

Playing with patterns, searching for strings

Citation for published version (APA): Woude, van der, J. C. S. P. (1987). *Playing with patterns, searching for strings*. (Computing science notes; Vol. 8713). Technische Universiteit Eindhoven.

Document status and date: Published: 01/01/1987

Document Version:

Publisher's PDF, also known as Version of Record (includes final page, issue and volume numbers)

Please check the document version of this publication:

• A submitted manuscript is the version of the article upon submission and before peer-review. There can be important differences between the submitted version and the official published version of record. People interested in the research are advised to contact the author for the final version of the publication, or visit the DOI to the publisher's website.

• The final author version and the galley proof are versions of the publication after peer review.

 The final published version features the final layout of the paper including the volume, issue and page numbers.

Link to publication

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- · Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
 You may freely distribute the URL identifying the publication in the public portal.

If the publication is distributed under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license above, please follow below link for the End User Agreement:

www.tue.nl/taverne

Take down policy

If you believe that this document breaches copyright please contact us at:

openaccess@tue.nl

providing details and we will investigate your claim.

Playing with patterns, searching for strings

Jaap van der Woude

-

87/13

.

•

june 87

COMPUTING SCIENCE NOTES

This is a series of notes of the Computing Science Section of the Department of Mathematics and Computing Science of Eindhoven University of Technology.

Since many of these notes are preliminary versions or may be published elsewhere, they have a limited distribution only and are not for review. Copies of these notes are available from the author or the editor.

Eindhoven University of Technology Department of Mathematics and Computing Science P.O. Box 513 5600 MB Eindhoven The Netherlands All rights reserved editor: F.A.J. van Neerven

Playing with patterns, searching for strings

Jaap van der Woude

june 87

§ 0

Pattern identification is a source of several instructive programming exercises. We shall present one such exercise that is especially interesting for problems dealing with periodicity. In particular it enables us to treat preprocessing and search in the Knuth-Morris-Pratt pattern search algorithm as a unity.

Some remarks on names and notations:

Let Σ be a fixed alphabet. A word over Σ , i.e. an element of the free monoid (Σ^* , Λ) generated by Σ , will be called a *string*. $\Sigma^* \setminus \{\Lambda\}$ will be denoted by Σ^+ , the nonempty strings. (In the sequel, capitals refer to strings, lowercase letters to naturals (including 0) or functions, unless stated otherwise.)

Let X be a string. In order to facilitate references to the length |X| of X and symbols occuring in X, we shall write $X (i: 0 \le i < N)$. Then |X| = N and X(i) is the i + 1th symbol of X.

If $n \le N$ we shall write $X \downarrow n$ for X ($i: 0 \le i < n$), the prefix of X with e length n.

A string X is called *periodic* if $X = P^m$ for some $P \in \Sigma^+$ and $m \ge 2$, where P^m is the concatenation of m copies of P. In that case, P as well as |P| is called a *period* of X. The period of X is the smallest period (if any).

For strings X and Y

 $X \le Y$ denotes "X is a prefix of Y" X < Y means $X \le Y \land X \ne Y$

So much for general remarks on names and notations.

First we shall state the basic problem (MPP) in its "historical" context, subsequently we shall give two applications

- Pattern-search (e.g. Knuth-Morris-Pratt)
- Periodicity-search (for all prefixes)

§ 1. The MPP problem

The problem we would like to consider evolved from the exercise below. We shall give some heuristics for this program evolution.

Exercise ("carrécheck") A string X is a *carré* if $X = P^2$ for some string $P \in \Sigma^+$. Derive a program to find every prefix of X that is a carré.

The formal specification : $i[N : int; \{N \ge 1\}$ $X (i: 0 \le i < N) : \underline{string};$ $i[c (i: 1 \le i \le N) : \underline{array of} \text{ bool};$ CARRÉCHECK $\{(A i: 1 \le i \le N : c(i) = car(X \downarrow i))\}$]|]|

where

(0) $\operatorname{car}(X) = (\mathbf{E}P : P \in \Sigma^+ : X = PP).$

A standard feature in programming methodology is: weaken in the postcondition by replacement of a constant by a variable.

If we consider the term in definition (0) : X = PP we might think of P as being a constant. Symmetry tells us to replace P by a variable twice. So X = PP might be "generalized" to $X = PE \land X = FP$ for variables E and F.

I.e. P is a pre- and postfix of X.

So for strings P and X we define

(1)
$$P \operatorname{pp} X \equiv (\mathbf{E} E, F : E, F \in \Sigma^+ : X = PE \land x = FP)$$

and note that $P pp X \wedge |X| = 2 * |P| \Rightarrow car(X)$.

The lack of reflexivity of pp, created by the domain in (1), seems unnatural, but it results from the definition of periodicity, i.e. the domain in definition (0).

Moreover, incorporating reflexivity of <u>pp</u> leads to additional non-triviality analysis (in <u>mpp</u> below for example).

