# Statically Typed String Sanitation Inside a Python (Technical Report)

Nathan Fulton Cyrus Omar Jonathan Aldrich

December 2014 CMU-ISR-14-112

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213

#### Abstract

This report contains supporting evidence for claims put forth and explained in the paper "Statically Typed String Sanitation Inside a Python" [1], including proofs of lemmas and theorems asserted in the paper, examples, additional discussion of the paper's technical content, and errata.

This work was supported by the National Security Agency lablet contract #H98230-14-C-0140.



## **Contents**

1	Tern	ninology and Notation	2
2	Regular Expressions		2
3	$\lambda_{RS}$ 3.1 3.2 3.3	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	2 3 3 6 10 10
4	Proo	ofs of Lemmas and Theorems About $\lambda_P$	11
5	Proo	ofs and Lemmas and Theorems About Translation	12
L	ist of	f Figures	
	1 2 3	Regular expressions over the alphabet $\Sigma$	14 14
	4	expressions	14 14
	5 6	Big step semantics for $\lambda_{RS}$	15 15
	7	Small step semantics for $\lambda_{RS}$ . Extends 6	16
	8		17 17
	10 11	Small step semantics for $\lambda_P$ (extends L-E rules)	18 19

## 1 Terminology and Notation

Theorems and lemmas appearing in [1] are numbered correspondingly, while supporting facts appearing only in the Technical Report are lettered.

## 2 Regular Expressions

The syntax of regular expressions over some alphabet  $\Sigma$  is shown in Figure 1.

**Assumption A** (Regular Expression Congruences). We assume regular expressions are implicitly identified up to the following congruences:

$$\epsilon \cdot r \equiv r$$

$$r \cdot \epsilon \equiv r$$

$$(r_1 \cdot r_2) \cdot r_3 \equiv r_1 \cdot (r_2 \cdot r_3)$$

$$r_1 + r_2 \equiv r_2 + r_1$$

$$(r_1 + r_2) + r_3 \equiv r_1 + (r_2 + r_3)$$

$$\epsilon^* = \epsilon$$

**Assumption B** (Properties of Regular Languages). We assume the following properties:

- 1. If  $s_1 \in \mathcal{L}\{r_1\}$  and  $s_2 \in \mathcal{L}\{r_2\}$  then  $s_1s_2 \in \mathcal{L}\{r_1 \cdot r_2\}$ .
- 2. For all strings s and regular expressions r, either  $s \in \mathcal{L}\{r\}$  or  $s \notin \mathcal{L}\{r\}$ .
- 3. Regular languages are closed under reversal.

#### $\lambda_{RS}$

The syntax of  $\lambda_{RS}$  is specified in Figure 2. The static semantics is specified in Figure 4.

#### 3.1 Head and Tail Operations

The following correctness conditions must hold for any definition of lhead(r) and ltail(r).

**Condition C** (Correctness of Head). *If*  $c_1s' \in \mathcal{L}\{r\}$ , then  $c_1 \in \mathcal{L}\{\text{lhead}(r)\}$ .

**Condition D** (Correctness of Tail). *If*  $c_1s' \in \mathcal{L}\{r\}$  *then*  $s' \in \mathcal{L}\{\text{ltail}(r)\}$ .

For example, we conjecture (but do not here prove) that the definitions below satisfy these conditions. Note that these are slightly amended relative to the published paper.

**Definition 1** (Definition of lhead(r)). We first define an auxiliary relation that determines the set of characters that the head might be, tracking the remainder of any sequences that appear:

$$\begin{aligned} \mathsf{Ihead}(\epsilon,\epsilon) &= \emptyset \\ \mathsf{Ihead}(\epsilon,r') &= \mathsf{Ihead}(r',\epsilon) \\ \mathsf{Ihead}(a,r') &= \{a\} \\ \mathsf{Ihead}(r_1 \cdot r_2,r') &= \mathsf{Ihead}(r_1,r_2 \cdot r') \\ \mathsf{Ihead}(r_1+r_2,r') &= \mathsf{Ihead}(r_1,r') \cup \mathsf{Ihead}(r_2,r') \\ \mathsf{Ihead}(r^*,r') &= \mathsf{Ihead}(r,\epsilon) \cup \mathsf{Ihead}(r',\epsilon) \end{aligned}$$

We define  $lhead(r) = a_1 + a_2 + ... + a_i$  iff  $lhead(r, \epsilon) = \{a_1, a_2, ..., a_i\}$ .

**Definition 2** (Brzozowski's Derivative). The *derivative of* r *with respect to* s is denoted by  $\delta_s(r)$  and is  $\delta_s(r) = \{t | st \in \mathcal{L}\{r\}\}.$ 

**Definition 3** (Definition of Itail(r)). If Ihead $(r, \epsilon) = \{a_1, a_2, ..., a_i\}$ , then we define Itail $(r) = \delta_{a_1}(r) + \delta_{a_2}(r) + ... + \delta_{a_i}(r)$ .

#### 3.2 Replacement

The following correctness condition must hold for any definition of  $lreplace(r, r_1, r_2)$ .

**Condition E** (Replacement Correctness). *If*  $s_1 \in \mathcal{L}\{r_1\}$  *and*  $s_2 \in \mathcal{L}\{r_2\}$  *then* 

$$\mathsf{replace}(r; s_1; s_2) \in \mathcal{L}\{\mathsf{lreplace}(r, r_1, r_2)\}$$

We do not give a particular definition for  $lreplace(r, r_1, r_2)$  here.

#### 3.3 Small Step Semantics of $\lambda_{RS}$

Figure 7 specifies a small-step operational semantics for  $\lambda_{RS}$ .

**Lemma F** (Canonical Forms). *If*  $\emptyset \vdash v : \sigma$  *then:* 

- 1. If  $\sigma = \operatorname{stringin}[r]$  then  $v = \operatorname{rstr}[s]$  and  $s \in \mathcal{L}\{r\}$ .
- 2. If  $\sigma = \sigma_1 \rightarrow \sigma_2$  then  $v = \lambda x.e'$ .

*Proof.* By inspection of the static and dynamic semantics.

**Lemma G** (Progress). If  $\emptyset \vdash e : \sigma$  either e = v for some v or  $e \mapsto e'$  for some e'.

*Proof.* The proof proceeds by rule induction on the derivation of  $\emptyset \vdash e : \sigma$ .

 $\lambda$  fragment. Cases SS-T-Var, SS-T-Abs, and SS-T-App are exactly as in a proof of progress for the simply typed lambda calculus.

**S-T-Stringin-I.** Suppose  $\emptyset \vdash \mathsf{rstr}[s]$  :  $\mathsf{stringin}[s]$ . Then  $e = \mathsf{rstr}[s]$ .

