

# The string search algorithm by Knuth, Morris and Pratt

Fabian Hellauer and Peter Lammich

March 17, 2025

## Abstract

The Knuth-Morris-Pratt algorithm[1] is often used to show that the problem of finding a string  $s$  in a text  $t$  can be solved deterministically in  $O(|s| + |t|)$  time. We use the Isabelle Refinement Framework[2] to formulate and verify the algorithm. Via refinement, we apply some optimisations and finally use the *Sepref* tool[3] to obtain executable code in *Imperative/HOL*.

## Contents

<b>1 Specification</b>	<b>3</b>
1.1 Sublist-predicate with a position check . . . . .	3
1.1.1 Definition . . . . .	3
1.1.2 Properties . . . . .	3
1.2 Sublist-check algorithms . . . . .	4
<b>2 Naive algorithm</b>	<b>4</b>
2.1 Invariants . . . . .	5
2.2 Algorithm . . . . .	5
2.3 Correctness . . . . .	5
<b>3 Knuth–Morris–Pratt algorithm</b>	<b>6</b>
3.1 Preliminaries: Borders of lists . . . . .	6
3.1.1 Properties . . . . .	6
3.1.2 Examples . . . . .	8
3.2 Main routine . . . . .	8
3.2.1 Invariants . . . . .	8
3.2.2 Algorithm . . . . .	8
3.2.3 Correctness . . . . .	9
3.2.4 Storing the $\mathfrak{f}$ -values . . . . .	10
3.3 Computing $\mathfrak{f}$ . . . . .	11
3.3.1 Invariants . . . . .	11

3.3.2	Algorithm . . . . .	11
3.3.3	Correctness . . . . .	11
3.3.4	Index shift . . . . .	13
3.4	Conflation . . . . .	14
<b>4</b>	<b>Refinement to Imperative/HOL</b>	<b>15</b>
4.1	Overall Correctness Theorem . . . . .	15
<b>5</b>	<b>Tests of Generated ML-Code</b>	<b>16</b>

```

theory KMP
imports Refine-Imperative-HOL.IICF
HOL-Library.Sublist
begin

declare len-greater-imp-nonempty[simp del] min-absorb2[simp]
no-notation Ref.update (<- := -> 62)

```

## 1 Specification

### 1.1 Sublist-predicate with a position check

#### 1.1.1 Definition

One could define

```
definition sublist-at' xs ys i ≡ take (length xs) (drop i ys) = xs
```

However, this doesn't handle out-of-bound indexes uniformly:

```

value[nbe] sublist-at' [] [a] 5
value[nbe] sublist-at' [a] [a] 5
value[nbe] sublist-at' [] [] 5

```

Instead, we use a recursive definition:

```

fun sublist-at :: 'a list ⇒ 'a list ⇒ nat ⇒ bool where
  sublist-at (x#xs) (y#ys) 0 ↔ x=y ∧ sublist-at xs ys 0 |
  sublist-at xs (y#ys) (Suc i) ↔ sublist-at xs ys i |
  sublist-at [] ys 0 ↔ True |
  sublist-at _ [] _ ↔ False

```

In the relevant cases, both definitions agree:

```
lemma i ≤ length ys ⇒ sublist-at xs ys i ↔ sublist-at' xs ys i
  ⟨proof⟩
```

However, the new definition has some reasonable properties:

#### 1.1.2 Properties

```
lemma sublist-lengths: sublist-at xs ys i ⇒ i + length xs ≤ length ys
  ⟨proof⟩
```

```
lemma Nil-is-sublist: sublist-at ([] :: 'x list) ys i ↔ i ≤ length ys
  ⟨proof⟩
```

Furthermore, we need:

```
lemma sublist-step[intro]:
  [| i + length xs < length ys; sublist-at xs ys i; ys!(i + length xs) = x |] ⇒ sublist-at
  (xs@[x]) ys i
```

$\langle proof \rangle$

**lemma** *all-positions-sublist*:

$\llbracket i + \text{length } xs \leq \text{length } ys; \forall jj < \text{length } xs. ys!(i+jj) = xs!jj \rrbracket \implies \text{sublist-at } xs \text{ } ys \text{ } i$   
 $\langle proof \rangle$

**lemma** *sublist-all-positions*:  $\text{sublist-at } xs \text{ } ys \text{ } i \implies \forall jj < \text{length } xs. ys!(i+jj) = xs!jj$   
 $\langle proof \rangle$

It also connects well to theory *HOL-Library.Sublist* (compare *sublist-def*):

**lemma** *sublist-at-altdef*:

$\text{sublist-at } xs \text{ } ys \text{ } i \longleftrightarrow (\exists ps \text{ } ss. ys = ps@xs@ss \wedge i = \text{length } ps)$   
 $\langle proof \rangle$

**corollary** *sublist-iff-sublist-at*:  $\text{Sublist.sublist } xs \text{ } ys \longleftrightarrow (\exists i. \text{sublist-at } xs \text{ } ys \text{ } i)$   
 $\langle proof \rangle$

## 1.2 Sublist-check algorithms

We use the Isabelle Refinement Framework (Theory *Refine-Monadic.Refine-Monadic*) to phrase the specification and the algorithm.