First we give a few properties of <u>pp</u>, the simple proofs are omitted. Let H, P and X be strings and h, $x \in \Sigma$, then

- (2) $H pp P \wedge P pp X \Rightarrow H pp X$
- $(3) \qquad H \operatorname{pp} X \wedge P \operatorname{pp} X \wedge |H| < |P| \Rightarrow H \operatorname{pp} P$
- (4) $Hh pp Xx = Hh < Xx \land H pp X \land h = x$

Property 4 will be used for prefixes of a fixed string X:

(4')
$$X \downarrow k+1 \operatorname{pp} X \downarrow n+1 = X \downarrow k \operatorname{pp} X \downarrow n \land X(k) = X(n)$$

$$(4'') \qquad \neg X \downarrow k + 1 \underline{pp} X \downarrow n + 1 \equiv \neg X \downarrow k \underline{pp} X \downarrow n \lor X(k) \neq X(n)$$

By 2, the transitivity of pp, we feel invited to consider maximal pre- and postfixes. So define :

(5)
$$P \operatorname{mpp} X = P \operatorname{pp} X \wedge (A H : H \operatorname{pp} X : H = P \vee H \operatorname{pp} P)$$

Note that the following are equivalent

$$\underline{a} P \underline{\text{mpp}} X$$

$$\underline{b} P \underline{\text{pp}} X \land (\mathbf{A} H : H \underline{\text{pp}} X : |H| \le |P|)$$

$$\underline{c} P \underline{\text{pp}} X \land (\mathbf{A} H : H \underline{\text{pp}} X : H \le P)$$

In section 6 we show that, given $P \mod X$, we have

$$car(X) \equiv |X| \mod (2*(|X| - |P|)) = 0$$

So indeed the generalization (weakening of carré) is fruitful. The carrécheck problem has evolved to the <u>MPP</u> problem:

Derive a program that calculates for every prefix of a given string its maximal pre- and postfix. The formal specification

```
\begin{split} &|[N: \text{ int }; \{N \ge 1\} \\ &X \ (i: 0 \le 1 < N) : \underline{\text{string }}; \\ &\|[f \ (i: 1 \le i \le N) : \underline{\text{array of }}[0..N); \\ & \text{MPP} \\ &\{(\mathbf{A} \ i: 1 \le i \le N : X \downarrow f \ (i) \underline{\text{mpp}} X \downarrow i)\} \\ ]\| \\ ]\| \end{split}
```

.

§ 2. Solution to the MPP problem

Forced by the postcondition of MPP :

 $\mathbf{R}_0 \quad (\mathbf{A} \ i : 1 \le i \le N : X \downarrow f(i) \underline{\mathrm{mpp}} X \downarrow i),$

we choose the following invariants :

 $P_0 \quad (A \ i : 1 \le i \le n : X \downarrow f(i) \underline{mpp} X \downarrow i)$

$$\mathbf{P}_1 \qquad 1 \le n \le N \,.$$

Approximation 0 (for MPP)

 $n := 1 ; f : (1) = 0 \{ P_0 \land P_1 \}$; do $n \neq N \rightarrow S_0 \{ X \downarrow k \text{ mpp } X \downarrow n + 1 \}$; $f : (n+1) = k \{ (P_0 \land P_1)_{n+1}^n \}$; $n := n + 1 \{ P_0 \land P_1 \}$ od $\{ P_0 \land P_1 \land n = N , \text{ hence } R_0 \}$

On cosmetical grounds, with 4' in mind, we consider a slightly different postcondition for S_0 :

$$R_1 \quad X \downarrow k + 1 \underline{\mathrm{mpp}} X \downarrow n + 1$$

By definition of <u>mpp</u> (version $5\underline{b}$) and by 4', R_1 equivales

$$X \downarrow k \operatorname{pp} X \downarrow n \land X(k) = X(n) \land (A j : X \downarrow j \operatorname{pp} X \downarrow n + 1 : j \le k + 1)$$

This leads us to a repetition for S_0 with guard $X(k) \neq X(n)$ and invariants :

 $Q_0 \quad (A \ j : X \downarrow j \text{ pp } X \downarrow n + 1 : j \le k + 1)$ $Q_1 \quad X \downarrow k \text{ pp } X \downarrow n \quad \land k \ge 0$ <u>Approximation 1 (for S_0)</u> $k := f(n) \{Q_0 \land Q_1, \text{ see the note below}\}$ $; \text{ do } X(k) \ne X(n) \land k \ne 0$ $\rightarrow \text{"decrease } k \text{ under invariance of } Q \text{ (and } P)\text{"}$ $\text{ od } \{Q_0 \land Q_1 \land (X(k) = X(n) \lor k = 0) \}$

The second conjunct in the guard is forced upon us by the wish to decrease k, leaving $k \ge 0$ invariant.

