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.

Throughout the following proofs, we abbreviate $s \in \mathcal{L}\{r\}$ by $s \in r$.

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

Proof. By structural induction on r.

If r=x for some character x then $s\in\mathcal{L}\{r\}$ implies s=x. By definition of lhead, lhead(s)= lhead $(s,\epsilon)=$ lhead $(x,\epsilon)=x$. Finally, $x\in\mathcal{L}\{x\}$.

If $r=r_1\cdot r_2$ then $s\in\mathcal{L}\{r\}$ implies $s=c_1s_1s_2$ such that $c_1s_1\in\mathcal{L}\{r_1\}$. By definition of Ihead, $|\operatorname{Ihead}(r_1\cdot r_2)|=|\operatorname{Ihead}(r_1,r_2)|=|\cdots|+|\operatorname{Ihead}(r_1,\epsilon)|=|\cdots|+|\operatorname{Ihead}(r_1)|$. Therefore, $r_1\subseteq \operatorname{Ihead}(r_1,r_2)$ and so it suffices to show that $s=c_1s_1s_2$ such that $c_1\in\operatorname{Ihead}(r_1)$, which holds by induction because $c_1s_1\in\mathcal{L}\{r_1\}$.

If $r=r_1+r_2$ then $s\in r$ implies either $s\in r_1$ or $s\in r_2$. Without loss of generality suppose $s\in r_1$. Then by induction $s=c_1\cdots c_m$ such that $c_1\in \operatorname{lhead}(r_1)$. By definition, $\operatorname{lhead}(r)=\operatorname{lhead}(r_1,\epsilon)+\operatorname{lhead}(r_2,\epsilon)=\operatorname{lhead}(r_1)+\cdots$. Therefore, $c_1\in \operatorname{lhead}(r)$.

If $r=r_1^*$ then $s\in r$ implies that $s=s_1\cdots s_k$ where $k\geq 0$ and $s_i\in \mathcal{L}\{r\}$. By definition of Ihead, $\mathsf{Ihead}(r_1^*)=\epsilon+\mathsf{Ihead}(r_1)$. If $s=\epsilon$, then $s\in \mathcal{L}\{\epsilon\}\subseteq \mathcal{L}\{\mathsf{Ihead}(r_1)\}$. If $s=s_1\cdots s_k$, then $s_1\in \mathcal{L}\{r\}$ and so by induction $s_1=c_1\cdots c_k$ such that $c_1\in \mathcal{L}\{\mathsf{Ihead}(r)\}$.

this is why we break out the c_1 from the s_1 , both here and going forward **Theorem B** (Correctness of Itail). *If* $s \in \mathcal{L}\{r\}$ *then* $s \in \mathcal{L}\{\text{Ihead}(r) \cdot \text{Itail}(r)\}$.

Proof. The proof proceeds by structural induction on r. In each case, we identify sublanguages of the head and tail of r such that s is in the concatenation of the the sublanguages. Throughout the proof, we use the fact that $\mathcal{L}\{r'\}\subseteq\mathcal{L}\{r\}$ and $\mathcal{L}\{q'\}\subseteq\mathcal{L}\{q\}$ implies $\mathcal{L}\{r'\cdot q'\}\subseteq\mathcal{L}\{r\cdot q\}$.

Case r = c for $c \in Sigma$. By definition, if $s \in \mathcal{L}\{c\}$ then s = c. By definition of lhead and ltail, $|\text{lhead}(c)| = |\text{lhead}(c, \epsilon)| = c$ and $\mathcal{L}\{\epsilon\} \subseteq \mathcal{L}\{|\text{ltail}(r)\}$. Therefore, $c \in \mathcal{L}\{c \cdot \epsilon\} \subseteq \mathcal{L}\{|\text{lhead}(r)| + |\text{ltail}(r)\}$.

