

Predictive Parsing

Monday, November 19, 2007
and Monday November 26, 2007
Reading: Appel 3.2

CS235 Languages and Automata

Department of Computer Science
Wellesley College

Overview

For some grammars, reading the first token (or first few tokens) of input is sufficient for determining which production to apply.

For such grammars, we can write a so-called **predictive parser**.

Today, we will see how to determine if a grammar is amenable to predictive parsing and how to build a predictive parser by hand.

We will also learn some techniques for handling grammars for which predictive parsers do not exist.

Running Example: SLiP--

As our running example, we will use SLiP--, a subset of Appel's straight-line programming language (SLiP).

The **abstract syntax** of SLiP-- is described by these SML datatypes:

```
datatype pgm = Pgm of Stm
and stm = Assign of string * exp
        | Print of exp
        | Seq of stm list
and exp = Id of string
        | Int of int
        | BinApp of exp * binop * exp
and binop = Add | Sub | Mul | Div
```

We will explore several versions of **concrete syntax** for SLiP--

Predictive Parsing

32/33-3

Our First Concrete Syntax for SLiP--

Productions for
Concrete Grammar

```
P → S EOF
S → ID(str) := E | print E | begin SL end
SL → % | S ; SL
E → ID(str) | INT(int) | ( E B E )
B → + | - | * | /
```

SML Data Types for
Abstract Grammar

```
datatype pgm = Pgm of Stm
and stm = Assign of string * exp
        | Print of exp
        | Seq of stm list
and exp = Id of string
        | Int of int
        | BinApp of exp * binop * exp
and binop = Add | Sub | Mul | Div
```

Predictive Parsing

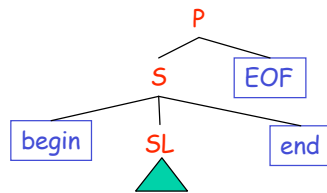
32/33-4

An Example SLiP-- Program

Chars: begin x := (3+4); print ((x-1)*(x+2)); end

Tokens: begin ID("x") := (INT(3) + INT(4)) ;
 print ((ID("x") - INT(1)) * (ID("x") +
 INT(2))) ; end EOF

Parse Tree:
 (see full tree
 on next slide)



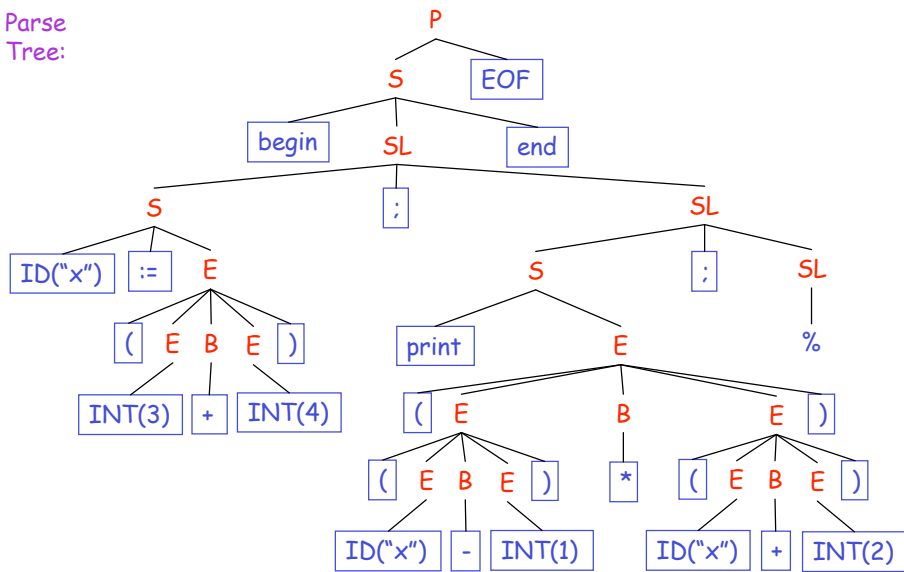
Abstract
 Syntax
 Tree (AST):

Pgm(Seq [Assign(BinApp(Int(3),Add,Int(4))),
 Print(BinApp(BinApp(Id("x"),Sub,Int(1)),
 Mul,
 BinApp(Id("x"),Add,Int(2)))]))

Predictive Parsing 32/33-5

An Example SLiP-- Program

Parse
 Tree:



Predictive Parsing 32/33-6

Predictive Parsing For SLiP--

```

P → S EOF
S → ID(str) := E | print E | begin SL end
SL → % | S ; SL
E → ID(str) | INT(int) | ( E B E )
B → + | - | * | /
    
```

Observe that:

- expressions (**E**) must begin with **ID(str)**, **INT(int)**, or **(**
- statements (**S**) must begin with **ID(str)**, **print**, or **begin**
- statement lists (**SL**) must begin with a statement (**S**) and so must begin with **ID(str)**, **prin**, or **begin** . They must end with **end** (a token that is not part of the **SL** tree but one immediately following it).
- programs (**P**) must begin with a statement (**S**) and so must begin with **ID(str)** , **print** , or **begin**

Predictive Parsing 32/33-7

Predictive Parsing Table for SLiP--

Can summarize observations on previous slide with a **predictive parsing table** of **variables** x **tokens** in which at most one production is valid per entry.

Empty slots in the table indicate parsing errors.

	ID(s)	INT(i)	(OP(b)	print	begin	end
P	P → S EOF				P → S EOF	P → S EOF	
S	S → ID(str) := E				S → print E	S → begin SL end	
SL	SL → S ; SL				SL → S ; SL	SL → S ; SL	SL → %
E	E → ID(str)	E → INT(num)	E → (E B E)				
B				B → OP(b)			

Predictive Parsing 32/33-8

Recursive Descent Parsing

From a predictive parsing table, it is possible to construct a **recursive descent parser** that parses tokens according to productions in the table.

Such a parser can eat (consume) or "peek" at the next token.

See ParserParens.sml for a recursive descent parser for SLiP-- written in SML.

Predictive Parsing 32/33-9

NULLABLE, FIRST, and FOLLOW

Parsing tables are constructed using the following notions:

Let t range over terminals, V and W range over variables, α range over terminals \cup variables, and γ range over sequences of terminals \cup variables.

- **NULLABLE(γ)** is true iff γ can derive the empty string (ϵ)
- **FIRST(γ)** is the set of terminals that can begin strings derived from γ .
- **FOLLOW(V)** is the set of terminals that can immediately follow V in some derivation.

Predictive Parsing 32/33-10

Computing NULLABLE For Variables

A variable V is NULLABLE iff

1. There is a production $V \rightarrow \epsilon$

OR

2. There is a production $V \rightarrow V_1 \dots V_n$
and each of V_1, \dots, V_n is NULLABLE

(Case 1 is really a special case of 2 with $n = 0$.)

In general, it is necessary to compute an **iterative fixed point** to determine nullability of a variable.

Example (from Appel 3.2)

$X \rightarrow a \mid Y$
$Y \rightarrow \epsilon \mid c$
$Z \rightarrow d \mid XYZ$

Predictive Parsing 32/33-11

Computing FIRST

$\text{FIRST}(t) = t$

$\text{FIRST}(V) = \cup \{ \text{FIRST}(\gamma) \mid V \rightarrow \gamma \text{ is a production for } V \}$

$\text{FIRST}(\alpha_0 \dots \alpha_i \dots \alpha_n) = \cup \{ \text{FIRST}(\alpha_i) \mid \alpha_0, \dots, \alpha_{i-1} \text{ are all nullable} \}$

Again, this is determined by an **iterative fixed point** computation

$X \rightarrow a \mid Y$
$Y \rightarrow \epsilon \mid c$
$Z \rightarrow d \mid XYZ$

Predictive Parsing 32/33-12

Computing FOLLOW

$$\text{FOLLOW}(V) =$$
$$\cup \{ \text{FIRST}(\alpha_j) \mid W \rightarrow \alpha_0 \dots \alpha_{i-1} V \alpha_{i+1} \dots \alpha_j \dots \alpha_n$$

is a production in the grammar
and $\alpha_{i+1}, \dots, \alpha_{j-1}$ are all nullable

$$\cup \cup \{ \text{FOLLOW}(W) \mid W \rightarrow \alpha_0 \dots \alpha_{i-1} V \alpha_{i+1} \dots \alpha_n$$

is a production in the grammar
and $\alpha_{i+1}, \dots, \alpha_n$ are all nullable