*s* for "searchword" / "searchlist", *t* for "text"

**definition** *kmp-SPEC*  $s \text{ } t = \text{SPEC} (\lambda$

$\text{None} \Rightarrow \nexists i. \text{sublist-at } s \text{ } t \text{ } i \mid$   
 $\text{Some } i \Rightarrow \text{sublist-at } s \text{ } t \text{ } i \wedge (\forall ii < i. \neg \text{sublist-at } s \text{ } t \text{ } ii))$

**lemma** *is-arg-min-id*:  $\text{is-arg-min id } P \text{ } i \longleftrightarrow P \text{ } i \wedge (\forall ii < i. \neg P \text{ } ii)$   
 $\langle proof \rangle$

**lemma** *kmp-result*:  $\text{kmp-SPEC } s \text{ } t =$

$\text{RETURN (if sublist } s \text{ } t \text{ then Some (LEAST } i. \text{sublist-at } s \text{ } t \text{ } i) \text{ else None)}$   
 $\langle proof \rangle$

**corollary** *weak-kmp-SPEC*:  $\text{kmp-SPEC } s \text{ } t \leq \text{SPEC } (\lambda pos. pos \neq \text{None} \longleftrightarrow \text{Sublist.sublist } s \text{ } t)$   
 $\langle proof \rangle$

**lemmas** *kmp-SPEC-altdefs* =

*kmp-SPEC-def*[folded *is-arg-min-id*]  
*kmp-SPEC-def*[folded *sublist-iff-sublist-at*]  
*kmp-result*

## 2 Naive algorithm

Since KMP is a direct advancement of the naive "test-all-starting-positions" approach, we provide it here for comparison:

## 2.1 Invariants

```

definition I-out-na s t  $\equiv \lambda(i,j,pos).$ 
   $(\forall ii < i. \neg \text{sublist-at } s t ii) \wedge$ 
   $(\text{case pos of None} \Rightarrow j = 0)$ 
   $| \text{Some } p \Rightarrow p=i \wedge \text{sublist-at } s t i)$ 
definition I-in-na s t i  $\equiv \lambda(j,pos).$ 
   $\text{case pos of None} \Rightarrow j < \text{length } s \wedge (\forall jj < j. t!(i+jj) = s!(jj))$ 
   $| \text{Some } p \Rightarrow \text{sublist-at } s t i)$ 

```

## 2.2 Algorithm

The following definition is taken from Helmut Seidl's lecture on algorithms and data structures[4] except that we

- output the identified position  $pos$  instead of just  $True$
- use  $pos$  as break-flag to support the abort within the loops
- rewrite  $i \leq \text{length } t - \text{length } s$  in the first while-condition to  $i + \text{length } s \leq \text{length } t$  to avoid having to use  $int$  for list indexes (or the additional precondition  $\text{length } s \leq \text{length } t$ )

```

definition naive-algorithm s t  $\equiv \text{do} \{$ 
  let  $i=0;$ 
  let  $j=0;$ 
  let  $pos=\text{None};$ 
   $(-,pos) \leftarrow \text{WHILEIT } (I\text{-out-na } s t) (\lambda(i,-,pos). i + \text{length } s \leq \text{length } t \wedge$ 
   $pos=\text{None}) (\lambda(i,j,pos). \text{do} \{$ 
     $(-,pos) \leftarrow \text{WHILEIT } (I\text{-in-na } s t i) (\lambda(j,pos). t!(i+j) = s!j \wedge pos=\text{None}) (\lambda(j,-).$ 
    do {
      let  $j=j+1;$ 
      if  $j=\text{length } s$  then RETURN  $(j,\text{Some } i)$  else RETURN  $(j,\text{None})$ 
    })  $(j,pos);$ 
    if  $pos=\text{None}$  then do {
      let  $i = i + 1;$ 
      let  $j = 0;$ 
      RETURN  $(i,j,\text{None})$ 
    } else RETURN  $(i,j,\text{Some } i)$ 
  })  $(i,j,pos);$ 

  RETURN  $pos$ 
}

