

CPEG: A Typed Tree Construction from Parsing Expression Grammars with Regex-Like Captures

Daisuke Yamaguchi

Graduate School of Engineering
Yokohama National University
Yokohama, Kanagawa, Japan
yamaguchi-daisuke-bf@ynu.jp

Kimio Kuramitsu

Department of Mathematical and Physical Sciences
Japan Women's University
Bunkyo-ku, Tokyo, Japan
kuramitsuk@fc.jwu.ac.jp

ABSTRACT

CPEG is an extended parsing expression grammar with regex-like capture annotation. Two annotations (capture and left-folding) allow a flexible construction of syntax trees from arbitrary parsing patterns. More importantly, CPEG is designed to guarantee structural constraints of syntax trees for any input strings. This reduces the amount of user code needed to check whether the intended elements exist.

To represent the structural constraints, we focus on regular expression types, a variant formalism of tree automata, which have been intensively studied in the context of XML schemas. Regular expression type is inferred from a given CPEG by the type inference that is formally developed in this paper. We prove the soundness and the uniqueness of the type inference. The type inference enables a CPEG to serve both as a syntactic specification of the input and a schematic specification of the output.

CCS CONCEPTS

• **Software and its engineering** → **Syntax; Parsers;**

KEYWORDS

Parsing, parsing expression grammars, type inference, regular expression types

ACM Reference Format:

Daisuke Yamaguchi and Kimio Kuramitsu. 2019. CPEG: A Typed Tree Construction from Parsing Expression Grammars with Regex-Like Captures. In *The 34th ACM/SIGAPP Symposium on Applied Computing (SAC '19)*, April 8–12, 2019, Limassol, Cyprus. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3297280.3297433>

1 INTRODUCTION

Regular expressions, or regexes, are a popular language tool that describe a complex text pattern and can match some set of strings. Its popularity, however, results not only from excellent pattern matching but also declarative data extraction, in which any sub-pattern within parentheses will be captured as numbered groups. Since the capture capability provides a straightforward means for extracting parsed data, many small parsers have been implemented simply with the help of regexes.

The aim of our study is to bring the regex-like capture capability into parsing expression grammars, or PEGs [4]. The reason for this aim is that PEGs are more expressive than regular expressions, which would enable us to extract some nested data (such as XML and JSON) that cannot be recognized by regexes. On the other hand, the capture capability would provide PEGs with a more straightforward means for implementing a parser than existing parser generators with action code.

Nez parser [14] is based on a declaratively extended parsing expression grammar with regex-like capture annotation. Two extended annotations (capture and left-folding) allow a flexible construction of syntax trees from arbitrary parsing patterns. We have demonstrated that Nez can parse many programming languages including Java, JavaScript, and Python.

While the declarative tree construction in Nez is convenient, it is still weak since the constructed trees are untyped. This means that they are treated as a common structure of tree data. To traverse their contents safely, the users need to check whether the intended elements exist. Parser generators such as Yacc [11] and ANTLR [18] can produce typed trees throughout their programmed action code, which seems more suitable for handling complex syntax trees. Similarly, typed trees, or trees whose structures are well guaranteed, are desirable in the declarative tree construction.

The main challenge of this paper is to provide typing rules for syntax trees that will be captured by a PEG. A critical issue is to infer a type of tree from a grammar, before constructing concrete trees. This indicates that a declarative grammar guarantees some structural constraints of parse trees of any input strings. This property could be a good foundation for further static binding with a programming language.

As the first attempt to infer a type of tree, we carefully designed CPEG to produce labeled unranked trees, which are equal to XML documents [7] (Note that CPEGs are a substantial subset of the Nez grammars). Here, we focus on *regular expression types*, or RETs, which have been studied in a foundation of XML schemas and a type system of tree automata [8]. The type inference for CPEGs developed in this paper shows that CPEGs become a schematic specification of the output and not only a syntactic specification of the input.

The remainder of this paper proceeds as follows. In Section 2, we describe our motivation using Nez grammar and RETs. In Section 3, we formally develop a CPEG as a string-to-tree transducer. In Section 4, we introduce regular expression types and define type inference rules for CPEG. In Section 5, we prove the soundness and the uniqueness of our type inference. In Section 6, we review related work, and Section 7 concludes this paper.

SAC '19, April 8–12, 2019, Limassol, Cyprus

© 2019 Copyright held by the owner/author(s). Publication rights licensed to ACM. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *The 34th ACM/SIGAPP Symposium on Applied Computing (SAC '19)*, April 8–12, 2019, Limassol, Cyprus. <https://doi.org/10.1145/3297280.3297433>.

| PEG | Type | Description |
|---------------------------------|--------------|--------------------------|
| ' ' | Primary | Matches text |
| [] | Primary | Matches character class |
| . | Primary | Any character |
| A | Primary | Nonterminal application |
| (e) | Primary | Grouping |
| e? | Unary suffix | Option |
| e* | Unary suffix | Zero-or-more repetitions |
| e+ | Unary suffix | One-or-more repetitions |
| &e | Unary prefix | And-predicate |
| !e | Unary prefix | Negation |
| e ₁ e ₂ | Binary | Sequencing |
| e ₁ / e ₂ | Binary | Prioritized Choice |

Table 1: PEG operators

2 MOTIVATING EXAMPLES

We describe a motivation of tree construction and its typing in the context of PEGs. Here we use Nez grammar [13], which is an open source implementation of CPEG.

2.1 Parsing Expression Grammars

Nez grammar is a PEG-based grammar specification language, whose constructs come from those of PEGs. Nez grammar is a set of syntax rules that are defined by a mapping from a non-terminal A to a parsing expression e :

$$A = e$$

Table 1 shows a list of PEG operators used in Nez which inherits the formal interpretation of PEGs [4]. This indicates that the string 'abc' exactly matches the same input, while [abc] matches one of these characters. The . operator matches any single character. The e?, e*, and e+ expressions behave as in common regular expressions, except that they are greedy and match until the longest position. The e₁ e₂ attempts two expressions e₁ and e₂ sequentially, backtracking to the starting position if either expression fails. The choice e₁ / e₂ first attempt e₁ and then attempt e₂ if e₁ fails. The expression &e attempts e without consuming any character. The expression !e succeeds if e fails but fails if e succeeds.