/

note

 P_0

$$\Rightarrow \{\text{instantiation at } n, \text{ def. mpp } (5\underline{b})\} \\ X \downarrow f(n) \text{ pp } X \downarrow n \land (A \ j : X \downarrow j \text{ pp } X \downarrow n : j \leq f(n)) \\ \Rightarrow \{\text{by } 4' : X \downarrow j + 1 \text{ pp } X \downarrow n + 1 \Rightarrow X \downarrow j \text{ pp } X \downarrow n \} \\ X \downarrow f(n) \text{ pp } X \downarrow n \land (A \ j : X \downarrow j + 1 \text{ pp } X \downarrow n + 1 : j \leq f(n)) \\ \Rightarrow \{\text{dummy change}\} \\ X \downarrow f(n) \text{ pp } X \downarrow n \land (A \ j : X \downarrow j \text{ pp } X \downarrow n + 1 : j \leq f(n) + 1) \\ = \{\text{def. } Q_0, Q_1\} \\ (Q_1 \land Q_0) \overset{k}{f}_{(n)} \end{cases}$$

Under assumption of $Q_0 \wedge Q_1 \wedge P_0 \wedge P_1$ and the guard, we study reduction of k. First with respect to Q_0 :

Let
$$j > 0$$
 then

$$X \downarrow j \text{ pp } X \downarrow n + 1$$

$$\Rightarrow \{Q_0; \text{ by } 4^n, X(k) \neq X(n) \Rightarrow \neg X \downarrow k + 1 \text{ pp } X \downarrow n + 1\}$$

$$X \downarrow j \text{ pp } X \downarrow n + 1 \land j \leq k + 1 \land j \neq k + 1$$

$$\Rightarrow \{j \geq 0, 4'\}$$

$$X \downarrow j - 1 \text{ pp } X \downarrow n \land j - 1 < k$$

$$\Rightarrow \{Q_1, 3\}$$

$$X \downarrow j - 1 \text{ pp } X \downarrow k$$

$$\Rightarrow \{P_0, 0 \neq k < n, \text{ def. mpp } 5a\}$$

$$j - 1 \leq f(k)$$

For j = 0, certainly we have $j \le f(k)+1$. Hence $Q_{0f(k)}^{k}$ holds.

With respect to Q_1 :

$$P_{0} \wedge Q_{1}$$

$$\Rightarrow \{k \neq 0 : P_{0} \text{ instantiated at } k\}$$

$$X \downarrow f(k) \text{ pp } X \downarrow k \land X \downarrow k \text{ pp } X \downarrow n$$

$$\Rightarrow \{2\}$$

$$X \downarrow f(k) \text{ pp } X \downarrow n$$

$$= \{\text{def. } Q_{1}, f(k) \ge 0\}$$

$$Q_{1} \overset{k}{f(k)}$$

This shows that "decrease k under invariance ..." is established by k := f(k). (Certainly $P_0 \wedge P_1$ is not affected.)

Because of the conjunct $k \neq 0$ in the guard, neither R₁ not the original postcondition of S₀ are met, but a mixture is :

for

and

$$Q_0 \wedge Q_1 \wedge X(k) = X(n) \implies R_1$$

$$Q_0 \wedge Q_1 \wedge X(k) \neq X(n) \wedge k = 0$$

 $\Rightarrow \{4^{"}\}$ $Q_0 \land Q_1 \land \neg X \downarrow k+1 pp X \downarrow n+1 \land k = 0$ $\Rightarrow \{X \downarrow 0 pp X \downarrow n+1, def. mpp\}$ $X \downarrow k mpp X \downarrow n+1 \land k = 0$

This proves the following solution for MPP :

$$n := 1 ; f : (1) = 0 \{ P_0 \land P_1 \}$$

; do $n \neq N$
 $\rightarrow | [k : int ;$
 $k := f(n) \{ Q_0 \land Q_1 \}$
; do $X(k) \neq X(n) \land k \neq 0$
 $\rightarrow k := f(k)$
od $\{ Q_0 \land Q_1 \land (X(k) = X(n) \lor k = 0) \}$
; if $X(k) = X(n) \rightarrow \{ X \downarrow k + 1 \text{ mpp } X \downarrow n + 1 \} k := k + 1$
 $[X(k) \neq X(n) \rightarrow \{ X \downarrow k + 1 \text{ mpp } X \downarrow n + 1 \} \text{ skip}$
fi $\{ X \downarrow k \text{ mpp } X \downarrow n + 1 \}$
; $f : (n+1) = k \{ (P_0 \land P_1)_{n+1}^n \}$
]|
; $n := n + 1$
od $\{ P_0 \land P_1 \land n = N \text{, hence } R_0 \}$

For the complexity of the algorithm, consider k to exist outside the innerblock $(P_2: k = f(n))$. A variant function that shows linearity is 2N - 2n + k.

