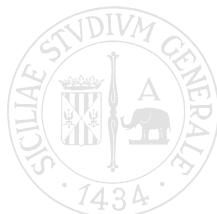# Efficient Online Abelian Pattern Matching in Strings by Simulating Reactive Multi-Automata

## Domenico Cantone and Simone Faro
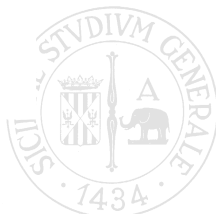
Department of Mathematics and Computer Science
**University of Catania (Italy)**

The Prague Stringology Conference 2014
Prague 1-4 August, 2014

**The Abelian Pattern Matching Problem**

Given a pattern p and a text t, the *abelian pattern matching* problem
(also known as *jumbled matching*) consists in finding all substrings of the
text t, whose characters have the same multiplicities as in p, so that they
could be converted into the input pattern just by permuting their
characters.

**The Abelian Pattern Matching Problem**

Given a pattern p and a text t, the *abelian pattern matching* problem (also known as *jumbled matching*) consists in finding all substrings of the text t, whose characters have the same multiplicities as in p, so that they could be converted into the input pattern just by permuting their characters.

---

when you listen, be silent but not be a tinsel

enlist

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The Abelian Pattern Matching Problem**

Given a pattern p and a text t, the *abelian pattern matching* problem
(also known as *jumbled matching*) consists in finding all substrings of the
text t, whose characters have the same multiplicities as in p, so that they
could be converted into the input pattern just by permuting their
characters.
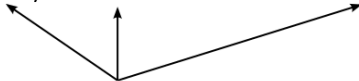
when you listen, be silent but not be a tinsel

enlist

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

## The counting filter

In the field of text processing and in computational biology, algorithms for abelian pattern matching are used as a filtering technique [Baeza-Yates.Navarro.2002]:

- $k$-mismatches [Grossi.Luccio.1989];
- $k$-differences [Jokinen.Tarhio.Ukkonen1996];
- inversions [Cantone.Cristofaro.Faro.2011];
- inversions and translocations [Cantone.Faro.Giaquinta.2011].

when you listen, be silent but not be a tinsel

enlist

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

## The counting filter

In the field of text processing and in computational biology, algorithms for abelian pattern matching are used as a filtering technique [Baeza-Yates.Navarro.2002]:

- $k$-mismatches [Grossi.Luccio.1989];
- $k$-differences [Jokinen.Tarhio.Ukkonen1996];
- inversions [Cantone.Cristofaro.Faro.2011];
- inversions and translocations [Cantone.Faro.Giaquinta.2011].

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

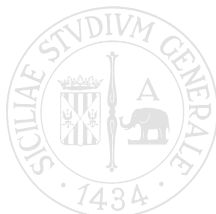The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The Parikh vector**

The *Parikh vector* of p is the vector of the multiplicities of the characters in p. More precisely, for each $c \in \Sigma$, we have

$$pv_\mathrm{p}[c] = |\{i : 0 \leq i < m \text{ and } \mathrm{p}[i] = c\}|\,.$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The Parikh vector**

The *Parikh vector* of p is the vector of the multiplicities of the characters in p. More precisely, for each $c \in \Sigma$, we have
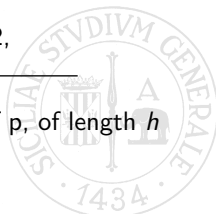
$$pv_p[c] = |\{i : 0 \le i < m \text{ and } p[i] = c\}|.$$

---

$$p = \text{enlist}$$
$$pv_p[e] = 1, \ pv_p[a] = 0, \ pv_p[t] = 1, \ pv_p[c] = 0$$

$$p = \text{stringology}$$
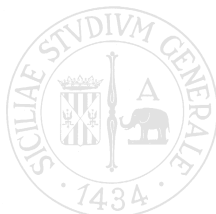$$pv_p[s] = 1, \ pv_p[o] = 2, \ pv_p[a] = 0, \ pv_p[g] = 2,$$

---

**Note**. The Parikh vector of the substring $p[i .. i + h - 1]$ of p, of length $h$ and starting at position $i$, will be denoted by $pv_{p(i,h)}$.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The Parikh vector**

In terms of Parikh vectors, the abelian pattern matching problem can be formally expressed as the problem of finding the set $\Gamma_{p,t}$ of positions in t, defined as

$$\Gamma_{p,t} = \{s : \ 0 \le s \le n - m \text{ and } pv_{t(s,m)} = pv_p\}.$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

## The naïve solution

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:1 | e:1 | i:1 | l:0 | n:1 | s:1 | t:1 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

