# Crochemore's String Matching Algorithm: Simplification and Extensions

Juha Kärkkäinen     Dominik Kempa     Simon J. Puglisi

University of Helsinki, Finland

Prague Stringology Conference 2013

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
Our contribution

# Outline

**Introduction**
Simplified Crochemore's algorithm
Generalizations

**Problem**
Existing solutions
Our contribution

## Problem

### Exact String Matching

Given strings $T[1..n]$ and $P[1..m]$ find all positions $i$ such that $T[i..i + m - 1] = P[1..m]$.

**Introduction**
Simplified Crochemore's algorithm
Generalizations

**Problem**
Existing solutions
Our contribution

## Problem

### Exact String Matching

Given strings $T[1..n]$ and $P[1..m]$ find all positions $i$ such that $T[i..i + m - 1] = P[1..m]$.

- Algorithms available for many models and variants, e.g., compressed, higher-dimentions, "jumbled", multiple patterns

**Introduction**
Simplified Crochemore's algorithm
Generalizations

**Problem**
Existing solutions
Our contribution

## Problem

### Exact String Matching

Given strings $T[1..n]$ and $P[1..m]$ find all positions $i$ such that $T[i..i + m - 1] = P[1..m]$.

- Algorithms available for many models and variants, e.g., compressed, higher-dimentions, "jumbled", multiple patterns

- http://www.dmi.unict.it/~faro/smart/ lists over 80 algorithms **for the basic variant**. Many of them are time-optimal: $\mathcal{O}(n + m)$,

**Introduction**
Simplified Crochemore's algorithm
Generalizations

**Problem**
Existing solutions
Our contribution

## Problem

### Exact String Matching

Given strings $T[1..n]$ and $P[1..m]$ find all positions $i$ such that
$T[i..i + m - 1] = P[1..m]$.

- Algorithms available for many models and variants, e.g., compressed, higher-dimentions, "jumbled", multiple patterns
- http://www.dmi.unict.it/~faro/smart/ lists over 80 algorithms **for the basic variant**. Many of them are time-optimal: $\mathcal{O}(n + m)$,
- but only few are also space-optimal, that is, use $\mathcal{O}(1)$ machine words in addition to the input strings

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
**Existing solutions**
Our contribution

# Time-space optimal string matching

| Algorithm | Prepr. | Key idea |
|---|---|---|
| Rytter [2003] | yes | lex-maximal suffix |
| Gasieniec et al. [1995] | yes | zooming |
| Crochemore [1992] | no | lex-maximal suffix |
| Crochemore and Perrin [1991] | yes | critical factorization |
| Galil, Seiferas [1981] | yes | perfect factorization |

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

## Our contribution

1. Simplified version of Crochemore's algorithm

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

# Our contribution

1. Simplified version of Crochemore's algorithm
   - much simpler shifting rule

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

## Our contribution

1. Simplified version of Crochemore's algorithm
   - much simpler shifting rule
   - simpler analysis and extensions

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

## Our contribution

1. Simplified version of Crochemore's algorithm
   - much simpler shifting rule
   - simpler analysis and extensions
   - still time-space-optimal

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

## Our contribution

① Simplified version of Crochemore's algorithm
  - much simpler shifting rule
  - simpler analysis and extensions
  - still time-space-optimal

② Generalizations

**Introduction**
Simplified Crochemore's algorithm
Generalizations

Problem
Existing solutions
**Our contribution**

# Our contribution

1. Simplified version of Crochemore's algorithm
   - much simpler shifting rule
   - simpler analysis and extensions
   - still time-space-optimal
2. Generalizations
   - *pattern matching* → *longest prefix matching*

# Our contribution

1. Simplified version of Crochemore's algorithm
   - much simpler shifting rule
   - simpler analysis and extensions
   - still time-space-optimal
2. Generalizations
   - *pattern matching* → *longest prefix matching*
   - *longest prefix matching* → *sparse longest prefix matching*

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
Simplified algorithm description

# Outline

Introduction
**Simplified Crochemore's algorithm**
Generalizations

**Original vs. simplified**
Basic tools
Simplified algorithm description

# Original vs. simplified

## Original

```
Algorithm Crochemore
1:  z ← 0; c ← 1
2:  (i, j, k, p) ← (0, 1, 1, 1)
3:  while i ≤ n do
4:      while z + c ≤ n and c ≤ m and T[z + c] = P[m] do
5:          c ← c + 1
6:      if z + c = n + 1 or c = m + 1 then
7:          report(z)
8:      if z + c = n + 1 then
9:          c ← c − 1
10:     (i, j, k, p) ← next(P, T[z + c], i, j, k, p)
11:     if P[1..i] = suf(pref(P[i + 1..c − 1], p)) then
12:         if j − i > p then
13:             z ← z + p
14:             c ← c − p + 1
15:             j ← j − p
16:         else
17:             z ← z + p
18:             c ← c − p + 1
19:             (i, j, k, p) ← (0, 1, 1, 1)
20:     else
21:         z ← z + max(i, min(c − i, j)) + 1
22:         c ← 1
23:         (i, j, k, p) ← (0, 1, 1, 1)
```

