

Parsing Algorithms

**CS 4447/CS 9545 -- Stephen Watt
University of Western Ontario**

The Big Picture

- Develop parsers based on grammars
- Figure out properties of the grammars
- Make tables that drive parsing engines

- Two essential ideas:
Derivations and *FIRST/FOLLOW sets*

Outline

- Grammars, parse trees and **derivations**.
- Recursive descent parsing
- Operator precedence parsing
- Predictive parsing
 - *FIRST* and *FOLLOW*
 - LL(1) parsing tables. LL(k) parsing.
- Left-most and right-most derivations
- **Shift-reduce** parsing
 - LR parsing automaton. LR(k) parsing.
 - LALR(k) parsing.

Example Grammar G1

- We have seen grammars already.

Here is an example.

[from *Modern Compiler Implementation in Java*, by Andrew W. Appel]

1. $S \rightarrow S \text{ “;” } S$
2. $S \rightarrow \text{id “:=” } E$
3. $S \rightarrow \text{“print” “(” } L \text{ “)”}$
4. $E \rightarrow \text{id}$
5. $E \rightarrow \text{num}$
6. $E \rightarrow E \text{ “+” } E$
7. $E \rightarrow \text{“(” } S \text{ “,” } E \text{ “)”}$
8. $L \rightarrow E$
9. $L \rightarrow L \text{ “,” } E$

Parse Trees

- A *parse tree* for a given grammar and input is a tree where each node corresponds to a grammar rule, and the leaves correspond to the input.

Example Parse Tree

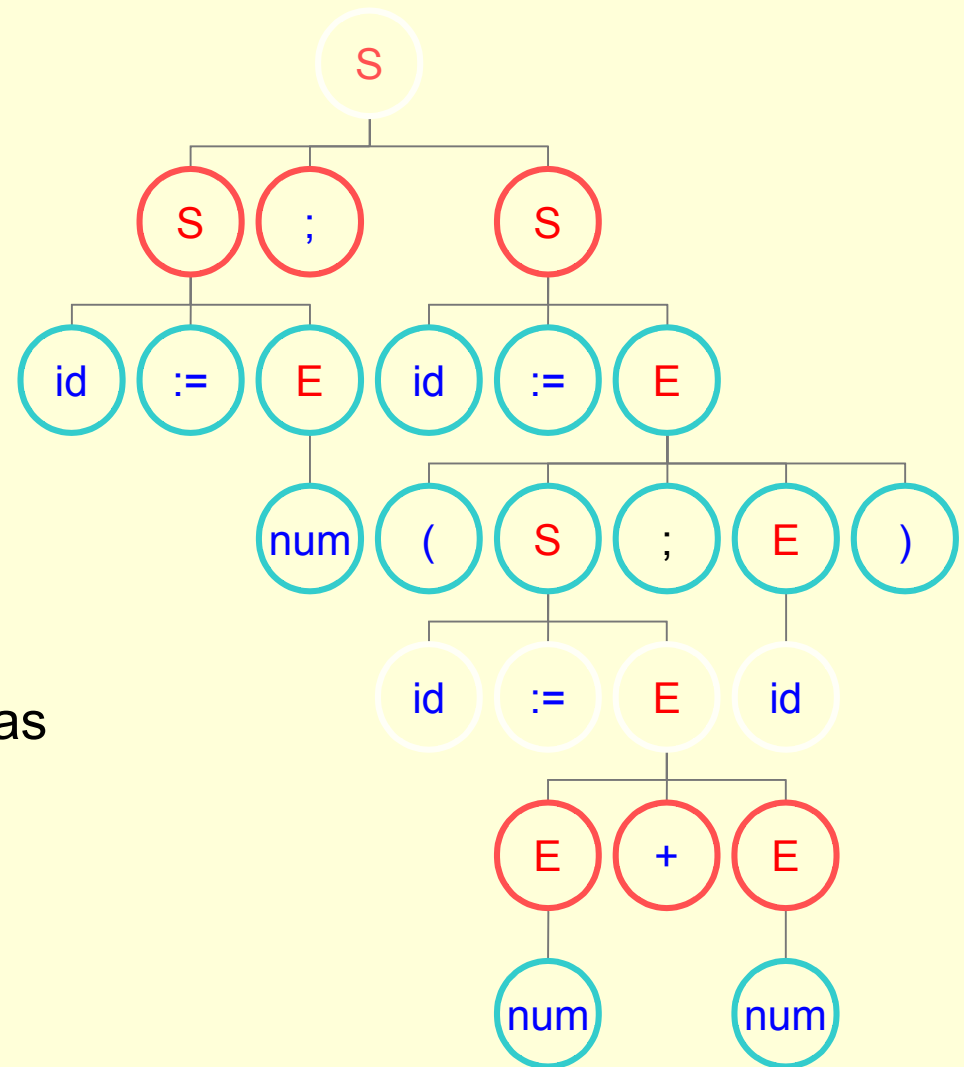
- Consider the following input:

```
a := 7;  
b := c + (d := 5 + 6, d)
```

