

# ***Practical and Optimal String Matching***

**Kimmo Fredriksson**

Department of Computer Science, University of Joensuu, Finland

**Szymon Grabowski**

Technical University of Łódź, Computer Engineering Department

# Problem Setting

---

- The classic string matching problem:

Given *text*  $T[1..n]$  and *pattern*  $P[1..m]$  over some finite alphabet  $\Sigma$  of size  $\sigma$ , find the occurrences of  $P$  in  $T$ .

- We focus on the case where  $m$  is relatively small

⇒

Bit-parallelism.

Vast number of algorithms exist. Some of the most well-known are (classics):

**Knuth-Morris-Prat:** The first  $O(n)$  worst case time algorithm.

**Boyer-Moore(-Horspool)-family:** Numerous variants, sublinear on average.

(bit-parallel:)

**Shift-or:**  $O(n)$  for  $m \leq w$  (Baeza-Yates & Gonnet, 1992).

**BNDM family:**  $O(n \log_{\sigma}(m)/m)$  on average for  $m \leq w$ .

SBNDM (Navarro, 2001; Peltola & Tarhio, 2003), LNDM (He & Fang, 2004), FNNDM (Holub & Durian, 2005).

- In practice, the best current algorithms for short patterns are the BNDM-family of algorithms (Navarro & Raffinot, 2000).

- We develop a novel pattern partitioning technique that allows us to use shift-or while skipping text characters.
- The algorithm has optimal  $O(n \log_{\sigma}(m)/m)$  average case running time if  $m \leq w$ .
- Very simple to implement, simple inner loop (comparable to plain shift-or)

⇒

very efficient in practice.

- $O(mn)$  worst case, but can be improved to  $O(n)$  without destroying the simplicity of the search algorithm.

# *Our algorithm: the idea*

---

The algorithm is based on the preprocessing / filtering / verification paradigm.

- The preprocessing phase generates  $q$  different alignements of the pattern, each containing only every  $q$ th pattern character.

I.e. we partition the pattern into  $q$  pieces.

- The filtering phase searches all the  $q$  pieces in parallel using shift-or algorithm, reading only every  $q$ th text character.
- If any of the  $q$  pieces match, then we invoke a verification algorithm.

- Given a pattern  $P$ , generate a set  $\mathcal{P} = \{P^0, \dots, P^{q-1}\}$  of  $q$  patterns as follows:

$$P^j[i] = P[j + iq], \quad j = 0 \dots q - 1, \quad i = 0 \dots \lfloor m/q \rfloor - 1.$$

- I.e. we generate  $q$  different alignments of the original pattern  $P$ , each alignment containing only every  $q$ th character.
- Each new pattern has length  $m' = \lfloor m/q \rfloor$ .
- The total length of the patterns is  $q \lfloor m/q \rfloor \leq m$ .
- For example, if  $P = \text{abcdef}$  and  $q = 3$ , then

$$P^0 = \text{ad}, \quad P^1 = \text{be} \quad \text{and} \quad P^2 = \text{cf}.$$

# Preprocessing: the rationale

---

Assume that  $P$  occurs at  $T[i..i + m - 1]$ .

⇒

$$P^j[h] = T[i + j + hq], \quad j = i \bmod q, \quad h = 0 \dots m' - 1.$$

⇒

(1) We can use the set  $\mathcal{P}$  as a filter for the pattern  $P$

(2) The filter needs to scan only every  $q$ th character of  $T$ .



