

1 Complete Parse Expressions: A Notation for Parse Trees

A parse tree can be unambiguously denoted using a symbolic form that I refer to as a *complete parse-expression*. Let G denote an arbitrary context-free grammar. Let A denote a nonterminal symbol in G and let α denote a string of symbols (terminals and/or nonterminals) such that:

$$A \Rightarrow^+ \alpha$$

Let α' denote a string that is derived from α by subscripting all nonterminals appearing in α . Under these assumptions we say that the term $A[\alpha']$ is a complete parse expression with respect to the grammar G .

1.1 An Example

Consider the following grammar containing the terminal symbols $\{c, d\}$ and the nonterminal symbols $\{A_list, B, D\}$:

A_list	$::=$	$A_list\ B \mid c$
B	$::=$	$B\ D \mid D$
D	$::=$	d

Among others, this grammar contains the following complete parse expressions:

$A_list[c]$
 $A_list[A_list_1\ B_1]$
 $A_list[A_list_1\ B_1\ B_1]$
 $A_list[A_list_1\ B_1\ B_2]$
 $A_list[cB_1\ B_2]$
 $A_list[cB_1\ d]$

 $B[D_1]$
 $B[B_1\ D_1]$
 $B[B_1\ D_1\ D_2]$
 $B[B_1\ D_1\ D_1]$
 $B[d]$

The following expressions are **NOT** complete parse expressions with respect to G :

Expression	Comment
$A_list[A_list_1]$	The length of the derivation must be at least 1.
$A_list[A_list_1\ B]$	The nonterminal B must be given a subscript.
$A_list[D_1]$	The derivation $A \xRightarrow{+} D$ is not possible in G .

1.2 An Example

In class conditions were given when the leading nonterminal of a complete parse expression can be dropped. Consider the following grammar containing the terminal symbols $\{c, d\}$ and the nonterminal symbols $\{A_list, B, D\}$:

A_list	$::=$	$A_list\ B \mid B$
B	$::=$	$c \mid D$
D	$::=$	d

Among others, this grammar contains the following parse expressions:

Parse-expression	Derived from
$\llbracket B_1 \rrbracket$	A
$\llbracket A_list_1\ B_1 \rrbracket$	A
$\llbracket A_list_1\ B_1\ B_1 \rrbracket$	A
$\llbracket A_list_1\ B_1\ B_2 \rrbracket$	A
$\llbracket c\ B_1\ B_2 \rrbracket$	A
$\llbracket c\ B_1\ d \rrbracket$	A
$\llbracket D_1 \rrbracket$	B
$\llbracket c \rrbracket$	B
$\llbracket d \rrbracket$	D

2 Matching

Parse expressions are useful because we will use them in equations to define the semantics of programming language constructs. In this framework, a program can be executed by using these equations to simplify a program until it can be simplified no further. This simplification is based on equational reasoning and requires matching as one of its computational steps. Therefore, to understand how a program can be simplified we must first understand matching in this context.

In the context of parse-expressions, a *match expression* is of the form:

$$\llbracket \alpha' \rrbracket \ll \llbracket \beta \rrbracket$$

In this case, α denotes a string that may contain terminals and subscripted nonterminals and β denotes a string that may only contain terminal symbols. It is also assumed that $\llbracket \alpha' \rrbracket$ and $\llbracket \beta \rrbracket$ are well-formed parse expressions.

In the context of matching, subscripted nonterminals occurring in parse expressions are treated as **variables** quantified over the set of all strings they can derive with respect to the given grammar. A match is said to **succeed** if the variables on the left-hand side of the match expression (i.e., the variables to the left of \ll) can be instantiated in such a manner so that the left and right sides of the match expression become syntactically equal.

The algorithm for solving a match expression of the form $\llbracket \alpha' \rrbracket \ll \llbracket \beta \rrbracket$ can be summarize as follows:

1. Check that $\llbracket \alpha' \rrbracket$ and $\llbracket \beta \rrbracket$ are well-formed.
2. Find values for variables (if possible) so that when the variables in α' are replaced with their values, the resulting parse-expression is $\llbracket \beta \rrbracket$.