Again, this is determined by an **iterative fixed point** computation

$$\begin{array}{l} X \rightarrow a \mid Y \\ Y \rightarrow \% \mid c \\ Z \rightarrow d \mid XYZ \end{array}$$

Predictive Parsing 32/33-13

Constructing Predictive Parsing Tables

To construct a predictive parsing table,
do the following for each production $V \rightarrow \gamma$:

- For each t in $\text{FIRST}(\gamma)$, enter $V \rightarrow \gamma$ in row V , column t .
- If $\text{NULLABLE}(\gamma)$, for each t in $\text{FOLLOW}(V)$, enter $V \rightarrow \gamma$ in row V , column t

$$\begin{array}{l} P \rightarrow S \text{ EOF} \\ S \rightarrow \text{ID}(\text{str}) := E \mid \text{print } E \mid \text{begin } SL \text{ end} \\ SL \rightarrow \% \mid S ; SL \\ E \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid (E B E) \\ B \rightarrow + \mid - \mid * \mid / \end{array}$$

Predictive Parsing 32/33-14

More Practice

$S' \rightarrow S \text{ EOF}$
 $S \rightarrow T \mid OS1$
 $T \rightarrow \% \mid 10T$

	NULLABLE	FIRST	FOLLOW
S'			
S			
T			

	0	1	EOF
S'			
S			
T			

Parsing not predictive since some table slots now have multiple entries!

Predictive Parsing 32/33-15

Adding Extra Lookahead

$S' \rightarrow S \text{ EOF}$
 $S \rightarrow T \mid OS1$
 $T \rightarrow \% \mid 10T$

Sometimes predictivity can be re-established by adding extra **lookahead**:

	0	10	11	1 EOF	EOF
S'					
S					
T					

Predictive Parsing 32/33-16

LL(k) Grammars

An **LL(k)** grammar is one that has a predictive parsing table with **k** symbols of lookahead.

- The SLiP-- grammar is LL(1).
- The S'/S/T grammar is LL(2) but not LL(1).

In LL,

- the first L means the tokens are consumed left-to-right.
- the second L means that the parse tree is constructed in the manner of a leftmost derivation.

SLiP-- Expressions with Prefix Syntax

Suppose we change Slip-- expressions to use prefix syntax:

$$E \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid \text{B } E E$$

Parsing is still predictive:

	ID(s)	INT(i)	OP(b)	print	begin	end
E	$E \rightarrow \text{ID}(\text{str})$	$E \rightarrow \text{INT}(\text{num})$	$E \rightarrow \text{B } E E$			
B			$B \rightarrow \text{OP}(b)$			

Postfix Syntax for Expressions

Suppose we change Slip-- expressions to use postfix syntax:

$$E \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid E E B$$

Parsing is no longer predictive since some table slots now have multiple entries:

	ID(s)	INT(i)	OP(b)	print	begin	end
E	$E \rightarrow \text{ID}(\text{str})$ $E \rightarrow E E B$	$E \rightarrow \text{INT}(\text{num})$ $E \rightarrow E E B$				
B			$B \rightarrow \text{OP}(b)$			

Postfix expressions are fundamentally not predictive (not LL(k) for any k) there's nothing we can do to parse them predictively.

But we'll soon see we *can* parse them with a shift/reduce parser.

Predictive Parsing 32/33-19

Infix Syntax for Expressions

Suppose we change Slip-- expressions to use infix syntax without required parens (but with optional ones)

$$E \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid E B E \mid (E)$$

Parsing is no longer predictive:

	ID(s)	INT(i)	OP(b)	(print	begin	end
E	$E \rightarrow \text{ID}(\text{str})$ $E \rightarrow E B E$	$E \rightarrow \text{INT}(\text{num})$ $E \rightarrow E B E$		$E \rightarrow (E)$			
B			$B \rightarrow \text{OP}(b)$				

This is not surprising: this grammar is ambiguous, and *no* ambiguous grammar can be uniquely parsed with *any* deterministic parsing algorithm.

