

# On Obtaining Knuth, Morris, and Pratt’s String Matcher by Partial Evaluation

Mads Sig Ager, Olivier Danvy, and Henning Korsholm Rohde

BRICS \*

Department of Computer Science

University of Aarhus †

## Abstract

We present the first formal proof that partial evaluation of a quadratic string matcher can yield the precise behaviour of Knuth, Morris, and Pratt’s linear string matcher.

Obtaining a KMP-like string matcher is a canonical example of partial evaluation: starting from the naive, quadratic program checking whether a pattern occurs in a text, one ensures that backtracking can be performed at partial-evaluation time (a binding-time shift that yields a staged string matcher); specializing the resulting staged program yields residual programs that do not back up on the text, à la KMP. We are not aware, however, of any formal proof that partial evaluation of a staged string matcher precisely yields the KMP string matcher, or in fact any other specific string matcher.

In this article, we present a staged string matcher and we formally prove that it performs the same sequence of comparisons between pattern and text as the KMP string matcher. To this end, we operationally specify each of the programming languages in which the matchers are written, and we formalize each sequence of comparisons with a trace semantics. We also state the (mild) conditions under which specializing the staged string matcher with respect to a pattern string provably yields a specialized string matcher whose size is proportional to the length of this pattern string and whose time complexity is proportional to the length of the text string. Finally, we show how tabulating one of the functions in this staged string matcher gives rise to the ‘next’ table of the original KMP algorithm.

The method scales for obtaining other linear string matchers, be they known or new.

---

\*Basic Research in Computer Science ([www.brics.dk](http://www.brics.dk)), funded by the Danish National Research Foundation.

†Ny Munkegade, Building 540  
DK-8000 Aarhus C, Denmark  
E-mail: {mads,danvy,hense}@brics.dk

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ASIA-PEPM’02, September 12-14, 2002, Aizu, Japan.  
Copyright 2002 ACM 1-58113-458-4/02/0009 ...\$5.00.

## Categories and Subject Descriptors

D.1.1 [Software]: Programming Techniques—*applicative (functional) programming*; D.1.2 [Software]: Programming Techniques—*automatic programming*; D.1.4 [Software]: Programming Techniques—*sequential programming*; D.2.4 [Software]: Software Engineering—*correctness proofs*; D.3.4 [Programming Languages]: Processors—*code generation*; F.2.2 [Logics and Meanings of Programs]: Analysis of Algorithms and Problem Complexity—*pattern matching*; F.3.2 [Logics and Meanings of Programs]: Semantics of Programming Languages—*operational semantics, partial evaluation, program analysis*; I.1.2 [Symbolic and Algebraic Manipulation]: Algorithms—*Analysis of algorithms*; I.1.1 [Computing Methodologies]: Symbolic and Algebraic Manipulation—*Simplification of expressions*.

## Keywords

Knuth-Morris-Pratt string matching, program specialization, data specialization, trace semantics.

## 1. Introduction

Obtaining Knuth, Morris, and Pratt’s linear string matcher out of a naive quadratic string matcher is a traditional exercise in partial evaluation:

$$\left\{ \begin{array}{l} \text{run } match \langle pat, txt \rangle = res \\ \text{run } PE \langle match, \langle pat, \_ \rangle \rangle = match_{\langle pat, \_ \rangle} \\ \text{run } match_{\langle pat, \_ \rangle} \langle \_ , txt \rangle = res \end{array} \right.$$

Given a static pattern, the partial evaluator should perform all backtracking statically to produce a specialized matcher that traverses the text in linear time.

Initially, the exercise was proposed by Futamura to illustrate Generalized Partial Computation, a form of partial evaluation that memoizes the result of dynamic tests when processing conditional branches [11]. Subsequently, Consel and Danvy pointed out that a binding-time improved (i.e., staged) quadratic string matcher could also be specialized into a linear string matcher, using a standard, Mix-style partial evaluator [8]. A number of publications followed, showing either a range of binding-time improved string matchers or presenting a range of partial evaluators integrating the binding-time improvement [2, 10, 12, 13, 16, 24, 25, 26].

After 15 years, however, we observe that

1. the KMP test, as it is called, appears to have had little impact, if any, on the development of algorithms outside the field of partial evaluation, and that

2. except for Grobauer and Lawall’s recent work [14], issues such as the precise characterization of time and space of specialized string matchers have not been addressed.

The goal of our work is to address the second item, with the hope to contribute to remedying the first one, in the long run.

## 1.1 This work

We relate the original KMP algorithm [19] to a staged quadratic string matcher that keeps one character of negative information (essentially Consel and Danvy’s original solution [8]; there are many ways to stage a string matcher [2, 14], and we show one in appendix). Our approach is semantic rather than algorithmic or intuitive:

- We formalize an imperative language similar to the one in which the KMP algorithm is traditionally specified, and we formalize the subset of Scheme in which the staged matcher is specified.
- We then present two trace semantics that account for the sequence of indices corresponding to the successive comparisons between characters in the pattern and in the text, and we show that the KMP algorithm and the staged matcher share the same trace.
- We analyze the binding times of the staged matcher using an off-the-shelf binding-time analysis (that of Similix [4, 5]), and we observe that the only dynamic comparisons are the ones between the static pattern and the dynamic text. Therefore, specializing this staged string matcher preserves its trace, given an offline program specializer (such as Similix’s) that (1) computes static operations at specialization time and (2) generates a residual program where dynamic operations do not disappear, are not duplicated, and are executed in the same order as in the source program. We also assess the size of residual programs: it is proportional to the size of the corresponding static patterns.

This correspondence and preservation of traces shows that a staged matcher that keeps one character of negative information corresponds to and specializes into (the second half of) the KMP algorithm, precisely. It also has two corollaries:

1. A staged matcher that does not keep track of negative information, as in Sørensen, Glück, and Jones’s work on positive supercompilation [26], does not give rise to the KMP algorithm. Instead, we observe that such a staged matcher gives rise to Morris and Pratt’s algorithm [6, Chapter 6], which is also linear but slightly less efficient.
2. A staged string matcher that keeps track of all the characters of negative information, as in Futamura’s Generalized Partial Computation [10, 12], Glück and Klimov’s supercompiler [13], and Jones, Gomard, and Sestoft’s textbook [16, Figure 12.3] does not give rise to the KMP algorithm either. The corresponding residual programs are slightly more efficient than the KMP algorithm, but their size is not linearly proportional to the length of the pattern. (Indeed, Grobauer and Lawall have shown that the size of these residual programs is bounded by  $|pat| \times |\Sigma|$ , where  $pat$  denotes the pattern and  $\Sigma$  denotes the alphabet [14].)

That said,

- (a) there is more to linear string matching than the KMP: for example, in their handbook on exact string matching [6], Charras and Lecroq list over 30 different algorithms; and
- (b) many naive string matchers exist that can be staged to yield a variety of linear string matchers, e.g., Boyer and Moore’s [2].

We observe that over half of the algorithms listed by Charras and Lecroq can be obtained as specialized versions of staged string matchers. Proving this observation can be done in the same manner as in the present article for the KMP. Furthermore, we can obtain new linear string matchers by exploring the variety of staged string matchers.

## 1.2 Overview

The rest of this article is organized as follows. In Section 2, we specify an operational semantics for the imperative language used by Knuth, Morris, and Pratt, and in Section 3, we specify an operational semantics for a subset of Scheme [17]. In each of these sections, we specify:

1. the abstract syntax of the language;
2. its expressible values;
3. its evaluation rules;
4. the string matcher;
5. the semantics of the string matcher;
6. an abstract semantics of the string matcher.

The point of the abstract semantics is to account for the sequence of comparisons between the pattern and the text.

In Section 4, we show that the imperative matcher and the functional matcher give rise to the same sequence of comparisons. In Section 5, we investigate the result of specializing the functional matcher with respect to a pattern string using program specialization and then using a simple form of data specialization. Section 6 concludes.

## 2. The KMP, Imperatively

In this section, we describe the imperative language in which the imperative string matcher is specified. The language is canonical, with constant and mutable identifiers and with immutable arrays. We then present the imperative string matcher and its meaning. Finally, we specify a trace semantics of the imperative matcher.

### 2.1 Abstract syntax

A program consists of statements  $s \in Stm$ , expressions  $e \in Exp$ , numerals  $num \in Num$ , constant identifiers  $c \in Cid$ , mutable identifiers  $x \in Mid$ , array identifiers  $a \in Aid$ , and operators  $opr \in Opr$ .

$$s ::= x := e \mid s; s \mid \text{if } e \text{ then } s \text{ else } s \text{ fi} \mid \text{while } e \text{ do } s \text{ od} \mid \text{return } e$$

$$e ::= num \mid x \mid c \mid a[e] \mid e \text{ opr } e \mid e \text{ and } e$$

$$opr ::= + \mid - \mid >= \mid < \mid !=$$

## 2.2 Expressible values

A value is an integer, a boolean, or a character in an alphabet:

$$Val = \mathbb{Z} + \mathbb{B} + \Sigma$$

## 2.3 Rules

In the following rules,  $e \in Exp$ ,  $v, v_1, v_2 \in Val$ ,  $s, s_1, s_2 \in Stm$ ,  $c \in Cid$ ,  $x \in Mid$ ,  $num \in Num$ ,  $opr \in Opr$ ,  $n \in \mathbb{Z}$ , and  $Unit = \{unit\}$ .