```

## 2.3 Correctness

The basic lemmas on  $\text{sublist-at}$  from the previous chapter together with *Refine-Monadic*.*Refine-Monadic*'s verification condition generator / solver suffice:

```
lemma  $s \neq [] \implies \text{naive-algorithm } s t \leq \text{kmp-SPEC } s t$ 
     $\langle \text{proof} \rangle$ 
```

Note that the precondition cannot be removed without an extra branch: If  $s = []$ , the inner while-condition accesses out-of-bound memory. This will apply to KMP, too.

### 3 Knuth–Morris–Pratt algorithm

Just like our templates[1][4], we first verify the main routine and discuss the computation of the auxiliary values  $f s$  only in a later section.

#### 3.1 Preliminaries: Borders of lists

```
definition border  $xs\ ys \longleftrightarrow \text{prefix } xs\ ys \wedge \text{suffix } xs\ ys$ 
definition strict-border  $xs\ ys \longleftrightarrow \text{border } xs\ ys \wedge \text{length } xs < \text{length } ys$ 
definition intrinsic-border  $ls \equiv \text{ARG-MAX } b. \text{strict-border } b\ ls$ 
```

##### 3.1.1 Properties

```
interpretation border-order: order border strict-border
     $\langle \text{proof} \rangle$ 
interpretation border-bot: order-bot Nil border strict-border
     $\langle \text{proof} \rangle$ 
```

```
lemma borderE[elim]:
    fixes  $xs\ ys :: \text{'a list}$ 
    assumes border  $xs\ ys$ 
    obtains prefix  $xs\ ys$  and suffix  $xs\ ys$ 
     $\langle \text{proof} \rangle$ 
```

```
lemma strict-borderE[elim]:
    fixes  $xs\ ys :: \text{'a list}$ 
    assumes strict-border  $xs\ ys$ 
    obtains border  $xs\ ys$  and length  $xs < \text{length } ys$ 
     $\langle \text{proof} \rangle$ 
```

```
lemma strict-border-simps[simp]:
    strict-border  $xs\ [] \longleftrightarrow \text{False}$ 
    strict-border  $[]\ (x \# xs) \longleftrightarrow \text{True}$ 
     $\langle \text{proof} \rangle$ 
```

```
lemma strict-border-prefix: strict-border  $xs\ ys \implies \text{strict-prefix } xs\ ys$ 
and strict-border-suffix: strict-border  $xs\ ys \implies \text{strict-suffix } xs\ ys$ 
and strict-border-imp-nonempty: strict-border  $xs\ ys \implies ys \neq []$ 
and strict-border-prefix-suffix: strict-border  $xs\ ys \longleftrightarrow \text{strict-prefix } xs\ ys \wedge \text{strict-suffix }$ 
 $xs\ ys$ 
     $\langle \text{proof} \rangle$ 
```

**lemma** *border-length-le*:  $\text{border } xs \ ys \implies \text{length } xs \leq \text{length } ys$   
 $\langle \text{proof} \rangle$

**lemma** *border-length-r-less* :  $\forall xs. \text{strict-border } xs \ ys \implies \text{length } xs < \text{length } ys$   
 $\langle \text{proof} \rangle$

**lemma** *border-positions*:  $\text{border } xs \ ys \implies \forall i < \text{length } xs. \ ys!i = ys!(\text{length } ys - \text{length } xs + i)$   
 $\langle \text{proof} \rangle$

**lemma** *all-positions-drop-length-take*:  $\llbracket i \leq \text{length } w; i \leq \text{length } x; \forall j < i. x ! j = w ! (\text{length } w + j - i) \rrbracket \implies \text{drop}(\text{length } w - i) \ w = \text{take } i \ x$   
 $\langle \text{proof} \rangle$

**lemma** *all-positions-suffix-take*:  $\llbracket i \leq \text{length } w; i \leq \text{length } x; \forall j < i. x ! j = w ! (\text{length } w + j - i) \rrbracket \implies \text{suffix}(\text{take } i \ x) \ w$   
 $\langle \text{proof} \rangle$

**lemma** *suffix-butlast*:  $\text{suffix } xs \ ys \implies \text{suffix}(\text{butlast } xs) \ (\text{butlast } ys)$   
 $\langle \text{proof} \rangle$

