# POSIX Lexing with Derivatives of Regular Expressions

Fahad Ausaf[1], Roy Dyckhoff[2], and Christian Urban[3]

[1] King's College London
`fahad.ausaf@icloud.com`
[2] University of St Andrews
`roy.dyckhoff@st-andrews.ac.uk`
[3] 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

**theory** *SizeBound*
  **imports** *Lexer*
**begin**

---

⋆ This paper is a revised and expanded version of [2]. Compared with that paper we give a second definition for POSIX values introduced by Okui Suzuki [10,11] 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 [13]. 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   Bit-Encodings

**datatype** *bit = Z | S*

**fun** *code :: val ⇒ 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 ⇒ val ⇒ val*
**where**
  *Stars__add v (Stars vs) = Stars (v # vs)*

**function**
  *decode′ :: bit list ⇒ rexp ⇒ (val * bit list)*
**where**
  *decode′ ds ZERO = (Void, [])*
| *decode′ ds ONE = (Void, ds)*
| *decode′ ds (CH 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* $\stackrel{def}{=}$ (*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** ⊨ *v : r*
  **shows** *decode* (*code v*) *r* = *Some v*
  **using** *assms* **unfolding** *decode__def*
  **by** (*smt append__Nil2 decode′__code old.prod.case*)

## 2   Annotated Regular Expressions

**datatype** *arexp* =
 *AZERO*
| *AONE bit list*
| *ACHAR bit list char*
| *ASEQ bit list arexp arexp*
| *AALTs bit list arexp list*
| *ASTAR bit list arexp*

**abbreviation**
  $AALT\ bs\ r1\ r2\ \overset{def}{=}\ AALTs\ bs\ [r1,\ r2]$

**fun** *asize* :: *arexp* ⇒ *nat* **where**
  *asize AZERO = 1*
| *asize (AONE cs) = 1*
| *asize (ACHAR cs c) = 1*
| *asize (AALTs cs rs) = Suc (sum__list (map asize rs))*
| *asize (ASEQ cs r1 r2) = Suc (asize r1 + asize r2)*
| *asize (ASTAR cs r) = Suc (asize r)*

**fun**
  *erase* :: *arexp* ⇒ *rexp*
**where**
  *erase AZERO = ZERO*
| *erase (AONE __) = ONE*
| *erase (ACHAR __ c) = CH c*
| *erase (AALTs __ []) = ZERO*
| *erase (AALTs __ [r]) = (erase r)*
| *erase (AALTs bs (r#rs)) = ALT (erase r) (erase (AALTs bs rs))*
| *erase (ASEQ __ r1 r2) = SEQ (erase r1) (erase r2)*
| *erase (ASTAR __ r) = STAR (erase r)*

**fun** *nonalt* :: *arexp* ⇒ *bool*
  **where**
  *nonalt (AALTs bs2 rs) = False*
| *nonalt r = True*

**fun** *good* :: *arexp* ⇒ *bool* **where**
  *good AZERO = False*
| *good (AONE cs) = True*
| *good (ACHAR cs c) = True*
| *good (AALTs cs []) = False*
| *good (AALTs cs [r]) = False*
| *good (AALTs cs (r1#r2#rs)) = ($\forall\,r' \in$ set (r1#r2#rs). good r' ∧ nonalt r')*
| *good (ASEQ __ AZERO __) = False*
| *good (ASEQ __ (AONE __) __) = False*
| *good (ASEQ __ __ AZERO) = False*
| *good (ASEQ cs r1 r2) = (good r1 ∧ good r2)*
| *good (ASTAR cs r) = True*

**fun** *fuse* :: *bit list* ⇒ *arexp* ⇒ *arexp* **where**
  *fuse bs AZERO = AZERO*
| *fuse bs (AONE cs) = AONE (bs @ cs)*
| *fuse bs (ACHAR cs c) = ACHAR (bs @ cs) c*
| *fuse bs (AALTs cs rs) = AALTs (bs @ cs) rs*
| *fuse bs (ASEQ cs r1 r2) = ASEQ (bs @ cs) r1 r2*
| *fuse bs (ASTAR cs r) = ASTAR (bs @ cs) r*

**lemma** *fuse___append*:
  **shows** *fuse (bs1 @ bs2) r = fuse bs1 (fuse bs2 r)*
  **apply**(*induct r*)
  **apply**(*auto*)
  **done**

**fun** *intern* :: *rexp* ⇒ *arexp* **where**
  *intern ZERO = AZERO*
| *intern ONE = AONE []*
| *intern (CH 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* ⇒ *val* ⇒ *bit list* **where**
  *retrieve (AONE bs) Void = bs*
| *retrieve (ACHAR bs c) (Char d) = bs*
| *retrieve (AALTs bs [r]) v = bs @ retrieve r v*
| *retrieve (AALTs bs (r#rs)) (Left v) = bs @ retrieve r v*
| *retrieve (AALTs bs (r#rs)) (Right v) = bs @ retrieve (AALTs [] rs) 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] @ retrieve r v @ retrieve (ASTAR [] r) (Stars vs)*

**fun**
 *bnullable* :: *arexp* ⇒ *bool*
**where**

$bnullable\ (AZERO) = False$
$|\ bnullable\ (AONE\ bs) = True$
$|\ bnullable\ (ACHAR\ bs\ c) = False$
$|\ bnullable\ (AALTs\ bs\ rs) = (\exists\, r \in set\ rs.\ bnullable\ r)$
$|\ bnullable\ (ASEQ\ bs\ r1\ r2) = (bnullable\ r1 \land bnullable\ r2)$
$|\ bnullable\ (ASTAR\ bs\ r) = True$

**fun**
  $bmkeps :: arexp \Rightarrow bit\ list$
**where**
  $bmkeps(AONE\ bs) = bs$
$|\ bmkeps(ASEQ\ bs\ r1\ r2) = bs\ @\ (bmkeps\ r1)\ @\ (bmkeps\ r2)$
$|\ bmkeps(AALTs\ bs\ [r]) = bs\ @\ (bmkeps\ r)$
$|\ bmkeps(AALTs\ bs\ (r\#rs)) = (if\ bnullable(r)\ then\ bs\ @\ (bmkeps\ r)\ else\ (bmkeps$
$(AALTs\ bs\ rs)))$
$|\ bmkeps(ASTAR\ bs\ r) = bs\ @\ [S]$

**fun**
 $bder :: char \Rightarrow arexp \Rightarrow arexp$
**where**
  $bder\ c\ (AZERO) = AZERO$
$|\ bder\ c\ (AONE\ bs) = AZERO$
$|\ bder\ c\ (ACHAR\ bs\ d) = (if\ c = d\ then\ AONE\ bs\ else\ AZERO)$
$|\ bder\ c\ (AALTs\ bs\ rs) = AALTs\ bs\ (map\ (bder\ c)\ rs)$
$|\ bder\ c\ (ASEQ\ bs\ r1\ r2) =$
    $(if\ bnullable\ r1$
      $then\ AALT\ bs\ (ASEQ\ []\ (bder\ c\ r1)\ r2)\ (fuse\ (bmkeps\ r1)\ (bder\ c\ r2))$
      $else\ ASEQ\ bs\ (bder\ c\ r1)\ r2)$
$|\ bder\ c\ (ASTAR\ bs\ r) = ASEQ\ bs\ (fuse\ [Z]\ (bder\ c\ r))\ (ASTAR\ []\ r)$

**fun**
  $bders :: arexp \Rightarrow string \Rightarrow arexp$
**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 rule*: *erase.induct*)
**apply**(*simp__all*)
**done**

**lemma** *erase__fuse*:
  **shows** *erase (fuse bs r) = erase r*
  **apply**(*induct r rule*: *erase.induct*)
  **apply**(*simp__all*)
  **done**

**thm** *Posix.induct*

**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 rule*: *erase.induct*)
  **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** $\forall\, v \in set\ vs. \models v : r \land code\ v = retrieve\ (intern\ r)\ v$
  **shows** *code (Stars vs) = retrieve (ASTAR [] (intern r)) (Stars vs)*
  **using** *assms*
  **apply**(*induct vs*)
  **apply**(*simp__all*)
  **done**


**lemma** *retrieve__fuse2*:
  **assumes** $\models v : (erase\ r)$
  **shows** *retrieve (fuse bs r) v = bs @ retrieve r v*
  **using** *assms*
  **apply**(*induct r arbitrary*: *v bs*)

       **apply**(*auto elim*: *Prf\_\_elims*)[*4*]
  **defer**
 **using** *retrieve\_\_encode\_\_STARS*
 **apply**(*auto elim*!: *Prf\_\_elims*)[*1*]
 **apply**(*case\_\_tac vs*)
  **apply**(*simp*)
 **apply**(*simp*)

 **apply**(*simp*)
 **apply**(*case\_\_tac x2a*)
  **apply**(*simp*)
 **apply**(*auto elim*!: *Prf\_\_elims*)[*1*]
 **apply**(*simp*)
  **apply**(*case\_\_tac list*)
  **apply**(*simp*)
 **apply**(*auto*)
 **apply**(*auto elim*!: *Prf\_\_elims*)[*1*]
 **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** *bnullable\_\_Hdbmkeps\_\_Hd*:
 **assumes** *bnullable a*
 **shows**  *bmkeps* (*AALTs bs* (*a* # *rs*)) = *bs* @ (*bmkeps a*)
 **using** *assms*
 **by** (*metis bmkeps.simps(3) bmkeps.simps(4) list.exhaust*)

**lemma** *r1*:
 **assumes** $\neg$ *bnullable a bnullable* (*AALTs bs rs*)
 **shows**  *bmkeps* (*AALTs bs* (*a* # *rs*)) = *bmkeps* (*AALTs bs rs*)
 **using** *assms*

**apply**(*induct rs*)
 **apply**(*auto*)
 **done**

**lemma** *r2*:
 **assumes** $x \in set\ rs\ bnullable\ x$
 **shows** *bnullable* (*AALTs bs rs*)
 **using** *assms*
 **apply**(*induct rs*)
 **apply**(*auto*)
 **done**

**lemma**  *r3*:
 **assumes** ¬ *bnullable r*
        $\exists\ x \in set\ rs.\ bnullable\ x$
 **shows** *retrieve* (*AALTs bs rs*) (*mkeps* (*erase* (*AALTs bs rs*))) =
       *retrieve* (*AALTs bs* (*r # rs*)) (*mkeps* (*erase* (*AALTs bs* (*r # rs*))))
 **using** *assms*
 **apply**(*induct rs arbitrary*: *r bs*)
 **apply**(*auto*)[*1*]
 **apply**(*auto*)
 **using** *bnullable___correctness* **apply** *blast*
 **apply**(*auto simp add*: *bnullable___correctness mkeps___nullable retrieve___fuse2*)
 **apply**(*subst retrieve___fuse2*[*symmetric*])
 **apply** (*smt bnullable.simps*(*4*) *bnullable___correctness erase.simps*(*5*) *erase.simps*(*6*)
*insert___iff list.exhaust list.set*(*2*) *mkeps.simps*(*3*) *mkeps___nullable*)
 **apply**(*simp*)
 **apply**(*case___tac bnullable a*)
 **apply** (*smt append___Nil2 bnullable.simps*(*4*) *bnullable___correctness erase.simps*(*5*)
*erase.simps*(*6*) *fuse.simps*(*4*) *insert___iff list.exhaust list.set*(*2*) *mkeps.simps*(*3*)
*mkeps___nullable retrieve___fuse2*)
 **apply**(*drule___tac x=a* **in** *meta___spec*)
 **apply**(*drule___tac x=bs* **in** *meta___spec*)
 **apply**(*drule meta___mp*)
 **apply**(*simp*)
 **apply**(*drule meta___mp*)
 **apply**(*auto*)
 **apply**(*subst retrieve___fuse2*[*symmetric*])
 **apply**(*case___tac rs*)
  **apply**(*simp*)
 **apply**(*auto*)[*1*]
    **apply** (*simp add*: *bnullable___correctness*)
 **apply** (*metis append___Nil2 bnullable___correctness erase___fuse fuse.simps*(*4*)
*list.set___intros*(*1*) *mkeps.simps*(*3*) *mkeps___nullable nullable.simps*(*4*) *r2*)
  **apply** (*simp add*: *bnullable___correctness*)

**apply** (*metis append___Nil2 bnullable___correctness erase.simps*(*6*) *erase___fuse fuse.simps*(*4*) *list.set___intros*(*2*) *mkeps.simps*(*3*) *mkeps___nullable r2*)
**apply**(*simp*)
**done**

**lemma** *t*:
  **assumes** $\forall\, r \in set\ rs.\ nullable\ (erase\ r) \longrightarrow bmkeps\ r = retrieve\ r\ (mkeps\ (erase\ r))$
          *nullable* (*erase* (*AALTs bs rs*))
  **shows**  *bmkeps* (*AALTs bs rs*) = *retrieve* (*AALTs bs rs*) (*mkeps* (*erase* (*AALTs bs rs*)))
  **using** *assms*
  **apply**(*induct rs arbitrary*: *bs*)
   **apply**(*simp*)
  **apply**(*auto simp add*: *bnullable___correctness*)
   **apply**(*case___tac rs*)
    **apply**(*auto simp add*: *bnullable___correctness*)[*2*]
   **apply**(*subst r1*)
    **apply**(*simp*)
   **apply**(*rule r2*)
    **apply**(*assumption*)
   **apply**(*simp*)
   **apply**(*drule___tac x=bs* **in** *meta___spec*)
   **apply**(*drule meta___mp*)
    **apply**(*auto*)[*1*]
   **prefer** *2*
  **apply**(*case___tac bnullable a*)
    **apply**(*subst bnullable___Hdbmkeps___Hd*)
     **apply** *blast*
    **apply**(*subgoal___tac nullable* (*erase a*))
   **prefer** *2*
  **using** *bnullable___correctness* **apply** *blast*
   **apply** (*metis* (*no___types, lifting*) *erase.simps*(*5*) *erase.simps*(*6*) *list.exhaust mkeps.simps*(*3*) *retrieve.simps*(*3*) *retrieve.simps*(*4*))
  **apply**(*subst r1*)
    **apply**(*simp*)
  **using** *r2* **apply** *blast*
  **apply**(*drule___tac x=bs* **in** *meta___spec*)
   **apply**(*drule meta___mp*)
    **apply**(*auto*)[*1*]
   **apply**(*simp*)
  **using** *r3* **apply** *blast*
  **apply**(*auto*)
  **using** *r3* **by** *blast*

**lemma** *bmkeps___retrieve*:
  **assumes** *nullable* (*erase r*)
  **shows** *bmkeps r = retrieve r* (*mkeps* (*erase r*))
  **using** *assms*
  **apply**(*induct r*)
        **apply**(*simp*)
      **apply**(*simp*)
     **apply**(*simp*)
   **apply**(*simp*)
   **defer**
  **apply**(*simp*)
  **apply**(*rule t*)
   **apply**(*auto*)
  **done**

**lemma** *bder___retrieve*:
  **assumes** $\models$ *v* : *der c* (*erase r*)
  **shows** *retrieve* (*bder c r*) *v = retrieve r* (*injval* (*erase r*) *c v*)
  **using** *assms*
  **apply**(*induct r arbitrary*: *v rule*: *erase.induct*)
        **apply**(*simp*)
         **apply**(*erule Prf___elims*)
        **apply**(*simp*)
        **apply**(*erule Prf___elims*)
        **apply**(*simp*)
     **apply**(*case___tac c = ca*)
      **apply**(*simp*)
      **apply**(*erule Prf___elims*)
      **apply**(*simp*)
     **apply**(*simp*)
      **apply**(*erule Prf___elims*)
  **apply**(*simp*)
     **apply**(*erule Prf___elims*)
    **apply**(*simp*)
   **apply**(*simp*)
  **apply**(*rename___tac* $r_1$ $r_2$ *rs v*)
   **apply**(*erule Prf___elims*)
    **apply**(*simp*)
   **apply**(*simp*)
   **apply**(*case___tac rs*)
    **apply**(*simp*)
   **apply**(*simp*)
  **apply** (*smt Prf___elims(3) injval.simps(2) injval.simps(3) retrieve.simps(4)*
*retrieve.simps(5) same___append___eq*)