**S-T-Concat**. Suppose  $\emptyset \vdash \mathsf{rconcat}(e_1; e_2) : \mathsf{stringin}[r_1 \cdot r_2]$  and  $\emptyset \vdash e_1 : \mathsf{stringin}[r_1]$  and  $\emptyset \vdash e_2 : \mathsf{stringin}[r_2]$ . By induction,  $e_1 \mapsto e_1'$  or  $e_1 = v_1$  and similarly,  $e_2 \mapsto e_2'$  or  $e_2 = v_2$ . If  $e_1$  steps, then SS-E-Concat-Left applies and so  $\mathsf{rconcat}(e_1; e_2) \mapsto \mathsf{rconcat}(e_1'; e_2)$ . Similarly, if  $e_2$  steps then e steps by SS-E-Concat-Right.

In the remaining case,  $e_1 = v_1$  and  $e_2 = v_2$ . But then it follows by Canonical Forms that  $e_1 = rstr[s_1]$  and  $e_2 = rstr[s_2]$ . Finally, by SS-E-Concat, rconcat( $rstr[s_1]$ ;  $rstr[s_2]$ )  $\mapsto rstr[s_1s_2]$ .

**S-T-Case.** Suppose  $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$  and  $\emptyset \vdash e_1 : \mathsf{stringin}[r]$ . By induction and Canonical Forms it follows that  $e_1 \mapsto e_1'$  or  $e_1 = \mathsf{rstr}[s]$ . In the former case, e steps by S-E-Case-Left. In the latter case, note that  $s = \epsilon$  or s = at for some string t. If  $s = \epsilon$  then e steps by S-E-Case- $\epsilon$ -Val, and if s = at the e steps by S-E-Case-Concat.

**S-T-Replace**. Suppose  $e = \text{rreplace}[r](e_1; e_2), \emptyset \vdash e : \text{stringin}[\texttt{lreplace}(r, r_1, r_2)]$  and:

$$\emptyset \vdash e_1 : \mathsf{stringin}[r_1]$$

$$\emptyset \vdash e_2 : \mathsf{stringin}[r_2]$$

By induction on (1),  $e_1 \mapsto e_1'$  or  $e_1 = v_1$  for some  $e_1'$ . If  $e_1 \mapsto e_1'$  then e steps by SS-E-Replace-Left. Similarly, if  $e_2$  steps then e steps by SS-E-Replace-Right. The only remaining case is where  $e_1 = v_1$  and also  $e_2 = v_2$ . By Canonical Forms,  $e_1 = \mathsf{rstr}[s_1]$  and  $e_2 = \mathsf{rstr}[s_2]$ . Therefore,  $e \mapsto \mathsf{rstr}[\mathsf{replace}(r;s_1;s_2)]$  by SS-E-Replace.

**S-T-SafeCoerce**. Suppose that  $\emptyset \vdash \mathsf{rcoerce}[r](e_1)$ :  $\mathsf{stringin}[r]$ . and  $\emptyset \vdash e_1$ :  $\mathsf{stringin}[r']$  for  $\mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}$ . By induction,  $e_1 = v_1$  or  $e_1 \mapsto e'_1$  for some  $e'_1$ . If  $e_1 \mapsto e'_1$  then e steps by SS-E-SafeCoerce-Step. Otherwise,  $e_1 = v$  and by Canonical Forms  $e_1 = \mathsf{rstr}[s]$ . In this case,  $e = \mathsf{rcoerce}[r](\mathsf{rstr}[s]) \mapsto \mathsf{rstr}[s]$  by SS-E-SafeCoerce.

**S-T-Check** Suppose that  $\emptyset \vdash \mathsf{rcheck}[r](e_0; x.e_1; e_2)$  :  $\mathsf{stringin}[r]$  and:

$$\emptyset \vdash e_0 : \mathsf{stringin}[r_0]$$

(4) 
$$\emptyset, x : \mathsf{stringin}[r] \vdash e_1 : \sigma$$

$$\emptyset \vdash e_2 : \sigma$$

By induction,  $e_0 \mapsto e_0'$  or  $e_0 = v$ . In the former case e steps by SS-E-Check-StepLeft. Otherwise,  $e_0 = \mathsf{rstr}[s]$  by Canonical Forms. By Lemma B part 2, either  $s \in \mathcal{L}\{r_0\}$  or  $s \notin \mathcal{L}\{r_0\}$ . In the former case e takes a step by SS-E-Check-Ok. In the latter case e takes a step by SS-E-Check-NotOk.

**Assumption H** (Substitution). If  $\Psi, x : \sigma' \vdash e : \sigma$  and  $\Psi \vdash e' : \sigma'$ , then  $\Psi \vdash [e'/x]e : \sigma$ .

**Lemma I** (Preservation for Small Step Semantics). If  $\emptyset \vdash e : \sigma$  and  $e \mapsto e'$  then  $\emptyset \vdash e' : \sigma$ .

*Proof.* By induction on the derivation of  $e \mapsto e'$  and  $\emptyset \vdash e : \sigma$ .

 $\lambda$  fragment. Cases SS-E-AppLeft, SS-E-AppRight, and SS-E-AppAbs are exactly as in a proof of type safety for the simply typed lambda calculus.

**SS-E-Concat-Left.** Suppose  $e = \mathsf{rconcat}(e_1; e_2) \mapsto \mathsf{rconcat}(e'_1; e_2)$  and  $e_1 \mapsto e'_1$ . The only rule that applies is S-T-Concat, so  $\emptyset \vdash e_1$ : stringin $[r_1]$  and  $\emptyset \vdash e_2$ : stringin $[r_2]$ . By induction,  $\emptyset \vdash e'_1$ : stringin $[r_1]$ . Therefore, by S-T-Concat,  $\emptyset \vdash \mathsf{rconcat}(e'_1; e_2)$ : stringin $[r_1r_2]$ .

**SS-E-Concat-Right**. Suppose  $e = \operatorname{rconcat}(e_1; e_2) \mapsto \operatorname{rconcat}(e_1; e_2')$  and  $e_2 \mapsto e_2'$ . The only rule that applies is S-T-Concat, so  $\emptyset \vdash e_1$ : stringin $[r_1]$  and  $\emptyset \vdash e_2$ : stringin $[r_2]$ . By induction,  $\emptyset \vdash e_2'$ : stringin $[r_2]$ . Therefore, by S-T-Concat,  $\emptyset \vdash \operatorname{rconcat}(e_1; e_2')$ : stringin $[r_1r_2]$ .