### 2.3.1 Auxiliary constructs

The language includes numeric operators and a comparison operator over characters:

$$\begin{aligned} number(num) &= n, \text{ if } num \text{ denotes } n \\ operate(+, n_1, n_2) &= n_1 + n_2 \\ operate(-, n_1, n_2) &= n_1 - n_2 \\ operate(>=, n_1, n_2) &= n_1 \geq n_2 \\ operate(<, n_1, n_2) &= n_1 < n_2 \\ operate(!=, c_1, c_2) &= c_1 \neq c_2 \end{aligned}$$

### 2.3.2 Stores

A store is a total function:

$$\sigma : Mid \rightarrow \mathbb{Z}$$

### 2.3.3 Constants

Constants are defined with a total function:

$$C : Cid \rightarrow \mathbb{Z}$$

### 2.3.4 Arrays

Arrays are defined with a partial function:

$$A : Aid \times \mathbb{N} \rightarrow \mathbb{Z} \cup \Sigma$$

where  $\mathbb{N}$  denotes the set of natural numbers including zero. Indexing arrays starts at zero, and indexing out of bounds is undefined.

### 2.3.5 Relations

The (big-step) evaluation relation for expressions reads as

$$A, C, \sigma \vdash e \rightarrow_I v$$

and the (small-step) evaluation relation for statements reads as

$$A, C \vdash \langle s, \sigma \rangle \rightarrow_I \langle r, \sigma' \rangle$$

where  $e \in Exp$ ,  $v \in Val$ ,  $s \in Stm$ , and  $r \in Stm + Unit + \mathbb{Z}$ . If  $r \in Stm$ , the computation of  $s$  is in progress. If  $r \in Unit$ , the computation of  $s$  completed normally. If  $r \in \mathbb{Z}$ , the computation of  $s$  aborted with a return.

We choose a big-step evaluation relation for expressions because we are not interested in intermediate evaluation steps. We choose a small-step evaluation relation for statements because we want to monitor the progress of imperative computations.

### 2.3.6 Expressions

$$\frac{n = number(num)}{A, C, \sigma \vdash num \rightarrow_I n}$$

$$\frac{n = \sigma(x)}{A, C, \sigma \vdash x \rightarrow_I n}$$

$$\frac{n = C(c)}{A, C, \sigma \vdash c \rightarrow_I n}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I n \quad v = A(a, n)}{A, C, \sigma \vdash a[e] \rightarrow_I v}$$

$$\frac{A, C, \sigma \vdash e_1 \rightarrow_I v_1 \quad A, C, \sigma \vdash e_2 \rightarrow_I v_2 \quad v = operate(opr, v_1, v_2)}{A, C, \sigma \vdash e_1 \text{ opr } e_2 \rightarrow_I v}$$

$$\frac{A, C, \sigma \vdash e_1 \rightarrow_I false}{A, C, \sigma \vdash e_1 \text{ and } e_2 \rightarrow_I false}$$

$$\frac{A, C, \sigma \vdash e_1 \rightarrow_I true \quad A, C, \sigma \vdash e_2 \rightarrow_I b}{A, C, \sigma \vdash e_1 \text{ and } e_2 \rightarrow_I b}$$

### 2.3.7 Statements

$$\frac{A, C, \sigma \vdash e \rightarrow_I n \quad \sigma' = \sigma[x \mapsto n]}{A, C \vdash \langle x := e, \sigma \rangle \rightarrow_I \langle unit, \sigma' \rangle}$$

$$\frac{A, C \vdash \langle s_1, \sigma \rangle \rightarrow_I \langle s'_1, \sigma' \rangle}{A, C \vdash \langle s_1 ; s_2, \sigma \rangle \rightarrow_I \langle s'_1 ; s_2, \sigma' \rangle}$$

$$\frac{A, C \vdash \langle s_1, \sigma \rangle \rightarrow_I \langle unit, \sigma' \rangle}{A, C \vdash \langle s_1 ; s_2, \sigma \rangle \rightarrow_I \langle s_2, \sigma' \rangle}$$

$$\frac{A, C \vdash \langle s_1, \sigma \rangle \rightarrow_I \langle n, \sigma' \rangle}{A, C \vdash \langle s_1 ; s_2, \sigma \rangle \rightarrow_I \langle n, \sigma' \rangle}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I true}{A, C \vdash \langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \rightarrow_I \langle s_1, \sigma \rangle}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I false}{A, C \vdash \langle \text{if } e \text{ then } s_1 \text{ else } s_2 \text{ fi}, \sigma \rangle \rightarrow_I \langle s_2, \sigma \rangle}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I false}{A, C \vdash \langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \rightarrow_I \langle unit, \sigma \rangle}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I true}{A, C \vdash \langle \text{while } e \text{ do } s \text{ od}, \sigma \rangle \rightarrow_I \langle s ; \text{while } e \text{ do } s \text{ od}, \sigma \rangle}$$

$$\frac{A, C, \sigma \vdash e \rightarrow_I n}{A, C \vdash \langle \text{return } e, \sigma \rangle \rightarrow_I \langle n, \sigma \rangle}$$

## 2.4 The string matcher

The KMP algorithm consists of two parts: the initialization of the next table and the actual string matching [19].

### 2.4.1 Initialization of the next table

The first part builds a next table for the pattern satisfying the following definition.

#### Definition 1 (Next table)

The next table is an array of indices with the same length as the pattern:  $next[j]$  is the largest  $i$  less than  $j$  such that  $pat[j - i] \dots pat[j - 1] = pat[0] \dots pat[i - 1]$  and  $pat[j] \neq pat[i]$ . If no such  $i$  exists then  $next[j]$  is  $-1$ .

The computation of the next table is described by the following pseudocode,<sup>1</sup> where we assume that `pat`, `txt`, `lpat`, and `ltxt` are given in an initial store  $\sigma$  in which `pat` denotes the pattern and `lpat` its length, and in which `txt` denotes the text and `ltxt` its length.

```

j := 0; t := -1; next[0] := -1;
while j < lpat - 1 do
  while t >= 0 and pat[j] != pat[t] do
    t := next[t]
  od;
  t := t+1; j := j+1;
  if pat[j] = pat[t]
  then next[j] := next[t]
  else next[j] := t
  fi
od

```

## 2.4.2 String matching

The second part traverses the text using the next table as described by the following program, which is written in the imperative language specified in Sections 2.1, 2.2, and 2.3. In this second part, `lpat` and `ltxt` are constant identifiers, `j` and `k` are mutable identifiers, and `pat` and `txt` are array identifiers. (`pat` denotes the pattern and `lpat` its length, and `txt` denotes the text and `ltxt` its length.)

```

j := 0; k := 0;
while j < lpat and k < ltxt do
  while j >= 0 and pat[j] != txt[k] do
    j := next[j]
  od;
  k := k+1;
  j := j+1
od;
if j >= lpat then return k-j else return -1 fi

```

In the rest of this article, we only consider the second part of the KMP algorithm and we refer to it as the *imperative matcher*. We state without proof that the imperative matcher accesses the pattern, the text, and the next table within their bounds.

## 2.5 Semantics of the imperative matcher

We now consider the meaning of the imperative matcher. What we are after is the sequence of indices corresponding to the successive comparisons between characters in the pattern and in the text. Because the imperative language is deterministic and the KMP algorithm is correct, this sequence exists and is unique. (It is also finite.)

### Definition 2 (Comparison)

An imperative comparison for the string matcher of Section 2.4 is a derivation tree of the form

$$\frac{E}{A, C, \sigma \vdash \text{pat}[j] \neq \text{txt}[k] \rightarrow_I b}$$

where  $E$  denotes another derivation tree.

<sup>1</sup> We write ‘pseudocode’ instead of ‘code’ because in the language of Sections 2.1, 2.2, and 2.3, arrays are immutable. We could easily extend the language to support mutable arrays, but doing so would clutter the rest of our development with side conditions to the effect that the next table is not updated in the second part of the KMP algorithm. We have therefore chosen to simplify the language.

### Definition 3 (Index)

The following function maps an imperative comparison into the corresponding pair of indices in the pattern and the text:

$$\text{index} \left( \frac{E}{A, C, \sigma \vdash \text{pat}[j] \neq \text{txt}[k] \rightarrow_I b} \right) = (\sigma(j), \sigma(k))$$

Let

$$\frac{S_1}{A, C \vdash \langle s_1, \sigma_1 \rangle \rightarrow_I \langle s_2, \sigma_2 \rangle},$$

$$\frac{S_2}{A, C \vdash \langle s_2, \sigma_2 \rangle \rightarrow_I \langle s_3, \sigma_3 \rangle},$$

$$\vdots$$

$$\frac{S_{n-1}}{A, C \vdash \langle s_{n-1}, \sigma_{n-1} \rangle \rightarrow_I \langle r, \sigma_n \rangle}$$

be a derivation of the imperative matcher, where the premises  $S_1, S_2, \dots, S_{n-1}$  are other derivation trees,  $A$  contains the pattern and the text,  $C$  contains their lengths,  $s_1$  is the imperative matcher, and  $\sigma_1$  is the initial state mapping all identifiers to zero.

Each premise might contain imperative comparisons. We want to build the sequence of indices corresponding to the successive comparisons between characters in the pattern and in the text. Applying the index function to each of the imperative comparisons in each premise gives such indices. We collect them in a sequence of non-empty sets of pairs of indices as follows.

### Definition 4 (Trace)

Let  $S_1, S_2, \dots, S_{n-1}$  be the premises of a derivation of the imperative matcher. Let  $c_i$  be the set of imperative comparisons in  $S_i$ , for  $0 < i < n$ . Let  $p_i = \{\text{index}(c) \mid c \in c_i\}$ , for  $0 < i < n$ . The imperative trace is the sequence  $\pi(p_1) \cdot \pi(p_2) \cdot \dots \cdot \pi(p_{n-1})$ , where

$$\pi(p) = \begin{cases} \varepsilon & \text{if } p = \emptyset \\ p & \text{otherwise} \end{cases}$$

and where  $\varepsilon$  is the neutral element for concatenation.