## Simplified

```
Algorithm Crochemore
1:  i ← 1; ℓ ← 0
2:  (s, p) ← (0, 0)
2:  while i ≤ n − m do
3:      while T[i + ℓ] = P[1 + ℓ] do
4:          update(ℓ, s, p)
6:      if ℓ = m then
7:          report(i)
8:      if 3p ≤ j and P[1..s] = P[p..p + s] then
9:          i ← i + p
10:         j ← j − p
11:     else
12:         i ← i + ⌊j/3⌋ + 1
13:         j ← 0
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
**Basic tools**
Simplified algorithm description
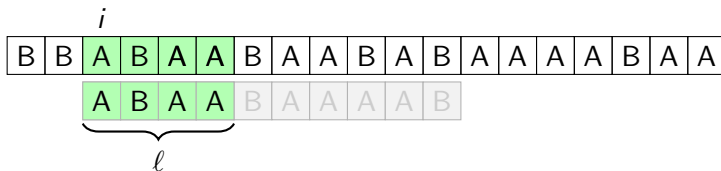
## Fundamental concept

### Definition

An integer $p$ is a period of string X if $p \leq |X|$ and X is a prefix of $X[1..p]^{\infty}$. The shortest period of X is denoted $per(X)$.

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
**Basic tools**
Simplified algorithm description

# Fundamental concept

### Definition

An integer $p$ is a period of string X if $p \leq |X|$ and X is a prefix of $X[1..p]^\infty$. The shortest period of X is denoted $per(X)$.

### Example

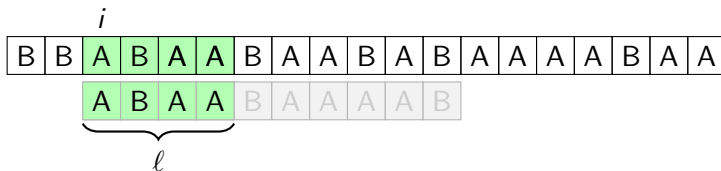Periods of $X = ABAABAAB$ are: 3, 6 and 8 thus $per(X) = 3$.

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



```
Algorithm Morris-Pratt
1: i ← 1; ℓ ← 0
2: while i ≤ n − m do
3:     while T[i + ℓ] = P[1 + ℓ] do
4:         ++ℓ
5:     if ℓ = m then report(i)
6:     Δ ← per[ℓ]            ▶ per[ℓ] = per(P[1..ℓ])
7:     ℓ′ ← ℓ − per[ℓ]
8:     i ← i + Δ; ℓ ← ℓ′
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



```
Algorithm Morris-Pratt
1: i ← 1; ℓ ← 0
2: while i ≤ n − m do
3:     while T[i + ℓ] = P[1 + ℓ] do
4:         ++ℓ
5:     if ℓ = m then report(i)
6:     Δ ← per[ℓ]           ▶ per[ℓ] = per(P[1..ℓ])
7:     ℓ' ← ℓ − per[ℓ]
8:     i ← i + Δ; ℓ ← ℓ'
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

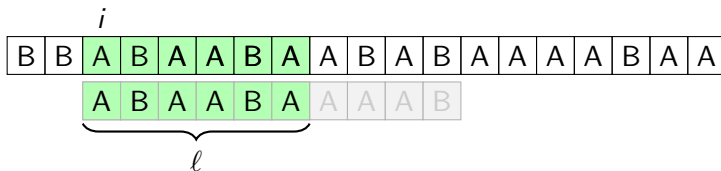# Basic principle



Algorithm Morris-Pratt
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** *report*$(i)$
6:    $\Delta \leftarrow per[\ell]$     ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

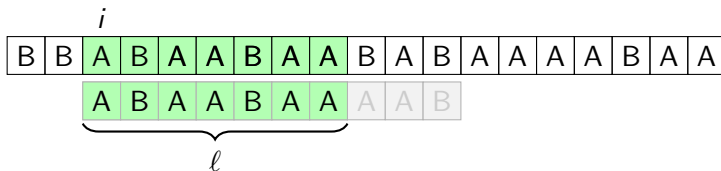# Basic principle



Algorithm Morris-Pratt
1: $i \leftarrow 1$; $\ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** *report*($i$)
6:    $\Delta \leftarrow per[\ell]$     ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

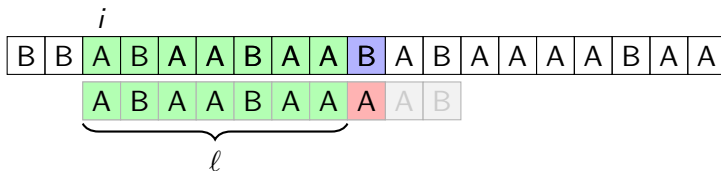# Basic principle



Algorithm Morris-Pratt
1: $i \leftarrow 1$; $\ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** *report*($i$)
6:    $\Delta \leftarrow per[\ell]$       ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



$i$

| B | B | A | B | A | A | B | A | A | B | A | B | A | A | A | A | B | A | A |

| A | B | A | A | B | A | A | A | A | B |

$\ell$

**Algorithm** `Morris-Pratt`
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** *report*$(i)$
6:    $\Delta \leftarrow per[\ell]$    ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta; \ell \leftarrow \ell'$
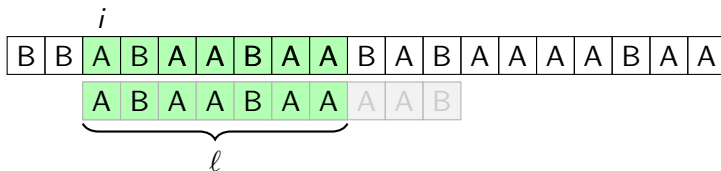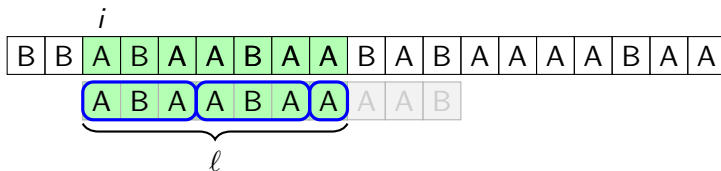
Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