**SS-E-Concat**. Suppose  $\operatorname{rconcat}(\operatorname{rstr}[s_1];\operatorname{rstr}[s_2]) \mapsto \operatorname{rstr}[s_1s_2]$ . The only applicable rule is S-T-Concat, so  $\emptyset \vdash \operatorname{rstr}[s_1]:\operatorname{stringin}[r_1]$  and  $\emptyset \vdash \operatorname{rstr}[s_2]:\operatorname{stringin}[r_2]$  and  $\emptyset \vdash \operatorname{rconcat}(\operatorname{rstr}[s_1];\operatorname{rstr}[s_2]):\operatorname{stringin}[r_1 \cdot r_2]$ . By Canonical Forms,  $s_1 \in \mathcal{L}\{r_1\}$  and  $s_2 \in \mathcal{L}\{r_2\}$  from which it follows by Lemma B that  $s_1s_2 \in \mathcal{L}\{r_1 \cdot r_2\}$ . Therefore,  $\emptyset \vdash \operatorname{rstr}[s_1s_2]:\operatorname{stringin}[r_1 \cdot r_2]$  by S-T-Rstr.

**S-E-Case-Left**. Suppose  $e \mapsto \mathsf{rstrcase}(e_1'; e_2; x, y.e_3)$  and  $\emptyset \vdash e : \sigma$  and  $e_1 \mapsto e_1'$ . The only rule that applies is S-T-Case, so:

$$\emptyset \vdash e_1 : \mathsf{stringin}[r]$$

$$\emptyset \vdash e_2 : \sigma$$

(8) 
$$\emptyset, x : \text{stringin}[\text{lhead}(r)], y : \text{stringin}[\text{ltail}(r)] \vdash e_3 : \sigma$$

By (6) and the assumption that  $e_1 \mapsto e_1'$ , it follows by induction that  $\emptyset \vdash e_1'$ : stringin[r]. This fact together with (7) and (8) implies by S-T-Case that  $\emptyset \vdash \mathsf{rstrcase}(e_1'; e_2; x, y.e_3) : \sigma$ .

**SS-E-Case-** $\epsilon$ **-Val**. Suppose  $\operatorname{rstrcase}(e_0; e_2; x, y.e_3) \mapsto e_2$ . The only rule that applies is S-T-Case, so  $\emptyset \vdash e_2 : \sigma$ .

**SS-E-Case-Concat**. Suppose that  $e = \mathsf{rstrcase}(\mathsf{rstr}[as]; e_2; x, y.e_3) \mapsto [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3$  and that  $\emptyset \vdash e : \sigma$ . The only rule that applies is S-T-Case so:

$$\emptyset \vdash \mathsf{rstr}[as] : \mathsf{stringin}[r]$$

$$\emptyset \vdash e_2 : \sigma$$

(11) 
$$\emptyset$$
,  $x$ : stringin[lhead $(r)$ ],  $y$ : stringin[ltail $(r)$ ]  $\vdash e_3 : \sigma$ 

We know that  $as \in \mathcal{L}\{r\}$  by Canonical Forms on (9) Therefore,  $a \in \mathcal{L}\{\mathsf{lhead}(r)\}$  by Condition C and  $s \in \mathcal{L}\{\mathsf{ltail}(r)\}$  by Condition D.

From these facts about a and s we know by S-T-Rstr that  $\emptyset \vdash \mathsf{rstr}[a] : \mathsf{stringin}[\mathsf{lhead}(r)]$  and  $\emptyset \vdash \mathsf{rstr}[s] : \mathsf{stringin}[\mathsf{ltail}(r)]$ . It follows by Assumption H that  $\emptyset \vdash [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3 : \sigma$ .\_\_\_\_\_

Cases SS-E-Replace-Left, SS-E-Replace-Right, SS-E-Check-StepLeft, SS-E-SafeCoerce-Step, SS-E-Check-StepRight.

#### Case SS-E-Replace.

Suppose  $e = \text{rreplace}[r](\text{rstr}[s_1]; \text{rstr}[s_2]) \mapsto \text{rstr}[\text{replace}(r; s_1; s_2)]$ . The only applicable rule is S-T-Replace, so

$$\emptyset \vdash \mathsf{rstr}[s_1] : \mathsf{stringin}[r_1]$$
  
 $\emptyset \vdash \mathsf{rstr}[s_2] : \mathsf{stringin}[r_2]$ 

Cyrus stopped here

todo

By conanical forms,  $s_1 \in \mathcal{L}\{r_1\}$  and  $s_2 \in \mathcal{L}\{r_2\}$ . Therefore,  $\texttt{lreplace}(r, s_1, s_2) \in \mathcal{L}\{\texttt{lreplace}(r, r_1, r_2)\}$  by Theorem E. It is finally derivable by S-T-Rstr that:

 $\emptyset \vdash \mathsf{rstr}[\mathsf{lreplace}(r, s_1, s_2)] : \mathsf{stringin}[\mathsf{lreplace}(r, r_1, r_2)].$ 

Case SS-E-SafeCoerce. Suppose that  $\mathsf{rcoerce}[r](\mathsf{rstr}[s_1]) \mapsto \mathsf{rstr}[s_1]$ . The only applicable rule is S-T-SafeCoerce, so  $\emptyset \vdash \mathsf{rcoerce}[r](s_1)$ :  $\mathsf{stringin}[r]$ . By Canonical Forms,  $s \in \mathcal{L}\{r\}$ . Therefore,  $\emptyset \vdash \mathsf{rstr}[s]$ :  $\mathsf{stringin}[r]$ .

**Case SS-E-Check-Ok.** Suppose  $\operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_1; e_2) \mapsto [\operatorname{rstr}[s]/x]e_1, s \in \mathcal{L}\{r\}, \text{ and } \emptyset \vdash \operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_1; e_2) : \sigma.$  By inversion of S-T-Check,  $x : \operatorname{stringin}[r] \vdash e_1 : \sigma.$  Note that  $s \in \mathcal{L}\{r\}$  implies that  $s : \operatorname{stringin}[r]$  by S-T-RStr. Therefore,  $\emptyset \vdash [\operatorname{rstr}[s]/x]e_1 : \sigma.$ 

**Case SS-E-Check-NotOk**. Suppose  $\operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_1; e_2) \mapsto e_2, s \notin \mathcal{L}\{r\}$ , and  $\emptyset \vdash \operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_1; e_2) : \sigma$ . The only applicable rule is S-T-Check, so  $\emptyset \vdash e_2 : \sigma$ .

**Theorem J** (Type Safety for small step semantics.). If  $\emptyset \vdash e : \sigma$  then either  $e \lor alor e \mapsto^* e'$  and  $\emptyset \vdash e' : \sigma$ .

*Proof.* Follows directly from progress and preservation.

#### 3.3.1 Semantic Correspondence between Big and Small Step Semantics for $\lambda_{RS}$

Before extending the previous theorem to the big step semantics, we first establish a correspondence between the big step semantics in Figure 7 and the small step semantics in Figure 5.