- This gives the *token* sequence:

```
id := num ; id := id +  
( id := num + num , id )
```

- Using example grammar G1, this has a parse tree shown on the right.



Derivations

- “Parsing” figures out a parse tree for an given input
- A “derivation” gives a rationale for a parse.

- Begin with the grammar’s start symbol and repeatedly replace non-terminals until only terminals remain.

Example Derivation

This derivation
justifies the parse
tree we showed.

S

S ; S

S ; id := E

id := E ; id := E

id := num ; id := E

id := num ; id := E + E

id := num ; id := E + (S, E)

id := num ; id := id + (S, E)

id := num ; id := id + (id := E, E)

id := num ; id := id + (id := E + E, E)

id := num ; id := id + (id := E + E, id)

id := num ; id := id + (id := num + E, id)

id := num ; id := id + (id := num + num, id)

Derivations and Parse Trees

- A node in a parse tree corresponds to the use of a rule in a derivation.
- A grammar is *ambiguous* if it can derive some sentence with two *different* parse trees.

E.g. $a + b * d$ can be derived two ways using the rules

$E \rightarrow \text{id}$ $E \rightarrow E \text{ "+" } E$ $E \rightarrow E \text{ "*" } E$

- Even for an unambiguous grammar, there is a choice of which non-terminal to replace in forming a derivation.

Two choices are

- Replace the leftmost non-terminal
- Replace the rightmost non-terminal

Recursive Descent Parsing

- Example for recursive descent parsing:

1. $S \rightarrow E$

2. $E \rightarrow T \text{ "+" } E$

3. $E \rightarrow T$

4. $T \rightarrow F \text{ "*" } T$

5. $T \rightarrow F$

6. $F \rightarrow P \text{ "^" } F$

7. $F \rightarrow P$

8. $P \rightarrow \text{id}$

9. $P \rightarrow \text{num}$

10. $P \rightarrow \text{"(" } E \text{ ")"}$

- Introduce one function for each non-terminal.

Recursive Descent Parsing (cont'd)

```
PT* S() { return E(); }
```

```
PT* E() { PT *pt = T();  
    if (peek("+")) { consume("+"); pt = mkPT(pt,E()); }  
    return pt; }
```

```
PT* T() { PT *pt = F();  
    if (peek("*")) { consume("*"); pt = mkPT(pt,F()); }  
    return pt; }
```

```
PT* F() { PT *pt = P();  
    if (peek("^")) { consume("^"); pt = mkPT(pt,P()); }  
    return pt; }
```

```
PT* P() { PT *pt;  
    if (peekDigit()) return new PT(Num());  
    if (peekLetter()) return new PT(Id());  
    consume("("); pt = E(); consume(")");  
    return pt; }
```

Recursive Descent Parsing -- Problems

- A slightly different grammar (G2) gives problems, though:

1. $S \rightarrow E$

2. $E \rightarrow E \text{ "+" } T$

3. $E \rightarrow T$

4. $T \rightarrow T \text{ "*" } F$

5. $T \rightarrow F$

6. $F \rightarrow P \text{ "^" } F$

7. $F \rightarrow P$

8. $P \rightarrow \text{id}$

9. $P \rightarrow \text{num}$

10. $P \rightarrow \text{"(" } E \text{ ")"}$

- This causes problems, e.g.:

- Do not know whether to use rule 2 or rule 3 parsing an E.
- Rule 2 gives an infinite recursion.