### The naïve solution

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you li<span style="color:red">sten, be</span> silent but not be a tinsel

| b:1 | e:2 | i:0 | l:0 | n:1 | s:1 | t:1 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:1 | e:2 | i:0 | l:0 | n:1 | s:1 | t:1 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

**Introduction**
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
**Previous Results**
Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

when you list**en, be si**lent but not be a tinsel

| b:1 | e:2 | i:1 | l:0 | n:1 | s:1 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \left|\{c\} \cap \{t[s]\}\right| + \left|\{c\} \cap \{t[s+m]\}\right|,$$

when you listen, be silent but not be a tinsel

| b:1 | e:1 | i:1 | l:1 | n:1 | s:1 | t:0 |
|-----|-----|-----|-----|-----|-----|-----|
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

Introduction     The Abelian Pattern Matching Problem
A New Algorithm Based on Reactive Multi-Automata     Previous Results
An Efficient Bit-Parallel Simulation     Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:1 | e:2 | i:1 | l:1 | n:0 | s:1 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
**Previous Results**
Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:0 | e:2 | i:1 | l:1 | n:1 | s:1 | t:0 |
|-----|-----|-----|-----|-----|-----|-----|
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
**Previous Results**
Reactive Automata

**The naïve solution**

For a pattern p of length $m$ and a text t of length $n$ over an alphabet $\Sigma$ of size $\sigma$, the *abelian pattern matching problem* can be solved in $O(n)$ time and $O(\sigma)$ space by using a naïve *prefix based approach* [Ejaz.2010].

For each position $s = 0, 1, \ldots, n - m - 1$ and character $c \in \Sigma$, we have

$$pv_{t(s+1,m)}[c] = pv_{t(s,m)}[c] - \big|\{c\} \cap \{t[s]\}\big| + \big|\{c\} \cap \{t[s+m]\}\big|,$$

---

when you listen, be silent but not be a tinsel

| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

---

Another efficient prefix-based approach, which uses less branch conditions, has been proposed in [Grabowski.Faro.Giaquinta.2011]

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a tinsel

b:0  e:0  i:0  l:0  n:0  s:0  t:0

b:0  e:1  i:1  l:1  n:1  s:1  t:1

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

| when you listen, be silent b̲u̲t̲ not be a tinsel |
|---|

|  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|
| b:1 | e:0 | i:0 | l:0 | n:0 | s:0 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a tinsel

   b:0   e:0   i:0   l:0   n:0   s:0   t:0

   b:0   e:1   i:1   l:1   n:1   s:1   t:1

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a tinsel

| b:1 | e:0 | i:0 | l:0 | n:0 | s:0 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a tinsel

| b:0 | e:0 | i:0 | l:0 | n:0 | s:0 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a tinsel

| b:0 | e:0 | i:0 | l:0 | n:0 | s:1 | t:0 |
|-----|-----|-----|-----|-----|-----|-----|
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.
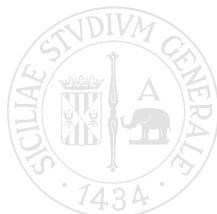
when you listen, be silent but not be a ti**ns**el

| b:0 | e:0 | i:0 | l:0 | n:1 | s:1 | t:0 |
|-----|-----|-----|-----|-----|-----|-----|
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

when you listen, be silent but not be a t**ins**el

| b:0 | e:0 | i:1 | l:0 | n:1 | s:1 | t:0 |
| b:0 | e:1 | i:1 | l:1 | n:1 | s:1 | t:1 |

**A suffix based solution**

A *suffix-based approach* to the problem has been proposed as an adaptation of the Horspool strategy [Ejaz.2010].

Characters are read from right to left. As soon as a frequency overflow occurs, the reading phase is stopped and a new alignment is attempted by sliding the window to the right. The resulting algorithm has an $O(nm)$ worst-case time complexity but performs well in practical cases.

---

when you listen, be silent but not be a tinsel

b:0   e:0   i:1   l:0   n:1   s:1   t:1

b:0   e:1   i:1   l:1   n:1   s:1   t:1

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation
The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**For more details ...**

A detailed analysis of the abelian pattern matching problem and of its
solutions is presented in [Ejaz.2010].

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
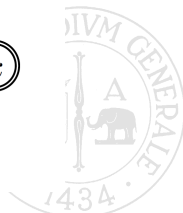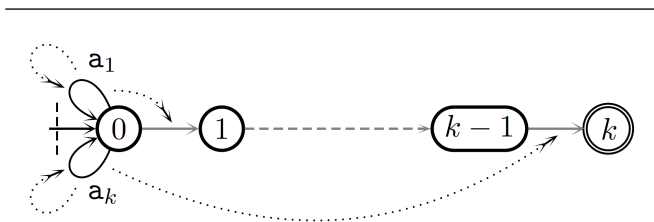Previous Results
Reactive Automata

**Reactive Automata**

A reactive automaton is an ordinary automaton extended with *reactive links* between its (ordinary) links. These can be of two types
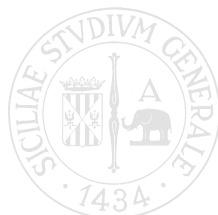
- *activation* reactive links;
- *deactivation* reactive links.