  **apply**(*simp*)
  **apply**(*case__tac nullable* (*erase r1*))
   **apply**(*simp*)
 **apply**(*erule Prf__elims*)
   **apply**(*subgoal__tac bnullable r1*)
  **prefer** *2*
  **using** *bnullable__correctness* **apply** *blast*
  **apply**(*simp*)
   **apply**(*erule Prf__elims*)
   **apply**(*simp*)
 **apply**(*subgoal__tac bnullable r1*)
  **prefer** *2*
  **using** *bnullable__correctness* **apply** *blast*
  **apply**(*simp*)
  **apply**(*simp add*: *retrieve__fuse2*)
  **apply**(*simp add*: *bmkeps__retrieve*)
 **apply**(*simp*)
 **apply**(*erule Prf__elims*)
 **apply**(*simp*)
  **using** *bnullable__correctness* **apply** *blast*
  **apply**(*rename__tac bs r v*)
  **apply**(*simp*)
  **apply**(*erule Prf__elims*)
    **apply**(*clarify*)
  **apply**(*erule Prf__elims*)
  **apply**(*clarify*)
  **apply**(*subst injval.simps*)
  **apply**(*simp del*: *retrieve.simps*)
  **apply**(*subst retrieve.simps*)
  **apply**(*subst retrieve.simps*)
  **apply**(*simp*)
  **apply**(*simp add*: *retrieve__fuse2*)
  **done**

 

**lemma** *MAIN__decode*:
  **assumes** $\models$ *v* : *ders s r*
  **shows** *Some* (*flex r id s v*) = *decode* (*retrieve* (*bders* (*intern r*) *s*) *v*) *r*
  **using** *assms*
**proof** (*induct s arbitrary*: *v rule*: *rev__induct*)
  **case** *Nil*
  **have** $\models$ *v* : *ders* [] *r* **by** *fact*
  **then have** $\models$ *v* : *r* **by** *simp*
  **then have** *Some v* = *decode* (*retrieve* (*intern r*) *v*) *r*

    **using** *decode__code retrieve__code* **by** *auto*
  **then show** *Some* (*flex r id* [] *v*) = *decode* (*retrieve* (*bders* (*intern r*) []) *v*) *r*
    **by** *simp*
**next**
  **case** (*snoc c s v*)
  **have** *IH*: $\bigwedge$*v.* $\models$ *v* : *ders s r* $\implies$
    *Some* (*flex r id s v*) = *decode* (*retrieve* (*bders* (*intern r*) *s*) *v*) *r* **by** *fact*
  **have** *asm*: $\models$ *v* : *ders* (*s* @ [*c*]) *r* **by** *fact*
  **then have** *asm2*: $\models$ *injval* (*ders s r*) *c v* : *ders s r*
    **by** (*simp add*: *Prf__injval ders__append*)
  **have** *Some* (*flex r id* (*s* @ [*c*]) *v*) = *Some* (*flex r id s* (*injval* (*ders s r*) *c v*))
    **by** (*simp add*: *flex__append*)
  **also have** ... = *decode* (*retrieve* (*bders* (*intern r*) *s*) (*injval* (*ders s r*) *c v*)) *r*
    **using** *asm2 IH* **by** *simp*
  **also have** ... = *decode* (*retrieve* (*bder c* (*bders* (*intern r*) *s*)) *v*) *r*
    **using** *asm* **by** (*simp__all add*: *bder__retrieve ders__append*)
  **finally show** *Some* (*flex r id* (*s* @ [*c*]) *v*) =
              *decode* (*retrieve* (*bders* (*intern r*) (*s* @ [*c*])) *v*) *r* **by** (*simp add*:
*bders__append*)
**qed**


**definition** *blex* **where**
 *blex a s* $\stackrel{def}{=}$ *if bnullable* (*bders a s*) *then Some* (*bmkeps* (*bders a s*)) *else None*


**definition** *blexer* **where**
 *blexer r s* $\stackrel{def}{=}$ *if bnullable* (*bders* (*intern r*) *s*) *then*
             *decode* (*bmkeps* (*bders* (*intern r*) *s*)) *r else None*

**lemma** *blexer__correctness*:
  **shows** *blexer r s* = *lexer r s*
**proof** −
  { **define** *bds* **where** *bds* $\stackrel{def}{=}$ *bders* (*intern r*) *s*

    **define** *ds* **where** *ds* $\stackrel{def}{=}$ *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**

**fun** *distinctBy* :: $'a$ *list* $\Rightarrow$ ($'a \Rightarrow$ $'b$) $\Rightarrow$ $'b$ *set* $\Rightarrow$ $'a$ *list*
  **where**
 *distinctBy* [] *f acc* = []
| *distinctBy* (*x#xs*) *f acc* =
   (*if* (*f x*) $\in$ *acc then distinctBy xs f acc*
   *else x* **#** (*distinctBy xs f* ({*f x*} $\cup$ *acc*)))

**fun** *flts* :: *arexp list* $\Rightarrow$ *arexp list*
  **where**
 *flts* [] = []
| *flts* (*AZERO* **#** *rs*) = *flts rs*
| *flts* ((*AALTs bs  rs1*) **#** *rs*) = (*map* (*fuse bs*) *rs1*) **@** *flts rs*
| *flts* (*r1* **#** *rs*) = *r1* **#** *flts rs*

**fun** *li* :: *bit list* $\Rightarrow$ *arexp list* $\Rightarrow$ *arexp*
  **where**
 *li __* [] = *AZERO*
| *li bs* [*a*] = *fuse bs a*
| *li bs as* = *AALTs bs as*

**fun** *bsimp__ASEQ* :: *bit list* $\Rightarrow$ *arexp* $\Rightarrow$ *arexp* $\Rightarrow$ *arexp*
  **where**
 *bsimp__ASEQ __ AZERO __* = *AZERO*
| *bsimp__ASEQ __ __ AZERO* = *AZERO*

| *bsimp___ASEQ bs1 (AONE bs2) r2 = fuse (bs1 @ bs2) r2*
| *bsimp___ASEQ bs1 r1 r2 = ASEQ  bs1 r1 r2*


**fun** *bsimp___AALTs :: bit list ⇒ arexp list ⇒ arexp*
  **where**
  *bsimp___AALTs ___ [] = AZERO*
| *bsimp___AALTs bs1 [r] = fuse bs1 r*
| *bsimp___AALTs bs1 rs = AALTs bs1 rs*


**fun** *bsimp :: arexp ⇒ arexp*
  **where**
  *bsimp (ASEQ bs1 r1 r2) = bsimp___ASEQ bs1 (bsimp r1) (bsimp r2)*
| *bsimp (AALTs bs1 rs) = bsimp___AALTs bs1 (distinctBy  (flts (map bsimp rs)) erase {} )*
| *bsimp r = r*




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

**definition** *blexer___simp* **where**
 *blexer___simp r s* $\overset{def}{=}$ *if bnullable (bders___simp (intern r) s) then*
          *decode (bmkeps (bders___simp (intern r) s)) r else None*

**export-code** *bders___simp* **in** *Scala* **module-name** *Example*

**lemma** *bders___simp___append*:
  **shows** *bders___simp r (s1 @ s2) = bders___simp (bders___simp r s1) s2*
  **apply**(*induct s1 arbitrary*: *r s2*)
   **apply**(*simp*)
  **apply**(*simp*)
  **done**

**lemma** *L___bsimp___ASEQ*:
  *L* (*SEQ* (*erase r1*) (*erase r2*)) = *L* (*erase* (*bsimp___ASEQ bs r1 r2*))
  **apply**(*induct bs r1 r2 rule*: *bsimp___ASEQ.induct*)
  **apply**(*simp___all*)
  **by** (*metis erase___fuse fuse.simps(4)*)

**lemma** *L___bsimp___AALTs*:
  *L* (*erase* (*AALTs bs rs*)) = *L* (*erase* (*bsimp___AALTs bs rs*))
  **apply**(*induct bs rs rule*: *bsimp___AALTs.induct*)
  **apply**(*simp___all add*: *erase___fuse*)
  **done**

**lemma** *L___erase___AALTs*:
  **shows** *L* (*erase* (*AALTs bs rs*)) = $\bigcup$ (*L* ' *erase* ' (*set rs*))
  **apply**(*induct rs*)
   **apply**(*simp*)
  **apply**(*simp*)
  **apply**(*case___tac rs*)
   **apply**(*simp*)
  **apply**(*simp*)
  **done**

**lemma** *L___erase___flts*:
  **shows** $\bigcup$ (*L* ' *erase* ' (*set* (*flts rs*))) = $\bigcup$ (*L* ' *erase* ' (*set rs*))
  **apply**(*induct rs rule*: *flts.induct*)
        **apply**(*simp___all*)
  **apply**(*auto*)
  **using** *L___erase___AALTs erase___fuse* **apply** *auto*[*1*]
  **by** (*simp add*: *L___erase___AALTs erase___fuse*)

**lemma** *L___erase___dB___acc*:
  **shows** ( $\bigcup$(*L* ' *acc*) $\cup$ ( $\bigcup$ (*L* ' *erase* ' (*set* (*distinctBy rs erase acc*) ) ) ))
= $\bigcup$(*L* ' *acc*) $\cup$ $\bigcup$ (*L* ' *erase* ' (*set rs*))
  **apply**(*induction rs arbitrary*: *acc*)
   **apply** *simp*
  **apply** *simp*
  **by** (*smt* (*z3*) *SUP___absorb UN___insert sup___assoc sup___commute*)

**lemma** *L___erase___dB*:
  **shows** ( $\bigcup$ (*L* ' *erase* ' (*set* (*distinctBy rs erase* {}) ) ) ) = $\bigcup$ (*L* ' *erase* '
(*set rs*))
  **by** (*metis L___erase___dB___acc Un___commute Union___image___empty*)

**lemma** *L___bsimp___erase*:

**shows** $L\ (erase\ r) = L\ (erase\ (bsimp\ r))$
**apply**($induct\ r$)
**apply**($simp$)
**apply**($simp$)
**apply**($simp$)
**apply**($auto\ simp\ add$: $Sequ\_\_def$)[1]
**apply**($subst\ L\_\_bsimp\_\_ASEQ[symmetric]$)
**apply**($auto\ simp\ add$: $Sequ\_\_def$)[1]
**apply**($subst\ (asm)\ \ L\_\_bsimp\_\_ASEQ[symmetric]$)
**apply**($auto\ simp\ add$: $Sequ\_\_def$)[1]
 **apply**($simp$)
 **apply**($subst\ L\_\_bsimp\_\_AALTs[symmetric]$)
 **defer**
 **apply**($simp$)
**apply**($subst\ (2)L\_\_erase\_\_AALTs$)
**apply**($subst\ L\_\_erase\_\_dB$)
**apply**($subst\ L\_\_erase\_\_flts$)
**apply**($auto$)
 **apply** ($simp\ add$: $L\_\_erase\_\_AALTs$)
 **using** $L\_\_erase\_\_AALTs$ **by** $blast$

**lemma** $bsimp\_\_ASEQ0$:
  **shows** $bsimp\_\_ASEQ\ bs\ r1\ AZERO = AZERO$
  **apply**($induct\ r1$)
  **apply**($auto$)
  **done**

**lemma** $bsimp\_\_ASEQ1$:
  **assumes** $r1 \neq AZERO\ r2 \neq AZERO\ \forall\ bs.\ r1 \neq AONE\ bs$
  **shows** $bsimp\_\_ASEQ\ bs\ r1\ r2 = ASEQ\ bs\ r1\ r2$
  **using** $assms$
  **apply**($induct\ bs\ r1\ r2\ rule$: $bsimp\_\_ASEQ.induct$)
  **apply**($auto$)
  **done**

**lemma** $bsimp\_\_ASEQ2$:
  **shows** $bsimp\_\_ASEQ\ bs\ (AONE\ bs1)\ r2 = fuse\ (bs\ @\ bs1)\ r2$
  **apply**($induct\ r2$)
  **apply**($auto$)
  **done**

**lemma** $L\_\_bders\_\_simp$:

**shows** $L$ (*erase* (*bders___simp r s*)) = $L$ (*erase* (*bders r s*))
**apply**(*induct s arbitrary*: *r rule*: *rev___induct*)
 **apply**(*simp*)
**apply**(*simp*)
**apply**(*simp add*: *ders___append*)
**apply**(*simp add*: *bders___simp___append*)
**apply**(*simp add*: *L___bsimp___erase*[*symmetric*])
**by** (*simp add*: *der___correctness*)


**lemma** *b2*:
 **assumes** *bnullable r*
 **shows** *bmkeps* (*fuse bs r*) = *bs* @ *bmkeps r*
  **by** (*simp add*: *assms bmkeps___retrieve bnullable___correctness erase___fuse*
*mkeps___nullable retrieve___fuse2*)


**lemma** *b4*:
 **shows** *bnullable* (*bders___simp r s*) = *bnullable* (*bders r s*)
 **by** (*metis L___bders___simp bnullable___correctness lexer.simps*(*1*) *lexer___correct___None*
*option.distinct*(*1*))


**lemma** *qq1*:
 **assumes** $\exists\, r \in$ *set rs. bnullable r*
 **shows** *bmkeps* (*AALTs bs* (*rs* @ *rs1*)) = *bmkeps* (*AALTs bs rs*)
 **using** *assms*
 **apply**(*induct rs arbitrary*: *rs1 bs*)
 **apply**(*simp*)
 **apply**(*simp*)
  **by** (*metis Nil___is___append___conv bmkeps.simps*(*4*) *neq___Nil___conv bnul-*
*lable___Hdbmkeps___Hd split___list___last*)

**lemma** *qq2*:
 **assumes** $\forall\, r \in$ *set rs.* $\neg$ *bnullable r* $\exists\, r \in$ *set rs1. bnullable r*
 **shows** *bmkeps* (*AALTs bs* (*rs* @ *rs1*)) = *bmkeps* (*AALTs bs rs1*)
 **using** *assms*
 **apply**(*induct rs arbitrary*: *rs1 bs*)
 **apply**(*simp*)
 **apply**(*simp*)
 **by** (*metis append___assoc in___set___conv___decomp r1 r2*)

**lemma** *qq3*:
 **shows** *bnullable* (*AALTs bs rs*) = ($\exists\, r \in$ *set rs. bnullable r*)
 **apply**(*induct rs arbitrary*: *bs*)

**apply**(*simp*)
**apply**(*simp*)
**done**

**fun** *nonnested* :: *arexp* ⇒ *bool*
  **where**
  *nonnested* (*AALTs bs2* []) = *True*
| *nonnested* (*AALTs bs2* ((*AALTs bs1 rs1*) # *rs2*)) = *False*
| *nonnested* (*AALTs bs2* (*r* # *rs2*)) = *nonnested* (*AALTs bs2 rs2*)
| *nonnested r* = *True*

**lemma**  *k0*:
  **shows** *flts* (*r* # *rs1*) = *flts* [*r*] @ *flts rs1*
  **apply**(*induct r arbitrary*: *rs1*)
   **apply**(*auto*)
  **done**

**lemma**  *k00*:
  **shows** *flts* (*rs1* @ *rs2*) = *flts rs1* @ *flts rs2*
  **apply**(*induct rs1 arbitrary*: *rs2*)
   **apply**(*auto*)
  **by** (*metis append.assoc k0*)

**lemma**  *k0a*:
  **shows** *flts* [*AALTs bs rs*] = *map* (*fuse bs*)  *rs*
  **apply**(*simp*)
  **done**

**lemma** *bsimp__AALTs__qq*:
  **assumes** *1 < length rs*
  **shows** *bsimp__AALTs bs rs* = *AALTs bs*  *rs*
  **using**  *assms*
  **apply**(*case__tac rs*)

