Statically Typed String Sanitation Inside a Python (Technical Report)

Nathan Fulton Cyrus Omar Jonathan Aldrich

December 2014 CMU-ISR-14-112

School of Computer Science Carnegie Mellon University Pittsburgh, PA 15213

Abstract

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

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



Contents

1	Terminology and Notation	2
2	Regular Expressions	2
3	Λ_{RS}	2
	.1 Head and Tail Operations	2
	Replacement	3
	S.3 Small Step Semantics of λ_{RS}	3
	3.3.1 The Security Theorem	6
4	Proofs of Lemmas and Theorems About λ_P	6
5	Proofs and Lemmas and Theorems About Translation	9
L	t of Figures	
	Regular expressions over the alphabet Σ	13
	Syntax of λ_{RS}	
	Syntax for the target language, λ_P , containing strings and statically constructed regular	
	expressions	13
	Typing rules for λ_{RS} . The typing context Ψ is standard	13
	Call-by-name small step Semantics for λ and its reflexive, transitive closure	14
	Small step semantics for λ_{RS} . Extends 5	14
	Typing rules for λ_P . The typing context Θ is standard	15
	Small step semantics for λ_P (extends L-E rules)	16
	1 ' '	

1 Terminology and Notation

Theorems and lemmas appearing in [1] are numbered correspondingly, while supporting facts appearing only in the Technical Report are lettered. Throughout this technical report, we use a small step semantics corresponding to the big step semantics given in [1].

2 Regular Expressions

The syntax of regular expressions over some alphabet Σ is shown in Figure 1.

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

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

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

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

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

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

$$\epsilon^* \equiv \epsilon$$

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

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

3 λ_{RS}

The syntax of λ_{RS} is specified in Figure 2. The static semantics is specified in Figure 4.

3.1 Head and Tail Operations

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

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

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

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

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

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

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

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

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

3.2 Replacement

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

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

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

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

3.3 Small Step Semantics of λ_{RS}

Figure 6 specifies a small-step operational semantics for λ_{RS} .

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

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

Proof. By inspection of the static and dynamic semantics.

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

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

 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 $e = \mathsf{rstr}[s]$.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

- **S-E-Concat-Left.** Suppose $e = \mathsf{rconcat}(e_1; e_2) \mapsto \mathsf{rconcat}(e'_1; e_2)$ and $e_1 \mapsto e'_1$. The only rule that applies is S-T-Concat, so $\emptyset \vdash e_1$: stringin $[r_1]$ and $\emptyset \vdash e_2$: stringin $[r_2]$. By induction, $\emptyset \vdash e'_1$: stringin $[r_1]$. Therefore, by S-T-Concat, $\emptyset \vdash \mathsf{rconcat}(e'_1; e_2)$: stringin $[r_1r_2]$.
- **S-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, $\emptyset \vdash \mathsf{rconcat}(e_1; e_2') : \mathsf{stringin}[r_1r_2]$.
- **S-E-Concat**. Suppose $\operatorname{rconcat}(\operatorname{rstr}[s_1];\operatorname{rstr}[s_2])\mapsto \operatorname{rstr}[s_1s_2]$. The only applicable rule is S-T-Concat, so $\emptyset \vdash \operatorname{rstr}[s_1]:\operatorname{stringin}[r_1]$ and $\emptyset \vdash \operatorname{rstr}[s_2]:\operatorname{stringin}[r_2]$ and $\emptyset \vdash \operatorname{rconcat}(\operatorname{rstr}[s_1];\operatorname{rstr}[s_2]):\operatorname{stringin}[r_1 \cdot r_2]$. By Canonical Forms, $s_1 \in \mathcal{L}\{r_1\}$ and $s_2 \in \mathcal{L}\{r_2\}$ from which it follows by Lemma B that $s_1s_2 \in \mathcal{L}\{r_1 \cdot r_2\}$. Therefore, $\emptyset \vdash \operatorname{rstr}[s_1s_2]:\operatorname{stringin}[r_1 \cdot r_2]$ by S-T-Rstr.
- **S-E-Case-Left**. Suppose $e \mapsto \mathsf{rstrcase}(e_1'; e_2; x, y.e_3)$ and $\emptyset \vdash e : \sigma$ and $e_1 \mapsto e_1'$. The only rule that applies is S-T-Case, so:

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

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

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

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

- **S-E-Case-** ϵ **-Val**. Suppose $\operatorname{rstrcase}(e_0; e_2; x, y.e_3) \mapsto e_2$. The only rule that applies is S-T-Case, so $\emptyset \vdash e_2 : \sigma$.
- **S-E-Case-Concat**. Suppose that $e = \mathsf{rstrcase}(\mathsf{rstr}[as]; e_2; x, y.e_3) \mapsto [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]e_3$ and that $\emptyset \vdash e : \sigma$. The only rule that applies is S-T-Case so:

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

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

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

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

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

Case S-E-Replace-Left. Suppose that $e = \text{rreplace}[r](e_1; e_2) \mapsto \text{rreplace}[r](e_1'; e_2)$ when $e_1 \mapsto e_1'$. The only rule that applies is S-T-Replace, so $\emptyset \vdash e : \text{stringin}[\text{lreplace}(r, r_1, r_2)]$ where:

 $\emptyset \vdash e_1 : \mathsf{stringin}[r_1]$ $\emptyset \vdash e_2 : \mathsf{stringin}[r_2]$

By induction, $\emptyset \vdash e_1'$: stringin $[r_1]$. Therefore, $\emptyset \vdash \mathsf{rreplace}[r](e_1'; e_2)$: stringin $[\mathsf{lreplace}(r, r_1, r_2)]$ by S-T-Replace.

