

# POSIX Lexing with Derivatives of Regular Expressions

Fahad Ausaf<sup>1</sup>, Roy Dyckhoff<sup>2</sup>, and Christian Urban<sup>3</sup>

<sup>1</sup> King's College London

fahad.ausaf@icloud.com

<sup>2</sup> University of St Andrews

roy.dyckhoff@st-andrews.ac.uk

<sup>3</sup> King's College London

christian.urban@kcl.ac.uk

**Abstract.** Brzozowski introduced the notion of derivatives for regular expressions. They can be used for a very simple regular expression matching algorithm. Sulzmann and Lu cleverly extended this algorithm in order to deal with POSIX matching, which is the underlying disambiguation strategy for regular expressions needed in lexers. Their algorithm generates POSIX values which encode the information of *how* a regular expression matches a string—that is, which part of the string is matched by which part of the regular expression. In this paper we give our inductive definition of what a POSIX value is and show (i) that such a value is unique (for given regular expression and string being matched) and (ii) that Sulzmann and Lu's algorithm always generates such a value (provided that the regular expression matches the string). We show that (iii) our inductive definition of a POSIX value is equivalent to an alternative definition by Okui and Suzuki which identifies POSIX values as least elements according to an ordering of values. We also prove the correctness of Sulzmann's bitcoded version of the POSIX matching algorithm and extend the results to additional constructors for regular expressions.

**Keywords:** POSIX matching, Derivatives of Regular Expressions, Isabelle/HOL

## 1 List prefixes, suffixes, and homeomorphic embedding

**theory** *Sublist*

**imports** *Main*

**begin**

---

\* This paper is a revised and expanded version of [3]. Compared with that paper we give a second definition for POSIX values introduced by Okui Suzuki [12,13] and prove that it is equivalent to our original one. This second definition is based on an ordering of values and very similar to, but not equivalent with, the definition given by Sulzmann and Lu [16]. The advantage of the definition based on the ordering is that it implements more directly the informal rules from the POSIX standard. We also prove Sulzmann & Lu's conjecture that their bitcoded version of the POSIX algorithm is correct. Furthermore we extend our results to additional constructors of regular expressions.

### 1.1 Prefix order on lists

**definition** *prefix* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool  
**where** *prefix* *xs* *ys*  $\longleftrightarrow (\exists z s. ys = xs @ z s)$

**definition** *strict\_prefix* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool  
**where** *strict\_prefix* *xs* *ys*  $\longleftrightarrow prefix\ xs\ ys \wedge xs \neq ys$

**interpretation** *prefix\_order*: *order prefix strict\_prefix*  
**by** *standard* (*auto simp: prefix\_def strict\_prefix\_def*)

**interpretation** *prefix\_bot*: *order\_bot Nil prefix strict\_prefix*  
**by** *standard* (*simp add: prefix\_def*)

**lemma** *prefixI* [*intro?*]: *ys = xs @ z s  $\implies prefix\ xs\ ys$*   
**unfolding** *prefix\_def* **by** *blast*

**lemma** *prefixE* [*elim?*]:  
**assumes** *prefix xs ys*  
**obtains** *z s* **where** *ys = xs @ z s*  
**using** *assms* **unfolding** *prefix\_def* **by** *blast*

**lemma** *strict\_prefixI'* [*intro?*]: *ys = xs @ z # z s  $\implies strict\_prefix\ xs\ ys$*   
**unfolding** *strict\_prefix\_def prefix\_def* **by** *blast*

**lemma** *strict\_prefixE'* [*elim?*]:  
**assumes** *strict\_prefix xs ys*  
**obtains** *z z s* **where** *ys = xs @ z # z s*

**proof** –

**from** (*strict\_prefix xs ys*) **obtain** *us* **where** *ys = xs @ us* **and** *xs  $\neq$  ys*  
**unfolding** *strict\_prefix\_def prefix\_def* **by** *blast*  
**with** *that* **show** *?thesis* **by** (*auto simp add: neq\_Nil\_conv*)

**qed**

**lemma** *strict\_prefixI* [*intro?*]: *prefix xs ys  $\implies xs \neq ys \implies strict\_prefix\ xs\ ys$*   
**by** (*fact prefix\_order.le\_neq\_trans*)

**lemma** *strict\_prefixE* [*elim?*]:  
**fixes** *xs ys* :: 'a list  
**assumes** *strict\_prefix xs ys*  
**obtains** *prefix xs ys* **and** *xs  $\neq$  ys*  
**using** *assms* **unfolding** *strict\_prefix\_def* **by** *blast*

### 1.2 Basic properties of prefixes

**theorem** *Nil\_prefix* [*simp*]: *prefix [] xs*

**by** (*fact prefix\_bot.bot\_least*)

**theorem** *prefix\_Nil* [*simp*]: (*prefix xs []*) = (*xs = []*)  
**by** (*fact prefix\_bot.bot\_unique*)

**lemma** *prefix\_snoc* [*simp*]: *prefix xs (ys @ [y])*  $\longleftrightarrow$  *xs = ys @ [y]*  $\vee$  *prefix xs ys*

**proof**

**assume** *prefix xs (ys @ [y])*

**then obtain** *zs* **where** *zs: ys @ [y] = xs @ zs ..*

**show** *xs = ys @ [y]*  $\vee$  *prefix xs ys*

**by** (*metis append\_Nil2 butlast\_append butlast\_snoc prefixI zs*)

**next**

**assume** *xs = ys @ [y]*  $\vee$  *prefix xs ys*

**then show** *prefix xs (ys @ [y])*

**by** (*metis prefix\_order.eq\_iff prefix\_order.order\_trans prefixI*)

**qed**

**lemma** *Cons\_prefix\_Cons* [*simp*]: *prefix (x # xs) (y # ys)* = (*x = y*  $\wedge$  *prefix xs ys*)  
**by** (*auto simp add: prefix\_def*)

**lemma** *prefix\_code* [*code*]:

*prefix [] xs*  $\longleftrightarrow$  *True*

*prefix (x # xs) []*  $\longleftrightarrow$  *False*

*prefix (x # xs) (y # ys)*  $\longleftrightarrow$  *x = y*  $\wedge$  *prefix xs ys*

**by** *simp\_all*

**lemma** *same\_prefix\_prefix* [*simp*]: *prefix (xs @ ys) (xs @ zs)* = *prefix ys zs*  
**by** (*induct xs simp\_all*)

**lemma** *same\_prefix\_nil* [*simp*]: *prefix (xs @ ys) xs* = (*ys = []*)  
**by** (*metis append\_Nil2 append\_self\_conv prefix\_order.eq\_iff prefixI*)

**lemma** *prefix\_prefix* [*simp*]: *prefix xs ys*  $\implies$  *prefix xs (ys @ zs)*  
**unfolding** *prefix\_def* **by** *fastforce*

**lemma** *append\_prefixD*: *prefix (xs @ ys) zs*  $\implies$  *prefix xs zs*  
**by** (*auto simp add: prefix\_def*)

**theorem** *prefix\_Cons*: *prefix xs (y # ys)* = (*xs = []*  $\vee$  ( $\exists zs. xs = y \# zs \wedge prefix zs ys$ ))  
**by** (*cases xs*) (*auto simp add: prefix\_def*)

**theorem** *prefix\_append*:

*prefix xs (ys @ zs)* = (*prefix xs ys*  $\vee$  ( $\exists us. xs = ys @ us \wedge prefix us zs$ ))

**apply** (*induct zs rule: rev\_induct*)

**apply** *force*  
**apply** (*simp flip: append\_assoc*)  
**apply** (*metis append\_eq\_append1*)  
**done**

**lemma** *append\_one\_prefix*:  
 $prefix\ xs\ ys \implies length\ xs < length\ ys \implies prefix\ (xs\ @\ [ys\ !\ length\ xs])\ ys$   
**proof** (*unfold prefix\_def*)  
**assume** *a1*:  $\exists\ zs.\ ys = xs\ @\ zs$   
**then obtain** *sk* :: 'a list **where** *sk*:  $ys = xs\ @\ sk$  **by** *fastforce*  
**assume** *a2*:  $length\ xs < length\ ys$   
**have** *f1*:  $\bigwedge v.\ ([::'a\ list]\ @\ v = v$  **using** *append\_Nil2* **by** *simp*  
**have**  $[] \neq sk$  **using** *a1 a2 sk less\_not\_refl* **by** *force*  
**hence**  $\exists v.\ xs\ @\ hd\ sk \# v = ys$  **using** *sk* **by** (*metis hd\_Cons\_tl*)  
**thus**  $\exists zs.\ ys = (xs\ @\ [ys\ !\ length\ xs])\ @\ zs$  **using** *f1* **by** *fastforce*  
**qed**

**theorem** *prefix\_length\_le*:  $prefix\ xs\ ys \implies length\ xs \leq length\ ys$   
**by** (*auto simp add: prefix\_def*)

**lemma** *prefix\_same\_cases*:  
 $prefix\ (xs_1::'a\ list)\ ys \implies prefix\ xs_2\ ys \implies prefix\ xs_1\ xs_2 \vee prefix\ xs_2\ xs_1$   
**unfolding** *prefix\_def* **by** (*force simp: append\_eq\_append\_conv2*)

**lemma** *prefix\_length\_prefix*:  
 $prefix\ ps\ xs \implies prefix\ qs\ xs \implies length\ ps \leq length\ qs \implies prefix\ ps\ qs$   
**by** (*auto simp: prefix\_def*) (*metis append\_Nil2 append\_eq\_append\_conv\_if*)

**lemma** *set\_mono\_prefix*:  $prefix\ xs\ ys \implies set\ xs \subseteq set\ ys$   
**by** (*auto simp add: prefix\_def*)

**lemma** *take\_is\_prefix*:  $prefix\ (take\ n\ xs)\ xs$   
**unfolding** *prefix\_def* **by** (*metis append\_take\_drop\_id*)

**lemma** *prefixeq\_butlast*:  $prefix\ (butlast\ xs)\ xs$   
**by** (*simp add: butlast\_conv\_take take\_is\_prefix*)

**lemma** *map\_mono\_prefix*:  $prefix\ xs\ ys \implies prefix\ (map\ f\ xs)\ (map\ f\ ys)$   
**by** (*auto simp: prefix\_def*)

**lemma** *filter\_mono\_prefix*:  $prefix\ xs\ ys \implies prefix\ (filter\ P\ xs)\ (filter\ P\ ys)$   
**by** (*auto simp: prefix\_def*)

**lemma** *sorted\_antimono\_prefix*:  $prefix\ xs\ ys \implies sorted\ ys \implies sorted\ xs$   
**by** (*metis sorted\_append prefix\_def*)

**lemma** *prefix\_length\_less*:  $strict\_prefix\ xs\ ys \implies length\ xs < length\ ys$   
**by** (*auto simp: strict\_prefix\_def prefix\_def*)

**lemma** *prefix\_snocD*:  $prefix\ (xs@[x])\ ys \implies strict\_prefix\ xs\ ys$   
**by** (*simp add: strict\_prefixI' prefix\_order.dual\_order.strict\_trans1*)

**lemma** *strict\_prefix\_simps* [*simp, code*]:  
 $strict\_prefix\ xs\ [] \longleftrightarrow False$   
 $strict\_prefix\ []\ (x\ \# \ xs) \longleftrightarrow True$   
 $strict\_prefix\ (x\ \# \ xs)\ (y\ \# \ ys) \longleftrightarrow x = y \wedge strict\_prefix\ xs\ ys$   
**by** (*simp\_all add: strict\_prefix\_def cong: conj\_cong*)

**lemma** *take\_strict\_prefix*:  $strict\_prefix\ xs\ ys \implies strict\_prefix\ (take\ n\ xs)\ ys$   
**proof** (*induct n arbitrary: xs ys*)  
**case** 0  
**then show** ?*case* **by** (*cases ys*) *simp\_all*  
**next**  
**case** (*Suc n*)  
**then show** ?*case* **by** (*metis prefix\_order.less\_trans strict\_prefixI take\_is\_prefix*)  
**qed**

**lemma** *not\_prefix\_cases*:  
**assumes** *pfx*:  $\neg prefix\ ps\ ls$   
**obtains**  
 (*c1*)  $ps \neq []$  **and**  $ls = []$   
 | (*c2*)  $a\ as\ x\ xs$  **where**  $ps = a\ \# \ as$  **and**  $ls = x\ \# \ xs$  **and**  $x = a$  **and**  $\neg prefix\ as\ xs$   
 | (*c3*)  $a\ as\ x\ xs$  **where**  $ps = a\ \# \ as$  **and**  $ls = x\ \# \ xs$  **and**  $x \neq a$   
**proof** (*cases ps*)  
**case** *Nil*  
**then show** ?*thesis* **using** *pfx* **by** *simp*  
**next**  
**case** (*Cons a as*)  
**note**  $c = \langle ps = a\ \# \ as \rangle$   
**show** ?*thesis*  
**proof** (*cases ls*)  
**case** *Nil* **then show** ?*thesis* **by** (*metis append\_Nil2 pfx c1 same\_prefix\_nil*)  
**next**  
**case** (*Cons x xs*)  
**show** ?*thesis*  
**proof** (*cases x = a*)  
**case** *True*  
**have**  $\neg prefix\ as\ xs$  **using** *pfx c Cons True* **by** *simp*  
**with**  $c\ Cons\ True$  **show** ?*thesis* **by** (*rule c2*)  
**next**

```

case False
with c Cons show ?thesis by (rule c3)
qed
qed
qed

```

```

lemma not_prefix_induct [consumes 1, case_names Nil Neq Eq]:
assumes np: ¬ prefix ps ls
and base: ∧x xs. P (x#xs) []
and r1: ∧x xs y ys. x ≠ y ⇒ P (x#xs) (y#ys)
and r2: ∧x xs y ys. [x = y; ¬ prefix xs ys; P xs ys] ⇒ P (x#xs) (y#ys)
shows P ps ls using np
proof (induct ls arbitrary: ps)
case Nil
then show ?case
by (auto simp: neq_Nil_conv elim!: not_prefix_cases intro!: base)
next
case (Cons y ys)
then have npfx: ¬ prefix ps (y # ys) by simp
then obtain x xs where pv: ps = x # xs
by (rule not_prefix_cases) auto
show ?case by (metis Cons.hyps Cons_prefix_Cons npfx pv r1 r2)
qed

```

### 1.3 Prefixes

```

primrec prefixes where
prefixes [] = [[]] |
prefixes (x#xs) = [] # map ((#) x) (prefixes xs)

```

```

lemma in_set_prefixes[simp]: xs ∈ set (prefixes ys) ⟷ prefix xs ys
proof (induct xs arbitrary: ys)
case Nil
then show ?case by (cases ys) auto
next
case (Cons a xs)
then show ?case by (cases ys) auto
qed

```

```

lemma length_prefixes[simp]: length (prefixes xs) = length xs + 1
by (induction xs) auto

```

```

lemma distinct_prefixes [intro]: distinct (prefixes xs)
by (induction xs) (auto simp: distinct_map)

```

```

lemma prefixes_snoc [simp]: prefixes (xs@[x]) = prefixes xs @ [xs@[x]]

```

**by** (*induction xs*) *auto*

**lemma** *prefixes\_not\_Nil* [*simp*]: *prefixes xs*  $\neq$  []  
**by** (*cases xs*) *auto*

**lemma** *hd\_prefixes* [*simp*]: *hd (prefixes xs)* = []  
**by** (*cases xs*) *simp\_all*

**lemma** *last\_prefixes* [*simp*]: *last (prefixes xs)* = *xs*  
**by** (*induction xs*) (*simp\_all add: last\_map*)

**lemma** *prefixes\_append*:  
*prefixes (xs @ ys)* = *prefixes xs @ map* ( $\lambda$ *ys'. xs @ ys'*) (*tl (prefixes ys)*)  
**proof** (*induction xs*)  
**case Nil**  
**thus** ?*case by* (*cases ys*) *auto*  
**qed** *simp\_all*

**lemma** *prefixes\_eq\_snoc*:  
*prefixes ys = xs @ [x]*  $\longleftrightarrow$   
(*ys = []*  $\wedge$  *xs = []*  $\vee$  ( $\exists$  *z zs. ys = zs@[z]*  $\wedge$  *xs = prefixes zs*))  $\wedge$  *x = ys*  
**by** (*cases ys rule: rev\_cases*) *auto*

**lemma** *prefixes\_tailrec* [*code*]:  
*prefixes xs = rev (snd (foldl* ( $\lambda$ (*acc1, acc2*) *x. (x#acc1, rev (x#acc1)#acc2*)) ([], [[]]  
*xs*))

**proof** –  
**have** *foldl* ( $\lambda$ (*acc1, acc2*) *x. (x#acc1, rev (x#acc1)#acc2*)) (*ys, rev ys # zs*) *xs =*  
(*rev xs @ ys, rev (map* ( $\lambda$ *as. rev ys @ as*) (*prefixes xs*)) @ *zs*) **for** *ys zs*  
**proof** (*induction xs arbitrary: ys zs*)  
**case** (*Cons x xs ys zs*)  
**from** *Cons.IH*[*of x # ys rev ys # zs*]  
**show** ?*case by* (*simp add: o\_def*)  
**qed** *simp\_all*  
**from this** [*of [] []*] **show** ?*thesis by simp*  
**qed**

**lemma** *set\_prefixes\_eq*: *set (prefixes xs)* = {*ys. prefix ys xs*}  
**by** *auto*

**lemma** *card\_set\_prefixes* [*simp*]: *card (set (prefixes xs))* = *Suc (length xs)*  
**by** (*subst distinct\_card*) *auto*

**lemma** *set\_prefixes\_append*:  
*set (prefixes (xs @ ys))* = *set (prefixes xs)*  $\cup$  {*xs @ ys' | ys'. ys'  $\in$  set (prefixes ys)*}

by (subst prefixes\_append, cases ys) auto

#### 1.4 Longest Common Prefix

**definition** Longest\_common\_prefix :: 'a list set  $\Rightarrow$  'a list **where**

Longest\_common\_prefix L = (ARG\_MAX length ps.  $\forall xs \in L. \text{prefix } ps \ xs$ )

**lemma** Longest\_common\_prefix\_ex:  $L \neq \{\}$   $\implies$

$\exists ps. (\forall xs \in L. \text{prefix } ps \ xs) \wedge (\forall qs. (\forall xs \in L. \text{prefix } qs \ xs) \longrightarrow \text{size } qs \leq \text{size } ps)$

(is  $\_ \implies \exists ps. ?P \ L \ ps$ )

**proof**(induction LEAST n.  $\exists xs \in L. n = \text{length } xs$  arbitrary: L)

**case** 0

**have**  $\square \in L$  **using** 0.hyps LeastI[of  $\lambda n. \exists xs \in L. n = \text{length } xs$ ]  $\langle L \neq \{\}$

**by** auto

**hence**  $?P \ L \ \square$  **by**(auto)

**thus** ?case ..

**next**

**case** (Suc n)

**let** ?EX =  $\lambda n. \exists xs \in L. n = \text{length } xs$

**obtain** x xs **where** xxs:  $x \# xs \in L$   $\text{size } xs = n$  **using** Suc.prem1 Suc.hyps(2)

**by**(metis LeastI\_ex[of ?EX] Suc\_length\_conv ex\_in\_conv)

**hence**  $\square \notin L$  **using** Suc.hyps(2) **by** auto

**show** ?case

**proof** (cases  $\forall xs \in L. \exists ys. xs = x \# ys$ )

**case** True

**let** ?L =  $\{ys. x \# ys \in L\}$

**have** 1: (LEAST n.  $\exists xs \in ?L. n = \text{length } xs$ ) = n

**using** xxs Suc.prem1 Suc.hyps(2) LeastIe[of ?EX]

**by** – (rule Least\_equality, fastforce+)

**have** 2:  $?L \neq \{\}$  **using**  $x \# xs \in L$  **by** auto

**from** Suc.hyps(1)[OF 1[symmetric] 2] **obtain** ps **where** IH:  $?P \ ?L \ ps \ ..$

{ **fix** qs

**assume**  $\forall qs. (\forall xa. x \# xa \in L \longrightarrow \text{prefix } qs \ xa) \longrightarrow \text{length } qs \leq \text{length } ps$

**and**  $\forall xs \in L. \text{prefix } qs \ xs$

**hence**  $\text{length } (tl \ qs) \leq \text{length } ps$

**by** (metis Cons\_prefix\_Cons hd\_Cons\_tl list.sel(2) Nil\_prefix)

**hence**  $\text{length } qs \leq \text{length } ps$  **by** auto

}

**hence**  $?P \ L \ (x \# ps)$  **using** True IH **by** auto

**thus** ?thesis ..

**next**

**case** False

**then obtain** y ys **where** yys:  $x \neq y \ y \# ys \in L$  **using**  $\langle \square \notin L \rangle$

**by** (auto) (metis list.exhaust)

**have**  $\forall qs. (\forall xs \in L. \text{prefix } qs \ xs) \longrightarrow qs = \square$  **using** yys  $x \# xs \in L$

**by** auto (metis Cons\_prefix\_Cons prefix\_Cons)



**hence** ?P L [] **by** auto  
**thus** ?thesis ..  
**qed**  
**qed**

**lemma** Longest\_common\_prefix\_unique:  $L \neq \{\}$   $\implies$   
 $\exists! ps. (\forall xs \in L. \text{prefix } ps \ xs) \wedge (\forall qs. (\forall xs \in L. \text{prefix } qs \ xs) \longrightarrow \text{size } qs \leq \text{size } ps)$   
**by**(rule ex\_ex11[OF Longest\_common\_prefix\_ex];  
 meson equals01 prefix\_length\_prefix\_prefix\_order.antisym)

**lemma** Longest\_common\_prefix\_eq:  
 $\llbracket L \neq \{\}; \forall xs \in L. \text{prefix } ps \ xs;$   
 $\forall qs. (\forall xs \in L. \text{prefix } qs \ xs) \longrightarrow \text{size } qs \leq \text{size } ps \rrbracket$   
 $\implies \text{Longest\_common\_prefix } L = ps$   
**unfolding** Longest\_common\_prefix\_def arg\_max\_def is\_arg\_max\_linorder  
**by**(rule some1\_equality[OF Longest\_common\_prefix\_unique]) auto

**lemma** Longest\_common\_prefix\_prefix:  
 $xs \in L \implies \text{prefix } (\text{Longest\_common\_prefix } L) \ xs$   
**unfolding** Longest\_common\_prefix\_def arg\_max\_def is\_arg\_max\_linorder  
**by**(rule someI2\_ex[OF Longest\_common\_prefix\_ex]) auto

**lemma** Longest\_common\_prefix\_longest:  
 $L \neq \{\} \implies \forall xs \in L. \text{prefix } ps \ xs \implies \text{length } ps \leq \text{length } (\text{Longest\_common\_prefix } L)$   
**unfolding** Longest\_common\_prefix\_def arg\_max\_def is\_arg\_max\_linorder  
**by**(rule someI2\_ex[OF Longest\_common\_prefix\_ex]) auto

**lemma** Longest\_common\_prefix\_max\_prefix:  
 $L \neq \{\} \implies \forall xs \in L. \text{prefix } ps \ xs \implies \text{prefix } ps \ (\text{Longest\_common\_prefix } L)$   
**by**(metis Longest\_common\_prefix\_prefix Longest\_common\_prefix\_longest  
 prefix\_length\_prefix\_ex\_in\_conv)

**lemma** Longest\_common\_prefix\_Nil:  $[] \in L \implies \text{Longest\_common\_prefix } L = []$   
**using** Longest\_common\_prefix\_prefix\_prefix\_Nil **by** blast

**lemma** Longest\_common\_prefix\_image\_Cons:  $L \neq \{\} \implies$   
 $\text{Longest\_common\_prefix } ((\#) \ x \ 'L) = x \# \text{Longest\_common\_prefix } L$   
**apply**(rule Longest\_common\_prefix\_eq)  
**apply**(simp)  
**apply** (simp add: Longest\_common\_prefix\_prefix)  
**apply** simp  
**by**(metis Longest\_common\_prefix\_longest[of L] Cons\_prefix\_Cons Nitpick.size\_list\_simp(2)  
 Suc\_le\_mono hd\_Cons\_tl order.strict\_implies\_order zero\_less\_Suc)

**lemma** Longest\_common\_prefix\_eq\_Cons: **assumes**  $L \neq \{\}$   $[] \notin L \ \forall xs \in L. \text{hd } xs = x$

**shows**  $\text{Longest\_common\_prefix } L = x \# \text{Longest\_common\_prefix } \{ys. x\#ys \in L\}$

**proof** –

**have**  $L = (\#) x \cdot \{ys. x\#ys \in L\}$  **using** *assms*(2,3)

**by** (*auto simp: image\_def*)(*metis hd\_Cons\_tl*)

**thus** *?thesis*

**by** (*metis Longest\_common\_prefix\_image\_Cons image\_is\_empty assms*(1))

**qed**

**lemma** *Longest\_common\_prefix\_eq\_Nil*:

$\llbracket x\#ys \in L; y\#zs \in L; x \neq y \rrbracket \implies \text{Longest\_common\_prefix } L = []$

**by** (*metis Longest\_common\_prefix\_prefix list.inject prefix\_Cons*)

**fun** *longest\_common\_prefix* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  'a list **where**

*longest\_common\_prefix* (x#xs) (y#ys) =

(if x=y then x # *longest\_common\_prefix* xs ys else []) |

*longest\_common\_prefix* \_ \_ = []

**lemma** *longest\_common\_prefix\_prefix1*:

*prefix* (*longest\_common\_prefix* xs ys) xs

**by**(*induction xs ys rule: longest\_common\_prefix.induct*) *auto*

**lemma** *longest\_common\_prefix\_prefix2*:

*prefix* (*longest\_common\_prefix* xs ys) ys

**by**(*induction xs ys rule: longest\_common\_prefix.induct*) *auto*

**lemma** *longest\_common\_prefix\_max\_prefix*:

$\llbracket \text{prefix } ps \ xs; \text{prefix } ps \ ys \rrbracket$

$\implies \text{prefix } ps \ (\text{longest\_common\_prefix } xs \ ys)$

**by**(*induction xs ys arbitrary: ps rule: longest\_common\_prefix.induct*)

(*auto simp: prefix\_Cons*)

## 1.5 Parallel lists

**definition** *parallel* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool (**infixl** || 50)

**where** (xs || ys) = ( $\neg \text{prefix } xs \ ys \wedge \neg \text{prefix } ys \ xs$ )

**lemma** *parallelI* [*intro*]:  $\neg \text{prefix } xs \ ys \implies \neg \text{prefix } ys \ xs \implies xs \ || \ ys$

**unfolding** *parallel\_def* **by** *blast*

**lemma** *parallelE* [*elim*]:

**assumes** xs || ys

**obtains**  $\neg \text{prefix } xs \ ys \wedge \neg \text{prefix } ys \ xs$

**using** *assms* **unfolding** *parallel\_def* **by** *blast*

**theorem** *prefix\_cases*:

**obtains** *prefix xs ys | strict\_prefix ys xs | xs || ys*  
**unfolding** *parallel\_def strict\_prefix\_def* **by** *blast*

**theorem** *parallel\_decomp*:

$xs \parallel ys \implies \exists as\ b\ bs\ c\ cs. b \neq c \wedge xs = as @ b \# bs \wedge ys = as @ c \# cs$

**proof** (*induct xs rule: rev\_induct*)

**case** *Nil*

**then have** *False* **by** *auto*

**then show** *?case ..*

**next**

**case** (*snoc x xs*)

**show** *?case*

**proof** (*rule prefix\_cases*)

**assume** *le: prefix xs ys*

**then obtain** *ys'* **where** *ys: ys = xs @ ys' ..*

**show** *?thesis*

**proof** (*cases ys'*)

**assume** *ys' = []*

**then show** *?thesis* **by** (*metis append\_Nil2 parallelE prefixI snoc.premys ys*)

**next**

**fix** *c cs* **assume** *ys': ys' = c # cs*

**have**  $x \neq c$  **using** *snoc.premys ys ys'* **by** *fastforce*

**thus**  $\exists as\ b\ bs\ c\ cs. b \neq c \wedge xs @ [x] = as @ b \# bs \wedge ys = as @ c \# cs$

**using** *ys ys'* **by** *blast*

**qed**

**next**

**assume** *strict\_prefix ys xs*

**then have** *prefix ys (xs @ [x])* **by** (*simp add: strict\_prefix\_def*)

**with** *snoc* **have** *False* **by** *blast*

**then show** *?thesis ..*

**next**

**assume**  $xs \parallel ys$

**with** *snoc* **obtain** *as b bs c cs* **where** *neq: (b::'a) ≠ c*

**and** *xs: xs = as @ b # bs* **and** *ys: ys = as @ c # cs*

**by** *blast*

**from** *xs* **have**  $xs @ [x] = as @ b \# (bs @ [x])$  **by** *simp*

**with** *neq ys* **show** *?thesis* **by** *blast*

**qed**

**qed**

**lemma** *parallel\_append*:  $a \parallel b \implies a @ c \parallel b @ d$

**apply** (*rule parallelI*)

**apply** (*erule parallelE, erule conjE,*

*induct rule: not\_prefix\_induct, simp+*)**+**

**done**

**lemma** *parallel\_appendI*:  $xs \parallel ys \implies x = xs @ xs' \implies y = ys @ ys' \implies x \parallel y$   
**by** (*simp add: parallel\_append*)

**lemma** *parallel\_commute*:  $a \parallel b \longleftrightarrow b \parallel a$   
**unfolding** *parallel\_def* **by** *auto*

## 1.6 Suffix order on lists

**definition** *suffix* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool  
**where** *suffix* xs ys = ( $\exists zs. ys = zs @ xs$ )

**definition** *strict\_suffix* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool  
**where** *strict\_suffix* xs ys  $\longleftrightarrow$  *suffix* xs ys  $\wedge$   $xs \neq ys$

**interpretation** *suffix\_order*: *order suffix strict\_suffix*  
**by** *standard* (*auto simp: suffix\_def strict\_suffix\_def*)

**interpretation** *suffix\_bot*: *order\_bot Nil suffix strict\_suffix*  
**by** *standard* (*simp add: suffix\_def*)

**lemma** *suffixI* [*intro?*]:  $ys = zs @ xs \implies \text{suffix } xs \text{ } ys$   
**unfolding** *suffix\_def* **by** *blast*

**lemma** *suffixE* [*elim?*]:  
**assumes** *suffix* xs ys  
**obtains** zs **where**  $ys = zs @ xs$   
**using** *assms* **unfolding** *suffix\_def* **by** *blast*

**lemma** *suffix\_tl* [*simp*]: *suffix* (tl xs) xs  
**by** (*induct* xs) (*auto simp: suffix\_def*)

**lemma** *strict\_suffix\_tl* [*simp*]:  $xs \neq [] \implies \text{strict\_suffix } (tl \text{ } xs) \text{ } xs$   
**by** (*induct* xs) (*auto simp: strict\_suffix\_def suffix\_def*)

**lemma** *Nil\_suffix* [*simp*]: *suffix* [] xs  
**by** (*simp add: suffix\_def*)

**lemma** *suffix\_Nil* [*simp*]: (*suffix* xs []) = ( $xs = []$ )  
**by** (*auto simp add: suffix\_def*)

**lemma** *suffix\_ConsI*: *suffix* xs ys  $\implies \text{suffix } xs \text{ } (y \# ys)$   
**by** (*auto simp add: suffix\_def*)

**lemma** *suffix\_ConsD*: *suffix* (x # xs) ys  $\implies \text{suffix } xs \text{ } ys$   
**by** (*auto simp add: suffix\_def*)

**lemma** *suffix\_appendI*:  $\text{suffix } xs \ ys \implies \text{suffix } xs \ (zs \ @ \ ys)$   
**by** (*auto simp add: suffix\_def*)

**lemma** *suffix\_appendD*:  $\text{suffix } (zs \ @ \ xs) \ ys \implies \text{suffix } xs \ ys$   
**by** (*auto simp add: suffix\_def*)

**lemma** *strict\_suffix\_set\_subset*:  $\text{strict\_suffix } xs \ ys \implies \text{set } xs \subseteq \text{set } ys$   
**by** (*auto simp: strict\_suffix\_def suffix\_def*)

**lemma** *set\_mono\_suffix*:  $\text{suffix } xs \ ys \implies \text{set } xs \subseteq \text{set } ys$   
**by** (*auto simp: suffix\_def*)

**lemma** *sorted\_antimono\_suffix*:  $\text{suffix } xs \ ys \implies \text{sorted } ys \implies \text{sorted } xs$   
**by** (*metis sorted\_append suffix\_def*)

**lemma** *suffix\_ConsD2*:  $\text{suffix } (x \ # \ xs) \ (y \ # \ ys) \implies \text{suffix } xs \ ys$   
**proof** –  
**assume** *suffix*  $(x \ # \ xs) \ (y \ # \ ys)$   
**then obtain** *zs* **where**  $y \ # \ ys = zs \ @ \ x \ # \ xs \ ..$   
**then show** *?thesis*  
**by** (*induct zs*) (*auto intro!: suffix\_appendI suffix\_ConsI*)  
**qed**

**lemma** *suffix\_to\_prefix* [*code*]:  $\text{suffix } xs \ ys \longleftrightarrow \text{prefix } (\text{rev } xs) \ (\text{rev } ys)$   
**proof**  
**assume** *suffix*  $xs \ ys$   
**then obtain** *zs* **where**  $ys = zs \ @ \ xs \ ..$   
**then have**  $\text{rev } ys = \text{rev } xs \ @ \ \text{rev } zs$  **by** *simp*  
**then show**  $\text{prefix } (\text{rev } xs) \ (\text{rev } ys) \ ..$   
**next**  
**assume**  $\text{prefix } (\text{rev } xs) \ (\text{rev } ys)$   
**then obtain** *zs* **where**  $\text{rev } ys = \text{rev } xs \ @ \ zs \ ..$   
**then have**  $\text{rev } (\text{rev } ys) = \text{rev } zs \ @ \ \text{rev } (\text{rev } xs)$  **by** *simp*  
**then have**  $ys = \text{rev } zs \ @ \ xs$  **by** *simp*  
**then show**  $\text{suffix } xs \ ys \ ..$   
**qed**