**Lemma K.** If  $e \Downarrow v$  and  $e \mapsto e'$  then  $e' \Downarrow v$ .

*Proof.* By induction on the structure of e.

Case  $e = e_1(e_2)$ . The only applicable rule is S-E-App, so  $e_1 \Downarrow \lambda x.e_3$  and  $e_2 \Downarrow v_2$  such that  $[v_2/x]e_3 \Downarrow v$ .

The term  $e = e_1(e_2)$  may step by three rules.

First, if e steps by L-E-AppLeft then  $e_1(e_2) \mapsto e_1'(e_2)$ . By induction,  $e_1' \Downarrow \lambda x.e_3$ . By S-E-App,  $e_1'(e_2) \Downarrow \lambda x.e_3$ .

Second, if e steps by L-E-AppRight then  $e_1(e_2) \mapsto e_1(e_2')$ . By induction,  $e_2' \Downarrow v_2$ . By S-E-App,  $e_1(e_2') \Downarrow \lambda x.e_3$ .

Third, if e steps by L-E-AppAbs then  $e=(\lambda x.e_3)(v_2)\mapsto [v_2/x]e_3$ . By induction,  $e=(\lambda x.e')v_2 \Downarrow v$ .

**Theorem L** (Semantic Correspondence for  $\lambda_{RS}$  (Part I)). If  $e \downarrow v$  then  $e \mapsto^* v$ .

*Proof.* We proceed by structural induction on e.

Case  $e = \lambda x.e_1$ . The only applicable rule is S-E-Abs, so  $v = \lambda x.e_1$ . Note that  $\lambda x.e_1 \mapsto^* \lambda x.e_1$  by RT-Refl.

Case  $e = e_1(e_2)$ . The only applicable rule is S-E-App. By inversion:

$$e_1 \Downarrow \lambda x. e'_1$$

$$e_2 \Downarrow v_2$$

$$[v_2/x]e'_1 \Downarrow v$$

From which it follows by induction that:

$$e_1 \mapsto^* \lambda x. e_1'$$

$$e_2 \mapsto^* v_2$$

$$[v_2/x]e_1' \mapsto^* v$$

If  $e_1 = \lambda x.e_1'$  and  $e_2 = v_2$  (henceforth the reflexive case) then  $e \mapsto [v_2/x]e_1'$  and the conclusion follows by RT-Trans.

If  $e_1 \mapsto \lambda x.e_1'$  then  $e_1(e_2) \mapsto (\lambda x.e_1')(e_2) \mapsto [e_2/x]e_1'$  and the conclusion follows by two applications of RT-Trans.

If  $e_1 \mapsto^k \lambda x.e_2'$  then  $e_1 \mapsto e'$  so  $e_1(e_2) \mapsto e'(e_2)$  and by Lemma K  $e'(e_2) \mapsto^* v$ . So  $e_1(e_2) \mapsto e'(e_2) \mapsto^* v$ ; it follows by RT-Trans that  $e_1(e_2) \mapsto v$ .

Note that the following rule is derivable by repeating applications of the left and right compatibility rules for application:

not modified below this line...

$$\frac{L^*\text{-APP}}{e_1 \mapsto^* e_1'} \quad e_2 \mapsto^* e_2'$$

$$e_1(e_2) \mapsto^* e_1'(e_2')$$

From these facts and L-AppAbs, we may establish that  $e_1(e_2) \mapsto^* (\lambda x. e_2)(v_2) \mapsto [v_2/x]e_2$ . Note that  $[v_2/x]e_2 \mapsto^* v$ , so by RT-Trans it follows that  $e = e_1(e_2) \mapsto^* v$ .

Case  $e = \mathsf{rstr}[s]$ . The only applicable rule is S-E-RStr, so  $v = \mathsf{rstr}[s]$ . By RT-Refl,  $\mathsf{rstr}[s] \mapsto^* \mathsf{rstr}[s]$ .

Case  $e = \mathsf{rconcat}(e_1; e_2)$ . The only applicable rule is S-E-Concat, so  $v = \mathsf{rstr}[s_1s_2]$ . By inversion,  $e_1 \Downarrow \mathsf{rstr}[s_1]$  and  $e_2 \Downarrow \mathsf{rstr}[s_2]$ . By induction,  $e_1 \mapsto^* \mathsf{rstr}[s_1]$  and  $e_2 \mapsto^* \mathsf{rstr}[s_2]$ . Note that the rule following is derivable:

$$\frac{\text{SS-E-Concat-LR*}}{e_1 \mapsto^* e_1'} \underbrace{e_2 \mapsto^* e_2'}_{\text{rconcat}(e_1;e_2) \mapsto^* \text{rconcat}(e_1';e_2')}$$

From these facts, it follows that  $\mathsf{rconcat}(e_1; e_2) \mapsto^* \mathsf{rconcat}(\mathsf{rstr}[s_1]; \mathsf{rstr}[s_2])$ . Finally,  $\mathsf{rconcat}(\mathsf{rstr}[s_1]; \mathsf{rstr}[s_2]) \mapsto^* \mathsf{rstr}[s_1s_2]$  by SS-E-Concat. By RT-Step, it follows that  $\mathsf{rconcat}(e_1; e_2) \mapsto^* \mathsf{rstr}[s_1s_2]$ .

Case  $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$ .

There are two subcases. For the first, suppose  $\mathsf{rstrcase}(e_1; e_2; x, y.e_3) \Downarrow v$  was finally derived by S-E-Case- $\epsilon$ . By inversion:

$$e_1 \Downarrow \mathsf{rstr}[\epsilon]$$
  
 $e_2 \Downarrow v$ 

from which it follows by induction that:

$$e_1 \mapsto^* \mathsf{rstr}[\epsilon]$$
$$e_2 \mapsto^* v$$

Note that the following rule is derivable:

SS-E-CASE-LR\*
$$e_1 \mapsto^* e_1' \qquad e_2 \mapsto^* e_2'$$

$$rstrcase(e_1; e_2; x, y.e_3) \mapsto^* rstrcase(e_1'; e_2'; x, y.e_3)$$

From these facts is follows that  $e \mapsto^* \mathsf{rstrcase}(\mathsf{rstr}[\epsilon]; v; x, y.e_3)$ . By S-E-Case- $\epsilon$ -Val and RT-Step it follows that  $e \mapsto^* v$ .

Now consider the other case where  $\mathsf{rstrcase}(e_1; e_2; x, y.e_3) \Downarrow v$  was finally derived by S-E-Case-Concat. By inversion,  $e_1 \Downarrow \mathsf{rstr}[as]$  and  $[\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3 \Downarrow v$ . From these facts it follows by induction that  $e_1 \mapsto^* \mathsf{rstr}[as]$  and  $[\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3 \mapsto^* v$ .