  **apply**(*simp*)
  **apply**(*case__tac list*)
  **apply**(*simp__all*)
  **done**


**lemma** *bbbbs1*:
  **shows** *nonalt r* $\vee$ ($\exists$ *bs rs. r* $=$ *AALTs bs rs*)
  **using** *nonalt.elims(3)* **by** *auto*


**lemma** *flts__append*:
  *flts* (*xs1* @ *xs2*) $=$ *flts xs1* @ *flts xs2*
  **apply**(*induct xs1  arbitrary*: *xs2  rule*: *rev__induct*)
  **apply**(*auto*)
  **apply**(*case__tac xs*)
  **apply**(*auto*)
  **apply**(*case__tac x*)
      **apply**(*auto*)
  **apply**(*case__tac x*)
      **apply**(*auto*)
  **done**

**fun** *nonazero* :: *arexp* $\Rightarrow$ *bool*
  **where**
  *nonazero AZERO* $=$ *False*
| *nonazero r* $=$ *True*


**lemma** *flts__single1*:
  **assumes** *nonalt r nonazero r*
  **shows** *flts* [*r*] $=$ [*r*]
  **using** *assms*
  **apply**(*induct r*)
  **apply**(*auto*)
  **done**


**lemma** *q3a*:
  **assumes** $\exists$ *r* $\in$ *set rs. bnullable r*

**shows** *bmkeps* (*AALTs bs* (*map* (*fuse bs1*) *rs*)) = *bmkeps* (*AALTs* (*bs@bs1*) *rs*)
  **using** *assms*
  **apply**(*induct rs arbitrary*: *bs bs1*)
   **apply**(*simp*)
  **apply**(*simp*)
  **apply**(*auto*)
    **apply** (*metis append___assoc b2 bnullable___correctness erase___fuse bnullable___Hdbmkeps___Hd*)
  **apply**(*case___tac bnullable a*)
  **apply** (*metis append.assoc b2 bnullable___correctness erase___fuse bnullable___Hdbmkeps___Hd*)
  **apply**(*case___tac rs*)
  **apply**(*simp*)
  **apply**(*simp*)
  **apply**(*auto*)[*1*]
   **apply** (*metis bnullable___correctness erase___fuse*)+
  **done**

**lemma** *qq4*:
  **assumes** ∃ *x*∈ *set list. bnullable x*
  **shows** ∃ *x*∈ *set* (*flts list*). *bnullable x*
  **using** *assms*
  **apply**(*induct list rule*: *flts.induct*)
        **apply**(*auto*)
  **by** (*metis UnCI bnullable___correctness erase___fuse imageI*)

**lemma** *qs3*:
  **assumes** ∃ *r* ∈ *set rs. bnullable r*
  **shows** *bmkeps* (*AALTs bs rs*) = *bmkeps* (*AALTs bs* (*flts rs*))
  **using** *assms*
  **apply**(*induct rs arbitrary*: *bs taking*: *size rule*: *measure___induct*)
  **apply**(*case___tac x*)
  **apply**(*simp*)
  **apply**(*simp*)
  **apply**(*case___tac a*)
      **apply**(*simp*)
       **apply** (*simp add*: *r1*)
      **apply**(*simp*)
      **apply** (*simp add*: *bnullable___Hdbmkeps___Hd*)
    **apply**(*simp*)
    **apply**(*case___tac flts list*)
      **apply**(*simp*)
 **apply** (*metis L___erase___AALTs L___erase___flts L___flat___Prf1 L___flat___Prf2 Prf___elims*(*1*) *bnullable___correctness erase.simps*(*4*) *mkeps___nullable r2*)

**apply**(*simp*)
 **apply** (*simp add*: *r1*)
 **prefer** *3*
 **apply**(*simp*)
 **apply** (*simp add*: *bnullable\_\_Hdbmkeps\_\_Hd*)
 **prefer** *2*
 **apply**(*simp*)
**apply**(*case\_\_tac* $\exists\, x \in set\ x52.\ bnullable\ x$)
**apply**(*case\_\_tac list*)
 **apply**(*simp*)
 **apply** (*metis b2 fuse.simps(4) q3a r2*)
 **apply**(*erule disjE*)
 **apply**(*subst qq1*)
  **apply**(*auto*)[*1*]
  **apply** (*metis bnullable\_\_correctness erase\_\_fuse*)
 **apply**(*simp*)
 **apply** (*metis b2 fuse.simps(4) q3a r2*)
 **apply**(*simp*)
 **apply**(*auto*)[*1*]
  **apply**(*subst qq1*)
   **apply** (*metis bnullable\_\_correctness erase\_\_fuse image\_\_eqI set\_\_map*)
  **apply** (*metis b2 fuse.simps(4) q3a r2*)
**apply**(*subst qq1*)
   **apply** (*metis bnullable\_\_correctness erase\_\_fuse image\_\_eqI set\_\_map*)
  **apply** (*metis b2 fuse.simps(4) q3a r2*)
 **apply**(*simp*)
 **apply**(*subst qq2*)
  **apply** (*metis bnullable\_\_correctness erase\_\_fuse imageE set\_\_map*)
 **prefer** *2*
**apply**(*case\_\_tac list*)
  **apply**(*simp*)
 **apply**(*simp*)
 **apply** (*simp add*: *qq4*)
**apply**(*simp*)
**apply**(*auto*)
 **apply**(*case\_\_tac list*)
  **apply**(*simp*)
 **apply**(*simp*)
 **apply** (*simp add*: *bnullable\_\_Hdbmkeps\_\_Hd*)
**apply**(*case\_\_tac bnullable* (*ASEQ x41 x42 x43*))
 **apply**(*case\_\_tac list*)
  **apply**(*simp*)
 **apply**(*simp*)
 **apply** (*simp add*: *bnullable\_\_Hdbmkeps\_\_Hd*)
**apply**(*simp*)

    **using** *qq4 r1 r2* **by** *auto*

  **lemma** *bder___fuse*:
    **shows** *bder c (fuse bs a) = fuse bs  (bder c a)*
    **apply**(*induct a arbitrary*: *bs c*)
       **apply**(*simp___all*)
    **done**

  **fun** *flts2* :: *char ⇒ arexp list ⇒ arexp list*
    **where**
  *flts2 ___ [] = []*
| *flts2 c (AZERO # rs) = flts2 c rs*
| *flts2 c (AONE ___ # rs) = flts2 c rs*
| *flts2 c (ACHAR bs d # rs) = (if c = d then (ACHAR bs d # flts2 c rs) else*
*flts2 c rs)*
| *flts2 c ((AALTs bs rs1) # rs) = (map (fuse bs) rs1) @ flts2 c rs*
| *flts2 c (ASEQ bs r1 r2 # rs) = (if (bnullable(r1) ∧ r2 = AZERO) then*
    *flts2 c rs*
    *else ASEQ bs r1 r2 # flts2 c rs)*
| *flts2 c (r1 # rs) = r1 # flts2 c rs*

  **lemma** *WQ1*:
    **assumes** *s ∈ L (der c r)*
    **shows** *s ∈ der c r → mkeps (ders s (der c r))*
    **using** *assms*
    **oops**

**lemma** *bder__bsimp__AALTs*:
  **shows** *bder c* (*bsimp__AALTs bs rs*) = *bsimp__AALTs bs* (*map* (*bder c*) *rs*)
  **apply**(*induct bs rs rule*: *bsimp__AALTs.induct*)
    **apply**(*simp*)
   **apply**(*simp*)
   **apply** (*simp add*: *bder__fuse*)
  **apply**(*simp*)
  **done**


**lemma**
  **assumes** *asize* (*bsimp a*) = *asize a   a* = *AALTs bs* [*AALTs bs2* [], *AZERO,*
*AONE bs3*]
  **shows** *bsimp a* = *a*
  **using** *assms*
  **apply**(*simp*)
  **oops**


**inductive** *rrewrite*:: *arexp* ⇒ *arexp* ⇒ *bool* (\_\_ ⇝ \_\_ [*99, 99*] *99*)
  **where**
   *ASEQ bs AZERO r2* ⇝ *AZERO*
  | *ASEQ bs r1 AZERO* ⇝ *AZERO*
  | *ASEQ bs* (*AONE bs1*) *r* ⇝ *fuse* (*bs@bs1*) *r*
  | *r1* ⇝ *r2* ⟹ *ASEQ bs r1 r3* ⇝ *ASEQ bs r2 r3*
  | *r3* ⇝ *r4* ⟹ *ASEQ bs r1 r3* ⇝ *ASEQ bs r1 r4*
  | *r* ⇝ *r′* ⟹ (*AALTs bs* (*rs1* @ [*r*] @ *rs2*)) ⇝ (*AALTs bs* (*rs1* @ [*r′*] @ *rs2*))

  | *AALTs bs* (*rsa@AZERO* # *rsb*) ⇝ *AALTs bs* (*rsa@rsb*)
  | *AALTs bs* (*rsa@*(*AALTs bs1 rs1*)# *rsb*) ⇝ *AALTs bs* (*rsa@*(*map* (*fuse bs1*)
*rs1*)*@rsb*)

  | *AALTs bs* (*map* (*fuse bs1*) *rs*) ⇝ *AALTs* (*bs@bs1*) *rs*

  | *AALTs* (*bs@bs1*) *rs* ⇝ *AALTs bs* (*map* (*fuse bs1*) *rs*)
  | *AALTs bs* [] ⇝ *AZERO*
  | *AALTs bs* [*r*] ⇝ *fuse bs r*

| *erase a1* = *erase a2* $\implies$ *AALTs bs* (*rsa*@[*a1*]@*rsb*@[*a2*]@*rsc*) $\rightsquigarrow$ *AALTs bs*
(*rsa*@[*a1*]@*rsb*@*rsc*)

**inductive** *rrewrites*:: *arexp* $\Rightarrow$ *arexp* $\Rightarrow$ *bool* (\_\_ $\rightsquigarrow*$ \_\_ [*100*, *100*] *100*)
  **where**
*rs1*[*intro*, *simp*]:*r* $\rightsquigarrow*$ *r*
| *rs2*[*intro*]: $[\![r1 \rightsquigarrow* r2; r2 \rightsquigarrow r3]\!] \implies r1 \rightsquigarrow* r3$

**inductive** *srewrites*:: *arexp list* $\Rightarrow$ *arexp list* $\Rightarrow$ *bool* ( \_\_ *s*$\rightsquigarrow*$ \_\_ [*100*, *100*]
*100*)
  **where**
*ss1*: [] *s*$\rightsquigarrow*$ []
|*ss2*: $[\![r \rightsquigarrow* r'; rs\ s\rightsquigarrow* rs']\!] \implies (r\#rs)\ s\rightsquigarrow* (r'\#rs')$

**lemma** *r\_\_in\_\_rstar* : *r1* $\rightsquigarrow$ *r2* $\implies$ *r1* $\rightsquigarrow*$ *r2*
  **using** *rrewrites.intros*(*1*) *rrewrites.intros*(*2*) **by** *blast*

**lemma** *real\_\_trans*:
  **assumes** *a1*: *r1* $\rightsquigarrow*$ *r2*  **and** *a2*: *r2* $\rightsquigarrow*$ *r3*
  **shows** *r1* $\rightsquigarrow*$ *r3*
  **using** *a2 a1*
  **apply**(*induct r2 r3 arbitrary*: *r1 rule*: *rrewrites.induct*)
   **apply**(*auto*)
  **done**

**lemma**  *many\_\_steps\_\_later*: $[\![r1 \rightsquigarrow r2; r2 \rightsquigarrow* r3\ ]\!] \implies r1 \rightsquigarrow* r3$
  **by** (*meson r\_\_in\_\_rstar real\_\_trans*)

**lemma** *contextrewrites1*: *r* $\rightsquigarrow*$ *r'* $\implies$ (*AALTs bs* (*r*#*rs*)) $\rightsquigarrow*$ (*AALTs bs*
(*r'*#*rs*))
  **apply**(*induct r r' rule*: *rrewrites.induct*)
   **apply** *simp*
  **by** (*metis append\_\_Cons append\_\_Nil rrewrite.intros*(*6*) *rs2*)

**lemma** *contextrewrites2*: *r* $\rightsquigarrow*$ *r'* $\implies$ (*AALTs bs* (*rs1*@[*r*]@*rs*)) $\rightsquigarrow*$ (*AALTs*
*bs* (*rs1*@[*r'*]@*rs*))
  **apply**(*induct r r' rule*: *rrewrites.induct*)
   **apply** *simp*

**using** *rrewrite.intros*(*6*) **by** *blast*

**lemma** *srewrites___alt*: *rs1 s⤳∗ rs2* $\implies$ (*AALTs bs* (*rs@rs1*)) ⤳∗ (*AALTs bs* (*rs@rs2*))

  **apply**(*induct rs1 rs2 arbitrary*: *bs rs rule*: *srewrites.induct*)
   **apply**(*rule rs1*)
  **apply**(*drule___tac x* = *bs* **in** *meta___spec*)
  **apply**(*drule___tac x* = *rsa@*[*r′*] **in** *meta___spec*)
  **apply** *simp*
  **apply**(*rule real___trans*)
   **prefer** *2*
   **apply**(*assumption*)
  **apply**(*drule contextrewrites2*)
  **apply** *auto*
  **done**

**corollary** *srewrites___alt1*: *rs1 s⤳∗ rs2* $\implies$ *AALTs bs rs1* ⤳∗ *AALTs bs rs2*
  **by** (*metis append.left___neutral srewrites___alt*)

**lemma** *star___seq*: *r1* ⤳∗ *r2* $\implies$ *ASEQ bs r1 r3* ⤳∗ *ASEQ bs r2 r3*
  **apply**(*induct r1 r2 arbitrary*: *r3 rule*: *rrewrites.induct*)
   **apply**(*rule rs1*)
  **apply**(*erule rrewrites.cases*)
   **apply**(*simp*)
   **apply**(*rule r___in___rstar*)
   **apply**(*rule rrewrite.intros*(*4*))
   **apply** *simp*
  **apply**(*rule rs2*)
   **apply**(*assumption*)
  **apply**(*rule rrewrite.intros*(*4*))
  **by** *assumption*

**lemma** *star___seq2*: *r3* ⤳∗ *r4* $\implies$ *ASEQ bs r1 r3* ⤳∗ *ASEQ bs r1 r4*
  **apply**(*induct r3 r4 arbitrary*: *r1 rule*: *rrewrites.induct*)
   **apply** *auto*
  **using** *rrewrite.intros*(*5*) **by** *blast*

**lemma** *continuous___rewrite*: [[*r1* ⤳∗ *AZERO*]] $\implies$ *ASEQ bs1 r1 r2* ⤳∗ *AZERO*

**apply**(*induction ra $\overset{def}{=}$ r1 rb $\overset{def}{=}$ AZERO arbitrary: bs1 r1 r2 rule: rrewrites.induct*)
  **apply** (*simp add: r__in__rstar rrewrite.intros(1)*)

  **by** (*meson rrewrite.intros(1) rrewrites.intros(2) star__seq*)

**lemma** *bsimp__aalts__simpcases*: *AONE bs $\rightsquigarrow*$ (bsimp (AONE bs))   AZERO $\rightsquigarrow*$ bsimp AZERO ACHAR bs c $\rightsquigarrow*$ (bsimp (ACHAR bs c))*
  **apply** (*simp add: rrewrites.intros(1)*)
  **apply** (*simp add: rrewrites.intros(1)*)
  **by** (*simp add: rrewrites.intros(1)*)

**lemma** *trivialbsimpsrewrites*: $[\![\bigwedge x.\ x \in\ set\ rs \implies x \rightsquigarrow* f\ x\ ]\!] \implies rs\ s\rightsquigarrow*$ (*map f rs*)