**lemma** *strict\_suffix\_to\_prefix* [*code*]:  $\text{strict\_suffix } xs \ ys \longleftrightarrow \text{strict\_prefix } (\text{rev } xs) \ (\text{rev } ys)$   
**by** (*auto simp: suffix\_to\_prefix strict\_suffix\_def strict\_prefix\_def*)

**lemma** *distinct\_suffix*:  $\text{distinct } ys \implies \text{suffix } xs \ ys \implies \text{distinct } xs$   
**by** (*clarsimp elim!: suffixE*)

**lemma** *map\_mono\_suffix*:  $\text{suffix } xs \ ys \implies \text{suffix } (\text{map } f \ xs) \ (\text{map } f \ ys)$   
**by** (*auto elim!*: *suffixE* *intro*: *suffixI*)

**lemma** *filter\_mono\_suffix*:  $\text{suffix } xs \ ys \implies \text{suffix } (\text{filter } P \ xs) \ (\text{filter } P \ ys)$   
**by** (*auto simp*: *suffix\_def*)

**lemma** *suffix\_drop*:  $\text{suffix } (\text{drop } n \ as) \ as$   
**unfolding** *suffix\_def* **by** (*rule exI* [**where**  $x = \text{take } n \ as$ ]) *simp*

**lemma** *suffix\_take*:  $\text{suffix } xs \ ys \implies ys = \text{take } (\text{length } ys - \text{length } xs) \ ys \ @ \ xs$   
**by** (*auto elim!*: *suffixE*)

**lemma** *strict\_suffix\_reflcp\_conv*:  $\text{strict\_suffix}^{\text{==}} = \text{suffix}$   
**by** (*intro ext*) (*auto simp*: *suffix\_def* *strict\_suffix\_def*)

**lemma** *suffix\_lists*:  $\text{suffix } xs \ ys \implies ys \in \text{lists } A \implies xs \in \text{lists } A$   
**unfolding** *suffix\_def* **by** *auto*

**lemma** *suffix\_snoc* [*simp*]:  $\text{suffix } xs \ (ys \ @ \ [y]) \longleftrightarrow xs = [] \vee (\exists zs. xs = zs \ @ \ [y] \wedge \text{suffix } zs \ ys)$   
**by** (*cases xs rule*: *rev\_cases*) (*auto simp*: *suffix\_def*)

**lemma** *snoc\_suffix\_snoc* [*simp*]:  $\text{suffix } (xs \ @ \ [x]) \ (ys \ @ \ [y]) = (x = y \wedge \text{suffix } xs \ ys)$   
**by** (*auto simp add*: *suffix\_def*)

**lemma** *same\_suffix\_suffix* [*simp*]:  $\text{suffix } (ys \ @ \ xs) \ (zs \ @ \ xs) = \text{suffix } ys \ zs$   
**by** (*simp add*: *suffix\_to\_prefix*)

**lemma** *same\_suffix\_nil* [*simp*]:  $\text{suffix } (ys \ @ \ xs) \ xs = (ys = [])$   
**by** (*simp add*: *suffix\_to\_prefix*)

**theorem** *suffix\_Cons*:  $\text{suffix } xs \ (y \ # \ ys) \longleftrightarrow xs = y \ # \ ys \vee \text{suffix } xs \ ys$   
**unfolding** *suffix\_def* **by** (*auto simp*: *Cons\_eq\_append\_conv*)

**theorem** *suffix\_append*:  
 $\text{suffix } xs \ (ys \ @ \ zs) \longleftrightarrow \text{suffix } xs \ zs \vee (\exists xs'. xs = xs' \ @ \ zs \wedge \text{suffix } xs' \ ys)$   
**by** (*auto simp*: *suffix\_def* *append\_eq\_append\_conv2*)

**theorem** *suffix\_length\_le*:  $\text{suffix } xs \ ys \implies \text{length } xs \leq \text{length } ys$   
**by** (*auto simp add*: *suffix\_def*)

**lemma** *suffix\_same\_cases*:  
 $\text{suffix } (xs_1 :: 'a \ \text{list}) \ ys \implies \text{suffix } xs_2 \ ys \implies \text{suffix } xs_1 \ xs_2 \vee \text{suffix } xs_2 \ xs_1$   
**unfolding** *suffix\_def* **by** (*force simp*: *append\_eq\_append\_conv2*)

```

lemma suffix_length_suffix:
  suffix ps xs  $\implies$  suffix qs xs  $\implies$  length ps  $\leq$  length qs  $\implies$  suffix ps qs
  by (auto simp: suffix_to_prefix intro: prefix_length_prefix)

lemma suffix_length_less: strict_suffix xs ys  $\implies$  length xs  $<$  length ys
  by (auto simp: strict_suffix_def suffix_def)

lemma suffix_ConsD': suffix (x#xs) ys  $\implies$  strict_suffix xs ys
  by (auto simp: strict_suffix_def suffix_def)

lemma drop_strict_suffix: strict_suffix xs ys  $\implies$  strict_suffix (drop n xs) ys
proof (induct n arbitrary: xs ys)
  case 0
  then show ?case by (cases ys) simp_all
next
  case (Suc n)
  then show ?case
    by (cases xs) (auto intro: Suc dest: suffix_ConsD' suffix_order.less_imp_le)
qed

lemma not_suffix_cases:
  assumes px:  $\neg$  suffix ps ls
  obtains
    (c1) ps  $\neq$  [] and ls = []
  | (c2) a as x xs where ps = as@[a] and ls = xs@[x] and x = a and  $\neg$  suffix as xs
  | (c3) a as x xs where ps = as@[a] and ls = xs@[x] and x  $\neq$  a
proof (cases ps rule: rev_cases)
  case Nil
  then show ?thesis using px by simp
next
  case (snoc as a)
  note c =  $\langle ps = as@[a] \rangle$ 
  show ?thesis
  proof (cases ls rule: rev_cases)
    case Nil then show ?thesis by (metis append_Nil2 px c1 same_suffix_nil)
  next
    case (snoc xs x)
    show ?thesis
    proof (cases x = a)
      case True
      have  $\neg$  suffix as xs using px c snoc True by simp
      with c snoc True show ?thesis by (rule c2)
    next
      case False
      with c snoc show ?thesis by (rule c3)

```

**qed**  
**qed**  
**qed**

**lemma** *not\_suffix\_induct* [*consumes 1, case\_names Nil Neq Eq*]:  
**assumes** *np*:  $\neg \text{suffix } ps \text{ } ls$   
**and** *base*:  $\bigwedge x \ xs. P (xs@[x]) \ []$   
**and** *r1*:  $\bigwedge x \ xs \ y \ ys. x \neq y \implies P (xs@[x]) (ys@[y])$   
**and** *r2*:  $\bigwedge x \ xs \ y \ ys. \llbracket x = y; \neg \text{suffix } xs \ ys; P \ xs \ ys \rrbracket \implies P (xs@[x]) (ys@[y])$   
**shows**  $P \ ps \ ls$  **using** *np*  
**proof** (*induct ls arbitrary: ps rule: rev\_induct*)  
**case** *Nil*  
**then show** *?case* **by** (*cases ps rule: rev\_cases*) (*auto intro: base*)  
**next**  
**case** (*snoc y ys ps*)  
**then have** *npfx*:  $\neg \text{suffix } ps \ (ys \ @ \ [y])$  **by** *simp*  
**then obtain** *x xs* **where** *pv*:  $ps = xs \ @ \ [x]$   
**by** (*rule not\_suffix\_cases*) *auto*  
**show** *?case* **by** (*metis snoc.hyps snoc\_suffix\_snoc npfx pv r1 r2*)  
**qed**

**lemma** *parallelD1*:  $x \parallel y \implies \neg \text{prefix } x \ y$   
**by** *blast*

**lemma** *parallelD2*:  $x \parallel y \implies \neg \text{prefix } y \ x$   
**by** *blast*

**lemma** *parallel\_Nil1* [*simp*]:  $\neg x \parallel []$   
**unfolding** *parallel\_def* **by** *simp*

**lemma** *parallel\_Nil2* [*simp*]:  $\neg [] \parallel x$   
**unfolding** *parallel\_def* **by** *simp*

**lemma** *Cons\_parallelI1*:  $a \neq b \implies a \# as \parallel b \# bs$   
**by** *auto*

**lemma** *Cons\_parallelI2*:  $\llbracket a = b; as \parallel bs \rrbracket \implies a \# as \parallel b \# bs$   
**by** (*metis Cons\_prefix\_Cons parallelE parallelI*)

**lemma** *not\_equal\_is\_parallel*:  
**assumes** *neq*:  $xs \neq ys$   
**and** *len*:  $\text{length } xs = \text{length } ys$   
**shows**  $xs \parallel ys$   
**using** *len neq*



```

proof (induct rule: list_induct2)
  case Nil
  then show ?case by simp
next
  case (Cons a as b bs)
  have ih: as ≠ bs ⇒ as || bs by fact
  show ?case
  proof (cases a = b)
    case True
    then have as ≠ bs using Cons by simp
    then show ?thesis by (rule Cons_parallelI2 [OF True ih])
  next
  case False
  then show ?thesis by (rule Cons_parallelI1)
qed
qed

```

## 1.7 Suffixes

**primrec** *suffixes* **where**

```

suffixes [] = [[]]
| suffixes (x#xs) = suffixes xs @ [x # xs]

```

**lemma** *in\_set\_suffixes* [simp]:  $xs \in \text{set } (\text{suffixes } ys) \longleftrightarrow \text{suffix } xs \text{ } ys$   
**by** (induction ys) (auto simp: suffix\_def Cons\_eq\_append\_conv)

**lemma** *distinct\_suffixes* [intro]: *distinct* (*suffixes* xs)  
**by** (induction xs) (auto simp: suffix\_def)

**lemma** *length\_suffixes* [simp]:  $\text{length } (\text{suffixes } xs) = \text{Suc } (\text{length } xs)$   
**by** (induction xs) auto

**lemma** *suffixes\_snoc* [simp]:  $\text{suffixes } (xs @ [x]) = [] \# \text{map } (\lambda ys. ys @ [x]) (\text{suffixes } xs)$   
**by** (induction xs) auto

**lemma** *suffixes\_not\_Nil* [simp]: *suffixes* xs ≠ []  
**by** (cases xs) auto

**lemma** *hd\_suffixes* [simp]:  $\text{hd } (\text{suffixes } xs) = []$   
**by** (induction xs) simp\_all

**lemma** *last\_suffixes* [simp]:  $\text{last } (\text{suffixes } xs) = xs$   
**by** (cases xs) simp\_all

**lemma** *suffixes\_append*:

```

suffixes (xs @ ys) = suffixes ys @ map ( $\lambda xs'. xs' @ ys$ ) (tl (suffixes xs))

```

```

proof (induction ys rule: rev_induct)
  case Nil
  thus ?case by (cases xs rule: rev_cases) auto
next
  case (snoc y ys)
  show ?case
  by (simp only: append.assoc [symmetric] suffixes_snoc snoc.IH) simp
qed

```

```

lemma suffixes_eq_snoc:
  suffixes ys = xs @ [x]  $\longleftrightarrow$ 
  (ys = []  $\wedge$  xs = []  $\vee$  ( $\exists z$  zs. ys = z#zs  $\wedge$  xs = suffixes zs))  $\wedge$  x = ys
  by (cases ys) auto

```

```

lemma suffixes_tailrec [code]:
  suffixes xs = rev (snd (foldl ( $\lambda$ (acc1, acc2) x. (x#acc1, (x#acc1)#acc2)) ([], []))
  (rev xs))
proof –
  have foldl ( $\lambda$ (acc1, acc2) x. (x#acc1, (x#acc1)#acc2)) (ys, ys # zs) (rev xs) =
    (xs @ ys, rev (map ( $\lambda$ as. as @ ys) (suffixes xs)) @ zs) for ys zs
  proof (induction xs arbitrary: ys zs)
    case (Cons x xs ys zs)
    from Cons.IH[of ys zs]
    show ?case by (simp add: o_def case_prod_unfold)
  qed simp_all
  from this [of [] []] show ?thesis by simp
qed

```

```

lemma set_suffixes_eq: set (suffixes xs) = {ys. suffix ys xs}
  by auto

```

```

lemma card_set_suffixes [simp]: card (set (suffixes xs)) = Suc (length xs)
  by (subst distinct_card) auto

```

```

lemma set_suffixes_append:
  set (suffixes (xs @ ys)) = set (suffixes ys)  $\cup$  {xs' @ ys | xs'. xs'  $\in$  set (suffixes xs)}
  by (subst suffixes_append, cases xs rule: rev_cases) auto

```

```

lemma suffixes_conv_prefixes: suffixes xs = map rev (prefixes (rev xs))
  by (induction xs) auto

```

```

lemma prefixes_conv_suffixes: prefixes xs = map rev (suffixes (rev xs))
  by (induction xs) auto

```

**lemma** *prefixes\_rev*:  $prefixes (rev\ xs) = map\ rev (suffixes\ xs)$   
**by** (*induction xs*) *auto*

**lemma** *suffixes\_rev*:  $suffixes (rev\ xs) = map\ rev (prefixes\ xs)$   
**by** (*induction xs*) *auto*

### 1.8 Homeomorphic embedding on lists

**inductive** *list\_emb* ::  $('a \Rightarrow 'a \Rightarrow bool) \Rightarrow 'a\ list \Rightarrow 'a\ list \Rightarrow bool$

**for**  $P :: ('a \Rightarrow 'a \Rightarrow bool)$

**where**

*list\_emb\_Nil* [*intro*, *simp*]:  $list\_emb\ P\ []\ ys$

| *list\_emb\_Cons* [*intro*]:  $list\_emb\ P\ xs\ ys \Longrightarrow list\_emb\ P\ xs\ (y\#\!ys)$

| *list\_emb\_Cons2* [*intro*]:  $P\ x\ y \Longrightarrow list\_emb\ P\ xs\ ys \Longrightarrow list\_emb\ P\ (x\#\!xs)\ (y\#\!ys)$

**lemma** *list\_emb\_mono*:

**assumes**  $\bigwedge x\ y. P\ x\ y \longrightarrow Q\ x\ y$

**shows**  $list\_emb\ P\ xs\ ys \longrightarrow list\_emb\ Q\ xs\ ys$

**proof**

**assume**  $list\_emb\ P\ xs\ ys$

**then show**  $list\_emb\ Q\ xs\ ys$  **by** (*induct*) (*auto simp: assms*)

**qed**

**lemma** *list\_emb\_Nil2* [*simp*]:

**assumes**  $list\_emb\ P\ xs\ []$  **shows**  $xs = []$

**using** *assms* **by** (*cases rule: list\_emb.cases*) *auto*

**lemma** *list\_emb\_refl*:

**assumes**  $\bigwedge x. x \in set\ xs \Longrightarrow P\ x\ x$

**shows**  $list\_emb\ P\ xs\ xs$

**using** *assms* **by** (*induct xs*) *auto*

**lemma** *list\_emb\_Cons\_Nil* [*simp*]:  $list\_emb\ P\ (x\#\!xs)\ [] = False$

**proof** –

{ **assume**  $list\_emb\ P\ (x\#\!xs)\ []$

**from** *list\_emb\_Nil2* [*OF this*] **have**  $False$  **by** *simp*

} **moreover** {

**assume**  $False$

**then have**  $list\_emb\ P\ (x\#\!xs)\ []$  **by** *simp*

} **ultimately show** *?thesis* **by** *blast*

**qed**

**lemma** *list\_emb\_append2* [*intro*]:  $list\_emb\ P\ xs\ ys \Longrightarrow list\_emb\ P\ xs\ (zs\ @\ ys)$

**by** (*induct zs*) *auto*

**lemma** *list\_emb\_prefix* [*intro*]:

**assumes**  $list\_emb\ P\ xs\ ys$  **shows**  $list\_emb\ P\ xs\ (ys\ @\ zs)$   
**using**  $assms$   
**by**  $(induct\ arbitrary:\ zs)\ auto$

**lemma**  $list\_emb\_ConsD$ :  
**assumes**  $list\_emb\ P\ (x\ \#\ xs)\ ys$   
**shows**  $\exists\ us\ v\ vs.\ ys = us\ @\ v\ \#\ vs \wedge P\ x\ v \wedge list\_emb\ P\ xs\ vs$   
**using**  $assms$   
**proof**  $(induct\ x\ \stackrel{def}{=} x\ \#\ xs\ ys\ arbitrary:\ x\ xs)$   
**case**  $list\_emb\_Cons$   
**then show**  $?case$  **by**  $(metis\ append\_Cons)$   
**next**  
**case**  $(list\_emb\_Cons2\ x\ y\ xs\ ys)$   
**then show**  $?case$  **by**  $blast$   
**qed**

**lemma**  $list\_emb\_appendD$ :  
**assumes**  $list\_emb\ P\ (xs\ @\ ys)\ zs$   
**shows**  $\exists\ us\ vs.\ zs = us\ @\ vs \wedge list\_emb\ P\ xs\ us \wedge list\_emb\ P\ ys\ vs$   
**using**  $assms$   
**proof**  $(induction\ xs\ arbitrary:\ ys\ zs)$   
**case**  $Nil$  **then show**  $?case$  **by**  $auto$   
**next**  
**case**  $(Cons\ x\ xs)$   
**then obtain**  $us\ v\ vs$  **where**  
 $zs:\ zs = us\ @\ v\ \#\ vs$  **and**  $p:\ P\ x\ v$  **and**  $lh:\ list\_emb\ P\ (xs\ @\ ys)\ vs$   
**by**  $(auto\ dest:\ list\_emb\_ConsD)$   
**obtain**  $sk_0 :: 'a\ list \Rightarrow 'a\ list \Rightarrow 'a\ list$  **and**  $sk_1 :: 'a\ list \Rightarrow 'a\ list \Rightarrow 'a\ list$  **where**  
 $sk:\ \forall\ x_0\ x_1.\ \neg\ list\_emb\ P\ (xs\ @\ x_0)\ x_1 \vee sk_0\ x_0\ x_1\ @\ sk_1\ x_0\ x_1 = x_1 \wedge list\_emb\ P\ xs$   
 $(sk_0\ x_0\ x_1) \wedge list\_emb\ P\ x_0\ (sk_1\ x_0\ x_1)$   
**using**  $Cons(1)$  **by**  $(metis\ (no\_types))$   
**hence**  $\forall\ x_2.\ list\_emb\ P\ (x\ \#\ xs)\ (x_2\ @\ v\ \#\ sk_0\ ys\ vs)$  **using**  $p\ lh$  **by**  $auto$   
**thus**  $?case$  **using**  $lh\ zs\ sk$  **by**  $(metis\ (no\_types)\ append\_Cons\ append\_assoc)$   
**qed**

**lemma**  $list\_emb\_strict\_suffix$ :  
**assumes**  $list\_emb\ P\ xs\ ys$  **and**  $strict\_suffix\ ys\ zs$   
**shows**  $list\_emb\ P\ xs\ zs$   
**using**  $assms(2)$  **and**  $list\_emb\_append2$   $[OF\ assms(1)]$  **by**  $(auto\ simp:\ strict\_suffix\_def\ suffix\_def)$

**lemma**  $list\_emb\_suffix$ :  
**assumes**  $list\_emb\ P\ xs\ ys$  **and**  $suffix\ ys\ zs$   
**shows**  $list\_emb\ P\ xs\ zs$   
**using**  $assms$  **and**  $list\_emb\_strict\_suffix$

**unfolding** *strict\_suffix\_reflclp\_conv*[*symmetric*] **by** *auto*

**lemma** *list\_emb\_length*: *list\_emb P xs ys*  $\implies$  *length xs*  $\leq$  *length ys*  
**by** (*induct rule: list\_emb.induct*) *auto*

**lemma** *list\_emb\_trans*:

**assumes**  $\bigwedge x y z. \llbracket x \in \text{set } xs; y \in \text{set } ys; z \in \text{set } zs; P x y; P y z \rrbracket \implies P x z$

**shows**  $\llbracket \text{list\_emb } P \text{ xs ys}; \text{list\_emb } P \text{ ys zs} \rrbracket \implies \text{list\_emb } P \text{ xs zs}$

**proof** –

**assume** *list\_emb P xs ys* **and** *list\_emb P ys zs*

**then show** *list\_emb P xs zs* **using** *assms*

**proof** (*induction arbitrary: zs*)

**case** *list\_emb\_Nil* **show** *?case* **by** *blast*

**next**

**case** (*list\_emb\_Cons xs ys y*)

**from** *list\_emb\_ConsD* [*OF* (*list\_emb P (y#ys) zs*)] **obtain** *us v vs*

**where** *zs*: *zs = us @ v # vs* **and**  $P = y v$  **and** *list\_emb P ys vs* **by** *blast*

**then have** *list\_emb P ys (v#vs)* **by** *blast*

**then have** *list\_emb P ys zs* **unfolding** *zs* **by** (*rule list\_emb\_append2*)

**from** *list\_emb\_Cons.IH* [*OF this*] **and** *list\_emb\_Cons.prem1* **show** *?case* **by** *auto*

**next**

**case** (*list\_emb\_Cons2 x y xs ys*)

**from** *list\_emb\_ConsD* [*OF* (*list\_emb P (y#ys) zs*)] **obtain** *us v vs*

**where** *zs*: *zs = us @ v # vs* **and**  $P y v$  **and** *list\_emb P ys vs* **by** *blast*

**with** *list\_emb\_Cons2* **have** *list\_emb P xs vs* **by** *auto*

**moreover have**  $P x v$

**proof** –

**from** *zs* **have**  $v \in \text{set } zs$  **by** *auto*

**moreover have**  $x \in \text{set } (x\#xs)$  **and**  $y \in \text{set } (y\#ys)$  **by** *simp\_all*

**ultimately show** *?thesis*

**using** ( $P x y$ ) **and** ( $P y v$ ) **and** *list\_emb\_Cons2*

**by** *blast*

**qed**

**ultimately have** *list\_emb P (x#xs) (v#vs)* **by** *blast*

**then show** *?case* **unfolding** *zs* **by** (*rule list\_emb\_append2*)

**qed**

**qed**

**lemma** *list\_emb\_set*:

**assumes** *list\_emb P xs ys* **and**  $x \in \text{set } xs$

**obtains** *y* **where**  $y \in \text{set } ys$  **and**  $P x y$

**using** *assms* **by** (*induct*) *auto*

**lemma** *list\_emb\_Cons\_iff1* [*simp*]:

**assumes**  $P x y$

**shows**  $list\_emb\ P\ (x\#\!xs)\ (y\#\!ys)\ \longleftrightarrow\ list\_emb\ P\ xs\ ys$   
**using** *assms* **by** (*subst list\_emb.simps*) (*auto dest: list\_emb\_ConsD*)

**lemma** *list\_emb\_Cons\_iff2* [*simp*]:  
**assumes**  $\neg P\ x\ y$   
**shows**  $list\_emb\ P\ (x\#\!xs)\ (y\#\!ys)\ \longleftrightarrow\ list\_emb\ P\ (x\#\!xs)\ ys$   
**using** *assms* **by** (*subst list\_emb.simps*) *auto*

**lemma** *list\_emb\_code* [*code*]:  
 $list\_emb\ P\ []\ ys\ \longleftrightarrow\ True$   
 $list\_emb\ P\ (x\#\!xs)\ []\ \longleftrightarrow\ False$   
 $list\_emb\ P\ (x\#\!xs)\ (y\#\!ys)\ \longleftrightarrow\ (if\ P\ x\ y\ then\ list\_emb\ P\ xs\ ys\ else\ list\_emb\ P\ (x\#\!xs)\ ys)$   
**by** *simp\_all*

### 1.9 Subsequences (special case of homeomorphic embedding)

**abbreviation** *subseq* :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool  
**where**  $subseq\ xs\ ys\ \stackrel{def}{=} list\_emb\ (=)\ xs\ ys$

**definition** *strict\_subseq* **where**  $strict\_subseq\ xs\ ys\ \longleftrightarrow\ xs\ \neq\ ys\ \wedge\ subseq\ xs\ ys$

**lemma** *subseq\_Cons2*:  $subseq\ xs\ ys\ \Longrightarrow\ subseq\ (x\#\!xs)\ (x\#\!ys)$  **by** *auto*

**lemma** *subseq\_same\_length*:  
**assumes**  $subseq\ xs\ ys$  **and**  $length\ xs = length\ ys$  **shows**  $xs = ys$   
**using** *assms* **by** (*induct*) (*auto dest: list\_emb\_length*)

**lemma** *not\_subseq\_length* [*simp*]:  $length\ ys < length\ xs\ \Longrightarrow\ \neg\ subseq\ xs\ ys$   
**by** (*metis list\_emb\_length linorder\_not\_less*)

**lemma** *subseq\_Cons'*:  $subseq\ (x\#\!xs)\ ys\ \Longrightarrow\ subseq\ xs\ ys$   
**by** (*induct xs, simp, blast dest: list\_emb\_ConsD*)

**lemma** *subseq\_Cons2'*:  
**assumes**  $subseq\ (x\#\!xs)\ (x\#\!ys)$  **shows**  $subseq\ xs\ ys$   
**using** *assms* **by** (*cases*) (*rule subseq\_Cons'*)

**lemma** *subseq\_Cons2\_neq*:  
**assumes**  $subseq\ (x\#\!xs)\ (y\#\!ys)$   
**shows**  $x\ \neq\ y\ \Longrightarrow\ subseq\ (x\#\!xs)\ ys$   
**using** *assms* **by** (*cases*) *auto*

**lemma** *subseq\_Cons2\_iff* [*simp*]:  
 $subseq\ (x\#\!xs)\ (y\#\!ys) = (if\ x = y\ then\ subseq\ xs\ ys\ else\ subseq\ (x\#\!xs)\ ys)$   
**by** *simp*

**lemma** *subseq\_append'*:  $subseq (zs @ xs) (zs @ ys) \longleftrightarrow subseq xs ys$   
**by** (*induct zs*) *simp\_all*

**interpretation** *subseq\_order*: *order subseq strict\_subseq*

**proof**

**fix** *xs ys* :: 'a list

{

**assume** *subseq xs ys* **and** *subseq ys xs*

**thus**  $xs = ys$

**proof** (*induct*)

**case** *list\_emb\_Nil*

**from** *list\_emb\_Nil2* [*OF this*] **show** ?*case* **by** *simp*

**next**

**case** *list\_emb\_Cons2*

**thus** ?*case* **by** *simp*

**next**

**case** *list\_emb\_Cons*

**hence** *False* **using** *subseq\_Cons'* **by** *fastforce*

**thus** ?*case* ..

**qed**

}

**thus**  $strict\_subseq\ xs\ ys \longleftrightarrow (subseq\ xs\ ys \wedge \neg subseq\ ys\ xs)$

**by** (*auto simp: strict\_subseq\_def*)

**qed** (*auto simp: list\_emb\_refl intro: list\_emb\_trans*)

**lemma** *in\_set\_subseqs* [*simp*]:  $xs \in set (subseqs\ ys) \longleftrightarrow subseq\ xs\ ys$

**proof**

**assume**  $xs \in set (subseqs\ ys)$

**thus** *subseq xs ys*

**by** (*induction ys arbitrary: xs*) (*auto simp: Let\_def*)

**next**

**have** [*simp*]:  $\square \in set (subseqs\ ys)$  **for** *ys* :: 'a list

**by** (*induction ys*) (*auto simp: Let\_def*)

**assume** *subseq xs ys*

**thus**  $xs \in set (subseqs\ ys)$

**by** (*induction xs ys rule: list\_emb.induct*) (*auto simp: Let\_def*)

**qed**

**lemma** *set\_subseqs\_eq*:  $set (subseqs\ ys) = \{xs.\ subseq\ xs\ ys\}$

**by** *auto*

**lemma** *subseq\_append\_le\_same\_iff*:  $subseq (xs @ ys) ys \longleftrightarrow xs = \square$

**by** (*auto dest: list\_emb\_length*)

**lemma** *subseq\_singleton\_left*:  $subseq [x] ys \longleftrightarrow x \in set\ ys$   
**by** (*fastforce dest: list\_emb\_ConsD split\_list\_last*)

**lemma** *list\_emb\_append\_mono*:  
 $\llbracket list\_emb\ P\ xs\ xs'; list\_emb\ P\ ys\ ys' \rrbracket \implies list\_emb\ P\ (xs@ys)\ (xs'@ys')$   
**by** (*induct rule: list\_emb.induct*) *auto*

**lemma** *prefix\_imp\_subseq* [*intro*]:  $prefix\ xs\ ys \implies subseq\ xs\ ys$   
**by** (*auto simp: prefix\_def*)

**lemma** *suffix\_imp\_subseq* [*intro*]:  $suffix\ xs\ ys \implies subseq\ xs\ ys$   
**by** (*auto simp: suffix\_def*)

### 1.10 Appending elements

**lemma** *subseq\_append* [*simp*]:  
 $subseq\ (xs\ @\ zs)\ (ys\ @\ zs) \longleftrightarrow subseq\ xs\ ys\ (is\ ?l = ?r)$   
**proof**  
{ **fix**  $xs'\ ys'\ xs\ ys\ zs :: 'a\ list$  **assume**  $subseq\ xs'\ ys'$   
**then have**  $xs' = xs\ @\ zs \wedge ys' = ys\ @\ zs \longrightarrow subseq\ xs\ ys$   
**proof** (*induct arbitrary: xs ys zs*)  
**case** *list\_emb\_Nil* **show** *?case* **by** *simp*  
**next**  
**case** (*list\_emb\_Cons*  $xs'\ ys'\ x$ )  
{ **assume**  $ys = []$  **then have** *?case* **using** *list\_emb\_Cons(1)* **by** *auto* }  
**moreover**  
{ **fix**  $us$  **assume**  $ys = x\ #\ us$   
**then have** *?case* **using** *list\_emb\_Cons(2)* **by** (*simp add: list\_emb.list\_emb\_Cons*)  
} }  
**ultimately show** *?case* **by** (*auto simp: Cons\_eq\_append\_conv*)  
**next**  
**case** (*list\_emb\_Cons2*  $x\ y\ xs'\ ys'$ )  
{ **assume**  $xs = []$  **then have** *?case* **using** *list\_emb\_Cons2(1)* **by** *auto* }  
**moreover**  
{ **fix**  $us\ vs$  **assume**  $xs = x\ #\ us\ ys = x\ #\ vs$  **then have** *?case* **using** *list\_emb\_Cons2* **by**  
*auto* }  
**moreover**  
{ **fix**  $us$  **assume**  $xs = x\ #\ us\ ys = []$  **then have** *?case* **using** *list\_emb\_Cons2(2)* **by**  
*bestsimp* }  
**ultimately show** *?case* **using**  $\langle (=)\ x\ y \rangle$  **by** (*auto simp: Cons\_eq\_append\_conv*)  
**qed** }  
**moreover assume** *?l*  
**ultimately show** *?r* **by** *blast*  
**next**  
**assume** *?r* **then show** *?l* **by** (*metis list\_emb\_append\_mono subseq\_order.order\_refl*)  
**qed**



**lemma** *subseq\_append\_iff*:  
 $subseq\ xs\ (ys\ @\ zs) \longleftrightarrow (\exists\ xs1\ xs2. xs = xs1\ @\ xs2 \wedge subseq\ xs1\ ys \wedge subseq\ xs2\ zs)$   
*(is ?lhs = ?rhs)*  
**proof**  
**assume** *?lhs thus ?rhs*  
**proof** (*induction xs ys @ zs arbitrary: ys zs rule: list\_emb.induct*)  
**case** (*list\_emb\_Cons xs ws y ys zs*)  
**from** *list\_emb\_Cons(2)[of tl ys zs]* **and** *list\_emb\_Cons(2)[of [] tl zs]* **and** *list\_emb\_Cons(1,3)*  
**show** *?case by (cases ys) auto*  
**next**  
**case** (*list\_emb\_Cons2 x y xs ws ys zs*)  
**from** *list\_emb\_Cons2(3)[of tl ys zs]* **and** *list\_emb\_Cons2(3)[of [] tl zs]*  
**and** *list\_emb\_Cons2(1,2,4)*  
**show** *?case by (cases ys) (auto simp: Cons\_eq\_append\_conv)*  
**qed** *auto*  
**qed** (*auto intro: list\_emb\_append\_mono*)

**lemma** *subseq\_appendE* [*case\_names append*]:  
**assumes** *subseq xs (ys @ zs)*  
**obtains** *xs1 xs2* **where**  $xs = xs1\ @\ xs2$  *subseq xs1 ys subseq xs2 zs*  
**using** *assms* **by** (*subst (asm) subseq\_append\_iff*) *auto*

**lemma** *subseq\_drop\_many*:  $subseq\ xs\ ys \implies subseq\ xs\ (zs\ @\ ys)$   
**by** (*induct zs*) *auto*

**lemma** *subseq\_rev\_drop\_many*:  $subseq\ xs\ ys \implies subseq\ xs\ (ys\ @\ zs)$   
**by** (*metis append\_Nil2 list\_emb\_Nil list\_emb\_append\_mono*)

### 1.11 Relation to standard list operations

**lemma** *subseq\_map*:  
**assumes** *subseq xs ys* **shows**  $subseq\ (map\ f\ xs)\ (map\ f\ ys)$   
**using** *assms* **by** (*induct*) *auto*

**lemma** *subseq\_filter\_left* [*simp*]:  $subseq\ (filter\ P\ xs)\ xs$   
**by** (*induct xs*) *auto*

**lemma** *subseq\_filter* [*simp*]:  
**assumes** *subseq xs ys* **shows**  $subseq\ (filter\ P\ xs)\ (filter\ P\ ys)$   
**using** *assms* **by** *induct auto*