By the first of these facts, it is derivable via SS-E-Case-LR\* that  $e\mapsto^* \mathrm{rstrcase}(e_1'; \mathrm{rstr}[as]; x, y.e_3)$ . SE-E-Case-Concat applies to this form, so by RT-Step we know  $e\mapsto^* [\mathrm{rstr}[a], \mathrm{rstr}[s]/x, y]e_3$ . Recall that  $[\mathrm{rstr}[a], \mathrm{rstr}[s]/x, y]e_3\mapsto^* v$ , so by RT-Trans we finally derive  $e\mapsto^* v$ .

Case  $e = \text{rreplace}[r](e_1; e_2)$ . There is only one applicable rule, so v = rstr[s] and by inversion it follows that:

$$e_1 \Downarrow \mathsf{rstr}[s_1]$$
  
 $e_2 \Downarrow \mathsf{rstr}[s_2]$ 

From which it follows by induction that:

$$e_1 \mapsto^* rstr[s_1]$$
  
 $e_2 \mapsto^* rstr[s_2]$ 

Furthermore, replace $(r; s_1; s_2) = s$  by induction. Note that the following rule is derivable:

$$\frac{\text{SS-E-Replace-LR*}}{e_1 \mapsto^* e_1' \qquad e_2 \mapsto^* e_2'} \\ \frac{e_1 \mapsto^* e_1' \qquad e_2 \mapsto^* replace[r](e_1'; e_2')}{\text{rreplace}[r](e_1'; e_2') \mapsto^* replace[r](e_1'; e_2')}$$

From these facts,  $rreplace[r](e_1; e_2) \mapsto^* rreplace[r](rstr[s_1]; rstr[s_2])$ .

Finally,  $rreplace[r](rstr[s_1]; rstr[s_2]) \mapsto replace(r; s_1; s_2)$ .

From these two facts we know via RT-Step that  $\operatorname{rreplace}[r](e_1; e_2) \mapsto^* \operatorname{rreplace}[r](e_1; e_2)$ . Recall that  $\operatorname{replace}(r; s_1; s_2) = s$ , from which the conclusion follows.

Case  $e = \text{rcoerce}[r](e_1)$ . In this case  $e \Downarrow v$  is only finally derivable via S-E-SafeCoerce. Therefore, v = rstr[s] and by inversion  $e_1 \Downarrow \text{rstr}[s]$ . By induction,  $e_1 \mapsto^* \text{rstr}[s]$ .

The following rule is derivable:

$$\frac{\text{SS-E-SafeCoerce-Step}}{e \mapsto^* e'} \frac{e \mapsto^* e'}{\text{rcoerce}[r](e) \mapsto^* \text{rcoerce}[r](e')}$$

Applying this rule at  $e_1 \mapsto^* \operatorname{rstr}[s]$  derives  $\operatorname{rcoerce}[r](e_1) \mapsto^* \operatorname{rcoerce}[r](\operatorname{rstr}[s])$ . In the final step,  $\operatorname{rcoerce}[r](\operatorname{rstr}[s]) \mapsto \operatorname{rstr}[s]$  by SS-E-SafeCoerce. From this fact, we may derive via RT-Trans that  $e \mapsto^* \operatorname{rstr}[s]$  as required.

**Case**  $e = \text{rcheck}[r](e_1; x.e_2; e_3).$ 

Note that the rule following is derivable:

$$\frac{\text{SS-E-Check-Step}}{e_1 \mapsto^* e_1' \qquad e_3 \mapsto^* e_3'} \\ \frac{e_1 \mapsto^* e_1' \qquad e_3 \mapsto^* \text{rcheck}[r](e_1'; x.e_2; e_3')}{\text{rcheck}[r](e_1'; x.e_2; e_3')}$$

There are two ways to finally derive  $e \downarrow v$ . In both cases,  $e_1 \downarrow \mathsf{rstr}[s]$  by inversion. Therefore, in both cases,  $e_1 \mapsto^* \mathsf{rstr}[s]$  by induction and so  $e \mapsto^* \mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_2; e_3)$  by SS-E-Check-Step.

Suppose  $e \Downarrow v$  is finally derived via SS-E-Check-Ok. By the facts mentioned above and SS-E-Check-Step,  $e \mapsto^* \operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_2; e_2)$ . Note that by inversion  $s \in \mathcal{L}\{r\}$ . Therefore, SS-E-Check-Ok applies and so by RT-Trans  $e \mapsto^* [\operatorname{rstr}[s]/x]e_1$ . By inversion,  $[\operatorname{rstr}[s]/x]e_1 \Downarrow v$ . Therefore, by induction and RT-Step  $e \mapsto^* v$  as required.

Suppose that  $e \Downarrow v$  is instead finally derived via SS-E-Check-NotOk. By inversion,  $e_3 \Downarrow v$  and by induction  $e_3 \mapsto^* v$ . From these facts at SS-E-Check-Step, it is derivable that  $e \mapsto^* \operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_2; v)$ .

Also by inversion,  $s \notin \mathcal{L}\{r\}$  and so SS-E-Check-NotOk applies. Therefore,  $\operatorname{rcheck}[r](\operatorname{rstr}[s]; x.e_2; v) \mapsto v$ .

The conclusion  $e \mapsto^* v$  follows from these facts by RT-Step.

**Theorem M** (Semantic Correspondence for  $\lambda_{RS}$  (Part II)). If  $\emptyset \vdash e : \sigma, e \mapsto^* v$  and v val then  $e \Downarrow v$ .

*Proof.* The proof proceeds by structural induction on e.

**Case**  $e = \text{concat}(e_1; e_2)$ . By inversion,  $\emptyset \vdash e_1 : \text{stringin}[r_1]$ . By Type Safety, Canonical Forms and Termination it follows that  $e_1 \mapsto^* \text{rstr}[s_1]$  for some  $s_1$ . By induction,  $e_1 \downarrow \text{rstr}[s_1]$ .

Similarly,  $e_2 \mapsto^* \operatorname{rstr}[s_2]$  and  $e_2 \Downarrow \operatorname{rstr}[s_2]$ .

Note that  $concat(e_1; e_2) \mapsto^* concat(rstr[s_1]; rstr[s_2]) \mapsto rstr[s_1s_2]$  by SS-E-Concat-LR\* and S-E-Concat. Therefore,  $e \mapsto^* rstr[s_1s_2]$  by RT-Step. So it suffices to show that  $e \Downarrow rstr[s_1s_2]$ .

Finally,  $e \Downarrow rstr[s_1s_2]$  follows via S-E-Concat from the facts that  $e_1 \Downarrow rstr[s_1]$  and  $e_2 \Downarrow rstr[s_2]$ . This completes the case.