  **apply**(*induction rs*)
   **apply** *simp*
   **apply**(*rule ss1*)
  **by** (*metis insert__iff list.simps(15) list.simps(9) srewrites.simps*)

**lemma** *bsimp__AALTsrewrites*: *AALTs bs1 rs $\rightsquigarrow*$ bsimp__AALTs bs1 rs*
  **apply**(*induction rs*)
  **apply** *simp*
   **apply**(*rule r__in__rstar*)
   **apply**(*simp add:  rrewrite.intros(11)*)
  **apply**(*case__tac rs = Nil*)
   **apply**(*simp*)
  **using** *rrewrite.intros(12)* **apply** *auto[1]*
  **apply**(*subgoal__tac length (a#rs) > 1*)
   **apply**(*simp add: bsimp__AALTs__qq*)
  **apply**(*simp*)
  **done**

**inductive** *frewrites*:: *arexp list $\Rightarrow$ arexp list $\Rightarrow$ bool ( ___ f$\rightsquigarrow*$ ___ [100, 100] 100)*
  **where**
*fs1*: *[] f$\rightsquigarrow*$ []*
*|fs2*: $[\![rs\ f\rightsquigarrow* rs']\!] \implies (AZERO\#rs)\ f\rightsquigarrow* rs'$
*|fs3*: $[\![rs\ f\rightsquigarrow* rs']\!] \implies ((AALTs\ bs\ rs1)\ \#\ rs)\ f\rightsquigarrow* ((map\ (fuse\ bs)\ rs1)\ @\ rs')$
*|fs4*: $[\![rs\ f\rightsquigarrow* rs';nonalt\ r;\ nonazero\ r]\!] \implies (r\#rs)\ f\rightsquigarrow* (r\#rs')$

**lemma** *flts__prepend*: ⟦*nonalt a; nonazero a*⟧ ⟹ *flts* (*a#rs*) = *a* # (*flts rs*)
  **by** (*metis append__Cons append__Nil flts__single1 k00*)

**lemma** *fltsfrewrites*: *rs f⤳∗* (*flts rs*)
  **apply**(*induction rs*)
  **apply** *simp*
  **apply**(*rule fs1*)

  **apply**(*case__tac a = AZERO*)


  **using** *fs2* **apply** *auto[1]*
  **apply**(*case__tac ∃ bs rs. a = AALTs bs rs*)
  **apply**(*erule exE*)+

  **apply** (*simp add: fs3*)
  **apply**(*subst flts__prepend*)
    **apply**(*rule nonalt.elims(2)*)
  **prefer** *2*
  **thm** *nonalt.elims*

       **apply** *blast*

  **using** *bbbbs1* **apply** *blast*
       **apply**(*simp add: nonalt.simps*)+

  **apply** (*meson nonazero.elims(3)*)

  **by** (*meson fs4 nonalt.elims(3) nonazero.elims(3)*)


**lemma** *rrewrite0away*: *AALTs bs* ( *AZERO* # *rsb*) ⤳ *AALTs bs rsb*
  **by** (*metis append__Nil rrewrite.intros(7)*)


**lemma** *frewritesaalts*:*rs f⤳∗ rs′* ⟹ (*AALTs bs* (*rs1@rs*)) ⤳∗ (*AALTs bs*
(*rs1@rs′*))
  **apply**(*induct rs rs′ arbitrary: bs rs1 rule:frewrites.induct*)
    **apply**(*rule rs1*)
    **apply**(*drule__tac x = bs* **in** *meta__spec*)
  **apply**(*drule__tac x = rs1* @ [*AZERO*] **in** *meta__spec*)
    **apply**(*rule real__trans*)
     **apply** *simp*

   **using** *r___in___rstar rrewrite.intros(7)* **apply** *presburger*
    **apply**(*drule___tac x = bsa* **in** *meta___spec*)
  **apply**(*drule___tac x = rs1a @ [AALTs bs rs1]* **in** *meta___spec*)
  **apply**(*rule real___trans*)
   **apply** *simp*
  **using** *r___in___rstar rrewrite.intros(8)* **apply** *presburger*
   **apply**(*drule___tac x = bs* **in** *meta___spec*)
  **apply**(*drule___tac x = rs1@[r]* **in** *meta___spec*)
   **apply**(*rule real___trans*)
  **apply** *simp*
  **apply** *auto*
  **done**

**lemma** *fltsrewrites*:   *AALTs bs1 rs* ⤳∗ *AALTs bs1 (flts rs)*
 **apply**(*induction rs*)
  **apply** *simp*
 **apply**(*case___tac a = AZERO*)
 **apply** (*metis append___Nil flts.simps(2) many___steps___later rrewrite.intros(7)*)


  **apply**(*case___tac* ∃ *bs2 rs2. a = AALTs bs2 rs2*)
  **apply**(*erule exE*)+
  **apply**(*simp add*: *flts.simps*)
  **prefer** *2*

 **apply**(*subst flts___prepend*)

   **apply** (*meson nonalt.elims(3)*)

   **apply** (*meson nonazero.elims(3)*)
  **apply**(*subgoal___tac (a#rs) f⤳∗ (a#flts rs)*)
  **apply** (*metis append___Nil frewritesaalts*)
  **apply** (*meson fltsfrewrites fs4 nonalt.elims(3) nonazero.elims(3)*)
  **by** (*metis append___Cons append___Nil fltsfrewrites frewritesaalts k00 k0a*)

**lemma** *alts___simpalts*: ⋀*bs1 rs.* (⋀*x. x* ∈ *set rs* ⟹ *x* ⤳∗ *bsimp x*) ⟹
*AALTs bs1 rs* ⤳∗ *AALTs bs1 (map bsimp rs)*
 **apply**(*subgoal___tac rs s⤳∗ (map bsimp rs)*)
  **prefer** *2*
 **using** *trivialbsimpsrewrites* **apply** *auto[1]*
 **using** *srewrites___alt1* **by** *auto*


**lemma** *threelistsappend*: *rsa@a#rsb = (rsa@[a])@rsb*

**apply** *auto*
**done**

**fun** *distinctByAcc* :: $'a\ list \Rightarrow ('a \Rightarrow 'b) \Rightarrow 'b\ set \Rightarrow 'b\ set$
  **where**
  *distinctByAcc* [] *f acc = acc*
| *distinctByAcc* (*x#xs*) *f acc* =
    (*if* (*f x*) $\in$ *acc then distinctByAcc xs f acc*
    *else* (*distinctByAcc xs f* ({*f x*} $\cup$ *acc*)))

**lemma** *dB__single__step*: *distinctBy* (*a#rs*) *f* {} = *a # distinctBy rs f* {*f a*}
  **apply** *simp*
  **done**

**lemma** *somewhereInside*: *r* $\in$ *set rs* $\Longrightarrow$ $\exists$ *rs1 rs2. rs = rs1@*[*r*]*@rs2*
  **using** *split__list* **by** *fastforce*

**lemma** *somewhereMapInside*: *f r* $\in$ *f ' set rs* $\Longrightarrow$ $\exists$ *rs1 rs2 a. rs = rs1@*[*a*]*@rs2*
$\wedge$ *f a = f r*
  **apply** *auto*
  **by** (*metis split__list*)

**lemma** *alts__dBrewrites__withFront*:  *AALTs bs* (*rsa @ rs*) $\rightsquigarrow$* *AALTs bs*
(*rsa @ distinctBy rs erase* (*erase ' set rsa*))
  **apply**(*induction rs arbitrary: rsa*)
   **apply** *simp*
  **apply**(*drule__tac x = rsa@*[*a*] **in** *meta__spec*)
  **apply**(*subst threelistsappend*)
  **apply**(*rule real__trans*)
  **apply** *simp*
  **apply**(*case__tac a* $\in$ *set rsa*)
   **apply** *simp*
   **apply**(*drule somewhereInside*)
   **apply**(*erule exE*)+
   **apply** *simp*
  **apply**(*subgoal__tac  AALTs bs*
        (*rs1 @*
        *a #*
        *rs2 @*
        *a #*
        *distinctBy rs erase*
        (*insert* (*erase a*)
         (*erase '*
         (*set rs1* $\cup$ *set rs2*))))) $\rightsquigarrow$ *AALTs bs* (*rs1@ a # rs2 @  distinctBy*
*rs erase*

        *(insert (erase a)*
         *(erase ʻ*
         *(set rs1 ∪ set rs2)))) )*
  **prefer** *2*
  **using** *rrewrite.intros(13)* **apply** *force*
  **using** *r___in___rstar* **apply** *force*
  **apply**(*subgoal___tac erase ʻ set (rsa @ [a]) = insert (erase a) (erase ʻ set rsa))*
  **prefer** *2*

  **apply** *auto[1]*
  **apply**(*case___tac erase a ∈ erase ʻset rsa*)

  **apply** *simp*
  **apply**(*subgoal___tac AALTs bs (rsa @ a # distinctBy rs erase (insert (erase a) (erase ʻ set rsa))) ⤳*
              *AALTs bs (rsa @ distinctBy rs erase (insert (erase a) (erase ʻ set rsa))))*
  **apply** *force*
  **apply** (*smt (verit, ccfv___threshold) append___Cons append___assoc append___self___conv2 r___in___rstar rrewrite.intros(13) same___append___eq somewhereMapInside*)
  **by** *force*

**lemma** *alts___dBrewrites*: *AALTs bs rs ⤳∗ AALTs bs (distinctBy rs erase {})*
  **apply**(*induction rs*)
  **apply** *simp*
  **apply** *simp*
  **using** *alts___dBrewrites___withFront*
  **by** (*metis append___Nil dB___single___step empty___set image___empty*)

**lemma** *bsimp___rewrite*: (*rrewrites r ( bsimp r)*)
  **apply**(*induction r rule: bsimp.induct*)
    **apply** *simp*
    **apply**(*case___tac bsimp r1 = AZERO*)
     **apply** *simp*
  **using** *continuous___rewrite* **apply** *blast*
    **apply**(*case___tac ∃ bs. bsimp r1 = AONE bs*)
     **apply**(*erule exE*)

     **apply** *simp*
     **apply**(*subst bsimp___ASEQ2*)
   **apply** (*meson real___trans rrewrite.intros(3) rrewrites.intros(2) star___seq star___seq2*)
      **apply** (*smt* (*verit, best*) *bsimp___ASEQ0 bsimp___ASEQ1 real___trans rrewrite.intros(2) rs2 star___seq star___seq2*)
   **defer**
 **using** *bsimp___aalts___simpcases(2)* **apply** *blast*
 **apply** *simp*
 **apply** *simp*
 **apply** *simp*

 **apply** *auto*


 **apply**(*subgoal___tac AALTs bs1 rs ⇝∗ AALTs bs1 (map bsimp rs)*)
  **apply**(*subgoal___tac AALTs bs1 (map bsimp rs) ⇝∗ AALTs bs1 (flts (map bsimp rs))*)
 **apply**(*subgoal___tac AALTs bs1 (flts (map bsimp rs)) ⇝∗ AALTs bs1 (distinctBy (flts (map bsimp rs)) erase {})*)
  **apply**(*subgoal___tac AALTs bs1 (distinctBy (flts (map bsimp rs)) erase {}) ⇝∗ bsimp___AALTs bs1 (distinctBy (flts (map bsimp rs)) erase {} )*)


   **apply** (*meson real___trans*)

 **apply** (*meson bsimp___AALTsrewrites*)

 **apply** (*meson alts___dBrewrites*)

 **using** *fltsrewrites* **apply** *auto[1]*

 **using** *alts___simpalts* **by** *force*


**lemma** *rewritenullable*: ⟦*r1 ⇝ r2; bnullable r1* ⟧ ⟹ *bnullable r2*
 **apply**(*induction r1 r2 rule: rrewrite.induct*)
     **apply**(*simp*)+
 **apply** (*metis bnullable___correctness erase___fuse*)
   **apply** *simp*
  **apply** *simp*
  **apply** *auto[1]*
  **apply** *auto[1]*
 **apply** *auto[4]*
  **apply** (*metis UnCI bnullable___correctness erase___fuse imageI*)

    **apply** (*metis bnullable__correctness erase__fuse*)
    **apply** (*metis bnullable__correctness erase__fuse*)

  **apply** (*metis bnullable__correctness erase.simps(5) erase__fuse*)


 **by** (*smt (z3) Un__iff bnullable__correctness insert__iff list.set(2) qq3 set__append*)

**lemma** *rewrite__non__nullable*: $\llbracket r1 \leadsto r2; \neg bnullable\ r1 \rrbracket \implies \neg bnullable\ r2$
 **apply**(*induction r1 r2 rule*: *rrewrite.induct*)
      **apply** *auto*
   **apply** (*metis bnullable__correctness erase__fuse*)+
 **done**


**lemma** *rewritesnullable*: $\llbracket r1 \leadsto* r2; bnullable\ r1 \rrbracket \implies bnullable\ r2$
 **apply**(*induction r1 r2 rule*: *rrewrites.induct*)
  **apply** *simp*
 **apply**(*rule rewritenullable*)
  **apply** *simp*
 **apply** *simp*
 **done**

**lemma** *nonbnullable__lists__concat*: $\llbracket \neg (\exists r0 \in set\ rs1.\ bnullable\ r0); \neg bnullable\ r; \neg (\exists r0 \in set\ rs2.\ bnullable\ r0) \rrbracket \implies$
$\neg(\exists r0 \in (set\ (rs1@[r]@rs2)).\ bnullable\ r0\ )$
 **apply** *simp*
 **apply** *blast*
 **done**




**lemma** *nomember__bnullable*: $\llbracket \neg (\exists r0 \in set\ rs1.\ bnullable\ r0); \neg bnullable\ r; \neg (\exists r0 \in set\ rs2.\ bnullable\ r0) \rrbracket$
$\implies \neg bnullable\ (AALTs\ bs\ (rs1\ @\ [r]\ @\ rs2))$
 **using** *nonbnullable__lists__concat qq3* **by** *presburger*

**lemma** *bnullable__segment*: $bnullable\ (AALTs\ bs\ (rs1@[r]@rs2)) \implies bnullable$
$(AALTs\ bs\ rs1) \lor bnullable\ (AALTs\ bs\ rs2) \lor bnullable\ r$
 **apply**(*case__tac $\exists r0 \in set\ rs1.\ bnullable\ r0$*)

 **using** *qq3* **apply** *blast*
 **apply**(*case__tac bnullable r*)

 **apply** *blast*

**apply**(*case__tac ∃ r0 ∈ set rs2.  bnullable r0*)

**using** *bnullable.simps(4)* **apply** *presburger*
**apply**(*subgoal__tac False*)

**apply** *blast*

**using** *nomember__bnullable* **by** *blast*


**lemma** *bnullablewhichbmkeps*: ⟦*bnullable  (AALTs bs (rs1@[r]@rs2)); ¬ bnullable (AALTs bs rs1); bnullable r* ⟧
⟹ *bmkeps (AALTs bs (rs1@[r]@rs2)) = bs @ (bmkeps r)*
  **using** *qq2 bnullable__Hdbmkeps__Hd* **by** *force*

**lemma** *rrewrite__nbnullable*: ⟦ *r1 ⤳ r2 ; ¬ bnullable r1* ⟧ ⟹ *¬bnullable r2*
  **apply**(*induction rule*: *rrewrite.induct*)
            **apply** *auto[1]*
           **apply** *auto[1]*
          **apply** *auto[1]*
          **apply** (*metis bnullable__correctness erase__fuse*)
         **apply** *auto[1]*
        **apply** *auto[1]*
       **apply** *auto[1]*
      **apply** *auto[1]*
     **apply** *auto[1]*
     **apply** (*metis bnullable__correctness erase__fuse*)
    **apply** *auto[1]*
    **apply** (*metis bnullable__correctness erase__fuse*)
   **apply** *auto[1]*
   **apply** (*metis bnullable__correctness erase__fuse*)
  **apply** *auto[1]*
  **apply** *auto[1]*