**lemma** *subseq\_conv\_nth*:  
 $subseq\ xs\ ys \longleftrightarrow (\exists\ N. xs = nth\ ys\ N)$  (**is** *?L = ?R*)

**proof**  
**assume** *?L*

```

then show ?R
proof (induct)
  case list_emb_Nil show ?case by (metis nth_empty)
next
  case (list_emb_Cons xs ys x)
  then obtain N where xs = nth ys N by blast
  then have xs = nth (x#ys) (Suc ` N)
    by (clarsimp simp add: nth_Cons inj_image_mem_iff)
  then show ?case by blast
next
  case (list_emb_Cons2 x y xs ys)
  then obtain N where xs = nth ys N by blast
  then have x#xs = nth (x#ys) (insert 0 (Suc ` N))
    by (clarsimp simp add: nth_Cons inj_image_mem_iff)
  moreover from list_emb_Cons2 have x = y by simp
  ultimately show ?case by blast
qed
next
assume ?R
  then obtain N where xs = nth ys N ..
  moreover have subseq (nth ys N) ys
  proof (induct ys arbitrary: N)
    case Nil show ?case by simp
  next
    case Cons then show ?case by (auto simp: nth_Cons)
  qed
  ultimately show ?L by simp
qed

```

## 1.12 Contiguous sublists

**definition** sublist :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool **where**  
 sublist xs ys = ( $\exists$  ps ss. ys = ps @ xs @ ss)

**definition** strict\_sublist :: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool **where**  
 strict\_sublist xs ys  $\longleftrightarrow$  sublist xs ys  $\wedge$  xs  $\neq$  ys

**interpretation** sublist\_order: order sublist strict\_sublist

**proof**

```

fix xs ys zs :: 'a list
assume sublist xs ys sublist ys zs
then obtain xs1 xs2 ys1 ys2 where ys = xs1 @ xs @ xs2 zs = ys1 @ ys @ ys2
  by (auto simp: sublist_def)
hence zs = (ys1 @ xs1) @ xs @ (xs2 @ ys2) by simp
thus sublist xs zs unfolding sublist_def by blast
next

```

```

fix xs ys :: 'a list
{
  assume sublist xs ys sublist ys xs
  then obtain as bs cs ds
  where xs: xs = as @ ys @ bs and ys: ys = cs @ xs @ ds
  by (auto simp: sublist_def)
  have xs = as @ cs @ xs @ ds @ bs by (subst xs, subst ys) auto
  also have length ... = length as + length cs + length xs + length bs + length ds
  by simp
  finally have as = [] bs = [] by simp_all
  with xs show xs = ys by simp
}
thus strict_sublist xs ys  $\longleftrightarrow$  (sublist xs ys  $\wedge$   $\neg$ sublist ys xs)
by (auto simp: strict_sublist_def)
qed (auto simp: strict_sublist_def sublist_def intro: exI[of _ []])

lemma sublist_Nil_left [simp, intro]: sublist [] ys
by (auto simp: sublist_def)

lemma sublist_Cons_Nil [simp]:  $\neg$ sublist (x#xs) []
by (auto simp: sublist_def)

lemma sublist_Nil_right [simp]: sublist xs []  $\longleftrightarrow$  xs = []
by (cases xs) auto

lemma sublist_appendI [simp, intro]: sublist xs (ps @ xs @ ss)
by (auto simp: sublist_def)

lemma sublist_append_leftI [simp, intro]: sublist xs (ps @ xs)
by (auto simp: sublist_def intro: exI[of _ []])

lemma sublist_append_rightI [simp, intro]: sublist xs (xs @ ss)
by (auto simp: sublist_def intro: exI[of _ []])

lemma sublist_altdef: sublist xs ys  $\longleftrightarrow$  ( $\exists$  ys'. prefix ys' ys  $\wedge$  suffix xs ys')
proof safe
  assume sublist xs ys
  then obtain ps ss where ys = ps @ xs @ ss by (auto simp: sublist_def)
  thus  $\exists$  ys'. prefix ys' ys  $\wedge$  suffix xs ys'
  by (intro exI[of _ ps @ xs] conjI suffix_appendI) auto
next
  fix ys'
  assume prefix ys' ys suffix xs ys'
  thus sublist xs ys by (auto simp: prefix_def suffix_def)
qed
    
```

**lemma** *sublist\_altdef'*:  $sublist\ xs\ ys \longleftrightarrow (\exists\ ys'.\ suffix\ ys'\ ys \wedge prefix\ xs\ ys')$

**proof** *safe*

**assume** *sublist xs ys*

**then obtain** *ps ss* **where**  $ys = ps @ xs @ ss$  **by** (*auto simp: sublist\_def*)

**thus**  $\exists\ ys'.\ suffix\ ys'\ ys \wedge prefix\ xs\ ys'$

**by** (*intro exI[of \_ xs @ ss] conjI suffixI*) *auto*

**next**

**fix** *ys'*

**assume** *suffix ys' ys prefix xs ys'*

**thus** *sublist xs ys* **by** (*auto simp: prefix\_def suffix\_def*)

**qed**

**lemma** *sublist\_Cons\_right*:  $sublist\ xs\ (y\ \#\ ys) \longleftrightarrow prefix\ xs\ (y\ \#\ ys) \vee sublist\ xs\ ys$

**by** (*auto simp: sublist\_def prefix\_def Cons\_eq\_append\_conv*)

**lemma** *sublist\_code* [*code*]:

*sublist [] ys*  $\longleftrightarrow$  *True*

*sublist (x # xs) []*  $\longleftrightarrow$  *False*

*sublist (x # xs) (y # ys)*  $\longleftrightarrow$  *prefix (x # xs) (y # ys)  $\vee$  sublist (x # xs) ys*

**by** (*simp\_all add: sublist\_Cons\_right*)

**lemma** *sublist\_append*:

*sublist xs (ys @ zs)*  $\longleftrightarrow$

*sublist xs ys  $\vee$  sublist xs zs  $\vee$  ( $\exists\ xs1\ xs2.\ xs = xs1 @ xs2 \wedge suffix\ xs1\ ys \wedge prefix\ xs2\ zs$ )*

**by** (*auto simp: sublist\_altdef prefix\_append suffix\_append*)

**primrec** *sublists* :: 'a list  $\Rightarrow$  'a list list **where**

*sublists []* = [[]]

| *sublists (x # xs)* = *sublists xs @ map ((#) x) (prefixes xs)*

**lemma** *in\_set\_sublists* [*simp*]:  $xs \in set\ (sublists\ ys) \longleftrightarrow sublist\ xs\ ys$

**by** (*induction ys arbitrary: xs*) (*auto simp: sublist\_Cons\_right prefix\_Cons*)

**lemma** *set\_sublists\_eq*:  $set\ (sublists\ xs) = \{ys.\ sublist\ ys\ xs\}$

**by** *auto*

**lemma** *length\_sublists* [*simp*]:  $length\ (sublists\ xs) = Suc\ (length\ xs * Suc\ (length\ xs)$

*div 2)*

**by** (*induction xs*) *simp\_all*

**lemma** *sublist\_length\_le*:  $sublist\ xs\ ys \Longrightarrow length\ xs \leq length\ ys$

**by** (*auto simp add: sublist\_def*)

**lemma** *set\_mono\_sublist*: *sublist xs ys*  $\implies$  *set xs*  $\subseteq$  *set ys*  
**by** (*auto simp add: sublist\_def*)

**lemma** *prefix\_imp\_sublist* [*simp, intro*]: *prefix xs ys*  $\implies$  *sublist xs ys*  
**by** (*auto simp: sublist\_def prefix\_def intro: exI[of \_ []]*)

**lemma** *suffix\_imp\_sublist* [*simp, intro*]: *suffix xs ys*  $\implies$  *sublist xs ys*  
**by** (*auto simp: sublist\_def suffix\_def intro: exI[of \_ []]*)

**lemma** *sublist\_take* [*simp, intro*]: *sublist (take n xs) xs*  
**by** (*rule prefix\_imp\_sublist*) (*simp\_all add: take\_is\_prefix*)

**lemma** *sublist\_drop* [*simp, intro*]: *sublist (drop n xs) xs*  
**by** (*rule suffix\_imp\_sublist*) (*simp\_all add: suffix\_drop*)

**lemma** *sublist\_tl* [*simp, intro*]: *sublist (tl xs) xs*  
**by** (*rule suffix\_imp\_sublist*) (*simp\_all add: suffix\_drop*)

**lemma** *sublist\_butlast* [*simp, intro*]: *sublist (butlast xs) xs*  
**by** (*rule prefix\_imp\_sublist*) (*simp\_all add: prefixeq\_butlast*)

**lemma** *sublist\_rev* [*simp*]: *sublist (rev xs) (rev ys) = sublist xs ys*  
**proof**

**assume** *sublist (rev xs) (rev ys)*  
**then obtain** *as bs* **where** *rev ys = as @ rev xs @ bs*  
**by** (*auto simp: sublist\_def*)  
**also have** *rev . . . = rev bs @ xs @ rev as* **by** *simp*  
**finally show** *sublist xs ys* **by** *simp*

**next**

**assume** *sublist xs ys*  
**then obtain** *as bs* **where** *ys = as @ xs @ bs*  
**by** (*auto simp: sublist\_def*)  
**also have** *rev . . . = rev bs @ rev xs @ rev as* **by** *simp*  
**finally show** *sublist (rev xs) (rev ys)* **by** *simp*

**qed**

**lemma** *sublist\_rev\_left*: *sublist (rev xs) ys = sublist xs (rev ys)*  
**by** (*subst sublist\_rev [symmetric]*) (*simp only: rev\_rev\_ident*)

**lemma** *sublist\_rev\_right*: *sublist xs (rev ys) = sublist (rev xs) ys*  
**by** (*subst sublist\_rev [symmetric]*) (*simp only: rev\_rev\_ident*)

**lemma** *snoc\_sublist\_snoc*:  
*sublist (xs @ [x]) (ys @ [y])*  $\longleftrightarrow$

$(x = y \wedge \text{suffix } xs \text{ } ys \vee \text{sublist } (xs @ [x]) \text{ } ys)$   
**by** (*subst* (1 2) *sublist\_rev* [*symmetric*])  
*(simp del: sublist\_rev add: sublist\_Cons\_right suffix\_to\_prefix)*

**lemma** *sublist\_snoc*:  
 $\text{sublist } xs \text{ } (ys @ [y]) \longleftrightarrow \text{suffix } xs \text{ } (ys @ [y]) \vee \text{sublist } xs \text{ } ys$   
**by** (*subst* (1 2) *sublist\_rev* [*symmetric*])  
*(simp del: sublist\_rev add: sublist\_Cons\_right suffix\_to\_prefix)*

**lemma** *sublist\_imp\_subseq* [*intro*]:  $\text{sublist } xs \text{ } ys \implies \text{subseq } xs \text{ } ys$   
**by** (*auto simp: sublist\_def*)

### 1.13 Parametricity

**context includes** *lifting\_syntax*

**begin**

**private lemma** *prefix\_primrec*:  
 $\text{prefix} = \text{rec\_list } (\lambda xs. \text{True}) (\lambda x \text{ } xs \text{ } xsa \text{ } ys. \text{case } ys \text{ of } [] \Rightarrow \text{False} \mid y \# ys \Rightarrow x = y \wedge xsa \text{ } ys)$   
**proof** (*intro ext, goal\_cases*)  
**case** (1 *xs ys*)  
**show** ?*case* **by** (*induction xs arbitrary: ys*) (*auto simp: prefix\_Cons split: list.splits*)  
**qed**

**private lemma** *sublist\_primrec*:  
 $\text{sublist} = (\lambda xs \text{ } ys. \text{rec\_list } (\lambda xs. xs = [])) (\lambda y \text{ } ys \text{ } ysa \text{ } xs. \text{prefix } xs \text{ } (y \# ys) \vee ysa \text{ } xs) \text{ } ys$   
**proof** (*intro ext, goal\_cases*)  
**case** (1 *xs ys*)  
**show** ?*case* **by** (*induction ys*) (*auto simp: sublist\_Cons\_right*)  
**qed**

**private lemma** *list\_emb\_primrec*:  
 $\text{list\_emb} = (\lambda uu \text{ } uua \text{ } uuaa. \text{rec\_list } (\lambda P \text{ } xs. \text{List.null } xs) (\lambda y \text{ } ys \text{ } ysa \text{ } P \text{ } xs. \text{case } xs \text{ of } [] \Rightarrow \text{True} \mid x \# xs \Rightarrow \text{if } P \text{ } x \text{ } y \text{ then } ysa \text{ } P \text{ } xs \text{ else } ysa \text{ } P \text{ } (x \# xs))) \text{ } uuaa \text{ } uu \text{ } uua$   
**proof** (*intro ext, goal\_cases*)  
**case** (1 *P xs ys*)  
**show** ?*case*  
**by** (*induction ys arbitrary: xs*)  
*(auto simp: list\_emb\_code List.null\_def split: list.splits)*  
**qed**

**lemma** *prefix\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique A*

**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ prefix\ prefix$   
**unfolding** *prefix\_primrec* **by** *transfer\_prover*

**lemma** *suffix\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ suffix\ suffix$   
**unfolding** *suffix\_to\_prefix* [*abs\_def*] **by** *transfer\_prover*

**lemma** *sublist\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ sublist\ sublist$   
**unfolding** *sublist\_primrec* **by** *transfer\_prover*

**lemma** *parallel\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ parallel\ parallel$   
**unfolding** *parallel\_def* **by** *transfer\_prover*

**lemma** *list\_emb\_transfer* [*transfer\_rule*]:  
 $((A\ ==\>\ A\ ==\>\ (=))\ ==\>\ list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))$   
*list\_emb* *list\_emb*  
**unfolding** *list\_emb\_primrec* **by** *transfer\_prover*

**lemma** *strict\_prefix\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ strict\_prefix\ strict\_prefix$   
**unfolding** *strict\_prefix\_def* **by** *transfer\_prover*

**lemma** *strict\_suffix\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ strict\_suffix\ strict\_suffix$   
**unfolding** *strict\_suffix\_def* **by** *transfer\_prover*

**lemma** *strict\_subseq\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ strict\_subseq\ strict\_subseq$   
**unfolding** *strict\_subseq\_def* **by** *transfer\_prover*

**lemma** *strict\_sublist\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique* *A*  
**shows**  $(list\_all2\ A\ ==\>\ list\_all2\ A\ ==\>\ (=))\ strict\_sublist\ strict\_sublist$   
**unfolding** *strict\_sublist\_def* **by** *transfer\_prover*

**lemma** *prefixes\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique A*  
**shows** (*list\_all2 A == => list\_all2 (list\_all2 A)*) *prefixes prefixes*  
**unfolding** *prefixes\_def* **by** *transfer\_prover*

**lemma** *suffixes\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique A*  
**shows** (*list\_all2 A == => list\_all2 (list\_all2 A)*) *suffixes suffixes*  
**unfolding** *suffixes\_def* **by** *transfer\_prover*

**lemma** *sublists\_transfer* [*transfer\_rule*]:  
**assumes** [*transfer\_rule*]: *bi\_unique A*  
**shows** (*list\_all2 A == => list\_all2 (list\_all2 A)*) *sublists sublists*  
**unfolding** *sublists\_def* **by** *transfer\_prover*

**end**

**end**

**theory** *Spec*  
**imports** *Main* *~/~/src/HOL/Library/Sublist*  
**begin**

## 2 Sequential Composition of Languages

**definition**

*Sequ* :: *string set*  $\Rightarrow$  *string set*  $\Rightarrow$  *string set* (*\_* ;; *\_* [*100,100*] *100*)

**where**

*A* ;; *B* = {*s1* @ *s2* | *s1 s2*. *s1*  $\in$  *A*  $\wedge$  *s2*  $\in$  *B*}

Two Simple Properties about Sequential Composition

**lemma** *Sequ\_empty\_string* [*simp*]:

**shows** *A* ;; {} = *A*

**and** {} ;; *A* = *A*

**by** (*simp\_all add: Sequ\_def*)

**lemma** *Sequ\_empty* [*simp*]:

**shows** *A* ;; {} = {}

**and** {} ;; *A* = {}

**by** (*simp\_all add: Sequ\_def*)

## 3 Semantic Derivative (Left Quotient) of Languages

**definition**

*Der* :: *char*  $\Rightarrow$  *string set*  $\Rightarrow$  *string set*



**where** $Der\ c\ A \stackrel{def}{=} \{s.\ c\ \#\ s \in A\}$ **definition** $Ders :: string \Rightarrow string\ set \Rightarrow string\ set$ **where** $Ders\ s\ A \stackrel{def}{=} \{s'.\ s\ @\ s' \in A\}$ **lemma** *Der\_null* [simp]:**shows**  $Der\ c\ \{\} = \{\}$ **unfolding** *Der\_def***by** *auto***lemma** *Der\_empty* [simp]:**shows**  $Der\ c\ \{\}\ = \{\}$ **unfolding** *Der\_def***by** *auto***lemma** *Der\_char* [simp]:**shows**  $Der\ c\ \{[d]\} = (if\ c = d\ then\ \{\}\ else\ \{\})$ **unfolding** *Der\_def***by** *auto***lemma** *Der\_union* [simp]:**shows**  $Der\ c\ (A \cup B) = Der\ c\ A \cup Der\ c\ B$ **unfolding** *Der\_def***by** *auto***lemma** *Der\_Sequ* [simp]:**shows**  $Der\ c\ (A ;; B) = (Der\ c\ A) ;; B \cup (if\ [] \in A\ then\ Der\ c\ B\ else\ \{\})$ **unfolding** *Der\_def Sequ\_def***by** (*auto simp add: Cons\_eq\_append\_conv*)

## 4 Kleene Star for Languages

**inductive-set** $Star :: string\ set \Rightarrow string\ set\ (-\star\ [101]\ 102)$ **for**  $A :: string\ set$ **where***start*[intro]:  $[] \in A\star$ *step*[intro]:  $[s1 \in A; s2 \in A\star] \Longrightarrow s1\ @\ s2 \in A\star$ **lemma** *Star\_cases*:

**shows**  $A^\star = \{\square\} \cup A ;; A^\star$   
**unfolding** *Sequ\_def*  
**by** (*auto*) (*metis Star.simps*)

**lemma** *Star\_decomp*:  
**assumes**  $c \# x \in A^\star$   
**shows**  $\exists s1\ s2. x = s1 @ s2 \wedge c \# s1 \in A \wedge s2 \in A^\star$   
**using** *assms*  
**by** (*induct*  $x \stackrel{def}{=} c \# x$  *rule: Star.induct*)  
(*auto simp add: append\_eq\_Cons\_conv*)

**lemma** *Star\_Der\_Sequ*:  
**shows**  $Der\ c\ (A^\star) \subseteq (Der\ c\ A) ;; A^\star$   
**unfolding** *Der\_def Sequ\_def*  
**by**(*auto simp add: Star\_decomp*)

**lemma** *Der\_star [simp]*:  
**shows**  $Der\ c\ (A^\star) = (Der\ c\ A) ;; A^\star$   
**proof** –  
**have**  $Der\ c\ (A^\star) = Der\ c\ (\{\square\} \cup A ;; A^\star)$   
**by** (*simp only: Star\_cases[symmetric]*)  
**also have**  $\dots = Der\ c\ (A ;; A^\star)$   
**by** (*simp only: Der\_union Der\_empty*) (*simp*)  
**also have**  $\dots = (Der\ c\ A) ;; A^\star \cup (\text{if } \square \in A \text{ then } Der\ c\ (A^\star) \text{ else } \{\})$   
**by** *simp*  
**also have**  $\dots = (Der\ c\ A) ;; A^\star$   
**using** *Star\_Der\_Sequ* **by** *auto*  
**finally show**  $Der\ c\ (A^\star) = (Der\ c\ A) ;; A^\star$ .  
**qed**

## 5 Regular Expressions

**datatype** *rexp* =  
*ZERO*  
| *ONE*  
| *CHAR char*  
| *SEQ rexp rexp*  
| *ALT rexp rexp*  
| *STAR rexp*

## 6 Semantics of Regular Expressions

**fun**

```

L :: rexp ⇒ string set
where
  | L (ZERO) = {}
  | L (ONE) = {[]}
  | L (CHAR c) = {[c]}
  | L (SEQ r1 r2) = (L r1) ;; (L r2)
  | L (ALT r1 r2) = (L r1) ∪ (L r2)
  | L (STAR r) = (L r)*

```

## 7 Nullable, Derivatives

```

fun
  nullable :: rexp ⇒ bool
where
  | nullable (ZERO) = False
  | nullable (ONE) = True
  | nullable (CHAR c) = False
  | nullable (ALT r1 r2) = (nullable r1 ∨ nullable r2)
  | nullable (SEQ r1 r2) = (nullable r1 ∧ nullable r2)
  | nullable (STAR r) = True

```

```

fun
  der :: char ⇒ rexp ⇒ rexp
where
  | der c (ZERO) = ZERO
  | der c (ONE) = ZERO
  | der c (CHAR d) = (if c = d then ONE else ZERO)
  | der c (ALT r1 r2) = ALT (der c r1) (der c r2)
  | der c (SEQ r1 r2) =
    (if nullable r1
     then ALT (SEQ (der c r1) r2) (der c r2)
     else SEQ (der c r1) r2)
  | der c (STAR r) = SEQ (der c r) (STAR r)

```

```

fun
  ders :: string ⇒ rexp ⇒ rexp
where
  | ders [] r = r
  | ders (c # s) r = ders s (der c r)

```

```

lemma nullable_correctness:
  shows nullable r ↔ [] ∈ (L r)
by (induct r) (auto simp add: Sequ_def)

```

**lemma** *der\_correctness*:  
**shows**  $L(\text{der } c \ r) = \text{Der } c \ (L \ r)$   
**by** (*induct r*) (*simp\_all add: nullable\_correctness*)

**lemma** *ders\_correctness*:  
**shows**  $L(\text{ders } s \ r) = \text{Ders } s \ (L \ r)$   
**by** (*induct s arbitrary: r*)  
*(simp\_all add: Ders\_def der\_correctness Der\_def)*

**lemma** *ders\_append*:  
**shows**  $\text{ders } (s1 \ @ \ s2) \ r = \text{ders } s2 \ (\text{ders } s1 \ r)$   
**apply**(*induct s1 arbitrary: s2 r*)  
**apply**(*auto*)  
**done**

## 8 Values

**datatype** *val* =  
*Void*  
| *Char char*  
| *Seq val val*  
| *Right val*  
| *Left val*  
| *Stars val list*

## 9 The string behind a value

**fun**  
*flat* :: *val*  $\Rightarrow$  *string*  
**where**  
*flat* (*Void*) = []  
| *flat* (*Char c*) = [c]  
| *flat* (*Left v*) = *flat v*  
| *flat* (*Right v*) = *flat v*  
| *flat* (*Seq v1 v2*) = (*flat v1*) @ (*flat v2*)  
| *flat* (*Stars []*) = []  
| *flat* (*Stars (v#vs)*) = (*flat v*) @ (*flat (Stars vs)*)

**abbreviation**  
*flats vs*  $\stackrel{\text{def}}{=} \text{concat } (\text{map } \text{flat } \text{vs})$

**lemma** *flat\_Stars* [*simp*]:  
*flat (Stars vs)* = *flats vs*  
**by** (*induct vs*) (*auto*)

**lemma** *Star\_concat*:

**assumes**  $\forall s \in \text{set } ss. s \in A$

**shows**  $\text{concat } ss \in A^*$

**using** *assms* **by** (*induct ss*) (*auto*)

**lemma** *Star\_cstring*:

**assumes**  $s \in A^*$

**shows**  $\exists ss. \text{concat } ss = s \wedge (\forall s \in \text{set } ss. s \in A \wedge s \neq [])$

**using** *assms*

**apply**(*induct rule: Star.induct*)

**apply**(*auto*)[*I*]

**apply**(*rule\_tac x=[] in exI*)

**apply**(*simp*)

**apply**(*erule exE*)

**apply**(*clarify*)

**apply**(*case\_tac sI = []*)

**apply**(*rule\_tac x=ss in exI*)

**apply**(*simp*)

**apply**(*rule\_tac x=sI#ss in exI*)

**apply**(*simp*)

**done**

## 10 Lexical Values

**inductive**

*Prf* :: *val*  $\Rightarrow$  *rexp*  $\Rightarrow$  *bool* ( $\models \_ : \_ [100, 100] 100$ )

**where**

$\llbracket \models v1 : r1 ; \models v2 : r2 \rrbracket \Longrightarrow \models \text{Seq } v1 \ v2 : \text{SEQ } r1 \ r2$

$\models v1 : r1 \Longrightarrow \models \text{Left } v1 : \text{ALT } r1 \ r2$

$\models v2 : r2 \Longrightarrow \models \text{Right } v2 : \text{ALT } r1 \ r2$

$\models \text{Void} : \text{ONE}$

$\models \text{Char } c : \text{CHAR } c$

$\models \forall v \in \text{set } vs. \models v : r \wedge \text{flat } v \neq [] \Longrightarrow \models \text{Stars } vs : \text{STAR } r$

**inductive-cases** *Prf\_elims*:

$\models v : \text{ZERO}$

$\models v : \text{SEQ } r1 \ r2$

$\models v : \text{ALT } r1 \ r2$

$\models v : \text{ONE}$

$\models v : \text{CHAR } c$

$\models vs : \text{STAR } r$

**lemma** *Prf\_Stars\_appendE*:

**assumes**  $\models \text{Stars } (vs1 @ vs2) : \text{STAR } r$

**shows**  $\models Stars\ vs1 : STAR\ r \wedge \models Stars\ vs2 : STAR\ r$   
**using** *assms*  
**by** (*auto intro: Prf.intros elim!: Prf.elims*)

**lemma** *Star\_cval*:

**assumes**  $\forall s \in set\ ss. \exists v. s = flat\ v \wedge \models v : r$   
**shows**  $\exists vs. flats\ vs = concat\ ss \wedge (\forall v \in set\ vs. \models v : r \wedge flat\ v \neq [])$   
**using** *assms*  
**apply**(*induct ss*)  
**apply**(*auto*)  
**apply**(*rule\_tac x=[] in exI*)  
**apply**(*simp*)  
**apply**(*case\_tac flat v = []*)  
**apply**(*rule\_tac x=vs in exI*)  
**apply**(*simp*)  
**apply**(*rule\_tac x=v#vs in exI*)  
**apply**(*simp*)  
**done**

**lemma** *L\_flat\_Prfl*:

**assumes**  $\models v : r$   
**shows**  $flat\ v \in L\ r$   
**using** *assms*  
**by** (*induct*) (*auto simp add: Sequ\_def Star\_concat*)

**lemma** *L\_flat\_Prfl2*:

**assumes**  $s \in L\ r$   
**shows**  $\exists v. \models v : r \wedge flat\ v = s$   
**using** *assms*  
**proof**(*induct r arbitrary: s*)  
**case** (*STAR r s*)  
**have** *IH*:  $\bigwedge s. s \in L\ r \implies \exists v. \models v : r \wedge flat\ v = s$  **by fact**  
**have**  $s \in L\ (STAR\ r)$  **by fact**  
**then obtain** *ss* **where**  $concat\ ss = s \wedge \forall s \in set\ ss. s \in L\ r \wedge s \neq []$   
**using** *Star\_cstring* **by auto**  
**then obtain** *vs* **where**  $flats\ vs = s \wedge \forall v \in set\ vs. \models v : r \wedge flat\ v \neq []$   
**using** *IH Star\_cval* **by metis**  
**then show**  $\exists v. \models v : STAR\ r \wedge flat\ v = s$   
**using** *Prf.intros(6) flat\_Stars* **by blast**  
**next**  
**case** (*SEQ r1 r2 s*)  
**then show**  $\exists v. \models v : SEQ\ r1\ r2 \wedge flat\ v = s$   
**unfolding** *Sequ\_def L.simps* **by** (*fastforce intro: Prf.intros*)

```

next
  case (ALT r1 r2 s)
  then show  $\exists v. \models v : ALT\ r1\ r2 \wedge flat\ v = s$ 
  unfolding L.simps by (fastforce intro: Prf.intros)
qed (auto intro: Prf.intros)

```

```

lemma L_flat_Prf:
  shows  $L(r) = \{flat\ v \mid v. \models v : r\}$ 
using L_flat_Prf1 L_flat_Prf2 by blast

```

## 11 Sets of Lexical Values

Shows that lexical values are finite for a given regex and string.

```

definition
  LV :: rexp  $\Rightarrow$  string  $\Rightarrow$  val set
where  $LV\ r\ s \stackrel{def}{=} \{v. \models v : r \wedge flat\ v = s\}$ 

```

```

lemma LV_simps:
  shows LV ZERO s = {}
  and LV ONE s = (if s = [] then {Void} else {})
  and LV (CHAR c) s = (if s = [c] then {Char c} else {})
  and LV (ALT r1 r2) s = Left ' LV r1 s  $\cup$  Right ' LV r2 s
unfolding LV_def
by (auto intro: Prf.intros elim: Prf.cases)

```

```

abbreviation
  Prefixes s  $\stackrel{def}{=} \{s'.\ prefix\ s'\ s\}$ 

```

```

abbreviation
  Suffixes s  $\stackrel{def}{=} \{s'.\ suffix\ s'\ s\}$ 

```

```

abbreviation
  SSuffixes s  $\stackrel{def}{=} \{s'.\ strict\_suffix\ s'\ s\}$ 

```

```

lemma Suffixes_cons [simp]:
  shows Suffixes (c # s) = Suffixes s  $\cup$  {c # s}
by (auto simp add: suffix_def Cons_eq_append_conv)

```

```

lemma finite_Suffixes:
  shows finite (Suffixes s)
by (induct s) (simp_all)

```

**lemma** *finite\_SSuffixes*:

**shows** *finite* (*SSuffixes* *s*)

**proof** –

**have** *SSuffixes* *s*  $\subseteq$  *Suffixes* *s*

**unfolding** *strict\_suffix\_def* *suffix\_def* **by** *auto*

**then show** *finite* (*SSuffixes* *s*)

**using** *finite\_Suffixes* *finite\_subset* **by** *blast*

**qed**

**lemma** *finite\_Prefixes*:

**shows** *finite* (*Prefixes* *s*)

**proof** –

**have** *finite* (*Suffixes* (*rev* *s*))

**by** (*rule* *finite\_Suffixes*)

**then have** *finite* (*rev* ‘ *Suffixes* (*rev* *s*)) **by** *simp*

**moreover**

**have** *rev* ‘ (*Suffixes* (*rev* *s*)) = *Prefixes* *s*

**unfolding** *suffix\_def* *prefix\_def* *image\_def*

**by** (*auto*)(*metis* *rev\_append* *rev\_rev\_ident*) +

**ultimately show** *finite* (*Prefixes* *s*) **by** *simp*

**qed**

**lemma** *LV\_STAR\_finite*:

**assumes**  $\forall s. \text{finite } (LV\ r\ s)$

**shows** *finite* (*LV* (*STAR* *r*) *s*)

**proof**(*induct* *s* *rule*: *length\_induct*)

**fix** *s*::*char* *list*

**assume**  $\forall s'. \text{length } s' < \text{length } s \longrightarrow \text{finite } (LV\ (STAR\ r)\ s')$

**then have** *IH*:  $\forall s' \in SSuffixes\ s. \text{finite } (LV\ (STAR\ r)\ s')$

**by** (*force* *simp* *add*: *strict\_suffix\_def* *suffix\_def*)

**define** *f* **where**  $f \stackrel{def}{=} \lambda(v, vs). Stars\ (v \# vs)$

**define** *S1* **where**  $S1 \stackrel{def}{=} \bigcup s' \in Prefixes\ s. LV\ r\ s'$

**define** *S2* **where**  $S2 \stackrel{def}{=} \bigcup s2 \in SSuffixes\ s. Stars\ -' (LV\ (STAR\ r)\ s2)$

**have** *finite* *S1* **using** *assms*

**unfolding** *S1\_def* **by** (*simp\_all* *add*: *finite\_Prefixes*)

**moreover**

**with** *IH* **have** *finite* *S2* **unfolding** *S2\_def*

**by** (*auto* *simp* *add*: *finite\_SSuffixes* *inj\_on\_def* *finite\_vimageI*)

**ultimately**

**have** *finite* ( $\{Stars\ []\} \cup f\ ' (S1 \times S2)$ ) **by** *simp*

**moreover**

**have** *LV* (*STAR* *r*) *s*  $\subseteq \{Stars\ []\} \cup f\ ' (S1 \times S2)$

**unfolding** *S1\_def* *S2\_def* *f\_def*

**unfolding** *LV\_def* *image\_def* *prefix\_def* *strict\_suffix\_def*



```

apply(auto)
apply(case_tac x)
apply(auto elim: Prf_elims)
apply(erule Prf_elims)
apply(auto)
apply(case_tac vs)
apply(auto intro: Prf.intros)
apply(rule exI)
apply(rule conjI)
apply(rule_tac x=flat a in exI)
apply(rule conjI)
apply(rule_tac x=flats list in exI)
apply(simp)
apply(blast)
apply(simp add: suffix_def)
using Prf.intros(6) by blast
ultimately
show finite (LV (STAR r) s) by (simp add: finite_subset)
qed

```

```

lemma LV_finite:
  shows finite (LV r s)
proof(induct r arbitrary: s)
  case (ZERO s)
  show finite (LV ZERO s) by (simp add: LV_simps)
next
  case (ONE s)
  show finite (LV ONE s) by (simp add: LV_simps)
next
  case (CHAR c s)
  show finite (LV (CHAR c) s) by (simp add: LV_simps)
next
  case (ALT r1 r2 s)
  then show finite (LV (ALT r1 r2) s) by (simp add: LV_simps)
next
  case (SEQ r1 r2 s)
  define f where  $f \stackrel{\text{def}}{=} \lambda(v1, v2). \text{Seq } v1 \ v2$ 
  define S1 where  $S1 \stackrel{\text{def}}{=} \bigcup s' \in \text{Prefixes } s. \text{LV } r1 \ s'$ 
  define S2 where  $S2 \stackrel{\text{def}}{=} \bigcup s' \in \text{Suffixes } s. \text{LV } r2 \ s'$ 
  have IHs:  $\bigwedge s. \text{finite (LV } r1 \ s) \wedge s. \text{finite (LV } r2 \ s)$  by fact+
  then have finite S1 finite S2 unfolding S1_def S2_def
  by (simp_all add: finite_Prefixes finite_Suffixes)
  moreover
  have LV (SEQ r1 r2) s  $\subseteq f' (S1 \times S2)$ 

```

```

unfolding f_def S1_def S2_def
unfolding LV_def image_def prefix_def suffix_def
apply (auto elim!: Prf_elims)
by (metis (mono_tags, lifting) mem_Collect_eq)
ultimately
show finite (LV (SEQ r1 r2) s)
by (simp add: finite_subset)
next
case (STAR r s)
then show finite (LV (STAR r) s) by (simp add: LV_STAR_finite)
qed

```

## 12 Our POSIX Definition

**inductive**

*Posix* :: *string*  $\Rightarrow$  *rexp*  $\Rightarrow$  *val*  $\Rightarrow$  *bool* ( $\_ \in \_ \rightarrow \_ [100, 100, 100] 100$ )

**where**

*Posix\_ONE*:  $\square \in ONE \rightarrow Void$

| *Posix\_CHAR*:  $[c] \in (CHAR\ c) \rightarrow (Char\ c)$

| *Posix\_ALT1*:  $s \in r1 \rightarrow v \Longrightarrow s \in (ALT\ r1\ r2) \rightarrow (Left\ v)$

| *Posix\_ALT2*:  $\llbracket s \in r2 \rightarrow v; s \notin L(r1) \rrbracket \Longrightarrow s \in (ALT\ r1\ r2) \rightarrow (Right\ v)$

| *Posix\_SEQ*:  $\llbracket s1 \in r1 \rightarrow v1; s2 \in r2 \rightarrow v2;$

$\neg(\exists s3\ s4. s3 \neq \square \wedge s3 @ s4 = s2 \wedge (s1 @ s3) \in L\ r1 \wedge s4 \in L\ r2) \rrbracket \Longrightarrow$   
 $(s1 @ s2) \in (SEQ\ r1\ r2) \rightarrow (Seq\ v1\ v2)$

| *Posix\_STAR1*:  $\llbracket s1 \in r \rightarrow v; s2 \in STAR\ r \rightarrow Stars\ vs; flat\ v \neq \square;$

$\neg(\exists s3\ s4. s3 \neq \square \wedge s3 @ s4 = s2 \wedge (s1 @ s3) \in L\ r \wedge s4 \in L\ (STAR\ r)) \rrbracket$   
 $\Longrightarrow (s1 @ s2) \in STAR\ r \rightarrow Stars\ (v \# vs)$

| *Posix\_STAR2*:  $\square \in STAR\ r \rightarrow Stars\ \square$

**inductive-cases** *Posix\_elims*:

$s \in ZERO \rightarrow v$

$s \in ONE \rightarrow v$

$s \in CHAR\ c \rightarrow v$

$s \in ALT\ r1\ r2 \rightarrow v$

$s \in SEQ\ r1\ r2 \rightarrow v$

$s \in STAR\ r \rightarrow v$

**lemma** *PosixI*:

**assumes**  $s \in r \rightarrow v$

**shows**  $s \in L\ r\ flat\ v = s$

**using** *assms*

**by** (*induct s r v rule: Posix.induct*)

(*auto simp add: Sequ\_def*)

Our Posix definition determines a unique value.