3. Make sure that a value assigned to a variable can actually be derived from that variable. That is if A_1 is a variable and v_1 is a value, then $A \xRightarrow{+} v_1$ must be possible within the grammar.

2.1 Example

In this example, we use a standard BNF grammar fragment describing a subset of mathematical expressions. We assume that *num* and *ident* are terminal symbols and have the usual meaning.

$\begin{aligned} E &::= E + T \mid T \\ T &::= T * F \mid F \\ F &::= \text{num} \mid \text{ident} \mid (E) \end{aligned}$
--

Match Expression	Result of Match
$\llbracket E_1 + T_1 \rrbracket \ll \llbracket 1 + 2 \rrbracket$	true – $E_1 \xRightarrow{*} 1$ and $T_1 \xRightarrow{*} 2$
$\llbracket E_1 + T_1 \rrbracket \ll \llbracket 1 + 1 \rrbracket$	true – $E_1 \xRightarrow{*} 1$ and $T_1 \xRightarrow{*} 1$
$\llbracket T_1 * F_1 \rrbracket \ll \llbracket 1 * x \rrbracket$	true – $T_1 \xRightarrow{*} 1$ and $F_1 \xRightarrow{*} x$
$\llbracket T_1 * F_1 \rrbracket \ll \llbracket x * x \rrbracket$	true – $T_1 \xRightarrow{*} x$ and $F_1 \xRightarrow{*} x$
$\llbracket E_1 + E_2 \rrbracket \ll \llbracket x + 2 \rrbracket$	false – $\llbracket E_1 + E_2 \rrbracket$ is not well-formed
$\llbracket T_1 + T_2 \rrbracket \ll \llbracket 1 + y \rrbracket$	true – $T_1 \xRightarrow{*} 1$ and $T_2 \xRightarrow{*} y$
$\llbracket T_1 + T_1 \rrbracket \ll \llbracket 1 + 2 \rrbracket$	false – T_1 cannot be instantiated to be both 1 and 2 at the same time.
$\llbracket T_1 * T_2 \rrbracket \ll \llbracket x * 2 \rrbracket$	false – $\llbracket T_1 * T_2 \rrbracket$ is not well-formed
$\llbracket F_1 * F_2 \rrbracket \ll \llbracket 1 * y \rrbracket$	true – $F_1 \xRightarrow{*} 1$ and $F_2 \xRightarrow{*} y$
$\llbracket F_1 * F_1 \rrbracket \ll \llbracket 1 * 2 \rrbracket$	false – F_1 cannot be instantiated to be both 1 and 2 at the same time.
$\llbracket E_1 + T_1 * F_1 \rrbracket \ll \llbracket 1 + 2 * 3 \rrbracket$	true – $E_1 \xRightarrow{*} 1$ and $T_1 \xRightarrow{*} 2$ and $F_1 \xRightarrow{*} 3$
$\llbracket T_1 + T_1 * F_1 \rrbracket \ll \llbracket 2 + 2 * 3 \rrbracket$	true – $T_1 \xRightarrow{*} 2$ and $F_1 \xRightarrow{*} 3$
$\llbracket T_1 + T_1 * F_1 \rrbracket \ll \llbracket 1 + 2 * 3 \rrbracket$	false – T_1 cannot be instantiated to be both 1 and 2 at the same time.
$\llbracket E_1 + F_1 * F_2 \rrbracket \ll \llbracket 1 + 2 * 3 \rrbracket$	true – $E_1 \xRightarrow{*} 1$ and $F_1 \xRightarrow{*} 2$ and $F_2 \xRightarrow{*} 3$
$\llbracket F_1 + F_2 * F_3 \rrbracket \ll \llbracket 1 + 2 * 3 \rrbracket$	true – $F_1 \xRightarrow{*} 1$ and $F_2 \xRightarrow{*} 2$ and $F_3 \xRightarrow{*} 3$
$\llbracket F_1 + F_2 * T_3 \rrbracket \ll \llbracket 1 + 2 * 3 \rrbracket$	false – $\llbracket F_1 + F_2 * T_3 \rrbracket$ is not well-formed.