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, and additional discussion of the paper's technical content.

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



1 Terminology and Notation

Theorems and lemmas appearing in [1] are numbered, while supporting facts appearing only in the Technical Report are lettered. Numbered items correspond to the numbering in [1].

2 λ_{RS}

This section presents proofs of lemmas and theorems about the type systems presented in [1], the accompanying paper. In addition, we provide some examples to help motivate and explain definitions.

To facilitate the type safety proof, we introduce a small step semantics for both λ_{RS} and λ_{P} . All theorems in this section are proven as stated in [1].

2.1 Head and Tail Operations

Definition 1 (Definition of lhead(r)). The relation lhead(r) = r' is defined in terms of the structure of r:

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

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)). The relation Itail(r) = r' is defined in terms of Ihead(r). Note that Ihead(r) = $a_1 + a_2 + ... + a_i$. We define Itail(r) = $\delta_{a_1}(r) + \delta_{a_2}(r) + ... + \delta_{a_i}(r) + \epsilon$.

Using these definitions of head and tail, we establish a correctness result upon which type soundness for the concatenation operator in λ_{RS} depends.

TR Lemma A (Leading characters are in the head). If $c_1 \cdot c_2 \cdot ... \cdot c_m \in \mathcal{L}\{r\}$, then $c_1 \in \mathsf{lhead}(r)$.

Proof. By structural induction on r...

TR Theorem B (Correctness of Itail). *If* $s \in \mathcal{L}\{r\}$ *then* $s \in \mathcal{L}\{\text{Ihead}(r)\} \cdot \mathcal{L}\{\text{Itail}(r)\}$.

Proof. The proof proceeds by structural induction on r. In each case, we demonstrate three facts:

- First, that the string $s = a \cdot b$ is a concatenation of two substrings (each of which may be ϵ).
- Second, $a \in \mathcal{L}\{\mathsf{lhead}(r)\}\ (\mathsf{or}, \mathsf{equivalently}, s \in \mathcal{L}\{\mathsf{lhead}(r)\}\ \mathsf{letting}\ b = \epsilon).$
- Third, there exists a $c \in \mathsf{Ihead}(r)$ such that $a = c \cdot a'$ and $a' \cdot b \in \mathcal{L}\{\delta_c(r)\}$, or else s = a.

Throughout, we consider ϵ as identity under all operators for expressions and \cdot for strings.

Case r = c for $c \in Sigma$. Note that if $s \in \mathcal{L}\{c\}$ then s = c. Note that $lhead(c) = lhead(c, \epsilon) = c$, and $c \in \mathcal{L}\{c\}$. So $c \in \mathcal{L}\{c\}$, which is sufficient since $\epsilon \in ltail(r)$.

Case $r = r_1 \cdot r_2$. Suppose $s \in \mathcal{L}\{r\}$. We may decompose s as $s_1 \cdot ... \cdot s_n$ such that $s_1 = c_1 \cdot ... \cdot c_m \in r_1$. Note that $c_1 \in \mathsf{lhead}(r_1)$ by A. Let $a' = c_2 \cdot ... \cdot c_m$. It suffices to show that $a' \cdot b \in \mathcal{L}\{\delta_c(r)\}$. Since $s = c_1 \cdot a' \cdot b$ and $s \in \mathcal{L}\{r\}$, it follows by the definition of derivative that $a' \cdot b \in \delta_{c_1}(r)$.

Case $r = r_1 + r_2$. Suppose $s \in \mathcal{L}\{r_1 + r_2\}$ so that $s \in \mathcal{L}\{r_1\}$ or $s \in \mathcal{L}\{r_2\}$. In either case, the result follows by induction.

Case $r = q^*$. Note that $s = \epsilon$ or else $s = s_1 \cdot s_2 \cdot ... s_n$ where each $s_i \in \mathcal{L}\{q\}^1$.

Suppose $s=\epsilon$. Note that $\operatorname{Ihead}(r)=\operatorname{Ihead}(q^*)=\operatorname{Ihead}(q^*,\epsilon)=\operatorname{Ihead}(\epsilon,\epsilon)+h'=\epsilon+h'$. Therefore, $s\in\operatorname{Ihead}(r)$.

In the other case, suppose $s = s_1 \cdot ... \cdot s_n$ where each $s_i \in \mathcal{L}\{q\}$. Let $a = s_1$ and $b = s_2 \cdot ... s_n$, so that $s = a \cdot b$. Note that $a = c_1 \cdot ... \cdot c_n$. We will show that $c_1 \in \mathcal{L}\{\mathsf{lhead}(q^*)\}$ and that $(c_2 \cdot ... c_n) \cdot b \in \delta_{c_1}(q^*)$.

The latter property is trivial. Note that $s = c_1 \cdot c_2 \cdot ... \cdot c_n \cdot b \in \mathcal{L}\{q^*\}$. The result follows immediately from the definition of derivative.