In general, PEGs can express all languages that can be expressed by deterministic context-free grammars (such as LALR and LL(k) grammars).

PEGs, on the other hand, provide no specification for the output of a parser, while the parser users require certain forms of syntax trees that contain all necessary information for further processing. Note that we can regard non-terminals as a tree constructor in a way that a labeled tree node A[...] is constructed from A = e. This approach is similar to that of derivation trees, resulting in redundantly nested trees. Besides, some forms of trees are not well constructed, as described in Section 2.2.2.

2.2 Tree Annotation

PEG only provides the syntactic matching capability while Nez provides two additional annotations (called capture and fold-capture) to construct complex syntax trees in a parser context.

2.2.1 Capture. The capture annotation {e #L} is a straightforward extension to specify a parsing expression e (with enclosing braces { }) that extracts its matched string as a tree node. The extracted node is labeled by the given #L, which is used to identify the type of tree nodes.

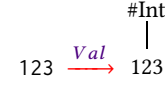
To start, let us consider a simple example, VAL rule, which recognizes a sequence of numeric characters. Here, Val is a capture version of VAL whose matched strings are constructed as a node labeled as #Int.

Let us start a token extraction.

$$\begin{aligned} VAL &= [0-9]^+ \\ Val &= \{ [0-9]^+ \#Int \} \end{aligned}$$

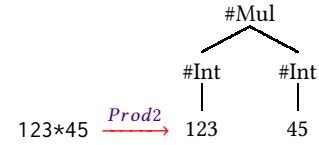
We use $x \xrightarrow{e} v$ to write that an expression e parses an input string x and then transforms it into a tree v.

Here, we show how Val parses the string 123:



The nested capture is naturally interpreted as a nested containment of trees that are constructed inside the braces. The following Prod2 parses from 123*45 to a tree depicted in the following manner:

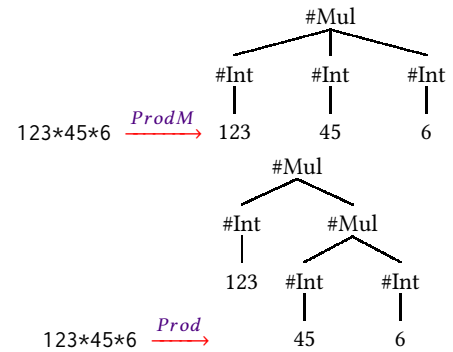
$$Prod2 = \{ Val '*' Val \#Mul \}$$



The Prod2 accepts only a single multiplication expression. Multiple multiplications such as 1*2*5 can be captured in two different forms.

$$\begin{aligned} ProdM &= \{ Val ('*' Val)^* \#Mul \} \\ Prod &= \{ Val ('*' Prod) \#Mul \} / Val \end{aligned}$$

The ProdM and Prod both accept the same inputs, while they produce different forms of trees. The ProdM uses a repetition which forms a variable-length list of trees. The Prod, on the other hand, uses a recursion which forms a recursively nested tree.



As shown, we can switch a form of trees by repetition and recursion.

2.2.2 *Fold-Capture*. In the previous section, we see that a recursion of capturing produces a nested form of recursive trees. It is important to note that the produced trees are always in the right-associative form since PEGs do not allow left recursions.

$$ProdL = \{ (ProdL / Val) '*' Val \#Mul \}$$

This example suggests that the capture does not well describe arbitrary forms of syntax trees, especially left-associative binary operators.

CPEG additionally provides the fold-capture annotation, denoted as $\wedge\{e \#L\}$, which constructs a tree containing the left-hand tree as its first subtree.

Here is an example of a fold-capture version of *Prod2*, where the first *Val* is factored outside and then folded inside by \wedge .

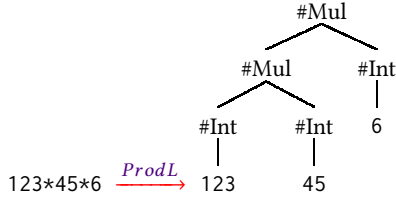
$$Prod2 = Val \wedge\{ '*' Val \#Mul \}$$

The \wedge operator is left-associative. The trees on the left hand are folded one after another.

$$((Val \wedge\{ '*' Val \#Mul \}) \dots) \wedge\{ '*' Val \#Mul \}$$

We use the repetition to denote this iterative folding. Finally, the left-associative version of *Prod* is described as follows:

$$ProdL = Val (\wedge\{ '*' Val \#Mul \})^*$$



In fact, left-recursion is a significant restriction of PEG. Although there is a known algorithm for eliminating any left-recursion from a grammar [19], this algorithm does not ensure the left-associative. Using fold-capture annotation, we can replace the left-recursion under the impression of keeping the associativity.

2.3 Typing with Regular Expression Types

The declarative tree annotations in Nez are convenient and powerful enough to express many types of syntax trees, ranging from XML and JSON to Java and JavaScript [14]. Trees that the parser users receive are untyped, which are the so-called *common trees* that are formed in a common structure.

To traverse the trees safely, the users need to check whether the traversed tree is in an intended structure (i.e. its label and arity), as an XML schema variation. Nevertheless, embedding the checking code in a traversal function makes the program fuzzier and prone to errors [20].

Besides, the traversal function should be implemented carefully so that the function is exhaustive—that is, some rule should be applied for all possible structured input trees. Static analysis techniques for checking exhaustiveness checking are based on a type of input value. For instance, XDuce [8], which is a statically typed XML processing language, provides an exhaustiveness checker that is found on a type called regular expression type [9]. Regular expression types are a variant expression of tree automata [2], which has been developed in the contexts of XML schema validations.

As our first attempt at the schematic variation on the trees and the static exhaustiveness checking, we propose typing rules that infer a regular expression type for a given grammar.

