

## ***LL(1) Parsing, pt. 1***

*CS 440: Programming Languages and Translators, Spring 2020*

### ***A. Parsing Using a Context-Free Grammar and Push-Down Automaton***

- Regular expressions, finite-state automata, and regular grammars work with the same set of languages.
- Similarly, context-free grammars and “push-down automata” both relate to context-free languages:
  - A context-free language is one generated by a context-free grammar. Context-free languages are also the ones recognizable using push-down automata.
  - A **Push-Down Automaton (PDA)** is a finite state machine augmented with a stack.
    - The state transition function maps a state, input symbol, and top-of-stack symbol to a new state and a stack operation (pop the stack, push a symbol onto the stack, or do nothing).
  - As with finite-state machines, there are deterministic and non-deterministic PDAs, but the non-deterministic PDAs can do things deterministic PDAs can't.
    - The set-of-states transformation used to map nondeterministic finite automata to equivalent deterministic finite automata doesn't work with PDAs.
    - (We can keep track of all the possible states we might be in, but we can't form the set of all possible stack configurations we might be in.)
  - For parsing, we much prefer the **deterministic PDAs (DPDAs)** because they don't require backtracking to execute the way a nondeterministic PDA can. A **deterministic CFL (DCFL)** is a language that can be parsed using a DPDA, and we use DCFLs to describe programming languages.
  - The simplest kind of DPDA has just one automaton state: Given the next input symbol and the top of the stack, gobble the input symbol and perform a stack operation.
  - It turns out that a language is LL(1) if it can be parsed by a DPDA with just one state.

### ***B. LL(1) Parsing, the Predict table, and First sets***

- We've been parsing LL(1) languages using recursive descent. The state of a recursive descent parser can be described using the current input symbol and a stack (which tracks the sequence of rule applications that got us where we're at). We have a nonterminal  $A$  that we're trying to expand and we look at the next input symbol to figure out which of the rules for  $A$  to apply ( $A \rightarrow \alpha_1$ ?  $A \rightarrow \alpha_2$ ? etc?)
  - As an example if we have rules  $A \rightarrow x \dots$  and  $A \rightarrow y \dots$  and the next input symbol is  $x$ , then we go with the  $A \rightarrow x \dots$  rule.
  - If no rules apply the next input symbol is  $x$  and there's no way for  $A$  to yield something that starts with  $x$  (i.e.,  $A \not\rightarrow^* x \dots$ ) then we have a parse error. If  $> 1$  rule applies, the language isn't LL(1).
- For a stack-based approach to LL(1) parsing, the “which rule do we apply?” decisions are stored in a **prediction table**:  $Predict(A, x)$  = the set of rules that apply if the current nonterminal to expand is  $A$  and the next input symbol is  $x$ . If  $Predict(A, x) = \emptyset$ , then we have a parse error. If  $Predict(A, x)$  contains  $> 1$  rule, then the language is not LL(1). (This is something we can check for when we build the prediction table.)

- The key to building the *Predict* table is the idea that for rule  $A \rightarrow \alpha$ , if from  $\alpha$  we can continue to something that begins with terminal symbol  $x$ , then it makes sense to try expanding  $A$  using  $A \rightarrow \alpha$  if the next input symbol is  $x$ .
- The set  $First(\alpha)$  holds the answers to the question of what terminal symbols we can find if we start with  $\alpha$  and look at all its possible derived yields.
- **Definition:**  $First(\alpha) = \{x \in T \mid \alpha \rightarrow^* x \beta \text{ for some } \beta \in (V \mid T)^*\}$ . I.e., what sentential forms  $x\dots$  can we get to if we start from  $\alpha$ ?
  - For all  $\beta$ ,  $\{x\} = First(x \beta)$ . I.e., if  $\alpha$  begins with a terminal symbol  $x$ , then certainly  $x \in First(\alpha)$ .
  - For all rules  $A \rightarrow \alpha$ ,  $First(A \beta) \supseteq First(\alpha \beta)$ . I.e., if we want to know the *First* set for  $A$ , we certainly need the *First* sets for all rhs's of rules for  $A$ .
  - There's a complicating factor: If  $A \rightarrow^* x\dots$ , then for all  $\beta$ ,  $A \beta \rightarrow^* x\dots \beta$ , so we don't need to know about  $First(\beta)$ , right? But if  $A \rightarrow^* \epsilon$ , then  $A \beta \rightarrow^* \epsilon \beta$ , so  $First(A \beta) \supseteq First(\beta)$ , so we do care about  $First(\beta)$ . So,
    - If  $A \not\rightarrow^* \epsilon$  then for all  $\beta$ ,  $First(A \beta) = First(A)$ .
    - If  $A \rightarrow^* \epsilon$  then for all  $\beta$ ,  $First(A \beta) \supseteq First(\beta)$ .
- **Definition 1:** For a nonterminal  $A$  and terminal  $x$ ,
  - $Predict(A, x) = \{\text{rule } A \rightarrow \alpha \mid (\text{For some } \beta), S \rightarrow^* \dots A \beta \text{ and } x \in First(A \beta)\}$
- If  $A \rightarrow^* \epsilon$  then in theory, we need to know  $First(\beta)$  for all possible  $\beta$  (of which there are an infinite number). But actually, we only need to worry about the  $A \beta$  we can get to from the start symbol. The situation is actually even better than that, but let's first look at the basic stack-based parsing algorithm and then go back to the details of calculating the *Predict* table using  $First(A)$  and a helper function  $Follow(A)$ .