At any step of the computation of a reactive automaton on a given input string $S$, states and links are distinguished as active and non-active.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
Reactive Automata

**Reactive Automata**

A reactive automaton is an ordinary automaton extended with *reactive links* between its (ordinary) links. These can be of two types

- *activation* reactive links;
- *deactivation* reactive links.

At any step of the computation of a reactive automaton on a given input string $S$, states and links are distinguished as active and non-active.
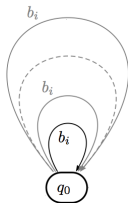
Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

The Abelian Pattern Matching Problem
Previous Results
**Reactive Automata**

**Reactive Automata**

A reactive automaton is an ordinary automaton extended with *reactive links* between its (ordinary) links. These can be of two types

- *activation* reactive links;
- *deactivation* reactive links.

At any step of the computation of a reactive automaton on a given input string $S$, states and links are distinguished as active and non-active.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

**Reactive Multi-Automata**

Reactive multi-automata extend reactive automata in that they allow the presence of multiple links labeled by a same character between any two states.
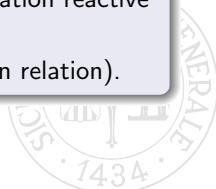
Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

### Reactive Multi-Automata

Reactive multi-automata extend reactive automata in that they allow the presence of multiple links labeled by a same character between any two states.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

### Definition (Reactive multi-automata)

Let $Q, \Sigma, L$ be finite sets of states, of characters, and of labels, respectively.

A *reactive multi-automaton* is a nonuple
$\mathcal{R} = \left( Q, \Sigma, L, q_0, \delta, \overline{\delta}, T^+, T^-, F \right)$, where

- $(Q, \Sigma, L, q_0, \delta, F)$ is a *multi-automaton* (called the *multi-automaton underlying* $\mathcal{R}$), with $q_0 \in Q$ (initial state), $F \subseteq Q$ (set of final states), and $\delta \subseteq Q \times \Sigma \times L \times Q$ (transition relation);
- $T^+, T^- \subseteq \delta \times \delta$ are the sets of activation and deactivation reactive links;
- $\overline{\delta} \subseteq \delta$ is the set of initially active links (initial transition relation).

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

**The Abelian Reactive Multi-Automaton**

Let p be a pattern of length $m$ over an alphabet $\Sigma$ and let $\langle b_0, b_1, \ldots, b_{k-1} \rangle$ be the sequence of the distinct characters occurring in p, ordered by their first occurrence. The *abelian reactive multi-automaton* (ARMA) for p is the reactive multi-automaton with $\varepsilon$-transitions

$$\mathcal{R} = \left( Q, \Sigma, L, q_0, \delta, \overline{\delta}, T^+, T^-, F \right)$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

### The Set of States of the Automaton

- $Q = \{q_0, q_1, \ldots, q_k, \omega\}$ is the set of states;
- $q_0$ is the initial state; $\omega$ is a special state called the *overflow state*;
- $F = \{q_k\}$ is the set of final states;

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## The Full Transition Relation

the transition relation $\delta$ of $\mathcal{R}$ and its subset $\overline{\delta} \subseteq \delta$ of the links initially active (initial transition relation) are defined as follows
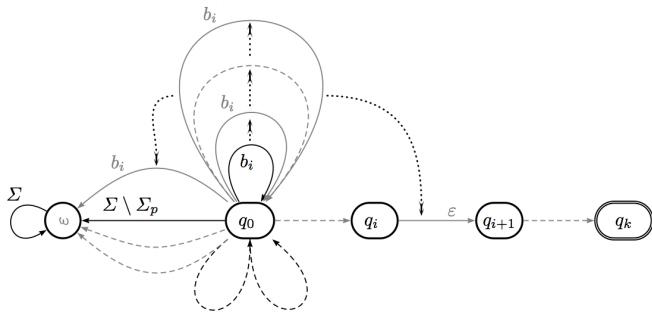
$$
\begin{aligned}
\delta = \quad & \{(q_i, \varepsilon, \ell_0, q_{i+1}) \mid 0 \leq i < k\} && (\varepsilon\text{-transitions}) \\
& \cup \{(q_0, p[i], \ell_i, q_0) \mid 0 \leq i < m\} && (\text{self-loops}) \\
& \cup \{(q_0, c, \ell_0, \omega) \mid c \in \Sigma\} && (\text{overflow transitions})
\end{aligned}
$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## The Initial Transition Relation