**lemma** *Posix\_determ*:

```

assumes  $s \in r \rightarrow v1\ s \in r \rightarrow v2$ 
shows  $v1 = v2$ 
using assms
proof (induct s r v1 arbitrary: v2 rule: Posix.induct)
  case (Posix_ONE v2)
    have  $\square \in ONE \rightarrow v2$  by fact
    then show  $Void = v2$  by cases auto
  next
    case (Posix_CHAR c v2)
    have  $[c] \in CHAR\ c \rightarrow v2$  by fact
    then show  $Char\ c = v2$  by cases auto
  next
    case (Posix_ALT1 s r1 v r2 v2)
    have  $s \in ALT\ r1\ r2 \rightarrow v2$  by fact
    moreover
    have  $s \in r1 \rightarrow v$  by fact
    then have  $s \in L\ r1$  by (simp add: Posix1)
    ultimately obtain  $v'$  where  $eq: v2 = Left\ v'\ s \in r1 \rightarrow v'$  by cases auto
    moreover
    have  $IH: \bigwedge v2. s \in r1 \rightarrow v2 \implies v = v2$  by fact
    ultimately have  $v = v'$  by simp
    then show  $Left\ v = v2$  using eq by simp
  next
    case (Posix_ALT2 s r2 v r1 v2)
    have  $s \in ALT\ r1\ r2 \rightarrow v2$  by fact
    moreover
    have  $s \notin L\ r1$  by fact
    ultimately obtain  $v'$  where  $eq: v2 = Right\ v'\ s \in r2 \rightarrow v'$ 
      by cases (auto simp add: Posix1)
    moreover
    have  $IH: \bigwedge v2. s \in r2 \rightarrow v2 \implies v = v2$  by fact
    ultimately have  $v = v'$  by simp
    then show  $Right\ v = v2$  using eq by simp
  next
    case (Posix_SEQ s1 r1 v1 s2 r2 v2 v')
    have  $(s1\ @\ s2) \in SEQ\ r1\ r2 \rightarrow v'$ 
       $s1 \in r1 \rightarrow v1\ s2 \in r2 \rightarrow v2$ 
       $\neg (\exists s3\ s4. s3 \neq \square \wedge s3\ @\ s4 = s2 \wedge s1\ @\ s3 \in L\ r1 \wedge s4 \in L\ r2)$  by fact+
    then obtain  $v1'\ v2'$  where  $v' = Seq\ v1'\ v2'\ s1 \in r1 \rightarrow v1'\ s2 \in r2 \rightarrow v2'$ 
    apply(cases) apply (auto simp add: append_eq_append_conv2)
    using Posix1(I) by fastforce+
    moreover
    have  $IHs: \bigwedge v1'. s1 \in r1 \rightarrow v1' \implies v1 = v1'$ 
       $\bigwedge v2'. s2 \in r2 \rightarrow v2' \implies v2 = v2'$  by fact+
    ultimately show  $Seq\ v1\ v2 = v'$  by simp

```

```

next
  case (Posix_STAR1 s1 r v s2 vs v2)
  have (s1 @ s2) ∈ STAR r → v2
    s1 ∈ r → v s2 ∈ STAR r → Stars vs flat v ≠ []
    ¬ (∃ s3 s4. s3 ≠ [] ∧ s3 @ s4 = s2 ∧ s1 @ s3 ∈ L r ∧ s4 ∈ L (STAR r)) by fact+
  then obtain v' vs' where v2 = Stars (v' # vs') s1 ∈ r → v' s2 ∈ (STAR r) → (Stars
vs')
  apply(cases) apply (auto simp add: append_eq_append_conv2)
  using Posix1(1) apply fastforce
  apply (metis Posix1(1) Posix_STAR1.hyps(6) append_Nil append_Nil2)
  using Posix1(2) by blast
  moreover
  have IHs: ∧v2. s1 ∈ r → v2 ⇒ v = v2
    ∧v2. s2 ∈ STAR r → v2 ⇒ Stars vs = v2 by fact+
  ultimately show Stars (v # vs) = v2 by auto
next
  case (Posix_STAR2 r v2)
  have [] ∈ STAR r → v2 by fact
  then show Stars [] = v2 by cases (auto simp add: Posix1)
qed

```

Our POSIX values are lexical values.

**lemma** *Posix\_LV*:

```

assumes s ∈ r → v
shows v ∈ LV r s
using assms unfolding LV_def
apply(induct rule: Posix.induct)
apply(auto simp add: intro!: Prf.intros elim!: Prf.elims)
done

```

**lemma** *Posix\_Prf*:

```

assumes s ∈ r → v
shows | = v : r
using assms Posix_LV LV_def
by simp

```

**end**

**theory** *Lexer*

**imports** *Spec*

**begin**

### 13 The Lexer Functions by Sulzmann and Lu

**fun**

```

mkeps :: rexp ⇒ val
where
  mkeps(ONE) = Void
  | mkeps(SEQ r1 r2) = Seq (mkeps r1) (mkeps r2)
  | mkeps(ALT r1 r2) = (if nullable(r1) then Left (mkeps r1) else Right (mkeps r2))
  | mkeps(STAR r) = Stars []

fun injval :: rexp ⇒ char ⇒ val ⇒ val
where
  injval (CHAR d) c Void = Char d
  | injval (ALT r1 r2) c (Left v1) = Left(injval r1 c v1)
  | injval (ALT r1 r2) c (Right v2) = Right(injval r2 c v2)
  | injval (SEQ r1 r2) c (Seq v1 v2) = Seq (injval r1 c v1) v2
  | injval (SEQ r1 r2) c (Left (Seq v1 v2)) = Seq (injval r1 c v1) v2
  | injval (SEQ r1 r2) c (Right v2) = Seq (mkeps r1) (injval r2 c v2)
  | injval (STAR r) c (Seq v (Stars vs)) = Stars ((injval r c v) # vs)

fun
  lexer :: rexp ⇒ string ⇒ val option
where
  lexer r [] = (if nullable r then Some(mkeps r) else None)
  | lexer r (c#s) = (case (lexer (der c r) s) of
    None ⇒ None
    | Some(v) ⇒ Some(injval r c v))

```

## 14 Mkeps, Injval Properties

```

lemma mkeps_nullable:
  assumes nullable(r)
  shows  $\models$  mkeps r : r
using assms
by (induct rule: nullable.induct)
    (auto intro: Prf.intros)

```

```

lemma mkeps_flat:
  assumes nullable(r)
  shows flat (mkeps r) = []
using assms
by (induct rule: nullable.induct) (auto)

```

```

lemma Prf_injval_flat:
  assumes  $\models$  v : der c r
  shows flat (injval r c v) = c # (flat v)
using assms
apply(induct c r arbitrary: v rule: der.induct)

```

```

apply(auto elim!: Prf_elims intro: mkeps_flat split: if_splits)
done

```

```

lemma Prf_injval:
  assumes  $\models v : \text{der } c \ r$ 
  shows  $\models (\text{injval } r \ c \ v) : r$ 
using assms
apply(induct r arbitrary: c v rule: rexp.induct)
apply(auto intro!: Prf.intros mkeps_nullable elim!: Prf_elims split: if_splits)
apply(simp add: Prf_injval_flat)
done

```

Mkeps and injval produce, or preserve, Posix values.

```

lemma Posix_mkeps:
  assumes nullable r
  shows  $\square \in r \rightarrow \text{mkeps } r$ 
using assms
apply(induct r rule: nullable.induct)
apply(auto intro: Posix.intros simp add: nullable_correctness Sequ_def)
apply(subst append.simps(1)[symmetric])
apply(rule Posix.intros)
apply(auto)
done

```

```

lemma Posix_injval:
  assumes  $s \in (\text{der } c \ r) \rightarrow v$ 
  shows  $(c \ \# \ s) \in r \rightarrow (\text{injval } r \ c \ v)$ 
using assms
proof(induct r arbitrary: s v rule: rexp.induct)
  case ZERO
  have  $s \in \text{der } c \ \text{ZERO} \rightarrow v$  by fact
  then have  $s \in \text{ZERO} \rightarrow v$  by simp
  then have False by cases
  then show  $(c \ \# \ s) \in \text{ZERO} \rightarrow (\text{injval } \text{ZERO} \ c \ v)$  by simp
next
  case ONE
  have  $s \in \text{der } c \ \text{ONE} \rightarrow v$  by fact
  then have  $s \in \text{ZERO} \rightarrow v$  by simp
  then have False by cases
  then show  $(c \ \# \ s) \in \text{ONE} \rightarrow (\text{injval } \text{ONE} \ c \ v)$  by simp
next
  case (CHAR d)
  consider  $(\text{eq}) \ c = d \mid (\text{ineq}) \ c \neq d$  by blast
  then show  $(c \ \# \ s) \in (\text{CHAR } d) \rightarrow (\text{injval } (\text{CHAR } d) \ c \ v)$ 
  proof (cases)
    case eq

```

```

have  $s \in \text{der } c \text{ (CHAR } d) \rightarrow v$  by fact
then have  $s \in \text{ONE} \rightarrow v$  using eq by simp
then have eqs:  $s = [] \wedge v = \text{Void}$  by cases simp
show  $(c \# s) \in \text{CHAR } d \rightarrow \text{injval (CHAR } d) c v$  using eq eqs
by (auto intro: Posix.intros)
next
case ineq
have  $s \in \text{der } c \text{ (CHAR } d) \rightarrow v$  by fact
then have  $s \in \text{ZERO} \rightarrow v$  using ineq by simp
then have False by cases
then show  $(c \# s) \in \text{CHAR } d \rightarrow \text{injval (CHAR } d) c v$  by simp
qed
next
case (ALT r1 r2)
have IH1:  $\bigwedge s v. s \in \text{der } c r1 \rightarrow v \implies (c \# s) \in r1 \rightarrow \text{injval } r1 c v$  by fact
have IH2:  $\bigwedge s v. s \in \text{der } c r2 \rightarrow v \implies (c \# s) \in r2 \rightarrow \text{injval } r2 c v$  by fact
have  $s \in \text{der } c \text{ (ALT } r1 r2) \rightarrow v$  by fact
then have  $s \in \text{ALT (der } c r1) \text{ (der } c r2) \rightarrow v$  by simp
then consider (left) v' where  $v = \text{Left } v' s \in \text{der } c r1 \rightarrow v'$ 
| (right) v' where  $v = \text{Right } v' s \notin L \text{ (der } c r1) s \in \text{der } c r2 \rightarrow v'$ 
by cases auto
then show  $(c \# s) \in \text{ALT } r1 r2 \rightarrow \text{injval (ALT } r1 r2) c v$ 
proof (cases)
case left
have  $s \in \text{der } c r1 \rightarrow v'$  by fact
then have  $(c \# s) \in r1 \rightarrow \text{injval } r1 c v'$  using IH1 by simp
then have  $(c \# s) \in \text{ALT } r1 r2 \rightarrow \text{injval (ALT } r1 r2) c (\text{Left } v')$  by (auto intro:
Posix.intros)
then show  $(c \# s) \in \text{ALT } r1 r2 \rightarrow \text{injval (ALT } r1 r2) c v$  using left by simp
next
case right
have  $s \notin L \text{ (der } c r1)$  by fact
then have  $c \# s \notin L r1$  by (simp add: der_correctness Der_def)
moreover
have  $s \in \text{der } c r2 \rightarrow v'$  by fact
then have  $(c \# s) \in r2 \rightarrow \text{injval } r2 c v'$  using IH2 by simp
ultimately have  $(c \# s) \in \text{ALT } r1 r2 \rightarrow \text{injval (ALT } r1 r2) c (\text{Right } v')$ 
by (auto intro: Posix.intros)
then show  $(c \# s) \in \text{ALT } r1 r2 \rightarrow \text{injval (ALT } r1 r2) c v$  using right by simp
qed
next
case (SEQ r1 r2)
have IH1:  $\bigwedge s v. s \in \text{der } c r1 \rightarrow v \implies (c \# s) \in r1 \rightarrow \text{injval } r1 c v$  by fact
have IH2:  $\bigwedge s v. s \in \text{der } c r2 \rightarrow v \implies (c \# s) \in r2 \rightarrow \text{injval } r2 c v$  by fact
have  $s \in \text{der } c \text{ (SEQ } r1 r2) \rightarrow v$  by fact

```

**then consider**

*(left\_nullable)*  $v1\ v2\ s1\ s2$  **where**

$v = \text{Left } (\text{Seq } v1\ v2)$   $s = s1\ @\ s2$

$s1 \in \text{der } c\ r1 \rightarrow v1\ s2 \in r2 \rightarrow v2\ \text{nullable } r1$

$\neg (\exists s3\ s4. s3 \neq [] \wedge s3\ @\ s4 = s2 \wedge s1\ @\ s3 \in L(\text{der } c\ r1) \wedge s4 \in L\ r2)$

| *(right\_nullable)*  $v1\ s1\ s2$  **where**

$v = \text{Right } v1\ s = s1\ @\ s2$

$s \in \text{der } c\ r2 \rightarrow v1\ \text{nullable } r1\ s1\ @\ s2 \notin L(\text{SEQ } (\text{der } c\ r1)\ r2)$

| *(not\_nullable)*  $v1\ v2\ s1\ s2$  **where**

$v = \text{Seq } v1\ v2\ s = s1\ @\ s2$

$s1 \in \text{der } c\ r1 \rightarrow v1\ s2 \in r2 \rightarrow v2\ \neg\text{nullable } r1$

$\neg (\exists s3\ s4. s3 \neq [] \wedge s3\ @\ s4 = s2 \wedge s1\ @\ s3 \in L(\text{der } c\ r1) \wedge s4 \in L\ r2)$

**by** (*force split: if\_splits elim!: Posix\_elims simp add: Sequ\_def der\_correctness*

*Der\_def*)

**then show**  $(c\ \# \ s) \in \text{SEQ } r1\ r2 \rightarrow \text{inval } (\text{SEQ } r1\ r2)\ c\ v$

**proof** (*cases*)

**case** *left\_nullable*

**have**  $s1 \in \text{der } c\ r1 \rightarrow v1$  **by fact**

**then have**  $(c\ \# \ s1) \in r1 \rightarrow \text{inval } r1\ c\ v1$  **using IH1 by simp**

**moreover**

**have**  $\neg (\exists s3\ s4. s3 \neq [] \wedge s3\ @\ s4 = s2 \wedge s1\ @\ s3 \in L(\text{der } c\ r1) \wedge s4 \in L\ r2)$  **by**

*fact*

**then have**  $\neg (\exists s3\ s4. s3 \neq [] \wedge s3\ @\ s4 = s2 \wedge (c\ \# \ s1)\ @\ s3 \in L\ r1 \wedge s4 \in L\ r2)$

**by** (*simp add: der\_correctness Der\_def*)

**ultimately have**  $((c\ \# \ s1)\ @\ s2) \in \text{SEQ } r1\ r2 \rightarrow \text{Seq } (\text{inval } r1\ c\ v1)\ v2$  **using**

*left\_nullable by (rule\_tac Posix.intros)*

**then show**  $(c\ \# \ s) \in \text{SEQ } r1\ r2 \rightarrow \text{inval } (\text{SEQ } r1\ r2)\ c\ v$  **using left\_nullable by**

*simp*

**next**

**case** *right\_nullable*

**have** *nullable r1 by fact*

**then have**  $[] \in r1 \rightarrow (\text{mkeps } r1)$  **by (rule Posix\_mkeps)**

**moreover**

**have**  $s \in \text{der } c\ r2 \rightarrow v1$  **by fact**

**then have**  $(c\ \# \ s) \in r2 \rightarrow (\text{inval } r2\ c\ v1)$  **using IH2 by simp**

**moreover**

**have**  $s1\ @\ s2 \notin L(\text{SEQ } (\text{der } c\ r1)\ r2)$  **by fact**

**then have**  $\neg (\exists s3\ s4. s3 \neq [] \wedge s3\ @\ s4 = c\ \# \ s \wedge []\ @\ s3 \in L\ r1 \wedge s4 \in L\ r2)$

**using** *right\_nullable*

**by**(*auto simp add: der\_correctness Der\_def append\_eq\_Cons\_conv Sequ\_def*)

**ultimately have**  $([]\ @\ (c\ \# \ s)) \in \text{SEQ } r1\ r2 \rightarrow \text{Seq } (\text{mkeps } r1)\ (\text{inval } r2\ c\ v1)$

**by**(*rule Posix.intros*)

**then show**  $(c\ \# \ s) \in \text{SEQ } r1\ r2 \rightarrow \text{inval } (\text{SEQ } r1\ r2)\ c\ v$  **using right\_nullable by**

*simp*

**next**



```

case not_nullable
  have  $s1 \in \text{der } c \ r1 \rightarrow v1$  by fact
  then have  $(c \# s1) \in r1 \rightarrow \text{inval } r1 \ c \ v1$  using IH1 by simp
  moreover
  have  $\neg (\exists s_3 \ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s2 \wedge s1 @ s_3 \in L(\text{der } c \ r1) \wedge s_4 \in L \ r2)$  by
fact
  then have  $\neg (\exists s_3 \ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s2 \wedge (c \# s1) @ s_3 \in L \ r1 \wedge s_4 \in L \ r2)$ 
by (simp add: der_correctness Der_def)
  ultimately have  $((c \# s1) @ s2) \in \text{SEQ } r1 \ r2 \rightarrow \text{Seq}(\text{inval } r1 \ c \ v1) \ v2$  using
not_nullable
  by (rule_tac Posix.intros) (simp_all)
  then show  $(c \# s) \in \text{SEQ } r1 \ r2 \rightarrow \text{inval}(\text{SEQ } r1 \ r2) \ c \ v$  using not_nullable by
simp
  qed
next
case (STAR r)
  have IH:  $\bigwedge s \ v. s \in \text{der } c \ r \rightarrow v \implies (c \# s) \in r \rightarrow \text{inval } r \ c \ v$  by fact
  have  $s \in \text{der } c \ (\text{STAR } r) \rightarrow v$  by fact
  then consider
    (cons)  $v1 \ vs \ s1 \ s2$  where
       $v = \text{Seq } v1 \ (\text{Stars } vs) \ s = s1 @ s2$ 
       $s1 \in \text{der } c \ r \rightarrow v1 \ s2 \in (\text{STAR } r) \rightarrow (\text{Stars } vs)$ 
       $\neg (\exists s_3 \ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s2 \wedge s1 @ s_3 \in L(\text{der } c \ r) \wedge s_4 \in L(\text{STAR } r))$ 
      apply(auto elim!: Posix_elims(1-5) simp add: der_correctness Der_def intro:
Posix.intros)
      apply(rotate_tac 3)
      apply(erule_tac Posix_elims(6))
      apply (simp add: Posix.intros(6))
      using Posix.intros(7) by blast
  then show  $(c \# s) \in \text{STAR } r \rightarrow \text{inval}(\text{STAR } r) \ c \ v$ 
  proof (cases)
  case cons
    have  $s1 \in \text{der } c \ r \rightarrow v1$  by fact
    then have  $(c \# s1) \in r \rightarrow \text{inval } r \ c \ v1$  using IH by simp
    moreover
    have  $s2 \in \text{STAR } r \rightarrow \text{Stars } vs$  by fact
    moreover
    have  $(c \# s1) \in r \rightarrow \text{inval } r \ c \ v1$  by fact
    then have flat  $(\text{inval } r \ c \ v1) = (c \# s1)$  by (rule Posix1)
    then have flat  $(\text{inval } r \ c \ v1) \neq []$  by simp
    moreover
    have  $\neg (\exists s_3 \ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s2 \wedge s1 @ s_3 \in L(\text{der } c \ r) \wedge s_4 \in L(\text{STAR } r))$  by fact
    then have  $\neg (\exists s_3 \ s_4. s_3 \neq [] \wedge s_3 @ s_4 = s2 \wedge (c \# s1) @ s_3 \in L \ r \wedge s_4 \in L(\text{STAR } r))$ 

```

```

    by (simp add: der_correctness Der_def)
  ultimately
  have ((c # s1) @ s2) ∈ STAR r → Stars (inval r c v1 # vs) by (rule Posix.intros)
  then show (c # s) ∈ STAR r → inval (STAR r) c v using cons by (simp)
qed
qed

```

## 15 Lexer Correctness

```

lemma lexer_correct_None:
  shows  $s \notin L r \longleftrightarrow \text{lexer } r s = \text{None}$ 
  apply (induct s arbitrary: r)
  apply (simp)
  apply (simp add: nullable_correctness)
  apply (simp)
  apply (drule_tac x=der a r in meta_spec)
  apply (auto)
  apply (auto simp add: der_correctness Der_def)
done

```

```

lemma lexer_correct_Some:
  shows  $s \in L r \longleftrightarrow (\exists v. \text{lexer } r s = \text{Some}(v) \wedge s \in r \rightarrow v)$ 
  apply (induct s arbitrary: r)
  apply (simp only: lexer.simps)
  apply (simp)
  apply (simp add: nullable_correctness Posix_mkeys)
  apply (drule_tac x=der a r in meta_spec)
  apply (simp (no_asm_use) add: der_correctness Der_def del: lexer.simps)
  apply (simp del: lexer.simps)
  apply (simp only: lexer.simps)
  apply (case_tac lexer (der a r) s = None)
  apply (auto)[1]
  apply (simp)
  apply (erule exE)
  apply (simp)
  apply (rule iffI)
  apply (simp add: Posix_inval)
  apply (simp add: PosixI(1))
done

```

```

lemma lexer_correctness:
  shows  $(\text{lexer } r s = \text{Some } v) \longleftrightarrow s \in r \rightarrow v$ 
  and  $(\text{lexer } r s = \text{None}) \longleftrightarrow \neg(\exists v. s \in r \rightarrow v)$ 
  using PosixI(1) Posix_determ lexer_correct_None lexer_correct_Some apply fastforce
  using PosixI(1) lexer_correct_None lexer_correct_Some by blast

```

```

fun flex :: rexp => (val => val) => string => (val => val)
  where
    flex r f [] = f
  | flex r f (c#s) = flex (der c r) (\v. f (injval r c v)) s

```

```

lemma flex_fun_apply:
  shows g (flex r f s v) = flex r (g o f) s v
  apply(induct s arbitrary: g f r v)
  apply(simp_all add: comp_def)
  by meson

```

```

lemma flex_append:
  shows flex r f (s1 @ s2) = flex (ders s1 r) (flex r f s1) s2
  apply(induct s1 arbitrary: s2 r f)
  apply(simp_all)
  done

```

```

lemma lexer_flex:
  shows lexer r s = (if nullable (ders s r)
                    then Some(flex r id s (mkeps (ders s r))) else None)
  apply(induct s arbitrary: r)
  apply(simp_all add: flex_fun_apply)
  done

```

**unused-thms**

```

end
theory Simplifying
  imports Lexer
begin

```

## 16 Lexer including simplifications

```

fun F_RIGHT where
  F_RIGHT f v = Right (f v)

```

```

fun F_LEFT where
  F_LEFT f v = Left (f v)

```

```

fun F_ALT where
  F_ALT f1 f2 (Right v) = Right (f2 v)

```

```
| F_ALT f1 f2 (Left v) = Left (f1 v)
| F_ALT fl f2 v = v
```

**fun F\_SEQ1 where**

```
F_SEQ1 f1 f2 v = Seq (f1 Void) (f2 v)
```

**fun F\_SEQ2 where**

```
F_SEQ2 f1 f2 v = Seq (f1 v) (f2 Void)
```

**fun F\_SEQ where**

```
F_SEQ f1 f2 (Seq v1 v2) = Seq (f1 v1) (f2 v2)
```

```
| F_SEQ fl f2 v = v
```

**fun simp\_ALT where**

```
simp_ALT (ZERO, f1) (r2, f2) = (r2, F_RIGHT f2)
```

```
| simp_ALT (r1, f1) (ZERO, f2) = (r1, F_LEFT f1)
```

```
| simp_ALT (r1, f1) (r2, f2) = (ALT r1 r2, F_ALT f1 f2)
```

**fun simp\_SEQ where**

```
simp_SEQ (ONE, f1) (r2, f2) = (r2, F_SEQ1 f1 f2)
```

```
| simp_SEQ (r1, f1) (ONE, f2) = (r1, F_SEQ2 f1 f2)
```

```
| simp_SEQ (ZERO, f1) (r2, f2) = (ZERO, undefined)
```

```
| simp_SEQ (r1, f1) (ZERO, f2) = (ZERO, undefined)
```

```
| simp_SEQ (r1, f1) (r2, f2) = (SEQ r1 r2, F_SEQ f1 f2)
```

**lemma simp\_SEQ\_simps[simp]:**

```
simp_SEQ p1 p2 = (if (fst p1 = ONE) then (fst p2, F_SEQ1 (snd p1) (snd p2))
```

```
  else (if (fst p2 = ONE) then (fst p1, F_SEQ2 (snd p1) (snd p2))
```

```
  else (if (fst p1 = ZERO) then (ZERO, undefined)
```

```
  else (if (fst p2 = ZERO) then (ZERO, undefined)
```

```
  else (SEQ (fst p1) (fst p2), F_SEQ (snd p1) (snd p2))))))
```

**by (induct p1 p2 rule: simp\_SEQ.induct) (auto)**

**lemma simp\_ALT\_simps[simp]:**

```
simp_ALT p1 p2 = (if (fst p1 = ZERO) then (fst p2, F_RIGHT (snd p2))
```

```
  else (if (fst p2 = ZERO) then (fst p1, F_LEFT (snd p1))
```

```
  else (ALT (fst p1) (fst p2), F_ALT (snd p1) (snd p2))))
```

**by (induct p1 p2 rule: simp\_ALT.induct) (auto)**

**fun**

```
simp :: rexp ⇒ rexp * (val ⇒ val)
```

**where**

```
simp (ALT r1 r2) = simp_ALT (simp r1) (simp r2)
```

```
| simp (SEQ r1 r2) = simp_SEQ (simp r1) (simp r2)
| simp r = (r, id)
```

**fun**

```
slexer :: rexp ⇒ string ⇒ val option
```

**where**

```
slexer r [] = (if nullable r then Some(mkeys r) else None)
| slexer r (c#s) = (let (rs, fr) = simp (der c r) in
  (case (slexer rs s) of
    None ⇒ None
  | Some(v) ⇒ Some(injval r c (fr v))))
```

**lemma** *slexer\_better\_simp*:

```
slexer r (c#s) = (case (slexer (fst (simp (der c r))) s) of
  None ⇒ None
  | Some(v) ⇒ Some(injval r c ((snd (simp (der c r))) v)))
```

**by** (*auto split: prod.split option.split*)

**lemma** *L\_fst\_simp*:

```
shows L(r) = L(fst (simp r))
```

**by** (*induct r*) (*auto*)

**lemma** *Posix\_simp*:

```
assumes s ∈ (fst (simp r)) → v
```

```
shows s ∈ r → ((snd (simp r)) v)
```

**using** *assms*

**proof**(*induct r arbitrary: s v rule: rexp.induct*)

```
case (ALT r1 r2 s v)
```

```
have IH1: ∧s v. s ∈ fst (simp r1) → v ⇒ s ∈ r1 → snd (simp r1) v by fact
```

```
have IH2: ∧s v. s ∈ fst (simp r2) → v ⇒ s ∈ r2 → snd (simp r2) v by fact
```

```
have as: s ∈ fst (simp (ALT r1 r2)) → v by fact
```

```
consider (ZERO_ZERO) fst (simp r1) = ZERO fst (simp r2) = ZERO
```

```
| (ZERO_NZERO) fst (simp r1) = ZERO fst (simp r2) ≠ ZERO
```

```
| (NZERO_ZERO) fst (simp r1) ≠ ZERO fst (simp r2) = ZERO
```

```
| (NZERO_NZERO) fst (simp r1) ≠ ZERO fst (simp r2) ≠ ZERO by auto
```

```
then show s ∈ ALT r1 r2 → snd (simp (ALT r1 r2)) v
```

```
proof(cases)
```

```
case (ZERO_ZERO)
```

```
with as have s ∈ ZERO → v by simp
```

```
then show s ∈ ALT r1 r2 → snd (simp (ALT r1 r2)) v by (rule Posix_elims(1))
```

```
next
```

```
case (ZERO_NZERO)
```

```
with as have s ∈ fst (simp r2) → v by simp
```

```

with IH2 have  $s \in r2 \rightarrow \text{snd}(\text{simp } r2) v$  by simp
moreover
from ZERO_NZERO have  $\text{fst}(\text{simp } r1) = \text{ZERO}$  by simp
then have  $L(\text{fst}(\text{simp } r1)) = \{\}$  by simp
then have  $L r1 = \{\}$  using L_fst_simp by simp
then have  $s \notin L r1$  by simp
ultimately have  $s \in \text{ALT } r1 r2 \rightarrow \text{Right}(\text{snd}(\text{simp } r2) v)$  by (rule Posix_ALT2)
then show  $s \in \text{ALT } r1 r2 \rightarrow \text{snd}(\text{simp}(\text{ALT } r1 r2)) v$ 
using ZERO_NZERO by simp
next
case (NZERO_ZERO)
with as have  $s \in \text{fst}(\text{simp } r1) \rightarrow v$  by simp
with IH1 have  $s \in r1 \rightarrow \text{snd}(\text{simp } r1) v$  by simp
then have  $s \in \text{ALT } r1 r2 \rightarrow \text{Left}(\text{snd}(\text{simp } r1) v)$  by (rule Posix_ALT1)
then show  $s \in \text{ALT } r1 r2 \rightarrow \text{snd}(\text{simp}(\text{ALT } r1 r2)) v$  using NZERO_ZERO by
simp
next
case (NZERO_NZERO)
with as have  $s \in \text{ALT}(\text{fst}(\text{simp } r1))(\text{fst}(\text{simp } r2)) \rightarrow v$  by simp
then consider  $(\text{Left}) v1$  where  $v = \text{Left } v1 s \in (\text{fst}(\text{simp } r1)) \rightarrow v1$ 
|  $(\text{Right}) v2$  where  $v = \text{Right } v2 s \in (\text{fst}(\text{simp } r2)) \rightarrow v2 s \notin L(\text{simp}$ 
 $r1)$ 
by (erule_tac Posix_elims(4))
then show  $s \in \text{ALT } r1 r2 \rightarrow \text{snd}(\text{simp}(\text{ALT } r1 r2)) v$ 
proof(cases)
case (Left)
then have  $v = \text{Left } v1 s \in r1 \rightarrow (\text{snd}(\text{simp } r1) v1)$  using IH1 by simp_all
then show  $s \in \text{ALT } r1 r2 \rightarrow \text{snd}(\text{simp}(\text{ALT } r1 r2)) v$  using NZERO_NZERO
by (simp_all add: Posix_ALT1)
next
case (Right)
then have  $v = \text{Right } v2 s \in r2 \rightarrow (\text{snd}(\text{simp } r2) v2) s \notin L r1$  using IH2 L_fst_simp
by simp_all
then show  $s \in \text{ALT } r1 r2 \rightarrow \text{snd}(\text{simp}(\text{ALT } r1 r2)) v$  using NZERO_NZERO
by (simp_all add: Posix_ALT2)
qed
qed
next
case (SEQ r1 r2 s v)
have IH1:  $\bigwedge s v. s \in \text{fst}(\text{simp } r1) \rightarrow v \implies s \in r1 \rightarrow \text{snd}(\text{simp } r1) v$  by fact
have IH2:  $\bigwedge s v. s \in \text{fst}(\text{simp } r2) \rightarrow v \implies s \in r2 \rightarrow \text{snd}(\text{simp } r2) v$  by fact
have as:  $s \in \text{fst}(\text{simp}(\text{SEQ } r1 r2)) \rightarrow v$  by fact
consider  $(\text{ONE\_ONE}) \text{fst}(\text{simp } r1) = \text{ONE} \text{fst}(\text{simp } r2) = \text{ONE}$ 
|  $(\text{ONE\_NONE}) \text{fst}(\text{simp } r1) = \text{ONE} \text{fst}(\text{simp } r2) \neq \text{ONE}$ 
|  $(\text{NONE\_ONE}) \text{fst}(\text{simp } r1) \neq \text{ONE} \text{fst}(\text{simp } r2) = \text{ONE}$ 