- We want to be able to *predict* which rule (which recursive function) to use, based on looking at the current input token.

Operator Precedence Parsing

- Each operator has left- and right- precedence. E.g.

100+101 200×201 301^300

- Group sub-expressions by binding highest numbers first.

A+B × C × D ^ E ^ F

A 100 +101 B 200×201 C 200×201 D 301^300 E 301^300 F

A 100 +101 B 200×201 C 200×201 D 301^300 (E 301^300 F)

A 100 +101 B 200×201 C 200×201 (D 301^300 (E 301^300 F))

A 100 +101 (B 200×201 C) 200×201 (D 301^300 (E 301^300 F))

A 100 +101 ((B 200×201 C) 200×201 (D 301^300 (E 301^300 F)))

A+((B × C) × (D ^ (E ^ F)))

- Works fine for infix expressions but not well for general CFL.

Predictive Parsing – FIRST sets

- We introduce the notion of “FIRST” sets that will be useful in predictive parsing.
- If α is a string of terminals and non-terminals, then $\text{FIRST}(\alpha)$ is the set of all terminals that may be the first symbol in a string derived from α .
- Eg1: For example grammar G1,

$$\text{FIRST}(S) = \{ \text{id}, \text{"print"} \}$$

- Eg2: For example grammar G2,

$$\text{FIRST}(T \text{"*"} F) = \{ \text{id}, \text{num}, \text{"("} \}$$

Predictive Parsing -- good vs bad grammars

- If two productions for the same LHS have RHS with intersecting FIRST sets, then the grammar cannot be parsed using predictive parsing.

E.g. with $E \rightarrow E \text{ "+" } T$ and $E \rightarrow T$
 $\text{FIRST}(E \text{ "+" } T) = \text{FIRST}(T) = \{ \text{id, num, "("} \}$

- To use predictive parsing, we need to formulate a different grammar for the same language.
- One technique is to eliminate left recursion:

E.g. replace $E \rightarrow E \text{ "+" } T$ and $E \rightarrow T$
with $E \rightarrow T E'$ $E' \rightarrow \text{ "+" } T E'$ $E' \rightarrow \epsilon$

The “nullable” property

- We say a non-terminal is “nullable” if it can derive the empty string.
- In the previous example E' is nullable.

FOLLOW sets

- The “FOLLOW” set for a non-terminal X is the set of terminals that can immediately follow X .
- The terminal t is in $\text{FOLLOW}(X)$ if there is a derivation containing Xt .
- This can occur if there is a derivation containing $X Y Z t$, if Y and Z are nullable.

Algorithm for FIRST, FOLLOW, nullable

for each symbol X

$\text{FIRST}[X] := \{\}, \text{FOLLOW}[X] := \{\}, \text{nullable}[X] := \text{false}$

for each terminal symbol t

$\text{FIRST}[t] := \{t\}$

repeat

for each production $X \rightarrow Y_1 Y_2 \dots Y_k$,

if all Y_i are nullable **then**

$\text{nullable}[X] := \text{true}$

if $Y_1..Y_{i-1}$ are nullable **then**

$\text{FIRST}[X] := \text{FIRST}[X] \cup \text{FIRST}[Y_i]$

if $Y_{i+1}..Y_k$ are all nullable **then**

$\text{FOLLOW}[Y_i] := \text{FOLLOW}[Y_i] \cup \text{FOLLOW}[X]$

if $Y_{i+1}..Y_{j-1}$ are all nullable **then**

$\text{FOLLOW}[Y_i] := \text{FOLLOW}[Y_i] \cup \text{FIRST}[Y_j]$

until FIRST, FOLLOW, nullable do not change

Example FIRST, FOLLOW, nullable

Example Grammar G3.