  **apply** (*metis bnullable__correctness erase__fuse*)

  **by** (*meson rewrite__non__nullable rrewrite.intros(13)*)


**lemma** *spillbmkepslistr*: *bnullable (AALTs bs1 rs1)*
    ⟹ *bmkeps (AALTs bs (AALTs bs1 rs1 # rsb)) = bmkeps (AALTs bs ( map (fuse bs1) rs1 @ rsb))*

**apply**(*subst bnullable__Hdbmkeps__Hd*)

  **apply** *simp*
**by** (*metis bmkeps.simps(3) k0a list.set__intros(1) qq1 qq4 qs3*)

**lemma** *third__segment__bnullable*: $[\![$*bnullable* (*AALTs bs* (*rs1@rs2@rs3*)); ¬*bnullable* (*AALTs bs rs1*); ¬*bnullable* (*AALTs bs rs2*)$]\!]$ $\implies$
*bnullable* (*AALTs bs rs3*)

  **by** (*metis append.left__neutral append__Cons bnullable.simps(1) bnullable__segment rrewrite.intros(7) rrewrite__nbnullable*)

**lemma** *third__segment__bmkeps*: $[\![$*bnullable* (*AALTs bs* (*rs1@rs2@rs3*)); ¬*bnullable* (*AALTs bs rs1*); ¬*bnullable* (*AALTs bs rs2*)$]\!]$ $\implies$
*bmkeps* (*AALTs bs* (*rs1@rs2@rs3*) ) = *bmkeps* (*AALTs bs rs3*)
  **apply**(*subgoal__tac bnullable* (*AALTs bs rs3*))
   **apply**(*subgoal__tac* $\forall\, r \in$ *set* (*rs1@rs2*). ¬*bnullable r*)
  **apply**(*subgoal__tac bmkeps* (*AALTs bs* (*rs1@rs2@rs3*)) = *bmkeps* (*AALTs bs* ((*rs1@rs2*)*@rs3*) ))
  **apply** (*metis qq2 qq3*)

  **apply** (*metis append.assoc*)

  **apply** (*metis append.assoc in__set__conv__decomp r2 third__segment__bnullable*)

  **using** *third__segment__bnullable* **by** *blast*

**lemma** *rewrite__bmkepsalt*: $[\![$*bnullable* (*AALTs bs* (*rsa @ AALTs bs1 rs1 #* *rsb*)); *bnullable* (*AALTs bs* (*rsa @ map* (*fuse bs1*) *rs1 @ rsb*))$]\!]$
    $\implies$ *bmkeps* (*AALTs bs* (*rsa @ AALTs bs1 rs1 # rsb*)) = *bmkeps* (*AALTs bs* (*rsa @ map* (*fuse bs1*) *rs1 @ rsb*))
  **apply**(*case__tac bnullable* (*AALTs bs rsa*))

  **using** *qq1* **apply** *force*
  **apply**(*case__tac bnullable* (*AALTs bs1 rs1*))
  **apply**(*subst qq2*)

  **using** *r2* **apply** *blast*

   **apply** (*metis list.set__intros(1)*)
 **apply** (*smt* (*verit, ccfv__threshold*) *append__eq__append__conv2 list.set__intros(1) qq2 qq3 rewritenullable rrewrite.intros(8) self__append__conv2 spillbmkepslistr*)

**thm** *qq1*
 **apply**(*subgoal___tac bmkeps  (AALTs bs (rsa @ AALTs bs1 rs1 # rsb)) = bmkeps (AALTs bs rsb) )*
  **prefer** *2*

 **apply** (*metis append___Cons append___Nil bnullable.simps(1) bnullable___segment rewritenullable rrewrite.intros(11) third___segment___bmkeps*)

 **by** (*metis bnullable.simps(4) rewrite___non___nullable rrewrite.intros(10) third___segment___bmkeps*)

**lemma** *rewrite___bmkeps*: $\llbracket$ *r1* $\rightsquigarrow$ *r2*; (*bnullable r1*)$\rrbracket$ $\Longrightarrow$ *bmkeps r1 = bmkeps r2*

 **apply**(*frule rewritenullable*)
 **apply** *simp*
 **apply**(*induction r1 r2 rule*: *rrewrite.induct*)
          **apply** *simp*
 **using** *bnullable.simps(1) bnullable.simps(5)* **apply** *blast*
      **apply** (*simp add*: *b2*)
     **apply** *simp*
      **apply** *simp*
 **apply**(*frule bnullable___segment*)
      **apply**(*case___tac bnullable (AALTs bs rs1)*)
 **using** *qq1* **apply** *force*
      **apply**(*case___tac bnullable r*)
 **using** *bnullablewhichbmkeps rewritenullable* **apply** *presburger*
      **apply**(*subgoal___tac bnullable (AALTs bs rs2)*)
 **apply**(*subgoal___tac* $\neg$ *bnullable r'*)
 **apply** (*simp add*: *qq2 r1*)

 **using** *rrewrite___nbnullable* **apply** *blast*

     **apply** *blast*
    **apply** (*simp add*: *flts___append qs3*)

 **apply** (*meson rewrite___bmkepsalt*)

 **using** *bnullable.simps(4) q3a* **apply** *blast*

 **apply** (*simp add*: *q3a*)

**using** *bnullable.simps(1)* **apply** *blast*

**apply** (*simp add: b2*)

**by** (*smt* (*z3*) *Un__iff bnullable__correctness erase.simps(5) qq1 qq2 qq3 set__append*)

**lemma** *rewrites__bmkeps*: $\llbracket$ (*r1* $\rightsquigarrow* $ *r2*); (*bnullable r1*)$\rrbracket$ $\implies$ *bmkeps r1* $=$ *bmkeps r2*
  **apply**(*induction r1 r2 rule*: *rrewrites.induct*)
   **apply** *simp*
  **apply**(*subgoal__tac bnullable r2*)
  **prefer** *2*
   **apply**(*metis rewritesnullable*)
  **apply**(*subgoal__tac bmkeps r1* $=$ *bmkeps r2*)
   **prefer** *2*
   **apply** *fastforce*
  **using** *rewrite__bmkeps* **by** *presburger*

**thm** *rrewrite.intros(12)*
**lemma** *alts__rewrite__front*: $r \rightsquigarrow r' \implies$ *AALTs bs* (*r* # *rs*) $\rightsquigarrow$ *AALTs bs* (*r'* # *rs*)
  **by** (*metis append__Cons append__Nil rrewrite.intros(6)*)

**lemma** *alt__rewrite__front*: $r \rightsquigarrow r' \implies$ *AALT bs r r2* $\rightsquigarrow$ *AALT bs r' r2*
  **using** *alts__rewrite__front* **by** *blast*

**lemma** *to__zero__in__alt*: *AALT bs* (*ASEQ* [] *AZERO r*) *r2* $\rightsquigarrow$ *AALT bs AZERO r2*
  **by** (*simp add: alts__rewrite__front rrewrite.intros(1)*)

**lemma** *alt__remove0__front*: *AALT bs AZERO r* $\rightsquigarrow$ *AALTs bs* [*r*]
  **by** (*simp add: rrewrite0away*)

**lemma** *alt__rewrites__back*: *r2* $\rightsquigarrow*$ *r2'* $\implies$*AALT bs r1 r2* $\rightsquigarrow*$ *AALT bs r1 r2'*
  **apply**(*induction r2 r2' arbitrary*: *bs rule*: *rrewrites.induct*)
   **apply** *simp*
  **by** (*meson rs1 rs2 srewrites__alt1 ss1 ss2*)

**lemma** *rewrite__fuse*: *r2* $\rightsquigarrow$ *r3* $\implies$ *fuse bs r2* $\rightsquigarrow*$ *fuse bs r3*
  **apply**(*induction r2 r3 arbitrary*: *bs rule*: *rrewrite.induct*)

  **apply** *auto*

   **apply** (*simp add*: *continuous__rewrite*)

   **apply** (*simp add*: *r__in__rstar rrewrite.intros*(*2*))

   **apply** (*metis fuse__append r__in__rstar rrewrite.intros*(*3*))

 **using** *r__in__rstar star__seq* **apply** *blast*

 **using** *r__in__rstar star__seq2* **apply** *blast*

 **using** *contextrewrites2 r__in__rstar* **apply** *auto*[*1*]

   **apply** (*simp add*: *r__in__rstar rrewrite.intros*(*7*))

 **using** *rrewrite.intros*(*8*) **apply** *auto*[*1*]

 **apply** (*metis append__assoc r__in__rstar rrewrite.intros*(*9*))

 **apply** (*metis append__assoc r__in__rstar rrewrite.intros*(*10*))

 **apply** (*simp add*: *r__in__rstar rrewrite.intros*(*11*))

 **apply** (*metis fuse__append r__in__rstar rrewrite.intros*(*12*))

 **using** *rrewrite.intros*(*13*) **by** *auto*


**lemma** *rewrites__fuse*: $r2 \rightsquigarrow* r2' \implies$ (*fuse bs1 r2*) $\rightsquigarrow*$ (*fuse bs1 r2'*)
 **apply**(*induction r2 r2' arbitrary*: *bs1 rule*: *rrewrites.induct*)
  **apply** *simp*
 **by** (*meson real__trans rewrite__fuse*)

**lemma** *bder__fuse__list*: *map* (*bder c* $\circ$ *fuse bs1*) *rs1* = *map* (*fuse bs1* $\circ$ *bder c*) *rs1*
 **apply**(*induction rs1*)
 **apply** *simp*
 **by** (*simp add*: *bder__fuse*)


**lemma** *rewrite__der__altmiddle*: *bder c* (*AALTs bs* (*rsa @ AALTs bs1 rs1 # rsb*)) $\rightsquigarrow*$ *bder c* (*AALTs bs* (*rsa @ map* (*fuse bs1*) *rs1 @ rsb*))

  **apply** *simp*
  **apply**(*simp add*: *bder___fuse___list*)
 **apply**(*rule many___steps___later*)
  **apply**(*subst rrewrite.intros*(*8*))
  **apply** *simp*

  **by** *fastforce*

**lemma** *lock___step___der___removal*:
  **shows** *erase a1 = erase a2* $\implies$
$$bder\ c\ (AALTs\ bs\ (rsa\ @\ [a1]\ @\ rsb\ @\ [a2]\ @\ rsc))$$
$\rightsquigarrow*$
$$bder\ c\ (AALTs\ bs\ (rsa\ @\ [a1]\ @\ rsb\ @\ rsc))$$
  **apply**(*simp*)

  **using** *rrewrite.intros*(*13*) **by** *auto*

**lemma** *rewrite___after___der*: *r1* $\rightsquigarrow$ *r2* $\implies$ (*bder c r1*) $\rightsquigarrow*$ (*bder c r2*)
  **apply**(*induction r1 r2 arbitrary*: *c rule*: *rrewrite.induct*)

      **apply** (*simp add*: *r___in___rstar rrewrite.intros*(*1*))
  **apply** *simp*

 **apply** (*meson contextrewrites1 r___in___rstar rrewrite.intros*(*11*) *rrewrite.intros*(*2*)
*rrewrite0away rs2*)
     **apply**(*simp*)
     **apply**(*rule many___steps___later*)
      **apply**(*rule to___zero___in___alt*)
     **apply**(*rule many___steps___later*)
 **apply**(*rule alt___remove0___front*)
     **apply**(*rule many___steps___later*)
      **apply**(*rule rrewrite.intros*(*12*))
  **using** *bder___fuse fuse___append rs1* **apply** *presburger*
     **apply**(*case___tac bnullable r1*)
  **prefer** *2*
     **apply**(*subgoal___tac ¬bnullable r2*)
      **prefer** *2*
  **using** *rewrite___non___nullable* **apply** *presburger*
     **apply** *simp+*

  **using** *star___seq* **apply** *auto*[*1*]
     **apply**(*subgoal___tac bnullable r2*)
     **apply** *simp+*
  **apply**(*subgoal___tac bmkeps r1 = bmkeps r2*)
  **prefer** *2*

**using** *rewrite__bmkeps* **apply** *auto[1]*
**using** *contextrewrites1 star__seq* **apply** *auto[1]*
**using** *rewritenullable* **apply** *auto[1]*
    **apply**(*case__tac bnullable r1*)
     **apply** *simp*
    **apply**(*subgoal__tac ASEQ* [] (*bder c r1*) *r3* $\rightsquigarrow$ *ASEQ* [] (*bder c r1*) *r4*)
     **prefer** *2*
**using** *rrewrite.intros(5)* **apply** *blast*
    **apply**(*rule many__steps__later*)
     **apply**(*rule alt__rewrite__front*)
     **apply** *assumption*
**apply** (*meson alt__rewrites__back rewrites__fuse*)

    **apply** (*simp add*: *r__in__rstar rrewrite.intros(5)*)

**using** *contextrewrites2* **apply** *force*

**using** *rrewrite.intros(7)* **apply** *force*

**using** *rewrite__der__altmiddle* **apply** *auto[1]*

**apply** (*metis bder.simps(4) bder__fuse__list map__map r__in__rstar rrewrite.intros(9)*)

**apply** (*metis List.map.compositionality bder.simps(4) bder__fuse__list r__in__rstar rrewrite.intros(10)*)

**apply** (*simp add*: *r__in__rstar rrewrite.intros(11)*)

**apply** (*metis bder.simps(4) bder__bsimp__AALTs bsimp__AALTs.simps(2) bsimp__AALTsrewrites*)

**using** *lock__step__der__removal* **by** *auto*

**lemma** *rewrites__after__der*:   *r1* $\rightsquigarrow*$ *r2* $\implies$ (*bder c r1*) $\rightsquigarrow*$ (*bder c r2*)
 **apply**(*induction r1 r2 rule*: *rrewrites.induct*)
  **apply**(*rule rs1*)
 **by** (*meson real__trans rewrite__after__der*)

**lemma** *central*: (*bders r s*) $\rightsquigarrow*$ (*bders__simp r s*)

**apply**(*induct s arbitrary*: *r rule*: *rev__induct*)

 **apply** *simp*
 **apply**(*subst bders__append*)
 **apply**(*subst bders__simp__append*)
 **by** (*metis bders.simps(1) bders.simps(2) bders__simp.simps(1) bders__simp.simps(2) bsimp__rewrite real__trans rewrites__after__der*)

**thm** *arexp.induct*

**lemma** *quasi__main*: *bnullable* (*bders r s*) $\Longrightarrow$ *bmkeps* (*bders r s*) = *bmkeps* (*bders__simp r s*)
 **using** *central rewrites__bmkeps* **by** *blast*

**theorem** *main__main*: *blexer r s* = *blexer__simp r s*
 **by** (*simp add*: *b4 blexer__def blexer__simp__def quasi__main*)

**theorem** *blexersimp__correctness*: *blexer__simp r s*= *lexer r s*
 **using** *blexer__correctness main__main* **by** *auto*

**unused-thms**

**end**

## 3   Introduction

This works builds on previous work by Ausaf and Urban using regular expression'd bit-coded derivatives to do lexing that is both fast and satisfied the POSIX specification. In their work, a bit-coded algorithm introduced by Sulzmann and Lu was formally verified in Isabelle, by a very clever use of flex function and retrieve to carefully mimic the way a value is built up by the injection funciton.
    In the previous work, Ausaf and Urban established the below equality:

**Lemma 1.** *If* $v : (r^{\downarrow})\backslash c$ *then retrieve* $(r\backslash\!\backslash c)$ $v$ = *retrieve* $r$ $(inj$ $(r^{\downarrow})$ $c$ $v)$.

    This lemma links the derivative of a bit-coded regular expression with the regular expression itself before the derivative.
    Brzozowski [3] introduced the notion of the *derivative* $r\backslash c$ 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 $s \in L(r \backslash c)$. 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 [12]. Another one in Isabelle/HOL is part of the work by Krauss and Nipkow [8]. And another one in Coq is given by Coquand and Siles [4].