**lemma** *positions-border*:  $\forall j < l. w!j = w!(\text{length } w - l + j) \implies \text{border}(\text{take } l \ w)$   
 $w$   
 $\langle \text{proof} \rangle$

**lemma** *positions-strict-border*:  $l < \text{length } w \implies \forall j < l. w!j = w!(\text{length } w - l + j) \implies \text{strict-border}(\text{take } l \ w) \ w$   
 $\langle \text{proof} \rangle$

**lemmas** *intrinsic-borderI* =  $\text{arg-max-natI}[OF - \text{border-length-r-less}, \text{folded intrinsic-border-def}]$

**lemmas** *intrinsic-borderI'* =  $\text{border-bot.bot.not-eq-extremum}[THEN \ iffD1, THEN \ intrinsic-borderI]$

**lemmas** *intrinsic-border-max* =  $\text{arg-max-nat-le}[OF - \text{border-length-r-less}, \text{folded intrinsic-border-def}]$

**lemma** *nonempty-is-arg-max-ib*:  $ys \neq [] \implies \text{is-arg-max length } (\lambda xs. \text{strict-border } xs \ ys) \ (\text{intrinsic-border } ys)$   
 $\langle \text{proof} \rangle$

**lemma** *intrinsic-border-less*:  $w \neq [] \implies \text{length } (\text{intrinsic-border } w) < \text{length } w$   
 $\langle \text{proof} \rangle$

**lemma** *intrinsic-border-take-less*:  $j > 0 \implies w \neq [] \implies \text{length } (\text{intrinsic-border } w) < \text{length } w$   
 $\langle \text{proof} \rangle$

$(take\ j\ w)) < length\ w$   
 $\langle proof \rangle$

### 3.1.2 Examples

**lemma** *border-example*:  $\{b. border\ b\ "aabaabaa"\} = \{"", "a", "aa", "aabaa", "aabaabaa"\}$   
 $\langle proof \rangle$   
 $(is\ \{b. border\ b\ ?l\} = \{\text{?take}0, \text{?take}1, \text{?take}2, \text{?take}5, \text{?l}\})$

**corollary** *strict-border-example*:  $\{b. strict-border\ b\ "aabaabaa"\} = \{"", "a", "aa", "aabaa"\}$   
 $\langle proof \rangle$   
 $(is\ ?l = ?r)$

**corollary** *intrinsic-border*  $"aabaabaa" = "aabaa"$   
 $\langle proof \rangle$

## 3.2 Main routine

The following is Seidl's "border"-table[4] (values shifted by 1 so we don't need *int*), or equivalently, "f" from Knuth's, Morris' and Pratt's paper[1] (with indexes starting at 0).

```
fun f :: 'a list ⇒ nat ⇒ nat where
  f s 0 = 0 — This increments the compare position while  $j = 0$  |
  f s j = length (intrinsic-border (take j s)) + 1
```

Note that we use their "next" only implicitly.

### 3.2.1 Invariants

**definition** *I-outer*  $s\ t \equiv \lambda(i,j,pos).$   
 $(\forall ii < i. \neg sublist-at\ s\ t\ ii) \wedge$   
 $(\text{case pos of None} \Rightarrow (\forall jj < j. t!(i+jj) = s!(jj)) \wedge j < length\ s$   
 $| \text{Some } p \Rightarrow p=i \wedge \text{sublist-at } s\ t\ i)$

For the inner loop, we can reuse *I-in-na*.

### 3.2.2 Algorithm

First, we use the non-evaluable function *f* directly:

```
definition kmp s t ≡ do {
  ASSERT (s ≠ []);
  let i=0;
  let j=0;
  let pos=None;
  (-,-,pos) ← WHILEIT (I-outer s t) (λ(i,j,pos). i + length s ≤ length t ∧ pos=None)
  (λ(i,j,pos). do {
```

```

ASSERT (i + length s ≤ length t);
(j,pos) ← WHILEIT (I-in-na s t i) (λ(j,pos). t!(i+j) = s!j ∧ pos=None)
(λ(j,pos). do {
    let j=j+1;
    if j=length s then RETURN (j,Some i) else RETURN (j,None)
}) (j,pos);
if pos=None then do {
    ASSERT (j < length s);
    let i = i + (j - ⌈ s j + 1);
    let j = max 0 (⌈ s j - 1); — max not necessary
    RETURN (i,j,None)
} else RETURN (i,j,Some i)
}) (i,j,pos);

RETURN pos
}

```

### 3.2.3 Correctness

**lemma**  $\mathfrak{f}\text{-eq-0-iff-}j\text{-eq-0}$  [simp]:  $\mathfrak{f} s j = 0 \longleftrightarrow j = 0$   
 $\langle proof \rangle$