```

```

    | (NONE_NONE) fst (simp r1) ≠ ONE fst (simp r2) ≠ ONE
    by auto
then show s ∈ SEQ r1 r2 → snd (simp (SEQ r1 r2)) v
proof(cases)
  case (ONE_ONE)
  with as have b: s ∈ ONE → v by simp
  from b have s ∈ r1 → snd (simp r1) v using IH1 ONE_ONE by simp
  moreover
  from b have c: s = [] v = Void using Posix_elims(2) by auto
  moreover
  have [] ∈ ONE → Void by (simp add: Posix_ONE)
  then have [] ∈ fst (simp r2) → Void using ONE_ONE by simp
  then have [] ∈ r2 → snd (simp r2) Void using IH2 by simp
  ultimately have ([] @ []) ∈ SEQ r1 r2 → Seq (snd (simp r1) Void) (snd (simp r2)
Void)
    using Posix_SEQ by blast
  then show s ∈ SEQ r1 r2 → snd (simp (SEQ r1 r2)) v using c ONE_ONE by simp
next
case (ONE_NONE)
with as have b: s ∈ fst (simp r2) → v by simp
from b have s ∈ r2 → snd (simp r2) v using IH2 ONE_NONE by simp
moreover
have [] ∈ ONE → Void by (simp add: Posix_ONE)
then have [] ∈ fst (simp r1) → Void using ONE_NONE by simp
then have [] ∈ r1 → snd (simp r1) Void using IH1 by simp
moreover
from ONE_NONE(1) have L (fst (simp r1)) = {[]} by simp
then have L r1 = {[]} by (simp add: L_fst_simp[symmetric])
ultimately have ([] @ s) ∈ SEQ r1 r2 → Seq (snd (simp r1) Void) (snd (simp r2)
v)
  by(rule_tac Posix_SEQ) auto
then show s ∈ SEQ r1 r2 → snd (simp (SEQ r1 r2)) v using ONE_NONE by simp
next
case (NONE_ONE)
with as have s ∈ fst (simp r1) → v by simp
with IH1 have s ∈ r1 → snd (simp r1) v by simp
moreover
have [] ∈ ONE → Void by (simp add: Posix_ONE)
then have [] ∈ fst (simp r2) → Void using NONE_ONE by simp
then have [] ∈ r2 → snd (simp r2) Void using IH2 by simp
ultimately have (s @ []) ∈ SEQ r1 r2 → Seq (snd (simp r1) v) (snd (simp r2)
Void)
  by(rule_tac Posix_SEQ) auto
then show s ∈ SEQ r1 r2 → snd (simp (SEQ r1 r2)) v using NONE_ONE by simp
next

```

```

case (NONE_NONE)
from as have 00: fst (simp r1) ≠ ZERO fst (simp r2) ≠ ZERO
  apply(auto)
  apply(smt Posix_elims(1) fst_conv)
  by (smt NONE_NONE(2) Posix_elims(1) fst1)
with NONE_NONE as have s ∈ SEQ (fst (simp r1)) (fst (simp r2)) → v by simp
then obtain s1 s2 v1 v2 where eqs: s = s1 @ s2 v = Seq v1 v2
  s1 ∈ (fst (simp r1)) → v1 s2 ∈ (fst (simp r2)) → v2
   $\neg (\exists s_3 s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L\ r1 \wedge s_4 \in L\ r2)$ 
  by (erule_tac Posix_elims(5)) (auto simp add: L_fst_simp[symmetric])
then have s1 ∈ r1 → (snd (simp r1) v1) s2 ∈ r2 → (snd (simp r2) v2)
  using IH1 IH2 by auto
then show s ∈ SEQ r1 r2 → snd (simp (SEQ r1 r2)) v using eqs NONE_NONE
00
  by(auto intro: Posix_SEQ)
qed
qed (simp_all)

```

```

lemma slexer_correctness:
shows slexer r s = lexer r s
proof(induct s arbitrary: r)
case Nil
show slexer r [] = lexer r [] by simp
next
case (Cons c s r)
have IH:  $\bigwedge r. slexer\ r\ s = lexer\ r\ s$  by fact
show slexer r (c # s) = lexer r (c # s)
proof (cases s ∈ L (der c r))
case True
  assume a1: s ∈ L (der c r)
  then obtain v1 where a2: lexer (der c r) s = Some v1 s ∈ der c r → v1
  using lexer_correct_Some by auto
  from a1 have s ∈ L (fst (simp (der c r))) using L_fst_simp[symmetric] by simp
  then obtain v2 where a3: lexer (fst (simp (der c r))) s = Some v2 s ∈ (fst (simp (der c r))) → v2
  using lexer_correct_Some by auto
  then have a4: slexer (fst (simp (der c r))) s = Some v2 using IH by simp
  from a3(2) have s ∈ der c r → (snd (simp (der c r))) v2 using Posix_simp by
simp
  with a2(2) have v1 = (snd (simp (der c r))) v2 using Posix_determ by simp
  with a2(1) a4 show slexer r (c # s) = lexer r (c # s) by (auto split: prod.split)
next
case False
  assume b1: s ∉ L (der c r)

```



```

then have lexer (der c r) s = None using lexer_correct_None by simp
moreover
from b1 have  $s \notin L (fst (simp (der c r)))$  using L_fst_simp[symmetric] by simp
then have lexer (fst (simp (der c r))) s = None using lexer_correct_None by simp
then have slexer (fst (simp (der c r))) s = None using IH by simp
ultimately show slexer r (c # s) = lexer r (c # s)
by (simp del: slexer.simps add: slexer_better_simp)

```

**qed**

**qed**

**end**

```

theory Positions
imports Spec Lexer
begin

```

## 17 Positions in Values

**fun**

*at* :: *val*  $\Rightarrow$  *nat list*  $\Rightarrow$  *val*

**where**

```

at v [] = v
| at (Left v) (0#ps) = at v ps
| at (Right v) (Suc 0#ps) = at v ps
| at (Seq v1 v2) (0#ps) = at v1 ps
| at (Seq v1 v2) (Suc 0#ps) = at v2 ps
| at (Stars vs) (n#ps) = at (nth vs n) ps

```

**fun** *Pos* :: *val*  $\Rightarrow$  (*nat list*) *set*

**where**

```

Pos (Void) = {}
| Pos (Char c) = {}
| Pos (Left v) = {}  $\cup$  {0#ps | ps. ps  $\in$  Pos v}
| Pos (Right v) = {}  $\cup$  {1#ps | ps. ps  $\in$  Pos v}
| Pos (Seq v1 v2) = {}  $\cup$  {0#ps | ps. ps  $\in$  Pos v1}  $\cup$  {1#ps | ps. ps  $\in$  Pos v2}
| Pos (Stars []) = {}
| Pos (Stars (v#vs)) = {}  $\cup$  {0#ps | ps. ps  $\in$  Pos v}  $\cup$  {Suc n#ps | n ps. n#ps  $\in$  Pos (Stars vs)}

```

**lemma** *Pos\_stars*:

*Pos* (*Stars* vs) = {}  $\cup$  ( $\bigcup n < \text{length vs. } \{n\#ps \mid ps. ps \in \text{Pos } (vs ! n)\}$ )

**apply**(*induct* vs)

**apply**(*auto simp add: insert\_ident less\_Suc\_eq\_0\_disj*)  
**done**

**lemma** *Pos\_empty*:  
**shows**  $\square \in \text{Pos } v$   
**by** (*induct v rule: Pos.induct*)(*auto*)

**abbreviation**  
 $\text{intlen } vs \stackrel{\text{def}}{=} \text{int } (\text{length } vs)$

**definition** *pflat\_len* ::  $\text{val} \Rightarrow \text{nat list} \Rightarrow \text{int}$   
**where**  
 $\text{pflat\_len } v \ p \stackrel{\text{def}}{=} (\text{if } p \in \text{Pos } v \text{ then intlen } (\text{flat } (\text{at } v \ p)) \text{ else } -1)$

**lemma** *pflat\_len\_simps*:  
**shows**  $\text{pflat\_len } (\text{Seq } v1 \ v2) \ (0\#p) = \text{pflat\_len } v1 \ p$   
**and**  $\text{pflat\_len } (\text{Seq } v1 \ v2) \ (\text{Suc } 0\#p) = \text{pflat\_len } v2 \ p$   
**and**  $\text{pflat\_len } (\text{Left } v) \ (0\#p) = \text{pflat\_len } v \ p$   
**and**  $\text{pflat\_len } (\text{Left } v) \ (\text{Suc } 0\#p) = -1$   
**and**  $\text{pflat\_len } (\text{Right } v) \ (\text{Suc } 0\#p) = \text{pflat\_len } v \ p$   
**and**  $\text{pflat\_len } (\text{Right } v) \ (0\#p) = -1$   
**and**  $\text{pflat\_len } (\text{Stars } (v\#vs)) \ (\text{Suc } n\#p) = \text{pflat\_len } (\text{Stars } vs) \ (n\#p)$   
**and**  $\text{pflat\_len } (\text{Stars } (v\#vs)) \ (0\#p) = \text{pflat\_len } v \ p$   
**and**  $\text{pflat\_len } v \ \square = \text{intlen } (\text{flat } v)$   
**by** (*auto simp add: pflat\_len\_def Pos\_empty*)

**lemma** *pflat\_len\_Stars\_simps*:  
**assumes**  $n < \text{length } vs$   
**shows**  $\text{pflat\_len } (\text{Stars } vs) \ (n\#p) = \text{pflat\_len } (vs!n) \ p$   
**using** *assms*  
**apply**(*induct vs arbitrary: n p*)  
**apply**(*auto simp add: less\_Suc\_eq\_0\_disj pflat\_len\_simps*)  
**done**

**lemma** *pflat\_len\_outside*:  
**assumes**  $p \notin \text{Pos } v1$   
**shows**  $\text{pflat\_len } v1 \ p = -1$   
**using** *assms* **by** (*simp add: pflat\_len\_def*)

## 18 Orderings

**definition** *prefix\_list*::  $'a \text{ list} \Rightarrow 'a \text{ list} \Rightarrow \text{bool}$  ( $-\ \sqsubseteq_{\text{pre}} -$  [60,59] 60)  
**where**

$$ps1 \sqsubseteq_{pre} ps2 \stackrel{def}{=} \exists ps'. ps1 @ps' = ps2$$

**definition** *sprefix\_list*:: 'a list  $\Rightarrow$  'a list  $\Rightarrow$  bool ( $\_ \sqsubseteq_{spre} \_ [60,59] 60$ )

**where**

$$ps1 \sqsubseteq_{spre} ps2 \stackrel{def}{=} ps1 \sqsubseteq_{pre} ps2 \wedge ps1 \neq ps2$$

**inductive** *lex\_list* :: nat list  $\Rightarrow$  nat list  $\Rightarrow$  bool ( $\_ \sqsubseteq_{lex} \_ [60,59] 60$ )

**where**

$$\square \sqsubseteq_{lex} (p\#ps)$$

$$| ps1 \sqsubseteq_{lex} ps2 \implies (p\#ps1) \sqsubseteq_{lex} (p\#ps2)$$

$$| p1 < p2 \implies (p1\#ps1) \sqsubseteq_{lex} (p2\#ps2)$$

**lemma** *lex\_irrfl*:

**fixes** *ps1 ps2* :: nat list

**assumes**  $ps1 \sqsubseteq_{lex} ps2$

**shows**  $ps1 \neq ps2$

**using** *assms*

**by** (*induct rule: lex\_list.induct*)(*auto*)

**lemma** *lex\_simps* [*simp*]:

**fixes** *xs ys* :: nat list

**shows**  $\square \sqsubseteq_{lex} ys \longleftrightarrow ys \neq \square$

**and**  $xs \sqsubseteq_{lex} \square \longleftrightarrow False$

**and**  $(x \# xs) \sqsubseteq_{lex} (y \# ys) \longleftrightarrow (x < y \vee (x = y \wedge xs \sqsubseteq_{lex} ys))$

**by** (*auto simp add: neq\_Nil\_conv elim: lex\_list.cases intro: lex\_list.intros*)

**lemma** *lex\_trans*:

**fixes** *ps1 ps2 ps3* :: nat list

**assumes**  $ps1 \sqsubseteq_{lex} ps2$   $ps2 \sqsubseteq_{lex} ps3$

**shows**  $ps1 \sqsubseteq_{lex} ps3$

**using** *assms*

**by** (*induct arbitrary: ps3 rule: lex\_list.induct*)

(*auto elim: lex\_list.cases*)

**lemma** *lex\_trichotomous*:

**fixes** *p q* :: nat list

**shows**  $p = q \vee p \sqsubseteq_{lex} q \vee q \sqsubseteq_{lex} p$

**apply**(*induct p arbitrary: q*)

**apply**(*auto elim: lex\_list.cases*)

**apply**(*case\_tac q*)

**apply**(*auto*)

**done**

## 19 POSIX Ordering of Values According to Okui & Suzuki

**definition** *PosOrd*::  $val \Rightarrow nat\ list \Rightarrow val \Rightarrow bool$  ( $-\sqsubset_{val} -$  [60, 60, 59] 60)

**where**

$$v1 \sqsubset_{val} p\ v2 \stackrel{def}{=} pflat\_len\ v1\ p > pflat\_len\ v2\ p \wedge \\ (\forall q \in Pos\ v1 \cup Pos\ v2. q \sqsubset_{lex} p \longrightarrow pflat\_len\ v1\ q = pflat\_len\ v2\ q)$$

**lemma** *PosOrd\_def2*:

**shows**  $v1 \sqsubset_{val} p\ v2 \longleftrightarrow$

$$pflat\_len\ v1\ p > pflat\_len\ v2\ p \wedge \\ (\forall q \in Pos\ v1. q \sqsubset_{lex} p \longrightarrow pflat\_len\ v1\ q = pflat\_len\ v2\ q) \wedge \\ (\forall q \in Pos\ v2. q \sqsubset_{lex} p \longrightarrow pflat\_len\ v1\ q = pflat\_len\ v2\ q)$$

**unfolding** *PosOrd\_def*

**apply**(*auto*)

**done**

**definition** *PosOrd\_ex*::  $val \Rightarrow val \Rightarrow bool$  ( $-\sqsubseteq_{val} -$  [60, 59] 60)

**where**

$$v1 : \sqsubseteq_{val} v2 \stackrel{def}{=} \exists p. v1 \sqsubset_{val} p\ v2$$

**definition** *PosOrd\_ex\_eq*::  $val \Rightarrow val \Rightarrow bool$  ( $-\sqsubseteq_{val} -$  [60, 59] 60)

**where**

$$v1 : \sqsubseteq_{val} v2 \stackrel{def}{=} v1 : \sqsubset_{val} v2 \vee v1 = v2$$

**lemma** *PosOrd\_trans*:

**assumes**  $v1 : \sqsubset_{val} v2\ v2 : \sqsubset_{val} v3$

**shows**  $v1 : \sqsubset_{val} v3$

**proof** –

**from** *assms* **obtain**  $p\ p'$

**where** *as*:  $v1 \sqsubset_{val} p\ v2\ v2 \sqsubset_{val} p'\ v3$  **unfolding** *PosOrd\_ex\_def* **by** *blast*

**then** **have** *pos*:  $p \in Pos\ v1\ p' \in Pos\ v2$  **unfolding** *PosOrd\_def* *pflat\_len\_def*

**by** (*smt not\_int\_zless\_negative*)+

**have**  $p = p' \vee p \sqsubset_{lex} p' \vee p' \sqsubset_{lex} p$

**by** (*rule lex\_trichotomous*)

**moreover**

{ **assume**  $p = p'$

**with** *as* **have**  $v1 \sqsubset_{val} p\ v3$  **unfolding** *PosOrd\_def* *pflat\_len\_def*

**by** (*smt Un\_iff*)

**then** **have**  $v1 : \sqsubset_{val} v3$  **unfolding** *PosOrd\_ex\_def* **by** *blast*

}

**moreover**

{ **assume**  $p \sqsubset_{lex} p'$

**with** *as* **have**  $v1 \sqsubset_{val} p\ v3$  **unfolding** *PosOrd\_def* *pflat\_len\_def*

```

    by (smt Un_iff lex_trans)
    then have v1 : $\sqsubseteq$ val v3 unfolding PosOrd_ex_def by blast
  }
  moreover
  { assume p'  $\sqsubseteq$  lex p
    with as have v1  $\sqsubseteq$ val p' v3 unfolding PosOrd_def
    by (smt Un_iff lex_trans pflat_len_def)
    then have v1 : $\sqsubseteq$ val v3 unfolding PosOrd_ex_def by blast
  }
  ultimately show v1 : $\sqsubseteq$ val v3 by blast
qed

```

```

lemma PosOrd_irrefl:
  assumes v : $\sqsubseteq$ val v
  shows False
using assms unfolding PosOrd_ex_def PosOrd_def
by auto

```

```

lemma PosOrd_assym:
  assumes v1 : $\sqsubseteq$ val v2
  shows  $\neg$ (v2 : $\sqsubseteq$ val v1)
using assms
using PosOrd_irrefl PosOrd_trans by blast

```

```

lemma PosOrd_ordering:
  shows ordering ( $\lambda$ v1 v2. v1 : $\sqsubseteq$ val v2) ( $\lambda$  v1 v2. v1 : $\sqsubseteq$ val v2)
  unfolding ordering_def PosOrd_ex_eq_def
  apply(auto)
  using PosOrd_irrefl apply blast
  using PosOrd_assym apply blast
  using PosOrd_trans by blast

```

```

lemma PosOrd_order:
  shows class.order ( $\lambda$ v1 v2. v1 : $\sqsubseteq$ val v2) ( $\lambda$  v1 v2. v1 : $\sqsubseteq$ val v2)
  using PosOrd_ordering
  apply(simp add: class.order_def class.preorder_def class.order_axioms_def)
  unfolding ordering_def
  by blast

```

```

lemma PosOrd_ex_eq2:
  shows v1 : $\sqsubseteq$ val v2  $\longleftrightarrow$  (v1 : $\sqsubseteq$ val v2  $\wedge$  v1  $\neq$  v2)
  using PosOrd_ordering

```

**unfolding** *ordering\_def*  
**by** *auto*

**lemma** *PosOrdeq\_trans*:  
**assumes**  $v1 : \sqsubseteq_{val} v2$   $v2 : \sqsubseteq_{val} v3$   
**shows**  $v1 : \sqsubseteq_{val} v3$   
**using** *assms PosOrd\_ordering*  
**unfolding** *ordering\_def*  
**by** *blast*

**lemma** *PosOrdeq\_antisym*:  
**assumes**  $v1 : \sqsubseteq_{val} v2$   $v2 : \sqsubseteq_{val} v1$   
**shows**  $v1 = v2$   
**using** *assms PosOrd\_ordering*  
**unfolding** *ordering\_def*  
**by** *blast*

**lemma** *PosOrdeq\_refl*:  
**shows**  $v : \sqsubseteq_{val} v$   
**unfolding** *PosOrd\_ex\_eq\_def*  
**by** *auto*

**lemma** *PosOrd\_shorterE*:  
**assumes**  $v1 : \sqsubseteq_{val} v2$   
**shows**  $length (flat v2) \leq length (flat v1)$   
**using** *assms unfolding PosOrd\_ex\_def PosOrd\_def*  
**apply** (*auto*)  
**apply** (*case\_tac p*)  
**apply** (*simp add: pflat\_len\_simps*)  
**apply** (*drule\_tac x= [] in bspec*)  
**apply** (*simp add: Pos\_empty*)  
**apply** (*simp add: pflat\_len\_simps*)  
**done**

**lemma** *PosOrd\_shorterI*:  
**assumes**  $length (flat v2) < length (flat v1)$   
**shows**  $v1 : \sqsubseteq_{val} v2$   
**unfolding** *PosOrd\_ex\_def PosOrd\_def pflat\_len\_def*  
**using** *assms Pos\_empty* **by** *force*

**lemma** *PosOrd\_spreI*:  
**assumes**  $flat v' \sqsubseteq_{spre} flat v$   
**shows**  $v : \sqsubseteq_{val} v'$   
**using** *assms*

```

apply(rule_tac PosOrd_shorter1)
unfolding prefix_list_def sprefix_list_def
by (metis append_Nil2 append_eq_conv_conj drop_all le_less_linear)

```

```

lemma pflat_len_inside:
  assumes pflat_len v2 p < pflat_len v1 p
  shows  $p \in \text{Pos } v1$ 
using assms
unfolding pflat_len_def
by (auto split: if_splits)

```

```

lemma PosOrd_Left_Right:
  assumes  $\text{flat } v1 = \text{flat } v2$ 
  shows  $\text{Left } v1 : \sqsubseteq \text{val } \text{Right } v2$ 
unfolding PosOrd_ex_def
apply(rule_tac x=[0] in ex1)
apply(auto simp add: PosOrd_def pflat_len_simps assms)
done

```

```

lemma PosOrd_LeftE:
  assumes  $\text{Left } v1 : \sqsubseteq \text{val } \text{Left } v2 \text{ flat } v1 = \text{flat } v2$ 
  shows  $v1 : \sqsubseteq \text{val } v2$ 
using assms
unfolding PosOrd_ex_def PosOrd_def2
apply(auto simp add: pflat_len_simps)
apply(frule pflat_len_inside)
apply(auto simp add: pflat_len_simps)
by (metis lex_simps(3) pflat_len_simps(3))

```

```

lemma PosOrd_LeftI:
  assumes  $v1 : \sqsubseteq \text{val } v2 \text{ flat } v1 = \text{flat } v2$ 
  shows  $\text{Left } v1 : \sqsubseteq \text{val } \text{Left } v2$ 
using assms
unfolding PosOrd_ex_def PosOrd_def2
apply(auto simp add: pflat_len_simps)
by (metis less_numeral_extra(3) lex_simps(3) pflat_len_simps(3))

```

```

lemma PosOrd_Left_eq:
  assumes  $\text{flat } v1 = \text{flat } v2$ 
  shows  $\text{Left } v1 : \sqsubseteq \text{val } \text{Left } v2 \longleftrightarrow v1 : \sqsubseteq \text{val } v2$ 
using assms PosOrd_LeftE PosOrd_LeftI
by blast

```

**lemma** *PosOrd\_RightE*:  
**assumes** *Right v1 : $\sqsubset$ val Right v2 flat v1 = flat v2*  
**shows** *v1 : $\sqsubset$ val v2*  
**using** *assms*  
**unfolding** *PosOrd\_ex\_def PosOrd\_def2*  
**apply**(*auto simp add: pflat\_len\_simps*)  
**apply**(*frule pflat\_len\_inside*)  
**apply**(*auto simp add: pflat\_len\_simps*)  
**by** (*metis lex\_simps(3) pflat\_len\_simps(5)*)

**lemma** *PosOrd\_RightI*:  
**assumes** *v1 : $\sqsubset$ val v2 flat v1 = flat v2*  
**shows** *Right v1 : $\sqsubset$ val Right v2*  
**using** *assms*  
**unfolding** *PosOrd\_ex\_def PosOrd\_def2*  
**apply**(*auto simp add: pflat\_len\_simps*)  
**by** (*metis lex\_simps(3) nat\_neq\_iff pflat\_len\_simps(5)*)

**lemma** *PosOrd\_Right\_eq*:  
**assumes** *flat v1 = flat v2*  
**shows** *Right v1 : $\sqsubset$ val Right v2  $\longleftrightarrow$  v1 : $\sqsubset$ val v2*  
**using** *assms PosOrd\_RightE PosOrd\_RightI*  
**by** *blast*

**lemma** *PosOrd\_SeqI1*:  
**assumes** *v1 : $\sqsubset$ val w1 flat (Seq v1 v2) = flat (Seq w1 w2)*  
**shows** *Seq v1 v2 : $\sqsubset$ val Seq w1 w2*  
**using** *assms(1)*  
**apply**(*subst (asm) PosOrd\_ex\_def*)  
**apply**(*subst (asm) PosOrd\_def*)  
**apply**(*clarify*)  
**apply**(*subst PosOrd\_ex\_def*)  
**apply**(*rule\_tac x=0#p in exI*)  
**apply**(*subst PosOrd\_def*)  
**apply**(*rule conjI*)  
**apply**(*simp add: pflat\_len\_simps*)  
**apply**(*rule ballI*)  
**apply**(*rule impI*)  
**apply**(*simp only: Pos\_simps*)  
**apply**(*auto*)[*I*]  
**apply**(*simp add: pflat\_len\_simps*)  
**apply**(*auto simp add: pflat\_len\_simps*)  
**using** *assms(2)*



```

apply(simp)
apply(metis length_append of_nat_add)
done

```

```

lemma PosOrd_SeqI2:
  assumes v2 : $\square$ val w2 flat v2 = flat w2
  shows Seq v v2 : $\square$ val Seq v w2
using assms(1)
apply(subst (asm) PosOrd_ex_def)
apply(subst (asm) PosOrd_def)
apply(clarify)
apply(subst PosOrd_ex_def)
apply(rule_tac x=Suc 0#p in exI)
apply(subst PosOrd_def)
apply(rule conjI)
apply(simp add: pflat_len_simps)
apply(rule ballI)
apply(rule impI)
apply(simp only: Pos_simps)
apply(auto)[1]
apply(simp add: pflat_len_simps)
using assms(2)
apply(simp)
apply(auto simp add: pflat_len_simps)
done

```

```

lemma PosOrd_Seq_eq:
  assumes flat v2 = flat w2
  shows (Seq v v2) : $\square$ val (Seq v w2)  $\longleftrightarrow$  v2 : $\square$ val w2
using assms
apply(auto)
prefer 2
apply(simp add: PosOrd_SeqI2)
apply(simp add: PosOrd_ex_def)
apply(auto)
apply(case_tac p)
apply(simp add: PosOrd_def pflat_len_simps)
apply(case_tac a)
apply(simp add: PosOrd_def pflat_len_simps)
apply(clarify)
apply(case_tac nat)
prefer 2
apply(simp add: PosOrd_def pflat_len_simps pflat_len_outside)
apply(rule_tac x=list in exI)
apply(auto simp add: PosOrd_def2 pflat_len_simps)

```

```

apply(smt Collect_disj_eq lex_list.intros(2) mem_Collect_eq pflat_len_simps(2))
apply(smt Collect_disj_eq lex_list.intros(2) mem_Collect_eq pflat_len_simps(2))
done

```

**lemma** *PosOrd\_StarsI*:

```

  assumes v1 : $\square$ val v2 flats (v1#vs1) = flats (v2#vs2)
  shows Stars (v1#vs1) : $\square$ val Stars (v2#vs2)
using assms(1)
apply(subst (asm) PosOrd_ex_def)
apply(subst (asm) PosOrd_def)
apply(clarify)
apply(subst PosOrd_ex_def)
apply(subst PosOrd_def)
apply(rule_tac x=0#p in exI)
apply(simp add: pflat_len_Stars_simps pflat_len_simps)
using assms(2)
apply(simp add: pflat_len_simps)
apply(auto simp add: pflat_len_Stars_simps pflat_len_simps)
by (metis length_append of_nat_add)

```

**lemma** *PosOrd\_StarsI2*:

```

  assumes Stars vs1 : $\square$ val Stars vs2 flats vs1 = flats vs2
  shows Stars (v#vs1) : $\square$ val Stars (v#vs2)
using assms(1)
apply(subst (asm) PosOrd_ex_def)
apply(subst (asm) PosOrd_def)
apply(clarify)
apply(subst PosOrd_ex_def)
apply(subst PosOrd_def)
apply(case_tac p)
apply(simp add: pflat_len_simps)
apply(rule_tac x=Suc a#list in exI)
apply(auto simp add: pflat_len_Stars_simps pflat_len_simps assms(2))
done

```

**lemma** *PosOrd\_Stars\_appendI*:

```

  assumes Stars vs1 : $\square$ val Stars vs2 flat (Stars vs1) = flat (Stars vs2)
  shows Stars (vs @ vs1) : $\square$ val Stars (vs @ vs2)
using assms
apply(induct vs)
apply(simp)
apply(simp add: PosOrd_StarsI2)
done

```

```

lemma PosOrd_StarsE2:
  assumes Stars (v # vs1) : $\square$ val Stars (v # vs2)
  shows Stars vs1 : $\square$ val Stars vs2
using assms
apply(subst (asm) PosOrd_ex_def)
apply(erule exE)
apply(case_tac p)
apply(simp)
apply(simp add: PosOrd_def pflat_len_simps)
apply(subst PosOrd_ex_def)
apply(rule_tac x=[] in exI)
apply(simp add: PosOrd_def pflat_len_simps Pos_empty)
apply(simp)
apply(case_tac a)
apply(clarify)
apply(auto simp add: pflat_len_simps PosOrd_def pflat_len_def split: if_splits)[1]
apply(clarify)
apply(simp add: PosOrd_ex_def)
apply(rule_tac x=nat#list in exI)
apply(auto simp add: PosOrd_def pflat_len_simps)[1]
apply(case_tac q)
apply(simp add: PosOrd_def pflat_len_simps)
apply(clarify)
apply(drule_tac x=Suc a # lista in bspec)
apply(simp)
apply(auto simp add: PosOrd_def pflat_len_simps)[1]
apply(case_tac q)
apply(simp add: PosOrd_def pflat_len_simps)
apply(clarify)
apply(drule_tac x=Suc a # lista in bspec)
apply(simp)
apply(auto simp add: PosOrd_def pflat_len_simps)[1]
done

```

```

lemma PosOrd_Stars_appendE:
  assumes Stars (vs @ vs1) : $\square$ val Stars (vs @ vs2)
  shows Stars vs1 : $\square$ val Stars vs2
using assms
apply(induct vs)
apply(simp)
apply(simp add: PosOrd_StarsE2)
done

```

```

lemma PosOrd_Stars_append_eq:

```

```

assumes flats vs1 = flats vs2
shows Stars (vs @ vs1) : $\sqsubseteq$ val Stars (vs @ vs2)  $\longleftrightarrow$  Stars vs1 : $\sqsubseteq$ val Stars vs2
using assms
apply(rule_tac iffI)
apply(erule PosOrd_Stars_appendE)
apply(rule PosOrd_Stars_appendI)
apply(auto)
done

```

```

lemma PosOrd_almost_trichotomous:
shows v1 : $\sqsubseteq$ val v2  $\vee$  v2 : $\sqsubseteq$ val v1  $\vee$  (length (flat v1) = length (flat v2))
apply(auto simp add: PosOrd_ex_def)
apply(auto simp add: PosOrd_def)
apply(rule_tac x=[] in exI)
apply(auto simp add: Pos_empty_pflat_len_simps)
apply(drule_tac x=[] in spec)
apply(auto simp add: Pos_empty_pflat_len_simps)
done

```

## 20 The Posix Value is smaller than any other Value

```

lemma Posix_PosOrd:
assumes s  $\in$  r  $\rightarrow$  v1 v2  $\in$  LV r s
shows v1 : $\sqsubseteq$ val v2
using assms
proof (induct arbitrary: v2 rule: Posix.induct)
case (Posix_ONE v)
have v  $\in$  LV ONE [] by fact
then have v = Void
by (simp add: LV_simps)
then show Void : $\sqsubseteq$ val v
by (simp add: PosOrd_ex_eq_def)
next
case (Posix_CHAR c v)
have v  $\in$  LV (CHAR c) [c] by fact
then have v = Char c
by (simp add: LV_simps)
then show Char c : $\sqsubseteq$ val v
by (simp add: PosOrd_ex_eq_def)
next
case (Posix_ALT1 s r1 v r2 v2)
have as1: s  $\in$  r1  $\rightarrow$  v by fact
have IH:  $\bigwedge$ v2. v2  $\in$  LV r1 s  $\implies$  v : $\sqsubseteq$ val v2 by fact
have v2  $\in$  LV (ALT r1 r2) s by fact
then have  $\models$  v2 : ALT r1 r2 flat v2 = s