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 [5] and the other is POSIX matching [1,9,10,13,14]. For example consider the string $xy$ and the regular expression $(x + y + xy)^\star$. 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 $^\star$ 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 $(r_{key} + r_{id})^\star$ 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 $(r_{key} + r_{id})^\star$ matches *iffoo*, respectively *if*, in exactly one 'iteration' of the star.

The Empty String Rule is for cases where, for example, the regular expression $(a^\star)^\star$ 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 [13] 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 $(x + (y + xy))^\star$ (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

$$\textit{Stars } [\textit{Left } (\textit{Char } x), \textit{Right } (\textit{Left } (\textit{Char } y))]$$

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, [7]). The value for matching *xy* in a single 'iteration', i.e. the POSIX value, would look as follows

$$\textit{Stars } [\textit{Seq } (\textit{Char } x) \ (\textit{Char } y)]$$

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 [13] 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 [5] 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 [13]); 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 [9] who found

that nearly all POSIX matching implementations are "buggy" [13, Page 203] and by Grathwohl et al [6, 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 [13]. We have proved the correctness of this algorithm according to our specification of what a POSIX value is (inspired by work of Vansummeren [14]). Sulzmann and Lu sketch in [13] an informal correctness proof: but to us it contains unfillable gaps.[4] 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 [10].

We extend our results to ??? Bitcoded version??

## 4  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 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:

$$r := \mathbf{0} \mid \mathbf{1} \mid c \mid r_1 + r_2 \mid r_1 \cdot r_2 \mid r^\star$$

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:

$$
\begin{aligned}
&\textit{(1)} & L(\mathbf{0}) &\stackrel{\text{def}}{=} \varnothing \\
&\textit{(2)} & L(\mathbf{1}) &\stackrel{\text{def}}{=} \{[]\} \\
&\textit{(3)} & L(c) &\stackrel{\text{def}}{=} \{[c]\} \\
&\textit{(4)} & L(r_1 \cdot r_2) &\stackrel{\text{def}}{=} L(r_1) \mathbin{@} L(r_2) \\
&\textit{(5)} & L(r_1 + r_2) &\stackrel{\text{def}}{=} L(r_1) \cup L(r_2) \\
&\textit{(6)} & L(r^\star) &\stackrel{\text{def}}{=} (L(r))\star
\end{aligned}
$$

---

[4] An extended version of [13] 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.

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

$$Der\ c\ A\ \stackrel{def}{=}\ \{s \mid c :: s \in A\}\ .$$

For semantic derivatives we have the following equations (for example mechanically proved in [8]):

$$
\begin{aligned}
Der\ c\ \varnothing &\stackrel{def}{=} \varnothing \\
Der\ c\ \{[]\} &\stackrel{def}{=} \varnothing \\
Der\ c\ \{[d]\} &\stackrel{def}{=} if\ c = d\ then\ \{[]\}\ else\ \varnothing \\
Der\ c\ (A \cup B) &\stackrel{def}{=} Der\ c\ A \cup Der\ c\ B \\
Der\ c\ (A\ @\ B) &\stackrel{def}{=} (Der\ c\ A\ @\ B) \cup (if\ [] \in A\ then\ Der\ c\ B\ else\ \varnothing) \\
Der\ c\ (A\star) &\stackrel{def}{=} Der\ c\ A\ @\ A\star
\end{aligned}
\tag{1}
$$

*Brzozowski's derivatives* of regular expressions [3] 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:

$$
\begin{aligned}
nullable\ (\mathbf{0}) &\stackrel{def}{=} False \\
nullable\ (\mathbf{1}) &\stackrel{def}{=} True \\
nullable\ (c) &\stackrel{def}{=} False \\
nullable\ (r_1 + r_2) &\stackrel{def}{=} nullable\ r_1 \lor nullable\ r_2 \\
nullable\ (r_1 \cdot r_2) &\stackrel{def}{=} nullable\ r_1 \land nullable\ r_2 \\
nullable\ (r^\star) &\stackrel{def}{=} True
\end{aligned}
$$

$$
\begin{aligned}
\mathbf{0} \backslash c &\stackrel{def}{=} \mathbf{0} \\
\mathbf{1} \backslash c &\stackrel{def}{=} \mathbf{0} \\
d \backslash c &\stackrel{def}{=} if\ c = d\ then\ \mathbf{1}\ else\ \mathbf{0} \\
(r_1 + r_2) \backslash c &\stackrel{def}{=} (r_1 \backslash c) + (r_2 \backslash c) \\
(r_1 \cdot r_2) \backslash c &\stackrel{def}{=} if\ nullable\ r_1\ then\ (r_1 \backslash c) \cdot r_2 + (r_2 \backslash c)\ else\ (r_1 \backslash c) \cdot r_2 \\
(r^\star) \backslash c &\stackrel{def}{=} (r \backslash c) \cdot r^\star
\end{aligned}
$$

We may extend this definition to give derivatives w.r.t. strings:

$$
\begin{aligned}
r \backslash [] &\stackrel{def}{=} r \\
r \backslash (c :: s) &\stackrel{def}{=} (r \backslash c) \backslash s
\end{aligned}
$$

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) $L(r \backslash c) = Der\ c\ (L(r))$.*

With this in place it is also very routine to prove that the regular expression matcher defined as

$$match\ r\ s \overset{def}{=} nullable\ (r \backslash s)$$

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 [13] is to append another phase to this algorithm in order to calculate a [lexical] value. We will explain the details next.

## 5   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 [13] 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

$$v := Empty \mid Char\ c \mid Left\ v \mid Right\ v \mid Seq\ v_1\ v_2 \mid Stars\ vs$$

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 [5], and Sulzmann and Lu for POSIX matching [13]). The string underlying a value can be calculated by the *flat* function, written $|\_|$ and defined as:

$$\begin{array}{llll}
|Empty| & \overset{def}{=} [] & |Seq\ v_1\ v_2| & \overset{def}{=} |v_1|\ @\ |v_2| \\
|Char\ c| & \overset{def}{=} [c] & |Stars\ []| & \overset{def}{=} [] \\
|Left\ v| & \overset{def}{=} |v| & |Stars\ (v::vs)| & \overset{def}{=} |v|\ @\ |Stars\ vs| \\
|Right\ v| & \overset{def}{=} |v| &
\end{array}$$

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:[5]

---

[5] Note that the rule for *Stars* differs from our earlier paper [2]. 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.

$$\overline{Empty : \mathbf{1}} \qquad\qquad \overline{Char\ c : c}$$

$$\frac{v_1 : r_1}{Left\ v_1 : r_1 + r_2} \qquad\qquad \frac{v_2 : r_1}{Right\ v_2 : r_2 + r_1}$$

$$\frac{v_1 : r_1 \qquad v_2 : r_2}{Seq\ v_1\ v_2 : r_1 \cdot r_2} \qquad \frac{\forall\, v \in vs.\ v : r \wedge |v| \neq []}{Stars\ vs : r^\star}$$

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 $^\star$ 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 $s$:[6]

$$LV\ r\ s \overset{def}{=} \{v \mid v : r \wedge |v| = s\}$$

The main property of $LV\ r\ s$ is that it is alway 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, $LV\ (\mathbf{1}^\star)\ []$ would contain infinitely many values, but according to our more restricted definition only a single value, namely $LV\ (\mathbf{1}^\star)$ $[] = \{Stars\ []\}$.

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:

---

[6] Okui and Suzuki refer to our lexical values as *canonical values* in [10]. The notion of *non-problematic values* by Cardelli and Frisch [5] is related, but not identical to our lexical values.

$$r_1 \xrightarrow{\_\backslash a} r_2 \xrightarrow{\_\backslash b} r_3 \xrightarrow{\_\backslash c} r_4 \; nullable$$

$$\big\downarrow mkeps$$

$$v_1 \xleftarrow[inj \; r_1 \; a]{} v_2 \xleftarrow[inj \; r_2 \; b]{} v_3 \xleftarrow[inj \; r_3 \; c]{} v_4$$
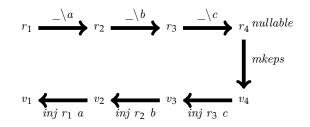
**Fig. 1.** The two phases of the algorithm by Sulzmann & Lu [13], matching the string $[a, b, c]$. The first phase (the arrows from left to right) is Brzozowski's matcher building successive derivatives. If the last regular expression is *nullable*, then the functions of the second phase are called (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.

$$
\begin{aligned}
mkeps \; \mathbf{1} &\stackrel{\text{def}}{=} Empty \\
mkeps \; (r_1 \cdot r_2) &\stackrel{\text{def}}{=} Seq \; (mkeps \; r_1) \; (mkeps \; r_2) \\
mkeps \; (r_1 + r_2) &\stackrel{\text{def}}{=} if \; nullable \; r_1 \; then \; Left \; (mkeps \; r_1) \; else \; Right \; (mkeps \; r_2) \\
mkeps \; (r^\star) &\stackrel{\text{def}}{=} Stars \; []
\end{aligned}
$$

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 [13] 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).

$$
\begin{array}{lll}
\textit{(1)} & \textit{inj d c (Empty)} & \stackrel{\text{def}}{=} \textit{Char d} \\
\textit{(2)} & \textit{inj } (r_1 + r_2) \textit{ c (Left } v_1) & \stackrel{\text{def}}{=} \textit{Left (inj } r_1 \textit{ c } v_1) \\
\textit{(3)} & \textit{inj } (r_1 + r_2) \textit{ c (Right } v_2) & \stackrel{\text{def}}{=} \textit{Right (inj } r_2 \textit{ c } v_2) \\
\textit{(4)} & \textit{inj } (r_1 \cdot r_2) \textit{ c (Seq } v_1 \textit{ } v_2) & \stackrel{\text{def}}{=} \textit{Seq (inj } r_1 \textit{ c } v_1) \textit{ } v_2 \\
\textit{(5)} & \textit{inj } (r_1 \cdot r_2) \textit{ c (Left (Seq } v_1 \textit{ } v_2)) & \stackrel{\text{def}}{=} \textit{Seq (inj } r_1 \textit{ c } v_1) \textit{ } v_2 \\
\textit{(6)} & \textit{inj } (r_1 \cdot r_2) \textit{ c (Right } v_2) & \stackrel{\text{def}}{=} \textit{Seq (mkeps } r_1) \textit{ (inj } r_2 \textit{ c } v_2) \\
\textit{(7)} & \textit{inj } (r^\star) \textit{ c (Seq } v \textit{ (Stars } vs)) & \stackrel{\text{def}}{=} \textit{Stars (inj r c } v :: vs)
\end{array}
$$

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:

$$
(r_1 \cdot r_2)\backslash c \stackrel{\text{def}}{=} \textit{ if nullable } r_1 \textit{ then } (r_1\backslash c) \cdot r_2 + (r_2\backslash c) \textit{ else } (r_1\backslash c) \cdot r_2
$$

Consider first the *else*-branch where the derivative is $(r_1\backslash c) \cdot r_2$. 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 $(r_1\backslash c) \cdot r_2 + (r_2\backslash c)$. 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* (*Stars vs*)—the derivative of a star is $(r\backslash c) \cdot r^\star$. 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.[7]

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 *mkeps*-value is always the empty string (given the regular expression is nullable since otherwise *mkeps* might not be defined).

**Lemma 2.**
*(1) If $v : r\backslash c$ then $|inj\ r\ c\ v| = c :: |v|$.*
*(2) If nullable r then $|mkeps\ r| = []$.*

---

[7] Sulzmann and Lu state this clause as *inj c c (Empty)* $\stackrel{\text{def}}{=}$ *Char c*, but our deviation is harmless.

*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.                                      □

Having defined the *mkeps* 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

$$
\begin{aligned}
lexer\ r\ [] \quad &\stackrel{\mathrm{def}}{=}\ if\ nullable\ r\ then\ Some\ (mkeps\ r)\ else\ None \\
lexer\ r\ (c::s) &\stackrel{\mathrm{def}}{=}\ case\ lexer\ (r\backslash c)\ s\ of \\
&\qquad None \Rightarrow None \\
&\qquad |\ Some\ v \Rightarrow Some\ (inj\ r\ c\ v)
\end{aligned}
$$

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 [13], 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 [14]. The relation we define is ternary and written as $(s,\ r) \to v$, 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 $(s,\ r) \to v$.

**Theorem 1.**
*(1) If $(s,\ r) \to v$ then $s \in L(r)$ and $|v| = s$.*
*(2) If $(s,\ r) \to v$ and $(s,\ r) \to v'$ then $v = v'$.*

*Proof.* Both by induction on the definition of $(s,\ r) \to v$. The second parts follows by a case analysis of $(s,\ r) \to v'$ and the first part.                      □

We claim that our $(s,\ r) \to v$ 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 $s_1\ @\ s_2$.

$$\frac{}{([], \mathbf{1}) \rightarrow Empty} P\mathbf{1} \qquad \frac{}{([c], c) \rightarrow Char\ c} Pc$$

$$\frac{(s, r_1) \rightarrow v}{(s, r_1 + r_2) \rightarrow Left\ v} P{+}L \qquad \frac{(s, r_2) \rightarrow v \qquad s \notin L(r_1)}{(s, r_1 + r_2) \rightarrow Right\ v} P{+}R$$

$$\frac{\begin{array}{c}(s_1, r_1) \rightarrow v_1 \qquad (s_2, r_2) \rightarrow v_2 \\ \nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r_1) \wedge s_4 \in L(r_2)\end{array}}{(s_1 @ s_2, r_1 \cdot r_2) \rightarrow Seq\ v_1\ v_2} PS$$

$$\frac{}{([], r^\star) \rightarrow Stars\ []} P[]$$

$$\frac{\begin{array}{c}(s_1, r) \rightarrow v \qquad (s_2, r^\star) \rightarrow Stars\ vs \qquad |v| \neq [] \\ \nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3 @ s_4 = s_2 \wedge s_1 @ s_3 \in L(r) \wedge s_4 \in L(r^\star)\end{array}}{(s_1 @ s_2, r^\star) \rightarrow Stars\ (v{::}vs)} P\star$$

**Fig. 2.** Our inductive definition of POSIX values.

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 *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 *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 3.** *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 4.** *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 5.** *If $(s, r\backslash c) \rightarrow v$ 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 $(s, r_1\backslash c) \rightarrow v'$; and $(b)$ $v = Right\ v'$, $s \notin L(r_1\backslash c)$ and $(s, r_2\backslash c) \rightarrow v'$. In $(a)$ we know $(s, r_1\backslash c) \rightarrow v'$, 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 $s \notin L(r_1\backslash c)$.

- 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 $(s_1, r_1\backslash c) \rightarrow v_1$ and $(s_2, r_2) \rightarrow v_2$ as well as

  $$\nexists s_3\ s_4.a.\ s_3 \neq [] \wedge s_3\ @\ s_4 = s_2 \wedge s_1\ @\ s_3 \in L(r_1\backslash c) \wedge s_4 \in L(r_2)$$

  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 $(s, r_2\backslash c) \rightarrow v_1$ and $s_1\ @\ s_2 \notin L((r_1\backslash c) \cdot r_2)$. 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 4 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 $r = r_1^{\star}$. 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 $(s_1, r_1\backslash c) \rightarrow v_1$ and the induction hypothesis). $\qquad\square$

With Lemma 5 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 lexer $r\ s = None$*
*(2) $s \in L(r)$ if and only if $\exists v.$ lexer $r\ s = Some\ v \wedge (s, r) \rightarrow v$*

*Proof.* By induction on *s* using Lemma 4 and 5.                    □

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) lexer r s = None if and only if $\nexists$ v.a. (s, r) → v*
*(2) lexer r s = Some v if and only if (s, r) → v*

This concludes our correctness proof. Note that we have not changed the algorithm of Sulzmann and Lu,[8] 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.