Using regular expression types, all trees that can be parsed from *ProdM* in the previous section have a type *ProdM*, which is defined as follows:

$$\begin{aligned} \text{type ProdM} &= \mathbf{Mul}[\mathbf{Val}, \mathbf{Val}^*] \\ \text{type Val} &= \mathbf{Int}[\mathbf{Empty}] \end{aligned}$$

Here, the type $\mathbf{Mul}[\mathbf{Val}, \mathbf{Val}^*]$ describes a tree with \mathbf{Mul} that has subtrees typed by the elements of $[\quad]$. The $'$, $'^*$, and \mathbf{Empty} inside respectively denotes *concatenation*, *repetition*, and *empty tree* as regular expression operators. This indicates that the \mathbf{Mul} tree has one or more subtrees that are typed with \mathbf{Val} , which is a type variable defined in the second line. The type $\mathbf{Int}[\mathbf{Empty}]$ represents no more subtrees and its type is \mathbf{Int} .

Regular expression types allow recursive type definition. All trees that can be parsed from *ProdL* are described in the following type:

$$\begin{aligned} \text{type ProdL} &= \mathbf{Prod}[\mathbf{ProdL}, \mathbf{Val}] \mid \mathbf{Val} \\ \text{type Val} &= \mathbf{Int}[\mathbf{Empty}] \end{aligned}$$

3 FORMAL DEFINITION OF CPEG

CPEG is a subset of Nez grammar [14] to highlight a focused extension to the construction of labeled unranked trees.

3.1 Grammar

A CPEG G is a 5-tuple $G = (N_G, \Sigma, P_G, e_s, S)$, where N_G is a finite set of non-terminal symbols, Σ is a finite set of terminal symbols, P_G is a finite set of *production rules*, e_s is a *start expression* and S is a finite set of label symbols.

Each production rule $r \in P_G$ is a pair (A, e) , that is written as $A \leftarrow e$, where $A \in N_G$ and e is an expression. For any $A \in N_G$, there is exactly one e such that $A \leftarrow e$. We regard P_G as a function from non-terminals to expressions.

The syntax of an expression e is defined in Figure 1.

| | | |
|-------------------------|------------------|----------------|
| $e ::= \varepsilon$ | | empty |
| a | $(a \in \Sigma)$ | terminal |
| A | $(A \in N_G)$ | nonterminal |
| $e_1 e_2$ | | sequence |
| e_1 / e_2 | | ordered choice |
| e^* | | repetition |
| $!e$ | | not-predicate |
| $\{e \#L\}$ | $(L \in S)$ | capture |
| $e_1 \wedge\{e_2 \#L\}$ | $(L \in S)$ | fold-capture |

Figure 1: Expressions in CPEG

The syntax of CPEGs is reasonably simple. Due to the following syntactic sugars (also defined in [4]), CPEG expresses all the parsing expressions that are constituted by PEG operators (also in Table 1).

$$\begin{aligned}
'abc' &= 'a' 'b' 'c' \\
[abc] &= 'a' / 'b' / 'c' \\
e^+ &= e e^* \\
e? &= e / \varepsilon \\
\&e &= !(e)
\end{aligned}$$

Similarly, the fold-capture annotation can be rewritten by the capture annotation, such as:

$$\begin{aligned}
e_1 \wedge \{e_2 \#L\} &= \{e_1 e_2 \#L\} \\
e_1(\wedge \{e_2 \#L\}) / (\wedge \{e_3 \#L'\}) &= \{e_1 e_2 \#L\} / \{e_1 e_3 \#L'\} \\
e_1(\wedge \{e_2 \#L\})? &= \{e_1 e_2 \#L\} / e_1
\end{aligned}$$

An important exception is the repetitive combination such as $e_1(\wedge \{e_2 \#L\})^*$, which can not be rewritten by the capture annotation. Thus, we focus only on the repetitive fold-capture as defined in Figure 1. Note that $e_1(\wedge \{e_2 \#L\})^*$ is a syntax sugar of $\{A e_2 \#L\}$ where $A = \{A e_2 \#L\} / e_1$. However, this syntactic sugar contains a left recursion that never terminates in PEGs.

Well-Formed Grammars. As will be discussed in Section 4, a regular expression type is inferred for each CPEG. Some inferred types, however, do not hold a *well-formedness* condition [9] without any syntactic restriction on CPEG. Here, we impose an additional restriction to the syntax of CPEG. We refer to the CPEG that holds the restriction as *well-formed grammar*.

Well-formed grammar is a CPEG that holds either of the following two conditions:

- recursive use of nonterminals occurs only in tail positions.
- if a sequence $A e$ exists in the production rules, where A is a recursively used nonterminal and e is a some expression (i.e., A is not in tail positions); then e has neither capture nor fold-capture as a subexpression.

For instance, the CPEG that has a production rule as the following:

$$A \leftarrow \{e_1 \#L_1\} A \{e_2 \#L_2\} / e$$

is not a well-formed grammar, because A is not in the tail position and the subsequent expression of A is a fold-capture.

Whereas, the CPEG that has a production rule as follows:

$$A \leftarrow \{e_1 \#L_1\} A / e$$

is a well-formed grammar, because A is in the tail position.

Furthermore, the CPEG which have a production rule as follows:

$$A \leftarrow \{e_1 \#L_1\} A a / e$$

is a well-formed grammar while A is not in the tail position, because the subsequent expression of A is a terminal a .

Naturally, the CPEG which have production rules:

$$A \leftarrow \{e_1 \#L_1\} B a / e$$

$$B \leftarrow \{e_2 \#L_2\} A / e'$$

is a well-formed grammar as well.

3.2 Formal Interpretation of CPEG

The semantics is defined as a relation from expressions to trees.

3.2.1 Tree. We start by defining a textual notation to denote a tree for convenience.

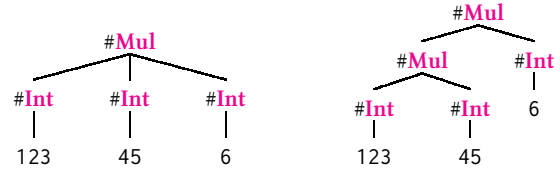
The tree of CPEG is a labeled unranked tree on \mathcal{S} and Σ^* , where the nodes are labeled with $L \in \mathcal{S}$ and the leaves hold a string $x \in \Sigma^*$.