`

§ 3. Pattern search

Let $P \in \Sigma^+$ be fixed, the pattern. Suppose we are interested in (all) occurrences of P in a string Z.

Put Y = PZ, then we are searching for numbers *n*, such that

 $n \ge 2 * |P| \land P pp Y \downarrow n$

In this setting the pattern might be represented by its length only.

Let X be a string, p a number $1 \le p < |X|$. As $X \downarrow p$ is a postfix of $X \downarrow n$ iff $p = n \lor X \downarrow p$ pp $X \downarrow n$, we define (the occurrence of the pattern as postfix of $X \downarrow n$):

(6)
$$O(n) \equiv p = n \lor X \downarrow p pp X \downarrow n$$

Let f be as in the MPP problem. It seems reasonable to hope for suitable O-information in the <u>mpp</u>-knowledge recorded in f. Indeed, for strings H, P and X we have

(7)
$$P \operatorname{mpp} X \Rightarrow (H \operatorname{pp} X \equiv H = P \lor H \operatorname{pp} P).$$

Property 7, which is closely linked to 3, follows easily from $5\underline{a}$, 2. It relates O(n) to f(n) as follows :

$$O(n)$$

$$= \{6\}$$

$$p = n \lor X \downarrow p \text{ pp } X \downarrow n$$

$$= \{X \downarrow f(n) \text{ mpp } X \downarrow n, 7 \text{ with } H, P, X := X \downarrow p, X \downarrow f(n), X \downarrow n\}$$

$$p = n \lor X \downarrow p = X \downarrow f(n) \lor X \downarrow p \text{ pp } X \downarrow f(n)$$

$$= \{X \downarrow p = X \downarrow f(n) \equiv p = f(n); 6\}$$

$$p = n \lor O(f(n))$$

As f(n) < n (nonreflexivity of pp), O(n) depends only on f(n) and O(i: $1 \le i < n$). This settles pattern search as a simple extension of the MPP problem, by adding invariant

$$P_3$$
 (A *i* : $1 \le i \le n$: $o(n) = O(n)$)

and initialization $o: (1) = (p = 1) \{P_3\}$ extra statement; $o: (n + 1) = (p = n + 1) \lor o(f(n + 1)) \{P_{3n+1}^n\}$ immediately following the innerblock.

§ 4. Knuth - Morris - Pratt

The pattern search presented in the previous section has a serious drawback : storage linear in the length of the given string (concatenated with the pattern). Indeed, the algorithm needs $f(i:1 \le i \le n)$ to calculate f(n+1) and $O(i:1 \le i \le n)$ and f(n+1) to calculate O(n+1).

As, for fixed p, we are interested in n such that $X \downarrow p pp X \downarrow n$ (instead of mpp !) the information recorded in f exceeds our needs : we might do with pre- and postfixes with lengths at most p. So we define : (P and p are not related !)

(8)
$$P \underline{\pi} \pi X \equiv P \underline{p} p X \land |P| < p$$

(9)
$$P \ \underline{\mu} \pi \pi \equiv P \ \underline{\pi} \pi X \land (A \ H : H \ \underline{\pi} \pi X : H = P \lor H \ \underline{\pi} \pi P)$$

The reader is urged to convince himself of the truth of the $\pi\pi$ -versions of 2, 3, $5^{a,b,c}$ (i.e. only (m)pp replaced by ((μ) $\pi\pi$). Property 4, however, has a slightly different $\pi\pi$ -version. We shall only provide the $\pi\pi$ -version of 4':

(10')
$$X \downarrow k \ \underline{\pi \pi} \ X \downarrow n \land X(k) = X(n)$$

$$\equiv (X \downarrow k + 1 \underline{\pi} \underline{\pi} X \downarrow n + 1 \land k$$

The $\mu\pi\pi$ problem is given by the postcondition

 ρ_0 (A $i: 1 \le i \le N: X \downarrow \phi(i) \mu \pi \pi X \downarrow i$)

For the solution of the $\mu\pi\pi$ problem we define invariants, (the obvious adaptations of the invariants for *MPP*)

 $\pi_0 \quad (\mathbf{A} \ i : 1 \le i \le n : X \downarrow \phi(i) \underline{\mu \pi \pi} X \downarrow i)$

$$\pi_1 \quad 1 \leq n \leq N$$

- $\Psi_0 \quad (\mathbf{A} \ j \ : X \downarrow j \ \underline{\pi\pi} X \downarrow n+1 \ : \ j \le k+1)$
- $\Psi_1 \quad X \downarrow k \ \underline{\pi\pi} \ X \downarrow n \ \land \ k \ge 0$

Except from the obvious adaptations, $\mu\pi\pi$ differs from MPP only in the case analysis in the innerblock : the drawback of 10' :