the transition relation $\delta$ of $\mathcal{R}$ and its subset $\overline{\delta} \subseteq \delta$ of the links initially active (initial transition relation) are defined as follows

$$\overline{\delta} = \quad \{(q_0, c, \ell_{\lambda(c)}, q_0) \mid c \in \Sigma_\mathsf{p}\}$$
$$\cup \{(q_0, c, \ell_0, \omega) \mid c \in \Sigma \setminus \Sigma_\mathsf{p}\}$$
$$\cup \{(\omega, c, \ell_0, \omega) \mid c \in \Sigma\}$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## The Set of Activation Reactive Links

The set $T^+$ of activation reactive links are defined as follows

$$
\begin{aligned}
T^+ = \ & \{((q_0, \mathsf{p}[\rho(b_i)], \ell_{\rho(b_i)}, q_0), (q_i, \varepsilon, \ell_0, q_{i+1})) \mid 0 \le i < k\} \\
& \cup \{((q_0, \mathsf{p}[\rho(b_i)], \ell_{\rho(b_i)}, q_0), (q_0, \mathsf{p}[\rho(b_i)], \ell_{\rho(b_i)}, \omega)) \mid 0 \le i < k\} \\
& \cup \{((q_0, \mathsf{p}[i], \ell_i, q_0), (q_0, \mathsf{p}[\nu(i)], \ell_{\nu(i)}, q_0)) \mid 0 \le i < m \text{ and } i \ne \rho(p_i)\}
\end{aligned}
$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
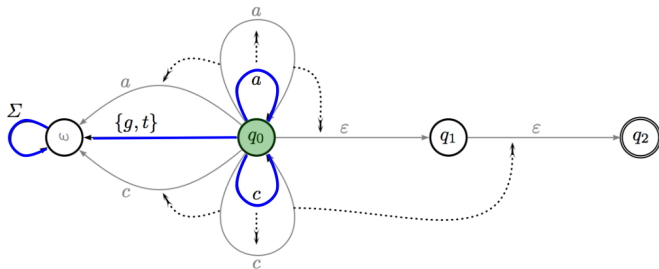The Algorithm

**The abelian reactive automaton for** *acca*

The complete abelian reactive automaton for the pattern $P = acca$ over
the DNA alphabet $\Sigma = \{a, c, g, t\}$. Standard transitions are represented
with solid lines while reactive links in $T^+$ are represented with dashed
lines. Reactive links in $T^-$ are not represented. Non active transitions
are represented in gray color.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automaton
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

# An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automaton
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automaton
The Abelian Reactive Multi-Automaton
The Algorithm

## An example

p = acca
t = gtcaaaccgtacgagtacgat...

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Reactive Multi-Automata
The Abelian Reactive Multi-Automaton
The Algorithm

# An example

p = acca
t = gtca<span style="color:red">aacc</span>gtacgagtacgat...

Introduction       Setting the Simulation
A New Algorithm Based on Reactive Multi-Automata     The Algorithm
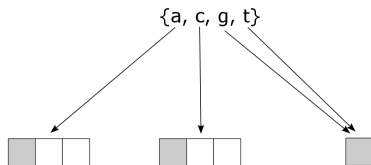An Efficient Bit-Parallel Simulation     Experimental Results

## A Bit-Parallel Simulation

The underlying idea is to associate a counter to each distinct character in p, plus a single 1-bit counter for the remaining characters of the alphabet which do not occur in p, maintaining them in the same computer word.

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## A Bit-Parallel Simulation

The underlying idea is to associate a counter to each distinct character in p, plus a single 1-bit counter for the remaining characters of the alphabet which do not occur in p, maintaining them in the same computer word.

The counter associated to the character $b_i$ in p is represented by a group of $l_i$ bits, where $l_i = \lceil \log(pv_p[b_i]) + 1 \rceil + 1$.
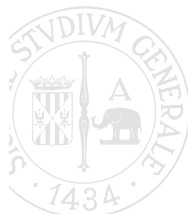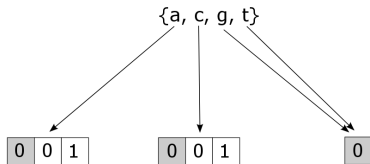


p = acca

{a, c, g, t}

## A Bit-Parallel Simulation

- Initially, the counter for $b_i$ is set to the value $2^{l_i} - pv_p[b_i] - 1$, so that its overflow bit is 0 and it remains so for up to $pv_p[b_i]$ increments.