### C. Stack-Based LL(1) Parsing using Predict

- First, so that we can distinguish between the first use of the start symbol and any subsequent use, we'll **modify the grammar** a bit:
  - 1: Add a fresh (unused) terminal symbol  $\$$  used to signal end-of-input.
  - 2: Add a new start symbol  $S'$  and a rule  $S' \rightarrow S \$$  where  $S$  is the nominal start symbol.

#### Basic algorithm:

```

Concatenate $ to end of input string and initialize the stack: push  $S'$  onto it.
while top of stack  $\neq \$$  {
    Let  $x$  be next character of remaining input.
    if  $x = \$$ , parse error (input ended early)
    Pop top of stack
    if it's a terminal  $y$ 
        if  $x = y$ , remove  $x$  from the input and continue loop else parse error
    else (the top of the stack was a nonterminal  $A$ )
        if  $Predict(A, x)$  says error (it's  $\emptyset$ ) then parse error (unexpected terminal  $x$ )
        else check  $Predict(A, x)$  for an entry  $A \rightarrow \alpha$  and push  $\alpha$  onto the stack
end

```

**if** (loop has ended, top of stack is \$ and) next the input character is \$  
    We parsed successfully!  
**else** parse error (leftover input)

***D. Example of Using Prediction table for LL(1) Parse***

- See Example 1.

**Example 1: Grammar of Balanced Parens plus variables  $x$  and  $y$** **Terminals** = ( )  $x$   $y$  \$**Rules**

Rule Nbr	Rule
0	$S' \rightarrow S \$$
1	$S \rightarrow P S$
2	$S \rightarrow x$
3	$P \rightarrow \backslash ( S \backslash )$
4	$P \rightarrow y$

**Predict( $A, s$ ) = Rule number (blank entry = error)**

Nonterminal	(	)	$x$	$y$
$S'$	0: $S' \rightarrow S \$$		2: $S \rightarrow x$	1: $S \rightarrow P S$
$S$	1: $S \rightarrow P S$		2: $S \rightarrow x$	1: $S \rightarrow P S$
$P$	3: $P \rightarrow \backslash ( S \backslash )$			4: $P \rightarrow y$

**Leftmost Derivation of (  $y x$  ) (  $x$  )  $x$  \$** $S' \$$  $\rightarrow S \$$  $\rightarrow P S \$$  $\rightarrow ( S ) S \$$  $\rightarrow ( P S ) S \$$  $\rightarrow ( y S ) S \$$  $\rightarrow ( y x ) P S \$$  $\rightarrow ( y x ) ( S ) S \$$  $\rightarrow ( y x ) ( x ) S \$$  $\rightarrow ( y x ) ( x ) x \$$

Trace of Parse of  $( y x ) ( x ) x \$$ 

Stack (Top to Left)	Rule nbr or = (terminals)	Input
$S'$	0	$( y x ) ( x ) x \$$
$S \$$	1	$( y x ) ( x ) x \$$
$P S \$$	3	$( y x ) ( x ) x \$$
$( S ) S \$$	=	$( y x ) ( x ) x \$$
$S ) S \$$	1	$y x ) ( x ) x \$$
$P S ) S \$$	4	$y x ) ( x ) x \$$
$y S ) S \$$	=	$y x ) ( x ) x \$$
$S ) S \$$	2	$x ) ( x ) x \$$
$x ) S \$$	=	$x ) ( x ) x \$$
$) S \$$	=	$) ( x ) x \$$
$S \$$	1	$( x ) x \$$
$P S \$$	3	$( x ) x \$$
$( S ) S \$$	=	$( x ) x \$$
$S ) S \$$	2	$x ) x \$$
$x ) S \$$	=	$x ) x \$$
$) S \$$	=	$) x \$$
$S \$$	2	$x \$$
$x \$$	=	$x \$$
$\$$	=	$\$$

NT	(	)	x	y
$S'$	0: $S' \rightarrow S \$$		2: $S \rightarrow x$	1: $S \rightarrow P S$
$S$	1: $S \rightarrow P S$		2: $S \rightarrow x$	1: $S \rightarrow P S$
$P$	3: $P \rightarrow \backslash ( S \backslash )$			4: $P \rightarrow y$

(End of example 1)

### *Activity Questions for Lecture 13*

1. Calculate the *First*(...) sets and the *Predict* table for a language without  $\epsilon$  rules, specifically,  
$$S \rightarrow A \$$$
$$A \rightarrow a B$$
$$B \rightarrow C D e$$
$$C \rightarrow c C \mid D$$
$$D \rightarrow d D$$
2. Use the *Predict* table from question 1 and calculate a trace of the parse of `a d a d c e e $` using the parsing algorithm.
3. Calculate the *First*(...) sets and the *Predict* table for the following language  
$$S' \rightarrow S \$$$
$$S \rightarrow a S b S$$
$$S \rightarrow c$$
4. Use the *Predict* table from question 2 and calculate a trace of the parse of `a a c b a c b c b c $`.