In Section 2.6, Lemma 1 shows that each of the premises in Definition 4 contains at most one imperative comparison. Therefore, for all  $i$ ,  $p_i$  is either empty or a singleton set. The imperative trace is thus a sequence of singleton sets, each of which corresponds to the successive comparisons of characters in `pat` and `txt`.

### Definition 5 (Program points)

The imperative program points  $\text{Match}_I$ ,  $\text{Compare}_I$  and  $\text{Shift}_I$  are defined as the following sets of configurations:

$$\begin{aligned} \text{Match}_I &= \{\langle P, \sigma \rangle \mid \sigma(j) \geq 0\} \\ \text{Compare}_I &= \{\langle W; P, \sigma \rangle \mid \sigma(j) \geq 0\} \\ \text{Shift}_I &= \{\langle j := \text{next}[j]; W; P, \sigma \rangle\} \end{aligned}$$

where

```

P = while j < lpat and k < ltxt do
  while j >= 0 and pat[j] != txt[k] do
    j := next[j]
  od;
  k := k+1;
  j := j+1
od;
if j >= lpat then return k-j else return -1 fi

```



$$\frac{\frac{n_1 = \sigma(j)}{A, C, \sigma \vdash j \rightarrow_I n_1} \quad \frac{n_2 = C(\text{lpat})}{A, C, \sigma \vdash \text{lpat} \rightarrow_I n_2}}{\text{operate}(>=, n_1, n_2) = \text{true}}}{A, C, \sigma \vdash j >= \text{lpat} \rightarrow_I \text{true}} \\ A, C \vdash \langle \text{if } j >= \text{lpat} \text{ then return } k-j \text{ else return } -1 \text{ fi}, \sigma \rangle \rightarrow_I \langle \text{return } k-j, \sigma \rangle$$

$$\frac{\frac{n_1 = \sigma(k)}{A, C, \sigma \vdash k \rightarrow_I n_1} \quad \frac{n_2 = \sigma(j)}{A, C, \sigma \vdash j \rightarrow_I n_2}}{n = \text{operate}(-, n_1, n_2)}}{A, C, \sigma \vdash k-j \rightarrow_I n} \\ A, C \vdash \langle \text{return } k-j, \sigma \rangle \rightarrow_I \langle n, \sigma \rangle$$

The abstract state corresponding to the initial configuration is  $(\text{match}, j, k)$ . Since  $j \geq C(\text{lpat})$ ,  $(\text{match}, j, k) \rightsquigarrow_I k-j$ , which corresponds to the result in the derivation. Furthermore, the derivation has no comparison subderivations.  $\square$

Since at most one comparison subderivation exists for each step in the derivation, the imperative trace of Definition 4 is a sequence of singleton sets.

### Definition 11 (Abstract trace)

An abstract imperative trace maps a sequence of abstract states to another sequence of abstract states:

$$\text{trace}_I : \text{States}_I^+ \rightarrow \text{States}_I^* \\ \text{trace}_I(s_1 \cdot s_2 \cdots s_n) = \pi(s_1) \cdot \pi(s_2) \cdots \pi(s_n)$$

where  $\pi(s_i) = s_i$  if  $s_i = (\text{compare}, j, k)$  and  $\pi(s_i) = \varepsilon$  otherwise.

The following corollary of Lemma 1 shows that abstract imperative traces represent imperative traces.

### Corollary 1 (Imperative traces are faithful)

Let  $\{(j_1, k_1)\} \cdot \{(j_2, k_2)\} \cdots \{(j_n, k_n)\}$  be the imperative trace for a derivation of the imperative matcher. Let  $(\text{compare}, j'_1, k'_1) \cdot (\text{compare}, j'_2, k'_2) \cdots (\text{compare}, j'_m, k'_m)$  be the abstract imperative trace for the imperative matcher. Then  $n = m$  and  $j_i = j'_i$  and  $k_i = k'_i$  for  $0 < i \leq n$ . In words, the abstract trace faithfully represents the imperative trace.

## 2.7 Summary

We have formally specified an imperative string matcher implementing the KMP algorithm, and we have given it a trace semantics accounting for the indices at which it successively compares characters in the pattern and in the text. In the next section, we turn to a functional string matcher and we treat it similarly.

## 3. The KMP, Functionally

In this section, we describe the functional language in which the functional string matcher is specified. The language is a first-order subset of Scheme (tail-recursive equations). We then present the functional string matcher and its meaning. Finally, we specify a trace semantics of the functional matcher.

### 3.1 Abstract syntax

A program consists of serious expressions  $e \in \text{Exp}$ , trivial expressions  $t \in \text{Triv}$ , operators  $\text{opr} \in \text{Opr}$ , numerals  $\text{num} \in \text{Num}$ , value identifiers  $x \in \text{Vid}$ , function identifiers  $f \in \text{Fid}$  and sequences of value identifiers  $\vec{x} \in \text{Vid}^*$ .

$$p ::= (\text{letrec } ([f_1 (\lambda(\vec{x}_1) e_1)] \dots [f_n (\lambda(\vec{x}_n) e_n)]) e) \\ e ::= t \mid (\text{if } t e_1 e_2) \mid (f t_1 \dots t_m) \\ t ::= \text{num} \mid x \mid (\text{opr } t_1 t_2) \\ \text{opr} ::= + \mid - \mid = \mid \text{eq?} \mid \text{string-ref}$$

### 3.2 Expressible values

A value is an integer, a boolean, a character, or a string:

$$\text{Val} = \mathbb{Z} + \mathbb{B} + \Sigma + \Sigma^*$$

### 3.3 Rules

#### 3.3.1 Auxiliary constructs

The language includes numeric operators, a comparison operator over characters and a string-indexing operator.

$$\text{number}(\text{num}) = n, \text{ if } \text{num} \text{ denotes } n$$

$$\text{operate}(+, n_1, n_2) = n_1 + n_2$$

$$\text{operate}(-, n_1, n_2) = n_1 - n_2$$

$$\text{operate}(=, n_1, n_2) = n_1 = n_2$$

$$\text{operate}(\text{eq?}, c_1, c_2) = c_1 = c_2$$

$$\text{operate}(\text{string-ref}, s, i) = c, \text{ if } c \text{ is the } i\text{'th character in } s.$$

Indexing strings starts at zero, and indexing out of bounds is undefined.

#### 3.3.2 Environments

Expressions are evaluated in a value environment  $\rho \in \text{Venv}$  and a function environment  $\theta \in \text{Fenv}$ :

$$\rho : \text{Vid} \rightarrow \mathbb{Z} + \Sigma^*$$

$$\theta : \text{Fid} \rightarrow \text{Vid}^* \times \text{Exp}$$

#### 3.3.3 Relations

The (big-step) evaluation relation for trivial expressions reads as

$$\rho \vdash t \rightarrow_F v$$

and the (small-step) evaluation relation for serious expressions reads as

$$\theta \vdash \langle e, \rho \rangle \rightarrow_F \langle r, \rho' \rangle$$

where  $\rho, \rho' \in \text{Venv}$ ,  $t \in \text{Triv}$ ,  $v \in \text{Val}$ ,  $\theta \in \text{Fenv}$ ,  $e \in \text{Exp}$ , and  $r \in \text{Exp} + \text{Val}$ .

We choose a big-step evaluation relation for trivial expressions because we are not interested in intermediate evaluation steps. We choose a small-step evaluation relation for serious expressions because we want to monitor the progress of computations.

#### 3.3.4 Programs

At the top level, a program is evaluated in an initial function environment  $\theta_0$  holding the predefined functions and an initial value environment  $\rho_0$  holding the predefined values. The initial configuration of a program

$$(\text{letrec } ([x_1 (\lambda(\vec{x}_1) e_1)] \dots [x_n (\lambda(\vec{x}_n) e_n)]) e)$$

is thus  $\langle e, \rho_0 \rangle$  in the function environment  $\theta$ :

$$\theta = \theta_0 \left[ \begin{array}{l} x_1 \mapsto \langle \vec{x}_1, e_1 \rangle, \\ \dots, \\ x_n \mapsto \langle \vec{x}_n, e_n \rangle \end{array} \right]$$

### 3.3.5 Trivial expressions

$$\frac{n = \text{number}(\text{num})}{\rho \vdash \text{num} \rightarrow_F n}$$

$$\frac{v = \rho(x)}{\rho \vdash x \rightarrow_F v}$$

$$\frac{\rho \vdash t_1 \rightarrow_F v_1 \quad \rho \vdash t_2 \rightarrow_F v_2 \quad v = \text{operate}(\text{opr}, v_1, v_2)}{\rho \vdash (\text{opr } t_1 t_2) \rightarrow_F v}$$

### 3.3.6 Serious expressions

$$\frac{\rho \vdash t \rightarrow_F \text{true}}{\theta \vdash \langle (\text{if } t e_1 e_2), \rho \rangle \rightarrow_F \langle e_1, \rho \rangle}$$

$$\frac{\rho \vdash t \rightarrow_F \text{false}}{\theta \vdash \langle (\text{if } t e_1 e_2), \rho \rangle \rightarrow_F \langle e_2, \rho \rangle}$$

$$\frac{\langle x_1 \dots x_m, e \rangle = \theta(f) \quad \rho \vdash t_1 \rightarrow_F v_1 \quad \dots \quad \rho \vdash t_m \rightarrow_F v_m}{\theta \vdash \langle (f t_1 \dots t_m), \rho \rangle \rightarrow_F \langle e, \rho[x_1 \mapsto v_1, \dots, x_m \mapsto v_m] \rangle}$$

## 3.4 The string matcher

We consider the following string matcher (motivated in appendix), which is written in the subset of Scheme specified in Sections 3.1, 3.2, and 3.3. The initial environment  $\rho_0$  binds `pat` and `lpat` to the pattern and its length, and `txt` and `ltxt` to the text and its length.

```
(letrec ([match
  (lambda (j k)
    (if (= j lpat)
        (- k j)
        (if (= k ltxt)
            -1
            (compare j k))))])
 [compare
  (lambda (j k)
    (if (eq? (string-ref pat j)
              (string-ref txt k))
        (match (+ j 1) (+ k 1))
        (if (= 0 j)
            (match 0 (+ k 1))
            (rematch j k 0 1))))])
 [rematch
  (lambda (j k jp kp)
    (if (= kp j)
        (if (eq? (string-ref pat jp)
                  (string-ref pat kp))
            (if (= jp 0)
                (match 0 (+ k 1))
                (rematch j k 0 (+ (- kp jp) 1)))
            (compare jp k))
        (if (eq? (string-ref pat jp)
                  (string-ref pat kp))
            (rematch j k (+ jp 1) (+ kp 1))
            (rematch j k 0 (+ (- kp jp) 1))))))])
 (match 0 0))
```

None of `pat`, `txt`, `lpat` and `ltxt` are bound in the program, and therefore they denote initial values throughout.