Let $\mathcal{T}_{\mathcal{S} \times \Sigma^*}$ be a set of trees on \mathcal{S} and Σ^* . The syntax of a tree $v \in \mathcal{T}_{\mathcal{S} \times \Sigma^*}$ is defined as follows.

$$\begin{array}{ll}
v ::= L[v] & (L \in \mathcal{S}) \quad \text{node} \\
v_1, v_2 & \quad \quad \quad \text{concatenation} \\
x & (x \in \Sigma^*) \quad \text{string}
\end{array}$$

Node $L[v]$ denotes the node that is labeled with a symbol L and has subtrees denoted by v . We use a concatenation operator $,$ to handle multiple subtrees. Since the order of subtrees is preserved, the concatenation $,$ is not commutative. The notation x denotes a string including empty string. We assume that a concatenation of two strings is equal to a single string. That is x, x' is equal to xx' . Additionally, we assume a concatenation of a node and a string is equal to the node. That is $L[v], x$ and $x, L[v]$ are equal to $L[v]$. This premise reads that if a tree $L[v], x$ or $x, L[v]$ is constructed, the tree can be regarded as $L[v]$. In other words, a string that is concatenated with a node can be ignored.

Here, we show the same trees with both pictorial and textual notations as follows.



$$\text{Mul}[\text{Int}[123], \text{Int}[45], \text{Int}[6]] \quad \text{Mul}[\text{Mul}[\text{Int}[123], \text{Int}[45]], \text{Int}[6]]$$

3.2.2 Operational Semantics. To formalize the semantics of a grammar G , we define a relation \Downarrow_G from (e, x) to (o, y) , where e is an expression of CPEG, x is an input string, y is an unconsumed string, and $o \in \mathcal{T}_{\mathcal{S} \times \Sigma^*} \cup \{f\}$ is an output. If o is in the form of a tree v , the matching succeeds and the tree v is constructed from an input string x . If $o = f$, the distinguished symbol f indicates *failure*.

For $((e, x), (o, y)) \in \Downarrow_G$, we write $e \Downarrow_y^x o$. This is read as “an expression e parses an input x and transforms it to an output o with an unconsumed string y ”.

Now supposing $v \in \mathcal{T}_{\mathcal{S} \times \Sigma^*}$, $a, b, c \in \Sigma$, $x, y, z \in \Sigma^*$, and ε is an empty string, \Downarrow_y^x is the smallest relation closed under the set of rules shown in Figure 2.

The matching interpretation in CPEG is the same as PEG [4], except for some tree construction. The capture and fold-capture are an explicit tree constructor as in (E-CAPTURE1) and (E-FOLD CAP1). All the trees constructed in the subexpressions are contained in a newly constructed tree.

The rules for fold-capture $e_1 \wedge^* \{e_2 \#L\}$ are (E-FOLD CAP1), (E-FOLD CAP2) and (E-FOLD CAP3). By (E-FOLD CAP1), a left-associative tree is derived from $e_1 \wedge^* \{e_2 \#L\}$, if the presuppositions of the rule are satisfied. The presuppositions are summarized in two conditions. The first condition is that the subexpression e_1 derives a tree v_1 from an input string $xy_1y_2 \dots y_nz$. Now we denote the unconsumed string by $y_1y_2 \dots y_nz$. The second condition is that subexpression e_2 derives a tree v_2 from the string $y_1y_2 \dots y_nz$. Additionally, if tree v_3 ,

$$\begin{array}{c}
\frac{}{\varepsilon \Downarrow_x^x \varepsilon} \text{ (E-EMPTY)} \qquad \frac{e_1 \Downarrow_y^{x_1 y} v_1 \quad e_2 \Downarrow_y^y f}{e_1 e_2 \Downarrow_{x_1 y}^{x_1 y} f} \text{ (E-SEQ3)} \qquad \frac{e \Downarrow_x^x f}{e^* \Downarrow_x^x \varepsilon} \text{ (E-REP2)} \\
\frac{}{a \Downarrow_x^{ax} a} \text{ (E-TERM1)} \qquad \frac{e_1 \Downarrow_y^{x_1 y} v_1}{e_1 / e_2 \Downarrow_y^{x_1 y} v_1} \text{ (E-ALT1)} \qquad \frac{e \Downarrow_y^{xy} v}{!e \Downarrow_{xy}^{xy} f} \text{ (E-NOT1)} \\
\frac{}{a \Downarrow_x^{bx} f \quad (a \neq b)} \text{ (E-TERM2)} \qquad \frac{e_1 \Downarrow_{xy}^{xy} f \quad e_2 \Downarrow_y^{x_2 y} v_2}{e_1 / e_2 \Downarrow_{xy}^{xy} v_2} \text{ (E-ALT2)} \qquad \frac{e \Downarrow_x^x f}{!e \Downarrow_x^x \varepsilon} \text{ (E-NOT2)} \\
\frac{P_G(A) \Downarrow_y^x v}{A \Downarrow_y^x v} \text{ (E-NT)} \qquad \frac{e_1 \Downarrow_x^{xy} f \quad e_2 \Downarrow_x^x f}{e_1 / e_2 \Downarrow_x^{xy} f} \text{ (E-ALT3)} \qquad \frac{e \Downarrow_y^{xy} v}{\{e \# L\} \Downarrow_y^{xy} L[v]} \text{ (E-CAPTURE1)} \\
\frac{e_1 \Downarrow_{x_2 y}^{x_1 x_2 y} v_1 \quad e_2 \Downarrow_y^{x_2 y} v_2}{e_1 e_2 \Downarrow_y^{x_1 x_2 y} v_1, v_2} \text{ (E-SEQ1)} \qquad \frac{e \Downarrow_x^x f \quad e^* \Downarrow_y^{x_2 y} v_2}{e^* \Downarrow_y^{x_1 x_2 y} v_1, v_2} \text{ (E-REP1)} \qquad \frac{e \Downarrow_x^x f}{\{e \# L\} \Downarrow_x^x f} \text{ (E-CAPTURE2)} \\
\frac{e_1 \Downarrow_x^x f}{e_1 \wedge^* \{e_2 \# L\} \Downarrow_x^x f} \text{ (E-FOLDCAP2)} \qquad \frac{e_1 \Downarrow_{y_1 y_2 \dots y_n z}^{x_1 y_1 y_2 \dots y_n z} v_1 \quad e_2 \Downarrow_{y_2 \dots y_n z}^{y_1 y_2 \dots y_n z} v_2 \quad \dots \quad e_2 \Downarrow_z^{y_n z} v_n \quad e_2 \Downarrow_z^z f}{e_1 \wedge^* \{e_2 \# L\} \Downarrow_z^{x_1 y_1 y_2 \dots y_n z} L[L[\dots L[L[v_1, v_2], v_3], \dots, v_{n-1}], v_n]} \text{ (E-FOLDCAP1)} \\
\frac{e_1 \Downarrow_x^x f}{e_1 \wedge^* \{e_2 \# L\} \Downarrow_x^x f} \text{ (E-FOLDCAP3)}
\end{array}$$