## Solutions to Activity Questions for Lecture 13

1. (Calculate First sets and LL prediction table for grammar with no  $\epsilon$  rules.)

Here's the reasoning behind the First set contents with each rule and its implications.

$S \rightarrow A \$$	$\text{First}(S) \supseteq \text{First}(A)$
$A \rightarrow a B$	$a \in \text{First}(A)$
$B \rightarrow C D e$	$\text{First}(B) \supseteq \text{First}(C)$
$C \rightarrow c \mid D$	$c \in \text{First}(C); \text{First}(C) \supseteq \text{First}(D)$
$D \rightarrow d$	$d \in \text{First}(C)$
$D \rightarrow A$	$\text{First}(D) \supseteq \text{First}(A)$

So we get

$$\text{First}(S) = \{a\}$$

$$\text{First}(A) = \{a\}$$

$$\text{First}(B) = \{c, d\}$$

$$\text{First}(C) = \{c, d\}$$

$$\text{First}(D) = \{a, d\}$$

**Prediction Table**

Nonterminal	a	c	d	e
$S$	0: $S \rightarrow A \$$			
$A$	1: $A \rightarrow a B$			
$B$		2: $B \rightarrow C D e$	2: $B \rightarrow C D e$	
$C$		3: $C \rightarrow c$	4: $C \rightarrow D$	5: $C \rightarrow D$
$D$	6: $D \rightarrow A$		5: $D \rightarrow d$	

2. (Trace parse)

**Trace of Parse of a d a d c e e \$**

Stack (Top to Left)	Rule # / = term?	Input
$S$	0	a d a d c e e \$
$A \$$	1	a d a d c e e \$
$a B \$$	=	a d a d c e e \$
$B \$$	2	d a d c e e \$
$C D e \$$	4	d a d c e e \$
$D D e \$$	5	d a d c e e \$
$d D e \$$	=	d a d c e e \$
$D e \$$	6	a d c e e \$
$A e \$$	1	a d c e e \$
$a B e \$$	=	a d c e e \$
$B e \$$	2	d c e e \$
$C D e e \$$	3	d c e e \$
$c D e e \$$	=	d c e e \$
$D e e \$$	5	d e e \$
$d e e \$$	=	d e e \$
$e e \$$	=	e e \$
$e \$$	=	e \$
$\$$	success!	\$

Rule nbr	Rule
0	$S \rightarrow A \$$
1	$A \rightarrow a B$
2	$B \rightarrow C D e$
3	$C \rightarrow c$
4	$C \rightarrow D$
5	$D \rightarrow d$
6	$D \rightarrow A$

3. (First sets and Predict table for LL(1) grammar without  $\epsilon$  rules)

Rule 0:  $S' \rightarrow S \$$  Implies  $First(S) \subseteq First(S')$

Rule 1:  $S \rightarrow a S b S$   $a \in First(S)$

Rule 2:  $S \rightarrow c$   $c \in First(S)$

So  $First(S) = First(S') = \{a, c\}$

**Predict table**

Nonterminal	a	b	c	\$
$S'$	0: $S' \rightarrow S \$$		0: $S' \rightarrow S \$$	
$S$	1: $S \rightarrow a S b S$		2: $S \rightarrow c$	

4. (Calculate trace of a parse)

**Trace of Parse of a a c b a c b c b c \$**

Stack (Top to Left)	Rule # / = term?	Input
$S'$	0	a a c b a c b c b c \$
$S \$$	1	a a c b a c b c b c \$
a S b S \$	=	a a c b a c b c b c \$
S b S \$	1	a c b a c b c b c \$
a S b S b S \$	=	a c b a c b c b c \$
S b S b S \$	2	c b a c b c b c \$
c b S b S \$	=	c b a c b c b c \$
b S b S \$	=	b a c b c b c \$
S b S \$	1	a c b c b c \$
a S b S b S \$	=	a c b c b c \$
S b S b S \$	2	c b c b c \$
c b S b S \$	=	c b c b c \$
b S b S \$	=	b c b c \$
S b S \$	2	c b c \$
c b S \$	=	c b c \$
b S \$	=	b c \$
S \$		c \$
c \$	=	c \$
\$	success!	\$