Case  $e = \text{rreplace}[r](e_1; e_2)$ . By inversion of S-T-Replace,  $\emptyset \vdash e_1 : \text{stringin}[r_1]$  for some  $r_1$ . It follows by Type Safety, Termination and Canonical Forms that  $e_1 \mapsto^* \text{rstr}[s_1]$ . By induction,  $e_1 \Downarrow \text{rstr}[s_1]$ .

Similarly,  $e_2 \mapsto^* \mathsf{rstr}[s_2]$  and  $e_2 \Downarrow \mathsf{rstr}[s_2]$ .

Note that  $e\mapsto^* \operatorname{rreplace}[r](\operatorname{rstr}[s_1];\operatorname{rstr}[s_2])\mapsto \operatorname{rstr}[\operatorname{replace}(r;s_1;s_2)]$  by SS-Replace-LR\* and SS-E-Replace. Therefore  $e\mapsto^* \operatorname{rstr}[\operatorname{replace}(r;s_1;s_2)]$  by RT-Step.

It suffices to show  $e \Downarrow \mathsf{rstr}[\mathsf{replace}(r; s_1; s_2)]$ , which follows by S-E-Replace from the facts that  $e_1 \Downarrow \mathsf{rstr}[s_1]$  and  $e_2 \Downarrow \mathsf{rstr}[s_2]$ .

Case  $e = \operatorname{rstrcase}(e_1; e_2; x.y.e_3)$ . By inversion,  $\emptyset \vdash e_1 : \operatorname{stringin}[r]$  and  $e_2 : \sigma$ . By Type Safety, Canonical Forms and Termination  $e_1 \mapsto^* \operatorname{stringin}[s_1]$  and by induction  $e_1 \Downarrow \operatorname{stringin}[s_1]$ . Similarly,  $e_2 \mapsto^* v_2$  and  $\emptyset \vdash e_2 \Downarrow v_2$ .

By SS-E-Case-LR\*,  $rstrcase(e_1; e_2; x, y.e_3) \mapsto^* rstrcase(v_1; v_2; x, y.e_3)$ .

Note that either  $s_1 = \epsilon$  or  $s_1 = as$  because we define strings as either empty or finite sequences of characters. We proceed by cases.

If  $s_1 = \epsilon$  then  $\mathsf{rstrcase}(v; v_2; x, y.e_3) \mapsto v_2$  by SS-E-Case- $\epsilon$ . Therefore, by RT-Step,  $e \mapsto^* v_2$ . Recall  $e_1 \Downarrow \mathsf{rstr}[\epsilon]$  and  $e_2 \Downarrow v_2$ , which is enough to establish by S-E-Case- $\epsilon$  that  $e \Downarrow v_2$ .

If  $s_1 = as$  instead, then  $\mathsf{rstrcase}(\mathsf{rstr}[s_1]; v_2; x, y.e_3) \mapsto [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3$  by SS-E-Case-Concat. Inversion of the typing relation satisfies the assumptions necessary to appeal to termination. Therefore,

$$[\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3 \mapsto^* v \text{ for } v \text{ val.}$$

It follows by RT-Step that  $e \mapsto^* v$ .

Note that the substitution does not change the structure of  $e_3$ . So by induction,  $[rstr[a], rstr[s]/x, y]e_3 \Downarrow v$ . Recall that  $e_1 \Downarrow rstr[s_1]$  and so by S-E-Case it follows that  $e \Downarrow [a, s/x, y]e_3 \Downarrow v$ .

#### 3.4 Extension of Safety for Small Step Semantics

**Theorem 4** (Type Safety). *If*  $\emptyset \vdash e : \sigma$  *and*  $e \Downarrow e'$  *then*  $\emptyset \vdash e' : \sigma$ .

*Proof.* If  $\emptyset \vdash e : \sigma$  then  $e \mapsto^* e'$ . Therefore,  $e \Downarrow e'$  by part 2 of the semantic correspondence theorem. Since  $\emptyset \vdash e : \sigma$  and  $e \mapsto^* e'$ , it follows that  $\emptyset \vdash e' : \sigma$  by type safety for the small step semantics.

#### 3.4.1 The Security Theorem

**Theorem 5** (Correctness of Input Sanitation for  $\lambda_{RS}$ ). If  $\emptyset \vdash e$ : stringin[r] and  $e \Downarrow rstr[s]$  then  $s \in \mathcal{L}\{r\}$ .

*Proof.* If  $\emptyset \vdash e$ : stringin[r] and  $e \Downarrow rstr[s]$  then  $\emptyset \vdash rstr[s]$ : stringin[r] by Type Safety. By Caonical Forms,  $s \in \mathcal{L}\{r\}$ .

The cases for coercion and checking are straightforward.

### 4 Proofs of Lemmas and Theorems About $\lambda_P$

This section follows the same structure as the safety proof for  $\lambda_{RS}$  – we prove type safety for a small-step semantics, prove a semantic correspondence, and then transfer the safety result to the big-step semantics in the paper.

Our proofs here are more sketchy than in the  $\lambda_{RS}$  section both because the proof method is already demonstrated, and because  $\lambda_P$  is a fairly straightforward extension of  $\lambda$  for which these results are already estbalished.

Lemma 6 (Canonical Forms for Target Language).

- If  $\emptyset \vdash \iota$ : regex then  $\iota \Downarrow \mathsf{rx}[r]$  such that r is a well-formed regular expression.
- If  $\emptyset \vdash \iota$ : string then  $\iota \Downarrow str[s]$ .

**Theorem 7** (Progress). *If*  $\emptyset \vdash \iota : \tau$  *either*  $\iota = \dot{v}$  *or*  $\iota \mapsto \iota'$  *for some*  $\iota'$ .

*Proof.* The proof proceeds by induction on the typing assumption. We elide contexts where context has obviously not changed.

Consider only the string and regex (non- $\lambda$ ) fragments of  $\lambda_P$ .

Regular expressions are always values.

For strings, the concatenation cases are trivial.

**P-T-Case**. Suppose  $\emptyset \vdash \mathsf{strcase}(\iota_1; \iota_2; x, y. \iota_3)$ . By inversion,  $\iota_1 : \mathsf{string}$  and so either  $\iota_1 \mapsto \iota_1'$  or by canonical forms,  $\iota_1 = \mathsf{str}[s_1]$ . Similarly,  $\iota_2 \mapsto \iota_2'$  or else  $\iota_2 = \mathsf{str}[s_2]$ . In the former cases, progress occurs via the compatibility rules. in the case where both are string values, progress occurs via the case concatenation rule.