```
Algorithm Morris-Pratt
1:  i ← 1; ℓ ← 0
2:  while i ≤ n − m do
3:      while T[i + ℓ] = P[1 + ℓ] do
4:          ++ℓ
5:      if ℓ = m then report(i)
6:      Δ ← per[ℓ]           ▶ per[ℓ] = per(P[1..ℓ])
7:      ℓ′ ← ℓ − per[ℓ]
8:      i ← i + Δ; ℓ ← ℓ′
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



**Algorithm** `Morris-Pratt`
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:   **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:     $++\ell$
5:   **if** $\ell = m$ **then** *report*$(i)$
6:   $\Delta \leftarrow per[\ell]$         ▶ $per[\ell] = per(P[1..\ell])$
7:   $\ell' \leftarrow \ell - per[\ell]$
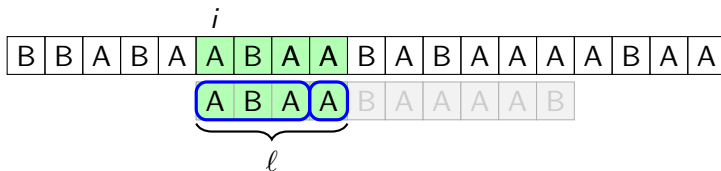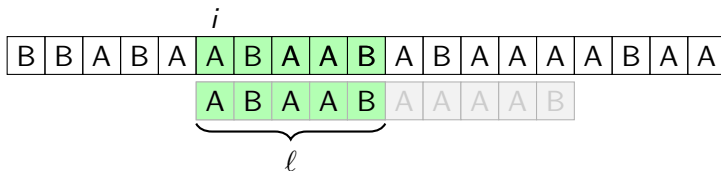8:   $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