Certainly,

$$\begin{split} & \Psi_0 \wedge \Psi_1 \wedge X(k) = X(n) \wedge k but
$$\begin{aligned} & \Psi_0 \wedge \Psi_1 \wedge X(k) = X(n) \wedge k = p - 1 \\ \Rightarrow & \{10^r, k + 1 = p\} \\ & X \downarrow p \ \underline{pp} X \downarrow n + 1 \end{aligned}$$
$$\Rightarrow & \{\pi_0, \text{ instantation at } p \ ; 2\} \\ & X \downarrow \phi(p) \ \underline{\mu}\pi\pi X \downarrow n + 1 \end{split}$$$$

This shows that the alternative statement following the inner repetition should be changed (for the $\mu\pi\pi$ problem) to

 $\{ \Psi_0 \land \Psi_1 \land (X(k) = X(n) \lor k = 0) \}$; if $X(k) = X(n) \land k$ $[] <math>X(k) = X(n) \land k = p - 1 \rightarrow \{X \downarrow \phi(p) \ \underline{\mu \pi \pi} X \downarrow n + 1\} k := \phi(p)$ [] $X(k) \neq X(n) \qquad \rightarrow \{k = 0 \land X \downarrow k \ \underline{\mu \pi \pi} X \downarrow n + 1\}$ skip fi { $X \downarrow k \ \underline{\mu \pi \pi} X \downarrow n + 1$ }

This and the change from f to ϕ make the solution of MPP to a solution of $\mu\pi\pi$. The code will reappear in the Knuth-Morris-Pratt pattern search algorithm, so we shall leave it with this. Note that for calculation of $\phi(n+1)$ only $\phi(n)$ and $\phi(i:1 \le i \le p)$ are needed : $\phi(n)$ for initializing k, $\phi(i:1 \le i \le p)$ in the inner repetition. I.e. $\mu\pi\pi$ needs storage proportional to the "pattern length".

Similar to § 3. we now transform $\mu\pi\pi$ to a pattern search algorithm : Knuth-Morris-Pratt. With respect to occurrence of $X \downarrow p$ as postfix of $X \downarrow n+1$, note that

$$O(n + 1)$$

$$= \{6\}$$

$$p = n + 1 \lor X \downarrow p \text{ pp } X \downarrow n + 1$$

$$= \{4'; \text{ def. } \underline{\pi}\pi\}$$

$$p = n + 1 \lor (X \downarrow p - 1 \underline{\pi}\underline{\pi} X \downarrow n \land X(p - 1) = X(n))$$

$$= \{\text{ def. of } \underline{\pi}\underline{\pi}: \Psi_{0p-1}^{k} = \text{ true }; \text{ def. } \Psi\}$$

$$p = n + 1 \lor (\Psi_{0} \land \Psi_{1} \land X(k) = X(n))_{p-1}^{k}$$

So calculation of O(n+1) depends only on the postcondition of the inner repetition and occurrence is to be signalled in the second alternative (the occurrence at n+1=p is not relevant of course).

We are now ready for the Knuth-Morris-Pratt algorithm.

As only $\phi(i: 1 \le i \le p)$ and $\phi(n)$ are needed to calculate $\phi(n+1)$, we have to distinguish between

preprocessing – "filling \$\phi" search – "signalling occurrences".

This separation of the two parts is inevitable, but an earlier separation is unnecessary, unelegant and confusing. In order to account for the reduced domain of ϕ we modify π_0 to π_0^1 , to "buffer" $\phi(n)$ we add π_2

 $\pi_0^{\perp} \quad (A \ i \ : \ 1 \le i \le p \ \underline{\min} n \ : \ X \downarrow \phi(i) \ \underline{\mu} \pi \pi X \downarrow i)$

$$\pi_2 \quad X \downarrow k \ \mu \pi \pi X \downarrow n$$

and we take k outside the inner repetition.

• • .

The Knuth-Morris-Pratt pattern search algorithm we derived :

$$I[k: int;
\phi(i: 1 \le i \le p): array of [0..p-1);
n, k:=1, 0; \phi:(1) = 0 {\pi_0^1 \land \pi_1 \land \pi_2};
; do n \neq N
\rightarrow {\pi_0^1 \land \pi_1 \land \pi_2 \land n \neq N, so \Psi_0 \land \Psi_1}
do X(k) \neq X(n) \land k \neq 0
\rightarrow k := \phi(k)
od {\Psi_0 \land \Psi_1 \land (X(k) = X(n) \lor k = 0)};
; if X(k) = X(n) \land k < p-1 \rightarrow k := k+1
[] X(k) = X(n) \land k = p-1 \rightarrow k := \phi(p); "MATCH"
[] X(k) \neq X(n) \rightarrow skip
fi {X \notherwise k + 1, so \pi_{2n+1}};
; if n
[] n \ge p \rightarrow skip
fi {(\pi_0^1 \land \pi_1 \land \pi_2)_{n+1}^n};
; n := n + 1
od {\pi_0^1 \land \pi_1 \land \pi_2 \land n = N, so \rightarrow 0}]
]]$$

The interested reader might want to try a direct approach via $\mu\pi\pi.$

· •

`

§ 5. Further remarks on pattern search

The second alternative statement in the algorithm above, distinguishing preprocessing and search, may also lead to a code with two (sequential) repetitions, one for each alternative. We chose for the form above to stress the uniformity: the difference between the parts is solely based upon coding, the genesis doesn't differentiate!