In the rest of this article, we refer to this string matcher as the *functional matcher*. We state without proof that the functional matcher accesses the pattern and the text within their bounds.

## 3.5 Semantics of the functional matcher

We now consider the meaning of the functional matcher. Again, what we are after is the sequence of indices corresponding to the successive comparisons between characters in the pattern and in the text.

### Definition 12 (Comparison)

A functional comparison for the string matcher of Section 3.4 is a derivation tree of the form

$$\frac{T}{\rho \vdash (\text{eq? (string-ref pat j) } \rightarrow_F b \text{ (string-ref txt k)})}$$

where  $T$  denotes another derivation tree.

### Definition 13 (Index)

The following function maps a functional comparison into the corresponding pair of indices in the pattern and the text:

$$\text{index} \left( \frac{T}{\rho \vdash (\text{eq? (string-ref pat j) } \rightarrow_F b \text{ (string-ref txt k)})} \right) = (\rho(j), \rho(k))$$

Let

$$\frac{\frac{\frac{E_1}{\theta \vdash \langle e_1, \rho_1 \rangle \rightarrow_F \langle e_2, \rho_2 \rangle}, \quad \frac{E_2}{\theta \vdash \langle e_2, \rho_2 \rangle \rightarrow_F \langle e_3, \rho_3 \rangle}, \quad \vdots, \quad E_{n-1}}{\theta \vdash \langle e_{n-1}, \rho_{n-1} \rangle \rightarrow_F \langle r, \rho_n \rangle}}$$

be a derivation of the functional matcher, where the premises  $E_1, E_2, \dots, E_{n-1}$  are other derivation trees,  $\theta$  is the initial function environment,  $e_1$  is the functional matcher, and  $\rho_1$  is a value environment mapping `pat`, `txt`, `lpat`, and `ltxt` to the pattern, the text, and their lengths, respectively, and all other value identifiers to zero.

Each premise might contain functional comparisons. We want to build the sequence of indices corresponding to the successive comparisons between characters in the pattern and in the text. Applying the index function to each of the functional comparisons in each premise gives such indices. We collect them in a sequence of non-empty sets of pairs of indices as follows.

### Definition 14 (Trace)

Let  $E_1, E_2, \dots, E_{n-1}$  be the premises of a derivation of the functional matcher. Let  $c_i$  be the set of functional comparisons in  $E_i$ , for  $0 < i < n$ . Let  $p_i = \{\text{index}(c) \mid c \in c_i\}$  for  $0 < i < n$ . The functional trace is the sequence  $\pi(p_1) \cdot \pi(p_2) \cdots \pi(p_{n-1})$ , where

$$\pi(p) = \begin{cases} \varepsilon & \text{if } p = \emptyset \\ p & \text{otherwise.} \end{cases}$$

In Section 3.6, Lemma 2 shows that each of the premises in Definition 14 contains at most one functional comparison. Therefore, for all  $i$ ,  $p_i$  is either empty or a singleton set. The functional trace is thus a sequence of singleton sets, each of which corresponds to the successive comparisons of characters in `pat` and `txt`.

### Definition 15 (Program points)

The functional program points  $Match_F$ ,  $Compare_F$  and  $Rematch_F$  are defined as the following sets of configurations:

$$\begin{aligned} Match_F &= \{\langle M, \rho \rangle\} \\ Compare_F &= \{\langle C, \rho \rangle\} \\ Rematch_F &= \{\langle R, \rho \rangle\} \end{aligned}$$

where  $M$  is the body of the match function,  $C$  is the body of the compare function, and  $R$  is the body of the rematch function.

The set of functional program points is defined as the sum

$$PP_F = Match_F + Compare_F + Rematch_F.$$

### 3.6 Abstract semantics

#### Definition 16 (Abstract states)

The set of abstract functional states is the sum of abstract functional final states and abstract functional intermediate states:

$$\begin{aligned} States_F &= States_F^{fn} + States_F^{int} \\ States_F^{fn} &= \{-1\} + \mathbb{N} \\ States_F^{int} &= (\{\text{match}, \text{compare}\} \times \mathbb{N} \times \mathbb{N}) + \\ &\quad (\text{rematch} \times \mathbb{N} \times \mathbb{N} \times \mathbb{N} \times \mathbb{N}) \end{aligned}$$

where  $\text{match}$ ,  $\text{compare}$  and  $\text{rematch}$  are injection tags.

#### Definition 17 (Program points and states)

We define the correspondence between abstract functional intermediate states and functional program points with the following relation  $\simeq_F \subseteq (States_F^{int} \times PP_F)$ :

$$\begin{aligned} (\text{match}, j, k) &\simeq_F \langle e, \rho \rangle \in Match_F \\ &\quad \text{if } \rho(j) = j \wedge \rho(k) = k \\ (\text{compare}, j, k) &\simeq_F \langle e, \rho \rangle \in Compare_F \\ &\quad \text{if } \rho(j) = j \wedge \rho(k) = k \\ (\text{rematch}, j, k, jp, kp) &\simeq_F \langle e, \rho \rangle \in Rematch_F \\ &\quad \text{if } \rho(j) = j \wedge \rho(k) = k \wedge \\ &\quad \quad \rho(jp) = jp \wedge \rho(kp) = kp \end{aligned}$$

#### Definition 18 (Abstract matcher)

Let  $pat, txt \in \Sigma^*$ . Then the abstract functional matcher is the following total function  $\rightsquigarrow_F \subseteq States_F^{int} \times States_F$ :

$$\begin{aligned} (\text{match}, j, k) &\rightsquigarrow_F \\ \left\{ \begin{array}{ll} (\text{compare}, j, k) & \text{if } j \neq |pat| \wedge k \neq |txt| \\ k-j & \text{if } j = |pat| \\ -1 & \text{otherwise} \end{array} \right. \\ (\text{compare}, j, k) &\rightsquigarrow_F \\ \left\{ \begin{array}{ll} (\text{match}, j+1, k+1) & \text{if } txt[k] = pat[j] \\ (\text{match}, 0, k+1) & \text{if } txt[k] \neq pat[j] \wedge j=0 \\ (\text{rematch}, j, k, 0, 1) & \text{otherwise} \end{array} \right. \\ (\text{rematch}, j, k, jp, kp) &\rightsquigarrow_F \\ \left\{ \begin{array}{ll} (\text{match}, 0, k+1) & \text{if } kp=j \wedge pat[jp]=pat[kp] \\ & \wedge jp=0 \\ (\text{rematch}, j, k, 0, kp-jp+1) & \text{if } kp=j \wedge pat[jp]=pat[kp] \\ & \wedge jp \neq 0 \\ (\text{compare}, jp, k) & \text{if } kp=j \wedge pat[jp] \neq pat[kp] \\ (\text{rematch}, j, k, jp+1, kp+1) & \text{if } kp \neq j \wedge pat[jp]=pat[kp] \\ (\text{rematch}, j, k, 0, kp-jp+1) & \text{otherwise} \end{array} \right. \end{aligned}$$

### Definition 19 (Last)

The function  $last_F$  yields the last element of a non-empty sequence of abstract states:

$$\begin{aligned} last_F &: States_F^+ \rightarrow States_F, \\ last_F(s_1 \cdot s_2 \cdot \dots \cdot s_n) &= s_n \end{aligned}$$

### Definition 20 (Abstract computations)

Let  $pat, txt \in \Sigma^*$  and let  $\rightsquigarrow_F$  be the corresponding abstract functional matcher. Then the set of abstract functional computations,  $AbsComp_F \subseteq States_F^+$ , is the least set closed under

- (1)  $(\text{match}, 0, 0) \in AbsComp_F$
- (2)  $S \in AbsComp_F \wedge last_F(S) \rightsquigarrow_F p \Rightarrow S \cdot p \in AbsComp_F$ .

$S$  is said to be complete iff  $last_F(S) \in \{-1\} + \mathbb{N}$ .

### Lemma 2 (Computations are faithful)

Abstract functional computations represent derivations. In other words:

1. A derivation sequence for the functional matcher starts with an initial derivation sequence containing no comparisons from the initial configuration to a program point  $P \in Match_F$  such that  $(\text{match}, 0, 0) \simeq_F P$ .
2. Whenever there is a derivation sequence from a functional program point  $P'$  to another functional program point or final result  $P''$ , the abstract functional computation represents this derivation in the sense that  $S \simeq_F P'$ ,  $S \rightsquigarrow_F^+ S'$ ,  $S' \simeq_F P''$ , where  $S$  and  $S'$  are the abstract states corresponding to  $P'$  and  $P''$ .

