

Compiler verification

Verified PEG Parser

From “TRX: A Formally Verified Parser Interpreter”
A. Koprowski and H. Binsztok, 2011

Neven Villani

2021-12-14

Motivation

- At time of writing, CompCert started at the AST: the parsing step is not verified
- But programming is not done directly with ASTs
- This is an attempt at filling this missing initial step

Choice of tools

- Parser extracted from Coq
- Grammar expressed as a PEG

Parsing Expression Grammar

Definitions, Advantages

Syntax

e	$::=$	ε	empty
		$[a]$	terminal
		A	nonterminal
		$e_1; e_2$	sequence
		e_1/e_2	prioritized choice
		e^*	greedy repetition
		$!e$	lookahead

Nonprimitives

$$\begin{aligned} [“a_1 \cdots a_n”] &:= [a_1]; \cdots ; [a_n] \\ [a - z] &:= [a] / \cdots / [z] \\ e^+ &:= e; e^* \\ e? &:= e / \varepsilon \\ \&e &:= !!e \end{aligned}$$

Examples

$$\{a^n b^n \mid n\}$$

$$X_{a,b} := ([a]; X_{a,b}; [b])?$$

$$S_{a,b} := X_{a,b}; ![\cdot]$$

palindromes

$$Y := ([a]; Y; [a]) / ([b]; Y; [b]) / a / b / \varepsilon$$

$$P := Y; ![\cdot]$$

$$\{a^n b^n c^n \mid n\}$$

$$S' := \&(X_{ab}; [c]^*); ([a]^*; X_{bc}); ![\cdot]$$

Semantics

Basic rules:

$$\overline{(\varepsilon, s) \rightsquigarrow \sqrt{s}^{()}}$$

$$\overline{([\cdot], []) \rightsquigarrow \perp}$$

$$\overline{([x], x :: s) \rightsquigarrow \sqrt{s}^x}$$

$$\frac{(e_1, s) \rightsquigarrow \perp}{(e_1; e_2, s) \rightsquigarrow \perp}$$

$$\frac{(e, s) \rightsquigarrow \sqrt{s'}^v \quad (e', s') \rightsquigarrow \sqrt{s''}^{v'}}{(e; e', s) \rightsquigarrow \sqrt{s''}^{(v, v')}}}$$

$$\frac{(e, s) \rightsquigarrow \perp}{(e^*, s) \rightsquigarrow \sqrt{s}^{\square}}$$

$$\frac{(e, s) \rightsquigarrow \sqrt{s'}^v \quad (e^*, s') \rightsquigarrow \sqrt{s''}^V}{(e^*, s) \rightsquigarrow \sqrt{s''}^{v::V}}$$

Semantics (cont)

Some more important rules:

$$\frac{(e, s) \rightsquigarrow \sqrt{s}^v}{(!e, s) \rightsquigarrow \perp}$$

$$\frac{(e, s) \rightsquigarrow \perp}{(!e, s) \rightsquigarrow \sqrt{s}^{(0)}}$$

$$\frac{(e, s) \rightsquigarrow \sqrt{s'}^v}{(e/e', s) \rightsquigarrow \sqrt{s'}^v}$$

$$\frac{(e, s) \rightsquigarrow \perp \quad (e', s) \rightsquigarrow r}{(e/e', s) \rightsquigarrow m + n + 1r}$$

Advantages

- No separate lexing step
- No large parsing table
- Unambiguous, fully deterministic
- More expressive than context-free

Contribution

TRX is proven correct.