Several people noticed the strong resemblance of those parts, but in the literature we searched in vain for a presentation or derivation (at all) of the algorithm that did justice to that resemblance. ([C 85] and [W 86] deserve some credit).

[Note that even in 1983 the preprocessing was said to be "complicated and difficult to understand" ([S 83] p. 242). As the two parts are almost identical such a statement is puzzling. Has it anything to do with the widespread chaotic algorithm presentation? (e.g. [KMP 77], [BM 77])].

In our opinion, exploitation of pre- and postfixes simplified the "derivation" of the algorithm such that it becomes within reach of every freshmen course.

We conclude the discussion of pattern search with a remark on the Boyer-Moore fast pattern search ([BM 77]). As this algorithm is slightly beyond the scope of this paper, we shall only hint at its relation with the MPP problem.

Consider $X \in \Sigma^*$ and pattern $X \downarrow p$. In the Knuth-Morris-Pratt pattern search we decided to build up pre- and postfixes bit by bit, but we could have been greedier :

To that end consider the (linear-search-like) invariant

PBM (A $i: X \downarrow p \underline{pp} X \downarrow i: i > n$)

As $X \downarrow p$ pp $X \downarrow n+1 \Rightarrow \{4'\} X(p-1) = X(n)$ we first check X(n) as a candidate for the end of a pattern occurrence.

Let $s = (\text{MAX } j : 0 \le j \le p \land X(n) = X(j) : j) \max_{n \ge 1} -1$

Then PBM_{n+p-s}^{n} holds.

In other words : the first candidate m to satisfy $X \downarrow p$ pp $X \downarrow m$ is m = n + p - s.

If s we can "leap further", if <math>s = p - 1 we check X(n-1), etcetera.

This requires knowledge of the occurrences of values and periodicities in the pattern.

The reader is challenged to give a "derivation" of the Boyer-Moore fast pattern search based on this early deviation of the MPP problem.

§ 6. Periodicity search

In section 1 we "generalized" the carrécheck exercise to the MPP problem, and we promised to show that a solution for carrécheck is found as soon as MPP is solved.

We shall keep our promise in the following way :

- we give a variant of carrécheck
- we proclaim an enrichment of MPP that solves the (variant) exercise
- we perform some string-mathematics to prove that the exercise is solved by that enrichment of MPP.

For fixed $m \ge 2$ consider the following postcondition

 $\mathbf{R} \qquad (\mathbf{A} \ i \ : \ 1 \le i \le N \ : \ c(i) = (\mathbf{E} \ P \ : \ P \in \Sigma^+ \ : \ X \ \downarrow \ i = P^m)$

 $\wedge (\mathbf{A} \ i : 1 \le i \le N : \operatorname{per}(i) = i \ \underline{\min}(\mathbf{MIN} \ p : p \ \operatorname{period} \ \operatorname{of} \ X \ \downarrow i : p))$

In case m = 2, the first conjunct of R is just the postcondition of carrécheck.

The second conjuct of R means :

per(i) is the period of $X \downarrow i$ if $X \downarrow i$ is periodic, otherwise per(i) = i.

Obviously we should extent invariant P for the MPP problem with a conjunct P_4 to get an invariant for the new problem.

 $P_4 = R_n^N$

Initialization of P_4 : ; c: (1) = false ; per: (1) = 1.

The outer repetition should contain an establishment of P_{4n+1} .

So, following " $f: (n + 1) = k \{P_{0n+1}^n\}$ " and before "n: = n + 1", we proclaim the statement list :

; c: $(n + 1) = (n + 1) \mod (m * (n + 1 - f (n + 1))) = 0$; if $(n + 1) \mod (n + 1 - f (n + 1)) = 0 \rightarrow \text{per}: (n + 1) = n + 1 - f (n + 1)$ [] $(n + 1) \mod (n + 1 - f (n + 1)) \neq 0 \rightarrow \text{per}: (n + 1) = n + 1$ fi {P_{4n+1}, see corollary 5 to follow }.

Indeed a minor adaptation, but it takes a proof!

The string-mathematics to follow is quite elementary, and has nothing to do with programming and - methodology. So we adopt a more conventional mathematical style, but (for the convenience of non-mathematicians) we still take small steps in the proofs.

The basic idea is to squeeze periodicity information out of pp or mpp knowledge.