What remains to be shown is that $c_1 \in \mathcal{L}\{q^*\}$. Note that $\mathsf{lhead}(q^*) = \epsilon + \mathsf{lhead}(q, \epsilon)$. Recall that $a = s_1$ where $s_1 \in \mathcal{L}\{q\}$. Also recall that $a = c_1 \cdot \ldots \cdot c_n$. Therefore, $c_1 \in \mathsf{lhead}(q)$ by A.

TR Example C (All the heads of all the tails can be more than one head and tail). $r \neq \text{lhead}(r) \cdot \text{ltail}(r)$.

Proof. A simple counter-example is ab + cd. Note that lhead(ab + cd) = a + c and ltail(ab + cd) = b + d. Therefore, $\{ad, bc\} \subset \mathcal{L}\{lhead(ab + cd) \cdot ltail(ab + cd)\}$ even though neither of these is in $\mathcal{L}\{r\}$.

Example C does not imply a counter-example to type soundness because $s \in \mathcal{L}\{r\} \implies s \in \mathsf{lhead}(r) \cdot \mathsf{ltail}(r)$ is the property required for soundness. Still, in a production implementation, it will make sense to massage the definitions of $\mathsf{lhead}(r)$ and $\mathsf{ltail}(r)$ so that type information is not unnecessarily lost during substring operations.

This is a general pattern in string operations: λ_{RS} simulates – within the type system – common operations on strings. If there is an operation for concatenating to strings, we define an operation for concatenating two regular expressions. If there is an operation for peeling off the first (n) characters of a string, then we define an operation for converting a regular expression r into a regular expression r' which recognizes any n^{th} suffix of a string in r.

It is important to note, however, that the type system need not *exactly* simulate the action of string operations. In the case of concatenation, we lose some information because more string values are possible – according the types – than are actually possible in the dynamic semantics. Soundness is not lost because the types are conservative in their approximation.

In the case of string replacement, there are *trivial* definitions of substitution (on strings) and replacement (on languages) which over-approximate the effect of a substitution. Closing these gaps in approximation is important, and motivates the string operations portion of this technical report.

¹intuitively, s is either empty or a finitary concatenation of strings matching q.

2.2 Some Corollaries About Substitution and Language Replacement

Definition 4 (subst). We consider several choices in the string operations section.

Definition 5 (Ireplace). We consider several choices in the string operations section.

Proposition 6 (Closure). If $\mathcal{L}\{r\}$, $\mathcal{L}\{r_1\}$ and $\mathcal{L}\{r_2\}$ are regular languages, then $\mathcal{L}\{\text{lreplace}(r, r_1, r_2)\}$ is also a regular language.

Proof. This result is proven for various formulations in the next section.

Proposition 7 (Substitution Correspondence). If $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ then $\mathrm{subst}(r;s_1;s_2) \in \mathcal{L}\{\mathrm{lreplace}(r,r_1,r_2)\}$.

Proof. This is exactly the correctness result proven for some pairs of subst and replace in the previous section. \Box

Lemma 8 (Properties of Regular Languages.).

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

Lemma 8 states some well-known properties about regular expressions.

Lemma 9. If $\emptyset \vdash e$: stringin[r] then r is a well-formed regular expression.

Proof. The proof proceeds by induction on the derivation of the premise. The only non-trivial cases (those which require more than an appeal to inversion) are S-T-Case, S-T-Replace and S-T-Concat.

In the S-T-Case case, note that lhead and Itail are total functions for well-formed regular expressions to well-formed regular expressions.

In the S-T-Concat case, note that Lemma 6 implies that if r_1 and r_2 are regular expressions then so is $r_1 \cdot r_2$.

In the S-T-Replace case, it suffices to show that $lreplace(r, r_1, r_2)$ is a regular expression assuming (inductively) that r, r_1 and r_2 are all regular expressions. This follows from the Closure proposition.

2.3 The Small Step Semantics

To prove type safety and the security theorems for the big step semantics, we first prove type safety for a small step semantics in Figure 7 and then extend this to the big step semantics in Figure 5 by proving a correspondence between the semantics.

TR Conjecture D. *if* $\emptyset \vdash e : \sigma$ *then* $e \mapsto^* v$ *such that* v val.

We do not develop the full proof here, but note that the simply typed lambda calculus terminates. For the string fragment, observe that the S-T- rules do not add any non-trivial binding structure because substitutions [e/x]e' may only occur in the special case where $e = \mathsf{rstr}[s]$, so that the length of the term never increases and the number of free variables strictly decreases. Therefore, the standard normalization argument proceeds without complication after fixing an evaluation order for the compatibility rules (all our other proofs are agnostic to evaluation order).