If \mathcal{G} is syntactically valid then:

$\text{well-formed}(\mathcal{G}) = \text{true} \implies \forall s.$

$(\text{parse}(\mathcal{G}, s) \text{ terminates})$

$\wedge (\forall v, s'. \text{parse}(\mathcal{G}, s) = \text{Ok}(v, s') \iff \exists m. (v_0, s) \rightsquigarrow \sqrt{v_{s'}})$

$\wedge (\text{parse}(\mathcal{G}, s) = \text{Err} \iff \exists m. (v_0, s) \rightsquigarrow \perp)$

(slightly redundant)

Appropriate characterization since \rightsquigarrow is fully deterministic

Limitations

- Exponential time (can be solved with memoization)

$$S := A; ![.]$$

$$A := ([a]; A; [b]) / ([a]; A; [c]) / \varepsilon$$

$$a^n c^n$$

- Nontermination

$$S := A; ![.]$$

$$A := A?; [a]$$

Termination

Well-formedness, Motivation, Limitations

Termination

- Source of nontermination is (mutual) left recursion
- Not so easy to check

$$A := B / (C; !D; A)$$

B can fail, C can succeed without consuming, D can fail

- Completeness is undecidable
- Well-formedness is easily checkable and implies completeness

Well-formedness

$e \in \mathbb{P}_0$ if e can accept without consuming input
WF fixpoint of well-formedness rules

$$\begin{array}{c} \overline{\varepsilon \in \text{WF}} \qquad \overline{[a] \in \text{WF}} \qquad \frac{e \in \text{WF}}{!e \in \text{WF}} \\[10pt] \frac{e \in \text{WF} \quad e \in \mathbb{P}_0 \Rightarrow e' \in \text{WF}}{e; e' \in \text{WF}} \qquad \frac{e \in \text{WF} \quad e' \in \text{WF}}{e/e' \in \text{WF}} \\[10pt] \frac{e \in \text{WF} \quad e \notin \mathbb{P}_0}{e^* \in \text{WF}} \end{array}$$

Well-formedness implies termination

All subexpressions of the grammar

$$E(\mathcal{G}) = \{e' \mid e' \sqsubseteq e_A, A \text{ nonterminal}\}$$

\mathcal{G} is well-formed if $E(\mathcal{G}) \subseteq \text{WF}$

This is an overapproximation

$A := !\varepsilon; A$ terminates but is rejected

Usability

Ergonomics, Performance, Improvements

Specification

Most parser interpreters provide a DSL to define the grammar.
Not the case of TRX.

- Grammar specified within Coq
- Correctness proof is easy
- Requires familiarity with Coq

Bootstrapping issues

- Parsing grammar description would require a type that contains itself
- Another solution: code generation (significantly more difficult)

Performance











tool	XML parser	Java parser
JAXP	2.3s. 	
JavaCC		23.0s. 
TRX-gen	5.1s. 	25.5s. 
TRX-int	40.0s. 	289.3s. 
TRX-cert	128.9s. 	662.4s. 
Mouse	206.4s. 	269.6s. 

Figure 8, page 20

Performance of TRX-cert on two examples

TRX-cert: certified version

TRX-int: prototype with similar functionality

TRX-gen: code generating parser

Performance diagnosis

- verification overhead – steps with only logical use
 - not significant
- extraction overhead – extracted code under-performs
 - conversions (probably optimized)
 - `List.rev` is quadratic
 - `Char` is a Coq natural
 - `Ascii` is 8 booleans
 - garbage collection spikes
- algorithmic overhead – suboptimal algorithm

Summary

Characterization










$\text{well-formed}(\mathcal{G}) = \text{true} \implies \forall s.$

$(\text{parse}(\mathcal{G}, s) \text{ terminates})$

$\wedge (\forall v, s'. \text{parse}(\mathcal{G}, s) = \text{Ok}(v, s') \iff \exists m. (v_0, s) \rightsquigarrow \sqrt{v_{s'}})$

$\wedge (\text{parse}(\mathcal{G}, s) = \text{Err} \iff \exists m. (v_0, s) \rightsquigarrow \perp)$

Performance

tool	XML parser	Java parser
JAXP	2.3s. 	
JavaCC		23.0s. 
TRX-gen	5.1s. 	25.5s. 
TRX-int	40.0s. 	289.3s. 
TRX-cert	128.9s. 	662.4s. 
Mouse	206.4s. 	269.6s. 