$i$

| B | B | A | B | A | A | B | A | A | B | A | B | A | A | A | A | B | A | A |

| A | B | A | A | B | A | A | A | A | B |

$\ell$

**Algorithm** Morris-Pratt
1: $i \leftarrow 1;\ \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** $report(i)$
6:    $\Delta \leftarrow per[\ell]$    ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta;\ \ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



```
Algorithm Morris-Pratt
1: i ← 1; ℓ ← 0
2: while i ≤ n − m do
3:     while T[i + ℓ] = P[1 + ℓ] do
4:         ++ℓ
5:     if ℓ = m then report(i)
6:     Δ ← per[ℓ]          ▶ per[ℓ] = per(P[1..ℓ])
7:     ℓ′ ← ℓ − per[ℓ]
8:     i ← i + Δ; ℓ ← ℓ′
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



**Algorithm** `Morris-Pratt`
1: $i \leftarrow 1$; $\ell \leftarrow 0$
2: **while** $i \le n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** $report(i)$
6:    $\Delta \leftarrow per[\ell]$     ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
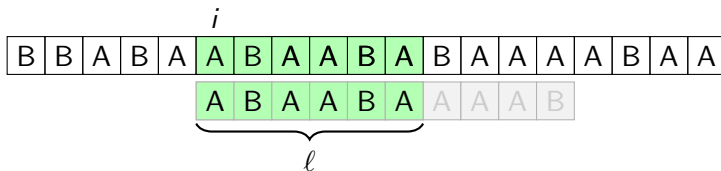8:    $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic principle



$i$

| B | B | A | B | A | A | B | A | A | B | A | B | A | A | A | A | B | A | A |

| A | B | A | A | B | A | A | A | A | B |

$\ell$

**Algorithm** `Morris-Pratt`
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$
5:    **if** $\ell = m$ **then** $report(i)$
6:    $\Delta \leftarrow per[\ell]$    ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

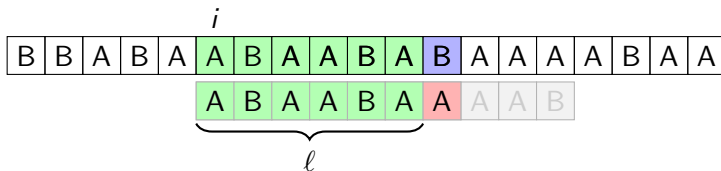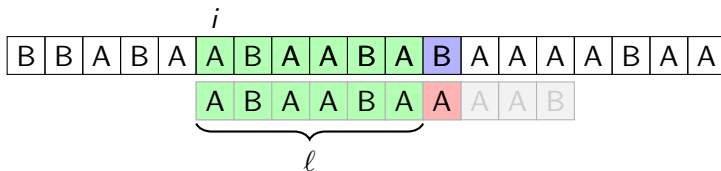# Basic principle



Algorithm Morris-Pratt
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \le n - m$ **do**
3:  **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:   $++\ell$
5:  **if** $\ell = m$ **then** *report*$(i)$
6:  $\Delta \leftarrow per[\ell]$       ▶ $per[\ell] = per(P[1..\ell])$
7:  $\ell' \leftarrow \ell - per[\ell]$
8:  $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

- Time: $\mathcal{O}(n + m)$ (each step increases the value of $i + \ell \le 2n$)

Introduction       Original vs. simplified
**Simplified Crochemore's algorithm**    Basic tools
Generalizations    **Simplified algorithm description**

# Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of $X$.

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

## Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of X.

### Example

```
A
A N A
A N A N A         X = B A N A N A
B A N A N A
N A
N A N A
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of $X$.

### Example

A
A N A
A N A N A                      $X = BA\boxed{NANA}$
B A N A N A
N A
$\boxed{NANA}$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of X.

### Example

A
A N A
A N A N A
B A N A N A
N A
N A N A

$$X = BA\,\boxed{NANA}$$
$$MaxSuf(X) = \quad NANA$$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

## Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of X.

### Definition

A pair of integers $(s, p)$ is called *MS-decomposition* of X if
$MaxSuf(X) = X[s..n]$ and $per(X[s..n]) = p$.

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

## Basic tools cont'd

### Definition

Let $MaxSuf(X)$ be the lex-maximal suffix of X.

### Definition

A pair of integers $(s, p)$ is called *MS-decomposition* of X if
$MaxSuf(X) = X[s..n]$ and $per(X[s..n]) = p$.

### Example

MS-decomposition of $X = BANANA$ is $(3, 2)$.

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- Only $\mathcal{O}(1)$ space is available $\Rightarrow$ no $per[1..n]$ precomputed

---

**Algorithm** Crochemore
1: $i \leftarrow 1; \ell \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:     **while** $\mathsf{T}[i + \ell] = \mathsf{P}[1 + \ell]$ **do**
4:         $++\ell$
5:     **if** $\ell = m$ **then** $report(i)$
6:     $\Delta \leftarrow per[\ell]$     ▶ $per[\ell] = per(\mathsf{P}[1..\ell])$
7:     $\ell' \leftarrow \ell - per[\ell]$
8:     $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

---

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- Only $\mathcal{O}(1)$ space is available $\Rightarrow$ no $per[1..n]$ precomputed

```
Algorithm Crochemore
1: i ← 1; ℓ ← 0
2: while i ≤ n − m do
3:     while T[i + ℓ] = P[1 + ℓ] do
4:         ++ℓ
5:     if ℓ = m then report(i)
6:     Δ ← per[ℓ]      ▶ per[ℓ] = per(P[1..ℓ])
7:     ℓ′ ← ℓ − per[ℓ]
8:     i ← i + Δ; ℓ ← ℓ′
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$

```
Algorithm Crochemore
1:  i ← 1; ℓ ← 0
2:  while i ≤ n − m do
3:      while T[i + ℓ] = P[1 + ℓ] do
4:          ++ℓ
5:      if ℓ = m then report(i)
6:      Δ ← per[ℓ]      ▶ per[ℓ] = per(P[1..ℓ])
7:      ℓ′ ← ℓ − per[ℓ]
8:      i ← i + Δ; ℓ ← ℓ′
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$