TR Lemma E (Canonical Forms). *Suppose v* val.

If
$$\emptyset \vdash v$$
: stringin[r] then $v = rstr[s]$.
If $\emptyset \vdash v : \sigma \to \sigma'$ then $v = \lambda x.e'$ for some e' .

Proof. By inspection of valuation and typing rules.

For the sake of completeness, we include a statement of the weaker lemma stated in the paper:

Lemma 10 (Canonical Forms for the String Fragment of λ_{RS}). If $\emptyset \vdash e$: stringin[r] and $e \Downarrow v$ then $v = \mathsf{rstr}[s]$.

Proof. This fact follows directly from Lemma E.

TR Lemma F (Progress of small step semantics.). *If* $\emptyset \vdash e : \sigma$ *either* e val or $e \mapsto e'$ *for some* e'.

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

 λ **fragment**. Cases SS-T-Var, SS-T-Abs, and SS-T-App are exactly as in a proof of type safety for the simply typed lambda calculus.

S-T-Stringin-I. Suppose $\emptyset \vdash \mathsf{rstr}[s] : \mathsf{stringin}[s]$. The $\mathsf{rstr}[s]$ val by SS-E-RStr.

S-T-Concat. Suppose $\emptyset \vdash \mathsf{rconcat}(e_1; e_2) : \mathsf{stringin}[s]$. By inversion and induction, $e_1 \mapsto e'_1$ or e_1 val and similarly for e_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 val and e_2 val. But then it follows by Canonical Forms that $e_1 = \mathsf{rstr}[s_1]$ and $e_2 = \mathsf{rstr}[s_2]$. Finally, by SS-E-Concat, $\mathsf{rconcat}(\mathsf{rstr}[s_1]; \mathsf{rstr}[s_2]) \mapsto \mathsf{rstr}[s_1s_2]$.

S-T-Case. Suppose $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$. By inversion,

 $\emptyset \vdash e_1$: stringin[r]. From this fact, 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- ϵ -Val, and if s = at the e steps by S-E-Case-Concat.

S-T-Replace. Suppose $e = \text{rreplace}[r](e_1; e_2)$ and $\emptyset \vdash e : \text{stringin}[r']$ where, by inversion of S-T-Replace,

- $\emptyset \vdash e_1 : \mathsf{stringin}[r_1]$
- (2) $\emptyset \vdash e_2 : \mathsf{stringin}[r_2]$
- (3) $lreplace(r, r_1, r_2) = r'$

By (1), inversion and induction e_1 val or $e_1 \mapsto e_1'$ 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 val and also e_2 val. But then by Canonical Forms, $e_1 = \mathsf{rstr}[s_1]$ and $e_2 = \mathsf{rstr}[s_2]$. Therefore, $e \mapsto \mathsf{rstr}[\mathsf{subst}(r;s_1;s_2)]$ by SS-E-Replace.

S-T-SafeCoerce. Suppose that $\emptyset \vdash \mathsf{rcoerce}[r](e_1)$: stringin[r]. By inversion of S-T-SafeCoerce, $\emptyset \vdash e_1$: stringin[r'] for $\mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}$.By induction, e_1 val 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 val 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-SafeCheck Suppose that $\emptyset \vdash \mathsf{rcheck}[r](e_0; x.e_1; e_2)$: $\mathsf{stringin}[r]$. By inversion of S-T-Check:

$$(4) \qquad \qquad \vdash e_0 : \mathsf{stringin}[r_0]$$

(5)
$$x : \mathsf{stringin}[r] \vdash e_1 : \sigma$$

(6)
$$\vdash e_2 : \sigma$$

By (6) and induction, $e_0 \mapsto e'_0$ or e_0 val. In the former case e steps by SS-E-Check-StepRight. Otherwise, $e_0 = \mathsf{rstr}[s]$ by Canonical Forms. By Lemma 8, 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.

TR Lemma G (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 $e \mapsto e'$.

 λ **fragment**. Cases SS-T-Var, SS-T-Abs, and SS-T-App 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 $e_1 \mapsto e_1'$. By inversion of S-T-Concat, $\emptyset \vdash e_1 : \mathsf{stringin}[r_1]$ where

 $\emptyset \vdash e : \mathsf{stringin}[r_1r_2]$. By induction, if $e_1 \mapsto e_1'$ then $\emptyset \vdash e_1' : \mathsf{stringin}[r_1]$. Therefore, $\emptyset \vdash \mathsf{rconcat}(e_1';e_2) : \mathsf{stringin}[r_1r_2]$.

SS-E-Concat-Right.