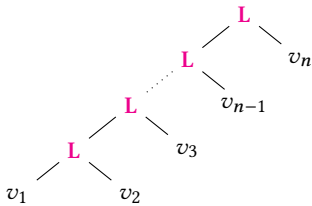
Figure 2: Rules for the relation \Downarrow_y^x

tree v_4, \dots and tree v_n are derived, by applying e_2 repeatedly to the unconsumed string until the out put goes failure, then the derived tree of $e_1 \wedge^* \{e_2 \# L\}$ is a left-associative tree:

$$L[L[\dots L[L[v_1, v_2], v_3], \dots, v_{n-1}], v_n].$$

$n-1$

The tree v_1 is contained at the left branch that is $n-1$ times nested from the root. The other tree v_i ($2 \leq i \leq n$) is stored at the right branch that is $n-i+1$ times nested from the root. The derived tree is also described pictorially as below.



By (E-FOLDCAP2), if the subexpression e_1 goes to failure, $e_1 \wedge^* \{e_2 \# L\}$ also goes to failure. By (E-FOLDCAP3), if the subexpression e_1 derives v_1 but the subexpression e_2 goes to failure, $e_1 \wedge^* \{e_2 \# L\}$ derives just v_1 . Note that v_1 is not nested in the node with L .

4 TYPING CPEG

4.1 Regular Expression Types

In this section, we review the syntax and the semantics of regular expression types (RETs) [9].

$$\begin{array}{c}
x : \text{Empty (S-EMPTY)} \qquad \frac{v_i : T \text{ for each } i}{v_1, \dots, v_n : T^*} \text{ (S-REP)} \\
\frac{v_1 : T_1 \quad v_2 : T_2}{v_1, v_2 : T_1, T_2} \text{ (S-SEQ)} \\
\frac{v : T_1}{v : T_1 | T_2} \text{ (S-OR1)} \qquad \frac{v : T}{L[v] : L[T]} \text{ (S-NODE)} \\
\frac{v : T_2}{v : T_1 | T_2} \text{ (S-OR2)} \qquad \frac{E(X) = T \quad v : T}{v : X} \text{ (S-VAR)}
\end{array}$$

Figure 3: Typing rules for trees

A RET T is inductively defined as follows:

| | | |
|---------|----------------------|----------------|
| $T ::=$ | Empty | empty sequence |
| | T_1, T_2 | concatenation |
| | $T_1 T_2$ | union |
| | T^* | repetition |
| | $L[T_1, \dots, T_n]$ | label |
| | X | type variable |

Here, the semantics of RETs is given by the relation $v : T$ ($v \in \mathcal{T}_{S \times \Sigma^*}$), read “the tree v has type T ”—the smallest relation closed under the set of typing rules in Figure 3.

The empty sequence Empty, the concatenation T_1, T_2 , the union $T_1 | T_2$, and the repetition T^* come from regular expressions, as its

name implies. The label $\mathbf{L}[T]$ represents a tree that contains a subtree that has a type T .

The X represents a type variable that binds an arbitrary RET. The bindings of type variables are given by a single global set E of type definitions of the following form:

$$\text{type } X = T$$

We regard E as a mapping from type variables to their bodies and write $E(X)$ for a reference to the mapped type from X in E .

As shown above, the type variables allow recursive types. The readers should note that the expressiveness of RETs correspond to some of the context-free grammars, despite the fact that its name is regular.

4.2 Subtyping

The subtype relation between two types is defined semantically: two types S, T are in the subtype relation $<$: if and only if $v : S$ implies $v : T$

$$S <: T \iff \forall v \in \mathcal{T}_{S \times \Sigma^*}. v : S \Rightarrow v : T$$

As Hosoya [9] indicated, the semantic notion of subtype relation over RETs immediately corresponds to the notion of inclusion relation on the set theory.

4.3 Type Inference for CPEG

CPEG intends to infer the types of trees from a grammar before constructing concrete trees.

To start, we consider that a tree v is derived from an input x with a CPEG $G = (N_G, \Sigma, P_G, e_s, \mathcal{S})$. Since the tree is constructed by the derivation $e_s \Downarrow_y^x v$, the type of G is regarded as the type of the start expression e_s .

Now, we define type inference rules for CPEG expressions. Let E be a single global set of type bindings and Γ be a type environment mapping from non-terminals to type variables. The mappings in Γ are denoted by $A : X_A$.

The inference rules are defined as a typing relation denoted by $\Gamma \vdash e : T \mid_\chi E$, which can be read “under a typing environment Γ , an expression e has a type T with a global set E ”. The relations are the smallest relation closed under the set of typing rules shown in Figure 4.

The χ is a set of type variables. The χ is used to store the type variables introduced in each subderivation, and $T \mid_\chi E$ ensure that the variables appearing in T are fresh, for each condition of χ prevent us from building a derivation in which the same variable is used as “fresh” in two different places. Since there is an infinite supply of type variable names, we can always find a way to satisfying the condition. These conventions come from Pierce’s textbook [21].