Cyrus stopped here **Case S-E-Replace-Right**. Suppose that $e = \text{rreplace}[r](e_1; e_2) \mapsto \text{rreplace}[r](e'_1; e_2)$ when $e_1 \mapsto e'_1$. The only rule that applies is S-T-Replace, so $\emptyset \vdash e : \text{stringin}[\text{lreplace}(r, r_1, r_2)]$ where:

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

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

By induction, $\emptyset \vdash e_1'$: stringin $[r_1]$. Therefore, $\emptyset \vdash \mathsf{rreplace}[r](r_1'; r_2)$: $\mathsf{stringin}[\mathsf{lreplace}(r, r_1, r_2)]$ by S-T-Replace.

Case S-E-Replace.

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

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

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 E. It is finally derivable by S-T-Rstr that:

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

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

Case S-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 S-E-Check-NotOk. Suppose $\mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) \mapsto e_2, s \not\in \mathcal{L}\{r\}, \text{ and } \emptyset \vdash \mathsf{rcheck}[r](\mathsf{rstr}[s]; x.e_1; e_2) : \sigma.$ The only applicable rule is S-T-Check, so $\emptyset \vdash e_2 : \sigma.$

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

Proof. Follows directly from progress and preservation.

3.3.1 The Security Theorem

Theorem 4 (Correctness of Input Sanitation for λ_{RS}). If $\emptyset \vdash e$: stringin[r] and $e \mapsto^* rstr[s]$ then $s \in \mathcal{L}\{r\}$.

4 Proofs of Lemmas and Theorems About λ_P

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

Lemma 5 (Canonical Forms for Target Language).

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

Theorem 6 (Progress). *If* $\emptyset \vdash \iota : \tau$ *either* $\iota = \dot{v}$ *or* $\iota \mapsto \iota'$ *for some* ι' .

Proof. The proof proceeds by induction on the typing assumption. Consider only the string and regex (non- λ) fragments of λ_P .

P-T-Case. Suppose $\emptyset \vdash \text{strcase}(\iota_1; \iota_2; x, y.\iota_3)$. By inversion, $\iota_1 : \text{string}$ and so either $\iota_1 \mapsto \iota'_1$ or by canonical forms, $\iota_1 = \text{str}[s_1]$. Similarly, $\iota_2 \mapsto \iota'_2$ or else $\iota_2 = \text{str}[s_2]$. In the former cases, progress occurs via the compatibility rules. in the case where both are string values, progress occurs via the case concatenation rule.

P-T-Replace. Suppose $\emptyset \vdash \text{replace}(\iota_1; \iota_2; \iota_3)$. By inversion, $\iota_1 : \text{regex}$ and so by canonical forms $\iota_1 = \text{rx}[r]$. By inversion, $\iota_2 : \text{string}$ and so by induction either $\iota_2 \mapsto \iota'_2$ or else $\iota_2 = \text{str}[s_2]$ for some string s_2 . Similarly, either ι_3 steps or else $\iota_3 = \text{str}[s_3]$. In case any steps occur, progress occurs. In the remaining case, PP-E-Replace applies and so progress occurs.

P-T-Check. Finally, suppose $\emptyset \vdash \iota_x \iota_1 \iota_2 \iota_3$. In case any of these step, then progress occurs. In the remaining cases, applications of inversion and canonical forms for each ι_x and ι_1 implies that the term at hand equals $rx[r]str[s]\iota_2\iota_3$, which evaluates to either ι_2 or ι_3 .

Lemma K (Substitution Lemma). *If* θ , $x : \tau \vdash \iota : \tau'$ and $\theta \vdash \iota' : \tau$ then $\theta \vdash [\iota'/x]\iota : \tau'$.

Theorem 7 (Preservation). *If* $\emptyset \vdash \iota : \tau$ *and* $\iota \mapsto \iota'$ *then* $\emptyset \vdash \iota' : \tau$.

Proof. The proof proceeds by induction of the derivations of $\emptyset \vdash \iota : \tau$ and $\iota \mapsto \iota'$. We treat only the non-lambda fragment.

Case PS-E-ConcatLeft. Suppose:

$$\iota = \mathsf{rconcat}(\iota_1; \iota_2) \mapsto \mathsf{rconcat}(\iota'_1; \iota_2)$$

 $\emptyset \vdash \iota : \mathsf{string}$
 $\iota \mapsto \iota'$

The only applicable typing rule is P-T-Concat, so $\emptyset \vdash \iota_1$: string and $\emptyset \vdash \iota_2$: string. By induction, $\emptyset \vdash \iota_1'$: string, so $\emptyset \vdash \mathsf{rconcat}(\iota_1'; \iota_2)$: string.

Case PS-E-ConcatRight

$$\begin{split} e &= \mathsf{rconcat}(e_1; e_2) \mapsto \mathsf{rconcat}(e_1; e_2') \\ \emptyset \vdash e : \mathsf{string} \\ \iota \mapsto \iota' \end{split}$$

The only applicable typing rule is P-T-Concat, so $\emptyset \vdash \iota_1$: string and $\emptyset \vdash \iota_2$: string. By induction, $\emptyset \vdash \iota_1'$: string, so $\emptyset \vdash \mathsf{rconcat}(\iota_1; \iota_2')$: string.