SS-E-Concat. Suppose $\emptyset \vdash \mathsf{rconcat}(\mathsf{rstr}[s_1]; \mathsf{rstr}[s_2]) : \mathsf{stringin}[r_1r_2] \text{ and } \mathsf{rconcat}(\mathsf{rstr}[s_1]; \mathsf{rstr}[s_2]) \mapsto \mathsf{rstr}[s_1s_2].$ Then by inversion $\emptyset \vdash \mathsf{rstr}[s_1] : \mathsf{stringin}[r_1]$ and similarly for $\mathsf{rstr}[s_2]$. Therefore, $s_1 \in \mathcal{L}\{r_1\}$ and

 $s_2 \in \mathcal{L}\{r_2\}$ from which it follows by Lemma 8 that $s_1s_2 \in \mathcal{L}\{r_1r_2\}$. Therefore, $\emptyset \vdash \mathsf{rstr}[s_1s_2] : \mathsf{stringin}[r_1r_2]$ by S-T-Rstr.

S-E-Case-Left. Suppose that $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$ and also that $e \mapsto \mathsf{rstrcase}(e_1'; e_2; x, y.e_3)$ and $\emptyset \vdash e : \mathsf{stringin}[r]$. By inversion of S-T-Case:

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

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

(9)
$$x : \text{stringin}[\text{lhead}(r)], y : \text{stringin}[\text{ltail}(r)] \vdash e_3 : \sigma$$

By (7) and the assumption that $e_1 \mapsto e_1'$, it follows by induction that $\emptyset \vdash e_1'$: stringin[r]. This fact together with (8) and (9) implies by S-T-Case that $\emptyset \vdash \mathsf{rstrcase}(e_1'; e_2; x, y.e_3) : \sigma$.

SS-E-Case-Right. We have that $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$, Suppose $e \mapsto \mathsf{rstrcase}(e_1; e_2'; x, y.e_3)$ and $\emptyset \vdash e : \mathsf{stringin}[r]$. By inversion of S-T-Case:

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

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

(12)
$$x : \text{stringin}[\text{lhead}(r)], y : \text{stringin}[\text{ltail}(r)] \vdash e_3 : \sigma$$

By (11) and the assumption that $e_2 \mapsto e_2'$, it follows by induction that $\emptyset \vdash e_2'$: stringin[r]. This fact together with (10) and (12) implies by S-T-Case that $\emptyset \vdash \mathsf{rstrcase}(e_1; e_2'; x, y.e_3) : \sigma$.

Similar to SS-E-Concat-Left.

SS-E-Case- ϵ **-Val**. Suppose:

$$e = \mathsf{rstrcase}(-; e_2; -)$$

 $\emptyset \vdash e : \sigma$
 $e \mapsto e_2$

By inversion of S-T-Case, $e_2 : \sigma$.

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 $\emptyset \vdash e : \sigma$. By inversion of S-T-Case:

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

$$\emptyset \vdash \mathsf{rstr}[e_2] : \sigma$$

(15)
$$x : \text{stringin}[\text{lhead}(r)], y : \text{stringin}[\text{ltail}(r)] \vdash e_3 : \sigma$$

We know that $as \in \mathcal{L}\{r\}$ by (13) and inversion of S-T-Rstr. Therefore, $a \in \mathcal{L}\{\text{lhead}(r)\}$ by definition of lhead. Furthermore, $|\text{tail}(r)| = ... |\delta_a r| ...$ by definition of ltail. Note that $s \in \mathcal{L}\{\delta_a r\}$ by definition of the derivative, and so $s \in \mathcal{L}\{|\text{tail}(r)|\}$

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 $\emptyset \vdash \mathsf{rstr}[s] : \mathsf{stringin}[\mathsf{lhead}(r)]$. It follows by (15) 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 $e = \text{rreplace}[r](\text{rstr}[s_1]; \text{rstr}[s_2]) \mapsto \text{rstr}[\text{subst}(r; s_1; s_2)]$. Assume $\emptyset \vdash e : \text{stringin}[r']$ for $r' = \texttt{lreplace}(r, r_1, r_2)$. Then by inversion of S-T-Replace:

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

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

from which follows that $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$. Therefore, subst $(r; s_1; s_2) \in \mathcal{L}\{\text{lreplace}(r, r_1, r_2)\}$ by Theorem 7. It is finally derivable by S-T-Rstr that:

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

Case SS-E-SafeCoerce. Suppose that $\mathsf{rcoerce}[r](s_1) \mapsto \mathsf{rstr}[s_1]$ and that $\emptyset \vdash \mathsf{rcoerce}[r](s_1) : \mathsf{stringin}[r]$. By inversion of S-T-SafeCoerce we know that $s \in \mathcal{L}\{r\}$. Therefore, $\emptyset \vdash s : \mathsf{stringin}[r]$.

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 $\mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) \mapsto e_2, s \not\in \mathcal{L}\{r\}, \text{ and } \emptyset \vdash \mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) : \sigma$. By inversion of S-T-Check, $\emptyset \vdash e_2 : \sigma$.