Predictive Parsing 32/33-20

Removing Ambiguity May not Help

Suppose we use an unambiguous infix grammar for arithmetic:

$E \rightarrow T \mid E + T \mid E - T$	<i>Expressions</i>
$T \rightarrow F \mid T * F \mid T / F$	<i>Terms</i>
$F \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid (E)$	<i>Factors</i>

Parsing is *still* not predictive due to **left recursion** in **E** and **T**:


	ID(s)	INT(i)	OP(b)	(print	begin	end
E	$E \rightarrow T$	$E \rightarrow T$		$E \rightarrow T$			
	$E \rightarrow E + T$	$E \rightarrow E + T$		$E \rightarrow E + T$			
	$E \rightarrow E - T$	$E \rightarrow E - T$		$E \rightarrow E - T$			
T	$T \rightarrow F$	$T \rightarrow F$		$T \rightarrow F$			
	$T \rightarrow T * F$	$T \rightarrow T * F$		$T \rightarrow T * F$			
	$T \rightarrow T / F$	$T \rightarrow T / F$		$T \rightarrow T / F$			
F	$F \rightarrow \text{ID}(\text{str})$	$F \rightarrow \text{INT}(\text{num})$		$F \rightarrow (E)$			

Predictive Parsing 32/33-21

Left Recursion Removal

Sometimes we can transform a grammar to remove left recursion (parse trees are transformed correspondingly).

$E \rightarrow T \mid E + T \mid E - T$
$T \rightarrow F \mid T * F \mid T / F$
$F \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid (E)$



$E \rightarrow T E'$
$E' \rightarrow \mid + T E' \mid - T E'$
$T \rightarrow F T'$
$T' \rightarrow \mid * F T' \mid / F T'$
$F \rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid (E)$

See Appel 3.2 for a general description of this transformation. You will see an example of this in PS5.

Predictive Parsing 32/33-22

The Transformed Grammar is Predictive!

$$\begin{aligned}
 E &\rightarrow T E' \\
 E' &\rightarrow \% \mid + T E' \mid - T E' \\
 T &\rightarrow F T' \\
 T' &\rightarrow \% \mid * F T' \mid / F T' \\
 F &\rightarrow \text{ID}(\text{str}) \mid \text{INT}(\text{int}) \mid (E)
 \end{aligned}$$

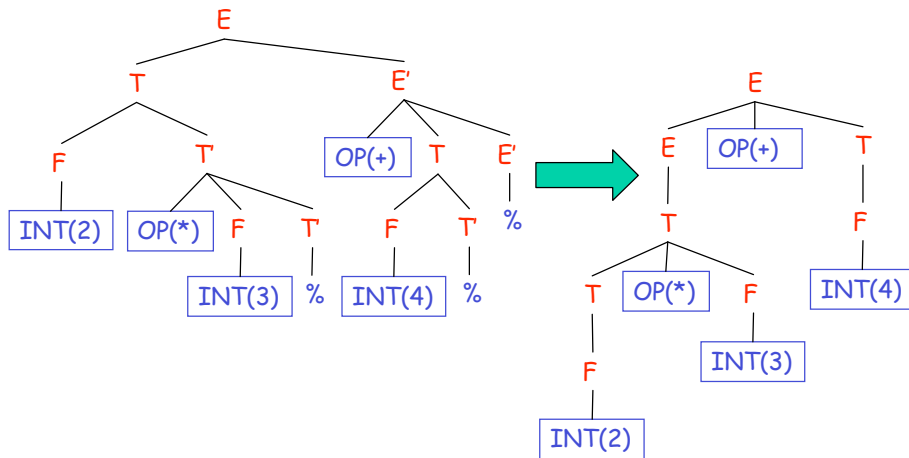
	ID(s)	INT(i)	+	*	()	:	EOF
E	$E \rightarrow T E'$	$E \rightarrow T E'$			$E \rightarrow T E'$			
E'			$E' \rightarrow + T E'$			$E' \rightarrow \%$	$E' \rightarrow \%$	$E' \rightarrow \%$
T	$T \rightarrow F T'$	$T \rightarrow F T'$			$T \rightarrow F T'$			
T'			$T' \rightarrow \%$	$T' \rightarrow * F T'$		$T' \rightarrow \%$	$T' \rightarrow \%$	$T' \rightarrow \%$
F	$F \rightarrow \text{ID}(\text{str})$	$F \rightarrow \text{INT}(\text{num})$			$F \rightarrow (E)$			

Predictive Parsing 32/33-23

Transforming Parse Trees

The parse tree from the transformed grammar can be transformed back to the untransformed grammar.

E.g. $2 * 3 + 4$



Predictive Parsing 32/33-24