The rules (T-NT1) and (T-NT2) are types for nonterminal symbols. The premise $A : T \notin \Gamma$ in (T-NT1) is simply an explicit reminder: first, we check if the relation $A : T$ is included in Γ ; if not, we add that relation to Γ and then try to calculate T for $P_G(A)$. In the rule (T-NT1), the sets E and χ are updated to $E \cup \{\text{type } X_A = T\}$ and $\chi \cup \{X_A\}$ respectively. The rules (T-SEQ), (T-ALT), and (T-FOLD CAP) update the global set and χ as well.

$$\begin{array}{c} \vdash \varepsilon : \text{Empty} \mid_\emptyset \emptyset \quad (\text{T-EMPTY}) \\ \vdash a : \text{Empty} \mid_\emptyset \emptyset \quad (\text{T-TERM}) \\ \frac{A : X_A \notin \Gamma \quad \Gamma, A : X_A \vdash P_G(A) : T \mid_\chi E \quad \{X_A\} \cap \chi = \emptyset}{\Gamma \vdash A : X_A \mid_{\chi \cup \{X_A\}} E \cup \{\text{type } X_A = T\}} \quad (\text{T-NT1}) \\ \frac{A : X \in \Gamma}{\Gamma \vdash A : X \mid_\emptyset \emptyset} \quad (\text{T-NT2}) \\ \frac{\Gamma \vdash e_1 : T_1 \mid_{\chi_1} E_1 \quad \Gamma \vdash e_2 : T_2 \mid_{\chi_2} E_2 \quad \chi_1 \cap \chi_2 = \emptyset}{\Gamma \vdash e_1 e_2 : T_1, T_2 \mid_{\chi_1 \cup \chi_2} E_1 \cup E_2} \quad (\text{T-SEQ}) \\ \frac{\Gamma \vdash e_1 : T_1 \mid_{\chi_1} E_1 \quad \Gamma \vdash e_2 : T_2 \mid_{\chi_2} E_2 \quad \chi_1 \cap \chi_2 = \emptyset}{\Gamma \vdash e_1 / e_2 : T_1 | T_2 \mid_{\chi_1 \cup \chi_2} E_1 \cup E_2} \quad (\text{T-ALT}) \\ \frac{\Gamma \vdash e : T \mid_\chi E}{\Gamma \vdash e * : T^* \mid_\chi E} \quad (\text{T-REP}) \\ \vdash !e : \text{Empty} \mid_\emptyset \emptyset \quad (\text{T-NOT}) \\ \frac{\Gamma \vdash e : T \mid_\chi E}{\Gamma \vdash \{e \# \mathbf{L}\} : \mathbf{L}[T] \mid_\chi E} \quad (\text{T-CAPTURE}) \\ \frac{\Gamma \vdash e_1 : T_1 \mid_{\chi_1} E_1 \quad \Gamma \vdash e_2 : T_2 \mid_{\chi_2} E_2 \quad \chi_1 \cap \chi_2 = \emptyset \quad \chi_1 \cap \{X\} = \emptyset \quad \{X\} \cap \chi_2 = \emptyset}{\Gamma \vdash e_1 \# \{e_2 \# \mathbf{L}\} : X \mid_{\chi_1 \cup \chi_2 \cup \{X\}} E_1 \cup E_2 \cup \{\text{type } X = \mathbf{L}[X, T_2] \mid T_1\}} \quad (\text{T-FOLD CAP}) \end{array}$$

Figure 4: Typing rules associated with a single global set E

Now, let us consider the typing of the following CPEG:

$$\begin{aligned} G &= (\{\text{Prod}, \text{Val}\}, \{0, 1, 2, \dots, 9, *\}, P_G, \text{Prod}, \{\text{Prod}, \text{Int}\}) \\ P_G &= \{\text{Prod} \leftarrow \text{Val} \# \{ * \text{Val} \# \text{Prod}\}, \\ &\quad \text{Val} \leftarrow \{[0 - 9] \# \text{Int}\}\} \end{aligned}$$

Note that $[0 - 9]$ is a derived form of $0/1/2/3/\dots/7/8/9$.

The derivation of typing is shown in Figure 5.

First, $\{[0 - 9] \# \text{Int}\}$ has the type $\text{Int}[\text{Empty}]$ by (T-CAPTURE), where both E and χ are empty. Although $[0 - 9]$ in this premise has some type, we omit this derivation for readability. We derive that Val has a type variable X_3 by (T-NT1). At this derivation, the empty global set is updated to E_2 as shown in the bottom of Figure 5. Moreover, the empty χ is updated to $\{X_3\}$. Note that X_3 is newly introduced in this step.

The next derivation $*\text{Val}$ is a little complex. Hence, we make a sub-derivation tree \mathcal{D} . The sub-derivation proceeds similarly from the upper left of the tree to the bottom. As a result of \mathcal{D} , the following is derived.

$$\text{Prod} : X_1 \vdash * \text{Val} : \text{Empty}, X_4 \mid_{\{X_4\}} E_3$$

where, $E_3 = \{\text{type } X_4 = \text{Int}[\text{Empty}]\}$.

$$\frac{\frac{\frac{\vdots}{\text{Prod} : X_1, \text{Val} : X_3 \vdash \{[0-9] \# \mathbf{Int}\} : \mathbf{Int}[\text{Empty}] \mid \emptyset \emptyset} \text{(T-CAPTURE)}}{\text{Prod} : X_1 \vdash \text{Val} : X_3 \mid_{\{X_3\}} E_2} \text{(T-NT1)}}{\text{Prod} : X_1 \vdash \text{Val} \overset{*}{\wedge} \{ * \text{Val} \# \mathbf{Prod} \} : X_2 \mid_{\{X_2, X_3, X_4\}} E_1} \mathcal{D} \text{(T-FOLDCAP)}}{\vdash \text{Prod} : X_1 \mid_{\{X_1, X_2, X_3, X_4\}} E} \text{(T-NT1)}$$

Where, $\mathcal{D} =$

$$\frac{\frac{\frac{\vdots}{\text{Prod} : X_1, \text{Val} : X_4 \vdash \{[0-9] \# \mathbf{Int}\} : \mathbf{Int}[\text{Empty}] \mid \emptyset \emptyset} \text{(T-CAPTURE)}}{\text{Prod} : X_1 \vdash \text{Val} : X_4 \mid_{\{X_4\}} E_3} \text{(T-NT1)}}{\text{Prod} : X_1 \vdash * : \text{Empty} \mid \emptyset \emptyset} \text{(T-TERM)}}{\text{Prod} : X_1 \vdash * \text{Val} : \text{Empty}, X_4 \mid_{\{X_4\}} E_3} \text{(T-SEQ)}$$

$E = \{\text{type } X_1 = X_2, \text{type } X_2 = \mathbf{Prod}[X_2, \text{Empty}, X_4] \mid X_3, \text{type } X_3 = \mathbf{Int}[\text{Empty}], \text{type } X_4 = \mathbf{Int}[\text{Empty}]\}$

$E_1 = \{\text{type } X_2 = \mathbf{Prod}[X_2, \text{Empty}, X_4] \mid X_3, \text{type } X_3 = \mathbf{Int}[\text{Empty}], \text{type } X_4 = \mathbf{Int}[\text{Empty}]\}$

$E_2 = \{\text{type } X_3 = \mathbf{Int}[\text{Empty}]\}$

$E_3 = \{\text{type } X_4 = \mathbf{Int}[\text{Empty}]\}$

Figure 5: A derivation tree that derives $\vdash \text{Prod} : X_1 \mid_{\{X_1, X_2, X_3, X_4\}} E$. This tree is omitted the freshness condition of type variables for simplicity.