TR Theorem H (Type Safety for small step semantics.). *If* $\emptyset \vdash e : \sigma$ *then either* e val or $e \mapsto^* e'$ and $\emptyset \vdash e' : \sigma$.

Proof. Follows directly from progress and preservation.

At this point the method for handling compatibility cases is clear; there-Tore, we elide these cases.

2.3.1 Semantic Correspondence between Big and Small Step Semantics for λ_{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.

TR Theorem I (Semantic Correspondence for λ_{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 $\lambda x.e_2 \mapsto^* \lambda x.e_2$ by RT-Refl.

Case $e = e_1(e_2)$. The only applicable rule is S-E-App. By inversion, we establish that the following:

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

$$e_2 \Downarrow v_2$$

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

From which it follows by induction that:

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

$$e_2 \mapsto^* v_2$$

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

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

$$\frac{e_1 \mapsto^* e'_1}{e_1(e_2) \mapsto^* e'_1(e'_2)}$$

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:

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

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- ϵ . By inversion:

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

 $e_2 \Downarrow v$

from which it follows by induction that:

$$e_1 \mapsto^* rstr[\epsilon]$$

 $e_2 \mapsto^* v$

Note that the following rule is derivable:

$$\frac{\text{SS-E-Case-LR*}}{e_1 \mapsto^* e_1' \qquad e_2 \mapsto^* e_2'} \\ \frac{e_1 \mapsto^* e_1' \qquad e_2 \mapsto^* e_2'}{\text{rstrcase}(e_1; e_2; x, y.e_3) \mapsto^* \text{rstrcase}(e_1'; e_2'; x, y.e_3)}$$

From these facts is follows that $e \mapsto^* \mathsf{rstrcase}(\mathsf{rstr}[\epsilon]; v; x, y.e_3)$. By S-E-Case- ϵ -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 $e \mapsto^* \operatorname{rstrcase}(e_1'; \operatorname{rstr}[as]; x, y.e_3)$. SE-E-Case-Concat applies to this form, so by RT-Step we know $e \mapsto^* [\operatorname{rstr}[a], \operatorname{rstr}[s]/x, y]e_3$. Recall that $[\operatorname{rstr}[a], \operatorname{rstr}[s]/x, y]e_3 \mapsto^* v$, 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, $subst(r; s_1; s_2) = s$ by induction. Note that the following rule is derivable:

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

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

Finally, $\operatorname{rreplace}[r](\operatorname{rstr}[s_1]; \operatorname{rstr}[s_2]) \mapsto \operatorname{subst}(r; s_1; s_2).$

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 $\operatorname{subst}(r; s_1; s_2) = s$, 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 $e_1 \mapsto^* \mathrm{rstr}[s]$ derives $\mathrm{rcoerce}[r](e_1) \mapsto^* \mathrm{rcoerce}[r](\mathrm{rstr}[s])$. In the final step, $\mathrm{rcoerce}[r](\mathrm{rstr}[s]) \mapsto \mathrm{rstr}[s]$ by SS-E-SafeCoerce. From this fact, we may derive via RT-Trans that $e \mapsto^* \mathrm{rstr}[s]$ as required.

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

Note that the rule following is derivable:

$$\frac{e_1 \mapsto^* e_1' \qquad e_3 \mapsto^* e_3'}{\mathsf{rcheck}[r](e_1; x.e_2; e_3) \mapsto^* \mathsf{rcheck}[r](e_1'; x.e_2; e_3')}$$

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, $e \mapsto^* \mathrm{rcheck}[r](\mathrm{rstr}[s]; x.e_2; e_2)$. Note that by inversion $s \in \mathcal{L}\{r\}$. Therefore, SS-E-Check-Ok applies and so by RT-Trans $e \mapsto^* [\mathrm{rstr}[s]/x]e_1$. By inversion, $[\mathrm{rstr}[s]/x]e_1 \Downarrow v$. 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.

TR Theorem J (Semantic Correspondence for λ_{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.

П

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

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

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

Finally, $e \Downarrow \mathsf{rstr}[s_1 s_2]$ follows via S-E-Concat from the facts that $e_1 \Downarrow \mathsf{rstr}[s_1]$ and $e_2 \Downarrow \mathsf{rstr}[s_2]$. 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, $e_2 \mapsto^* \operatorname{rstr}[s_2]$ and $e_2 \Downarrow \operatorname{rstr}[s_2]$.

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

It suffices to show $e \Downarrow \mathsf{rstr}[\mathsf{subst}(r; s_1; s_2)]$, 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- ϵ . 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- ϵ 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 $e \Downarrow [a, s/x, y]e_3 \Downarrow v$.

2.4 Extension of Safety for Small Step Semantics

Theorem 11 (Type Safety). *If* $\emptyset \vdash e : \sigma$ *then* $e \Downarrow v$ *and* $\emptyset \vdash v : \sigma$.