**P-T-Replace.** Suppose  $\emptyset \vdash \text{replace}(\iota_1; \iota_2; \iota_3)$ . By inversion,  $\iota_1 : \text{regex}$  and so by canonical forms  $\iota_1 = \text{rx}[r]$ . By inversion,  $\iota_2 : \text{string}$  and so by induction either  $\iota_2 \mapsto \iota'_2$  or else  $\iota_2 = \text{str}[s_2]$  for some string  $s_2$ . Similarly, either  $\iota_3$  steps or else  $\iota_3 = \text{str}[s_3]$ . In case any steps occur, progress occurs. In the remaining case, PP-E-Replace applies and so progress occurs.

**P-T-Check**. Finally, suppose  $\emptyset \vdash \iota_x \iota_1 \iota_2 \iota_3$ . In case any of these step, then progress occurs. In the remaining cases, applications of inversion and canonical forms for each  $\iota_x$  and  $\iota_1$  implies that the term at hand equals  $rx[r]str[s]\iota_2\iota_3$ , which evaluates to either  $\iota_2$  or  $\iota_3$ .

**Theorem 8** (Preservation). *If*  $\emptyset \vdash \iota : \tau \text{ and } \iota \mapsto \iota' \text{ then } \emptyset \vdash \iota' : \tau$ .

Proof.

**Theorem 9** (Safety for Small-Step Semantics). *If*  $\iota \Downarrow \dot{v}$  *and*  $\iota \mapsto \iota'$  *then*  $\iota' \Downarrow \dot{v}$ .

*Proof.* A direct result from progress and preservation.

We only prove semantic correspondence in one direction; again, whereas  $\lambda_{RS}$  proofs were detailed, here we provide a less verbosy proof.

**Theorem 10** (Semantic Correspondence). *If*  $\iota \mapsto^* \iota'$  *then*  $\iota \Downarrow \iota'$ .

Fill in if there's time. Otherwise there's nothing interesting going on.

*Proof.* By induction on the structure of  $\iota$ .

**Theorem 11** (Safety for  $\lambda_P$ ). *If*  $\emptyset \vdash \iota : \tau$  *then*  $\iota \Downarrow \dot{v}$  *and*  $\emptyset \vdash \dot{v} : \tau$ .

*Proof.* If  $\emptyset \vdash \iota : \tau$  then  $\iota \mapsto^* \iota'$ . Therefore,  $\iota \Downarrow \iota'$  by part 2 of the semantic correspondence theorem. Since  $\emptyset \vdash \iota : \tau$  and  $\iota \mapsto^* \iota'$ , it follows that  $\emptyset \vdash \iota' : \tau$  by type safety for the small step semantics.

We can also define canonical forms for regular expressions and strings in the usual way:

#### 5 Proofs and Lemmas and Theorems About Translation

**Theorem 12** (Translation Correctness). *If*  $\Theta \vdash e : \sigma$  *then there exists an*  $\iota$  *such that*  $\llbracket e \rrbracket = \iota$  *and*  $\llbracket \Theta \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$ . *Furthermore,*  $e \Downarrow v$  *and*  $\iota \Downarrow \dot{v}$  *such that*  $\llbracket v \rrbracket = \dot{v}$ .

*Proof.* We present a proof by induction on the derivation that  $\Theta \vdash e : \sigma$ . we write  $e \leadsto \iota$  as shorthand for the final property.

Case  $e = \mathsf{rstr}[s]$ . Suppose  $\Theta \vdash \mathsf{rstr}[s] : \sigma$ .

By examination the syntactic structure of conclusions in the relation S-T, we know this is true just in case  $\sigma = \text{stringin}[r]$  for some r such that  $s \in \mathcal{L}\{r\}$ ; and of course, there is always such an r.

There are no free variables in rstr[s], so we might as well proceed from the fact that  $\emptyset \vdash rstr[s]$ : stringin[r].

By definition of the translation ( $[\![\cdot]\!]$ ) the following statements hold:

- [rstr[s]] = strings
- [stringin[r]] = string
- $[\![\emptyset]\!] = \emptyset$

Note that  $\emptyset \vdash \text{string } s$ : string by P-T-Str. Recall that contexts are standard and, in particular, can be weakened. So since  $\llbracket \Theta \rrbracket$  is either a weakening of  $\emptyset$  or  $\emptyset$  itself,  $\llbracket \Theta \rrbracket \vdash \text{str}[s]$ : string by weakening.

Summarily, strings is a term of  $\lambda_P$  such that  $\llbracket \Theta \rrbracket \vdash \mathsf{string} s : \llbracket \sigma \rrbracket$ 

It remains to be shown that there exist  $v, \dot{v}$  such that  $\mathsf{rstr}[s] \Downarrow v$ ,  $\mathsf{string} s \Downarrow \dot{v}$ , and  $[\![v]\!] = \dot{v}$ . But this is immediate because each term evaluates to itself and we have already established the equality.

Case  $e = \text{rconcat}(e_1; e_2)$ . This case is an obvious appeal to induction.

Case  $e = \text{rstrcase}(e_1; e_2; x, y.e_3)$ . This case relies on our definition of context translation.

Suppose  $\Psi \vdash \mathsf{rstrcase}(e_1; e_2; x, y.e_3) : \sigma$ . By inversion of the typing relation it follows that  $\Psi \vdash e_1 : \mathsf{stringin}[r], \Psi \vdash e_2 : \sigma \text{ and } \Psi, x : \mathsf{stringin}[\mathsf{lhead}(r)], y : \mathsf{stringin}[\mathsf{ltail}(r)] \vdash e_3 : \sigma$ .

By induction, there exists an  $\iota_1$  such that  $\llbracket e_1 \rrbracket = \iota_1$ ,  $\llbracket \Psi \rrbracket \vdash \iota_1 : \llbracket \sigma \rrbracket$ , and  $e_1 \leadsto \iota_1$ . Similarly for  $e_2$  and some  $\iota_2$ .

By canonical forms,  $e_1 \Downarrow \mathsf{rstr}[s]$  and so  $\iota_1 \Downarrow \mathsf{str}[s]$  by  $\leadsto$ .

Choose  $\iota = \operatorname{concat}(\iota_1; \iota_2)x, y.\iota_3$  and note that by the properties established via induction,  $\llbracket e \rrbracket = \iota$  and  $\llbracket \Psi \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$ .

Suppose  $s = \epsilon$ . Then  $e \Downarrow v$  where  $e_2 \Downarrow v$  and  $\iota \Downarrow \dot{v}$  where  $\iota_2 \Downarrow \dot{v}$ . But recall that  $e_2 \leadsto v_2$  and so  $\llbracket v \rrbracket = \dot{v}$ .

Suppose otherwise that s=at for some character a and string t. Then  $e \Downarrow v$  where  $[a,t/x,y]e_3 \Downarrow v$ . Similarly,  $\iota \Downarrow \dot{v}$  where  $[a,t/x,y]\iota_3 \Downarrow \dot{v}$ 

**Theorem 13** (Correctness of Input Sanitation for Translated Terms). *If*  $\llbracket e \rrbracket = \iota$  *and*  $\emptyset \vdash e : \mathsf{stringin}[r]$  *then*  $\iota \Downarrow \mathsf{str}[s]$  *for*  $s \in \mathcal{L}\{r\}$ .