**lemma**  $j\text{-le-}\mathfrak{f}\text{-le}$ :  $j \leq \text{length } s \implies \mathfrak{f} s j \leq j$   
 $\langle proof \rangle$

**lemma**  $j\text{-le-}\mathfrak{f}\text{-le}'$ :  $0 < j \implies j \leq \text{length } s \implies \mathfrak{f} s j - 1 < j$   
 $\langle proof \rangle$

**lemma**  $\mathfrak{f}\text{-le}$ :  $s \neq [] \implies \mathfrak{f} s j - 1 < \text{length } s$   
 $\langle proof \rangle$

**lemma** *reuse-matches*:  
**assumes**  $j\text{-le}$ :  $j \leq \text{length } s$   
**and** *old-matches*:  $\forall jj < j. t ! (i + jj) = s ! jj$   
**shows**  $\forall jj < \mathfrak{f} s j - 1. t ! (i + (j - \mathfrak{f} s j + 1) + jj) = s ! jj$   
 $(\text{is } \forall jj < ?j'. t ! (?i' + jj) = s ! jj)$   
 $\langle proof \rangle$

**theorem** *shift-safe*:  
**assumes**  
 $\forall ii < i. \neg \text{sublist-at } s t ii$   
 $t!(i+j) \neq s!j$  **and**  
[simp]:  $j < \text{length } s$  **and**  
*matches*:  $\forall jj < j. t!(i+jj) = s!jj$   
**defines**  
*assignment*:  $i' \equiv i + (j - \mathfrak{f} s j + 1)$   
**shows**

$\forall ii < i'. \neg \text{sublist-at } s t ii$   
 $\langle \text{proof} \rangle$

**lemma** *kmp-correct*:  $s \neq []$   
 $\implies \text{kmp } s t \leq \text{kmp-SPEC } s t$   
 $\langle \text{proof} \rangle$

### 3.2.4 Storing the $f$ -values

We refine the algorithm to compute the  $f$ -values only once at the start:

**definition** *compute-fs-SPEC* :: 'a list  $\Rightarrow$  nat list nres **where**  
 $\text{compute-fs-SPEC } s \equiv \text{SPEC } (\lambda fs. \text{length } fs = \text{length } s + 1 \wedge (\forall j \leq \text{length } s. fs!j = f s j))$

**definition** *kmp1*  $s t \equiv \text{do } \{$   
 $\quad \text{ASSERT } (s \neq []);$   
 $\quad \text{let } i=0;$   
 $\quad \text{let } j=0;$   
 $\quad \text{let } pos=\text{None};$   
 $\quad fs \leftarrow \text{compute-fs-SPEC } (\text{butlast } s); \text{ — At the last char, we abort instead.}$   
 $\quad (-, -, pos) \leftarrow \text{WHILEIT } (I\text{-outer } s t) (\lambda(i, j, pos). i + \text{length } s \leq \text{length } t \wedge pos=\text{None})$   
 $\quad (\lambda(i, j, pos). \text{do } \{$   
 $\quad \quad \text{ASSERT } (i + \text{length } s \leq \text{length } t);$   
 $\quad \quad (j, pos) \leftarrow \text{WHILEIT } (I\text{-in-na } s t i) (\lambda(j, pos). t!(i+j) = s!j \wedge pos=\text{None})$   
 $\quad \quad (\lambda(j, pos). \text{do } \{$   
 $\quad \quad \text{let } j=j+1;$   
 $\quad \quad \text{if } j=\text{length } s \text{ then RETURN } (j, \text{Some } i) \text{ else RETURN } (j, \text{None})$   
 $\quad \quad \}) (j, pos);$   
 $\quad \quad \text{if } pos=\text{None} \text{ then do } \{$   
 $\quad \quad \quad \text{ASSERT } (j < \text{length } fs);$   
 $\quad \quad \quad \text{let } i = i + (j - fs!j + 1);$   
 $\quad \quad \quad \text{let } j = \max 0 (fs!j - 1); \text{ — max not necessary}$   
 $\quad \quad \quad \text{RETURN } (i, j, \text{None})$   
 $\quad \quad \} \text{ else RETURN } (i, j, \text{Some } i)$   
 $\quad \}) (i, j, pos);$   
 $\quad \text{RETURN } pos$   
 $\}$

**lemma** *f-butlast[simp]*:  $j < \text{length } s \implies f(\text{butlast } s) j = f s j$   
 $\langle \text{proof} \rangle$