Proof. If $\emptyset \vdash e : \sigma$ then $e \mapsto^* v$ for v val by termination. Therefore, $e \Downarrow v$ by part 2 of the semantic correspondence theorem.

Since $\emptyset \vdash e : \sigma$ and $e \mapsto^* v$, it follows that $\emptyset \vdash v : \sigma$ by type safety for the small step semantics.

The cases for coercion and checking are straightforward.

2.4.1 The Security Theorem

Theorem 12 (Correctness of Input Sanitation for λ_{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 inversion of S-T-Rstr, $s \in \mathcal{L}\{r\}$.

3 Proofs of Lemmas and Theorems About λ_P

Theorem 13 (Safety for λ_P). *If* $\emptyset \vdash \iota : \tau$ *then* $\iota \Downarrow \dot{\upsilon}$ *and* $\emptyset \vdash \dot{\upsilon} : \tau$.

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

Lemma 14 (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 $\iota \Downarrow \mathsf{str}[s]$.

4 Proofs and Lemmas and Theorems About Translation

Theorem 15 (Translation Correctness). *If* $\Theta \vdash e : \sigma$ *then there exists an* ι *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 $\mathsf{rstr}[s]$, so we might as well proceed from the fact that $\emptyset \vdash \mathsf{rstr}[s]$: $\mathsf{stringin}[r]$.

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 λ_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 ι_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 ι_2 .

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

Choose $\iota = \operatorname{concat}(\iota_1; \iota_2)x, y.\iota_3$

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 16 (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 15 and the rules given, $\iota \Downarrow \mathsf{str}[s]$ implies that $e \Downarrow \mathsf{rstr}[s]$. Theorem 12 together with the assumption that e is well-typed implies that $s \in \mathcal{L}\{r\}$.

5 String Substitution and Language Replacement

5.1 The Trivial Definition

5.2 An Automaton Construction

Insert Automaton stuff...

5.3 Toward a Precise Definition

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 Σ .

and note that by the properties established via induction, $\llbracket e \rrbracket = \iota$ and $\llbracket \Psi \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$.

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

Figure 2: Syntax of λ_{RS} .

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

Figure 3: Syntax for the target language, λ_P , containing strings and statically constructed regular expressions.

$$\begin{array}{c|c} \Psi \vdash e : \sigma & \Psi ::= \emptyset \mid \Psi, x : \sigma \\ \hline S\text{-T-VAR} \\ x : \sigma \in \Psi \\ \overline{\Psi} \vdash x : \sigma & \overline{\Psi}, x : \sigma_1 \vdash e : \sigma_2 \\ \hline \Psi \vdash \lambda x.e : \sigma_1 \to \sigma_2 & \underline{\Psi} \vdash e_1 : \sigma_2 \to \sigma \quad \Psi \vdash e_2 : \sigma_2 \\ \hline \Psi \vdash e_1 : \text{stringin}[r_1] & \underline{\Psi} \vdash e_2 : \text{stringin}[r_2] \\ \hline \Psi \vdash r \text{concat}(e_1; e_2) : \text{stringin}[r_1 : r_2] \\ \hline \\ S\text{-T-CASE} \\ \underline{\Psi} \vdash e_1 : \text{stringin}[r] & \underline{\Psi} \vdash e_2 : \sigma \quad \Psi, x : \text{stringin}[\text{lhead}(r)], y : \text{stringin}[\text{ltail}(r)] \vdash e_3 : \sigma \\ \hline \underline{\Psi} \vdash r \text{stricase}(e_1; e_2; x, y.e_3) : \sigma \\ \hline \\ S\text{-T-Replace} \\ \underline{\Psi} \vdash e_1 : \text{stringin}[r_1] & \underline{\Psi} \vdash e_2 : \text{stringin}[r_2] \\ \underline{\Psi} \vdash r \text{replace}(r, r_1, r_2) = r' \\ \hline \underline{\Psi} \vdash r \text{replace}[r](e_1; e_2) : \text{stringin}[r'] & \underline{\Psi} \vdash e : \text{stringin}[r'] \\ \hline \underline{\Psi} \vdash r \text{coerce}[r](e) : \text{stringin}[r] \\ \underline{\Psi} \vdash r \text{coerce}[r](e) : \text{stringin}[r] \\ \underline{\Psi} \vdash r \text{check}[r](e_0; x.e_1; e_2) : \sigma \\ \hline \end{array}$$

Figure 4: Typing rules for λ_{RS} . The typing context Ψ is standard.

Figure 5: Big step semantics for λ_{RS} .