# Preprocessing: the rationale

$P$ 

a	b	c	d	e	f
---	---	---	---	---	---

$T$ 

	x	x	a	b	c	d	e	f	x	x	x	
--	---	---	---	---	---	---	---	---	---	---	---	--

$P^0$ 

a	d
---	---

$P^1$ 

b	e
---	---

$P^2$ 

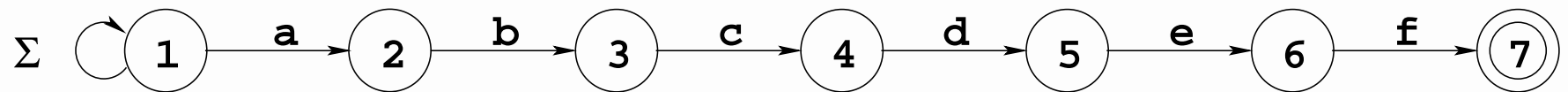
c	f
---	---

$P'$ 

a	d	b	e	c	f
---	---	---	---	---	---

# Prelude to filtering: Shift-or algorithm

- The algorithm is based on a non-deterministic automaton. The automaton for  $P = \text{'abcdef}$  is:



- The transitions are encoded in a table  $B$  of bit-masks: For  $1 \leq i \leq m$ , the mask  $B[c]$  has the  $i$ th bit set to 0, iff  $P[i] = c$ .
- The bit-vector  $D$  has one bit per state in the automaton, the  $i$ th bit of the vector is set to 0, iff the state  $i$  is active (initially all bits are 1).
- It can be shown that the automaton can be simulated as:

$$D \leftarrow (D \ll 1) | B[T[i]]$$

## ***Prelude to filtering: Shift-or algorithm***

---

- If after the simulation step, the  $m$ th bit of  $D$  is zero, then  $P$  occurs at  $T[i - m + 1 \dots i]$ .

Can be detected as

$$(D \& mm) \neq mm$$

where  $mm$  has only the  $m$ th bit set.

- Clearly each step of the automaton is simulated in time  $O(\lceil m/w \rceil)$ , which leads to  $O(n \lceil m/w \rceil)$  total time.

- The whole set  $\mathcal{P}$  of patterns can be searched simultaneously using the Shift-or algorithm (Baeza-Yates & Gonnet, 1992).
- All the patterns are preprocessed together, as if they were concatenated:  
For  $P = abcdef$ , we effectively preprocess a pattern  $P' = P^0 P^1 P^2 = adbecf$ .
- If the pattern  $P^j$  matches, then the  $(j + 1)m'$ -th bit in  $D$  is zero. This can be detected as

$$(D \& mm) \neq mm$$

where  $mm$  has every  $(j + 1)m'$ -th bit set to 1.

# Filtering: the simplicity illustrated

- Plain shift-or search:

```
1       $D \leftarrow \sim 0; i \leftarrow 0$ 
2      while  $i < n$  do
3           $D \leftarrow (D \ll 1) | B[T[i]]$ 
4          if  $(D \& mm) \neq mm$  then report match
5           $i \leftarrow i + 1$ 
```

- Our shift-or search:

```
1       $D \leftarrow \sim 0; i \leftarrow 0$ 
2      while  $i < n$  do
3           $D \leftarrow ((D \& \sim mm) \ll 1) | B[T[i]]$ 
4          if  $(D \& mm) \neq mm$  then Verify
5           $i \leftarrow i + q$ 
```

- If any of the pattern pieces in  $\mathcal{P}$  match, we verify if the original pattern  $P$  matches (with the corresponding alignment).
- Can be done by brute force algorithm, with  $O(m)$  worst case cost.

- The filtering time is  $O(n/q)$ .
- Assuming that each character occurs with probability  $1/\sigma$ , the probability that  $P^j$  occurs in a given text position is  $(1/\sigma)^{\lfloor m/q \rfloor}$ .

⇒

The verification cost is on average at most  $O(mn/\sigma^{m/q})$

- We select  $q$  so that  $mn/\sigma^{m/q} \leq n/q$ , i.e.  $q = O(m/\log_\sigma m)$ .

⇒

Total average time is  $O(n \log_\sigma m/m)$ , which is optimal.

- If  $qm' > w$ , we must use several computer words  
⇒ Asymptotic running time becomes  $O(n \log_{\sigma}(m)/w)$ .
- The trick in (Peltola & Tarhio, 2003) to make BNDM work with  $m > w$  can be applied to our algorithm too.  
⇒ Omitting the details, we obtain  $O(n \log_{\sigma/h}(m)/m)$  average time where  $h = \lfloor (m-1)/w \rfloor + 1$ .
- Not optimal anymore.