**lemma** *kmp1-refine*:  $\text{kmp1 } s t \leq \text{kmp } s t$   
 $\langle \text{proof} \rangle$

Next, an algorithm that satisfies *compute-fs-SPEC*:

### 3.3 Computing $\mathfrak{f}$

#### 3.3.1 Invariants

```

definition I-out-cb s ≡ λ(fs,i,j).
  length s + 1 = length fs ∧
  (forall jj < j. fs!jj = f s jj) ∧
  fs!(j-1) = i ∧
  0 < j
definition I-in-cb s j ≡ λi.
  if j=1 then i=0 — first iteration
  else
    strict-border (take (i-1) s) (take (j-1) s) ∧
    fs j ≤ i + 1

```

#### 3.3.2 Algorithm

Again, we follow Seidl[4], p.582. Apart from the +1-shift, we make another modification: Instead of directly setting  $\mathfrak{f}s ! 1$ , we let the first loop-iteration (if there is one) do that for us. This allows us to remove the precondition  $s \neq []$ , as the index bounds are respected even in that corner case.

```

definition compute-fs :: 'a list ⇒ nat list nres where
  compute-fs s = do {
    let fs=replicate (length s + 1) 0; — only the first 0 is needed
    let i=0;
    let j=1;
    (fs, -, -) ← WHILEIT (I-out-cb s) (λ(fs,-,j). j < length fs) (λ(fs,i,j). do {
      i ← WHILEIT (I-in-cb s j) (λi. i>0 ∧ s!(i-1) ≠ s!(j-1)) (λi. do {
        ASSERT (i-1 < length fs);
        let i=fs!(i-1);
        RETURN i
      }) i;
      let i=i+1;
      ASSERT (j < length fs);
      let fs=fs[j:=i];
      let j=j+1;
      RETURN (fs,i,j)
    }) (fs,i,j);
    RETURN fs
  }

```

#### 3.3.3 Correctness

```

lemma take-length-ib[simp]:
  assumes 0 < j j ≤ length s
  shows take (length (intrinsic-border (take j s))) s = intrinsic-border (take j s)
  ⟨proof⟩

```

**lemma** *ib-singleton*[simp]: *intrinsic-border* [z] = []  
*(proof)*

**lemma** *border-butlast*: *border* xs ys  $\implies$  *border* (*butlast* xs) (*butlast* ys)  
*(proof)*

**corollary** *strict-border-butlast*: xs  $\neq$  []  $\implies$  *strict-border* xs ys  $\implies$  *strict-border* (*butlast* xs) (*butlast* ys)  
*(proof)*

**lemma** *border-take-lengths*: i  $\leq$  *length* s  $\implies$  *border* (*take* i s) (*take* j s)  $\implies$  i  $\leq$  j  
*(proof)*

**lemma** *border-step*: *border* xs ys  $\longleftrightarrow$  *border* (xs@[ys!length xs]) (ys@[ys!length xs])  
*(proof)*

**corollary** *strict-border-step*: *strict-border* xs ys  $\longleftrightarrow$  *strict-border* (xs@[ys!length xs]) (ys@[ys!length xs])  
*(proof)*

**lemma** *ib-butlast*: *length* w  $\geq$  2  $\implies$  *length* (*intrinsic-border* w)  $\leq$  *length* (*intrinsic-border* (*butlast* w)) + 1  
*(proof)*

**corollary**  $\mathfrak{f}$ -*Suc*: *Suc* i  $\leq$  *length* w  $\implies$   $\mathfrak{f}$  w (*Suc* i)  $\leq$   $\mathfrak{f}$  w i + 1  
*(proof)*

**lemma**  $\mathfrak{f}$ -*step-bound*:  
**assumes** j  $\leq$  *length* w  
**shows**  $\mathfrak{f}$  w j  $\leq$   $\mathfrak{f}$  w (j-1) + 1  
*(proof)*

**lemma** *border-take-f*: *border* (*take* ( $\mathfrak{f}$  s i - 1) s) (*take* i s)  
*(proof)*

**corollary**  $\mathfrak{f}$ -*strict-borderI*: y =  $\mathfrak{f}$  s (i-1)  $\implies$  *strict-border* (*take* (i-1) s) (*take* (j-1) s)  $\implies$  *strict-border* (*take* (y-1) s) (*take* (j-1) s)  
*(proof)*

**corollary** *strict-border-take-f*: 0 < i  $\implies$  i  $\leq$  *length* s  $\implies$  *strict-border* (*take* ( $\mathfrak{f}$  s i - 1) s) (*take* i s)  
*(proof)*