Formally stated:

Let  $\langle e_0, \rho_0 \rangle$  be the initial configuration for a derivation of the functional matcher in the function environment  $\theta$ . Then there is a derivation sequence

$$\begin{aligned} \theta \vdash \langle e_0, \rho_0 \rangle \rightarrow_F \langle e_1, \rho_1 \rangle, \dots, \theta \vdash \langle e_n, \rho_n \rangle \rightarrow_F \langle e', \rho' \rangle \\ \langle e_i, \rho_i \rangle \notin PP_F, \langle e', \rho' \rangle \in Match_F, (\text{match}, 0, 0) \simeq_F \langle e', \rho' \rangle, \\ \text{and no comparison subderivations exist.} \\ \theta \vdash \langle e, \rho \rangle \rightarrow_F \langle e_1, \rho_1 \rangle, \dots, \theta \vdash \langle e_n, \rho_n \rangle \rightarrow_F \langle e', \rho' \rangle \\ \langle e_i, \rho_i \rangle \notin PP_I, \langle e, \rho \rangle \in Match_F, \langle e', \rho' \rangle \in PP_F \Rightarrow \\ (\text{match}, j, k) \simeq_F \langle e, \rho \rangle, (\text{match}, j, k) \rightsquigarrow_I m', m' \simeq_F \langle e', \rho' \rangle, \\ \text{and no comparison subderivations exist.} \\ \theta \vdash \langle e, \rho \rangle \rightarrow_F \langle e_1, \rho_1 \rangle, \dots, \theta \vdash \langle e_n, \rho_n \rangle \rightarrow_F \langle n, \rho' \rangle \\ \langle e_i, \rho_i \rangle \notin PP_I, \langle e, \rho \rangle \in Match_F, n \in \mathbb{N} \Rightarrow (\text{match}, j, k) \simeq_F \\ \langle e, \rho \rangle, (\text{match}, j, k) \rightsquigarrow_I n, \text{ and no comparison subderivations exist.} \\ \theta \vdash \langle e, \rho \rangle \rightarrow_F \langle e_1, \rho_1 \rangle, \dots, \theta \vdash \langle e_n, \rho_n \rangle \rightarrow_F \langle e', \rho' \rangle \\ \langle e_i, \rho_i \rangle \notin PP_I, \langle e, \rho \rangle \in Compare_F, \langle e', \rho' \rangle \in PP_F \Rightarrow \\ (\text{compare}, j, k) \simeq_F \langle e, \rho \rangle, (\text{compare}, j, k) \rightsquigarrow_I m', m' \simeq_F \langle e', \rho' \rangle, \\ \text{and exactly one comparison subderivation exists, using environment } \rho. \\ \theta \vdash \langle e, \rho \rangle \rightarrow_F \langle e_1, \rho_1 \rangle, \dots, \theta \vdash \langle e_n, \rho_n \rangle \rightarrow_F \langle e', \rho' \rangle \\ \langle e_i, \rho_i \rangle \notin PP_I, \langle e, \rho \rangle \in Rematch_F, \langle e', \rho' \rangle \in PP_F \Rightarrow \\ (\text{rematch}, j, k, jp, kp) \simeq_F \langle e, \rho \rangle, (\text{rematch}, j, k, jp, kp) \rightsquigarrow_I m', \\ m' \simeq_F \langle e', \rho' \rangle, \text{ and no comparison subderivations exist.} \end{aligned}$$

PROOF. We show one case for a derivation of the body of match in a value environment  $\rho$  where  $\rho(j) \neq \rho(1pat)$  and  $\rho(k) \neq \rho(1txt)$ . The other cases are similar.

The derivation tree is

$$\begin{array}{c}
\frac{n_1 = \rho(j) \quad n_2 = \rho(\text{lpat}) \quad \text{operate}(=, n_1, n_2) = \text{false}}{\rho \vdash j \rightarrow_F n_1 \quad \rho \vdash \text{lpat} \rightarrow_F n_2} \\
\hline
\rho \vdash (= j \text{ lpat}) \rightarrow_F \text{false} \\
\hline
\theta \vdash \left\langle \left( \begin{array}{l} \text{(if } (= j \text{ lpat}) \\ \text{(- k j}) \\ \text{(if } (= k \text{ ltxt}) \text{ -1 (compare j k))} \end{array} \right), \rho \right\rangle \\
\rightarrow_F \left\langle \text{(if } (= k \text{ ltxt}) \text{ -1 (compare j k)), } \rho \right\rangle \\
\hline
\frac{n_1 = \rho(k) \quad n_2 = \rho(\text{ltxt}) \quad \text{operate}(=, n_1, n_2) = \text{false}}{\rho \vdash k \rightarrow_F n_1 \quad \rho \vdash \text{ltxt} \rightarrow_F n_2} \\
\hline
\rho \vdash (= k \text{ ltxt}) \rightarrow_F \text{false} \\
\hline
\theta \vdash \left\langle \text{(if } (= k \text{ ltxt}) \text{ -1 (compare j k)), } \rho \right\rangle \\
\rightarrow_F \left\langle \text{(compare j k), } \rho \right\rangle \\
\hline
\frac{n_1 = \rho(j) \quad n_2 = \rho(k) \quad \theta(\text{compare}) = \langle j \cdot k, \mathbf{C} \rangle}{\rho \vdash j \rightarrow_F n_1 \quad \rho \vdash k \rightarrow_F n_2} \\
\hline
\theta \vdash \langle \text{(compare j k), } \rho \rangle \rightarrow_F \langle \mathbf{C}, \rho[j \mapsto n_1, k \mapsto n_2] \rangle
\end{array}$$

where  $\mathbf{C}$  denotes the body of the `compare` function, as in Definition 15.

The abstract state corresponding to the initial configuration is  $(\text{match}, j, k)$ , where  $j = \rho(j)$  and  $k = \rho(k)$ . Since  $j \neq \rho(\text{lpat})$  and  $k \neq \rho(\text{ltxt})$ ,  $(\text{match}, j, k) \rightsquigarrow_F (\text{compare}, j, k)$ . Since  $n_1 = \rho(j) = j$  and  $n_2 = \rho(k) = k$ ,  $(\text{compare}, j, k)$  corresponds to the final configuration in the derivation.  $\square$

Since at most one comparison subderivation exists for each step in the derivation, the functional trace of Definition 14 is a sequence of singleton sets.

### Definition 21 (Abstract trace)

An abstract functional trace maps a sequence of abstract states to another sequence of abstract states:

$$\begin{array}{l}
\text{trace}_F : \text{States}_F^+ \rightarrow \text{States}_F^* \\
\text{trace}_F(s_1 \cdot s_2 \cdot \dots \cdot s_n) = \pi(s_1) \cdot \pi(s_2) \cdot \dots \cdot \pi(s_n)
\end{array}$$

where  $\pi(s_i) = s_i$  if  $s_i = (\text{compare}, j, k)$  and  $\pi(s_i) = \varepsilon$  otherwise.

The following corollary of Lemma 2 shows that abstract functional traces represent functional traces.

### Corollary 2 (Functional traces are faithful)

Let  $\{(j_1, k_1)\} \cdot \{(j_2, k_2)\} \cdot \dots \cdot \{(j_n, k_n)\}$  be the functional trace for a derivation of the functional matcher. Let  $(\text{compare}, j'_1, k'_1) \cdot (\text{compare}, j'_2, k'_2) \cdot \dots \cdot (\text{compare}, j'_m, k'_m)$  be the abstract trace for the functional matcher. Then  $n = m$  and  $j_i = j'_i$  and  $k_i = k'_i$  for  $0 < i \leq n$ . In words, the abstract trace faithfully represents the functional trace.

### Lemma 3 (Invariants)

Let  $\text{pat}, \text{txt} \in \Sigma^*$  and  $\text{AbsComp}_F$  be the corresponding set of abstract functional computations. Then for all  $s_1 \cdot s_2 \cdot \dots \cdot s_n \in \text{AbsComp}_F$ , the following conditions, whose conclusions we call invariants, are satisfied:

- If  $s_i = (\text{match}, j, k)$  then
  - (m1)  $0 \leq j \leq |\text{pat}|$
  - (m2)  $k \leq |\text{txt}|$

- If  $s_i = (\text{compare}, j, k)$  then

$$\begin{array}{l}
(c1) \quad 0 \leq j < |\text{pat}| \\
(c2) \quad k < |\text{txt}|
\end{array}$$

- If  $s_i = (\text{rematch}, j, k, jp, kp)$  then

$$\begin{array}{l}
(r1) \quad 0 < j < |\text{pat}| \\
(r2) \quad k < |\text{txt}| \\
(r3) \quad 0 \leq jp < kp \leq j \\
(r4) \quad \text{pat}[0] \cdot \dots \cdot \text{pat}[jp-1] = \text{pat}[kp-jp] \cdot \dots \cdot \text{pat}[kp-1] \\
(r5) \quad \forall \underline{k} \in [1, kp-jp-1]. \\
\quad \neg(\text{pat}[0] \cdot \dots \cdot \text{pat}[j-\underline{k}-1] = \text{pat}[\underline{k}] \cdot \dots \cdot \text{pat}[j-1] \wedge \\
\quad \text{pat}[j-\underline{k}] \neq \text{pat}[j])
\end{array}$$

PROOF. See extended version [1].  $\square$

The key connection between the abstract functional matcher and the abstract imperative matcher is stated in the following remark. The remark shows how to interpret Invariant (r5) in terms of the next table.

**Remark 1** We notice that for any  $j$  and  $0 \leq a \leq b$ , if  $\forall \underline{k} \in [a, b]. \neg(\text{pat}[0] \cdot \dots \cdot \text{pat}[j-\underline{k}-1] = \text{pat}[\underline{k}] \cdot \dots \cdot \text{pat}[j-1] \wedge \text{pat}[j-\underline{k}] \neq \text{pat}[j])$ , then by Definition 1  $\text{next}[j]$  cannot occur in the interval  $[j-b, j-a]$ .