Case PS-E-Concat Let $e = \text{rconcat}(\text{rstr}[s_1]; \text{rstr}[s_2]) \mapsto \text{rstr}[s_1s_2]$. The only rule that applies is P-T-Concat, so $\emptyset \vdash e$: string. By canonical forms, $\emptyset \text{rstr}[s_1s_2]$: string.

Case PS-E-CaseLeft Let $\iota = \mathsf{rstrcase}(\iota_1; \iota_2; x, y.\iota_3) \mapsto \mathsf{rstrcase}(\iota'_1; \iota_2; x, y.\iota_3)$ when $\iota_1 \mapsto \iota'_1$. The only rule that applies is P-T-Case, so $\emptyset \vdash \iota : \tau$ where:

$$\emptyset \vdash \iota_1: \mathsf{string}$$
 $\emptyset \vdash \iota_2: au$ $\emptyset, x: \mathsf{string}, y: \mathsf{string} \vdash \iota_3: au$

By induction, $\emptyset \vdash \iota'_1$: string. By P-T-Case, $\emptyset \vdash \mathsf{rstrcase}(\iota_1; \iota_2; x, y.\iota_3) : \tau$.

Case PS-E-CaseEpsilon Let $\iota = \mathsf{rstrcase}(\mathsf{rstr}[\epsilon]; \iota_2; x, y. \iota_3) \mapsto \iota_2$. The only rule that applies is P-T-Case, so $\emptyset \vdash \iota : \tau$ where $\iota_2 : \tau$.

Case PS-E-Case Let $\iota = \mathsf{rstrcase}(\mathsf{rstr}[as]; \iota_2; x, y.\iota_3) \mapsto [\mathsf{rstr}[a], \mathsf{rstr}[s]/x, y]\iota_3$ The only rule that applies is P-T-Case, so $\emptyset \vdash \iota : \tau$ where:

$$\emptyset \vdash \iota_1: \mathsf{string}$$
 $\emptyset \vdash \iota_2: au$ $\emptyset, x: \mathsf{string}, y: \mathsf{string} \vdash \iota_3: au$

The result follows by the substitution lemma.

Case PS-E-ReplaceLeft Let $\iota = \mathsf{rreplace}[\iota_1](\iota_2; \iota_3) \mapsto \mathsf{rreplace}[\iota'_1](\iota_2; \iota_3)$ where $\iota_1 \mapsto \iota'_1$. The applicable typing rule is P-T-Replace, so $\emptyset \vdash \iota$: string where:

$$\emptyset \vdash \iota_1 : \mathsf{regex}$$

 $\emptyset \vdash \iota_2 : \mathsf{string}$
 $\emptyset \vdash \iota_3 : \mathsf{string}$

By induction, $\emptyset \vdash \iota_1'$: regex. Therefore, $\emptyset \vdash \mathsf{rreplace}[\iota_1'](\iota_2; \iota_3)$.

Case PS-E-ReplaceMid Let $\iota = \text{rreplace}[\iota_1](\iota_2; \iota_3) \mapsto \text{rreplace}[\iota_1](\iota_2'; \iota_3)$ where $\iota_2 \mapsto \iota_2'$. The applicable typing rule is P-T-Replace, so $\emptyset \vdash \iota$: string where:

$$\emptyset \vdash \iota_1 : \mathsf{regex}$$

 $\emptyset \vdash \iota_2 : \mathsf{string}$
 $\emptyset \vdash \iota_3 : \mathsf{string}$

By induction, $\emptyset \vdash \iota_2'$: string. Therefore, $\emptyset \vdash \mathsf{rreplace}[\iota_1](\iota_2'; \iota_3)$.

Case PS-E-ReplaceRight Let $\iota = \text{rreplace}[\iota_1](\iota_2; \iota_3) \mapsto \text{rreplace}[\iota_1](\iota_2; \iota_3')$ where $\iota_3 \mapsto \iota_3'$. The applicable typing rule is P-T-Replace, so $\emptyset \vdash \iota$: string where:

$$\emptyset \vdash \iota_1 : \mathsf{regex}$$

 $\emptyset \vdash \iota_2 : \mathsf{string}$
 $\emptyset \vdash \iota_3 : \mathsf{string}$

By induction, $\emptyset \vdash \iota_3'$: string. Therefore, $\emptyset \vdash \mathsf{rreplace}[\iota_1](\iota_2; \iota_3')$.

Case PS-E-Replace Let $\iota = \mathsf{rreplace}[\mathsf{rx}[r]](\mathsf{rstr}[s_2]; \mathsf{rstr}[s_3]) \mapsto \mathsf{rstr}[\mathsf{lreplace}(r, s_2, s_3)]$. The applicable typing rule is P-T-Replace, so $\emptyset \vdash \iota$: string. The result follows by canonical forms.