**Algorithm** Crochemore
1: $i \leftarrow 1$; $\ell \leftarrow 0$; $(s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:     **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:         $++\ell$
5:     **if** $\ell = m$ **then** *report*$(i)$
6:     $\Delta \leftarrow per[\ell]$     ▶ $per[\ell] = per(P[1..\ell])$
7:     $\ell' \leftarrow \ell - per[\ell]$
8:     $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$

**Algorithm** Crochemore
1: $i \leftarrow 1; \ell \leftarrow 0; (s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:     **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:         $++\ell$; update$(s, p)$          $\blacktriangleright \mathcal{O}(1)$ (amortized)
5:     **if** $\ell = m$ **then** *report*(i)
6:     $\Delta \leftarrow per[\ell]$     $\blacktriangleright per[\ell] = per(P[1..\ell])$
7:     $\ell' \leftarrow \ell - per[\ell]$
8:     $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$
- Recall that $p = per(MaxSuf(P[1..\ell]))$

```
Algorithm Crochemore
1:  i ← 1; ℓ ← 0; (s, p) ← (0, 0)
2:  while i ≤ n − m do
3:      while T[i + ℓ] = P[1 + ℓ] do
4:          ++ℓ; update(s, p)         ▶ O(1) (amortized)
5:      if ℓ = m then report(i)
6:      Δ ← per[ℓ]       ▶ per[ℓ] = per(P[1..ℓ])
7:      ℓ' ← ℓ − per[ℓ]
8:      i ← i + Δ; ℓ ← ℓ'
```

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$
- Recall that $p = per(MaxSuf(P[1..\ell]))$
    - Lemma: If $per(P[1..\ell]) \leq \ell/3$ then $per(P[1..\ell]) = p$

---

**Algorithm** Crochemore
1: $i \leftarrow 1; \ell \leftarrow 0; (s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell$; update$(s, p)$     ▶$\mathcal{O}(1)$ (amortized)
5:    **if** $\ell = m$ **then** report$(i)$
6:    $\Delta \leftarrow per[\ell]$   ▶ $per[\ell] = per(P[1..\ell])$
7:    $\ell' \leftarrow \ell - per[\ell]$
8:    $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

---

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$
- Recall that $p = per(MaxSuf(P[1..\ell]))$
  - ▶ Lemma: If $per(P[1..\ell]) \leq \ell/3$ then $per(P[1..\ell]) = p$

---

**Algorithm** Crochemore
1: $i \leftarrow 1;\ \ell \leftarrow 0;\ (s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell;\ update(s, p)$    ▶$\mathcal{O}(1)$ (amortized)
5:    **if** $\ell = m$ **then** $report(i)$
6:    **if** $per(P[1..\ell]) \leq \ell/3$ **then**  ▶$\mathcal{O}(1)$ (amortized)
7:       $\Delta \leftarrow p;\ \ell' \leftarrow \ell - p$    ▶ $per(P[1..\ell]) = p$
8:    $i \leftarrow i + \Delta;\ \ell \leftarrow \ell'$

---

Introduction
**Simplified Crochemore's algorithm**
Generalizations

Original vs. simplified
Basic tools
**Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$
- Recall that $p = per(MaxSuf(P[1..\ell]))$
  - ▶ Lemma: If $per(P[1..\ell]) \leq \ell/3$ then $per(P[1..\ell]) = p$

---

**Algorithm** Crochemore
1: $i \leftarrow 1; \ell \leftarrow 0;$ $(s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell;$ update$(s, p)$     ▶ $\mathcal{O}(1)$ (amortized)
5:    **if** $\ell = m$ **then** report$(i)$
6:    **if** $per(P[1..\ell]) \leq \ell/3$ **then**  ▶ $\mathcal{O}(1)$ (amortized)
7:       $\Delta \leftarrow p; \ell' \leftarrow \ell - p$    ▶ $per(P[1..\ell]) = p$
8:    **else**
9:       $\Delta \leftarrow \lfloor \ell/3 \rfloor + 1$      ▶ safe shift
10:      $(\ell', s, p) \leftarrow (0, 0, 0)$
11:   $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

---

Introduction    Original vs. simplified
**Simplified Crochemore's algorithm**    Basic tools
Generalizations    **Simplified algorithm description**

# Simplified algorithm overview

- We keep MS-decomposition $(s, p)$ of matching prefix $P[1..\ell]$
- Recall that $p = per(MaxSuf(P[1..\ell]))$
  - ▶ Lemma: If $per(P[1..\ell]) \leq \ell/3$ then $per(P[1..\ell]) = p$

