

# 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:

$$\begin{array}{ll} \text{FUN} ::= \text{fun } f_0 \ ( \text{PRM} ) : \tau = \text{EXP} & \text{INV} ::= \text{FID} \ ( \text{ARG} ) \\ \text{PRM} ::= \text{PID} : \tau \mid \text{PRM} , \text{PID} : \tau & \text{ARG} ::= \text{EXP} \mid \text{ARG} , \text{EXP} \\ \text{EXP} ::= \lceil \text{T} \rfloor \mid \text{PID} \mid \text{INV} \mid \text{IFE} \mid \text{OPX} & \text{OPR} ::= + \mid * \mid < \mid == \\ \text{OPX} ::= ( \text{EXP} \text{ OPR } \text{ EXP} ) & \text{PID} ::= p_1 \mid \dots \mid p_k \\ \text{IFE} ::= \text{if } \text{EXP} \{ \text{EXP} \} \text{ else } \{ \text{EXP} \} & \text{FID} ::= f_0 \mid \dots \mid f_n \end{array}$$

**Type universe.** We assume a finite set  $\mathbb{T}$  (size  $d$ ) with at least  $\mathbb{B}, \mathbb{N}$ , and an ambient global context  $\Gamma$  of named functions  $f_i : (\tau_1, \dots, \tau_m) \rightarrow \tau$ .

## Static Semantics

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

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

where the operator typing function  $\delta_{\text{OPR}}$  is:

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

## 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  $f_0 : \vec{\tau} \rightarrow \dot{\tau}$ .

$$\begin{array}{c} \frac{\langle \vec{\tau}, \dot{\tau} \rangle \in \mathbb{T}^{0..k} \times \mathbb{T} \quad \vec{\tau}_{0..|\vec{\tau}|} \in \vec{\tau}}{\left( S_{\Gamma} \rightarrow \text{fun } f_0 ( \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} } \text{ FUN}_{\varphi} \\ \frac{\text{EXP}[\tau, \pi] \in V_{\Gamma} \quad \Gamma \vdash f_{\_} : (\tau_1, \dots, \tau_m) \rightarrow \tau}{( \text{EXP}[\tau, \pi] \rightarrow f_{\_} ( \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 f_0 ( \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}\} \lceil \mathbb{T} \lrcorner_{\varphi}}{( \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 CNF ( $G'_{\Gamma}$ ), pruning unreachable/unproductive nonterminals. In this case, the blow-up is close to linear in  $|G_{\Gamma}|$ .

## Finite Language Intersection

Context-free languages are closed under intersection with regular languages. Constructively, given CNF productions  $W \rightarrow XZ$  and  $\alpha = \langle Q, \Sigma, \delta, q_{\alpha}, F \rangle$ , we build synthetic nonterminals  $pWr$  (for all  $p, r \in Q$ ) and then add:

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

If  $\alpha$  is acyclic, the  $\mathcal{L}(\alpha_{\cap})$  admits efficient enumeration and exact sampling.

## Generalized Regular Expressions

Finite slices of a CFL are finite and therefore regular. We use GREs (union, intersection, concatenation, star, complement) as a compact algebra for propagating regular constraints during parsing and decoding.

$$\begin{array}{ll} \mathcal{L}(\emptyset) = \emptyset & \mathcal{L}(R^*) = \{\varepsilon\} \cup \mathcal{L}(R \cdot R^*) \\ \mathcal{L}(\varepsilon) = \{\varepsilon\} & \mathcal{L}(R \vee S) = \mathcal{L}(R) \cup \mathcal{L}(S) \\ \mathcal{L}(a) = \{a\} & \mathcal{L}(R \wedge S) = \mathcal{L}(R) \cap \mathcal{L}(S) \\ \mathcal{L}(R \cdot S) = \mathcal{L}(R) \times \mathcal{L}(S) & \mathcal{L}(\neg R) = \Sigma^* \setminus \mathcal{L}(R) \end{array}$$

When computing the closure of a parse chart, each chart cell carries a vector of GREs, one per nonterminal, representing exactly which terminal strings can fill that span while respecting the constraints.

## Brzozowski Differentiation

Brzozowski derivatives provide a rewrite-based procedure for incremental LTR decoding:

$$\partial_a L = \{b \in \Sigma^* \mid ab \in L\}.$$

Enables sampling from intersections without materializing product automata.

$$\begin{array}{ll} \partial_a(\emptyset) = \emptyset & \delta(\emptyset) = \emptyset \\ \partial_a(\varepsilon) = \emptyset & \delta(\varepsilon) = \varepsilon \\ \partial_a(a) = \varepsilon & \delta(a) = \emptyset \\ \partial_a(b) = \emptyset \text{ for each } a \neq b & \delta(R^*) = \varepsilon \\ \partial_a(R^*) = (\partial_x R) \cdot R^* & \delta(\neg R) = \varepsilon \text{ if } \delta(R) = \emptyset \\ \partial_a(\neg R) = \neg \partial_a R & \delta(\neg R) = \emptyset \text{ if } \delta(R) = \varepsilon \\ \partial_a(R \cdot S) = (\partial_a R) \cdot S \vee \delta(R) \cdot \partial_a S & \delta(R \cdot S) = \delta(R) \wedge \delta(S) \\ \partial_a(R \vee S) = \partial_a R \vee \partial_a S & \delta(R \vee S) = \delta(R) \vee \delta(S) \\ \partial_a(R \wedge S) = \partial_a R \wedge \partial_a S & \delta(R \wedge S) = \delta(R) \wedge \delta(S) \end{array}$$

## Takeaways & Next Steps

- A fixed-parameter tractable reduction from syntax-directed type system to a CFG with soundness and completeness guarantees
- Intersection with an acyclic regular grammar yields a regular expression
- Brzozowski derivatives enable incremental autoregressive decoding
- TODO
- TODO
- Future: lazy CNF materialization, quotienting by syntactic symmetries (e.g.,  $\alpha$ -equivalence / invariants), richer typing (subtyping, polymorphism, substructural constraints).