## 6   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 [13]: they introduced an ordering for values and identified POSIX values as the maximal elements. An extended version of [13] 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 [10,11] 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*

$$v \stackrel{def}{=} Stars\ [Seq\ (Char\ x)\ (Char\ y),\ Char\ z]$$

At position *[0,1]* of this value is the subvalue *Char y* and at position *[1]* the subvalue *Char z*. At the 'root' position, or empty list *[]*, is the whole value *v*. Positions such as *[0,1,0]* or *[2]* are outside of *v*. If it exists, the subvalue of *v* at a position *p*, written $v\downarrow_p$, can be recursively defined by

---

[8] All deviations we introduced are harmless.

$$v\!\downarrow_{[]} \quad \overset{def}{=} \quad v$$
$$Left\ v\!\downarrow_{0::ps} \quad \overset{def}{=} \quad v\!\downarrow_{ps}$$
$$Right\ v\!\downarrow_{1::ps} \quad \overset{def}{=} \quad v\!\downarrow_{ps}$$
$$Seq\ v_1\ v_2\!\downarrow_{0::ps} \quad \overset{def}{=} \quad v_1\!\downarrow_{ps}$$
$$Seq\ v_1\ v_2\!\downarrow_{1::ps} \quad \overset{def}{=} \quad v_2\!\downarrow_{ps}$$
$$Stars\ vs\!\downarrow_{n::ps} \quad \overset{def}{=} \quad vs_{[n]}\!\downarrow_{ps}$$

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 *Pos v*, is given by

$$Pos\ (Empty) \quad \overset{def}{=} \quad \{[]\}$$
$$Pos\ (Char\ c) \quad \overset{def}{=} \quad \{[]\}$$
$$Pos\ (Left\ v) \quad \overset{def}{=} \quad \{[]\} \cup \{0::ps \mid ps \in Pos\ v\}$$
$$Pos\ (Right\ v) \quad \overset{def}{=} \quad \{[]\} \cup \{1::ps \mid ps \in Pos\ v\}$$
$$Pos\ (Seq\ v_1\ v_2) \quad \overset{def}{=} \quad \{[]\} \cup \{0::ps \mid ps \in Pos\ v_1\} \cup \{1::ps \mid ps \in Pos\ v_2\}$$
$$Pos\ (Stars\ vs) \quad \overset{def}{=} \quad \{[]\} \cup (\bigcup n < len\ vs\ \{n::ps \mid ps \in Pos\ vs_{[n]}\})$$

whereby *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 $w$:

$$v \overset{def}{=} Stars\ [Seq\ (Char\ x)\ (Char\ y),\ Char\ z]$$
$$w \overset{def}{=} Stars\ [Char\ x,\ Char\ y,\ Char\ z]$$

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

$$\|v\|_p \overset{def}{=} if\ p \in Pos\ v\ then\ len\ |v\!\downarrow_p|\ else\ -1$$

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

$$v \overset{def}{=} Left\ (Char\ x) \qquad and \qquad w \overset{def}{=} Right\ (Char\ x)$$

Both values match $x$. At position $[0]$ the norm of $v$ is $1$ (the subvalue matches $x$), but the norm of $w$ is $-1$ (the position is outside $w$ according to how we defined the 'inside' positions of *Left-* and *Right*-values). Of course at position $[1]$, the norms $\|v\|_{[1]}$ and $\|w\|_{[1]}$ 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 $[1]$ and thus also $v$ will be preferred over $w$. The lexicographic ordering of positions, written $\_ \prec_{lex} \_$, can be conveniently formalised by three inference rules

$$\frac{}{[] \prec_{lex} p :: ps} \qquad \frac{p_1 < p_2}{p_1 :: ps_1 \prec_{lex} p_2 :: ps_2} \qquad \frac{ps_1 \prec_{lex} ps_2}{p :: ps_1 \prec_{lex} p :: ps_2}$$

With the norm and lexicographic order in place, we can state the key definition of Okui and Suzuki [10]: 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

$$v_1 \prec_p v_2 \stackrel{def}{=} \begin{cases} (i) & \|v_2\|_p < \|v_1\|_p \quad \text{and} \\ (ii) & \forall\, q \in Pos\ v_1 \cup Pos\ v_2.\ q \prec_{lex} p \longrightarrow \|v_1\|_q = \|v_2\|_q \end{cases}$$

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 [10]. 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:

$$v_1 \prec v_2 \stackrel{def}{=} \exists\, p.\ v_1 \prec_p v_2$$
$$v_1 \preccurlyeq v_2 \stackrel{def}{=} v_1 \prec v_2 \vee v_1 = v_2$$

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 $\_ \preccurlyeq \_$ by Okui and Suzuki.

**Lemma 6 (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 $\prec_{lex}$ is trichotomous, we need to consider three cases, namely $p = q$, $p \prec_{lex} q$ and $q \prec_{lex} p$. 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 $p' \in Pos\ v_1 \cup Pos\ v_3$ with $p' \prec_{lex} p$ that $\|v_1\|_{p'} = \|v_3\|_{p'}$ holds. Suppose $p' \in Pos\ v_1$, then we can infer from the first assumption that $\|v_1\|_{p'} = \|v_2\|_{p'}$. But this means that $p'$ must be in $Pos\ v_2$ too (the norm cannot be $-1$ given $p' \in Pos\ v_1$). 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.                                                                              □

The proof for $\preccurlyeq$ is similar and omitted. It is also straightforward to show that $\prec$ and $\preccurlyeq$ are partial orders. Okui and Suzuki furthermore show that they are linear orderings for lexical values [10] 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 $v_1 \preccurlyeq v_2$.*

*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 $v_2 \in LV$ **1** $[]$ must also be of the form $v_2 = Empty$. Therefore we have $v_1 \preccurlyeq v_2$. 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 $v_1 \preccurlyeq v_2$ 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 $w_1 \preccurlyeq w_2$. Because $w_1$ and $w_2$ have the same underlying string $s$, we can conclude with *Left* $w_1 \preccurlyeq$ *Left* $w_2$ using Proposition 5*(2)*.
- Case $P+R$ with $(s, r_1 + r_2) \rightarrow$ *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 $v_1 \preccurlyeq v_2$.
- Case $PS$ with $(s_1 @ s_2, r_1 \cdot r_2) \rightarrow$ *Seq* $w_1$ $w_2$: We can assume $v_2 =$ *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 $v_1 \preccurlyeq v_2$ 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 $w_1 \preccurlyeq u_1$ and $w_2 \preccurlyeq u_2$. By Proposition 5*(4,5)* we can again infer $v_1 \preccurlyeq v_2$.

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 3 we also have $v_P \in LV\ r\ s$. By Theorem 3 we therefore have $v_P \preccurlyeq v_1$. 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 $LV\ r\ s \neq \varnothing$ then $\exists! vmin.\ vmin \in LV\ r\ s \wedge (\forall\, v \in LV\ r\ s.\ vmin \preccurlyeq v)$.*

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 5. This provides an independent confirmation that our ternary relation formalises the informal POSIX rules.

## 7  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

*lexer r s = (if nullable (r\s) then Some (flex r id s (mkeps (r\s))) else None)*

## 8  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 [13]. 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

$$\mathbf{0} + r \Rightarrow r \qquad r + \mathbf{0} \Rightarrow r \qquad \mathbf{1} \cdot r \Rightarrow r \qquad r \cdot \mathbf{1} \Rightarrow r \tag{2}$$

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 [13] 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:

$$F_{Right}\ f\ v \quad\stackrel{\text{def}}{=}\quad Right\ (f\ v)$$

$$F_{Left}\ f\ v \quad\stackrel{\text{def}}{=}\quad Left\ (f\ v)$$

$$F_{Alt}\ f_1\ f_2\ (Right\ v) \quad\stackrel{\text{def}}{=}\quad Right\ (f_2\ v)$$

$$F_{Alt}\ f_1\ f_2\ (Left\ v) \quad\stackrel{\text{def}}{=}\quad Left\ (f_1\ v)$$

$$F_{Seq1}\ f_1\ f_2\ v \quad\stackrel{\text{def}}{=}\quad Seq\ (f_1\ ())\ (f_2\ v)$$

$$F_{Seq2}\ f_1\ f_2\ v \quad\stackrel{\text{def}}{=}\quad Seq\ (f_1\ v)\ (f_2\ ())$$

$$F_{Seq}\ f_1\ f_2\ (Seq\ v_1\ v_2) \quad\stackrel{\text{def}}{=}\quad Seq\ (f_1\ v_1)\ (f_2\ v_2)$$

$$simp_{Alt}\ (\mathbf{0},\ \_)\ (r_2, f_2) \quad\stackrel{\text{def}}{=}\quad (r_2,\ F_{Right}\ f_2)$$

$$simp_{Alt}\ (r_1, f_1)\ (\mathbf{0},\ \_) \quad\stackrel{\text{def}}{=}\quad (r_1,\ F_{Left}\ f_1)$$

$$simp_{Alt}\ (r_1, f_1)\ (r_2, f_2) \quad\stackrel{\text{def}}{=}\quad (r_1 + r_2,\ F_{Alt}\ f_1\ f_2)$$

$$simp_{Seq}\ (\mathbf{1}, f_1)\ (r_2, f_2) \quad\stackrel{\text{def}}{=}\quad (r_2,\ F_{Seq1}\ f_1\ f_2)$$

$$simp_{Seq}\ (r_1, f_1)\ (\mathbf{1}, f_2) \quad\stackrel{\text{def}}{=}\quad (r_1,\ F_{Seq2}\ f_1\ f_2)$$

$$simp_{Seq}\ (r_1, f_1)\ (r_2, f_2) \quad\stackrel{\text{def}}{=}\quad (r_1 \cdot r_2,\ F_{Seq}\ f_1\ f_2)$$

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

$$simp\ (r_1 + r_2) \quad\stackrel{\text{def}}{=}\quad simp_{Alt}\ (simp\ r_1)\ (simp\ r_2)$$

$$simp\ (r_1 \cdot r_2) \quad\stackrel{\text{def}}{=}\quad simp_{Seq}\ (simp\ r_1)\ (simp\ r_2)$$

$$simp\ r \quad\stackrel{\text{def}}{=}\quad (r,\ id)$$

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:

$$lexer^+\ r\ [] \quad\stackrel{\text{def}}{=}\quad if\ nullable\ r\ then\ Some\ (mkeps\ r)\ else\ None$$

$$lexer^+\ r\ (c::s) \quad\stackrel{\text{def}}{=}\quad let\ (r_s, f_r) = simp\ (r\backslash c)\ in$$
$$case\ lexer^+\ r_s\ s\ of$$
$$None \Rightarrow None$$
$$|\ Some\ v \Rightarrow Some\ (inj\ r\ c\ (f_r\ v))$$

In the second clause we first calculate the derivative $r\backslash c$ and then simpli

text    *Incremental calculation of the value. To simplify the proof we first define the function @{const flex} which calculates the "iterated" injection function. With this we can rewrite the lexer as \begin{center} @{thm lexer__flex} \end{center} \begin{center} \begin{tabular}{lcl} @{thm (lhs) code.simps(1)} & $\dn$ & @{thm (rhs) code.simps(1)}\\ @{thm (lhs) code.simps(2)} & $\dn$ & @{thm (rhs) code.simps(2)}\\ @{thm (lhs) code.simps(3)} & $\dn$*

& @{thm (rhs) code.simps(3)}\\ @{thm (lhs) code.simps(4)} & $\dn$ &
@{thm (rhs) code.simps(4)}\\ @{thm (lhs) code.simps(5)[of $v_1$ $v_2$]} & $\dn$
& @{thm (rhs) code.simps(5)[of $v_1$ $v_2$]}\\ @{thm (lhs) code.simps(6)} &
$\dn$ & @{thm (rhs) code.simps(6)}\\ @{thm (lhs) code.simps(7)} & $\dn$
& @{thm (rhs) code.simps(7)} \end{tabular} \end{center} \begin{center}
\begin{tabular}{lcl} @{term areg} & $::=$ & @{term AZERO}\\ & $\mid$
& @{term AONE bs}\\ & $\mid$ & @{term ACHAR bs c}\\ & $\mid$
& @{term AALT bs r1 r2}\\ & $\mid$ & @{term ASEQ bs $r_1$ $r_2$}\\ &
$\mid$ & @{term ASTAR bs r} \end{tabular} \end{center} \begin{center}
\begin{tabular}{lcl} @{thm (lhs) intern.simps(1)} & $\dn$ & @{thm (rhs)
intern.simps(1)}\\ @{thm (lhs) intern.simps(2)} & $\dn$ & @{thm (rhs) in-
tern.simps(2)}\\ @{thm (lhs) intern.simps(3)} & $\dn$ & @{thm (rhs) in-
tern.simps(3)}\\ @{thm (lhs) intern.simps(4)[of $r_1$ $r_2$]} & $\dn$ & @{thm
(rhs) intern.simps(4)[of $r_1$ $r_2$]}\\ @{thm (lhs) intern.simps(5)[of $r_1$ $r_2$]} &
$\dn$ & @{thm (rhs) intern.simps(5)[of $r_1$ $r_2$]}\\ @{thm (lhs) intern.simps(6)}
& $\dn$ & @{thm (rhs) intern.simps(6)}\\ \end{tabular} \end{center} \begin{center}
\begin{tabular}{lcl} @{thm (lhs) erase.simps(1)} & $\dn$ & @{thm (rhs)
erase.simps(1)}\\ @{thm (lhs) erase.simps(2)[of bs]} & $\dn$ & @{thm (rhs)
erase.simps(2)[of bs]}\\ @{thm (lhs) erase.simps(3)[of bs]} & $\dn$ & @{thm
(rhs) erase.simps(3)[of bs]}\\ @{thm (lhs) erase.simps(4)[of bs $r_1$ $r_2$]} & $\dn$
& @{thm (rhs) erase.simps(4)[of bs $r_1$ $r_2$]}\\ @{thm (lhs) erase.simps(5)[of
bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs) erase.simps(5)[of bs $r_1$ $r_2$]}\\ @{thm
(lhs) erase.simps(6)[of bs]} & $\dn$ & @{thm (rhs) erase.simps(6)[of bs]}\\
\end{tabular} \end{center} Some simple facts about erase \begin{lemma}\mbox{}\\
@{thm erase__bder}\\ @{thm erase__intern} \end{lemma} \begin{center}
\begin{tabular}{lcl} @{thm (lhs) bnullable.simps(1)} & $\dn$ & @{thm (rhs)
bnullable.simps(1)}\\ @{thm (lhs) bnullable.simps(2)} & $\dn$ & @{thm (rhs)
bnullable.simps(2)}\\ @{thm (lhs) bnullable.simps(3)} & $\dn$ & @{thm (rhs)
bnullable.simps(3)}\\ @{thm (lhs) bnullable.simps(4)[of bs $r_1$ $r_2$]} & $\dn$ &
@{thm (rhs) bnullable.simps(4)[of bs $r_1$ $r_2$]}\\ @{thm (lhs) bnullable.simps(5)[of
bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs) bnullable.simps(5)[of bs $r_1$ $r_2$]}\\ @{thm
(lhs) bnullable.simps(6)} & $\dn$ & @{thm (rhs) bnullable.simps(6)}\medskip\\
% \end{tabular} % \end{center} % \begin{center} % \begin{tabular}{lcl}
@{thm (lhs) bder.simps(1)} & $\dn$ & @{thm (rhs) bder.simps(1)}\\ @{thm
(lhs) bder.simps(2)} & $\dn$ & @{thm (rhs) bder.simps(2)}\\ @{thm (lhs)
bder.simps(3)} & $\dn$ & @{thm (rhs) bder.simps(3)}\\ @{thm (lhs) bder.simps(4)[of
bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs) bder.simps(4)[of bs $r_1$ $r_2$]}\\ @{thm (lhs)
bder.simps(5)[of bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs) bder.simps(5)[of bs $r_1$
$r_2$]}\\ @{thm (lhs) bder.simps(6)} & $\dn$ & @{thm (rhs) bder.simps(6)}
\end{tabular} \end{center} \begin{center} \begin{tabular}{lcl} @{thm (lhs)
bmkeps.simps(1)} & $\dn$ & @{thm (rhs) bmkeps.simps(1)}\\ @{thm (lhs)
bmkeps.simps(2)[of bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs) bmkeps.simps(2)[of bs
$r_1$ $r_2$]}\\ @{thm (lhs) bmkeps.simps(3)[of bs $r_1$ $r_2$]} & $\dn$ & @{thm (rhs)
bmkeps.simps(3)[of bs $r_1$ $r_2$]}\\ @{thm (lhs) bmkeps.simps(4)} & $\dn$ &
@{thm (rhs) bmkeps.simps(4)}\medskip\\ \end{tabular} \end{center} @{thm