---

**Algorithm** Crochemore

1: $i \leftarrow 1;\ \ell \leftarrow 0;\ (s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:       $++\ell;$ update$(s, p)$      ▶ $\mathcal{O}(1)$ (amortized)
5:    **if** $\ell = m$ **then** report$(i)$
6:    **if** $per(P[1..\ell]) \leq \ell/3$ **then**  ▶ $\mathcal{O}(1)$ (amortized)
7:       $\Delta \leftarrow p;\ \ell' \leftarrow \ell - p$    ▶ $per(P[1..\ell]) = p$
8:    **else**
9:       $\Delta \leftarrow \lfloor \ell/3 \rfloor + 1$       ▶ safe shift
10:      $(\ell', s, p) \leftarrow (0, 0, 0)$
11:    $i \leftarrow i + \Delta;\ \ell \leftarrow \ell'$

---

- Time: $\mathcal{O}(n + m)$ (each step increases the value of $3i + \ell \leq 4n$)

**Introduction**
**Simplified Crochemore's algorithm**
**Generalizations**

Motivation
Longest prefix matching
Sparse longest prefix matching

# Outline

Introduction
Simplified Crochemore's algorithm
**Generalizations**

**Motivation**
Longest prefix matching
Sparse longest prefix matching

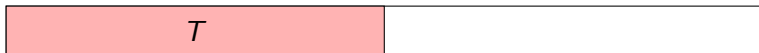## Motivation

- Scan-based lightweight LZ77 factorization:

Introduction
Simplified Crochemore's algorithm
**Generalizations**

**Motivation**
Longest prefix matching
Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
  ▶ text is stored on disk

Introduction
Simplified Crochemore's algorithm
**Generalizations**
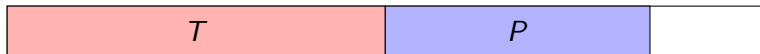**Motivation**
Longest prefix matching
Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
  - ▶ text is stored on disk
  - ▶ a prefix $T$ of text was already processed

| $T$ | |
|---|---|

Introduction
Simplified Crochemore's algorithm
**Generalizations**

**Motivation**
Longest prefix matching
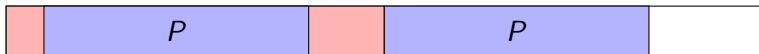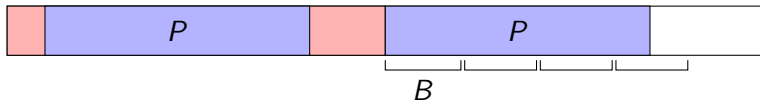Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
  - ▶ text is stored on disk
  - ▶ a prefix $T$ of text was already processed
  - ▶ goal: the length of the longest $P$ following $T$ that occurs twice in $TP$

| $T$ | $P$ | |
|:---:|:---:|:---:|

Introduction
Simplified Crochemore's algorithm
**Generalizations**

**Motivation**
Longest prefix matching
Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
    - ▶ text is stored on disk
    - ▶ a prefix $T$ of text was already processed
    - ▶ goal: the length of the longest $P$ following $T$ that occurs twice in $TP$

Introduction
Simplified Crochemore's algorithm
**Generalizations**
**Motivation**
Longest prefix matching
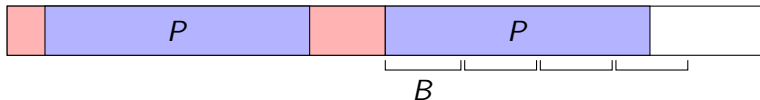Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
  - ▶ text is stored on disk
  - ▶ a prefix $T$ of text was already processed
  - ▶ goal: the length of the longest $P$ following $T$ that occurs twice in $TP$

- $P$ can have length $\gg B$ (where $B =$ memory budget)

Introduction
Simplified Crochemore's algorithm
Generalizations

Motivation
Longest prefix matching
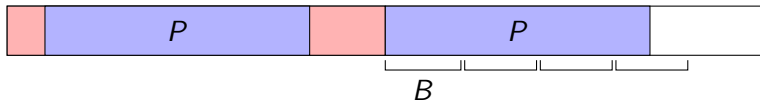Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
  - ▶ text is stored on disk
  - ▶ a prefix $T$ of text was already processed
  - ▶ goal: the length of the longest $P$ following $T$ that occurs twice in $TP$
- $P$ can have length $\gg B$ (where $B$ = memory budget)
- to keep the algorithm lightweight, we cannot index text

Introduction
Simplified Crochemore's algorithm
**Generalizations**
**Motivation**
Longest prefix matching
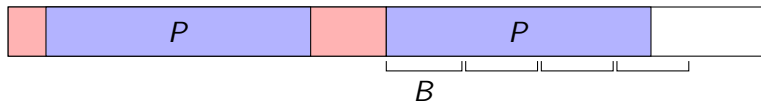Sparse longest prefix matching