$Z \rightarrow d$ $Y \rightarrow \epsilon$ $X \rightarrow Y$

$Z \rightarrow X Y Z$ $Y \rightarrow c$ $X \rightarrow a$

	nullable	FIRST	FOLLOW
X	false		
Y	false		
Z	false		

	nullable	FIRST	FOLLOW
X	false	a	c d
Y	true	c	d
Z	false	d	

	nullable	FIRST	FOLLOW
X	true	a c	a c d
Y	true	c	a c d
Z	false	a c d	

Predictive Parsing Tables

- Rows: Non-terminals
- Columns: Terminals
- Entries: Productions

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

Enter production $X \rightarrow \alpha$ in row X , column t for each t in $\text{FIRST}(\alpha)$.

If α is nullable, enter the productions in row X , column t for each t in $\text{FOLLOW}(X)$.

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

	nullable	FIRST	FOLLOW
X	true	a c	a c d
Y	true	c	a c d
Z	false	a c d	

Example of Predictive Parsing

Initial grammar

$S \rightarrow E$

$E \rightarrow E \text{ "+" } T \quad E \rightarrow T$

$T \rightarrow T \text{ "*" } F \quad T \rightarrow F$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow \text{"(" } E \text{ ")"}$

Modified grammar

$S \rightarrow E \$$

$E \rightarrow T E' \quad E' \rightarrow \quad E' \rightarrow \text{"+" } T E'$

$T \rightarrow F T' \quad T' \rightarrow \quad T' \rightarrow \text{"*" } F T'$

$F \rightarrow \text{id}$

$F \rightarrow \text{num}$

$F \rightarrow \text{"(" } E \text{ ")"}$

	Nullable	FIRST	FOLLOW
S	False	(id num	
E	False	(id num) \$
E'	True	+) \$
T	False	(id num) + \$
T'	True	*) + \$
F	False	(id num) * + \$

Example of Predictive Parsing (contd)

$S \rightarrow E \$$
 $E \rightarrow T E'$ $E' \rightarrow$ $E' \rightarrow "+" T E'$
 $T \rightarrow F T'$ $T' \rightarrow$ $T' \rightarrow "*" F T'$
 $F \rightarrow \text{id}$
 $F \rightarrow \text{num}$
 $F \rightarrow "(" E ")"$

	Nullable	FIRST	FOLLOW
S	False	(id num	
E	False	(id num) \$
E'	True	+) \$
T	False	(id num) + \$
T'	True	*) + \$
F	False	(id num) * + \$

	+	*	id	num	()	\$
S			$S \rightarrow E \$$	$S \rightarrow E \$$	$S \rightarrow E \$$		
E			$E \rightarrow T E'$	$E \rightarrow T E'$	$E \rightarrow T E'$		
E'	$E' \rightarrow "+" T E'$					$E' \rightarrow$	$E' \rightarrow$
T			$T \rightarrow F T'$	$T \rightarrow F T'$	$T \rightarrow F T'$		
T'	$T' \rightarrow$	$T' \rightarrow "*" F T'$				$T' \rightarrow$	$T' \rightarrow$
F			$F \rightarrow \text{id}$	$F \rightarrow \text{num}$	$F \rightarrow "(" E ")"$		

LL(k) Grammars

- The predictive parser we built makes use of one look ahead token.
 - We say the grammar is LL(1).
 - LL stands for “Left to right parse, Leftmost derivation”
- If k look ahead tokens are needed, then we say the grammar is LL(k).
 - For $k > 1$, the columns are the possible sequences of k tokens, and the tables become large.
- There is a better way...

LR Parsing

- LL parsing always uses a grammar rule for the *left-most* non-terminal.
- If we aren't so eager, we can apply grammar rules to other non-terminals
- This allows us to decide about the “hard” non-terminals later.
- We keep a stack of unfinished work.
- Using the *right-most* derivation leads to LR parsing.

LR Parsing

- Parser state consists of a *stack* and *input*.
- First k tokens of the unused input is the “*lookahead*”
- Based on what is on the top of the stack and the lookahead, the parser decides whether to
 - Shift =
 1. consume the first input token
 2. push it to the top of the stack
 - Reduce =
 1. choose a grammar rule $X \rightarrow A B C$
 2. pop C, B, A from the stack
 3. push X onto the stack