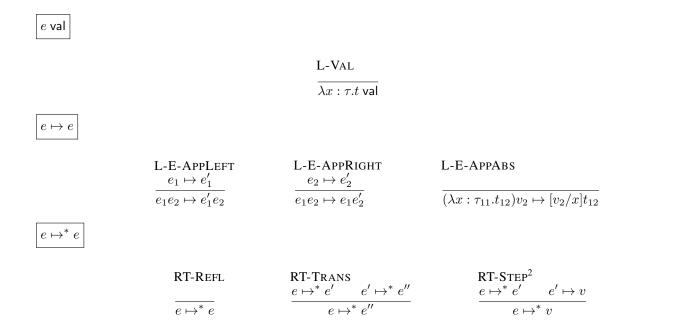


Figure 6: Call-by-name small step Semantics for λ and its reflexive, transitive closure.

Figure 7: Small step semantics for λ_{RS} . Extends 6.

Figure 8: Typing rules for λ_P . The typing context Θ is standard.

 $\iota \Downarrow \dot{v}$

P-E-ABS P-E-APP P-E-STR P-E-Rx $\frac{\iota_1 \Downarrow \lambda x. \iota_3}{\iota_2 \Downarrow \dot{v}_2} \frac{\iota_2 \Downarrow \dot{v}_2}{[\dot{v}_2/x]\iota_3 \Downarrow \dot{v}_3}$ $\lambda x.e \Downarrow \lambda x.e$ $\mathsf{str}[s] \Downarrow \mathsf{str}[s]$ $rx[r] \Downarrow rx[r]$ $\iota_2 \Downarrow \dot{v_2}$ P-E-Case- ϵ P-E-CONCAT P-E-CASE-CONCAT $\iota_1 \Downarrow \mathsf{str}[as] \qquad [\mathsf{str}[a], \mathsf{str}[s]/x, y] \iota_3 \Downarrow \dot{v}$ $\iota_1 \Downarrow \mathsf{str}[s_1]$ $\iota_2 \Downarrow \mathsf{str}[s_2]$ $\iota_1 \Downarrow \mathsf{str}[\epsilon]$ $concat(\iota_1; \iota_2) \Downarrow str[s_1s_2]$ $strcase(\iota_1; \iota_2; x, y.\iota_3) \Downarrow \dot{v_2}$ $strcase(\iota_1; \iota_2; x, y.\iota_3) \downarrow \dot{v}$ P-E-REPLACE $\iota_1 \Downarrow \mathsf{rx}[r] \qquad \iota_2 \Downarrow \mathsf{str}[s_2] \qquad \iota_3 \Downarrow \mathsf{str}[s_3] \qquad \mathsf{subst}(r; s_2; s_3) = s$ $replace(\iota_1; \iota_2; \iota_3) \Downarrow str[s]$ P-E-CHECK-OK P-E-CHECK-NOTOK $\iota_r \Downarrow \mathsf{rx}[r] \qquad \iota \Downarrow \mathsf{str}[s] \qquad s \not\in \mathcal{L}\{r\} \qquad \iota_2 \Downarrow \dot{v_2}$ $\iota \Downarrow \mathsf{str}[s] \qquad s \in \mathcal{L}\{r\} \qquad \iota_1 \Downarrow \dot{v_1}$ $\iota_r \Downarrow \mathsf{rx}[r]$ $\mathsf{check}(\iota_r; \iota; \iota_1; \iota_2) \Downarrow \dot{v_1}$ $\mathsf{check}(\iota_r; \iota; \iota_1; \iota_2) \Downarrow \dot{v_2}$