Case $r = r_1 \cdot r_2$. Suppose $s \in \mathcal{L}\{r\}$. We may decompose s as $s_1 \cdots s_n$ such that $s_1 = c_1 s_1' \in \mathcal{L}\{r_1\}$. Note that $c_1 \in \mathsf{lhead}(r_1)$ by A. Therefore, $c_1 \in \mathsf{lhead}(r_1 \cdot r_2)$ by definition of lhead. It suffices to show that $s_1' s_2 \in \mathcal{L}\{\delta_{c_1}(r)\}$. Since $s = c_1 s_1' s_2 \in \mathcal{L}\{r\}$, it follows by the definition of derivative that $s_1' s_2 \in \delta_{c_1}(r)$. Therefore, $s = c_1 s_1' s_2 \in \mathcal{L}\{c_1 \cdot \delta_{c_1}\} \subseteq \mathcal{L}\{\mathsf{lhead}(r) \cdot \mathsf{ltail}(r)\}$.

Case $r=r_1+r_2$. Suppose $s\in \mathcal{L}\{r_1+r_2\}$ so that $s\in \mathcal{L}\{r_1\}$ or $s\in \mathcal{L}\{r_2\}$. Consider without loss of generality the case where $s\in \mathcal{L}\{r_1\}$. By induction, $s\in \mathcal{L}\{\operatorname{lhead}(r_1)\cdot\operatorname{ltail}(r_1)\}$ Note that $\operatorname{lhead}(r_1)\subseteq\operatorname{lhead}(r)$ and $\operatorname{ltail}(r_1)\subseteq\operatorname{ltail}(r)$. Therefore, $\operatorname{lhead}(r_1)\cdot\operatorname{ltail}(r_2)\subseteq\operatorname{lhead}(r_1)\cdot\operatorname{ltail}(r_2)$.

Case $r = q^*$. Either s is the empty string or else the k^{th} unwinding (i.e., $s = s_1 \cdots s_k$ where k > 0 and $s_i \in \mathcal{L}\{q\}$).

Suppose $s=\epsilon$. By definition of Ihead, Ihead(r)= Ihead $(q^*)=$ Ihead $(q^*,\epsilon)=$ Ihead $(\epsilon,\epsilon)+\cdots=\epsilon+\cdots$. By definition of Itail, Itail $(r)=\epsilon\cdots$. From these facts we conclude that $s=\epsilon\in\mathcal{L}\{\epsilon\cdot\epsilon\}\subseteq\mathcal{L}\{\mathsf{Ihead}(r)\cdot\mathsf{Itail}(r)\}$.

Suppose instead that $s=s_1\cdots s_k$ for k<0 where $s_i\in\mathcal{L}\{q\}$. If $q=\epsilon$ then $s=\epsilon$ and the conclusion follows for the same reason as the epsilon subcase of the character case. Otherwise, s is not the empty string and so $s=c_1s_1s_2$ where s_1 and s_2 might be empty. By the lemma above, $c_1\in \mathsf{lhead}(s)$. By the definition of derivative, $s_1s_2\in\delta_{c_1}(s)$. Therefore, $s\in\mathcal{L}\{\mathsf{lhead}(s)\cdot\delta_{c_1}(s)\}\subseteq\mathcal{L}\{\mathsf{lhead}(r)\cdot\mathsf{ltail}(r)\}$, where the set inclusion follows by the definition of lhead and ltail.

Example C (Decomposition is not Reversible). $r \neq \mathsf{lhead}(r) \cdot \mathsf{ltail}(r)$.

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

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.

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.

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

Assumption D (Properties of Regular Languages.).

- 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.

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

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.

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

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

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 8 (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 and the preservation lemma.

Lemma F (Progress of small step semantics.). If $\emptyset \vdash e : \sigma$ either e = v 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 progress for the simply typed lambda calculus.

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

S-T-Concat. Suppose $\emptyset \vdash \mathsf{rconcat}(e_1; e_2) : \mathsf{stringin}[r] \text{ and } \emptyset e_1 : \mathsf{stringin}[r_1] \text{ and } \emptyset \vdash e_2 : \mathsf{stringin}[r_2].$ By induction, $e_1 \mapsto e_1'$ or $e_1 = v_1$ 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, $e_1 = v_1$ and $e_2 = v_2$. 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 $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$ and $\emptyset \vdash e_1 : \mathsf{stringin}[r]$. By induction and Lemma 8 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 $e = \text{rreplace}[r](e_1; e_2), \emptyset \vdash e : \text{stringin}[r']$ and:

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

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

(3)
$$lreplace(r, r_1, r_2) = r'$$

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, $e_1 = \operatorname{rstr}[s_1]$ and $e_2 = \operatorname{rstr}[s_2]$. Therefore, $e \mapsto \operatorname{rstr}[\operatorname{subst}(r;s_1;s_2)]$ 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-SafeCheck Suppose that $\emptyset \vdash \mathsf{rcheck}[r](e_0; x.e_1; e_2)$: $\mathsf{stringin}[r]$ and:

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

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

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

By induction, $e_0 \mapsto e_0'$ or $e_0 = v$. In the former case e steps by SS-E-Check-StepRight. Otherwise, $e_0 = \mathsf{rstr}[s]$ by Canonical Forms. By Lemma D 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.

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

Lemma H (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 $\emptyset \vdash e : \sigma$.

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

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

SS-E-Concat. Suppose rconcat(rstr[s_1]; rstr[s_2]) \mapsto rstr[s_1s_2]. The only applicable rule is S-T-Concat, so $\emptyset \vdash \text{rstr}[s_1]$: stringin[r_1] and $\emptyset \vdash \text{rstr}[s_2]$: stringin[r_2]. By Lemma $\ref{lem:space}$?, $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ from which it follows by Lemma D that $s_1s_2 \in \mathcal{L}\{r_1r_2\}$. Therefore, $\emptyset \vdash \text{rstr}[s_1s_2]$: stringin[r_1r_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 $e_1 \to e_1'$. The only rule that applies is S-T-Case, so:

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

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

(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- ϵ **-Val**. Suppose:

$$e = \mathsf{rstrcase}(e_0; e_2; x, y.e_3)$$

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

The only rule htat applies is S-T-Case, so $e_2 : \sigma$.

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:

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

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

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

We know that $as \in \mathcal{L}\{r\}$ by canonical forms on (10) Therefore, $a \in \mathcal{L}\{\mathsf{lhead}(r)\}$ by definition of lhead. Furthermore, $\mathsf{ltail}(r) = \ldots + \delta_a r + \ldots$ by definition of ltail. Note that $s \in \mathcal{L}\{\delta_a r\}$ by definition of the derivative, and so $s \in \mathcal{L}\{\mathsf{ltail}(r)\}$

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 (G) 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.

todo

Case SS-E-Replace.

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

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

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

By conanical forms, $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$. Therefore, $lreplace(r, s_1, s_2) \in \mathcal{L}\{lreplace(r, r_1, r_2)\}$ by Theorem 7. 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 I (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.

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.

Lemma J. 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 K (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 $\lambda x.e_1 \mapsto^* \lambda x.e_1$ 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 J $e'(e_2) \mapsto^* v$. 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...

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

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

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

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:

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

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

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, $\operatorname{rreplace}[r](e_1; e_2) \mapsto^* \operatorname{rreplace}[r](\operatorname{rstr}[s_1]; \operatorname{rstr}[s_2]).$

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.

Theorem L (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.

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

Note that $e\mapsto^* \operatorname{rreplace}[r](\operatorname{rstr}[s_1];\operatorname{rstr}[s_2])\mapsto \operatorname{rstr}[\operatorname{subst}(r;s_1;s_2)]$ 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 9 (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.

The cases for coercion and checking are straightforward.

2.4.1 The Security Theorem

Theorem 10 (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 Caonical Forms, $s \in \mathcal{L}\{r\}$.

3 Proofs of Lemmas and Theorems About λ_P

Theorem 11 (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 12 (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 13 (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 $e = \mathsf{rstrcase}(e_1; e_2; x, y.e_3)$. 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 $\iota = \operatorname{concat}(\iota_1; \iota_2)x, y.\iota_3$ and note that by the properties established via induction, $\llbracket e \rrbracket = \iota$ and $\llbracket \Psi \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$.

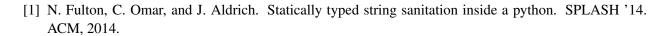
Suppose $s = \epsilon$. Then $e \Downarrow v$ where $e_2 \Downarrow v$ and $\iota \Downarrow \dot{v}$ where $\iota_2 \Downarrow \dot{v}$. But recall that $e_2 \leadsto v_2$ and so $\llbracket v \rrbracket = \dot{v}$.

