

A Word Sampler for Well-Typed Functions

Breandan Considine

Syntactic Terms

Consider a simply-typed, first-order functional programming language with function calls, conditionals, and binary operators:

$$\begin{array}{ll} \text{FUN} ::= \text{fun } f0 \ (\text{PRM}) : \mathbb{T} = \text{EXP} & \text{INV} ::= \text{FID} (\text{ARG}) \\ \text{PRM} ::= \text{PID} : \mathbb{T} \mid \text{PRM} , \text{ PID} : \mathbb{T} & \text{ARG} ::= \text{EXP} \mid \text{ARG} , \text{ EXP} \\ \text{EXP} ::= \sqcap \mathbb{T} \sqcup \mid \text{PID} \mid \text{INV} \mid \text{IFE} \mid \text{OPX} & \text{OPR} ::= + \mid * \mid < \mid == \\ \text{OPX} ::= (\text{EXP} \text{ OPR } \text{ EXP}) & \text{PID} ::= p1 \mid \dots \mid pk \\ \text{IFE} ::= \text{if } \text{EXP} \{ \text{EXP} \} \text{ else } \{ \text{EXP} \} & \text{FID} ::= f0 \mid \dots \mid fn \end{array}$$

Type universe. We assume a finite universe \mathbb{T} with at two base types \mathbb{B}, \mathbb{N} , and an ambient global context Γ of named functions $f_0 : (\tau_1, \dots, \tau_m) \rightarrow \tau$.

Static Semantics

Typing judgements are standard; we highlight just a few of them below:

$$\begin{array}{c} \frac{\Gamma \vdash e_c : \mathbb{B} \quad \Gamma \vdash e_{\top} : \tau \quad \Gamma \vdash e_{\perp} : \tau}{\Gamma \vdash \text{if } e_c \{ e_{\top} \} \text{ else } \{ e_{\perp} \} : \tau} \text{ IFE} \\ \frac{\Gamma \vdash f_0 : (\tau_1, \dots, \tau_m) \rightarrow \tau \quad \Gamma \vdash e_i : \tau_i \ \forall i \in [1, m]}{\Gamma \vdash f_0 (e_1 , \dots , e_m) : \tau} \text{ INV} \\ \frac{\delta_{\text{OPR}}(\odot, \tau, \tau') = \hat{\tau} \quad \Gamma \vdash e_1 : \tau \quad \Gamma \vdash e_2 : \tau'}{\Gamma \vdash (e_1 \odot e_2) : \hat{\tau}} \text{ OPX} \end{array}$$

where the infix operator typing function δ_{OPR} is defined as follows:

$$\delta_{\text{OPR}}(\odot, \tau, \tau') = \begin{cases} \mathbb{B} \quad \odot = <, \tau = \tau' = \mathbb{B} \\ \mathbb{N} \quad \odot \in \{+, *\}, \tau = \tau' = \mathbb{N} \\ \mathbb{B} \quad \odot = ==, \tau = \tau' \end{cases}$$

Embedding the Type Checker

Typing derivations are compiled by decorating nonterminals with a pair, $\text{EXP}[\cdot, \cdot]$, carrying the type annotation, $e : \tau$, and type signature $f0 : \vec{\tau} \rightarrow \dot{\tau}$.

$$\begin{array}{c} \langle \vec{\tau}, \dot{\tau} \rangle \in \mathbb{T}^{0..k} \times \mathbb{T} \quad \vec{\tau}_{0..|\vec{\tau}|} \in \vec{\tau} \quad \text{FUN}_{\varphi} \\ \left(S_{\Gamma} \rightarrow \text{fun } f0 (\vec{\tau}_{i=1}^{|vec{\tau}|} (p_i : \vec{\tau}_i)) : \dot{\tau} = \text{EXP}[\dot{\tau}, \vec{\tau} \rightarrow \dot{\tau}] \right) \in P_{\Gamma} \\ \frac{\text{EXP}[\tau, \pi] \in V_{\Gamma} \quad \Gamma \vdash f_0 : (\tau_1, \dots, \tau_m) \rightarrow \tau}{(\text{EXP}[\tau, \pi] \rightarrow f_0 (\vec{\tau}_{i=1}^m \text{EXP}[\tau_i, \pi])) \in P_{\Gamma}} \text{ INV}_{\varphi} \\ \frac{\text{EXP}[\tau, \vec{\tau} \rightarrow \dot{\tau}] \in V_{\Gamma} \quad \tau = \dot{\tau} \quad \vec{\tau}_{0..|\vec{\tau}|} \in \vec{\tau}}{(\text{EXP}[\tau, \vec{\tau} \rightarrow \dot{\tau}] \rightarrow f0 (\vec{\tau}_{i=1}^{|vec{\tau}|} \text{EXP}[\vec{\tau}_i, \vec{\tau} \rightarrow \dot{\tau}])) \in P_{\Gamma}} \text{ REC}_{\varphi} \\ \frac{\text{EXP}[\tau, \pi] \in V_{\Gamma} \quad \tau = \tau' \quad \tau, \tau' \in \mathbb{T}}{(\text{EXP}[\tau, \pi] \rightarrow \text{if } \text{EXP}[\mathbb{B}, \pi] \{ \text{EXP}[\tau, \pi] \} \text{ else } \{ \text{EXP}[\tau', \pi] \}) \in P_{\Gamma}} \text{ IFE}_{\varphi} \\ \frac{\text{EXP}[\hat{\tau}, \pi] \in V_{\Gamma} \quad \delta_{\text{OPR}}(\odot, \tau, \tau') = \hat{\tau} \quad \odot \in \{==, <, +, *\}}{(\text{EXP}[\hat{\tau}, \pi] \rightarrow (\text{EXP}[\tau, \pi] \odot \text{EXP}[\tau', \pi])) \in P_{\Gamma}} \text{ OPX}_{\varphi} \\ \frac{\text{EXP}[\tau, \vec{\tau} \rightarrow \dot{\tau}] \in V_{\Gamma} \quad \exists \vec{\tau}_i = \tau \text{ PID}_{\varphi} \quad \text{EXP}[\tau, \pi] \in V_{\Gamma} \quad __ : \tau \in \{\mathbb{B}, \mathbb{N}\}}{(\text{EXP}[\tau, \vec{\tau} \rightarrow \dot{\tau}] \rightarrow \text{pi}) \in P_{\Gamma} \quad (\text{EXP}[\tau, \pi] \rightarrow _) \in P_{\Gamma}} \text{ PID}_{\varphi} \end{array}$$