Now, we can go back to the main derivation tree. The following is derived by (T-FOLDCAP).

$$\text{Prod} : X_1 \vdash \text{Val} \overset{*}{\wedge} \{ * \text{Val} \# \mathbf{Prod} \} : X_2 \mid_{\{X_2, X_3, X_4\}} E_1$$

where, $E_1 = \{\text{type } X_2 = \mathbf{Prod}[X_2, \text{Empty}, X_4] \mid X_3, \text{type } X_3 = \mathbf{Int}[\text{Empty}], \text{type } X_4 = \mathbf{Int}[\text{Empty}]\}$.

Finally, applying (T-NT1), X_1 and the global set E :

$$\left\{ \begin{array}{ll} \text{type } X_1 = X_2, & \text{type } X_2 = \mathbf{Prod}[X_2, \text{Empty}, X_4] \mid X_3, \\ \text{type } X_3 = \mathbf{Int}[\text{Empty}], & \text{type } X_4 = \mathbf{Int}[\text{Empty}] \end{array} \right\}$$

are derived.

5 PROPERTIES OF TYPING RULES FOR CPEG

In this section, we show type soundness and uniqueness of types property of the typing rules.

First, we have shown the *uniqueness of types* property, which means that the inferred type is always unique for a given CPEG. Furthermore, this property says that the derivation tree is also deterministic.

THEOREM 5.1 (UNIQUENESS OF TYPES). *In a given typing context Γ , an expression e has uniquely one type with type variables all in the domain of a global set of types E . Moreover, there is just one derivation of this typing built from the typing rule that derives the typing relation.*

PROOF. The proof goes by the structural induction on e . \square

Next, we have shown type soundness. Intuitively, type soundness states that if a CPEG has type T under the type inference, then the trees derived from the CPEG are typed by the type T .

THEOREM 5.2 (SOUNDNESS). *Let G be a CPEG such that $G = (N_G, \Sigma, P_G, e_s, S)$. Let E be a global set of regular expression types.*

$$\forall v \in \{v \mid \exists x, y, e_s \Downarrow_y^x v\}. \Gamma \vdash e_s : T \mid_{\mathcal{X}} E \Rightarrow v : T$$

PROOF. By induction on derivation of $e \Downarrow_y^x v$. The (E-EMPTY), (E-TERM1), (E-NOT2) cases are immediate by the rules (T-EMPTY), (T-TERM), (T-NOT), (S-EMPTY). For the other cases, we will discuss the (E-CAPTURE1) case and the proof for the remaining rules proceed in the same manner.

Case E-CAPTURE1: The start expression is $\{e \# \mathbf{L}\}$ and the derived tree is $\mathbf{L}[v]$ such that $e \Downarrow_y^x v$. Let the following be an induction hypothesis.

$$\Gamma \vdash e : T \mid_{\mathcal{X}} E \Rightarrow v : T$$

The only typing rule for $\{e \# \mathbf{L}\}$ is the rule (T-CAPTURE). By inversion of (T-CAPTURE),

$$\Gamma \vdash e : T \mid_{\mathcal{X}} E$$

Using the induction hypothesis, v has type T . Finally, $\mathbf{L}[v] : \mathbf{L}[T]$ by (S-NODE). \square

6 RELATED WORK

Here we review parsing with an emphasis on tree construction technique.

PEGs have been gaining popularity among many language developers since it was presented by B. Ford in 2004. In the context of PEGs, a parser generator is a standard approach to the development of parsers [6]. PEGs, as well as other CFG-based grammars, can well describe the syntactic pattern of the input while providing a very poor specification of the output. Consequently, the construction of syntax trees is mostly left to the embedded action code, which is usually written in a target programming language. However, the use of action code reduces a good property of declarative grammar specification. Moreover, their additional code generation and compilation process are cumbersome compared to regex-style matching and capturing.

Parser combinators such as Parsec [15], FParsec [24], and Scala Parser Combinators [17] provide a more integrated means for writing a recursive descent parser. Although parser combinators are not always based on some grammar formalism, many PEG parsers have been implemented by the combinators. In parser combinators, the resulting syntax trees can be well-typed. The tree construction, however, relies on code fragments that directly manipulates parsing results. Consequently, the specification of a parser is still hard to maintain [1, 12].

More recently, declarative parsing (no action code) has been focused in many grammar formalisms, since action code makes it difficult to maintain a parser specification and reduces grammar reusability [1, 12]. In the contexts of PEGs, LPeg [10] is implemented as a PEG-based pattern matching tool that provides a grouped capture like regex. Moreover, Nez grammar [13, 14], an ancestor of CPEG, provides a structured capture that can construct complex syntax trees. These parse data require no action code, but they are untyped.

For the parser users, types are significant. This viewpoint has made another attempt to parser generations from data types. Notably, PADS/ML [16] describes DDC^α [3, 16] and then generates a parser from the data specification. The “parser from types” approach can follow various syntax patterns including programming languages. However, formally specifying the provided types safety would be more challenging [20].

Finally, a type system for grammar is new. As a starting point, we use RETs, which have intensively been studied in the context of XML schemas and tree automata [8, 9]. We consider that RETs are a straightforward type representation of syntax trees, and RETs can make a theoretical bridge between declarative parsing, tree automata, and programming language design. Indeed, binding RETs with ML and OCaml has been reported in [5, 22].

7 CONCLUSION

Regular expressions, or regexes, have had great success both as pattern matching and as a library tool to develop small parsers. However, the absence of recursive patterns results in very limited parser applications. Since PEGs are more powerful than regular expressions, PEGs with regex-like captures could make it much easier to integrate a full-fledged parser into programs.