Let $D, Y \in \Sigma^+$ with D pp Y. Then there are $E, F \in \Sigma^+$ such that $Y = DE \land Y = FD$. Lemmata 1 and 2 tell us about (almost) periodicity of Y. (they are well-known, e.g. see [L 79] Ch. 11.5). Much more can be said if D mpp Y, some of which is done in 3,4,5. Lemma 1 Let $D \in \Sigma^*$ and $E, F \in \Sigma^+$ such that DE = FD. Then there are $L \in \Sigma^*, K \in \Sigma^+$ and $n \ge 0$ with $\underline{0} \quad D = F^n L$ and L < F (hence $DE = FD = F^{n+1} L$) $\underline{1} \quad E = KL$ and F = LK. Proof Certainly, there are $L \in \Sigma^*$ and $n \ge 0$ such that $D = F^n L$ and |L| < |F|.

Then $F^{n}LE = DE = FD = F^{n+1}L$, so LE = FL.

As |L| < |F| = |E|, there are K, $K^1 \in \Sigma^+$ with LK = F and $E = K^1 L$.

Hence, $L(K^{1}L) = LE = FL = (LK)L$ and it follows that $K^{1} = K$.

Lemma 2 Let $D, F \in \Sigma^*$ with DF = FD. Then there is a $P \in \Sigma^*$ such that $D, F \in \{P^m \mid m \ge 0\}$. If $D, F \in \Sigma^+$ it follows that DF is periodic with period at most gcd(|D|, |F|).

Proof By induction to the length of DF.

If $D = \Lambda$ or $F = \Lambda$ the existence of P is obvious.

Let D, $F \in \Sigma^+$. By 1, there are K, L, n such that $D = F^n L$, F = KL and F = LK.

Hence KL = LK.

As $D \neq \Lambda$, |KL| = |F| < |DF|, so by induction there is a $P \in \Sigma^*$ (as $K \in \Sigma^+$ even $P \in \Sigma^+$) such that $K, L \in \{P^m \mid m \ge 0\}$.

Consequently, $D, F \in \{P^m \mid m \ge 0\}$ which proves the first part.

If D, $F \in \Sigma^+$ then $DF \in \{P^{m+2} \mid m \ge 0\}$ while $P \in \Sigma^+$.

Note that |P| divides |D| and |F|.

Lemma 3 Let Y = FD and $D \mod Y$. Then F is not periodic.

.

Proof By definition of (m)pp, $F \neq \Lambda$. So, by 1, there are L, n with $D = F^n L$ and L < F. Suppose F is periodic, say $F = Q^m$ for some $Q \in \Sigma^+$, $m \ge 2$.

Then QF = FQ, hence QL < QF = FQ. As also L < F < FQ, it follows that FQ has both Land QL as prefix. Since |L| < |QL| we have L < QL and, equivalently, $Q^m Q^{mn-1} L < Q^m Q^{mn} L$.

As $Y = FD = Q^m Q^{mn} L$, it follows that $Q^m Q^{mn-1} L pp Y$. However, since $m \ge 2$, $|D| = |Q^{mn} L| < |Q^m Q^{mn-1} L|$ which contradicts D mpp Y. This falsifies periodicity of F.

Π

Π

Lemma 4 Let Y = FD, $D \mod Y$. Let $P \mod Y$ and $|P| \ge |F|$, then there is a $k \ge 0$ such that $D = F^k P$. (I.e. all pre-postfixes of Y with lenght $\ge |F|$ are known).

Proof As $D \mod Y$, $F \neq \Lambda$, so there are $n \ge 0$ and L < F with $D = F^n L$. Because $|L| < |F| \le |P|$ and P is a postfix of D, L is a postfix of P. Hence there are a $k: 0 \le k \le n$ and H < F with $D = F^k HP$. Let HG = F then $P = GF^{n-k-1}L$. As $|H| + |G| = |F| \le |P| = |G| + |F^{n-k-1}L|$, $|H| \le |F^{n-k-1}L|$. Since $H < F^{n-k}$ and $F^{n-k-1}L < F^{n-k}$ it follows that $H \le F^{n-k-1}L$, so $GH \le GF^{n-k-1}L = P$. On the other hand $HG = F \le P$, so GH = HG. As H < F, and, by 3, F is not periodic it follows from 2 that $H = \Lambda$, which shows $D = F^k P$.

Corollary 5 Let D mpp Y, say Y = FD. Let m ≥ 2, then
(E C :: Y = C^m) iff |Y| mod (m * |F|) = 0.
In particular, Y is a carré iff |Y| mod 2|F| = 0, and Y is periodic iff D ≠ A and |Y| mod |F| = 0.
Proof By 1, Y = Fⁿ⁺¹L and L < F, so the if-part is obvious.