**lemma**  $\mathfrak{f}$ -*is-max*: j  $\leq$  *length* s  $\implies$  *strict-border* b (*take* j s)  $\implies$   $\mathfrak{f}$  s j  $\geq$  *length* b + 1  
*(proof)*

**theorem** *skipping-ok*:  
**assumes** j-bounds[simp]: 1 < j j  $\leq$  *length* s

```

and mismatch:  $s!(i-1) \neq s!(j-1)$ 
and greater-checked:  $\mathfrak{f} s j \leq i + 1$ 
and strict-border (take  $(i-1)$   $s$ ) (take  $(j-1)$   $s$ )
shows  $\mathfrak{f} s j \leq \mathfrak{f} s (i-1) + 1$ 
⟨proof⟩

```

```

lemma extend-border:
assumes  $j \leq \text{length } s$ 
assumes  $s!(i-1) = s!(j-1)$ 
assumes strict-border (take  $(i-1)$   $s$ ) (take  $(j-1)$   $s$ )
assumes  $\mathfrak{f} s j \leq i + 1$ 
shows  $\mathfrak{f} s j = i + 1$ 
⟨proof⟩

```

```

lemma compute- $\mathfrak{f}s$ -correct: compute- $\mathfrak{f}s$   $s \leq \text{compute-}\mathfrak{f}s\text{-SPEC } s$ 
⟨proof⟩

```

### 3.3.4 Index shift

To avoid inefficiencies, we refine *compute- $\mathfrak{f}s$*  to take  $s$  instead of *butlast*  $s$  (it still only uses *butlast*  $s$ ).

```

definition compute-butlast- $\mathfrak{f}s$  :: 'a list  $\Rightarrow$  nat list nres where
  compute-butlast- $\mathfrak{f}s$   $s = \text{do } \{$ 
    let  $\mathfrak{f}s = \text{replicate} (\text{length } s) 0;$ 
    let  $i = 0;$ 
    let  $j = 1;$ 
     $(\mathfrak{f}s, -, -) \leftarrow \text{WHILEIT } (I\text{-out-cb } (\text{butlast } s)) (\lambda(b, i, j). j < \text{length } b) (\lambda(\mathfrak{f}s, i, j). \text{do } \{$ 
      ASSERT ( $j < \text{length } \mathfrak{f}s$ );
       $i \leftarrow \text{WHILEIT } (I\text{-in-cb } (\text{butlast } s) j) (\lambda i. i > 0 \wedge s!(i-1) \neq s!(j-1)) (\lambda i. \text{do } \{$ 
        ASSERT ( $i-1 < \text{length } \mathfrak{f}s$ );
        let  $i = \mathfrak{f}s!(i-1);$ 
        RETURN  $i$ 
       $\}) i;$ 
      let  $i = i + 1;$ 
      ASSERT ( $j < \text{length } \mathfrak{f}s$ );
      let  $\mathfrak{f}s = \mathfrak{f}s[j := i];$ 
      let  $j = j + 1;$ 
      RETURN  $(\mathfrak{f}s, i, j)$ 
     $\}) (\mathfrak{f}s, i, j);$ 
  RETURN  $\mathfrak{f}s$ 
}

```

```

lemma compute- $\mathfrak{f}s$ -inner-bounds:
assumes  $I\text{-out-cb } s (\mathfrak{f}s, ix, j)$ 
assumes  $j < \text{length } \mathfrak{f}s$ 
assumes  $I\text{-in-cb } s j i$ 
shows  $i-1 < \text{length } s j-1 < \text{length } s$ 
⟨proof⟩

```

```

lemma compute-butlast-fs-refine[refine]:
  assumes (s,s') ∈ br butlast ((≠) [])
  shows compute-butlast-fs s ≤ ⇲ Id (compute-fs-SPEC s')
⟨proof⟩

```

### 3.4 Conflation

We replace *compute-fs-SPEC* with *compute-butlast-fs*

```

definition kmp2 s t ≡ do {
  ASSERT (s ≠ []);
  let i=0;
  let j=0;
  let pos=None;
  fs ← compute-butlast-fs s;
  (-,-,pos) ← WHILEIT (I-outer s t) (λ(i,j,pos). i + length s ≤ length t ∧ pos=None)
  (λ(i,j,pos). do {
    ASSERT (i + length s ≤ length t ∧ pos=None);
    (j,pos) ← WHILEIT (I-in-na s t i) (λ(j,pos). t!(i+j) = s!j ∧ pos=None)
    (λ(j,pos). do {
      let j=j+1;
      if j=length s then RETURN (j,Some i) else RETURN (j,None)
    }) (j,pos);
    if pos=None then do {
      ASSERT (j < length fs);
      let i = i + (j - fs!j + 1);
      let j = max 0 (fs!j - 1); — max not necessary
      RETURN (i,j,None)
    } else RETURN (i,j,Some i)
  }) (i,j,pos);
  RETURN pos
}