Case PS-E-CheckLeft Let $\iota = \mathsf{rcheck}[\iota_x](\iota_1; \iota_2; \iota_3) \mapsto \mathsf{rcheck}[\iota_x'](\iota_1; \iota_2; \iota_3)$ where $\iota_x \mapsto \iota_x'$. The applicable typing rule is P-T-Check, so $\emptyset \vdash \iota : \tau$ where:

$$\emptyset \vdash \iota_x : \mathsf{regex}$$

$$\emptyset \vdash \iota_1 : \mathsf{string}$$

$$\emptyset \vdash \iota_2 : \tau$$

$$\emptyset \vdash \iota_3 : \tau$$

By induction, ι_x : regex. Therefore, \emptyset rcheck $[\iota'_x](\iota_1; \iota_2; \iota_3) : \tau$.

Case PS-E-CheckRight Let $\iota = \mathsf{rcheck}[\iota_x](\iota_1; \iota_2; \iota_3) \mapsto \mathsf{rcheck}[\iota_x](\iota_1'; \iota_2; \iota_3)$ where $\iota_1 \mapsto \iota_1'$. The applicable typing rule is P-T-Check, so $\emptyset \vdash \iota : \tau$ where:

$$\emptyset \vdash \iota_x : \mathsf{regex}$$

$$\emptyset \vdash \iota_1 : \mathsf{string}$$

$$\emptyset \vdash \iota_2 : \tau$$

$$\emptyset \vdash \iota_3 : \tau$$

By induction, ι'_1 : string. Therefore, \emptyset rcheck $[\iota_x](\iota'_1; \iota_2; \iota_3) : \tau$.

Case PS-E-Check-Ok Let $\iota = \mathsf{rcheck}[\mathsf{rx}[r]](\mathsf{rstr}[s]; \iota_2; \iota_3) \mapsto \iota_2 \text{ and } s \in \mathcal{L}\{r\}$. The applicable typing rule is P-T-Check, so $\emptyset \vdash \iota : \tau$ where $\emptyset \vdash \iota_2 : \tau$.

Case PS-E-Check-NotOk Let $\iota = \mathsf{rcheck}[\mathsf{rx}[r]](\mathsf{rstr}[s]; \iota_2; \iota_3) \mapsto \iota_3$ where $s \notin \mathcal{L}\{r\}$. The applicable typing rule is P-T-Check, so $\emptyset \vdash \iota : \tau$ where $\emptyset \vdash \iota_3 : \tau$.

5 Proofs and Lemmas and Theorems About Translation

Theorem 8 (Type-Preserving Compilation). *If* $\Psi \vdash e : \sigma$ *then* $\llbracket \Psi \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$

Proof. By induction on the typing relation.

Case S-T-Abs Suppose $\Psi \vdash \lambda x.e' : \sigma_1 \to \sigma_2$. By induction, $\llbracket \Psi x : \sigma_1 \rrbracket \vdash \llbracket e' \rrbracket : \llbracket \sigma_1 \rrbracket$. By P-T-Abs, $\llbracket \Psi \rrbracket \vdash \llbracket \lambda x.e' \rrbracket : \llbracket \sigma_1 \to \sigma_2 \rrbracket$.

Case S-T-App Suppose $\Psi \vdash e_1(e_2) : \sigma$. By induction, $\llbracket \Psi \rrbracket \vdash \llbracket e_1 \rrbracket : \llbracket \sigma_2 \to \sigma \rrbracket$ and $\llbracket \Psi \rrbracket \vdash \llbracket e_2 \rrbracket : \llbracket \sigma_2 \rrbracket$. By the definition of type translation, $\llbracket \sigma_2 \to \sigma \rrbracket = \llbracket \sigma_2 \rrbracket \to \llbracket \sigma \rrbracket$. Therefore, $\llbracket \Psi \rrbracket \vdash \llbracket e_1(e_2) \rrbracket : \llbracket \sigma \rrbracket$.

Case S-T-StringIn-I. Suppose $\Psi \vdash \mathsf{rstr}[s] : \mathsf{stringin}[r]$ where $s \in \mathcal{L}\{r\}$. By the definition of term translation, $[\![\mathsf{rstr}[s]]\!] = \mathsf{str}[s]$ and $[\![\mathsf{stringin}[r]]\!] = \mathsf{string}$. By P-T-String, $\Theta \vdash \mathsf{str}[s] : \mathsf{string}$.

Case S-T-Concat. Suppose $\Psi \vdash \mathsf{rconcat}(e_1; e_2) : \mathsf{stringin}[r_1 \cdot r_2]$. By induction, $\llbracket \Psi \rrbracket \vdash \llbracket e_1 \rrbracket : \llbracket \mathsf{stringin}[r_1] \rrbracket$ and $\llbracket \Psi \rrbracket \vdash \llbracket e_2 \rrbracket : \llbracket \mathsf{stringin}[r_2] \rrbracket$. The result follows by P-T-Concat and the definition of term and type translation.

Case S-T-Case. Suppose $\Psi \vdash \mathsf{rstrcase}(e_1; e_2; x, y.e_3) : \sigma$. By induction, $\llbracket \Psi \rrbracket \vdash \llbracket e_1 \rrbracket : \llbracket \mathsf{stringin}[r] \rrbracket$, $\llbracket \Psi \rrbracket \vdash \llbracket e_2 \rrbracket : \llbracket \sigma \rrbracket$, and $\llbracket \Psi, x : \mathsf{stringin}[\mathsf{lhead}(r)], y : \mathsf{stringin}[\mathsf{ltail}(r)] \rrbracket \vdash \llbracket e_3 \rrbracket : \llbracket \sigma \rrbracket$. By the definition of context and type translation, $\llbracket \Psi, x : \mathsf{stringin}[\mathsf{lhead}(r)], y : \mathsf{stringin}[\mathsf{ltail}(r)] \rrbracket = \llbracket \Psi \rrbracket, x : \mathsf{string}, y : \mathsf{string}$. The result follows by P-T-Case.