## Motivation

- Scan-based lightweight LZ77 factorization:
    - ▶ text is stored on disk
    - ▶ a prefix $T$ of text was already processed
    - ▶ goal: the length of the longest $P$ following $T$ that occurs twice in $TP$
- $P$ can have length $\gg B$ (where $B$ = memory budget)
- to keep the algorithm lightweight, we cannot index text
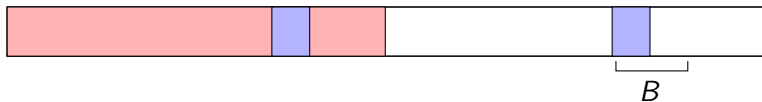- solution: search each block using string matching

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
**Longest prefix matching**
Sparse longest prefix matching

# Longest prefix matching

- if the block under consideration contains the phrase endpoint we need to compute its length

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
**Longest prefix matching**
Sparse longest prefix matching

# Longest prefix matching

- if the block under consideration contains the phrase endpoint we need to compute its length

Introduction
Simplified Crochemore's algorithm
**Generalizations**

Motivation
**Longest prefix matching**
Sparse longest prefix matching

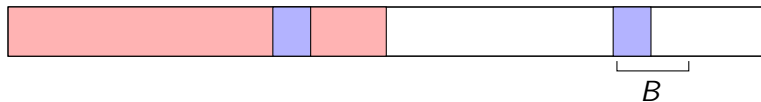# Longest prefix matching

- if the block under consideration contains the phrase endpoint we need to compute its length
  - ▶ *longest prefix matching* problem

Introduction
Simplified Crochemore's algorithm
Generalizations

Motivation
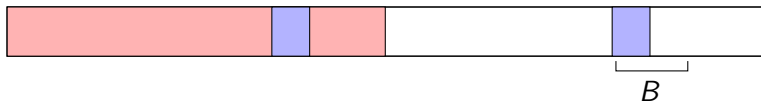**Longest prefix matching**
Sparse longest prefix matching

# Longest prefix matching

- if the block under consideration contains the phrase endpoint we need to compute its length
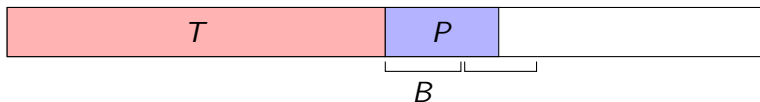  - ▶ *longest prefix matching* problem



**Algorithm** Morris-Pratt
1: $i \leftarrow 1$; $\ell \leftarrow 0$; $\ell_{max} \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:     **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:         $++\ell$
5:     **if** $\ell > \ell_{max}$ **then**
6:         $\ell_{max} \leftarrow \ell$
7:     $\Delta \leftarrow per[\ell]$; $\ell' \leftarrow \ell - per[\ell]$
8:     $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Sparse longest prefix matching

- searching for occurrences of prefixes of $2^{nd}$ block should be restricted to positions where occurrences of the $1^{st}$ block end

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
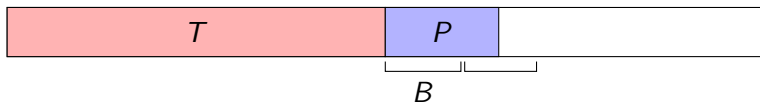**Sparse longest prefix matching**

# Sparse longest prefix matching

- searching for occurrences of prefixes of $2^{nd}$ block should be restricted to positions where occurrences of the $1^{st}$ block end
  - ▶ **sparse** *longest prefix matching* problem

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
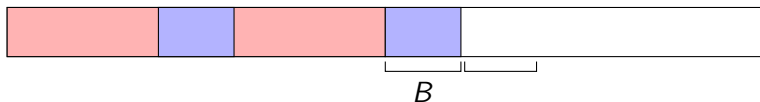**Sparse longest prefix matching**

## Sparse longest prefix matching

- searching for occurrences of prefixes of $2^{nd}$ block should be restricted to positions where occurrences of the $1^{st}$ block end
  - ▶ **sparse** *longest prefix matching* problem

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Sparse longest prefix matching

- searching for occurrences of prefixes of $2^{nd}$ block should be restricted to positions where occurrences of the $1^{st}$ block end
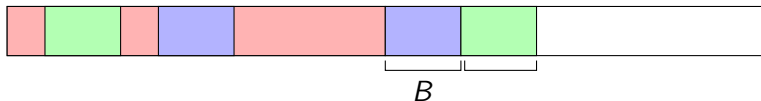  - ▶ **sparse** *longest prefix matching* problem