[mode=IfThen] bder__retrieve}   By induction on ⟨r⟩   \begin{theorem}[Main Lemma]\mbox{}\\  @{thm [mode=IfThen]  MAIN__decode}  \end{theorem} \noindent Definition of the bitcoded lexer @{thm blexer__def}   \begin{theorem} @{thm blexer__correctness}  \end{theorem}

   section *Optimisations*

   text   *Derivatives as calculated by* \Brz'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* @{term ALT ZERO r}, @{term ALT r ZERO}, @{term SEQ ONE r} *and* @{term SEQ r ONE} *to* @{term r}. *These simplifications can speed up the algorithms considerably, as noted in* \cite{Sulzmann2014}. *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* \begin{equation}\label{Simpl} \begin{array}{lcllcllcllcl} @{term ALT ZERO r} & ⟨⇒⟩ & @{term r} \hspace{8mm}%\\ @{term ALT r ZERO} & ⟨⇒⟩ & @{term r} \hspace{8mm}%\\ @{term SEQ ONE r}  & ⟨⇒⟩ & @{term r} \hspace{8mm}%\\ @{term SEQ r ONE}  & ⟨⇒⟩ & @{term r} \end{array} \end{equation}   \noindent *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,* \emph{not} *for the original (unsimplified) derivative regular expression. Sulzmann and Lu* \cite{Sulzmann2014} *overcome this obstacle by not just calculating a simplified regular expression, but also calculating a* \emph{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{center} \begin{tabular}{lcl} @{thm (lhs) F__RIGHT.simps(1)} & $\dn$ & ⟨Right (f v)⟩\\ @{thm (lhs) F__LEFT.simps(1)} & $\dn$ & ⟨Left (f v)⟩\\  @{thm (lhs) F__ALT.simps(1)} & $\dn$ & ⟨Right $(f_2$ v)⟩\\ @{thm (lhs) F__ALT.simps(2)} & $\dn$ & ⟨Left $(f_1$ v)⟩\\  @{thm (lhs) F__SEQ1.simps(1)} & $\dn$ & ⟨Seq $(f_1$ ()) $(f_2$ v)⟩\\ @{thm (lhs) F__SEQ2.simps(1)} & $\dn$ & ⟨Seq $(f_1$ v) $(f_2$ ())⟩\\ @{thm (lhs) F__SEQ.simps(1)} & $\dn$ & ⟨Seq $(f_1$ $v_1)$ $(f_2$ $v_2)$⟩\medskip\\ %\end{tabular} % %\begin{tabular}{lcl} @{term simp__ALT (ZERO, DUMMY) $(r_2, f_2)$} & $\dn$ & @{term $(r_2,$ F__RIGHT $f_2)$}\\ @{term simp__ALT $(r_1, f_1)$ (ZERO, DUMMY)} & $\dn$ & @{term $(r_1,$ F__LEFT $f_1)$}\\ @{term simp__ALT $(r_1, f_1)$ $(r_2, f_2)$} & $\dn$ & @{term (ALT $r_1$ $r_2,$ F__ALT $f_1$ $f_2)$}\\ @{term simp__SEQ (ONE, $f_1)$ $(r_2, f_2)$} & $\dn$ & @{term $(r_2,$ F__SEQ1 $f_1$ $f_2)$}\\ @{term simp__SEQ $(r_1, f_1)$ (ONE, $f_2)$} & $\dn$ & @{term $(r_1,$ F__SEQ2 $f_1$ $f_2)$}\\ @{term simp__SEQ $(r_1, f_1)$ $(r_2, f_2)$} & $\dn$ & @{term (SEQ $r_1$ $r_2,$ F__SEQ $f_1$ $f_2)$}\\ \end{tabular} \end{center}   \noindent *The functions* ⟨$simp_{Alt}$⟩ *and* ⟨$simp_{Seq}$⟩ *encode the simplification rules in* \eqref{Simpl} *and compose the rectification functions (simplifications can occur deep inside the regular expression). The main simplification function is then* \begin{center} \begin{tabular}{lcl}

@{term simp (ALT $r_1$ $r_2$)} & $\dn$ & @{term simp__ALT (simp $r_1$) (simp $r_2$)}\\ @{term simp (SEQ $r_1$ $r_2$)} & $\dn$ & @{term simp__SEQ (simp $r_1$) (simp $r_2$)}\\ @{term simp r} & $\dn$ & @{term (r, id)}\\ \end{tabular} \end{center}   \noindent where @{term id} stands for the identity function. The function @{const 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{center} \begin{tabular}{lcl} @{thm (lhs) slexer.simps(1)} & $\dn$ & @{thm (rhs) slexer.simps(1)}\\ @{thm (lhs) slexer.simps(2)} & $\dn$ & ⟨let ($r_s$, $f_r$) = simp (r ⟩$\backslash$⟨ c⟩ in⟩\\ & & ⟨case⟩ @{term slexer $r_s$ s} ⟨of⟩\\ & & \phantom{$|$} @{term None} ⟨⇒⟩ @{term None}\\ & & $|$ @{term Some v} ⟨⇒⟩ ⟨Some (inj r c ($f_r$ v))⟩ \end{tabular} \end{center}   \noindent In the second clause we first calculate the derivative @{term der c r} 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 @{term c} into the value @{term v}, we need to rectify @{term v} (that is construct @{term $f_r$ v}). Before we can establish the correctness of @{term slexer}, we need to show that simplification preserves the language and simplification preserves our POSIX relation once the value is rectified (recall @{const simp} generates a (regular expression, rectification function) pair): \begin{lemma}\mbox{}\smallskip\\\label{slexeraux} \begin{tabular}{ll} (1) & @{thm L__fst__simp[symmetric]}\\ (2) & @{thm[mode=IfThen] Posix__simp} \end{tabular} \end{lemma} \begin{proof} Both are by induction on ⟨r⟩. There is no interesting case for the first statement. For the second statement, of interest are the @{term r = ALT $r_1$ $r_2$} and @{term r = SEQ $r_1$ $r_2$} cases. In each case we have to analyse four subcases whether @{term fst (simp $r_1$)} and @{term fst (simp $r_2$)} equals @{const ZERO} (respectively @{const ONE}). For example for @{term r = ALT $r_1$ $r_2$}, consider the subcase @{term fst (simp $r_1$) = ZERO} and @{term fst (simp $r_2$) ≠ ZERO}. By assumption we know @{term s ∈ fst (simp (ALT $r_1$ $r_2$)) → v}. From this we can infer @{term s ∈ fst (simp $r_2$) → v} and by IH also (∗) @{term s ∈ $r_2$ → (snd (simp $r_2$) v)}. Given @{term fst (simp $r_1$) = ZERO} we know @{term L (fst (simp $r_1$)) = {}}. By the first statement @{term L $r_1$} is the empty set, meaning (∗∗) @{term s ∉ L $r_1$}. Taking (∗) and (∗∗) together gives by the \mbox{⟨P+R⟩−rule} @{term s ∈ ALT $r_1$ $r_2$ → Right (snd (simp $r_2$) v)}. In turn this gives @{term s ∈ ALT $r_1$ $r_2$ → snd (simp (ALT $r_1$ $r_2$)) v} as we need to show. The other cases are similar.\qed \end{proof}   \noindent We can now prove relatively straightforwardly that the optimised lexer produces the expected result:  \begin{theorem} @{thm slexer__correctness} \end{theorem}   \begin{proof} By induction on @{term s} generalising over @{term r}. The case @{term []} is trivial. For the cons−case suppose the string is of the form @{term c # s}. By induction hypothesis we know @{term slexer r s = lexer r s} holds for all @{term r} (in particular for @{term r} being the derivative @{term der c r}). Let @{term $r_s$} be the simplified derivative regular expression, that is @{term fst (simp (der c

$r))\}$, and @{term $f_r$} be the rectification function, that is @{term snd (simp (der c r))}. We distinguish the cases whether $(*)$ @{term $s \in L$ (der c r)} or not. In the first case we have by Theorem~\ref{lexercorrect}(2) a value @{term v} so that @{term lexer (der c r) s = Some v} and @{term $s \in$ der $c\ r \to v$} hold. By Lemma~\ref{slexeraux}(1) we can also infer from~$(*)$ that @{term $s \in L\ r_s$} holds. Hence we know by Theorem~\ref{lexercorrect}(2) that there exists a @{term $v'$} with @{term lexer $r_s$ s = Some v'} and @{term $s \in r_s \to v'$}. From the latter we know by Lemma~\ref{slexeraux}(2) that @{term $s \in$ der c r $\to (f_r\ v')$} holds. By the uniqueness of the POSIX relation (Theorem~\ref{posixdeterm}) we can infer that @{term v} is equal to @{term $f_r\ v'$}———that is the rectification function applied to @{term $v'$} produces the original @{term v}. Now the case follows by the definitions of @{const lexer} and @{const slexer}. In the second case where @{term $s \notin L$ (der c r)} we have that @{term lexer (der c r) s = None} by Theorem~\ref{lexercorrect}(1). We also know by Lemma~\ref{slexeraux}(1) that @{term $s \notin L\ r_s$}. Hence @{term lexer $r_s$ s = None} by Theorem~\ref{lexercorrect}(1) and by IH then also @{term slexer $r_s$ s = None}. With this we can conclude in this case too.\qed \end{proof}  fy 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 $(s, fst\ (simp\ r)) \to v$ then $(s, r) \to snd\ (simp\ r)\ v$.

*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 **0** (respectively **1**). For example for $r = r_1 + r_2$, consider the subcase *fst* (*simp* $r_1$) = **0** and *fst* (*simp* $r_2$) $\neq$ **0**. By assumption we know $(s, fst\ (simp\ (r_1 + r_2))) \to v$. From this we can infer $(s, fst\ (simp\ r_2)) \to v$ and by IH also $(*)\ (s, r_2) \to snd\ (simp\ r_2)\ v$. Given *fst* (*simp* $r_1$) = **0** we know $L(fst\ (simp\ r_1)) = \varnothing$. 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 $(s, r_1 + r_2) \to Right\ (snd\ (simp\ r_2)\ v)$. In turn this gives $(s, r_1 + r_2) \to snd\ (simp\ (r_1 + r_2))\ v$ as we need to show. The other cases are similar.  □

We can now prove relatively straightforwardly that the optimised lexer produces the expected result:

**Theorem 5.** $lexer^+\ r\ s = lexer\ r\ s$

*Proof.* By induction on $s$ generalising over $r$. The case $[]$ is trivial. For the cons-case suppose the string 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 $r\backslash c$).
Let $r_s$ be the simplified derivative regular expression, that is $fst$ $(simp$ $(r\backslash c))$,
and $f_r$ be the rectification function, that is $snd$ $(simp$ $(r\backslash c))$. We distinguish the
cases whether (*) $s \in L(r\backslash c)$ or not. In the first case we have by Theorem 2(2)
a value $v$ so that $lexer$ $(r\backslash c)$ $s = Some$ $v$ and $(s, r\backslash c) \rightarrow v$ 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 $(s, r\backslash c) \rightarrow f_r$ $v'$ 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 $s \notin L(r\backslash c)$ we have that $lexer$ $(r\backslash c)$ $s = None$ 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.                                              □

## 9   HERE

**Lemma 8.** *If* $v : (r^\downarrow)\backslash c$ *then* $retrieve$ $(r\backslash\!\backslash c)$ $v = retrieve$ $r$ $(inj$ $(r^\downarrow)$ $c$ $v)$.

*Proof.* By induction on the definition of $r^\downarrow$. The cases for rule 1) and 2) are
straightforward as $\mathbf{0}\backslash c$ and $\mathbf{1}\backslash c$ are both equal to $\mathbf{0}$. This means $v : \mathbf{0}$ cannot
hold. Similarly in case of rule 3) where $r$ is of the form *ACHAR* $d$ with $c = d$.
Then by assumption we know $v : \mathbf{1}$, which implies $v = Empty$. The equation
follows by simplification of left- and right-hand side. In case $c \neq d$ we have again
$v : \mathbf{0}$, which cannot hold.

   For rule 4a) we have again $v : \mathbf{0}$. The property holds by IH for rule 4b). The
induction hypothesis is

$$retrieve\ (r\backslash\!\backslash c)\ v = retrieve\ r\ (inj\ (r^\downarrow)\ c\ v)$$

which is what left- and right-hand side simplify to. The slightly more inter-
esting case is for 4c). By assumption we have $v : ((r_1^\downarrow)\backslash c)$ $+$ $(((AALTs$ $bs$
$(r_2 :: rs))^\downarrow)\backslash c)$. This means we have either (*) $v1 : (r_1^\downarrow)\backslash c$ with $v = Left$ $v1$
or (**) $v2 : ((AALTs$ $bs$ $(r_2 :: rs))^\downarrow)\backslash c$ with $v = Right$ $v2$. The former case is
straightforward by simplification. The second case is ...TBD.

   Rule 5) TBD.

   Finally for rule 6) the reasoning is as follows: By assumption we have $v :$
$((r^\downarrow)\backslash c)$ $\cdot$ $(r^\downarrow)^\star$. This means we also have $v = Seq$ $v1$ $v2$, $v1 : (r^\downarrow)\backslash c$ and $v2 =$
*Stars vs.* We want to prove

$$retrieve\ (ASEQ\ bs\ (fuse\ [Z]\ (r\backslash\!\backslash c))\ (ASTAR\ []\ r))\ v \tag{3}$$

$$= retrieve\ (ASTAR\ bs\ r)\ (inj\ ((r^\downarrow)^\star)\ c\ v) \tag{4}$$

The right-hand side *inj*-expression is equal to *Stars* (*inj* ($r^{\downarrow}$) *c v1* :: *vs*), which means the *retrieve*-expression simplifies to

$$bs \ @ \ [Z] \ @ \ retrieve \ r \ (inj \ (r^{\downarrow}) \ c \ v1) \ @ \ retrieve \ (ASTAR \ [] \ r) \ (Stars \ vs)$$

The left-hand side (3) above simplifies to

$$bs \ @ \ retrieve \ (fuse \ [Z] \ (r \backslash\backslash c)) \ v1 \ @ \ retrieve \ (ASTAR \ [] \ r) \ (Stars \ vs)$$

We can move out the *fuse* [Z] and then use the IH to show that left-hand side and right-hand side are equal. This completes the proof.

## 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.
2. 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.
3. J. A. Brzozowski. Derivatives of Regular Expressions. *Journal of the ACM*, 11(4):481–494, 1964.
4. 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.
5. 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.
6. 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.
7. 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.
8. A. Krauss and T. Nipkow. Proof Pearl: Regular Expression Equivalence and Relation Algebra. *Journal of Automated Reasoning*, 49:95–106, 2012.
9. C. Kuklewicz. Regex Posix. https://wiki.haskell.org/Regex_Posix.
10. 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.
11. S. Okui and T. Suzuki. Disambiguation in Regular Expression Matching via Position Automata with Augmented Transitions. Technical report, University of Aizu, 2013.
12. S. Owens and K. Slind. Adapting Functional Programs to Higher Order Logic. *Higher-Order and Symbolic Computation*, 21(4):377–409, 2008.
13. 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.
14. S. Vansummeren. Type Inference for Unique Pattern Matching. *ACM Transactions on Programming Languages and Systems*, 28(3):389–428, 2006.