Case S-T-Replace. Suppose $\Psi \vdash \mathsf{rreplace}[r](e_1; e_2) : \mathsf{stringin}[\mathsf{lreplace}(r, r_1, r_2)]$. Note that $[\mathsf{rreplace}[r](e_1; e_2)] = \mathsf{replace}(\mathsf{rx}[r]; [\![e_1]\!]; [\![e_2]\!])$. Let $\Theta = [\![\Psi]\!]$. By induction, $\Theta \vdash [\![e_1]\!]: [\![\mathsf{stringin}[r_1]\!]]$ and $\Theta \vdash [\![e_2]\!]: [\![\mathsf{stringin}[r_2]\!]]$ from which it follows by the definition of type translation that $\Theta \vdash [\![e_1]\!]: \mathsf{string}$ and $\Theta \vdash [\![e_2]\!]: \mathsf{string}$. So P-T-Replace, $\Theta \vdash \mathsf{replace}(\mathsf{rx}[r]; [\![e_1]\!]; [\![e_2]\!]): \mathsf{string}$

Case S-T-SafeCoerce. Suppose $\Psi \vdash \mathsf{rcoerce}[r](e) : \mathsf{stringin}[r] \text{ where } \Psi \vdash e : \mathsf{stringin}[r'] \text{ and } \mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}. \text{ Let } \llbracket\Psi\rrbracket = \Theta. \text{ Note that } \llbracket\mathsf{rcoerce}[r](e)\rrbracket = \llbracket e \rrbracket. \text{ By induction, } \Theta \vdash \llbracket e \rrbracket : \llbracket\mathsf{stringin}[r'] \rrbracket.$

Case S-T-Check. Suppose $\Psi \vdash \mathsf{rcheck}[r](e_0; x.e_1; e_2) : \sigma \text{ where } \Psi \vdash e_0 : \mathsf{stringin}[r], \Psi, x : \mathsf{stringin}[r] \vdash e_1 : \sigma, \text{ and } \Psi \vdash e_2 : \sigma.$

Note that $\llbracket \operatorname{rcheck}[r](e_0; x.e_1; e_2) \rrbracket = \operatorname{check}(\operatorname{rx}[r]\llbracket e_0 \rrbracket \lambda x.\llbracket e_1 \rrbracket (\llbracket e_0 \rrbracket) \llbracket e_2 \rrbracket ; ; ; L) \operatorname{et} \llbracket \Psi \rrbracket = \Theta.$ Note that $\llbracket r \rrbracket = \operatorname{rx}[r]$ and by canonical forms $\Theta \vdash \operatorname{rx}[r] : \operatorname{regex}.$ By induction, $\Theta \vdash \llbracket e_0 \rrbracket : \llbracket \operatorname{stringin}[r] \rrbracket$ and so $\Theta \vdash \llbracket e_0 \rrbracket : \operatorname{string}.$ By context translation, $\llbracket \Psi, x : \operatorname{stringin}[r] \rrbracket = \Theta, x : \operatorname{string}$ and so by induction $\Theta, x : \operatorname{string} \vdash \llbracket e_1 \rrbracket : \llbracket \sigma \rrbracket.$ Therefore, $\Theta \vdash \lambda x.\llbracket e_1 \rrbracket (\llbracket e_0 \rrbracket) : \llbracket \sigma \rrbracket.$ By induction, $\Theta \vdash \llbracket e_0 \rrbracket : \llbracket \sigma \rrbracket.$ From these three facts, it follows that $\Theta \vdash \operatorname{check}(\operatorname{rx}[r] \llbracket e_0 \rrbracket \lambda x.\llbracket e_1 \rrbracket (\llbracket e_0 \rrbracket) \llbracket e_2 \rrbracket ; ; \llbracket ; \rrbracket \sigma).$

Theorem 9 (Translation Correctness). If $\Psi \vdash e : \sigma$ and $e \mapsto e'$ then $[\![e]\!] \mapsto^* [\![e']\!]$.

Proof. By induction on the typing and evaluation derivations.

Theorem 10 (Translation Correctness). *If* $\Psi \vdash e : \sigma$ *then there exists an* ι *such that* $\llbracket e \rrbracket = \iota$ *and* $\llbracket \Psi \rrbracket \vdash \iota : \llbracket \sigma \rrbracket$. *Furthermore, if* $e \mapsto^* v$ *then* $\iota \mapsto^* \dot{v}$ *such that* $\llbracket v \rrbracket = \dot{v}$.

Proof. We present a proof by induction on the structure of e. We write $e \leadsto \iota$ as shorthand for the final property.

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

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

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

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

$$[rstr[s]] = str[s]$$

$$[stringin[r]] = string$$

$$[\![\emptyset]\!] = \emptyset$$

Note that $\emptyset \vdash \mathsf{str}[s]$: 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, $\llbracket \Theta \rrbracket \vdash \mathsf{str}[s]$: string by weakening.

Summarily, $\mathsf{str}[s]$ is a term of λ_P such that $[\![\Theta]\!] \vdash \mathsf{str}[s] : [\![\sigma]\!]$

This proof needs to be changed to use only the small-step semantics.