```

```

by(auto simp add: LV_def prefix_list_def)
then consider
  (Left)  $v3$  where  $v2 = \text{Left } v3 \models v3 : r1 \text{ flat } v3 = s$ 
| (Right)  $v3$  where  $v2 = \text{Right } v3 \models v3 : r2 \text{ flat } v3 = s$ 
by (auto elim: Prf.cases)
then show  $\text{Left } v : \sqsubseteq_{\text{val}} v2$ 
proof(cases)
  case (Left  $v3$ )
  have  $v3 \in LV \ r1 \ s$  using Left(2,3)
  by (auto simp add: LV_def prefix_list_def)
  with IH have  $v : \sqsubseteq_{\text{val}} v3$  by simp
  moreover
  have  $\text{flat } v3 = \text{flat } v$  using as1 Left(3)
  by (simp add: Posix1(2))
  ultimately have  $\text{Left } v : \sqsubseteq_{\text{val}} \text{Left } v3$ 
  by (simp add: PosOrd_ex_eq_def PosOrd_Left_eq)
  then show  $\text{Left } v : \sqsubseteq_{\text{val}} v2$  unfolding Left .
next
  case (Right  $v3$ )
  have  $\text{flat } v3 = \text{flat } v$  using as1 Right(3)
  by (simp add: Posix1(2))
  then have  $\text{Left } v : \sqsubseteq_{\text{val}} \text{Right } v3$ 
  unfolding PosOrd_ex_eq_def
  by (simp add: PosOrd_Left_Right)
  then show  $\text{Left } v : \sqsubseteq_{\text{val}} v2$  unfolding Right .
qed
next
  case (Posix_ALT2  $s \ r2 \ v \ r1 \ v2$ )
  have  $as1 : s \in r2 \rightarrow v$  by fact
  have  $as2 : s \notin L \ r1$  by fact
  have  $IH : \bigwedge v2. v2 \in LV \ r2 \ s \implies v : \sqsubseteq_{\text{val}} v2$  by fact
  have  $v2 \in LV \ (ALT \ r1 \ r2) \ s$  by fact
  then have  $\models v2 : ALT \ r1 \ r2 \ \text{flat } v2 = s$ 
  by(auto simp add: LV_def prefix_list_def)
  then consider
    (Left)  $v3$  where  $v2 = \text{Left } v3 \models v3 : r1 \ \text{flat } v3 = s$ 
  | (Right)  $v3$  where  $v2 = \text{Right } v3 \models v3 : r2 \ \text{flat } v3 = s$ 
  by (auto elim: Prf.cases)
  then show  $\text{Right } v : \sqsubseteq_{\text{val}} v2$ 
  proof (cases)
    case (Right  $v3$ )
    have  $v3 \in LV \ r2 \ s$  using Right(2,3)
    by (auto simp add: LV_def prefix_list_def)
    with IH have  $v : \sqsubseteq_{\text{val}} v3$  by simp
    moreover

```

```

have flat v3 = flat v using as1 Right(3)
  by (simp add: Posix1(2))
ultimately have Right v :  $\sqsubseteq$ val Right v3
  by (auto simp add: PosOrd_ex_eq_def PosOrd_Right1)
then show Right v :  $\sqsubseteq$ val v2 unfolding Right .
next
case (Left v3)
have v3  $\in$  LV r1 s using Left(2,3) as2
  by (auto simp add: LV_def prefix_list_def)
then have flat v3 = flat v  $\wedge$   $\models$  v3 : r1 using as1 Left(3)
  by (simp add: Posix1(2) LV_def)
then have False using as1 as2 Left
  by (auto simp add: Posix1(2) L_flat_Prfl)
then show Right v :  $\sqsubseteq$ val v2 by simp
qed
next
case (Posix_SEQ s1 r1 v1 s2 r2 v2 v3)
have s1  $\in$  r1  $\rightarrow$  v1 s2  $\in$  r2  $\rightarrow$  v2 by fact+
then have as1: s1 = flat v1 s2 = flat v2 by (simp_all add: Posix1(2))
have IH1:  $\bigwedge$ v3. v3  $\in$  LV r1 s1  $\implies$  v1 :  $\sqsubseteq$ val v3 by fact
have IH2:  $\bigwedge$ v3. v3  $\in$  LV r2 s2  $\implies$  v2 :  $\sqsubseteq$ val v3 by fact
have cond:  $\neg$  ( $\exists$  s3 s4. s3  $\neq$   $\square$   $\wedge$  s3 @ s4 = s2  $\wedge$  s1 @ s3  $\in$  L r1  $\wedge$  s4  $\in$  L r2) by fact
have v3  $\in$  LV (SEQ r1 r2) (s1 @ s2) by fact
then obtain v3a v3b where eqs:
  v3 = Seq v3a v3b  $\models$  v3a : r1  $\models$  v3b : r2
  flat v3a @ flat v3b = s1 @ s2
  by (force simp add: prefix_list_def LV_def elim: Prf.cases)
with cond have flat v3a  $\sqsubseteq$ pre s1 unfolding prefix_list_def
  by (smt L_flat_Prfl append_eq_append_conv2 append_self_conv)
then have flat v3a  $\sqsubseteq$ spre s1  $\vee$  (flat v3a = s1  $\wedge$  flat v3b = s2) using eqs
  by (simp add: spre_list_def append_eq_conv_conj)
then have q2: v1 :  $\sqsubseteq$ val v3a  $\vee$  (flat v3a = s1  $\wedge$  flat v3b = s2)
  using PosOrd_spre1 as1(1) eqs by blast
then have v1 :  $\sqsubseteq$ val v3a  $\vee$  (v3a  $\in$  LV r1 s1  $\wedge$  v3b  $\in$  LV r2 s2) using eqs(2,3)
  by (auto simp add: LV_def)
then have v1 :  $\sqsubseteq$ val v3a  $\vee$  (v1 :  $\sqsubseteq$ val v3a  $\wedge$  v2 :  $\sqsubseteq$ val v3b) using IH1 IH2 by blast
then have Seq v1 v2 :  $\sqsubseteq$ val Seq v3a v3b using eqs q2 as1
  unfolding PosOrd_ex_eq_def by (auto simp add: PosOrd_SeqI1 PosOrd_Seq_eq)
then show Seq v1 v2 :  $\sqsubseteq$ val v3 unfolding eqs by blast
next
case (Posix_STAR1 s1 r v s2 vs v3)
have s1  $\in$  r  $\rightarrow$  v s2  $\in$  STAR r  $\rightarrow$  Stars vs by fact+
then have as1: s1 = flat v s2 = flat (Stars vs) by (auto dest: Posix1(2))
have IH1:  $\bigwedge$ v3. v3  $\in$  LV r s1  $\implies$  v :  $\sqsubseteq$ val v3 by fact
have IH2:  $\bigwedge$ v3. v3  $\in$  LV (STAR r) s2  $\implies$  Stars vs :  $\sqsubseteq$ val v3 by fact

```

```

have cond:  $\neg (\exists s_3 s_4. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L r \wedge s_4 \in L (STAR r))$ 
by fact
have cond2:  $flat\ v \neq []$  by fact
have  $v_3 \in LV (STAR r) (s_1 @ s_2)$  by fact
then consider
   $(NonEmpty)\ v_3a\ vs_3$  where  $v_3 = Stars\ (v_3a \# vs_3)$ 
   $\models v_3a : r \models Stars\ vs_3 : STAR\ r$ 
   $flat\ (Stars\ (v_3a \# vs_3)) = s_1 @ s_2$ 
   $| (Empty)\ v_3 = Stars\ []$ 
unfolding LV_def
apply(auto)
apply(erule Prf.cases)
apply(auto)
apply(case_tac vs)
apply(auto intro: Prf.intros)
done
then show  $Stars\ (v \# vs) : \sqsubseteq_{val}\ v_3$ 
proof (cases)
  case  $(NonEmpty\ v_3a\ vs_3)$ 
  have  $flat\ (Stars\ (v_3a \# vs_3)) = s_1 @ s_2$  using NonEmpty(4) .
  with cond have  $flat\ v_3a \sqsubseteq_{pre}\ s_1$  using NonEmpty(2,3)
  unfolding prefix_list_def
  by (smt L_flat_Prfl append_Nil2 append_eq_append_conv2 flat.simps(7))
  then have  $flat\ v_3a \sqsubseteq_{spre}\ s_1 \vee (flat\ v_3a = s_1 \wedge flat\ (Stars\ vs_3) = s_2)$  using
NonEmpty(4)
  by (simp add: sprefix_list_def append_eq_conv_conj)
  then have  $q_2: v : \sqsubseteq_{val}\ v_3a \vee (flat\ v_3a = s_1 \wedge flat\ (Stars\ vs_3) = s_2)$ 
  using PosOrd_spreI as1(1) NonEmpty(4) by blast
  then have  $v : \sqsubseteq_{val}\ v_3a \vee (v_3a \in LV\ r\ s_1 \wedge Stars\ vs_3 \in LV\ (STAR\ r)\ s_2)$ 
  using NonEmpty(2,3) by (auto simp add: LV_def)
  then have  $v : \sqsubseteq_{val}\ v_3a \vee (v : \sqsubseteq_{val}\ v_3a \wedge Stars\ vs : \sqsubseteq_{val}\ Stars\ vs_3)$  using IH1 IH2
by blast
  then have  $v : \sqsubseteq_{val}\ v_3a \vee (v = v_3a \wedge Stars\ vs : \sqsubseteq_{val}\ Stars\ vs_3)$ 
  unfolding PosOrd_ex_eq_def by auto
  then have  $Stars\ (v \# vs) : \sqsubseteq_{val}\ Stars\ (v_3a \# vs_3)$  using NonEmpty(4) q2 as1
  unfolding PosOrd_ex_eq_def
  using PosOrd_StarsI PosOrd_StarsI2 by auto
  then show  $Stars\ (v \# vs) : \sqsubseteq_{val}\ v_3$  unfolding NonEmpty by blast
next
  case Empty
  have  $v_3 = Stars\ []$  by fact
  then show  $Stars\ (v \# vs) : \sqsubseteq_{val}\ v_3$ 
  unfolding PosOrd_ex_eq_def using cond2
  by (simp add: PosOrd_shorter1)
qed

```

```

next
  case (Posix_STAR2 r v2)
  have  $v2 \in LV (STAR\ r)$  by fact
  then have  $v2 = Stars$  by
    unfolding LV_def by (auto elim: Prf.cases)
  then show  $Stars \sqsubseteq \sqsubseteq val\ v2$ 
  by (simp add: PosOrd_ex_eq_def)
qed

```

```

lemma Posix_PosOrd_reverse:
  assumes  $s \in r \rightarrow v1$ 
  shows  $\neg(\exists v2 \in LV\ r\ s.\ v2 \sqsubseteq \sqsubseteq val\ v1)$ 
using assms
by (metis Posix_PosOrd_less_irrefl PosOrd_def
  PosOrd_ex_eq_def PosOrd_ex_def PosOrd_trans)

```

```

lemma PosOrd_Posix:
  assumes  $v1 \in LV\ r\ s \forall v2 \in LV\ r\ s.\ \neg v2 \sqsubseteq \sqsubseteq val\ v1$ 
  shows  $s \in r \rightarrow v1$ 
proof –
  have  $s \in L\ r$  using assms(1) unfolding LV_def
  using L_flat_Prfl by blast
  then obtain vposix where  $vp: s \in r \rightarrow vposix$ 
  using lexer_correct_Some by blast
  with assms(1) have  $vposix \sqsubseteq \sqsubseteq val\ v1$  by (simp add: Posix_PosOrd)
  then have  $vposix = v1 \vee vposix \sqsubseteq \sqsubseteq val\ v1$  unfolding PosOrd_ex_eq2 by auto
  moreover
  { assume  $vposix \sqsubseteq \sqsubseteq val\ v1$ 
    moreover
    have  $vposix \in LV\ r\ s$  using vp
    using Posix_LV by blast
    ultimately have False using assms(2) by blast
  }
  ultimately show  $s \in r \rightarrow v1$  using vp by blast
qed

```

```

lemma Least_existence:
  assumes  $LV\ r\ s \neq \{\}$ 
  shows  $\exists vmin \in LV\ r\ s.\ \forall v \in LV\ r\ s.\ vmin \sqsubseteq \sqsubseteq val\ v$ 
proof –
  from assms
  obtain vposix where  $s \in r \rightarrow vposix$ 
  unfolding LV_def
  using L_flat_Prfl lexer_correct_Some by blast

```



```

then have  $\forall v \in LV\ r\ s. vposix : \sqsubseteq val\ v$ 
by (simp add: Posix_PosOrd)
then show  $\exists vmin \in LV\ r\ s. \forall v \in LV\ r\ s. vmin : \sqsubseteq val\ v$ 
using Posix_LV  $\langle s \in r \rightarrow vposix \rangle$  by blast
qed
    
```

```

lemma Least_existence1:
  assumes  $LV\ r\ s \neq \{\}$ 
  shows  $\exists !vmin \in LV\ r\ s. \forall v \in LV\ r\ s. vmin : \sqsubseteq val\ v$ 
using Least_existence[OF assms] assms
using PosOrdeq_antisym by blast
    
```

```

lemma Least_existence2:
  assumes  $LV\ r\ s \neq \{\}$ 
  shows  $\exists !vmin \in LV\ r\ s. lexer\ r\ s = Some\ vmin \wedge (\forall v \in LV\ r\ s. vmin : \sqsubseteq val\ v)$ 
using Least_existence[OF assms] assms
using PosOrdeq_antisym
using PosOrd_Posix PosOrd_ex_eq2 lexer_correctness(1) by auto
    
```

```

lemma Least_existence1_pre:
  assumes  $LV\ r\ s \neq \{\}$ 
  shows  $\exists !vmin \in LV\ r\ s. \forall v \in (LV\ r\ s \cup \{v'. flat\ v' \sqsubseteq spre\ s\}). vmin : \sqsubseteq val\ v$ 
using Least_existence[OF assms] assms
apply –
apply (erule bexE)
apply (rule_tac a=vmin in ex1I)
apply (auto)[I]
apply (metis PosOrd_Posix PosOrd_ex_eq2 PosOrd_spre1 PosOrdeq_antisym Posix1(2))
apply (auto)[I]
apply (simp add: PosOrdeq_antisym)
done
    
```

```

lemma
  shows partial_order_on UNIV  $\{(v1, v2). v1 : \sqsubseteq val\ v2\}$ 
apply (simp add: partial_order_on_def)
apply (simp add: preorder_on_def refl_on_def)
apply (simp add: PosOrdeq_refl)
apply (auto)
apply (rule transI)
apply (auto intro: PosOrdeq_trans)[I]
apply (rule antisymI)
apply (simp add: PosOrdeq_antisym)
done
    
```

```

lemma
  wf {(v1, v2). v1 : $\square$ val v2  $\wedge$  v1  $\in$  LV r s  $\wedge$  v2  $\in$  LV r s}
apply(rule finite_acyclic_wf)
prefer 2
apply(simp add: acyclic_def)
apply(induct_tac rule: trancl.induct)
apply(auto)[1]
oops

```

```

unused-thms

```

```

end

```

```

theory Sulzmann
  imports Lexer
begin

```

## 21 Bit-Encodings

```

datatype bit = Z | S

```

```

fun
  code :: val  $\Rightarrow$  bit list
where
  code Void = []
| code (Char c) = []
| code (Left v) = Z # (code v)
| code (Right v) = S # (code v)
| code (Seq v1 v2) = (code v1) @ (code v2)
| code (Stars []) = [S]
| code (Stars (v # vs)) = (Z # code v) @ code (Stars vs)

```

```

fun
  Stars_add :: val  $\Rightarrow$  val  $\Rightarrow$  val
where
  Stars_add v (Stars vs) = Stars (v # vs)

```

```

function
  decode' :: bit list  $\Rightarrow$  rexp  $\Rightarrow$  (val * bit list)
where
  decode' ds ZERO = (Void, [])
| decode' ds ONE = (Void, ds)
| decode' ds (CHAR d) = (Char d, ds)

```

```

| decode' [] (ALT r1 r2) = (Void, [])
| decode' (Z # ds) (ALT r1 r2) = (let (v, ds') = decode' ds r1 in (Left v, ds'))
| decode' (S # ds) (ALT r1 r2) = (let (v, ds') = decode' ds r2 in (Right v, ds'))
| decode' ds (SEQ r1 r2) = (let (v1, ds') = decode' ds r1 in
                           let (v2, ds'') = decode' ds' r2 in (Seq v1 v2, ds''))
| decode' [] (STAR r) = (Void, [])
| decode' (S # ds) (STAR r) = (Stars [], ds)
| decode' (Z # ds) (STAR r) = (let (v, ds') = decode' ds r in
                              let (vs, ds'') = decode' ds' (STAR r)
                              in (Stars_add v vs, ds''))

```

**by** pat\_completeness auto

**lemma** decode'\_smaller:

```

  assumes decode'_dom (ds, r)
  shows length (snd (decode' ds r)) ≤ length ds

```

**using** assms

**apply**(induct ds r)

**apply**(auto simp add: decode'.psimps split: prod.split)

**using** dual\_order.trans **apply** blast

**by** (meson dual\_order.trans le\_SucI)

**termination** decode'

```

apply(relation inv_image (measure(%cs. size cs) <*>lex*> measure(%s. size s)) (%(ds,r).
(r,ds)))
apply(auto dest!: decode'_smaller)
by (metis less_Suc_eq_le snd_conv)

```

**definition**

```

decode :: bit list ⇒ rexp ⇒ val option

```

**where**

```

decode ds r  $\stackrel{\text{def}}{=} (let (v, ds') = decode' ds r
                    in (if ds' = [] then Some v else None))$ 

```

**lemma** decode'\_code\_Stars:

```

  assumes ∀ v ∈ set vs. ⊨ v : r ∧ (∀ x. decode' (code v @ x) r = (v, x)) ∧ flat v ≠ []
  shows decode' (code (Stars vs) @ ds) (STAR r) = (Stars vs, ds)

```

**using** assms

**apply**(induct vs)

**apply**(auto)

**done**

**lemma** decode'\_code:

```

  assumes ⊨ v : r
  shows decode' ((code v) @ ds) r = (v, ds)

```

**using** assms

```

apply(induct v r arbitrary: ds)
apply(auto)
using decode'_code_Stars by blast

```

**lemma** *decode\_code*:

```

assumes  $\models v : r$ 
shows decode (code v) r = Some v
using assms unfolding decode_def
by (smt append_Nil2 decode'_code old.prod.case)

```

**datatype** *arexp* =

```

  AZERO
| AONE bit list
| ACHAR bit list char
| ASEQ bit list arexp arexp
| AALT bit list arexp arexp
| ASTAR bit list arexp

```

**fun** *fuse* :: *bit list*  $\Rightarrow$  *arexp*  $\Rightarrow$  *arexp* **where**

```

  fuse bs AZERO = AZERO
| fuse bs (AONE cs) = AONE (bs @ cs)
| fuse bs (ACHAR cs c) = ACHAR (bs @ cs) c
| fuse bs (AALT cs r1 r2) = AALT (bs @ cs) r1 r2
| fuse bs (ASEQ cs r1 r2) = ASEQ (bs @ cs) r1 r2
| fuse bs (ASTAR cs r) = ASTAR (bs @ cs) r

```

**fun** *intern* :: *rexp*  $\Rightarrow$  *arexp* **where**

```

  intern ZERO = AZERO
| intern ONE = AONE []
| intern (CHAR c) = ACHAR [] c
| intern (ALT r1 r2) = AALT [] (fuse [Z] (intern r1))
   (fuse [S] (intern r2))
| intern (SEQ r1 r2) = ASEQ [] (intern r1) (intern r2)
| intern (STAR r) = ASTAR [] (intern r)

```

**fun** *retrieve* :: *arexp*  $\Rightarrow$  *val*  $\Rightarrow$  *bit list* **where**

```

  retrieve (AONE bs) Void = bs
| retrieve (ACHAR bs c) (Char d) = bs
| retrieve (AALT bs r1 r2) (Left v) = bs @ retrieve r1 v
| retrieve (AALT bs r1 r2) (Right v) = bs @ retrieve r2 v
| retrieve (ASEQ bs r1 r2) (Seq v1 v2) = bs @ retrieve r1 v1 @ retrieve r2 v2
| retrieve (ASTAR bs r) (Stars []) = bs @ [S]
| retrieve (ASTAR bs r) (Stars (v#vs)) =

```

$bs @ [Z] @ \text{retrieve } r \vee @ \text{retrieve } (ASTAR [] r)$  (*Stars vs*)

**fun**

$erase :: arexp \Rightarrow rexp$

**where**

$erase AZERO = ZERO$   
 $| erase (AONE \_) = ONE$   
 $| erase (ACHAR \_ c) = CHAR c$   
 $| erase (AALT \_ r1 r2) = ALT (erase r1) (erase r2)$   
 $| erase (ASEQ \_ r1 r2) = SEQ (erase r1) (erase r2)$   
 $| erase (ASTAR \_ r) = STAR (erase r)$

**fun**

$bnullable :: arexp \Rightarrow bool$

**where**

$bnullable (AZERO) = False$   
 $| bnullable (AONE bs) = True$   
 $| bnullable (ACHAR bs c) = False$   
 $| bnullable (AALT bs r1 r2) = (bnullable r1 \vee bnullable r2)$   
 $| bnullable (ASEQ bs r1 r2) = (bnullable r1 \wedge bnullable r2)$   
 $| bnullable (ASTAR bs r) = True$

**fun**

$bmkeps :: arexp \Rightarrow \text{bit list}$

**where**

$bmkeps(AONE bs) = bs$   
 $| bmkeps(ASEQ bs r1 r2) = bs @ (bmkeps r1) @ (bmkeps r2)$   
 $| bmkeps(AALT bs r1 r2) = (\text{if } bnullable(r1) \text{ then } bs @ (bmkeps r1) \text{ else } bs @ (bmkeps r2))$   
 $| bmkeps(ASTAR bs r) = bs @ [S]$

**fun**

$bder :: char \Rightarrow arexp \Rightarrow arexp$

**where**

$bder c (AZERO) = AZERO$   
 $| bder c (AONE bs) = AZERO$   
 $| bder c (ACHAR bs d) = (\text{if } c = d \text{ then } AONE bs \text{ else } AZERO)$   
 $| bder c (AALT bs r1 r2) = AALT bs (bder c r1) (bder c r2)$   
 $| bder c (ASEQ bs r1 r2) =$   
      $(\text{if } bnullable r1$   
          $\text{then } AALT bs (ASEQ [] (bder c r1) r2) (\text{fuse } (bmkeps r1) (bder c r2))$   
          $\text{else } ASEQ bs (bder c r1) r2)$   
 $| bder c (ASTAR bs r) = ASEQ bs (\text{fuse } [Z] (bder c r)) (ASTAR [] r)$

```

fun
  bders :: arexp ⇒ string ⇒ arexp
where
  bders r [] = r
  | bders r (c#s) = bders (bder c r) s

lemma bders_append:
  bders r (s1 @ s2) = bders (bders r s1) s2
apply(induct s1 arbitrary: r s2)
apply(simp_all)
done

lemma bnullable_correctness:
  shows nullable (erase r) = bnullable r
apply(induct r)
apply(simp_all)
done

lemma erase_fuse:
  shows erase (fuse bs r) = erase r
apply(induct r)
apply(simp_all)
done

lemma erase_intern[simp]:
  shows erase (intern r) = r
apply(induct r)
apply(simp_all add: erase_fuse)
done

lemma erase_bder[simp]:
  shows erase (bder a r) = der a (erase r)
apply(induct r)
apply(simp_all add: erase_fuse bnullable_correctness)
done

lemma erase_bders[simp]:
  shows erase (bders r s) = ders s (erase r)
apply(induct s arbitrary: r)
apply(simp_all)
done

lemma retrieve_encode_STARS:
  assumes ∀ v ∈ set vs. ⊨ v : r ∧ code v = retrieve (intern r) v

```

```

shows code (Stars vs) = retrieve (ASTAR [] (intern r)) (Stars vs)
using assms
apply(induct vs)
apply(simp_all)
done

```

```

lemma retrieve_fuse2:
  assumes  $\models v : (\text{erase } r)$ 
  shows retrieve (fuse bs r) v = bs @ retrieve r v
  using assms
  apply(induct r arbitrary: v bs)
  using retrieve_encode_STARS
  apply(auto elim!: Prf_elims)
  apply(case_tac vs)
  apply(simp)
  apply(simp)
  done

```

```

lemma retrieve_fuse:
  assumes  $\models v : r$ 
  shows retrieve (fuse bs (intern r)) v = bs @ retrieve (intern r) v
  using assms
  by (simp_all add: retrieve_fuse2)

```

```

lemma retrieve_code:
  assumes  $\models v : r$ 
  shows code v = retrieve (intern r) v
  using assms
  apply(induct v r)
  apply(simp_all add: retrieve_fuse retrieve_encode_STARS)
  done

```

```

lemma bmkeys_retrieve:
  assumes nullable (erase r)
  shows bmkeys r = retrieve r (mkeys (erase r))
  using assms
  apply(induct r)
  apply(auto simp add: nullable_correctness)
  done

```

```

lemma bder_retrieve:
  assumes  $\models v : \text{der } c (\text{erase } r)$ 
  shows retrieve (bder c r) v = retrieve r (inval (erase r) c v)

```

```

using assms
apply(induct r arbitrary: v)
apply(auto elim!: Prf_elims simp add: retrieve_fuse2 bnullable_correctness bmkeys_retrieve)
done

```

**lemma** *MAIN\_decode:*

```

assumes  $\models v : \text{ders } s \ r$ 
shows  $\text{Some } (\text{flex } r \ \text{id } s \ v) = \text{decode } (\text{retrieve } (\text{bders } (\text{intern } r) \ s) \ v) \ r$ 
using assms
proof (induct s arbitrary: v rule: rev_induct)
  case Nil
    have  $\models v : \text{ders } [] \ r$  by fact
    then have  $\models v : r$  by simp
    then have  $\text{Some } v = \text{decode } (\text{retrieve } (\text{intern } r) \ v) \ r$ 
      using decode_code retrieve_code by auto
    then show  $\text{Some } (\text{flex } r \ \text{id } [] \ v) = \text{decode } (\text{retrieve } (\text{bders } (\text{intern } r) \ []) \ v) \ r$ 
      by simp
  next
    case (snoc c s v)
    have IH:  $\bigwedge v. \models v : \text{ders } s \ r \implies$ 
       $\text{Some } (\text{flex } r \ \text{id } s \ v) = \text{decode } (\text{retrieve } (\text{bders } (\text{intern } r) \ s) \ v) \ r$  by fact
    have asm:  $\models v : \text{ders } (s \ @ \ [c]) \ r$  by fact
    then have asm2:  $\models \text{injval } (\text{ders } s \ r) \ c \ v : \text{ders } s \ r$ 
      by(simp add: Prf_injval ders_append)
    have  $\text{Some } (\text{flex } r \ \text{id } (s \ @ \ [c]) \ v) = \text{Some } (\text{flex } r \ \text{id } s \ (\text{injval } (\text{ders } s \ r) \ c \ v))$ 
      by (simp add: flex_append)
    also have  $\dots = \text{decode } (\text{retrieve } (\text{bders } (\text{intern } r) \ s) \ (\text{injval } (\text{ders } s \ r) \ c \ v)) \ r$ 
      using asm2 IH by simp
    also have  $\dots = \text{decode } (\text{retrieve } (\text{bder } c \ (\text{bders } (\text{intern } r) \ s)) \ v) \ r$ 
      using asm by(simp_all add: bder_retrieve ders_append)
    finally show  $\text{Some } (\text{flex } r \ \text{id } (s \ @ \ [c]) \ v) =$ 
       $\text{decode } (\text{retrieve } (\text{bders } (\text{intern } r) \ (s \ @ \ [c])) \ v) \ r$  by (simp add: bders_append)
  qed

```

**definition** *blexer where*

```

blexer r s  $\stackrel{\text{def}}{=} \text{if } \text{bnullable } (\text{bders } (\text{intern } r) \ s) \ \text{then}$ 
   $\text{decode } (\text{bmkeys } (\text{bders } (\text{intern } r) \ s)) \ r$  else None

```

**lemma** *blexer\_correctness:*

```

shows blexer r s = lexer r s

```

**proof** –

```

{ define bds where bds  $\stackrel{\text{def}}{=} \text{bders } (\text{intern } r) \ s$ 
  define ds where ds  $\stackrel{\text{def}}{=} \text{ders } s \ r$ 
  assume asm: nullable ds

```



```

have era: erase bds = ds
  unfolding ds_def bds_def by simp
have mke:  $\models$  mkeps ds : ds
  using asm by (simp add: mkeps_nullable)
have decode (bmkeps bds) r = decode (retrieve bds (mkeps ds)) r
  using bmkeps_retrieve
  using asm era by (simp add: bmkeps_retrieve)
also have ... = Some (flex r id s (mkeps ds))
  using mke by (simp_all add: MAIN_decode ds_def bds_def)
finally have decode (bmkeps bds) r = Some (flex r id s (mkeps ds))
  unfolding bds_def ds_def .
}
then show blexer r s = lexer r s
  unfolding blexer_def lexer_flex
  apply(subst bnullable_correctness[symmetric])
  apply(simp)
done
qed

```

**end**

## 22 Introduction

Brzozowski [4] introduced the notion of the *derivative*  $r \setminus c$  of a regular expression  $r$  w.r.t. a character  $c$ , and showed that it gave a simple solution to the problem of matching a string  $s$  with a regular expression  $r$ : if the derivative of  $r$  w.r.t. (in succession) all the characters of the string matches the empty string, then  $r$  matches  $s$  (and *vice versa*). The derivative has the property (which may almost be regarded as its specification) that, for every string  $s$  and regular expression  $r$  and character  $c$ , one has  $cs \in L(r)$  if and only if  $s \in L(r \setminus c)$ . The beauty of Brzozowski's derivatives is that they are neatly expressible in any functional language, and easily definable and reasoned about in theorem provers—the definitions just consist of inductive datatypes and simple recursive functions. A mechanised correctness proof of Brzozowski's matcher in for example HOL4 has been mentioned by Owens and Slind [14]. Another one in Isabelle/HOL is part of the work by Krauss and Nipkow [9]. And another one in Coq is given by Coquand and Siles [5].

If a regular expression matches a string, then in general there is more than one way of how the string is matched. There are two commonly used disambiguation strategies to generate a unique answer: one is called GREEDY matching [6] and the other is POSIX matching [1, 10, 12, 16, 17]. For example consider the string  $xy$  and the regular expression  $(x + y + xy)^*$ . Either the string can be matched in two 'iterations' by the single letter-regular expressions  $x$  and  $y$ , or directly in one iteration by  $xy$ . The first case corresponds to GREEDY matching, which first matches with the left-most symbol and only matches the next symbol in case of a mismatch (this is greedy in the sense of preferring instant

gratification to delayed repletion). The second case is POSIX matching, which prefers the longest match.

In the context of lexing, where an input string needs to be split up into a sequence of tokens, POSIX is the more natural disambiguation strategy for what programmers consider basic syntactic building blocks in their programs. These building blocks are often specified by some regular expressions, say  $r_{key}$  and  $r_{id}$  for recognising keywords and identifiers, respectively. There are a few underlying (informal) rules behind tokenising a string in a POSIX [1] fashion:

- *The Longest Match Rule* (or “*Maximal Munch Rule*”): The longest initial substring matched by any regular expression is taken as next token.
- *Priority Rule*: For a particular longest initial substring, the first (leftmost) regular expression that can match determines the token.
- *Star Rule*: A subexpression repeated by  $*$  shall not match an empty string unless this is the only match for the repetition.
- *Empty String Rule*: An empty string shall be considered to be longer than no match at all.

Consider for example a regular expression  $r_{key}$  for recognising keywords such as *if*, *then* and so on; and  $r_{id}$  recognising identifiers (say, a single character followed by characters or numbers). Then we can form the regular expression  $(r_{key} + r_{id})^*$  and use POSIX matching to tokenise strings, say *iffoo* and *if*. For *iffoo* we obtain by the Longest Match Rule a single identifier token, not a keyword followed by an identifier. For *if* we obtain by the Priority Rule a keyword token, not an identifier token—even if  $r_{id}$  matches also. By the Star Rule we know  $(r_{key} + r_{id})^*$  matches *iffoo*, respectively *if*, in exactly one ‘iteration’ of the star. The Empty String Rule is for cases where, for example, the regular expression  $(a^*)^*$  matches against the string *bc*. Then the longest initial matched substring is the empty string, which is matched by both the whole regular expression and the parenthesised subexpression.

One limitation of Brzozowski’s matcher is that it only generates a YES/NO answer for whether a string is being matched by a regular expression. Sulzmann and Lu [16] extended this matcher to allow generation not just of a YES/NO answer but of an actual matching, called a [lexical] *value*. Assuming a regular expression matches a string, values encode the information of *how* the string is matched by the regular expression—that is, which part of the string is matched by which part of the regular expression. For this consider again the string *xy* and the regular expression  $(x + (y + xy))^*$  (this time fully parenthesised). We can view this regular expression as tree and if the string *xy* is matched by two Star ‘iterations’, then the *x* is matched by the left-most alternative in this tree and the *y* by the right-left alternative. This suggests to record this matching as

$$\text{Stars} [\text{Left} (\text{Char } x), \text{Right} (\text{Left} (\text{Char } y))]$$

where *Stars*, *Left*, *Right* and *Char* are constructors for values. *Stars* records how many iterations were used; *Left*, respectively *Right*, which alternative is used. This ‘tree view’ leads naturally to the idea that regular expressions act as types and values as inhabiting those types (see, for example, [8]). The value for matching *xy* in a single ‘iteration’, i.e. the POSIX value, would look as follows

$$\text{Stars } [Seq (Char\ x) (Char\ y)]$$

where *Stars* has only a single-element list for the single iteration and *Seq* indicates that  $xy$  is matched by a sequence regular expression.