Indeed, if for some  $\underline{k}$  and some  $j$ ,  $(\text{pat}[0] \cdot \dots \cdot \text{pat}[j-\underline{k}-1] = \text{pat}[\underline{k}] \cdot \dots \cdot \text{pat}[j-1]$  and  $\text{pat}[j-\underline{k}] \neq \text{pat}[j])$ , then  $j-\underline{k}$  is a candidate for  $\text{next}[j]$ . Therefore the negation of the condition gives us that  $j-\underline{k}$  is not a candidate for  $\text{next}[j]$ .

## 3.7 Summary

We have formally specified a functional string matcher, and we have given it a trace semantics accounting for the indices at which it successively compares characters in the pattern and in the text. In the next section, we show that for any given pattern and text, the traces of the imperative matcher and of the functional matcher coincide.

## 4. Extensional correspondence between imperative and functional matchers

### Definition 22 (Correspondence)

We define the correspondence between imperative and functional states with the relation  $\simeq_{\subseteq} \text{States}_I \times \text{States}_F$ :

$$\begin{array}{l}
-1 \simeq -1 \\
n \simeq n', \quad \text{if } n = n' \\
(\text{match}, j, k) \simeq (\text{match}, j', k') \quad \text{if } j = j' \wedge k = k' \\
(\text{compare}, j, k) \simeq (\text{compare}, j', k') \quad \text{if } j = j' \wedge k = k' \\
(\text{shift}, j, k) \simeq (\text{rematch}, j', k', jp, kp) \quad \text{if } j = j' \wedge k = k'
\end{array}$$

We define  $\simeq^*_{\subseteq} \text{States}_I^* \times \text{States}_F^*$  such that for any sequences  $S = s_1 \cdot s_2 \cdot \dots \cdot s_p \in \text{States}_I^+$  and  $S' = s'_1 \cdot s'_2 \cdot \dots \cdot s'_q \in \text{States}_F^+$ ,  $S \simeq^* S'$  iff  $p = q$  and  $s_i \simeq s'_i$  for all  $0 < i \leq p$ . We make  $\simeq^*$  hold for empty sequences.

### Definition 23 (Synchronization)

Synchronization is a relation  $\text{sync} \subseteq \text{States}_I^+ \times \text{States}_F^+$  defined as

$$\text{sync}(S, S') \text{ iff } \text{trace}_I(S) \simeq^* \text{trace}_F(S') \wedge \text{last}_I(S) \simeq \text{last}_F(S')$$

### Theorem 1 (Abstract equivalence)

For any given pattern and text, there is a unique complete abstract imperative computation  $S$  and a unique complete abstract functional computation  $S'$ , and these two abstract computations are synchronized, i.e.,  $\text{sync}(S, S')$  holds.

PROOF. See extended version [1].  $\square$

### Corollary 3 (Equivalence)

Given the next table for the pattern string, the imperative matcher of Section 2.4 and the functional matcher of Section 3.4 give rise to the same sequence of comparisons between pattern and text, and yield the same result.

PROOF. By Corollary 1 and Corollary 2, an abstract trace corresponds to an actual trace; by Theorem 1, the abstract computations are synchronized; and by Lemma 1 and Lemma 2, the abstract computations represent the derivations. In particular, the result of the abstract computations is the same.  $\square$

## 5. Intensional correspondence between imperative and functional matchers

We now turn to specializing the functional string matcher with respect to given patterns. First we use partial evaluation (i.e., program specialization), and next we consider a simple form of data specialization. We first show that the size of the specialized programs is linear in the size of the pattern, and that the specialized programs run in time linear in the size of the text. We next show that the specialized data coincides with the next table of the KMP.

This section is more informal and makes a somewhat liberal use of partial-evaluation terminology [22].

### 5.1 Program specialization

Figure 1 displays a binding-time annotated version of the functional matcher. Formal parameters are tagged with “s” (for “static”) or “d” (for “dynamic”) depending on whether they only denote values that depend on data available at partial-evaluation time or whether they denote values that may depend on data available at run time. In addition, dynamic conditional expressions, dynamic tests, and dynamic additions and subtractions are boxed. All the other parts in the source program are static and will be evaluated at partial-evaluation time. All the dynamic parts will be reconstructed, giving rise to the residual program.

A partial evaluator such as Similix [4, 5] is designed to preserve dynamic computations and their order. In the present case, the dynamic tests are among the dynamic computations. They are guaranteed to occur in specialized programs in the same order as in the source program. Therefore, by construction, Similix generates programs that traverse the text in the same order as the functional matcher and thus the KMP algorithm.

For example, we have specialized the functional matcher with respect to the pattern “abac” (without post-unfolding). The resulting residual program is displayed in Figure 2, after lambda-dropping [9] and renaming (the character following the “[” in the subscripts, is the next character in the pattern to be matched against the text—an intuitive notation suggested by Grobauer and Lawall [14]). The specialized string matcher traverses the text linearly and compares characters in the text and literal characters from the pattern. In their

```
(define (main pats txtd)
  (let ([lpats (string-length pat)]
        [ltxtd (string-length txt)])
    (letrec
      ([match
        (lambda (js kd)
          (if (= j lpat)
              (boxed -) k j
              (boxed if) (boxed (=) k ltxt)
              (compare j k)))]
       [compare
        (lambda (js kd)
          (boxed if) (boxed eq?) (string-ref pat j)
                    (string-ref txt k))
          (match (+ j 1) (boxed +) k 1))
          (if (= 0 j)
              (match 0 (boxed +) k 1)
              (rematch j k 0 1)))]
       [rematch
        (lambda (js kd jps kps)
          (if (= kp j)
              (if (eq? (string-ref pat jp)
                       (string-ref pat kp))
                  (if (= jp 0)
                      (match 0 (boxed +) k 1)
                      (rematch j k 0 (+ (- kp jp) 1)))
                  (compare jp k))
              (if (eq? (string-ref pat jp)
                       (string-ref pat kp))
                  (rematch j k (+ jp 1) (+ kp 1))
                  (rematch j k 0 (+ (- kp jp) 1)))))]
          (match 0 0)))]
      (match 0 0))))
```

Figure 1: The binding-time annotated functional matcher

article [19, page 330], Knuth, Morris and Pratt display a similar program where the next table has been “compiled” into the control flow. We come back to this point at the end of Section 5.2.

In their revisitation of partial evaluation of pattern matching in strings [14], Grobauer and Lawall analyzed the size and complexity of the residual code produced by Similix, measured in terms of the number of residual tests. They showed that the size of a residual program is linear in the length of the pattern, and that the time complexity is linear in the length of the text. In the same manner, we can show that Similix yields residual programs that are linear in the length of the patterns, and whose time complexity is linear in the length of the text.

Similix is a polyvariant program-point specializer that builds mutually recursive specialized versions of source program points (by default: conditional expressions with dynamic tests). Each source program point is specialized with respect to a set of static values. The corresponding residual program point is indexed with this set. If a source program point is met again with the same set of static values, a residual call to the corresponding residual program point is generated.

### Proposition 1

Specializing the functional matcher of Figure 1 with respect to a pattern yields a residual program whose size is linear in the length of the pattern.

```

(define (main-abac txt)
  (let ([ltx (string-length txt)])
    (define (match|abac k)
      (if (= k ltx) -1 (compare|abac k)))
    (define (compare|abac k)
      (if (eq? #\a (string-ref txt k))
          (matcha|bac (+ k 1))
          (match|abac (+ k 1))))
    (define (matcha|bac k)
      (if (= k ltx) -1 (comparea|bac k)))
    (define (comparea|bac k)
      (if (eq? #\b (string-ref txt k))
          (matchab|ac (+ k 1))
          (compare|abac k)))
    (define (matchab|ac k)
      (if (= k ltx) -1 (compareab|ac k)))
    (define (compareab|ac k)
      (if (eq? #\a (string-ref txt k))
          (matchaba|c (+ k 1))
          (match|abac (+ k 1))))
    (define (matchaba|c k)
      (if (= k ltx) -1 (compareaba|c k)))
    (define (compareaba|c k)
      (if (eq? #\c (string-ref txt k))
          (- (+ k 1) 4)
          (comparea|bac k)))
    (match|abac 0)))

```

- For all `txt`, evaluating `(main-abac txt)` yields the same result as evaluating `(main "abac" txt)`.
- For all `k`, evaluating `(match|abac k)` in the scope of `ltx` yields the same result as evaluating `(match 0 k)` in the scope of `lpat` and `ltx`, where `lpat` denotes the length of `pat` and `ltx` denotes the length of `txt`.
- For all `k`, evaluating `(matcha|bac k)` in the scope of `ltx` yields the same result as evaluating `(match 1 k)` in the scope of `lpat` and `ltx`.
- For all `k`, evaluating `(matchab|ac k)` in the scope of `ltx` yields the same result as evaluating `(match 2 k)` in the scope of `lpat` and `ltx`.
- For all `k`, evaluating `(matchaba|c k)` in the scope of `ltx` yields the same result as evaluating `(match 3 k)` in the scope of `lpat` and `ltx`.