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 14 (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 13 and the rules given, $\iota \Downarrow \mathsf{str}[s]$ implies that $e \Downarrow \mathsf{rstr}[s]$. Theorem 10 together with the assumption that e is well-typed implies that $s \in \mathcal{L}\{r\}$.

References



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

$$a \in \Sigma$$

Figure 1: Regular expressions over the alphabet Σ .

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

Figure 2: Syntax of λ_{RS} .

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

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

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

$$\begin{array}{c} \textbf{S-E-ABS} \\ \textbf{S-E-ABS} \\ \hline \lambda x.e \Downarrow \lambda x.e \end{array} \qquad \begin{array}{c} \textbf{S-E-APP} \\ e_1 \Downarrow \lambda x.e_3 & e_2 \Downarrow v_2 & [v_2/x]e_3 \Downarrow v \\ \hline & e_1(e_2) \Downarrow v \end{array} \qquad \begin{array}{c} \textbf{S-E-RSTR} \\ \textbf{rstr}[s] \Downarrow \textbf{rstr}[s] \end{array} \qquad \begin{array}{c} \textbf{S-E-CONCAT} \\ e_1 \Downarrow \textbf{rstr}[s_1] & e_2 \Downarrow \textbf{rstr}[s_2] \\ \hline & \textbf{rconcat}(e_1;e_2) \Downarrow \textbf{rstr}[s_1s_2] \end{array} \\ \\ \textbf{S-E-CASE-} \\ \hline e_1 \Downarrow \textbf{rstr}[e] & e_2 \Downarrow v_2 \\ \hline & \textbf{rstrcase}(e_1;e_2;x,y.e_3) \Downarrow v_2 \end{array} \qquad \begin{array}{c} \textbf{S-E-CASE-CONCAT} \\ e_1 \Downarrow \textbf{rstr}[s] & \textbf{rstr}[s]/x,y]e_3 \Downarrow v_3 \\ \hline & \textbf{rstrcase}(e_1;e_2;x,y.e_3) \Downarrow v_3 \end{array} \\ \\ \textbf{S-E-REPLACE} \\ \hline e_1 \Downarrow \textbf{rstr}[s_1] & e_2 \Downarrow \textbf{rstr}[s_2] & \textbf{subst}(r;s_1;s_2) = s \\ \hline & \textbf{rreplace}[r](e_1;e_2) \Downarrow \textbf{rstr}[s] \end{array} \qquad \begin{array}{c} \textbf{S-E-SAFECOERCE} \\ e \Downarrow \textbf{rstr}[s] \\ \hline & \textbf{rcoerce}[r](e) \Downarrow \textbf{rstr}[s] \end{array} \\ \\ \textbf{S-E-CHECK-OK} \\ \hline & \textbf{e} \Downarrow \textbf{rstr}[s] & s \not\in \mathcal{L}\{r\} & e_2 \Downarrow v \\ \hline & \textbf{rcheck}[r](e;x.e_1;e_2) \Downarrow v \end{array}$$

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

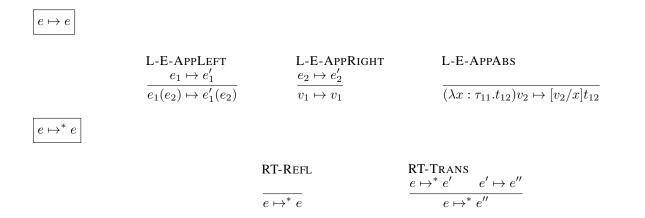


Figure 6: Call-by-name small step Semantics for λ and its reflexive, transitive closure.