- The overflow bit gets set only when the $(pv_p[b_i] + 1)$-st occurrence of $b_i$ is encountered in the text window.

p = acca

{a, c, g, t}

| 0 | 0 | 1 |

| 0 | 0 | 1 |

| 0 |

## A Bit-Parallel Simulation

- Initially, the counter for $b_i$ is set to the value $2^{l_i} - pv_p[b_i] - 1$, so that its overflow bit is 0 and it remains so for up to $pv_p[b_i]$ increments.

- The overflow bit gets set only when the $(pv_p[b_i] + 1)$-st occurrence of $b_i$ is encountered in the text window.
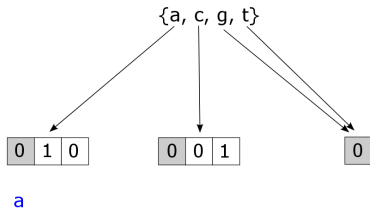
Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## A Bit-Parallel Simulation

- Initially, the counter for $b_i$ is set to the value $2^{l_i} - pv_p[b_i] - 1$, so that its overflow bit is 0 and it remains so for up to $pv_p[b_i]$ increments.

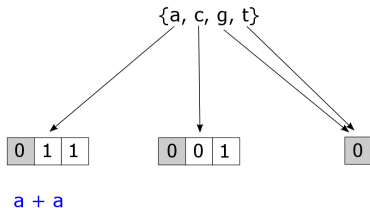- The overflow bit gets set only when the $(pv_p[b_i] + 1)$-st occurrence of $b_i$ is encountered in the text window.

## A Bit-Parallel Simulation

- Initially, the counter for $b_i$ is set to the value $2^{l_i} - pv_p[b_i] - 1$, so that its overflow bit is 0 and it remains so for up to $pv_p[b_i]$ increments.

- The overflow bit gets set only when the $(pv_p[b_i] + 1)$-st occurrence of $b_i$ is encountered in the text window.
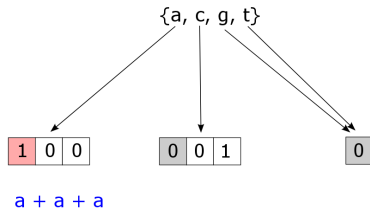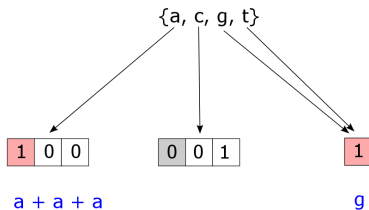
p = acca

{a, c, g, t}

| 1 | 0 | 0 |

| 0 | 0 | 1 |

| 0 |

a + a + a

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## A Bit-Parallel Simulation

- Likewise, the 1-bit counter reserved for all the characters not occurring in p is initially null and it gets set as soon as any character not in p is encountered in the text window.

p = acca

{a, c, g, t}

| 1 | 0 | 0 |

| 0 | 0 | 1 |

| 1 |

a + a + a                                  g

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

### The Preprocessing Phase

for each distinct character $b_i$ occurring in p we compute, a bit mask
$M[b_i]$ of $l + 1$ bits is computed, where

$$l = \sum_{i=0}^{k-1} l_i \qquad \text{and} \qquad M[b_i] = 1 \ll \Big( \sum_{j=0}^{i-1} l_j \Big).$$

The bit mask $M[b_i]$ is then used to increment the counter in $D$
associated to the character $b_i$.

---

p = acca

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
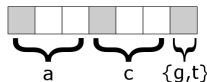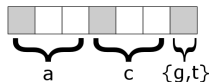The Algorithm
Experimental Results

## The Preprocessing Phase

for each distinct character $b_i$ occurring in p we compute, a bit mask $M[b_i]$ of $l + 1$ bits is computed, where

$$l = \sum_{i=0}^{k-1} l_i \qquad \text{and} \qquad M[b_i] = 1 \ll \Big( \sum_{j=0}^{i-1} l_j \Big).$$

The bit mask $M[b_i]$ is then used to increment the counter in $D$ associated to the character $b_i$.

p = acca

| M[a] | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
|------|---|---|---|---|---|---|---|
| M[c] | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| M[g] | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

a    c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
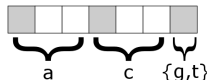The Algorithm
Experimental Results

### The Preprocessing Phase

Two additional bit masks are used:

- the bit mask $I$, which contains the initial values for each counter

$$I = \sum_{i=0}^{k-1} \left[ \left( 2^{l_i} - pv_p[b_i] - 1 \right) \ll \sum_{j=0}^{i-1} l_j \right]$$