Sulzmann and Lu give a simple algorithm to calculate a value that appears to be the value associated with POSIX matching. The challenge then is to specify that value, in an algorithm-independent fashion, and to show that Sulzmann and Lu’s derivative-based algorithm does indeed calculate a value that is correct according to the specification. The answer given by Sulzmann and Lu [16] is to define a relation (called an “order relation”) on the set of values of  $r$ , and to show that (once a string to be matched is chosen) there is a maximum element and that it is computed by their derivative-based algorithm. This proof idea is inspired by work of Frisch and Cardelli [6] on a GREEDY regular expression matching algorithm. However, we were not able to establish transitivity and totality for the “order relation” by Sulzmann and Lu. There are some inherent problems with their approach (of which some of the proofs are not published in [16]); perhaps more importantly, we give in this paper a simple inductive (and algorithm-independent) definition of what we call being a *POSIX value* for a regular expression  $r$  and a string  $s$ ; we show that the algorithm by Sulzmann and Lu computes such a value and that such a value is unique. Our proofs are both done by hand and checked in Isabelle/HOL. The experience of doing our proofs has been that this mechanical checking was absolutely essential: this subject area has hidden snares. This was also noted by Kuklewicz [10] who found that nearly all POSIX matching implementations are “buggy” [16, Page 203] and by Grathwohl et al [7, Page 36] who wrote:

*“The POSIX strategy is more complicated than the greedy because of the dependence on information about the length of matched strings in the various subexpressions.”*

**Contributions:** We have implemented in Isabelle/HOL the derivative-based regular expression matching algorithm of Sulzmann and Lu [16]. We have proved the correctness of this algorithm according to our specification of what a POSIX value is (inspired by work of Vansummeren [17]). Sulzmann and Lu sketch in [16] an informal correctness proof: but to us it contains unfillable gaps.<sup>4</sup> Our specification of a POSIX value consists of a simple inductive definition that given a string and a regular expression uniquely determines this value. We also show that our definition is equivalent to an ordering of values based on positions by Okui and Suzuki [12].

We extend our results to ??? Bitcoded version??

## 23 Preliminaries

Strings in Isabelle/HOL are lists of characters with the empty string being represented by the empty list, written  $[]$ , and list-cons being written as  $\_ :: \_$ . Often we use the usual

<sup>4</sup> An extended version of [16] is available at the website of its first author; this extended version already includes remarks in the appendix that their informal proof contains gaps, and possible fixes are not fully worked out.

bracket notation for lists also for strings; for example a string consisting of just a single character  $c$  is written  $[c]$ . We use the usual definitions for *prefixes* and *strict prefixes* of strings. By using the type *char* for characters we have a supply of finitely many characters roughly corresponding to the ASCII character set. Regular expressions are defined as usual as the elements of the following inductive datatype:

$$r := \mathbf{0} \mid \mathbf{1} \mid c \mid r_1 + r_2 \mid r_1 \cdot r_2 \mid r^*$$

where  $\mathbf{0}$  stands for the regular expression that does not match any string,  $\mathbf{1}$  for the regular expression that matches only the empty string and  $c$  for matching a character literal. The language of a regular expression is also defined as usual by the recursive function  $L$  with the six clauses:

$$\begin{aligned} (1) \quad L(\mathbf{0}) &\stackrel{\text{def}}{=} \emptyset \\ (2) \quad L(\mathbf{1}) &\stackrel{\text{def}}{=} \{\emptyset\} \\ (3) \quad L(c) &\stackrel{\text{def}}{=} \{[c]\} \\ (4) \quad L(r_1 \cdot r_2) &\stackrel{\text{def}}{=} L(r_1) @ L(r_2) \\ (5) \quad L(r_1 + r_2) &\stackrel{\text{def}}{=} L(r_1) \cup L(r_2) \\ (6) \quad L(r^*) &\stackrel{\text{def}}{=} (L(r))^* \end{aligned}$$

In clause (4) we use the operation  $_ @ _$  for the concatenation of two languages (it is also list-append for strings). We use the star-notation for regular expressions and for languages (in the last clause above). The star for languages is defined inductively by two clauses: (i) the empty string being in the star of a language and (ii) if  $s_1$  is in a language and  $s_2$  in the star of this language, then also  $s_1 @ s_2$  is in the star of this language. It will also be convenient to use the following notion of a *semantic derivative* (or *left quotient*) of a language defined as

$$\text{Der } c A \stackrel{\text{def}}{=} \{s \mid c :: s \in A\}.$$

For semantic derivatives we have the following equations (for example mechanically proved in [9]):

$$\begin{aligned} \text{Der } c \emptyset &\stackrel{\text{def}}{=} \emptyset \\ \text{Der } c \{\emptyset\} &\stackrel{\text{def}}{=} \emptyset \\ \text{Der } c \{[d]\} &\stackrel{\text{def}}{=} \text{if } c = d \text{ then } \{\emptyset\} \text{ else } \emptyset \\ \text{Der } c (A \cup B) &\stackrel{\text{def}}{=} \text{Der } c A \cup \text{Der } c B \\ \text{Der } c (A @ B) &\stackrel{\text{def}}{=} (\text{Der } c A @ B) \cup (\text{if } \emptyset \in A \text{ then } \text{Der } c B \text{ else } \emptyset) \\ \text{Der } c (A^*) &\stackrel{\text{def}}{=} \text{Der } c A @ A^* \end{aligned} \tag{1}$$

*Brzozowski's derivatives* of regular expressions [4] can be easily defined by two recursive functions: the first is from regular expressions to booleans (implementing a test when a regular expression can match the empty string), and the second takes a regular expression and a character to a (derivative) regular expression:

$nullable(\mathbf{0})$	$\stackrel{\text{def}}{=} False$
$nullable(\mathbf{1})$	$\stackrel{\text{def}}{=} True$
$nullable(c)$	$\stackrel{\text{def}}{=} False$
$nullable(r_1 + r_2)$	$\stackrel{\text{def}}{=} nullable\ r_1 \vee nullable\ r_2$
$nullable(r_1 \cdot r_2)$	$\stackrel{\text{def}}{=} nullable\ r_1 \wedge nullable\ r_2$
$nullable(r^*)$	$\stackrel{\text{def}}{=} True$
$\mathbf{0} \setminus c$	$\stackrel{\text{def}}{=} \mathbf{0}$
$\mathbf{1} \setminus c$	$\stackrel{\text{def}}{=} \mathbf{0}$
$d \setminus c$	$\stackrel{\text{def}}{=} \text{if } c = d \text{ then } \mathbf{1} \text{ else } \mathbf{0}$
$(r_1 + r_2) \setminus c$	$\stackrel{\text{def}}{=} (r_1 \setminus c) + (r_2 \setminus c)$
$(r_1 \cdot r_2) \setminus c$	$\stackrel{\text{def}}{=} \text{if } nullable\ r_1 \text{ then } (r_1 \setminus c) \cdot r_2 + (r_2 \setminus c) \text{ else } (r_1 \setminus c) \cdot r_2$
$(r^*) \setminus c$	$\stackrel{\text{def}}{=} (r \setminus c) \cdot r^*$

We may extend this definition to give derivatives w.r.t. strings:

$$\begin{aligned} r \setminus [] &\stackrel{\text{def}}{=} r \\ r \setminus (c :: s) &\stackrel{\text{def}}{=} (r \setminus c) \setminus s \end{aligned}$$

Given the equations in (1), it is a relatively easy exercise in mechanical reasoning to establish that

**Proposition 1.**

- (1)  $nullable\ r$  if and only if  $[] \in L(r)$ , and
- (2)  $L(r \setminus c) = Der\ c\ (L(r))$ .

With this in place it is also very routine to prove that the regular expression matcher defined as

$$match\ r\ s \stackrel{\text{def}}{=} nullable\ (r \setminus s)$$

gives a positive answer if and only if  $s \in L(r)$ . Consequently, this regular expression matching algorithm satisfies the usual specification for regular expression matching. While the matcher above calculates a provably correct YES/NO answer for whether a regular expression matches a string or not, the novel idea of Sulzmann and Lu [16] is to append another phase to this algorithm in order to calculate a [lexical] value. We will explain the details next.

## 24 POSIX Regular Expression Matching

There have been many previous works that use values for encoding *how* a regular expression matches a string. The clever idea by Sulzmann and Lu [16] is to define a function on values that mirrors (but inverts) the construction of the derivative on regular expressions. *Values* are defined as the inductive datatype

$$v := \text{Empty} \mid \text{Char } c \mid \text{Left } v \mid \text{Right } v \mid \text{Seq } v_1 v_2 \mid \text{Stars } vs$$

where we use  $vs$  to stand for a list of values. (This is similar to the approach taken by Frisch and Cardelli for GREEDY matching [6], and Sulzmann and Lu for POSIX matching [16]). The string underlying a value can be calculated by the *flat* function, written  $|\_$  and defined as:

$$\begin{array}{ll} |\text{Empty}| \stackrel{\text{def}}{=} [] & |\text{Seq } v_1 v_2| \stackrel{\text{def}}{=} |v_1| @ |v_2| \\ |\text{Char } c| \stackrel{\text{def}}{=} [c] & |\text{Stars } []| \stackrel{\text{def}}{=} [] \\ |\text{Left } v| \stackrel{\text{def}}{=} |v| & |\text{Stars } (v :: vs)| \stackrel{\text{def}}{=} |v| @ |\text{Stars } vs| \\ |\text{Right } v| \stackrel{\text{def}}{=} |v| & \end{array}$$

We will sometimes refer to the underlying string of a value as *flattened value*. We will also overload our notation and use  $|vs|$  for flattening a list of values and concatenating the resulting strings.

Sulzmann and Lu define inductively an *inhabitation relation* that associates values to regular expressions. We define this relation as follows:<sup>5</sup>

$$\begin{array}{c} \overline{\text{Empty} : \mathbf{1}} \\ \frac{v_1 : r_1}{\text{Left } v_1 : r_1 + r_2} \\ \frac{v_1 : r_1 \quad v_2 : r_2}{\text{Seq } v_1 v_2 : r_1 \cdot r_2} \end{array} \qquad \begin{array}{c} \overline{\text{Char } c : c} \\ \frac{v_2 : r_1}{\text{Right } v_2 : r_2 + r_1} \\ \frac{\forall v \in vs. v : r \wedge |v| \neq []}{\text{Stars } vs : r^*} \end{array}$$

where in the clause for *Stars* we use the notation  $v \in vs$  for indicating that  $v$  is a member in the list  $vs$ . We require in this rule that every value in  $vs$  flattens to a non-empty string. The idea is that *Stars*-values satisfy the informal Star Rule (see Introduction) where the  $*$  does not match the empty string unless this is the only match for the repetition. Note also that no values are associated with the regular expression  $\mathbf{0}$ , and that the only value associated with the regular expression  $\mathbf{1}$  is *Empty*. It is routine to establish how values “inhabiting” a regular expression correspond to the language of a regular expression, namely

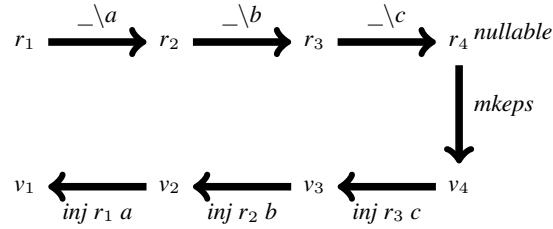
**Proposition 2.**  $L(r) = \{|v| \mid v : r\}$

Given a regular expression  $r$  and a string  $s$ , we define the set of all *Lexical Values* inhabited by  $r$  with the underlying string being  $s$ :<sup>6</sup>

$$LV r s \stackrel{\text{def}}{=} \{v \mid v : r \wedge |v| = s\}$$

<sup>5</sup> Note that the rule for *Stars* differs from our earlier paper [3]. There we used the original definition by Sulzmann and Lu which does not require that the values  $v \in vs$  flatten to a non-empty string. The reason for introducing the more restricted version of lexical values is convenience later on when reasoning about an ordering relation for values.

<sup>6</sup> Okui and Suzuki refer to our lexical values as *canonical values* in [12]. The notion of *non-problematic values* by Cardelli and Frisch [6] is related, but not identical to our lexical values.



**Fig. 1.** The two phases of the algorithm by Sulzmann & Lu [16], matching the string  $[a, b, c]$ . The first phase (the arrows from left to right) is Brzozowski’s matcher building successive derivatives. If the last regular expression is *nullable*, then the functions of the second phase are called (the top-down and right-to-left arrows): first *mkeps* calculates a value  $v_4$  witnessing how the empty string has been recognised by  $r_4$ . After that the function *inj* “injects back” the characters of the string into the values.

The main property of  $LV\ r\ s$  is that it is always finite.

**Proposition 3.** *finite* ( $LV\ r\ s$ )

This finiteness property does not hold in general if we remove the side-condition about  $|v| \neq []$  in the *Stars*-rule above. For example using Sulzmann and Lu’s less restrictive definition,  $LV\ (\mathbf{1}^*)\ []$  would contain infinitely many values, but according to our more restricted definition only a single value, namely  $LV\ (\mathbf{1}^*)\ [] = \{\text{Stars}\ []\}$ .

If a regular expression  $r$  matches a string  $s$ , then generally the set  $LV\ r\ s$  is not just a singleton set. In case of POSIX matching the problem is to calculate the unique lexical value that satisfies the (informal) POSIX rules from the Introduction. Graphically the POSIX value calculation algorithm by Sulzmann and Lu can be illustrated by the picture in Figure 1 where the path from the left to the right involving *derivatives/nullable* is the first phase of the algorithm (calculating successive Brzozowski’s derivatives) and *mkeps/inj*, the path from right to left, the second phase. This picture shows the steps required when a regular expression, say  $r_1$ , matches the string  $[a, b, c]$ . We first build the three derivatives (according to  $a, b$  and  $c$ ). We then use *nullable* to find out whether the resulting derivative regular expression  $r_4$  can match the empty string. If yes, we call the function *mkeps* that produces a value  $v_4$  for how  $r_4$  can match the empty string (taking into account the POSIX constraints in case there are several ways). This function is defined by the clauses:

$$\begin{aligned}
 \textit{mkeps}\ \mathbf{1} &\stackrel{\text{def}}{=} \textit{Empty} \\
 \textit{mkeps}\ (r_1 \cdot r_2) &\stackrel{\text{def}}{=} \textit{Seq}\ (\textit{mkeps}\ r_1)\ (\textit{mkeps}\ r_2) \\
 \textit{mkeps}\ (r_1 + r_2) &\stackrel{\text{def}}{=} \textit{if nullable}\ r_1\ \textit{then Left}\ (\textit{mkeps}\ r_1)\ \textit{else Right}\ (\textit{mkeps}\ r_2) \\
 \textit{mkeps}\ (r^*) &\stackrel{\text{def}}{=} \textit{Stars}\ []
 \end{aligned}$$

Note that this function needs only to be partially defined, namely only for regular expressions that are nullable. In case *nullable* fails, the string  $[a, b, c]$  cannot be matched

by  $r_1$  and the null value *None* is returned. Note also how this function makes some subtle choices leading to a POSIX value: for example if an alternative regular expression, say  $r_1 + r_2$ , can match the empty string and furthermore  $r_1$  can match the empty string, then we return a *Left*-value. The *Right*-value will only be returned if  $r_1$  cannot match the empty string.

The most interesting idea from Sulzmann and Lu [16] is the construction of a value for how  $r_1$  can match the string  $[a, b, c]$  from the value how the last derivative,  $r_4$  in Fig. 1, can match the empty string. Sulzmann and Lu achieve this by stepwise “injecting back” the characters into the values thus inverting the operation of building derivatives, but on the level of values. The corresponding function, called *inj*, takes three arguments, a regular expression, a character and a value. For example in the first (or right-most) *inj*-step in Fig. 1 the regular expression  $r_3$ , the character  $c$  from the last derivative step and  $v_4$ , which is the value corresponding to the derivative regular expression  $r_4$ . The result is the new value  $v_3$ . The final result of the algorithm is the value  $v_1$ . The *inj* function is defined by recursion on regular expressions and by analysing the shape of values (corresponding to the derivative regular expressions).

$$\begin{array}{ll}
(1) & \text{inj } d \ c \ (\text{Empty}) \stackrel{\text{def}}{=} \text{Char } d \\
(2) & \text{inj } (r_1 + r_2) \ c \ (\text{Left } v_1) \stackrel{\text{def}}{=} \text{Left } (\text{inj } r_1 \ c \ v_1) \\
(3) & \text{inj } (r_1 + r_2) \ c \ (\text{Right } v_2) \stackrel{\text{def}}{=} \text{Right } (\text{inj } r_2 \ c \ v_2) \\
(4) & \text{inj } (r_1 \cdot r_2) \ c \ (\text{Seq } v_1 \ v_2) \stackrel{\text{def}}{=} \text{Seq } (\text{inj } r_1 \ c \ v_1) \ v_2 \\
(5) & \text{inj } (r_1 \cdot r_2) \ c \ (\text{Left } (\text{Seq } v_1 \ v_2)) \stackrel{\text{def}}{=} \text{Seq } (\text{inj } r_1 \ c \ v_1) \ v_2 \\
(6) & \text{inj } (r_1 \cdot r_2) \ c \ (\text{Right } v_2) \stackrel{\text{def}}{=} \text{Seq } (\text{mkeps } r_1) \ (\text{inj } r_2 \ c \ v_2) \\
(7) & \text{inj } (r^*) \ c \ (\text{Seq } v \ (\text{Stars } vs)) \stackrel{\text{def}}{=} \text{Stars } (\text{inj } r \ c \ v :: vs)
\end{array}$$

To better understand what is going on in this definition it might be instructive to look first at the three sequence cases (clauses (4) – (6)). In each case we need to construct an “injected value” for  $r_1 \cdot r_2$ . This must be a value of the form *Seq*  $\_ \_$ . Recall the clause of the *derivative*-function for sequence regular expressions:

$$(r_1 \cdot r_2) \setminus c \stackrel{\text{def}}{=} \text{if nullable } r_1 \text{ then } (r_1 \setminus c) \cdot r_2 + (r_2 \setminus c) \text{ else } (r_1 \setminus c) \cdot r_2$$

Consider first the *else*-branch where the derivative is  $(r_1 \setminus c) \cdot r_2$ . The corresponding value must therefore be of the form *Seq*  $v_1 \ v_2$ , which matches the left-hand side in clause (4) of *inj*. In the *if*-branch the derivative is an alternative, namely  $(r_1 \setminus c) \cdot r_2 + (r_2 \setminus c)$ . This means we either have to consider a *Left*- or *Right*-value. In case of the *Left*-value we know further it must be a value for a sequence regular expression. Therefore the pattern we match in the clause (5) is *Left* (*Seq*  $v_1 \ v_2$ ), while in (6) it is just *Right*  $v_2$ . One more interesting point is in the right-hand side of clause (6): since in this case the regular expression  $r_1$  does not “contribute” to matching the string, that means it only matches the empty string, we need to call *mkeps* in order to construct a value for how  $r_1$  can match this empty string. A similar argument applies for why we can expect in the left-hand side of clause (7) that the value is of the form *Seq*  $v \ (\text{Stars } vs)$ —the derivative of a star is  $(r \setminus c) \cdot r^*$ . Finally, the reason for why we can ignore the second argument in clause (1) of *inj* is that it will only ever be called in cases where  $c = d$ , but the usual



linearity restrictions in patterns do not allow us to build this constraint explicitly into our function definition.<sup>7</sup>

The idea of the *inj*-function to “inject” a character, say *c*, into a value can be made precise by the first part of the following lemma, which shows that the underlying string of an injected value has a prepended character *c*; the second part shows that the underlying string of an *mkeys*-value is always the empty string (given the regular expression is nullable since otherwise *mkeys* might not be defined).

**Lemma 1.**

- (1) If  $v : r \setminus c$  then  $|inj\ r\ c\ v| = c :: |v|$ .
- (2) If nullable *r* then  $|mkeys\ r| = []$ .

*Proof.* Both properties are by routine inductions: the first one can, for example, be proved by induction over the definition of *derivatives*; the second by an induction on *r*. There are no interesting cases.  $\square$

Having defined the *mkeys* and *inj* function we can extend Brzozowski’s matcher so that a value is constructed (assuming the regular expression matches the string). The clauses of the Sulzmann and Lu lexer are

$$\begin{aligned} \text{lexer } r\ [] &\stackrel{\text{def}}{=} \text{if nullable } r \text{ then Some (mkeys } r) \text{ else None} \\ \text{lexer } r\ (c :: s) &\stackrel{\text{def}}{=} \text{case lexer } (r \setminus c) \text{ of} \\ &\quad \text{None} \Rightarrow \text{None} \\ &\quad | \text{Some } v \Rightarrow \text{Some (inj } r\ c\ v) \end{aligned}$$

If the regular expression does not match the string, *None* is returned. If the regular expression *does* match the string, then *Some* value is returned. One important virtue of this algorithm is that it can be implemented with ease in any functional programming language and also in Isabelle/HOL. In the remaining part of this section we prove that this algorithm is correct.

The well-known idea of POSIX matching is informally defined by some rules such as the Longest Match and Priority Rules (see Introduction); as correctly argued in [16], this needs formal specification. Sulzmann and Lu define an “ordering relation” between values and argue that there is a maximum value, as given by the derivative-based algorithm. In contrast, we shall introduce a simple inductive definition that specifies directly what a *POSIX value* is, incorporating the POSIX-specific choices into the side-conditions of our rules. Our definition is inspired by the matching relation given by Vansummeren [17]. The relation we define is ternary and written as  $(s, r) \rightarrow v$ , relating strings, regular expressions and values; the inductive rules are given in Figure 2. We can prove that given a string *s* and regular expression *r*, the POSIX value *v* is uniquely determined by  $(s, r) \rightarrow v$ .

**Theorem 1.**

- (1) If  $(s, r) \rightarrow v$  then  $s \in L(r)$  and  $|v| = s$ .
- (2) If  $(s, r) \rightarrow v$  and  $(s, r) \rightarrow v'$  then  $v = v'$ .

<sup>7</sup> Sulzmann and Lu state this clause as  $inj\ c\ c\ (\text{Empty}) \stackrel{\text{def}}{=} \text{Char } c$ , but our deviation is harmless.

$$\begin{array}{c}
\frac{}{([\ ], \mathbf{1}) \rightarrow \text{Empty}} P\mathbf{1} \quad \frac{}{([c], c) \rightarrow \text{Char } c} Pc \\
\frac{(s, r_1) \rightarrow v}{(s, r_1 + r_2) \rightarrow \text{Left } v} P+L \quad \frac{(s, r_2) \rightarrow v \quad s \notin L(r_1)}{(s, r_1 + r_2) \rightarrow \text{Right } v} P+R \\
\frac{\begin{array}{l} (s_1, r_1) \rightarrow v_1 \quad (s_2, r_2) \rightarrow v_2 \\ \nexists s_3 s_4.a. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2) \end{array}}{(s_1 @ s_2, r_1 \cdot r_2) \rightarrow \text{Seq } v_1 v_2} PS \\
\frac{}{([\ ], r^*) \rightarrow \text{Stars } []} P[] \\
\frac{\begin{array}{l} (s_1, r) \rightarrow v \quad (s_2, r^*) \rightarrow \text{Stars } vs \quad |v| \neq [] \\ \nexists s_3 s_4.a. s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r) \wedge s_4 \in L(r^*) \end{array}}{(s_1 @ s_2, r^*) \rightarrow \text{Stars } (v :: vs)} P\star
\end{array}$$

**Fig. 2.** Our inductive definition of POSIX values.

*Proof.* Both by induction on the definition of  $(s, r) \rightarrow v$ . The second part follows by a case analysis of  $(s, r) \rightarrow v'$  and the first part.  $\square$

We claim that our  $(s, r) \rightarrow v$  relation captures the idea behind the four informal POSIX rules shown in the Introduction: Consider for example the rules  $P+L$  and  $P+R$  where the POSIX value for a string and an alternative regular expression, that is  $(s, r_1 + r_2)$ , is specified—it is always a *Left*-value, *except* when the string to be matched is not in the language of  $r_1$ ; only then it is a *Right*-value (see the side-condition in  $P+R$ ). Interesting is also the rule for sequence regular expressions ( $PS$ ). The first two premises state that  $v_1$  and  $v_2$  are the POSIX values for  $(s_1, r_1)$  and  $(s_2, r_2)$  respectively. Consider now the third premise and note that the POSIX value of this rule should match the string  $s_1 @ s_2$ . According to the Longest Match Rule, we want that the  $s_1$  is the longest initial split of  $s_1 @ s_2$  such that  $s_2$  is still recognised by  $r_2$ . Let us assume, contrary to the third premise, that there *exist* an  $s_3$  and  $s_4$  such that  $s_2$  can be split up into a non-empty string  $s_3$  and a possibly empty string  $s_4$ . Moreover the longer string  $s_1 @ s_3$  can be matched by  $r_1$  and the shorter  $s_4$  can still be matched by  $r_2$ . In this case  $s_1$  would *not* be the longest initial split of  $s_1 @ s_2$  and therefore  $\text{Seq } v_1 v_2$  cannot be a POSIX value for  $(s_1 @ s_2, r_1 \cdot r_2)$ . The main point is that our side-condition ensures the Longest Match Rule is satisfied.

A similar condition is imposed on the POSIX value in the  $P\star$ -rule. Also there we want that  $s_1$  is the longest initial split of  $s_1 @ s_2$  and furthermore the corresponding value  $v$  cannot be flattened to the empty string. In effect, we require that in each “iteration” of the star, some non-empty substring needs to be “chipped” away; only in case of the empty string we accept  $\text{Stars } []$  as the POSIX value. Indeed we can show that our POSIX values are lexical values which exclude those *Stars* that contain subvalues that flatten to the empty string.

**Lemma 2.** *If  $(s, r) \rightarrow v$  then  $v \in LV r s$ .*

*Proof.* By routine induction on  $(s, r) \rightarrow v$ . □

Next is the lemma that shows the function *mkeps* calculates the POSIX value for the empty string and a nullable regular expression.

**Lemma 3.** *If nullable  $r$  then  $([], r) \rightarrow mkeps\ r$ .*

*Proof.* By routine induction on  $r$ . □

The central lemma for our POSIX relation is that the *inj*-function preserves POSIX values.

**Lemma 4.** *If  $(s, r \setminus c) \rightarrow v$  then  $(c :: s, r) \rightarrow inj\ r\ c\ v$ .*

*Proof.* By induction on  $r$ . We explain two cases.

- Case  $r = r_1 + r_2$ . There are two subcases, namely (a)  $v = Left\ v'$  and  $(s, r_1 \setminus c) \rightarrow v'$ ; and (b)  $v = Right\ v'$ ,  $s \notin L(r_1 \setminus c)$  and  $(s, r_2 \setminus c) \rightarrow v'$ . In (a) we know  $(s, r_1 \setminus c) \rightarrow v'$ , from which we can infer  $(c :: s, r_1) \rightarrow inj\ r_1\ c\ v'$  by induction hypothesis and hence  $(c :: s, r_1 + r_2) \rightarrow inj\ (r_1 + r_2)\ c\ (Left\ v')$  as needed. Similarly in subcase (b) where, however, in addition we have to use Proposition 1(2) in order to infer  $c :: s \notin L(r_1)$  from  $s \notin L(r_1 \setminus c)$ .
- Case  $r = r_1 \cdot r_2$ . There are three subcases:
  - (a)  $v = Left\ (Seq\ v_1\ v_2)$  and nullable  $r_1$
  - (b)  $v = Right\ v_1$  and nullable  $r_1$
  - (c)  $v = Seq\ v_1\ v_2$  and  $\neg$  nullable  $r_1$

For (a) we know  $(s_1, r_1 \setminus c) \rightarrow v_1$  and  $(s_2, r_2) \rightarrow v_2$  as well as

$$\nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1 \setminus c) \wedge s_4 \in L(r_2)$$

From the latter we can infer by Proposition 1(2):

$$\nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge c :: s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2)$$

We can use the induction hypothesis for  $r_1$  to obtain  $(c :: s_1, r_1) \rightarrow inj\ r_1\ c\ v_1$ . Putting this all together allows us to infer  $(c :: s_1 @ s_2, r_1 \cdot r_2) \rightarrow Seq\ (inj\ r_1\ c\ v_1)\ v_2$ . The case (c) is similar.

For (b) we know  $(s, r_2 \setminus c) \rightarrow v_1$  and  $s_1 @ s_2 \notin L((r_1 \setminus c) \cdot r_2)$ . From the former we have  $(c :: s, r_2) \rightarrow inj\ r_2\ c\ v_1$  by induction hypothesis for  $r_2$ . From the latter we can infer

$$\nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3 @ s_4 = c :: s \wedge s_3 \in L(r_1) \wedge s_4 \in L(r_2)$$

By Lemma 3 we know  $([], r_1) \rightarrow mkeps\ r_1$  holds. Putting this all together, we can conclude with  $(c :: s, r_1 \cdot r_2) \rightarrow Seq\ (mkeps\ r_1)\ (inj\ r_2\ c\ v_1)$ , as required.

Finally suppose  $r = r_1^*$ . This case is very similar to the sequence case, except that we need to also ensure that  $|inj\ r_1\ c\ v_1| \neq []$ . This follows from  $(c :: s_1, r_1) \rightarrow inj\ r_1\ c\ v_1$  (which in turn follows from  $(s_1, r_1 \setminus c) \rightarrow v_1$  and the induction hypothesis). □

With Lemma 4 in place, it is completely routine to establish that the Sulzmann and Lu lexer satisfies our specification (returning the null value *None* iff the string is not in the language of the regular expression, and returning a unique POSIX value iff the string *is* in the language):

**Theorem 2.**

- (1)  $s \notin L(r)$  if and only if  $\text{lexer } r s = \text{None}$
- (2)  $s \in L(r)$  if and only if  $\exists v. \text{lexer } r s = \text{Some } v \wedge (s, r) \rightarrow v$

*Proof.* By induction on  $s$  using Lemma 3 and 4. □

In (2) we further know by Theorem 1 that the value returned by the lexer must be unique. A simple corollary of our two theorems is:

**Corollary 1.**

- (1)  $\text{lexer } r s = \text{None}$  if and only if  $\nexists v.a. (s, r) \rightarrow v$
- (2)  $\text{lexer } r s = \text{Some } v$  if and only if  $(s, r) \rightarrow v$

This concludes our correctness proof. Note that we have not changed the algorithm of Sulzmann and Lu,<sup>8</sup> but introduced our own specification for what a correct result—a POSIX value—should be. In the next section we show that our specification coincides with another one given by Okui and Suzuki using a different technique.

## 25 Ordering of Values according to Okui and Suzuki

While in the previous section we have defined POSIX values directly in terms of a ternary relation (see inference rules in Figure 2), Sulzmann and Lu took a different approach in [16]: they introduced an ordering for values and identified POSIX values as the maximal elements. An extended version of [16] is available at the website of its first author; this includes more details of their proofs, but which are evidently not in final form yet. Unfortunately, we were not able to verify claims that their ordering has properties such as being transitive or having maximal elements.

Okui and Suzuki [12,13] described another ordering of values, which they use to establish the correctness of their automata-based algorithm for POSIX matching. Their ordering resembles some aspects of the one given by Sulzmann and Lu, but overall is quite different. To begin with, Okui and Suzuki identify POSIX values as minimal, rather than maximal, elements in their ordering. A more substantial difference is that the ordering by Okui and Suzuki uses *positions* in order to identify and compare subvalues. Positions are lists of natural numbers. This allows them to quite naturally formalise the Longest Match and Priority rules of the informal POSIX standard. Consider for example the value  $v$

$$v \stackrel{\text{def}}{=} \text{Stars } [\text{Seq } (\text{Char } x) (\text{Char } y), \text{Char } z]$$

---

<sup>8</sup> All deviations we introduced are harmless.

At position  $[0, I]$  of this value is the subvalue *Char*  $y$  and at position  $[I]$  the subvalue *Char*  $z$ . At the ‘root’ position, or empty list  $[]$ , is the whole value  $v$ . Positions such as  $[0, I, 0]$  or  $[2]$  are outside of  $v$ . If it exists, the subvalue of  $v$  at a position  $p$ , written  $v|_p$ , can be recursively defined by

$$\begin{aligned}
 v|_[] &\stackrel{\text{def}}{=} v \\
 \text{Left } v|_{0::ps} &\stackrel{\text{def}}{=} v|_ps \\
 \text{Right } v|_{I::ps} &\stackrel{\text{def}}{=} v|_ps \\
 \text{Seq } v_1 v_2|_{0::ps} &\stackrel{\text{def}}{=} v_1|_ps \\
 \text{Seq } v_1 v_2|_{I::ps} &\stackrel{\text{def}}{=} v_2|_ps \\
 \text{Stars } vs|_{n::ps} &\stackrel{\text{def}}{=} vs_{[n]}|_ps
 \end{aligned}$$

In the last clause we use Isabelle’s notation  $vs_{[n]}$  for the  $n$ th element in a list. The set of positions inside a value  $v$ , written  $\text{Pos } v$ , is given by

$$\begin{aligned}
 \text{Pos } (\text{Empty}) &\stackrel{\text{def}}{=} \{[]\} \\
 \text{Pos } (\text{Char } c) &\stackrel{\text{def}}{=} \{[]\} \\
 \text{Pos } (\text{Left } v) &\stackrel{\text{def}}{=} \{[]\} \cup \{0::ps \mid ps \in \text{Pos } v\} \\
 \text{Pos } (\text{Right } v) &\stackrel{\text{def}}{=} \{[]\} \cup \{I::ps \mid ps \in \text{Pos } v\} \\
 \text{Pos } (\text{Seq } v_1 v_2) &\stackrel{\text{def}}{=} \{[]\} \cup \{0::ps \mid ps \in \text{Pos } v_1\} \cup \{I::ps \mid ps \in \text{Pos } v_2\} \\
 \text{Pos } (\text{Stars } vs) &\stackrel{\text{def}}{=} \{[]\} \cup (\bigcup n < \text{len } vs \{n::ps \mid ps \in \text{Pos } vs_{[n]}\})
 \end{aligned}$$

whereby  $\text{len}$  in the last clause stands for the length of a list. Clearly for every position inside a value there exists a subvalue at that position.