*Proof.* By Theorem 12 and the rules given,  $\iota \Downarrow \mathsf{str}[s]$  implies that  $e \Downarrow \mathsf{rstr}[s]$ . Theorem 5 together with the assumption that e is well-typed implies that  $s \in \mathcal{L}\{r\}$ .

## References

[1] N. Fulton, C. Omar, and J. Aldrich. Statically typed string sanitation inside a python. SPLASH '14. ACM, 2014.

$$r ::= \epsilon \mid . \mid a \mid r \cdot r \mid r + r \mid r *$$
  $a \in \Sigma$ 

**Figure 1:** Regular expressions over the alphabet  $\Sigma$ .

$$\begin{array}{lll} \sigma & ::= & \sigma \rightarrow \sigma \mid \mathsf{stringin}[r] & \mathsf{source types} \\ e & ::= & x \mid v & \mathsf{source terms} \\ & \mid & \mathsf{rconcat}(e;e) \mid \mathsf{rstrcase}(e;e;x,y.e) & s \in \Sigma^* \\ & \mid & \mathsf{rreplace}[r](e;e) \mid \mathsf{rcoerce}[r](e) \mid \mathsf{rcheck}[r](e;x.e;e) \\ \\ v & ::= & \lambda x.e \mid \mathsf{rstr}[s] & \mathsf{source values} \end{array}$$

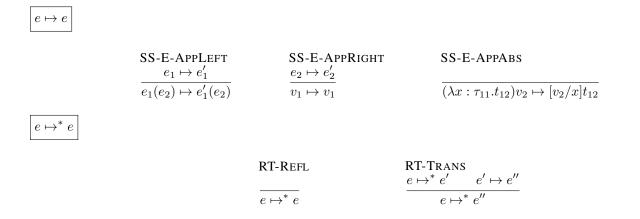
**Figure 2:** Syntax of  $\lambda_{RS}$ .

$$\tau \ \ \, ::= \tau \rightarrow \tau \ | \ \, \text{string} \ | \ \, \text{regex} \qquad \qquad \text{target types}$$
 
$$\iota \ \ \, ::= x \ | \ \, \dot{v} \\ \ \, | \ \, \text{concat}(\iota;\iota) \ | \ \, \text{strcase}(\iota;\iota;x,y.\iota) \\ \ \, | \ \, \text{rx}[r] \ | \ \, \text{replace}(\iota;\iota;\iota) \ | \ \, \text{check}(\iota;\iota;\iota)$$
 
$$\dot{v} \ \ \, ::= \lambda x.\iota \ | \ \, \text{str}[s] \ | \ \, \text{rx}[r]$$
 
$$\text{target values}$$

**Figure 3:** Syntax for the target language,  $\lambda_P$ , containing strings and statically constructed regular expressions.

**Figure 4:** Typing rules for  $\lambda_{RS}$ . The typing context  $\Psi$  is standard.

**Figure 5:** Big step semantics for  $\lambda_{RS}$ .



**Figure 6:** Call-by-name small step Semantics for  $\lambda$  and its reflexive, transitive closure.

**Figure 7:** Small step semantics for  $\lambda_{RS}$ . Extends 6.

**Figure 8:** Typing rules for  $\lambda_P$ . The typing context  $\Theta$  is standard.

$$\begin{array}{lll} \text{P-E-CONCAT} & \text{P-E-CASE-}\epsilon & \text{P-E-CASE-CONCAT} \\ \underline{\iota_1 \Downarrow \mathsf{str}[s_1]} & \iota_2 \Downarrow \mathsf{str}[s_2] & \underline{\iota_1 \Downarrow \mathsf{str}[\epsilon]} & \iota_2 \Downarrow \dot{v_2} \\ \hline{\mathsf{concat}}(\iota_1; \iota_2) \Downarrow \mathsf{str}[s_1s_2] & \underline{\mathsf{strcase}}(\iota_1; \iota_2; x, y. \iota_3) \Downarrow \dot{v_2} & \underline{\mathsf{strcase}}(\iota_1; \iota_2; x, y. \iota_3) \Downarrow \dot{v} \end{array}$$

$$\begin{array}{lll} & \text{P-E-Replace} \\ & \iota_1 \Downarrow \mathsf{rx}[r] & \iota_2 \Downarrow \mathsf{str}[s_2] & \iota_3 \Downarrow \mathsf{str}[s_3] \\ & \text{replace}(\iota_1; \iota_2; \iota_3) \Downarrow \mathsf{str}[\mathsf{replace}(r; s_2; s_3)] \end{array} & \begin{array}{lll} & \text{P-E-CHECK-OK} \\ & \iota_x \Downarrow \mathsf{rx}[r] & \iota \Downarrow \mathsf{str}[s] & s \in \mathcal{L}\{r\} & \iota_1 \Downarrow \dot{v_1} \\ & & \text{check}(\iota_x; \iota; \iota_1; \iota_2) \Downarrow \dot{v_1} \end{array}$$

$$\frac{\text{P-E-CHECK-NOTOK}}{\iota_x \Downarrow \mathsf{rx}[r]} \frac{\iota_v \Downarrow \mathsf{str}[s] \qquad s \not\in \mathcal{L}\{r\} \qquad \iota_2 \Downarrow \dot{v_2}}{\mathsf{check}(\iota_x; \iota; \iota_1; \iota_2) \Downarrow \dot{v_2}}$$

**Figure 9:** Big step semantics for  $\lambda_P$ 

$$\iota \Downarrow \dot{v}$$

**Figure 10:** Small step semantics for  $\lambda_P$  (extends L-E rules)

$$\begin{array}{c|c} & TR-T-STRING & TR-T-ARROW & \hline \llbracket \sigma \rrbracket = \tau_1 & \llbracket \sigma_2 \rrbracket = \tau_2 \\ \hline \llbracket \Psi \rrbracket = \Theta & & TR-T-CONTEXT-EMP & TR-T-CONTEXT-EXT & \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau_1 & \llbracket \sigma_2 \rrbracket = \tau_2 \\ \hline \llbracket \Psi \rrbracket = \theta & & TR-T-CONTEXT-EMP & TR-T-CONTEXT-EXT & \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Theta & \llbracket \sigma_1 \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \Psi \end{bmatrix} = \Psi_1 & \Psi_1 & \Psi_2 \end{bmatrix} = \Psi_2 & \Psi_2 & \Psi_3 \end{bmatrix} = \Psi_3 & \Psi_3$$

**Figure 11:** Translation from source terms (e) to target terms  $(\iota)$ .