 $e \mapsto e$ (Continues figure 6) $\frac{\textit{SS-E-CONCAT-LEFT}}{e_1 \mapsto e_1'} \\ \frac{e_1 \mapsto e_1'}{\mathsf{rconcat}(e_1; e_2) \mapsto \mathsf{rconcat}(e_1'; e_2)}$ SS-E-CONCAT-RIGHT $\frac{e_2 \mapsto e_2'}{\mathsf{rconcat}(v_1; e_2) \mapsto \mathsf{rconcat}(v_1; e_2')}$ SS-E-CASE-LEFT SS-E-CONCAT $\frac{e_1 \mapsto e_1'}{\mathsf{rstrcase}(e_1; e_2; x, y.e_3) \mapsto \mathsf{rstrcase}(e_1'; e_2; x, y.e_3)}$ $rconcat(rstr[s_1]; rstr[s_2]) \mapsto rstr[s_1s_2]$ SS-E-CASE- ϵ -VAL SS-E-CASE-CONCAT $\overline{\mathsf{rstrcase}(\mathsf{rstr}[as]; e_2; x, y.e_3) \mapsto [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3}$ $\mathsf{rstrcase}(\mathsf{rstr}[\epsilon]; e_2; x.y.e_3) \mapsto e_2$ SS-E-REPLACE-LEFT SS-E-REPLACE-RIGHT $\frac{e_1 \mapsto e_1'}{\mathsf{rreplace}[r](v_1; e_2) \mapsto \mathsf{rreplace}[r](v_1'; e_2)}$ $\frac{e_2 \mapsto e_2'}{\mathsf{rreplace}[r](e_1; e_2) \mapsto \mathsf{rreplace}[r](e_1; e_2')}$ $\begin{array}{c} {\rm SS\text{-}E\text{-}SafeCoerce\text{-}Step} \\ e \mapsto e' \end{array}$ SS-E-REPLACE $\frac{}{\mathsf{rcoerce}[r](e) \mapsto \mathsf{rcoerce}[r](e')}$ $\overline{\mathsf{rreplace}[r](\mathsf{rstr}[s_1];\mathsf{rstr}[s_2])} \mapsto \mathsf{rstr}[\mathsf{subst}(r;s_1;s_2)]$ SS-E-SAFECOERCE SS-E-CHECK-STEPLEFT $\overline{\mathsf{rcoerce}[r](\mathsf{rstr}[s]) \mapsto \mathsf{rstr}[s]} \qquad \overline{\mathsf{rcheck}[r](e; x.e_1; e_2) \mapsto \mathsf{rcheck}[r](e'; x.e_1; e_2)}$ SS-E-CHECK-OK SS-E-CHECK-NOTOK $\frac{s \in \mathcal{L}\{r\}}{\mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) \mapsto [\mathsf{rstr}[s]/x]e_1}$ $s \not\in \mathcal{L}\{r\}$ $\frac{c \not\vdash z_1, c_2}{\mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) \mapsto e_2}$

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

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

 $\iota \Downarrow \dot{v}$

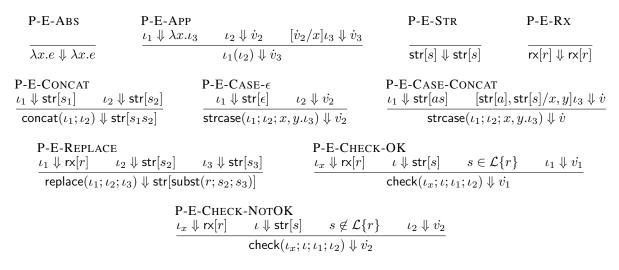


Figure 9: Big step semantics for λ_P

 $\iota \Downarrow \dot{v}$

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 $\frac{\llbracket \Psi \rrbracket = \Theta \qquad \llbracket \sigma \rrbracket = \tau}{\llbracket \Psi, x : \sigma \rrbracket = \Theta, x : \tau}$ $\boxed{\llbracket\emptyset\rrbracket=\emptyset}$ $\llbracket e \rrbracket = \iota$ $\begin{array}{ll} \text{TR-ABS} & \text{TR-APP} & \text{TR-CASE} \\ \underline{\llbracket e \rrbracket = \iota} & \underline{\llbracket e_1 \rrbracket = \iota_1} & \underline{\llbracket e_2 \rrbracket = \iota_2} & \underline{\llbracket e_1 \rrbracket = \iota_1} & \underline{\llbracket e_2 \rrbracket = \iota_2} \\ \underline{\llbracket \lambda x. e \rrbracket = \lambda x. \iota} & \underline{\llbracket e_1 (e_2) \rrbracket = \iota_1 (\iota_2)} & \underline{\llbracket rstrcase(e_1; e_2; x, y. e_3) \rrbracket = strcase(\iota_1; \iota_2; x, y. \iota_3)} \end{array}$ 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) .