## ***Linear worst case time***

---

- The worst case running time is  $O(mn)$ .
- Use any  $O(n)$  worst case time algorithm for the verifications, and do the verifications incrementally, saving the search state of the worst case algorithm after each verification.

⇒

'Standard trick', worst case becomes  $O(n)$ .

- Not a real problem: if verification time is a problem, then the filter does not work well, and can use the linear time algorithm instead.

# Implementation

- In modern pipelined CPUs branching is costly.
  - ⇒  
Unroll  $U$  times (i.e. repeat inline  $U$  times) the code  
 $D \leftarrow (D \ll 1) | B[T[i]]$ .
  - ⇒  
The bit positions indicating the occurrences will overflow
  - ⇒  
Reserve  $U - 1$  extra bits per pattern to avoid interference.
  - ⇒  
 $q(U - 1 + \lfloor m/q \rfloor)$  bits in total.
- Verification is done only every  $U$ th step, for those (at most  $U$ ) alignments that could match.
- Much faster in practice.

# *Experimental results*

---

- Implementation in C, compiled using `icc 8.1` with full optimizations, run in a 2.4GHz Pentium 4 ( $w = 32$ ), with 512MB RAM, running Linux 2.4.20-8.
- 100 patterns were randomly extracted from the text.
- Each pattern was then searched for separately.
- We report the average speed in megabytes per second.
- Our data: real DNA and protein data, English natural language and random ASCII text ( $\sigma = 96$ ).

# *Experimental results*

---

We compared against:

**BNDM:** (Navarro & Raffinot, 2000), competitive only for random ASCII.

**SBNDM:** Simplified version of BNDM (Peltola & Tarhio, 2003), competitive only for random ASCII.

**BMH, BMHS:** Boyer-Moore-Horspool, and the Sunday variant of BMH.

Not competitive on any data (results omitted).

Our algorithms:

**AOSO:** Our basic algorithm...

**FAOSO:** ...with loop-unrolling.

# Experiments: DNA

$m, q$	AOSO	FAOSO	BNDM	SBNDM
4, 2	321	<b>503</b>	181	210
8, 2	539	<b>763</b>	312	357
12, 3	702	<b>941</b>	438	492
16, 3	1029	<b>1229</b>	567	598
20, 4	1079	<b>1341</b>	750	804
24, 4	1229	<b>1525</b>	1106	1164
28, 5	1427	<b>1638</b>	1106	1164

## *Experiments: proteins*

$m, q$	AOSO	FAOSO	BNDM	SBNDM
4, 2	580	<b>909</b>	415	512
8, 4	944	<b>1267</b>	642	678
12, 4	1120	<b>1376</b>	816	926
16, 4	1120	<b>1459</b>	963	1025
20, 4	1235	<b>1376</b>	1175	1204
24, 5	1267	<b>1338</b>	1235	1302
28, 6	<b>1302</b>	<b>1302</b>	<b>1302</b>	<b>1302</b>

## *Experiments: natural language*

$m, q$	AOSO	FAOSO	BNDM	SBNDM
4, 2	579	884	368	476
8, 4	1034	1262	685	778
12, 4	1144	1279	797	845
16, 5	1200	1389	831	944
20, 6	1279	1389	1013	1092

## *Experiments: random ASCII*

$m, q$	AOSO	FAOSO	BNDM	SBNDM
4, 2	599	952	633	<b>1053</b>
8, 4	1124	<b>1333</b>	1064	1220
12, 4	1250	<b>1389</b>	1299	1282
16, 4	1351	<b>1389</b>	1351	<b>1389</b>
20, 6	1449	<b>1471</b>	1370	1429



- Very simple to implement.
- Very efficient in practice.
- Optimal for short patterns ( $m \leq w$ ).
- The techniques can be adapted for several other algorithms as well, e.g.
  - Shift-add (for Hamming distance):  
 $O(n(k + \log_{\sigma}(m))/m)$  average time.
- Any algorithm for multiple string matching can be used in place of Shift-or.