It remains to be shown that there exist v, \dot{v} such that $\mathsf{rstr}[s] \mapsto^* v$, $\mathsf{str}[s] \mapsto^* \dot{v}$, and $[v] = \dot{v}$. But this is immediate because each term is already a value and s = s.

Case $e = \text{rconcat}(e_1; e_2)$. The applicable typing rule is S-T-Concat, so $\Psi \vdash \text{rconcat}(e_1; e_2)$: stringin $[r_1 \cdot r_2]$ where $\Psi \vdash e_1$: stringin $[r_1]$ and $\Psi \vdash e_2$: stringin $[r_2]$.

By induction, $e_1 \leadsto \iota_1$ and $e_2 \leadsto \iota_2$. Therefore, $\llbracket \Psi \rrbracket \vdash \mathsf{concat}(\iota_1; \iota_2)$ by P-T-Concat.

By canonical forms, $e_1 \mapsto^* \mathsf{rstr}[s_1]$ where by induction $\iota_1 \mapsto^* \mathsf{str}[s_1]$. Similarly, $e_2 \mapsto^* \mathsf{rstr}[s_2]$ and $\iota_2 \mapsto^* \mathsf{str}[s_2]$. Therefore, $e \mapsto^* \mathsf{rstr}[s_1s_2]$ by S-E-Concat at last, and $\mathsf{concat}(\iota_1; \iota_2) \mapsto^* \mathsf{str}[s_1s_2]$ by P-E-Concat at last. Note that $\lceil \mathsf{rstr}[s_1s_2] \rceil = \mathsf{str}[s_1s_2]$.

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 $e_1 \mapsto \iota_1$.

By canonical forms, $e_1 \mapsto^* \mathsf{rstr}[s]$. Therefore, $\iota_1 \mapsto *\mathsf{str}[s]$ because $e_1 \rightsquigarrow \iota_1$.

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

To prove the evaluation correspondence, we consider two cases for the value of s.

Suppose $s = \epsilon$. Then $e \mapsto^* v$ where $e_2 \mapsto^* v$, from which it follows that $\iota \mapsto^* \dot{v}$ where $\iota_2 \mapsto^* \dot{v}$. But recall that $e_2 \rightsquigarrow v_2$ and so $\llbracket v \rrbracket = \dot{v}$.

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

Case $e = \operatorname{rreplace}[r](e_1; e_2)$. There is only one applicable typing rule, so suppose $\Psi \vdash \operatorname{rreplace}[r](e_1; e_2)$: stringin[lreplace (r, e_1, e_2)]. Let $\Theta = \llbracket \Psi \rrbracket$. Note that $\llbracket \operatorname{rreplace}[r](e_1; e_2) \rrbracket = \operatorname{replace}(\operatorname{rx}[r]; \iota_1; \iota_2)$ where by induction $\llbracket e_1 \rrbracket = \iota_1$ and $\llbracket e_2 \rrbracket = \iota_2$ such that $\Theta \vdash \iota_1$ and $\Theta \vdash \iota_2$. It follows by P-T-Replace that $\Theta \vdash \operatorname{replace}(\operatorname{rx}[r]; \iota_1; \iota_2)$: string. Finally, note that $\llbracket \operatorname{stringin}[\operatorname{lreplace}(r, e_1, e_2)] \rrbracket = \operatorname{string}$.

For evaluation correspondence, note that $[\![rstr[lreplace(r,s_1,s_2)]\!] = rstr[lreplace(r,s_1,s_2)]\!]$ and so it suffices to show that $[\![rstr[lreplace(r,s_1,s_2)]\!]] \mapsto rstr[r]s_1s_2$. Note that $[\![lreplace(r,e_1,e_2) \mapsto rstr[lreplace(r,s_1,s_2)]\!]$ where $[\![e_1 \mapsto rstr[s_1]\!]]$, $[\![e_2 \mapsto rstr[s_2]\!]]$, $[\![r \mapsto rstr[s_2]\!]]$, $[\![r \mapsto rstr[s_2]\!]]$, and $[\![rstr[s_2]\!]]$. So by S-E-Replace, the sufficient condition holds.

Case e = rcoerce[r](e'). The only applicable tpying rule is S-T-SafeCoerce, so suppose $\Psi \vdash \text{rcoerce}[r](e')$: stringin[r] where $\Psi \vdash e'$: stringin[r'] and $\mathcal{L}\{r'\} \subseteq \mathcal{L}\{r\}$. By induction, $e' \leadsto \iota$ for some ι . Therefore, $\llbracket \text{rcoerce}[r](e') \rrbracket = \iota$ by Tr-SafeCoerce.

For evaluation correspondence, note that $e \mapsto^* v$ where $e' \mapsto^* v$. The result follows by induction because $e' \rightsquigarrow \iota$.

Case $e = \mathsf{rcheck}[r](e_1; x.e_2; e_3)$. The applicable typing rule is S-T-Check, so $\Psi \vdash e : \sigma$ where $\Psi \vdash e_1 : \mathsf{stringin}[r], \Psi, x : \mathsf{stringin}[r] \vdash e_2 : \sigma$, and $\Psi \vdash e_3 : \sigma$. By induction and a corresponding substitution principle there exists $\iota_1, \iota_2, \iota_3$ such that $e_1 \leadsto \iota_1, e_2 \leadsto \iota_2$ in context $\Psi, s : \mathsf{stringin}[r]$, and $e_3 \leadsto \iota_3$. Choose $\iota = \mathsf{check}(\mathsf{rx}[r]; \iota_1; \lambda x.\iota_2; \iota_3)$. The result follows by induction.