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

## Sparse longest prefix matching

- searching for occurrences of prefixes of 2$^{\text{nd}}$ block should be restricted to positions where occurrences of the 1$^{\text{st}}$ block end
  - ▶ **sparse** *longest prefix matching* problem

Introduction
Simplified Crochemore's algorithm
**Generalizations**

Motivation
Longest prefix matching
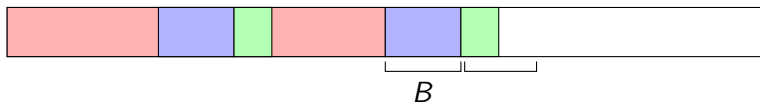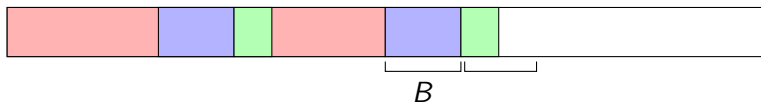**Sparse longest prefix matching**

# Sparse longest prefix matching

- searching for occurrences of prefixes of $2^{nd}$ block should be restricted to positions where occurrences of the $1^{st}$ block end
  - ▶ **sparse** *longest prefix matching* problem



**Algorithm** `Morris-Pratt`
1: $i \leftarrow 1; \ell \leftarrow 0; \ell_{max} \leftarrow 0$
2: **while** $i \leq n - m$ **do**
3:     **if** $i \in \mathcal{A}$ **then**  ▶ $\mathcal{A} =$ considered positions (ascending)
4:         **while** $T[i + \ell] = P[1 + \ell]$ **do**
5:             $++\ell$
6:     **if** $\ell > \ell_{max}$ **then**
7:         $\ell_{max} \leftarrow \ell$
8:     $\Delta \leftarrow per[\ell]; \ell' \leftarrow \ell - per[\ell]$
9:     $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Streaming issue

- in basic form of Crochemore's algorithm text is not streamed

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Streaming issue

- in basic form of Crochemore's algorithm text is not streamed

---

**Algorithm** Crochemore
1: $i \leftarrow 1$; $\ell \leftarrow 0$; $(s, p) \leftarrow (0, 0)$
2: **while** $i \leq n - m$ **do**
3:     **while** $T[i + \ell] = P[1 + \ell]$ **do**
4:         $++\ell$; update$(s, p)$
6:     **if** $\ell = m$ **then** report$(i)$
7:     **if** per$(P[1..\ell]) \leq \ell/3$
8:         $\Delta \leftarrow p$; $\ell' \leftarrow \ell - p$
9:     **else**
10:        $\Delta \leftarrow \lfloor \ell/3 \rfloor + 1$    ▶ This can cause the value
11:        $(\ell', s, p) \leftarrow (0, 0, 0)$  ▶ of $i + \ell$ to decrease
12:    $i \leftarrow i + \Delta$; $\ell \leftarrow \ell'$

---

Introduction
Simplified Crochemore's algorithm
**Generalizations**

Motivation
Longest prefix matching
**Sparse longest prefix matching**

## Streaming issue

- in basic form of Crochemore's algorithm text is not streamed

> **Algorithm** Crochemore
> 1: $i \leftarrow 1; \ell \leftarrow 0; (s, p) \leftarrow (0, 0)$
> 2: **while** $i \leq n - m$ **do**
> 3:    **while** $T[i + \ell] = P[1 + \ell]$ **do**
> 4:       $++\ell$; update$(s, p)$
> 6:    **if** $\ell = m$ **then** $report(i)$
> 7:    **if** $per(P[1..\ell]) \leq \ell/3$
> 8:       $\Delta \leftarrow p; \ell' \leftarrow \ell - p$
> 9:    **else**
> 10:      $\Delta \leftarrow \lfloor \ell/3 \rfloor + 1$   ▶ This can cause the value
> 11:      $(\ell', s, p) \leftarrow (0, 0, 0)$ ▶ of $i + \ell$ to decrease
> 12:    $i \leftarrow i + \Delta; \ell \leftarrow \ell'$

- Observation: every streamed text symbol is compared to pattern, i.e., pattern stores recent text history

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Conclusion

- minimalistic version of Crochemore's algorithm retaining all its key features

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

## Conclusion

- minimalistic version of Crochemore's algorithm retaining all its key features
- after several modifications very useful in lightweight computation of LZ77 parsing

Introduction
Simplified Crochemore's algorithm
**Generalizations**
Motivation
Longest prefix matching
**Sparse longest prefix matching**

# Conclusion

- minimalistic version of Crochemore's algorithm retaining all its key features
- after several modifications very useful in lightweight computation of LZ77 parsing
- prefix matching problems could be of independent interest

Introduction  Motivation
Simplified Crochemore's algorithm  Longest prefix matching
**Generalizations**  **Sparse longest prefix matching**

## Conclusion

# Thank you!