```

Using *compute-butlast-fs-refine* (it has attribute *refine*), the proof is trivial:

```

lemma kmp2-refine: kmp2 s t ≤ kmp1 s t
⟨proof⟩

```

```

lemma kmp2-correct: s ≠ []
  ⇒ kmp2 s t ≤ kmp-SPEC s t
⟨proof⟩

```

For convenience, we also remove the precondition:

```

definition kmp3 s t ≡ do {
  if s=[] then RETURN (Some 0) else kmp2 s t
}

```

```

lemma kmp3-correct: kmp3 s t ≤ kmp-SPEC s t
⟨proof⟩

```

## 4 Refinement to Imperative/HOL

```

lemma eq-id-param:  $((=), (=)) \in Id \rightarrow Id \rightarrow Id$   $\langle proof \rangle$ 

lemmas in-bounds-aux = compute-fs-inner-bounds[of butlast s for s, simplified]

sepref-definition compute-butlast-fs-impl is compute-butlast-fs ::  $(arl-assn id-assn)^k$   

 $\rightarrow_a array-assn nat-assn$   

 $\langle proof \rangle$ 

declare compute-butlast-fs-impl.refine[sepref-fr-rules]

sepref-register compute-fs

lemma kmp-inner-in-bound:  

assumes  $i + length s \leq length t$   

assumes I-in-na s t i (j, None)  

shows  $i + j < length t$   $j < length s$   

 $\langle proof \rangle$ 

sepref-definition kmp-impl is uncurry kmp3 ::  $(arl-assn id-assn)^k *_a (arl-assn id-assn)^k \rightarrow_a option-assn nat-assn$   

 $\langle proof \rangle$ 

export-code kmp-impl in SML-imp module-name KMP

lemma kmp3-correct':  

 $(uncurry kmp3, uncurry kmp-SPEC) \in Id \times_r Id \rightarrow_f \langle Id \rangle nres-rel$   

 $\langle proof \rangle$ 

lemmas kmp-impl-correct' = kmp-impl.refine[FCOMP kmp3-correct']

```

### 4.1 Overall Correctness Theorem

The following theorem relates the final Imperative HOL algorithm to its specification, using, beyond basic HOL concepts

- Hoare triples for Imperative/HOL, provided by the Separation Logic Framework for Imperative/HOL (Theory *Separation-Logic-Imperative-HOL.Sep-Main*);
- The assertion *arl-assn* to specify array-lists, which we use to represent the input strings of the algorithm;
- The *sublist-at* function that we defined in section 1.

```

theorem kmp-impl-correct:  

 $< arl-assn id-assn s si * arl-assn id-assn t ti >$   

 $kmp-impl si ti$ 

```

```

<λr. arl-assn id-assn s si * arl-assn id-assn t ti * ↑(
  case r of None ⇒ ∉ i. sublist-at s t i
  | Some i ⇒ sublist-at s t i ∧ ( ∀ ii < i. ¬ sublist-at s t ii)
)>t
⟨proof⟩

```

```
definition kmp-string-impl ≡ kmp-impl :: (char array × nat) ⇒ -
```

## 5 Tests of Generated ML-Code

```
⟨ML⟩
```

```
end
```

## References

- [1] D. E. Knuth, J. James H. Morris, and V. R. Pratt. Fast pattern matching in strings. *SIAM Journal on Computing*, 6(2):323–350, 1977.
- [2] P. Lammich. Refinement for monadic programs. *Archive of Formal Proofs*, Jan. 2012. [http://isa-afp.org/entries/Refine\\_Monadic.html](http://isa-afp.org/entries/Refine_Monadic.html), Formal proof development.
- [3] P. Lammich. The imperative refinement framework. *Archive of Formal Proofs*, Aug. 2016. [http://isa-afp.org/entries/Refine\\_Imperative\\_HOL.html](http://isa-afp.org/entries/Refine_Imperative_HOL.html), Formal proof development.
- [4] H. Seidl. Grundlagen: Algorithmen und Datenstrukturen. <http://www2.in.tum.de/hp/file?fid=1347>, German lecture notes, 2016.