Theorem 11 (Correctness of Input Sanitation for Translated Terms). *If* $\llbracket e \rrbracket = \iota$ *and* $\emptyset \vdash e : \mathsf{stringin}[r]$ *then* $\iota \mapsto^* \mathsf{str}[s]$ *for* $s \in \mathcal{L}\{r\}$.

Proof. By 4, $e \mapsto^* \mathrm{rstr}[s]$ where $\emptyset \vdash \mathrm{rstr}[s]$: stringin[r]. Therefore, $s \in \mathcal{L}\{r\}$. Note that $\llbracket \cdot \rrbracket$ is a function and $\llbracket \mathrm{rstr}[s] \rrbracket = \mathrm{str}[s]$; therefore, by theorem 10, $\iota \mapsto^* \mathrm{str}[s]$.

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 .

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

Figure 2: Syntax of λ_{RS} .

$$\tau \ \ \, ::= \tau \rightarrow \tau \ | \ \, \text{string} \ | \ \, \text{regex} \qquad \qquad \text{target types}$$

$$\iota \ \ \, ::= x \ | \ \, \dot{v} \qquad \qquad \qquad \text{target terms}$$

$$| \ \, \text{concat}(\iota;\iota) \ | \ \, \text{strcase}(\iota;\iota;x,y.\iota)$$

$$| \ \, \text{rx}[r] \ | \ \, \text{replace}(\iota;\iota;\iota) \ | \ \, \text{check}(\iota;\iota;\iota;\iota)$$

$$\dot{v} \ \ \, ::= \lambda x.\iota \ | \ \, \text{str}[s] \ | \ \, \text{rx}[r]$$

$$\text{target values}$$

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

$$\begin{array}{c|c} \Psi \vdash e : \sigma & \Psi ::= \emptyset \mid \Psi, x : \sigma \\ \hline S\text{-T-VAR} \\ x : \sigma \in \Psi \\ \overline{\Psi} \vdash x : \sigma & \underline{\Psi, x : \sigma_1 \vdash e : \sigma_2} \\ \overline{\Psi} \vdash \lambda x.e : \sigma_1 \rightarrow \sigma_2 & \underline{\Psi \vdash e_1 : \sigma_2 \rightarrow \sigma} \quad \underline{\Psi \vdash e_2 : \sigma_2} \\ \overline{\Psi \vdash e_1 : stringin[r_1]} \quad \underline{\Psi \vdash e_2 : stringin[r_2]} \\ \hline & \underline{\Psi \vdash rconcat(e_1; e_2) : stringin[r_1]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r_1]} \quad \underline{\Psi \vdash e_2 : stringin[r_1]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r_1]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r_1]} \\ \hline & \underline{\Psi \vdash e_2 : \sigma} \quad \underline{\Psi, x : stringin[lhead(r)], y : stringin[ltail(r)] \vdash e_3 : \sigma} \\ \hline & \underline{\Psi \vdash rstrcase(e_1; e_2; x, y.e_3) : \sigma} \\ \hline \\ \hline & \underline{\Psi \vdash e_1 : stringin[r_1]} \quad \underline{\Psi \vdash e_2 : stringin[r_2]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r_1]} \quad \underline{\Psi \vdash e_2 : stringin[r_2]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r]} \quad \underline{\Psi \vdash e_2 : stringin[r_2]} \\ \hline & \underline{\Psi \vdash e_1 : stringin[r]} \quad \underline{\Psi \vdash e_2 : stringin[r]} \\ \hline & \underline{\Psi \vdash rcoerce[r](e_1; e_2) : stringin[r]} \\ \hline & \underline{\Psi \vdash rcoerce[r](e_1; e_2) : stringin[r]} \\ \hline & \underline{\Psi \vdash rcheck[r](e_0; x.e_1; e_2) : \sigma} \\ \hline \end{array}$$

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

$$\begin{array}{c} \textbf{SS-E-APPLEFT} \\ & \underbrace{\frac{e_1 \mapsto e_1'}{e_1(e_2) \mapsto e_1'(e_2)}} \\ \hline \\ e \mapsto^* e \\ \hline \\ \hline \\ \textbf{RT-Refl} \\ \hline \\ & \underbrace{\frac{e_2 \mapsto e_2'}{v_1 \mapsto v_1}} \\ \hline \\ \hline \\ \textbf{RT-Trans} \\ \underbrace{\frac{e \mapsto^* e'}{e \mapsto^* e'}} \\ \hline \\ \hline \\ \hline \\ e \mapsto^* e'' \\ \hline \\ \hline \\ \hline \\ e \mapsto^* e'' \\ \hline \\ \hline \\ \hline \\ e \mapsto^* e'' \\ \hline \\ \hline \\ \hline \end{array} \right]$$

Figure 5: Call-by-name small step Semantics for λ and its reflexive, transitive closure.