Figure 2: Result of specializing the functional matcher wrt. "abac"

PROOF (INFORMAL). The only functions for which residual code is generated are `main`, `match` and `compare`. The first one, `main`, is the goal function, but it contains no memoization points, so only one residual `main` function is generated. There is exactly one memoization point—a dynamic conditional expression—in each of the functions `match` and `compare`. The only static data available at the two memoization points are bound to `j`, `pat`, and `lpat`. The only piece of static data that varies is the value of `j`, i.e.,  $j$ , and since  $0 \leq j < |pat|$  at the memoization points (because of the invariants of Lemma 3 in Section 4, and the fact that the memoization point in `match` is only reached if  $j \neq |pat|$ ), at most  $|pat|$  variants of the two memoization points can be generated. The number of residual functions is therefore linear in the size of the pattern. In addition, the size of each function is bounded by a small constant, as can be seen if one writes the BNF of residual programs [21].  $\square$

## Proposition 2

*Specializing the functional matcher of Figure 1 with respect to a pattern yields a residual program whose time complexity is linear in the length of the text.*

PROOF (INFORMAL). As proven by Knuth, Morris and Pratt, the KMP algorithm performs a number of comparisons between characters in the pattern and in the text, that is linear in the length of the text [19]. Corollary 3 shows that the functional matcher performs the exact same sequence of comparisons between characters in the pattern and in the text as the KMP algorithm. All comparisons are performed in the `compare` function, and exactly one comparison is performed at each call to `compare`. The number of calls to `compare` is therefore linear in the length of the text, and since the `match` function either terminates or calls `compare`, the number of calls to `match` is bounded by the number of calls to `compare`. By Proposition 1, residual code is only generated for the functions `main`, `compare`, and `match`. The time complexity of each of the functions `main`, `compare`, and `match` is easily seen to be bounded by a small constant. Since `main` is only called once and the number of calls to `compare` and `match` is linear in the length of the text, the time complexity of the residual program is linear in the length of the text.  $\square$

## 5.2 Data specialization

In Section 3.6, Remark 1 connects the `rematch` function in the functional matcher and the next table of the KMP algorithm. In this section, we revisit this connection and show how to actually derive the KMP algorithm with a next table from the functional matcher using a simple form of data specialization [3, 7, 18, 20]. To this end, we first restate the functional matcher.

In the functional matcher, all functions are tail recursive, i.e., they iteratively call themselves or each other. In particular, `rematch` completes either by calling `match` or by calling `compare`. The two actual parameters to `match` are `0`, a literal, and an increment over `k`, which is available in the scope of `match`. The two actual parameters to `compare` are `jp`, which has been computed in the course of `rematch`, and `k`, which is available in the scope of `compare`.

To make it possible to tabulate the `rematch` function, we modify the functional matcher so that it is no longer tail recursive. Instead of having `rematch` call `match` or `compare`, tail recursively, we make it return a value on which to call `match` or `compare`. We set this value to be that of `jp` (a natural number) or `-1`. Correspondingly, instead of having `compare` call `rematch` tail recursively, we make it dispatch on the result of `rematch` to call `match` or `compare`, tail recursively. The result is displayed in Figure 3.

In the proof of Theorem 1 [1], we show that when `rematch` terminates by calling `compare`, `jp` is equal to `next[j]` in the KMP algorithm. We also show that when `match` is called from `rematch`, the value `next[j]` in the KMP algorithm is `-1`. We only call `rematch` from `compare`, and only with  $0 \leq j < |pat|$ , `jp` = `0`, and `kp` = `1`. Therefore calling the new `rematch` function is equivalent to a lookup in the next table in the KMP algorithm. In particular, tabulating the  $|pat|$  input values of `rematch` corresponding to all  $j$  between `0` and  $|pat| - 1$  yields the next table as used in the KMP algorithm.

This simple data specialization yields a string matcher that traverses the text linearly, matching it against the pattern, and looking up the next index into the pattern in the

```

(define (main pat txt)
  (let ([lpat (string-length pat)]
        [ltx (string-length txt)])
    (letrec
      ([match
       (lambda (j k)
         (if (= j lpat)
             (- k j)
             (if (= k ltx)
                 -1
                 (compare j k))))])
      [compare
       (lambda (j k)
         (if (eq? (string-ref pat j)
                  (string-ref txt k))
             (match (+ j 1) (+ k 1))
             (if (= 0 j)
                 (match 0 (+ k 1))
                 (let ([next (rematch j 0 1)])
                   (if (= next -1)
                       (match 0 (+ k 1))
                       (compare next k))))))])
      [rematch
       (lambda (j jp kp)
         (if (= kp j)
             (if (eq? (string-ref pat jp)
                      (string-ref pat kp))
                 (if (= jp 0)
                     -1
                     (rematch j 0 (+ (- kp jp) 1)))
                 jp)
             (if (eq? (string-ref pat jp)
                      (string-ref pat kp))
                 (rematch j (+ jp 1) (+ kp 1))
                 (rematch j 0 (+ (- kp jp) 1))))))])
      (match 0 0))))))

```

Figure 3: Variation on the functional matcher

next table in case of mismatch. In other words, data specialization of the functional matcher yields the KMP algorithm.

In particular, specializing the string matcher of Figure 3 (or its tabulated version) with respect to a pattern would compile the corresponding next table into the control flow of the residual program. The result would coincide with the compiled code in Knuth, Morris and Pratt’s article [19, page 330].

## 6. Conclusion and issues

We have presented the first formal proof that partial evaluation can precisely yield the KMP, both extensionally (trace semantics, synchronization) and intensionally (size of specialized programs, relation to the next table, actual derivation of the KMP algorithm). We have shown that the key to obtaining the KMP out of a naive, quadratic string matcher is not only to keep backtracking under static control, but also to maintain exactly one character of negative information, as in Consel and Danvy’s original solution. Together with Grobauer and Lawall’s complexity proofs about the size and time complexity of residual programs, the buildup of Corollary 3 paves the way to relating the effect of staged string matchers with independently known string matchers, e.g., Boyer and Moore’s [2].

Our work has led us to consider a family of KMP algorithms in relation with the following family of staged string matchers:

- A staged string matcher that does not keep track of negative information gives rise not to Knuth, Morris, and Pratt’s next table, but to their  $f$  function [19, page 327], i.e., to Morris and Pratt’s algorithm [6, Chapter 6]. Tabulating this function yields an array of the same size as the pattern.
- A staged string matcher that keeps track of one character of negative information corresponds to Knuth, Morris, and Pratt’s algorithm and next table.
- A staged string matcher that keeps track of a limited number of characters of negative information gives rise to a KMP-like algorithm. The corresponding residual programs are more efficient, but they are also bigger.
- A staged string matcher that keeps track of all the characters of negative information also gives rise to a KMP-like algorithm. The corresponding residual programs are even more efficient, but they are also even bigger. Grobauer and Lawall have shown that the size of these residual programs is bounded by  $|pat| \times |\Sigma|$ , where  $|\Sigma|$  is the size of the alphabet [14].

It is however our conjecture that for string matchers that keep track of two or more characters of negative information, a tighter upper bound on the size is twice the length of the pattern, i.e.,  $2|pat|$ . This conjecture holds for short patterns.

Let us conclude on two points: obtaining *efficient* string matchers by partial evaluation of a naive string matcher and obtaining them *efficiently*.

The essence of obtaining efficient string matchers by partial evaluation of a naive string matcher is to ensure that backtracking in the naive matcher is static. One can then either stage the naive matcher and use a simple partial evaluator, or keep the naive matcher unstaged and use a sophisticated partial evaluator. What matters is that backtracking is carried out at specialization time and that dynamic computations are preserved in specialized programs.

The size of residual programs provides a lower bound to the time complexity of specialization. For example, looking at the KMP, the size of a residual program is proportional to the size of the pattern if only positive information is kept. At best, a general-purpose partial evaluator could thus proceed in time linear in  $|pat|$ , i.e.,  $O(|pat|)$ , as in the first pass of the KMP algorithm. However, evaluating the static parts of the source program at specialization time, as driven by the static control flow of the source program, does not seem like an optimal strategy, even discounting the complexity of binding-time analysis. For example, the data specialization in Section 5.2 works in time quadratic in  $|pat|$ , i.e.,  $O(|pat|^2)$ , to construct the next table. On the other hand, such an efficient treatment could be one of the bullets in a partial evaluator’s gun [23, Section 11], i.e., a treatment that is not generally applicable but has a dramatic effect occasionally. For example, proving the conjecture above could lead to such a bullet.

**Acknowledgments.** We are grateful to Torben Amtoft, Julia Lawall, Karoline Malmkjær, Jan Midtgaard, Mikkel Nygaard, and the anonymous reviewers for a variety of comments. Special thanks to Andrzej Filinski for further comments that led us to reshape this article.

This work is supported by the ESPRIT Working Group APPSEM (<http://www.md.chalmers.se/Cs/Research/Semantics/APPSEM/>).

## APPENDIX

### Staging a quadratic string matcher

Figure 4 displays a naive, quadratic string matcher that successively checks whether the pattern `pat` is a prefix of one of the successive suffixes of the text `txt`. The `main` function initializes the indices `j` and `k` with which to access `pat` and `txt`. The `match` function checks whether the matching is finished (either with a success or with a failure), or whether one more comparison is needed. The `compare` function carries out this comparison. Either it continues to match the rest of `pat` with the rest of the current suffix of `txt` or it starts to match `pat` and the next suffix of `txt`.

Figure 5 displays a staged version of the quadratic string matcher. Instead of matching `pat` and the next suffix of `txt`, this version uses a `rematch` function and a `recompare` function to first match `pat` and a prefix of a suffix of `pat`, which we know to be equal to the corresponding segment in `txt`. Eventually, the `rematch` function resumes matching the rest of the pattern and the rest of `txt`. As a result, the staged string matcher does not backtrack on `txt`.

In partial-evaluation jargon, the string matcher of Figure 5 uses positive information about the text. A piece of negative information is also available, namely the latest character having provoked a mismatch. Figure 6 displays a staged version of the quadratic string matcher that exploits this negative information. Rather than blindly resuming the `compare` function, the `rematch` function first checks whether the character having caused the latest mismatch could cause a new mismatch, thereby avoiding one access to the text.

To simplify the formal development, we inline `recompare` in `rematch` and lambda-lift `rematch` to the same lexical level as `match` and `compare` [9, 15]. The resulting string matcher is displayed in Figure 7 and in Section 3.4.

There are of course many ways to stage a string matcher. The one we have chosen is easy to derive and easy to reason about.