$p = acca$

| M[a] | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
|------|---|---|---|---|---|---|---|

| M[c] | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
|------|---|---|---|---|---|---|---|

| M[g] | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
|------|---|---|---|---|---|---|---|

| I | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

$$\underbrace{\phantom{XXX}}_{a} \underbrace{\phantom{XXX}}_{c} \underbrace{\phantom{X}}_{\{g,t\}}$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

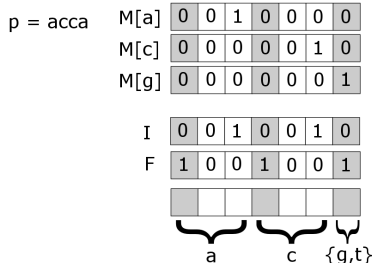Setting the Simulation
The Algorithm
Experimental Results

### The Preprocessing Phase

Two additional bit masks are used:

- the bit mask $F$, whose bits set are exactly the overflow bits.

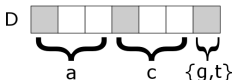$$F = \sum_{i=0}^{k-1} \left[ 1 \ll \left( \sum_{j=0}^{i} l_j - 1 \right) \right]$$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

p = acca
t = gtcaaaccgtacgagtacgat...

D

a      c      {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

### The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
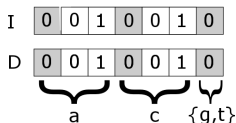
---

p = acca
t = gtcaaaccgtacgagtacgat...

| I | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

| D | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

    a       c     {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l+1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
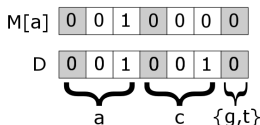
p = acca
t = gtcaaaccgtacgagtacgat...

M[a] | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

D | 0 | 0 | 1 | 0 | 0 | 1 | 0 |

a        c        {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
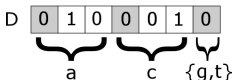
---

p = acca
t = gtcaaaccgtacgagtacgat...

D | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

  a  c {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
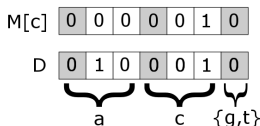
p = acca
t = gtcaaaccgtacgagtacgat...

| M[c] | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
|------|---|---|---|---|---|---|---|

| D | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

a     c     {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

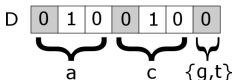Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l+1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

---

p = acca
t = gtcaaaccgtacgagtacgat...

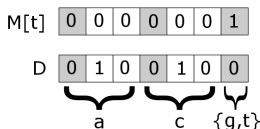$D$ | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

a · · · c · · {g,t}

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $l$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[\text{t}[j]]$.

---

p = acca
t = gtcaaaccgtacgagtacgat...

M[t]  | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

D     | 0 | 1 | 0 | 0 | 1 | 0 | 0 |

    a      c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

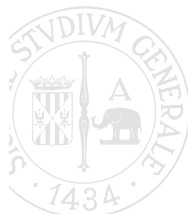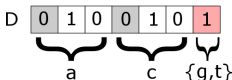- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

```
p = acca
t = gtcaaaccgtacgagtacgat...
    └──┘
```

D | 0 | 1 | 0 | 0 | 1 | 0 | 1

      a      c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[\text{t}[j]]$.
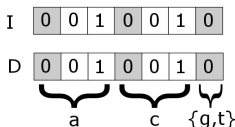
p = acca
t = gtcaaaccgtacgagtacgat...

I | 0 | 0 | 1 | 0 | 0 | 1 | 0 |

D | 0 | 0 | 1 | 0 | 0 | 1 | 0 |

a      c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
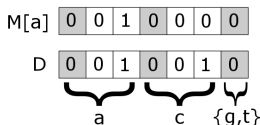
---

p = acca
t = gtcaaaccgtacgagtacgat...

| M[a] | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
|------|---|---|---|---|---|---|---|

| D | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

a     c     {g,t}

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
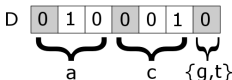
p = acca
t = gtcaaaccgtacgagtacgat...

D $\boxed{0}$ $\boxed{1}$ $\boxed{0}$ $\boxed{0}$ $\boxed{0}$ $\boxed{1}$ $\boxed{0}$

$\underbrace{\qquad}_{a}$ $\underbrace{\qquad\qquad}_{c}$ $\underbrace{\quad}_{\{g,t\}}$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
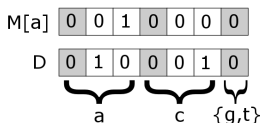
p = acca
t = gtcaaaccgtacgagtacgat...

M[a] | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

D | 0 | 1 | 0 | 0 | 0 | 1 | 0 |

a          c       {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
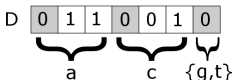
---

p = acca
t = gtcaaaccgtacgagtacgat...