CPEGs are a formally developed extension of PEGs with regex-like captures. Two annotations (capture and fold-capture) allow a flexible construction of complex syntax trees. More importantly, a CPEG tree is a company with regular expression types that are a foundation of XML schema and tree automata. A regular expression type for a given CPEG is inferred syntactically. We present a formal definition of the type inference and we proved its soundness and uniqueness of types property.

Our attempt to type system for grammar is new. There are several interesting issues that remain unexplored. The future direction is that we will implement the CPEG based parser and investigate a practical aspect of CPEG. We intend to implement the CPEG based parser in F# and integrate our type inference into F# type provider [23].

REFERENCES

- [1] Michael D. Adams. 2013. Principled Parsing for Indentation-sensitive Languages: Revisiting Landin’s Offside Rule. In *Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL ’13)*. ACM, New York, NY, USA, 511–522. <https://doi.org/10.1145/2429069.2429129>
- [2] H. Comon, M. Dauchet, R. Gilleron, C. Löding, F. Jacquemard, D. Lugiez, S. Tison, and M. Tommasi. 2007. Tree Automata Techniques and Applications. Available on: <http://www.grappa.univ-lille3.fr/tata>. release October, 12th 2007.
- [3] Kathleen Fisher, Yitzhak Mandelbaum, and David Walker. 2006. The Next 700 Data Description Languages. *SIGPLAN Not.* 41, 1 (Jan. 2006), 2–15. <https://doi.org/10.1145/1111320.1111039>
- [4] Bryan Ford. 2004. Parsing Expression Grammars: A Recognition-based Syntactic Foundation. In *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL ’04)*. ACM, New York, NY, USA, 111–122. <https://doi.org/10.1145/964001.964011>
- [5] Alain Frisch. 2006. OCaml + XDuce. In *Proceedings of the Eleventh ACM SIGPLAN International Conference on Functional Programming (ICFP ’06)*. ACM, New York, NY, USA, 192–200. <https://doi.org/10.1145/1159803.1159829>
- [6] Robert Grimm. 2006. Better Extensibility Through Modular Syntax. In *Proceedings of the 27th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI ’06)*. ACM, New York, NY, USA, 38–51. <https://doi.org/10.1145/1133981.1133987>
- [7] Haruo Hosoya. 2010. *Foundations of XML Processing: The Tree-Automata Approach* (1st ed.). Cambridge University Press, New York, NY, USA.
- [8] Haruo Hosoya and Benjamin C. Pierce. 2003. XDuce: A Statically Typed XML Processing Language. *ACM Trans. Internet Technol.* 3, 2 (May 2003), 117–148. <https://doi.org/10.1145/767193.767195>
- [9] Haruo Hosoya, Jérôme Vouillon, and Benjamin C. Pierce. 2005. Regular Expression Types for XML. *ACM Trans. Program. Lang. Syst.* 27, 1 (Jan. 2005), 46–90. <https://doi.org/10.1145/1053468.1053470>
- [10] Roberto Ierusalimsky. 2009. A Text Pattern-matching Tool Based on Parsing Expression Grammars. *Softw. Pract. Exper.* 39, 3 (March 2009), 221–258. <https://doi.org/10.1002/spe.v39:3>
- [11] Stephen C Johnson. 1975. *Yacc: Yet another compiler-compiler*. Vol. 32. ell Laboratories Murray Hill, NJ.
- [12] Paul Klint, Tijs van der Storm, and Jurgen Vinju. 2010. On the Impact of DSL Tools on the Maintainability of Language Implementations. In *Proceedings of the Tenth Workshop on Language Descriptions, Tools and Applications (LDTA ’10)*. ACM, New York, NY, USA, Article 10, 9 pages. <https://doi.org/10.1145/1868281.1868291>
- [13] Kimio Kuramitsu. 2015. Nez open grammar repository. <https://github.com/nez-peg/nez-grammar>.
- [14] Kimio Kuramitsu. 2016. Nez: Practical Open Grammar Language. In *Proceedings of the 2016 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software (Onward! 2016)*. ACM, New York, NY, USA, 29–42. <https://doi.org/10.1145/2986012.2986019>
- [15] Daan Leijen and Erik Meijer. 2002. Parsec: Direct style monadic parser combinators for the real world. (2002).
- [16] Yitzhak Mandelbaum, Kathleen Fisher, David Walker, Mary Fernandez, and Artem Gleyzer. 2007. PADS/ML: A Functional Data Description Language. In *Proceedings of the 34th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL ’07)*. ACM, New York, NY, USA, 77–83. <https://doi.org/10.1145/1190216.1190231>
- [17] Adriaan Moors, Frank Piessens, and Martin Odersky. 2008. Parser Combinators in Scala.
- [18] Terence Parr and Kathleen Fisher. 2011. LL(*): The Foundation of the ANTLR Parser Generator. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI ’11)*. ACM, New York, NY, USA, 425–436. <https://doi.org/10.1145/1993498.1993548>
- [19] Terence Parr, Sam Harwell, and Kathleen Fisher. 2014. Adaptive LL(*) Parsing: The Power of Dynamic Analysis. In *Proceedings of the 2014 ACM International Conference on Object Oriented Programming Systems Languages & Applications (OOPSLA ’14)*. ACM, New York, NY, USA, 579–598. <https://doi.org/10.1145/2660193.2660202>
- [20] Tomas Petricek, Gustavo Guerra, and Don Syme. 2016. Types from Data: Making Structured Data First-class Citizens in F#. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI ’16)*. ACM, New York, NY, USA, 477–490. <https://doi.org/10.1145/2908080.2908115>
- [21] Benjamin C. Pierce. 2002. *Types and Programming Languages* (1st ed.). The MIT Press.
- [22] Martin Sulzmann and Kenny Zhuo Ming Lu. 2006. A Type-Safe Embedding of XDuce into ML. *Electron. Notes Theor. Comput. Sci.* 148, 2 (March 2006), 239–264. <https://doi.org/10.1016/j.entcs.2005.11.047>
- [23] Don Syme, Keith Battocchi, et al. 2012. Strongly-typed language support for internet-scale information sources. (2012).

- [24] Stephan Tolksdorf. 2013. FParsec—a parser combinator library for F#.