## References

- [1] Mads Sig Ager, Olivier Danvy, and Henning Korsholm Rohde. On obtaining Knuth, Morris, and Pratt's string matcher by partial evaluation. Technical Report BRICS RS-02-32, Department of Computer Science, University of Aarhus, Aarhus, Denmark, July 2002.
- [2] Torben Amtoft, Charles Consel, Olivier Danvy, and Karoline Malmkjær. The abstraction and instantiation of string-matching programs. Technical Report BRICS RS-01-12, DAIMI, Department of Computer Science, University of Aarhus, Aarhus, Denmark, April 2001. To appear in Neil Jones's Festschrift.
- [3] Guntis J. Barzdins and Mikhail A. Bulyonkov. Mixed computation and translation: Linearisation and decomposition of compilers. Preprint 791, Computing Centre of Siberian Division of USSR Academy of Sciences, Novosibirsk, Siberia, 1988.
- [4] Anders Bondorf. Similix 5.1 manual. Technical report, DIKU, Computer Science Department, University of Copenhagen, Copenhagen, Denmark, May 1993. Included in the Similix 5.1 distribution.
- [5] Anders Bondorf and Olivier Danvy. Automatic autoprojection of recursive equations with global variables and abstract data types. *Science of Computer Programming*, 16:151–195, 1991.
- [6] Christian Charras and Thierry Lecroq. Exact string matching algorithms. <http://www-igm.univ-mlv.fr/~lecroq/string/>, 1997.
- [7] Sandrine Chirokoff, Charles Consel, and Renaud Marlet. Combining program and data specialization. *Higher-Order and Symbolic Computation*, 12(4):309–335, 1999.
- [8] Charles Consel and Olivier Danvy. Partial evaluation of pattern matching in strings. *Information Processing Letters*, 30(2):79–86, January 1989.
- [9] Olivier Danvy and Ulrik P. Schultz. Lambda-dropping: Transforming recursive equations into programs with block structure. *Theoretical Computer Science*, 248(1-2):243–287, 2000.
- [10] Yoshihiko Futamura, Zenjiro Konishi, and Robert Glück. Program transformation system based on generalized partial computation. *New Generation Computing*, 20(1):75–99, 2002.
- [11] Yoshihiko Futamura and Kenroku Nogi. Generalized partial computation. In Dines Bjørner, Andrei P. Ershov, and Neil D. Jones, editors, *Partial Evaluation and Mixed Computation*, pages 133–151. North-Holland, 1988.
- [12] Yoshihiko Futamura, Kenroku Nogi, and Akihiko Takano. Essence of generalized partial computation. *Theoretical Computer Science*, 90(1):61–79, 1991.
- [13] Robert Glück and Andrei Klimov. Occam's razor in metacomputation: the notion of a perfect process tree. In Patrick Cousot, Moreno Falaschi, Gilberto Filé, and Antoine Rauzy, editors, *Proceedings of the Third International Workshop on Static Analysis WSA'93*, number 724 in Lecture Notes in Computer Science, pages 112–123, Padova, Italy, September 1993. Springer-Verlag.
- [14] Bernd Grobauer and Julia L. Lawall. Partial evaluation of pattern matching in strings, revisited. *Nordic Journal of Computing*, 8(4):437–462, 2002.
- [15] Thomas Johnsson. Lambda lifting: Transforming programs to recursive equations. In Jean-Pierre Jouannaud, editor, *Functional Programming Languages and Computer Architecture*, number 201 in Lecture Notes in Computer Science, pages 190–203, Nancy, France, September 1985. Springer-Verlag.
- [16] Neil D. Jones, Carsten K. Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice-Hall International, London, UK, 1993. Available online at <http://www.dina.kvl.dk/~sestoft/pebook/>.
- [17] Richard Kelsey, William Clinger, and Jonathan Rees, editors. Revised<sup>5</sup> report on the algorithmic language Scheme. *Higher-Order and Symbolic Computation*, 11(1):7–105, 1998.
- [18] Todd B. Knoblock and Erik Ruf. Data specialization. In *Proceedings of the ACM SIGPLAN'96 Conference on Programming Languages Design and Implementation*, SIGPLAN Notices, Vol. 31, No 5,

```

(define (main pat txt)
  (let ([lpat (string-length pat)] [ltxt (string-length txt)])
    (letrec ([match (lambda (j k)
                  (if (= j lpat)
                      (- k j)
                      (if (= k ltxt)
                          -1
                          (compare j k))))]
      [compare (lambda (j k)
                 (if (eq? (string-ref pat j) (string-ref txt k))
                     (match (+ j 1) (+ k 1))
                     (match 0 (+ (- k j) 1))))])
      (match 0 0))))

```

Figure 4: The naive, quadratic functional matcher

```

(define (main pat txt)
  (let ([lpat (string-length pat)] [ltxt (string-length txt)])
    (letrec ([match (lambda (j k)
                  (if (= j lpat)
                      (- k j)
                      (if (= k ltxt)
                          -1
                          (compare j k))))]
      [compare (lambda (j k)
                 (if (eq? (string-ref pat j) (string-ref txt k))
                     (match (+ j 1) (+ k 1))
                     (if (= 0 j)
                         (match 0 (+ k 1))
                         (letrec ([rematch (lambda (jp jk)
                                             (if (= jk j)
                                                 (compare jp k)
                                                 (recompare jp jk)))]
                           [recompare (lambda (jp jk)
                                         (if (eq? (string-ref pat jp) (string-ref pat jk))
                                             (rematch (+ jp 1) (+ jk 1))
                                             (rematch 0 (+ (- jk jp) 1)))))]
                         (rematch 0 1)))))]
      (match 0 0))))

```

Figure 5: The functional matcher with positive information

```

(define (main pat txt)
  (let ([lpat (string-length pat)] [ltxt (string-length txt)])
    (letrec ([match (lambda (j k)
                  (if (= j lpat)
                      (- k j)
                      (if (= k ltxt)
                          -1
                          (compare j k))))]
      [compare (lambda (j k)
                 (if (eq? (string-ref pat j) (string-ref txt k))
                     (match (+ j 1) (+ k 1))
                     (if (= 0 j)
                         (match 0 (+ k 1))
                         (letrec ([rematch (lambda (jp kp)
                                             (if (= kp j)
                                                 (if (eq? (string-ref pat jp) (string-ref pat kp))
                                                     (if (= jp 0)
                                                         (match 0 (+ k 1))
                                                         (rematch 0 (+ (- kp jp) 1)))
                                                     (compare jp k))
                                                 (recompare jp kp)))]
                           [recompare (lambda (jp kp)
                                         (if (eq? (string-ref pat jp) (string-ref pat kp))
                                             (rematch (+ jp 1) (+ kp 1))
                                             (rematch 0 (+ (- kp jp) 1)))))]
                         (rematch 0 1)))))]
      (match 0 0))))

```

Figure 6: The functional matcher with positive information and one character of negative information

```

(define (main pat txt)
  (let ([lpat (string-length pat)] [ltxt (string-length txt)])
    (letrec ([match (lambda (j k)
                  (if (= j lpat)
                      (- k j)
                      (if (= k ltxt)
                          -1
                          (compare j k))))])
      [compare (lambda (j k)
                 (if (eq? (string-ref pat j) (string-ref txt k))
                     (match (+ j 1) (+ k 1))
                     (if (= 0 j)
                         (match 0 (+ k 1))
                         (rematch j k 0 1))))])
      [rematch (lambda (j k jp kp)
                 (if (= kp j)
                     (if (eq? (string-ref pat jp) (string-ref pat kp))
                         (if (= jp 0)
                             (match 0 (+ k 1))
                             (rematch j k 0 (+ (- kp jp) 1)))
                         (compare jp k))
                     (if (eq? (string-ref pat jp) (string-ref pat kp))
                         (rematch j k (+ jp 1) (+ kp 1))
                         (rematch j k 0 (+ (- kp jp) 1))))))])
      (match 0 0))))

```

Figure 7: The functional matcher with positive information and one character of negative information (final version)

- pages 215–225. ACM Press, June 1996.
- [19] Donald E. Knuth, James H. Morris, and Vaughan R. Pratt. Fast pattern matching in strings. *SIAM Journal on Computing*, 6(2):323–350, 1977.
- [20] Karoline Malmkjær. Program and data specialization: Principles, applications, and self-application. Master’s thesis, DIKU, Computer Science Department, University of Copenhagen, August 1989.
- [21] Karoline Malmkjær. *Abstract Interpretation of Partial-Evaluation Algorithms*. PhD thesis, Department of Computing and Information Sciences, Kansas State University, Manhattan, Kansas, March 1993.
- [22] Torben Æ. Mogensen. Glossary for partial evaluation and related topics. *Higher-Order and Symbolic Computation*, 13(4):355–368, 2000.
- [23] Simon Peyton Jones and André Santos. A transformation-based optimiser for Haskell. *Science of Computer Programming*, 32(1-3):3–47, 1998.
- [24] Christian Queinnec and Jean-Marie Geffroy. Partial evaluation applied to pattern matching with intelligent backtrack. In *Proceedings of the Second International Workshop on Static Analysis WSA’92*, volume 81-82 of *Bigre Journal*, pages 109–117, Bordeaux, France, September 1992. IRISA, Rennes, France.
- [25] Donald A. Smith. Partial evaluation of pattern matching in constraint logic programming languages. In Paul Hudak and Neil D. Jones, editors, *Proceedings of the ACM SIGPLAN Symposium on Partial Evaluation and Semantics-Based Program Manipulation*, SIGPLAN Notices, Vol. 26, No 9, pages 62–71, New Haven, Connecticut, June 1991. ACM Press.
- [26] Morten Heine Sørensen, Robert Glück, and Neil D. Jones. A positive supercompiler. *Journal of Functional Programming*, 6(6):811–838, 1996.