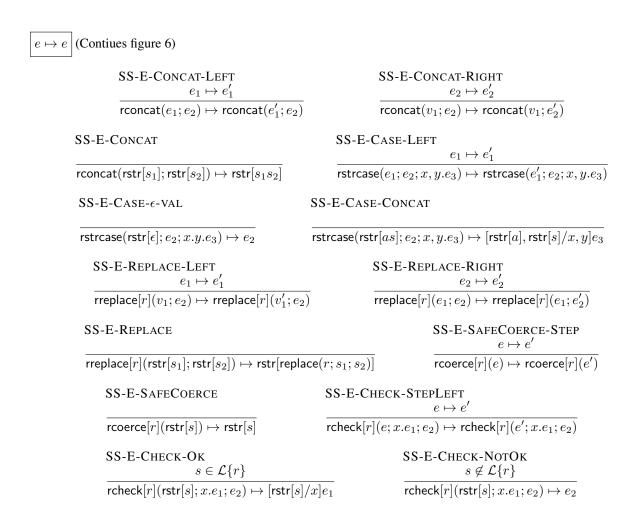


Figure 6: Small step semantics for λ_{RS} . Extends 5.

Figure 7: Typing rules for λ_P . The typing context Θ is standard.

$$\iota \mapsto \iota$$

$$\frac{\mathsf{PS-E-ConcatRight}}{\mathsf{concat}(\iota_1; \iota_2) \mapsto \mathsf{concat}(\iota_1'; \iota_2)} \qquad \frac{\mathsf{PS-E-ConcatRight}}{\mathsf{concat}(\iota_1; \iota_2) \mapsto \mathsf{concat}(\iota_1; \iota_2')} \qquad \frac{\mathsf{PS-E-Concat}}{\mathsf{concat}(\mathsf{str}[s_1]; \mathsf{str}[s_2]) \mapsto \mathsf{str}[s_1 s_2]} \\ \frac{\mathsf{PS-E-CaseLeft}}{\mathsf{strcase}(\iota_1; \iota_2; x, y. \iota_3) \mapsto \mathsf{strcase}(\iota_1'; \iota_2; x, y. \iota_3)} \qquad \frac{\mathsf{PS-E-Case-Epsilon}}{\mathsf{strcase}(\iota_1; \iota_2; x, y. \iota_3) \mapsto \mathsf{str}[s_1 s_2]} \\ \frac{\mathsf{PS-E-Case}}{\mathsf{strcase}(\mathsf{str}[as]; \iota_2; x, y. \iota_3) \mapsto \mathsf{str}[as]} \qquad \frac{\mathsf{PS-E-ReplaceLeft}}{\mathsf{replace}(\iota_1; \iota_2; \iota_3) \mapsto \mathsf{replace}(\iota_1'; \iota_2; \iota_3)} \\ \frac{\mathsf{PS-E-ReplaceMid}}{\mathsf{replace}(\mathsf{rx}[r]; \iota_2; \iota_3) \mapsto \mathsf{replace}(\mathsf{rx}[r]; \iota_2'; \iota_3)} \qquad \frac{\mathsf{PS-E-ReplaceRight}}{\mathsf{replace}(\mathsf{rx}[r]; \mathsf{str}[s_2]; \iota_3) \mapsto \mathsf{replace}(\mathsf{rx}[r]; \mathsf{str}[s_2]; \iota_3)} \\ \frac{\mathsf{PS-E-Replace}}{\mathsf{replace}(\mathsf{rx}[r]; \mathsf{str}[s_2]; \mathsf{str}[s_3]) \mapsto \mathsf{str}[\mathsf{replace}(r; s_2; s_3)]} \qquad \frac{\mathsf{PS-E-CheckLeft}}{\mathsf{replace}(\mathsf{rx}[\iota_1; \iota_2) \mapsto \mathsf{rcheck}[\iota_1](\iota_1; \iota_1; \iota_2)} \\ \frac{\mathsf{PS-E-Check-Notook}}{\mathsf{rcheck}[\mathsf{rx}[r]](\iota_1; \iota_1; \iota_2) \mapsto \mathsf{rcheck}[\mathsf{rx}[r]](\mathsf{str}[s]; \iota_1; \iota_2) \mapsto \iota_2} \\ \frac{\mathsf{PS-E-Check-Notook}}{\mathsf{rcheck}[\mathsf{rx}[r]](\mathsf{str}[s]; \iota_1; \iota_2) \mapsto \iota_2} \\ \\ \frac{\mathsf{PS-E-Check-Notook}}{\mathsf{rcheck}[\mathsf{rx}[r]](\mathsf{str}[s]; \iota_1; \iota_2) \mapsto \iota_2}$$

Figure 8: Small step semantics for λ_P (extends L-E rules)

$$\begin{array}{c|c} & T_{R}\text{-T-STRING} & T_{R}\text{-T-STRING} \\ \hline \llbracket \Psi \rrbracket = \Theta & & T_{R}\text{-T-CONTEXT-EMP} \\ \hline \llbracket \Psi \rrbracket = \Theta & & T_{R}\text{-T-CONTEXT-EMP} \\ \hline \llbracket \Psi \rrbracket = \theta & & T_{R}\text{-T-CONTEXT-EMP} \\ \hline \llbracket \Psi \rrbracket = \theta & & \llbracket \Psi \rrbracket = \Theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta & \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta \\ \hline \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta \\ \hline \llbracket \Psi \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \theta \rrbracket = \tau \\ \hline \llbracket \Psi \rrbracket = \tau \\ \hline$$

Figure 9: Translation from source terms (e) to target terms (ι) .