Let $Y = C^m$; note that since $m \ge 2$, $C^{m-1} pp Y$. As $|D| \ge |C^{m-1}|$, $|F| \le |C| \le |C^{m-1}|$, so by $\underline{4}$, $D = F^k C^{m-1}$. Hence $F^{k+1} C^{m-1} = Y = CC^{m-1}$ and $C = F^{k+1}$, so $Y = F^{m*(k+1)}$.

The remark on Y being a carré is an instantation for m = 2.

Finally, as $Y \neq F$, $|Y| \mod |F| = 0 = (E m : m \ge 2 : |Y| \mod (m * |F|) = 0)$.

Note that if Y is periodic, |Y| - |D| is the period.

Acknowledgements

The carré problem and related other problems were "en vogue" in the environment of the Eindhoven University of Technology. So, inevitably, I was contaminated too.

I would like to thank all colleagues that contributed to the genesis of this paper by comments and interest.

•

References

[BM 77]	Boyer, R.S. and Moore J.S.,	
	A fast string searching algorithm,	
	CACM <u>20</u> , 762 - 772 (1977).	
[C 85]	Chengdian, C.,	
	A derivation of the Knuth-Morris-Pratt pat-	
	tern matching program,	
	EUT-report 85-Wsk-02,	
	Eindhoven University of Technology (1985).	
[KMP 77]	Knuth, D.E., Morris, J.H. and Pratt, V.R.,	
	Fast pattern matching in strings,	
	SIAM J. Comput. 6, 323 - 350 (1977).	
[L 79]	Lallement, G.,	
	Semigroups and combinatorial applications,	
	Wiley-interscience, (1979).	
[S 83]	Sedgewick, R.,	
	Algorithms,	
	Addison - Wesley, (1983).	
[W 86]	Wiltink, G.,	
	Knuth-Morris-Pratt, private communication,	
	(1986).	

In this series appeared :

<u>No.</u>	Author(s)	<u>Title</u>
85/01	R.H. Mak	The formal specification and derivation of CMOS-circuits
85/02	W.M.C.J. van Overveld	On arithmetic operations with M-out-of-N-codes
85/03	W.J.M. Lemmens	Use of a computer for evaluation of flow films
85/04	T. Verhoeff H.M.J.L. Schols	Delay insensitive directed trace structures satisfy the foam rubber wrapper postulate
86/01	R. Koymans	Specifying message passing and real-time systems
86/02	G.A. Bussing K.M. van Hee M. Voorhoeve	ELISA, A language for formal specifications of information systems
86/03	Rob Hoogerwoord	Some reflections on the implementation of trace structures
86/04	G.J. Houben J. Paredaens K.M. van Hee	The partition of an information system in several parallel systems
86/05	Jan L.G. Dietz Kees M. van Hee	A framework for the conceptual modeling of discrete dynamic systems
86/06	Tom Verhoeff	Nondeterminism and divergence created by concealment in CSP
86/07	R. Gerth L. Shira	On proving communication closedness of distributed layers

86/08	R. Koymans R.K. Shyamasundar W.P. de Roever R. Gerth S. Arum Kumar	Compositional semantics for real-time distributed computing (Inf. & Control 1987)
86/09	C. Huizing R. Gerth W.P. de Roever	Full abstraction of a real-time denotational semantics for an OCCAM-like language
86/10	J. Hooman	A compositional proof theory for real-time distributed message passing
86/11	W.P. de Roever	Questions to Robin Milner - A responders commentary (IFIP86)
86/12	A. Boucher R. Gerth	A timed failures model for extended communicating processes
86/13	R. Gerth W.P. de Roever	Proving monitors revisited: a first step towards verifying object oriented systems (Fund. Informatica IX-4)
86/14	R. Koymans	Specifying passing systems requires extending temporal logic
87/01	R. Gerth	On the existence of a sound and complete axiomatizations of the monitor concept
87/02	Simon J. Klaver Chris F.M. Verberne	Federatieve Databases
87/03	G.J. Houben J. Paredaens	A formal approach to distributed information systems
87/04	T. Verhoeff	Delay-insensitive codes - An overview
87/05	R. Kuiper	Enforcing non-determinism via linear time temporal logic specification

s

	·	- 17 -
87/06	R. Koymans	Temporele logica specificatie van message passing en real-time systemen (in Dutch)
87/07	R. Koymans	Specifying message passing and real-time systems with real-time temporal logic
87/08	H.M.J.L. Schols	The maximum number of states after projection
87/09	J. Kalisvaart L.R.A. Kessener W.J.M. Lemmens M.L.P van Lierop F.J. Peters H.M.M. van de Wetering	Language extensions to study structures for raster graphics
87/10	T. Verhoeff	Three families of maximally nondeterministic automata
87/11	P. Lemmens	Eldorado ins and outs. Specifications of a data base management toolkit according to the functional model
87/12	K.M. van Hee A. Lapinski	OR and AI approaches to decision support systems
87/13	J. van der Woude	Playing with patterns, searching for strings