Figure 9: Big step semantics for λ_P

16

 $\iota \Downarrow \dot{v}$

$$\begin{array}{c} \text{SP-E-ABS} & \text{SP-E-APP} \\ \hline \lambda x.e \Downarrow \lambda x.e \end{array} & \begin{array}{c} \iota_1 \Downarrow \lambda x.\iota_3 & \iota_2 \Downarrow \dot{v}_2 & [\dot{v}_2/x]\iota_3 \Downarrow \dot{v}_3 \\ \hline \lambda x.e \Downarrow \lambda x.e \end{array} & \begin{array}{c} \text{SP-E-STR} & \text{SP-E-RX} \\ \hline \lambda x.e \Downarrow \lambda x.e \end{array} & \begin{array}{c} \iota_1 \Downarrow \lambda x.\iota_3 & \iota_2 \Downarrow \dot{v}_2 & [\dot{v}_2/x]\iota_3 \Downarrow \dot{v}_3 \\ \hline \lambda x.e \Downarrow \lambda x.e \end{array} & \begin{array}{c} I_1 \Downarrow \lambda x.\iota_3 & \iota_2 \Downarrow \dot{v}_2 & \text{str}[s] \Downarrow \text{str}[s] \Downarrow \text{str}[s] \\ \hline \lambda_1 \Downarrow \text{str}[s_1] & \iota_2 \Downarrow \text{str}[s_2] & \begin{array}{c} \text{SP-E-CASE-CONCAT} \\ \iota_1 \Downarrow \text{str}[s_1] & \iota_2 \Downarrow \text{str}[s_2] & \\ \hline \text{concat}(\iota_1; \iota_2) \Downarrow \text{str}[s_1 s_2] \end{array} & \begin{array}{c} \text{SP-E-CASE-CONCAT} \\ I_1 \Downarrow \text{str}[s] & \text{str}[s] \text{str}[s], \text{str}[s]/x, y]\iota_3 \Downarrow \dot{v} \\ \hline \text{str}[s] & \text{str}[s] & \text{str}[s] & \text{str}[s] \text{str}[s]/x, y]\iota_3 \Downarrow \dot{v} \\ \hline \\ \text{SP-E-REPLACE} \\ I_1 \Downarrow \text{rx}[r] & I_2 \Downarrow \text{str}[s_2] & I_3 \Downarrow \text{str}[s_3] & \text{subst}(r; s_2; s_3) = s \\ \hline \text{replace}(\iota_1; \iota_2; \iota_3) \Downarrow \text{str}[s] \\ \hline \\ \text{SP-E-CHECK-NOTOK} \\ I_2 \Downarrow \text{rx}[r] & I \Downarrow \text{str}[s] & s \not\in \mathcal{L}\{r\} & I_2 \Downarrow \dot{v}_2 \\ \hline \\ \text{check}(\iota_r; \iota; \iota_1; \iota_2) \Downarrow \dot{v}_1 & \text{check}(\iota_r; \iota; \iota_1; \iota_2) \Downarrow \dot{v}_2 \\ \hline \end{array}$$

Figure 10: Small step semantics for λ_P

 $\llbracket \sigma \rrbracket = \tau$ $\begin{aligned} & \text{TR-T-ARROW} \\ & \underline{\llbracket \sigma_1 \rrbracket} = \tau_1 & \underline{\llbracket \sigma_2 \rrbracket} = \tau_2 \\ & \underline{\llbracket \sigma_1 \to \sigma_2 \rrbracket} = \tau_1 \to \tau_2 \end{aligned}$ TR-T-STRING $\overline{\|\text{stringin}[r]\|} = \text{string}$ $\llbracket \Psi \rrbracket = \Theta$ TR-T-CONTEXT-EXT TR-T-CONTEXT-EMP $\frac{\llbracket \Psi \rrbracket = \Theta \qquad \llbracket \sigma \rrbracket = \tau}{\llbracket \Psi, x : \sigma \rrbracket = \Theta, x : \tau}$ $\boxed{\llbracket\emptyset\rrbracket=\emptyset}$ $\llbracket e \rrbracket = \iota$ $\begin{array}{ll} \text{TR-ABS} & \text{TR-APP} \\ \underline{\llbracket e \rrbracket = \iota} \\ \overline{\llbracket \lambda x. e \rrbracket = \lambda x. \iota} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_2 \rrbracket = \iota_2}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_2}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_2}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_2 \rrbracket = \iota_2}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1 \rrbracket = \iota_1}{\llbracket e_1 \rrbracket = \iota_1}} & \underline{\frac{\mathbb{I} e_1$ TR-VAR $[\![x]\!]=x$ $\begin{array}{c|c} \text{TR-STRING} & \text{TR-CONCAT} & \text{TR-SUBST} \\ \hline \llbracket e_1 \rrbracket = \iota_1 & \llbracket e_2 \rrbracket = \iota_2 & \llbracket e_1 \rrbracket = \iota_1 & \llbracket e_2 \rrbracket = \iota_2 \\ \hline \llbracket \mathsf{rstr}[s] \rrbracket = \mathsf{str}[s] & \hline \llbracket \mathsf{rconcat}(e_1; e_2) \rrbracket = \mathsf{concat}(\iota_1; \iota_2) & \hline \llbracket \mathsf{rreplace}[r](e_1; e_2) \rrbracket = \mathsf{replace}(\mathsf{rx}[r]; \iota_1; \iota_2) \\ \hline \end{array}$ TR-STRING $\begin{array}{lll} \text{TR-SAFECOERCE} & & \text{TR-CHECK} \\ & \llbracket e \rrbracket = \iota & & \llbracket e \rrbracket = \iota & \llbracket e_1 \rrbracket = \iota_1 & \llbracket e_2 \rrbracket = \iota_2 \\ & \llbracket \text{rcoerce}[r'](e) \rrbracket = \iota & & \hline{\llbracket \text{rcheck}[r](e; x.e_1; e_2) \rrbracket = \text{check}(\text{rx}[r]; \iota; (\lambda x.\iota_1)(\iota); \iota_2)} \end{array}$

Figure 11: Translation from source terms (e) to target terms (ι) .