D  | 0 | 1 | 1 | 0 | 0 | 1 | 0 |

a    c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
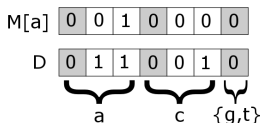
---

p = acca
t = gtcaaaccgtacgagtacgat...

| $M[a]$ | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|

| $D$ | 0 | 1 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

a     c     {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character t[$j$] of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[\text{t}[j]]$.
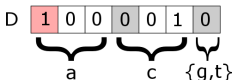
p = acca
t = gtcaaaccgtacgagtacgat...

D  | 1 | 0 | 0 | 0 | 0 | 1 | 0 |

a      c      {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

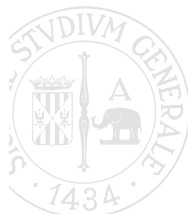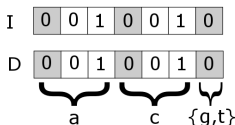- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

---

p = acca
t = gtcaaaccgtacgagtacgat...

| I | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

| D | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
|---|---|---|---|---|---|---|---|

a    c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $l$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
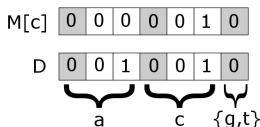
---

p = acca
t = gtcaaaccgtacgagtacgat...

M[c] | 0 | 0 | 0 | 0 | 0 | 1 | 0

D | 0 | 0 | 1 | 0 | 0 | 1 | 0

a    c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

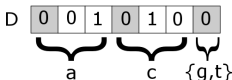- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

p = acca
t = gtcaaaccgtacgagtacgat...

D  | 0 | 0 | 1 | 0 | 1 | 0 | 0 |

a    c    {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
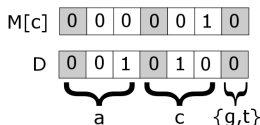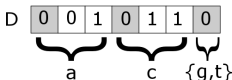
---

p = acca
t = gtcaaaccgtacgagtacgat...

| M[c] | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
|------|---|---|---|---|---|---|---|

| D | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
|---|---|---|---|---|---|---|---|

      a        c     {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

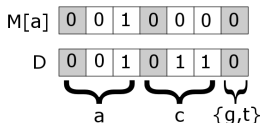- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

---

```
p = acca
t = gtcaaaccgtacgagtacgat...
```

$D$ | 0 | 0 | 1 | 0 | 1 | 1 | 0 |

a          c      {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.

- Then, during the attempt, the window is read character by character, proceeding from right to left.

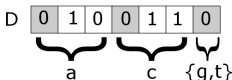- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

---

$p = \text{acca}$
$t = \text{gtca}\underbrace{\text{aacc}}\text{gtacgagtacgat...}$

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $M[a]$ | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $D$ | 0 | 0 | 1 | 0 | 1 | 1 | 0 |

$\underbrace{\phantom{00}}_{a} \quad \underbrace{\phantom{000}}_{c} \quad \underbrace{\phantom{0}}_{\{g,t\}}$

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
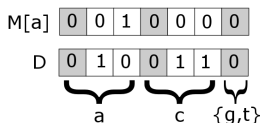
p = acca
t = gtcaaaccgtacgagtacgat...

D  | 0 | 1 | 0 | 0 | 1 | 1 | 0 |

a        c      {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
The Algorithm
Experimental Results

## The Searching Phase

- At the beginning of each attempt, a bit mask $D$ of $l+1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.
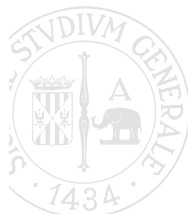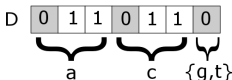
p = acca
t = gtcaaaccgtacgagtacgat...

M[a]  | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

D  | 0 | 1 | 0 | 0 | 1 | 1 | 0 |

a        c      {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

Setting the Simulation
**The Algorithm**
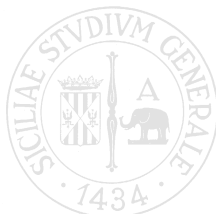Experimental Results

**The Searching Phase**

- At the beginning of each attempt, a bit mask $D$ of $l + 1$ bits is initialized to $I$.
- Then, during the attempt, the window is read character by character, proceeding from right to left.
- When reading the character $t[j]$ of the text, the bit mask $D$ is updated accordingly by setting it to $D + M[t[j]]$.

---

p = acca
t = gtcaaaccgtacgagtacgat...