Finally, normalize to Chomsky Normal Form (CNF), pruning unreachable and unproductive nonterminals. Expansion will be close to linear in $|G_{\Gamma}|$.

Finite Language Intersection

Context-free languages are closed under intersection with regular languages. Constructively, given CNF productions of the form $W \rightarrow XZ$ or $W \rightarrow a$, and $\alpha = \langle Q, \Sigma, \delta, q_0, F \rangle$, we build synthetic nonterminals pWr , then add:

- Binary rules $pWr \rightarrow pXq \ qZr$ for each $W \rightarrow XZ$ and $p, q, r \in Q$,
- Unit rules $pWq \rightarrow a$ for each $W \rightarrow a$ with $\delta(p, a) = q$, and
- Start rules $S \rightarrow q_a Sq_{\omega}$ for each $q_{\omega} \in F$.

This naïve construction (Salomaa, 1973) can be significantly improved via a semiring dynamic programming algorithm that avoids useless productions. When α is acyclic, $\mathcal{L}(\alpha_{\cap})$ admits efficient enumeration and exact sampling.

Autoregressive Decoding

Finite slices of a CFL are finite and therefore regular. We use star free regular expressions as a compact algebra for propagating regular constraints during parsing and decoding. Let $e : E$ be an expression defined by the grammar:

$$e \rightarrow \emptyset \mid \varepsilon \mid \Sigma \mid e \cdot e \mid e \vee e \mid e \wedge e$$

where ε is the empty symbol. We interpret these expressions as denoting regular languages, where $\mathcal{L}(x) \circ \mathcal{L}(z) := \{a \cdot b \mid a \in \mathcal{L}(x) \wedge b \in \mathcal{L}(z)\}$:

$$\begin{array}{ll} \mathcal{L}(\emptyset) = \emptyset & \mathcal{L}(x \cdot z) = \mathcal{L}(x) \circ \mathcal{L}(z)^1 \\ \mathcal{L}(\varepsilon) = \{\varepsilon\} & \mathcal{L}(x \vee z) = \mathcal{L}(x) \cup \mathcal{L}(z) \\ \mathcal{L}(a) = \{a\} & \mathcal{L}(x \wedge z) = \mathcal{L}(x) \cap \mathcal{L}(z) \end{array}$$

Brzozowski (1962) provides a rewrite-based procedure which supports both recognition and generation. He defines a quotient $\partial_a(L) = \{b \mid ab \in L\}$:

$$\begin{array}{ll} \partial_a(\emptyset) = \emptyset & \delta(\emptyset) = \emptyset \\ \partial_a(\varepsilon) = \emptyset & \delta(\varepsilon) = \varepsilon \\ \partial_a(b) = \begin{cases} \varepsilon & \text{if } a = b \\ \emptyset & \text{if } a \neq b \end{cases} & \delta(a) = \emptyset \\ \partial_a(x \cdot z) = (\partial_a x) \cdot z \vee \delta(x) \cdot \partial_a z & \delta(x \cdot z) = \delta(x) \wedge \delta(z) \\ \partial_a(x \vee z) = \partial_a x \vee \partial_a z & \delta(x \vee z) = \delta(x) \vee \delta(z) \\ \partial_a(x \wedge z) = \partial_a x \wedge \partial_a z & \delta(x \wedge z) = \delta(x) \wedge \delta(z) \end{array}$$

Now, for any nonempty (ε, \wedge) -free regex, e , to witness $\sigma \in \mathcal{L}(e)$:

$$\begin{array}{ll} \text{follow}(e) : E \rightarrow 2^{\Sigma} = \begin{cases} \{e\} & \text{if } e \in \Sigma \\ \text{follow}(x) & \text{if } e = x \cdot z \\ \text{follow}(x) \cup \text{follow}(z) & \text{if } e = x \vee z \end{cases} \\ \text{choose}(e) : E \rightarrow \Sigma^+ = \begin{cases} e & \text{if } e \in \Sigma \\ (s \leftrightarrow \text{follow}(e)) \cdot \text{choose}(\partial_s e) & \text{if } e = x \cdot z \\ \text{choose}(e' \leftrightarrow \{x, z\}) & \text{if } e = x \vee z \end{cases} \end{array}$$

This enables LTR decoding without materializing the product automaton.

Takeaways & Next Steps

- A fixed-parameter $\langle k, |\mathbb{T}| \rangle$ tractable embedding from a syntax-directed type system to a CFG with soundness and completeness guarantees.
- Intersection with an acyclic FSA yields a regular expression
- Brzozowski derivatives enable incremental autoregressive decoding
- Future work: lazy CNF materialization, quotienting (e.g., α -equivalence), richer typing (subtyping, polymorphism, substructural constraints).