To help understanding the ordering of Okui and Suzuki, consider again the earlier value  $v$  and compare it with the following  $w$ :

$$\begin{aligned}
 v &\stackrel{\text{def}}{=} \text{Stars } [\text{Seq } (\text{Char } x) (\text{Char } y), \text{Char } z] \\
 w &\stackrel{\text{def}}{=} \text{Stars } [\text{Char } x, \text{Char } y, \text{Char } z]
 \end{aligned}$$

Both values match the string  $xyz$ , that means if we flatten these values at their respective root position, we obtain  $xyz$ . However, at position  $[0]$ ,  $v$  matches  $xy$  whereas  $w$  matches only the shorter  $x$ . So according to the Longest Match Rule, we should prefer  $v$ , rather than  $w$  as POSIX value for string  $xyz$  (and corresponding regular expression). In order to formalise this idea, Okui and Suzuki introduce a measure for subvalues at position  $p$ , called the *norm* of  $v$  at position  $p$ . We can define this measure in Isabelle as an integer as follows

$$\|v\|_p \stackrel{\text{def}}{=} \text{if } p \in \text{Pos } v \text{ then } \text{len } |v|_p \text{ else } -1$$

where we take the length of the flattened value at position  $p$ , provided the position is inside  $v$ ; if not, then the norm is  $-1$ . The default for outside positions is crucial for the POSIX requirement of preferring a *Left*-value over a *Right*-value (if they can match the same string—see the Priority Rule from the Introduction). For this consider

$$v \stackrel{\text{def}}{=} \text{Left}(\text{Char } x) \quad \text{and} \quad w \stackrel{\text{def}}{=} \text{Right}(\text{Char } x)$$

Both values match  $x$ . At position  $[0]$  the norm of  $v$  is  $I$  (the subvalue matches  $x$ ), but the norm of  $w$  is  $-I$  (the position is outside  $w$  according to how we defined the ‘inside’ positions of *Left*- and *Right*-values). Of course at position  $[I]$ , the norms  $\|v\|_{[I]}$  and  $\|w\|_{[I]}$  are reversed, but the point is that subvalues will be analysed according to lexicographically ordered positions. According to this ordering, the position  $[0]$  takes precedence over  $[I]$  and thus also  $v$  will be preferred over  $w$ . The lexicographic ordering of positions, written  $\_ \prec_{\text{lex}} \_$ , can be conveniently formalised by three inference rules

$$\frac{}{\Box \prec_{\text{lex}} p :: ps} \quad \frac{p_1 < p_2}{p_1 :: ps_1 \prec_{\text{lex}} p_2 :: ps_2} \quad \frac{ps_1 \prec_{\text{lex}} ps_2}{p :: ps_1 \prec_{\text{lex}} p :: ps_2}$$

With the norm and lexicographic order in place, we can state the key definition of Okui and Suzuki [12]: a value  $v_1$  is *smaller at position  $p$*  than  $v_2$ , written  $v_1 \prec_p v_2$ , if and only if (i) the norm at position  $p$  is greater in  $v_1$  (that is the string  $|v_1 \downarrow_p|$  is longer than  $|v_2 \downarrow_p|$ ) and (ii) all subvalues at positions that are inside  $v_1$  or  $v_2$  and that are lexicographically smaller than  $p$ , we have the same norm, namely

$$v_1 \prec_p v_2 \stackrel{\text{def}}{=} \begin{cases} (i) & \|v_2\|_p < \|v_1\|_p \quad \text{and} \\ (ii) & \forall q \in \text{Pos } v_1 \cup \text{Pos } v_2. q \prec_{\text{lex}} p \longrightarrow \|v_1\|_q = \|v_2\|_q \end{cases}$$

The position  $p$  in this definition acts as the *first distinct position* of  $v_1$  and  $v_2$ , where both values match strings of different length [12]. Since at  $p$  the values  $v_1$  and  $v_2$  match different strings, the ordering is irreflexive. Derived from the definition above are the following two orderings:

$$v_1 \prec v_2 \stackrel{\text{def}}{=} \exists p. v_1 \prec_p v_2$$

$$v_1 \preceq v_2 \stackrel{\text{def}}{=} v_1 \prec v_2 \vee v_1 = v_2$$

While we encountered a number of obstacles for establishing properties like transitivity for the ordering of Sulzmann and Lu (and which we failed to overcome), it is relatively straightforward to establish this property for the orderings  $\_ \prec \_$  and  $\_ \preceq \_$  by Okui and Suzuki.

**Lemma 5 (Transitivity).** *If  $v_1 \prec v_2$  and  $v_2 \prec v_3$  then  $v_1 \prec v_3$ .*

*Proof.* From the assumption we obtain two positions  $p$  and  $q$ , where the values  $v_1$  and  $v_2$  (respectively  $v_2$  and  $v_3$ ) are ‘distinct’. Since  $\prec_{\text{lex}}$  is trichotomous, we need to consider three cases, namely  $p = q$ ,  $p \prec_{\text{lex}} q$  and  $q \prec_{\text{lex}} p$ . Let us look at the first case. Clearly  $\|v_2\|_p < \|v_1\|_p$  and  $\|v_3\|_p < \|v_2\|_p$  imply  $\|v_3\|_p < \|v_1\|_p$ . It remains to show that for a  $p' \in \text{Pos } v_1 \cup \text{Pos } v_3$  with  $p' \prec_{\text{lex}} p$  that  $\|v_1\|_{p'} = \|v_3\|_{p'}$  holds. Suppose  $p' \in \text{Pos } v_1$ , then we can infer from the first assumption that  $\|v_1\|_{p'} = \|v_2\|_{p'}$ . But this means that  $p'$  must be in  $\text{Pos } v_2$  too (the norm cannot be  $-I$  given  $p' \in \text{Pos } v_1$ ). Hence we can use the second assumption and infer  $\|v_2\|_{p'} = \|v_3\|_{p'}$ , which concludes this case with  $v_1 \prec v_3$ . The reasoning in the other cases is similar.  $\square$

The proof for  $\preceq$  is similar and omitted. It is also straightforward to show that  $\prec$  and  $\preceq$  are partial orders. Okui and Suzuki furthermore show that they are linear orderings for lexical values [12] of a given regular expression and given string, but we have not formalised this in Isabelle. It is not essential for our results. What we are going to show below is that for a given  $r$  and  $s$ , the orderings have a unique minimal element on the set  $LV\ r\ s$ , which is the POSIX value we defined in the previous section. We start with two properties that show how the length of a flattened value relates to the  $\prec$ -ordering.

**Proposition 4.**

- (1) *If  $v_1 \prec v_2$  then  $len\ |v_2| \leq len\ |v_1|$ .*
- (2) *If  $len\ |v_2| < len\ |v_1|$  then  $v_1 \prec v_2$ .*

Both properties follow from the definition of the ordering. Note that (2) entails that a value, say  $v_2$ , whose underlying string is a strict prefix of another flattened value, say  $v_1$ , then  $v_1$  must be smaller than  $v_2$ . For our proofs it will be useful to have the following properties—in each case the underlying strings of the compared values are the same:

**Proposition 5.**

- (1) *If  $|v_1| = |v_2|$  then  $Left\ v_1 \prec Right\ v_2$ .*
- (2) *If  $|v_1| = |v_2|$  then  $Left\ v_1 \prec Left\ v_2$  iff  $v_1 \prec v_2$*
- (3) *If  $|v_1| = |v_2|$  then  $Right\ v_1 \prec Right\ v_2$  iff  $v_1 \prec v_2$*
- (4) *If  $|v_2| = |w_2|$  then  $Seq\ v\ v_2 \prec Seq\ v\ w_2$  iff  $v_2 \prec w_2$*
- (5) *If  $|v_1| @ |v_2| = |w_1| @ |w_2|$  and  $v_1 \prec w_1$  then  $Seq\ v_1\ v_2 \prec Seq\ w_1\ w_2$*
- (6) *If  $|vs_1| = |vs_2|$  then  $Stars\ (vs\ @\ vs_1) \prec Stars\ (vs\ @\ vs_2)$  iff  $Stars\ vs_1 \prec Stars\ vs_2$*
- (7) *If  $|v_1 :: vs_1| = |v_2 :: vs_2|$  and  $v_1 \prec v_2$  then  $Stars\ (v_1 :: vs_1) \prec Stars\ (v_2 :: vs_2)$*

One might prefer that statements (4) and (5) (respectively (6) and (7)) are combined into a single *iff*-statement (like the ones for *Left* and *Right*). Unfortunately this cannot be done easily: such a single statement would require an additional assumption about the two values  $Seq\ v_1\ v_2$  and  $Seq\ w_1\ w_2$  being inhabited by the same regular expression. The complexity of the proofs involved seems to not justify such a ‘cleaner’ single statement. The statements given are just the properties that allow us to establish our theorems without any difficulty. The proofs for Proposition 5 are routine.

Next we establish how Okui and Suzuki’s orderings relate to our definition of POSIX values. Given a *POSIX* value  $v_1$  for  $r$  and  $s$ , then any other lexical value  $v_2$  in  $LV\ r\ s$  is greater or equal than  $v_1$ , namely:

**Theorem 3.** *If  $(s, r) \rightarrow v_1$  and  $v_2 \in LV\ r\ s$  then  $v_1 \preceq v_2$ .*

*Proof.* By induction on our POSIX rules. By Theorem 1 and the definition of *LV*, it is clear that  $v_1$  and  $v_2$  have the same underlying string  $s$ . The three base cases are straightforward: for example for  $v_1 = Empty$ , we have that  $v_2 \in LV\ \mathbf{1}\ \square$  must also be of the form  $v_2 = Empty$ . Therefore we have  $v_1 \preceq v_2$ . The inductive cases for  $r$  being of the form  $r_1 + r_2$  and  $r_1 \cdot r_2$  are as follows:

- Case  $P+L$  with  $(s, r_1 + r_2) \rightarrow Left\ w_1$ : In this case the value  $v_2$  is either of the form  $Left\ w_2$  or  $Right\ w_2$ . In the latter case we can immediately conclude with

$v_1 \preceq v_2$  since a *Left*-value with the same underlying string  $s$  is always smaller than a *Right*-value by Proposition 5(1). In the former case we have  $w_2 \in LV r_1 s$  and can use the induction hypothesis to infer  $w_1 \preceq w_2$ . Because  $w_1$  and  $w_2$  have the same underlying string  $s$ , we can conclude with *Left*  $w_1 \preceq w_2$  using Proposition 5(2).

- Case  $P+R$  with  $(s, r_1 + r_2) \rightarrow \text{Right } w_1$ : This case similar to the previous case, except that we additionally know  $s \notin L(r_1)$ . This is needed when  $v_2$  is of the form *Left*  $w_2$ . Since  $|v_2| = |w_2| = s$  and  $w_2 : r_1$ , we can derive a contradiction for  $s \notin L(r_1)$  using Proposition 2. So also in this case  $v_1 \preceq v_2$ .
- Case  $PS$  with  $(s_1 @ s_2, r_1 \cdot r_2) \rightarrow \text{Seq } w_1 w_2$ : We can assume  $v_2 = \text{Seq } u_1 u_2$  with  $u_1 : r_1$  and  $u_2 : r_2$ . We have  $s_1 @ s_2 = |u_1| @ |u_2|$ . By the side-condition of the  $PS$ -rule we know that either  $s_1 = |u_1|$  or that  $|u_1|$  is a strict prefix of  $s_1$ . In the latter case we can infer  $w_1 \prec u_1$  by Proposition 4(2) and from this  $v_1 \preceq v_2$  by Proposition 5(5) (as noted above  $v_1$  and  $v_2$  must have the same underlying string). In the former case we know  $u_1 \in LV r_1 s_1$  and  $u_2 \in LV r_2 s_2$ . With this we can use the induction hypotheses to infer  $w_1 \preceq u_1$  and  $w_2 \preceq u_2$ . By Proposition 5(4,5) we can again infer  $v_1 \preceq v_2$ .

The case for  $P\star$  is similar to the  $PS$ -case and omitted. □

This theorem shows that our *POSIX* value for a regular expression  $r$  and string  $s$  is in fact a minimal element of the values in  $LV r s$ . By Proposition 4(2) we also know that any value in  $LV r s'$ , with  $s'$  being a strict prefix, cannot be smaller than  $v_1$ . The next theorem shows the opposite—namely any minimal element in  $LV r s$  must be a *POSIX* value. This can be established by induction on  $r$ , but the proof can be drastically simplified by using the fact from the previous section about the existence of a *POSIX* value whenever a string  $s \in L(r)$ .

**Theorem 4.** *If  $v_1 \in LV r s$  and  $\forall v_2 \in LV r s. v_2 \not\prec v_1$  then  $(s, r) \rightarrow v_1$ .*

*Proof.* If  $v_1 \in LV r s$  then  $s \in L(r)$  by Proposition 2. Hence by Theorem 2(2) there exists a *POSIX* value  $v_P$  with  $(s, r) \rightarrow v_P$  and by Lemma 2 we also have  $v_P \in LV r s$ . By Theorem 3 we therefore have  $v_P \preceq v_1$ . If  $v_P = v_1$  then we are done. Otherwise we have  $v_P \prec v_1$ , which however contradicts the second assumption about  $v_1$  being the smallest element in  $LV r s$ . So we are done in this case too. □

From this we can also show that if  $LV r s$  is non-empty (or equivalently  $s \in L(r)$ ) then it has a unique minimal element:

**Corollary 2.** *If  $LV r s \neq \emptyset$  then  $\exists! v_{min}. v_{min} \in LV r s \wedge (\forall v \in LV r s. v_{min} \preceq v)$ .*

To sum up, we have shown that the (unique) minimal elements of the ordering by Okui and Suzuki are exactly the *POSIX* values we defined inductively in Section 24. This provides an independent confirmation that our ternary relation formalise the informal *POSIX* rules.



## 26 Bitcoded Lexing

Incremental calculation of the value. To simplify the proof we first define the function *flex* which calculates the “iterated” injection function. With this we can rewrite the lexer as

$$\text{lexer } r \ s = (\text{if nullable } (r \setminus s) \text{ then Some } (\text{flex } r \ \text{id } s \ (\text{mkeps } (r \setminus s))) \text{ else None})$$

$$\begin{aligned} \text{code } (\text{Empty}) & \stackrel{\text{def}}{=} [] \\ \text{code } (\text{Char } c) & \stackrel{\text{def}}{=} [] \\ \text{code } (\text{Left } v) & \stackrel{\text{def}}{=} Z :: \text{code } v \\ \text{code } (\text{Right } v) & \stackrel{\text{def}}{=} S :: \text{code } v \\ \text{code } (\text{Seq } v_1 \ v_2) & \stackrel{\text{def}}{=} \text{code } v_1 \ @ \ \text{code } v_2 \\ \text{code } (\text{Stars } []) & \stackrel{\text{def}}{=} [S] \\ \text{code } (\text{Stars } (v :: vs)) & \stackrel{\text{def}}{=} Z :: \text{code } v \ @ \ \text{code } (\text{Stars } vs) \end{aligned}$$

$$\begin{aligned} \text{areg} ::= & \text{AZERO} \\ & | \text{AONE } bs \\ & | \text{ACHAR } bs \ c \\ & | \text{AALT } bs \ r_1 \ r_2 \\ & | \text{ASEQ } bs \ r_1 \ r_2 \\ & | \text{ASTAR } bs \ r \end{aligned}$$

$$\begin{aligned} (\mathbf{0})^\uparrow & \stackrel{\text{def}}{=} \text{AZERO} \\ (\mathbf{1})^\uparrow & \stackrel{\text{def}}{=} \text{AONE } [] \\ (c)^\uparrow & \stackrel{\text{def}}{=} \text{ACHAR } [] \ c \\ (r_1 + r_2)^\uparrow & \stackrel{\text{def}}{=} \text{AALT } [] \ (\text{fuse } [Z] \ (r_1^\uparrow)) \ (\text{fuse } [S] \ (r_2^\uparrow)) \\ (r_1 \cdot r_2)^\uparrow & \stackrel{\text{def}}{=} \text{ASEQ } [] \ (r_1^\uparrow) \ (r_2^\uparrow) \\ (r^*)^\uparrow & \stackrel{\text{def}}{=} \text{ASTAR } [] \ (r^\uparrow) \end{aligned}$$

$$\begin{aligned} \text{AZERO}^\downarrow & \stackrel{\text{def}}{=} \mathbf{0} \\ (\text{AONE } bs)^\downarrow & \stackrel{\text{def}}{=} \mathbf{1} \\ (\text{ACHAR } bs \ c)^\downarrow & \stackrel{\text{def}}{=} c \\ (\text{AALT } bs \ r_1 \ r_2)^\downarrow & \stackrel{\text{def}}{=} (r_1^\downarrow) + (r_2^\downarrow) \\ (\text{ASEQ } bs \ r_1 \ r_2)^\downarrow & \stackrel{\text{def}}{=} (r_1^\downarrow) \cdot (r_2^\downarrow) \\ (\text{ASTAR } bs \ r)^\downarrow & \stackrel{\text{def}}{=} (r^\downarrow)^* \end{aligned}$$

Some simple facts about erase

**Lemma 6.**

$$\begin{aligned} (r \setminus a)^\downarrow &= (r^\downarrow) \setminus a \\ (r^\uparrow)^\downarrow &= r \end{aligned}$$

$nullable_b AZERO$	$\stackrel{\text{def}}{=} False$
$nullable_b (AONE bs)$	$\stackrel{\text{def}}{=} True$
$nullable_b (ACHAR bs c)$	$\stackrel{\text{def}}{=} False$
$nullable_b (AALT bs r_1 r_2)$	$\stackrel{\text{def}}{=} nullable_b r_1 \vee nullable_b r_2$
$nullable_b (ASEQ bs r_1 r_2)$	$\stackrel{\text{def}}{=} nullable_b r_1 \wedge nullable_b r_2$
$nullable_b (ASTAR bs r)$	$\stackrel{\text{def}}{=} True$
$AZERO \ll c$	$\stackrel{\text{def}}{=} AZERO$
$AONE bs \ll c$	$\stackrel{\text{def}}{=} AZERO$
$ACHAR bs d \ll c$	$\stackrel{\text{def}}{=} \text{if } c = d \text{ then } AONE bs \text{ else } AZERO$
$AALT r_1 r_2 r2.0 \ll bs$	$\stackrel{\text{def}}{=} AALT r_1 (r_2 \ll bs) (r2.0 \ll bs)$
$ASEQ r_1 r_2 r2.0 \ll bs$	$\stackrel{\text{def}}{=} \text{if } nullable_b r_2 \text{ then } AALT r_1 (ASEQ [] (r_2 \ll bs) r2.0) (\text{fuse } (mkeps_b r_2) (r2.0 \ll bs)) \text{ else } AONE bs$
$ASTAR bs r \ll c$	$\stackrel{\text{def}}{=} ASEQ bs (\text{fuse } [Z] (r \ll c)) (ASTAR [] r)$
$mkeps_b (AONE bs)$	$\stackrel{\text{def}}{=} bs$
$mkeps_b (ASEQ bs r_1 r_2)$	$\stackrel{\text{def}}{=} bs @ mkeps_b r_1 @ mkeps_b r_2$
$mkeps_b (AALT bs r_1 r_2)$	$\stackrel{\text{def}}{=} \text{if } nullable_b r_1 \text{ then } bs @ mkeps_b r_1 \text{ else } bs @ mkeps_b r_2$
$mkeps_b (ASTAR bs r)$	$\stackrel{\text{def}}{=} bs @ [S]$

If  $v : (r^\perp) \ll c$  then  $\text{retrieve } (r \ll c) v = \text{retrieve } r (\text{inj } (r^\perp) c v)$ .  
By induction on  $r$

**Theorem 5 (Main Lemma).**

If  $v : r \ll s$  then  $\text{Some } (\text{flex } r \text{ id } s v) = \text{decode } (\text{retrieve } (r^\uparrow \ll s) v) r$ .

Definition of the bitcoded lexer

$\text{lexer}_b r s \stackrel{\text{def}}{=} \text{if } nullable_b (r^\uparrow \ll s) \text{ then } \text{decode } (mkeps_b (r^\uparrow \ll s)) r \text{ else } None$

**Theorem 6.**  $\text{lexer}_b r s = \text{lexer } r s$

## 27 Optimisations

Derivatives as calculated by Brzozowski’s method are usually more complex regular expressions than the initial one; the result is that the derivative-based matching and lexing algorithms are often abysmally slow. However, various optimisations are possible, such as the simplifications of  $\mathbf{0} + r$ ,  $r + \mathbf{0}$ ,  $\mathbf{1} \cdot r$  and  $r \cdot \mathbf{1}$  to  $r$ . These simplifications can speed up the algorithms considerably, as noted in [16]. One of the advantages of having a simple specification and correctness proof is that the latter can be refined to prove the correctness of such simplification steps. While the simplification of regular expressions according to rules like

$$\mathbf{0} + r \Rightarrow r \quad r + \mathbf{0} \Rightarrow r \quad \mathbf{1} \cdot r \Rightarrow r \quad r \cdot \mathbf{1} \Rightarrow r \quad (2)$$

is well understood, there is an obstacle with the POSIX value calculation algorithm by Sulzmann and Lu: if we build a derivative regular expression and then simplify it, we will calculate a POSIX value for this simplified derivative regular expression, *not* for the original (unsimplified) derivative regular expression. Sulzmann and Lu [16] overcome this obstacle by not just calculating a simplified regular expression, but also calculating a *rectification function* that “repairs” the incorrect value.

The rectification functions can be (slightly clumsily) implemented in Isabelle/HOL as follows using some auxiliary functions:

$$\begin{aligned} F_{Right} f v &\stackrel{\text{def}}{=} Right (f v) \\ F_{Left} f v &\stackrel{\text{def}}{=} Left (f v) \\ F_{Alt} f_1 f_2 (Right v) &\stackrel{\text{def}}{=} Right (f_2 v) \\ F_{Alt} f_1 f_2 (Left v) &\stackrel{\text{def}}{=} Left (f_1 v) \\ F_{Seq1} f_1 f_2 v &\stackrel{\text{def}}{=} Seq (f_1 ()) (f_2 v) \\ F_{Seq2} f_1 f_2 v &\stackrel{\text{def}}{=} Seq (f_1 v) (f_2 ()) \\ F_{Seq} f_1 f_2 (Seq v_1 v_2) &\stackrel{\text{def}}{=} Seq (f_1 v_1) (f_2 v_2) \\ simp_{Alt} (\mathbf{0}, \_) (r_2, f_2) &\stackrel{\text{def}}{=} (r_2, F_{Right} f_2) \\ simp_{Alt} (r_1, f_1) (\mathbf{0}, \_) &\stackrel{\text{def}}{=} (r_1, F_{Left} f_1) \\ simp_{Alt} (r_1, f_1) (r_2, f_2) &\stackrel{\text{def}}{=} (r_1 + r_2, F_{Alt} f_1 f_2) \\ simp_{Seq} (\mathbf{1}, f_1) (r_2, f_2) &\stackrel{\text{def}}{=} (r_2, F_{Seq1} f_1 f_2) \\ simp_{Seq} (r_1, f_1) (\mathbf{1}, f_2) &\stackrel{\text{def}}{=} (r_1, F_{Seq2} f_1 f_2) \\ simp_{Seq} (r_1, f_1) (r_2, f_2) &\stackrel{\text{def}}{=} (r_1 \cdot r_2, F_{Seq} f_1 f_2) \end{aligned}$$

The functions  $simp_{Alt}$  and  $simp_{Seq}$  encode the simplification rules in (2) and compose the rectification functions (simplifications can occur deep inside the regular expression). The main simplification function is then

$$\begin{aligned} simp (r_1 + r_2) &\stackrel{\text{def}}{=} simp_{Alt} (simp r_1) (simp r_2) \\ simp (r_1 \cdot r_2) &\stackrel{\text{def}}{=} simp_{Seq} (simp r_1) (simp r_2) \\ simp r &\stackrel{\text{def}}{=} (r, id) \end{aligned}$$

where  $id$  stands for the identity function. The function  $simp$  returns a simplified regular expression and a corresponding rectification function. Note that we do not simplify under stars: this seems to slow down the algorithm, rather than speed it up. The optimised lexer is then given by the clauses:

$$\begin{aligned} lexer^+ r [] &\stackrel{\text{def}}{=} \text{if nullable } r \text{ then Some (mkeps } r) \text{ else None} \\ lexer^+ r (c :: s) &\stackrel{\text{def}}{=} \text{let } (r_s, f_r) = simp (r \setminus c) \text{ in} \\ &\quad \text{case } lexer^+ r_s \text{ of} \\ &\quad \quad \text{None} \Rightarrow \text{None} \\ &\quad \quad | \text{Some } v \Rightarrow \text{Some (inj } r \text{ c (} f_r v)) \end{aligned}$$

In the second clause we first calculate the derivative  $r \setminus c$  and then simplify the result. This gives us a simplified derivative  $r_s$  and a rectification function  $f_r$ . The lexer is then recursively called with the simplified derivative, but before we inject the character  $c$  into the value  $v$ , we need to rectify  $v$  (that is construct  $f_r v$ ). Before we can establish the correctness of  $lexer^+$ , we need to show that simplification preserves the language and simplification preserves our POSIX relation once the value is rectified (recall  $simp$  generates a (regular expression, rectification function) pair):

**Lemma 7.**

- (1)  $L(fst (simp r)) = L(r)$
- (2) If  $(s, fst (simp r)) \rightarrow v$  then  $(s, r) \rightarrow snd (simp r) v$ .

*Proof.* Both are by induction on  $r$ . There is no interesting case for the first statement. For the second statement, of interest are the  $r = r_1 + r_2$  and  $r = r_1 \cdot r_2$  cases. In each case we have to analyse four subcases whether  $fst (simp r_1)$  and  $fst (simp r_2)$  equals  $\mathbf{0}$  (respectively  $\mathbf{1}$ ). For example for  $r = r_1 + r_2$ , consider the subcase  $fst (simp r_1) = \mathbf{0}$  and  $fst (simp r_2) \neq \mathbf{0}$ . By assumption we know  $(s, fst (simp (r_1 + r_2))) \rightarrow v$ . From this we can infer  $(s, fst (simp r_2)) \rightarrow v$  and by IH also  $(*) (s, r_2) \rightarrow snd (simp r_2) v$ . Given  $fst (simp r_1) = \mathbf{0}$  we know  $L(fst (simp r_1)) = \emptyset$ . By the first statement  $L(r_1)$  is the empty set, meaning  $(**) s \notin L(r_1)$ . Taking  $(*)$  and  $(**)$  together gives by the  $P+R$ -rule  $(s, r_1 + r_2) \rightarrow Right (snd (simp r_2) v)$ . In turn this gives  $(s, r_1 + r_2) \rightarrow snd (simp (r_1 + r_2)) v$  as we need to show. The other cases are similar.  $\square$

We can now prove relatively straightforwardly that the optimised lexer produces the expected result:

**Theorem 7.**  $lexer^+ r s = lexer r s$

*Proof.* By induction on  $s$  generalising over  $r$ . The case  $\square$  is trivial. For the cons-case suppose the string  $s$  is of the form  $c :: s$ . By induction hypothesis we know  $lexer^+ r s = lexer r s$  holds for all  $r$  (in particular for  $r$  being the derivative  $r \setminus c$ ). Let  $r_s$  be the simplified derivative regular expression, that is  $fst (simp (r \setminus c))$ , and  $f_r$  be the rectification function, that is  $snd (simp (r \setminus c))$ . We distinguish the cases whether  $(*) s \in L(r \setminus c)$  or not. In the first case we have by Theorem 2(2) a value  $v$  so that  $lexer (r \setminus c) s = Some v$  and  $(s, r \setminus c) \rightarrow v$  hold. By Lemma 7(1) we can also infer from  $(*)$  that  $s \in L(r_s)$  holds. Hence we know by Theorem 2(2) that there exists a  $v'$  with  $lexer r_s s = Some v'$  and  $(s, r_s) \rightarrow v'$ . From the latter we know by Lemma 7(2) that  $(s, r \setminus c) \rightarrow f_r v'$  holds. By the uniqueness of the POSIX relation (Theorem 1) we can infer that  $v$  is equal to  $f_r v'$ —that is the rectification function applied to  $v'$  produces the original  $v$ . Now the case follows by the definitions of  $lexer$  and  $lexer^+$ .

In the second case where  $s \notin L(r \setminus c)$  we have that  $lexer (r \setminus c) s = None$  by Theorem 2(1). We also know by Lemma 7(1) that  $s \notin L(r_s)$ . Hence  $lexer r_s s = None$  by Theorem 2(1) and by IH then also  $lexer^+ r_s s = None$ . With this we can conclude in this case too.  $\square$

## 28 Conclusion

We have implemented the POSIX value calculation algorithm introduced by Sulzmann and Lu [16]. Our implementation is nearly identical to the original and all modifications we introduced are harmless (like our char-clause for *inj*). We have proved this algorithm to be correct, but correct according to our own specification of what POSIX values are. Our specification (inspired from work by Vansummeren [17]) appears to be much simpler than in [16] and our proofs are nearly always straightforward. We have attempted to formalise the original proof by Sulzmann and Lu [16], but we believe it contains unfillable gaps. In the online version of [16], the authors already acknowledge some small problems, but our experience suggests that there are more serious problems.

Having proved the correctness of the POSIX lexing algorithm in [16], which lessons have we learned? Well, this is a perfect example for the importance of the *right* definitions. We have (on and off) explored mechanisations as soon as first versions of [16] appeared, but have made little progress with turning the relatively detailed proof sketch in [16] into a formalisable proof. Having seen [17] and adapted the POSIX definition given there for the algorithm by Sulzmann and Lu made all the difference: the proofs, as said, are nearly straightforward. The question remains whether the original proof idea of [16], potentially using our result as a stepping stone, can be made to work? Alas, we really do not know despite considerable effort.

Closely related to our work is an automata-based lexer formalised by Nipkow [11]. This lexer also splits up strings into longest initial substrings, but Nipkow's algorithm is not completely computational. The algorithm by Sulzmann and Lu, in contrast, can be implemented with ease in any functional language. A bespoke lexer for the Imp-language is formalised in Coq as part of the Software Foundations book by Pierce et al [15]. The disadvantage of such bespoke lexers is that they do not generalise easily to more advanced features. Our formalisation is available from the Archive of Formal Proofs [2] under <http://www.isa-afp.org/entries/Posix-Lexing.shtml>.

**Acknowledgements:** We are very grateful to Martin Sulzmann for his comments on our work and moreover for patiently explaining to us the details in [16]. We also received very helpful comments from James Cheney and anonymous referees.

## References

1. The Open Group Base Specification Issue 6 IEEE Std 1003.1 2004 Edition, 2004. [http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd\\_chap09.html](http://pubs.opengroup.org/onlinepubs/009695399/basedefs/xbd_chap09.html).
2. F. Ausaf, R. Dyckhoff, and C. Urban. POSIX Lexing with Derivatives of Regular Expressions. *Archive of Formal Proofs*, 2016. <http://www.isa-afp.org/entries/Posix-Lexing.shtml>, Formal proof development.
3. F. Ausaf, R. Dyckhoff, and C. Urban. POSIX Lexing with Derivatives of Regular Expressions (Proof Pearl). In *Proc. of the 7th International Conference on Interactive Theorem Proving (ITP)*, volume 9807 of *LNCS*, pages 69–86, 2016.
4. J. A. Brzozowski. Derivatives of Regular Expressions. *Journal of the ACM*, 11(4):481–494, 1964.
5. T. Coquand and V. Siles. A Decision Procedure for Regular Expression Equivalence in Type Theory. In *Proc. of the 1st International Conference on Certified Programs and Proofs (CPP)*, volume 7086 of *LNCS*, pages 119–134, 2011.

6. A. Frisch and L. Cardelli. Greedy Regular Expression Matching. In *Proc. of the 31st International Conference on Automata, Languages and Programming (ICALP)*, volume 3142 of *LNCS*, pages 618–629, 2004.
7. N. B. B. Grathwohl, F. Henglein, and U. T. Rasmussen. A Crash-Course in Regular Expression Parsing and Regular Expressions as Types. Technical report, University of Copenhagen, 2014.
8. H. Hosoya, J. Vouillon, and B. C. Pierce. Regular Expression Types for XML. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 27(1):46–90, 2005.
9. A. Krauss and T. Nipkow. Proof Pearl: Regular Expression Equivalence and Relation Algebra. *Journal of Automated Reasoning*, 49:95–106, 2012.
10. C. Kuklewicz. Regex Posix. [https://wiki.haskell.org/Regex\\_Posix](https://wiki.haskell.org/Regex_Posix).
11. T. Nipkow. Verified Lexical Analysis. In *Proc. of the 11th International Conference on Theorem Proving in Higher Order Logics (TPHOLs)*, volume 1479 of *LNCS*, pages 1–15, 1998.
12. S. Okui and T. Suzuki. Disambiguation in Regular Expression Matching via Position Automata with Augmented Transitions. In *Proc. of the 15th International Conference on Implementation and Application of Automata (CIAA)*, volume 6482 of *LNCS*, pages 231–240, 2010.
13. S. Okui and T. Suzuki. Disambiguation in Regular Expression Matching via Position Automata with Augmented Transitions. Technical report, University of Aizu, 2013.
14. S. Owens and K. Slind. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation*, 21(4):377–409, 2008.
15. B. C. Pierce, C. Casinghino, M. Gaboardi, M. Greenberg, C. Hrițcu, V. Sjöberg, and B. Yorgey. *Software Foundations*. Electronic textbook, 2015. <http://www.cis.upenn.edu/~bcpierce/sf>.
16. M. Sulzmann and K. Lu. POSIX Regular Expression Parsing with Derivatives. In *Proc. of the 12th International Conference on Functional and Logic Programming (FLOPS)*, volume 8475 of *LNCS*, pages 203–220, 2014.
17. S. Vansummeren. Type Inference for Unique Pattern Matching. *ACM Transactions on Programming Languages and Systems*, 28(3):389–428, 2006.