D | 0 | 1 | 1 | 0 | 1 | 1 | 0 |

a          c          {g,t}

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

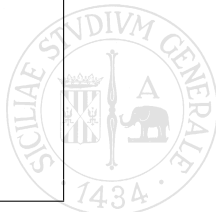Setting the Simulation
The Algorithm
Experimental Results

**The Searching Phase**

When $l + 1 > w$, we must content ourselves to maintain the counters only for a proper selection $\Sigma'_p$ of the set of characters occurring in p. In this case, when a match relative to the characters in $\Sigma'_p$ is reported, an additional verification phase must be run, in order to discard possible *false positives*.

Introduction
Setting the Simulation

A New Algorithm Based on Reactive Multi-Automata
The Algorithm

An Efficient Bit-Parallel Simulation
Experimental Results

```
BAM(p, m, t, n, Σ)
 1    for each c ∈ Σ do M[c] ← pv_p[c] ← 0
 2    I ← F ← sh ← 0
 3    for i ← 0 to m − 1 do pv_p[p[i]] ← pv_p[p[i]] + 1
 4    for each c ∈ Σ do
 5        if pv_p[c] > 0 then
 6            M[c] ← M[c] | (1 ≪ sh)
 7            I ← I | (((1 ≪ log m) − pv_p[c] − 1) ≪ sh)
 8            F ← F | (1 ≪ (sh + log m))
 9            sh ← sh + log m + 1
10    F ← F | (1 ≪ sh)
11    for each c ∈ Σ do
12        if pv_p[c] = 0 then M[c] ← M[c] | (1 ≪ sh)
13    s ← 0
14    while s ≤ n − m do
15        D ← I; j ← s + m − 1
16        while j ≥ s do
17            D ← D + M[t[j]]
18            if (D & F) then break
19            j ← j − 1
20        if j < s then
21            OUTPUT(s)
22            s ← s + 1
23        else s ← j + 1
```

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation

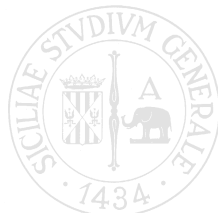Setting the Simulation
The Algorithm
Experimental Results

**Experimental Results**

In this section we evaluate the performance of the bit-parallel simulation BAM described in the previous section and compare it with some standard solutions known in literature.

We compare the performances of the following algorithms:

- The prefix based algorithm due to Grabowsky *et al.* (GFG);
- The algorithm using the suffix based approach (SBA),
- The Bit-parallel Abelian Matcher (BAM) described in this paper.

Introduction     Setting the Simulation
A New Algorithm Based on Reactive Multi-Automata     The Algorithm
An Efficient Bit-Parallel Simulation     Experimental Results
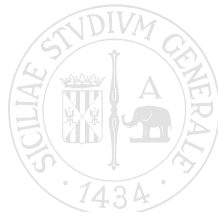
**Experimental Results**

In this section we evaluate the performance of the bit-parallel simulation BAM described in the previous section and compare it with some standard solutions known in literature.

We compare the performances of the following algorithms:

- The prefix based algorithm due to Grabowsky *et al.* $(\text{GFG})$;
- The algorithm using the suffix based approach $(\text{SBA})$,
- The Bit-parallel Abelian Matcher $(\text{BAM})$ described in this paper.

We use the following text buffers

- a genome sequence (alphabet of 4 characters);
- a protein sequence (alphabet of 20 characters);

Introduction
A New Algorithm Based on Reactive Multi-Automata
An Efficient Bit-Parallel Simulation
Setting the Simulation
The Algorithm
Experimental Results

## Experimental Results

| $m$ | GFG | SBA | BAM | | $m$ | GFG | SBA | BAM |
|---|---|---|---|---|---|---|---|---|
| 2 | **23.56** | 39.20 | 27.03 | | 2 | 23.08 | **18.07** | 12.51 |
| 4 | **23.56** | 33.27 | 23.17 | | 4 | 23.00 | **15.39** | 10.36 |
| 8 | 23.54 | 27.54 | **19.01** | | 8 | 22.96 | 13.67 | **9.40** |
| 16 | 23.49 | 24.05 | **16.21** | | 16 | 23.03 | 11.91 | **8.44** |
| 32 | 23.52 | 23.78 | **15.63** | | 32 | 23.04 | 9.58 | **7.16** |
| 64 | 23.50 | 25.33 | **16.12** | | 64 | 23.01 | 8.46 | **6.64** |
| 128 | 23.57 | 28.74 | **17.69** | | 128 | 22.97 | 7.82 | **6.49**$^{*}$ |
| 256 | 23.53 | 33.14 | **19.63** | | 256 | 22.96 | 7.84 | **7.69**$^{*}$ |

Tabella: Experimental results on a genome sequence (on the left) and a on a protein sequence (on the right). An asterisk symbol ($^{*}$) indicates those runs where false positives have been detected. All best results have been boldfaced.