

Bottom-Up Parsing

Lecture 8
(From slides by G. Necula & R. Bodik)

2/12/09

Prof. Hilfinger CS164 Lecture 8

1

Administrivia

- **Test I** during class on 10 March.
- Notes updated (at last)

2/12/09

Prof. Hilfinger CS164 Lecture 8

2

Bottom-Up Parsing

- We've been looking at general context-free parsing.
- It comes at a price, measured in overheads, so in practice, we design programming languages to be parsed by less general but faster means, like top-down recursive descent.
- Deterministic bottom-up parsing is more general than top-down parsing, and just as efficient.
- Most common form is *LR parsing*
 - L means that tokens are read left to right
 - R means that it constructs a rightmost derivation

2/12/09

Prof. Hilfinger CS164 Lecture 8

3

An Introductory Example

- LR parsers don't need left-factored grammars and can also handle left-recursive grammars

- Consider the following grammar:

$$E \rightarrow E + (E) \mid \text{int}$$

- Why is this not LL(1)?

- Consider the string: `int + (int) + (int)`

2/12/09

Prof. Hilfinger CS164 Lecture 8

4

The Idea

- LR parsing *reduces* a string to the start symbol by inverting productions:

`sent` ← input string of terminals

while `sent` ≠ `S`:

- Identify first β in `sent` such that $A \rightarrow \beta$ is a production and $S \rightarrow^* \alpha A \gamma \rightarrow \alpha \beta \gamma = \text{sent}$
- Replace β by A in `sent` (so $\alpha A \gamma$ becomes new `sent`)
- Such $\alpha \beta$'s are called *handles*

2/12/09

Prof. Hilfinger CS164 Lecture 8

5

A Bottom-up Parse in Detail (1)

`int + (int) + (int)`

`int + (int) + (int)`

2/12/09

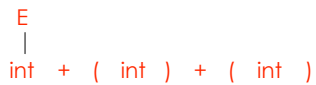
Prof. Hilfinger CS164 Lecture 8

6

A Bottom-up Parse in Detail (2)

int + (int) + (int)
E + (int) + (int)

(handles in red)



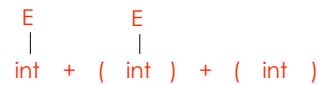
2/12/09

Prof. Hilfinger CS164 Lecture 8

7

A Bottom-up Parse in Detail (3)

int + (int) + (int)
E + (int) + (int)
E + (E) + (int)



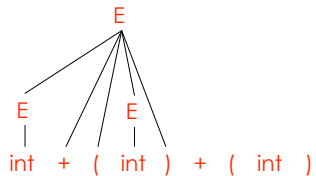
2/12/09

Prof. Hilfinger CS164 Lecture 8

8

A Bottom-up Parse in Detail (4)

int + (int) + (int)
E + (int) + (int)
E + (E) + (int)
E + (int)



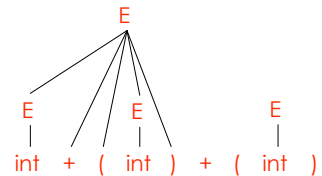
2/12/09

Prof. Hilfinger CS164 Lecture 8

9

A Bottom-up Parse in Detail (5)

int + (int) + (int)
E + (int) + (int)
E + (E) + (int)
E + (int)
E + (E)



2/12/09

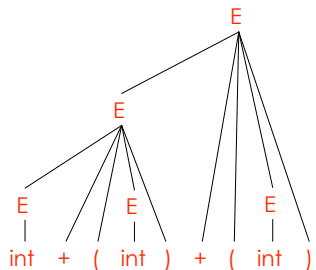
Prof. Hilfinger CS164 Lecture 8

10

A Bottom-up Parse in Detail (6)

↑
int + (int) + (int)
E + (int) + (int)
E + (E) + (int)
E + (int)
E + (E)
E

A reverse rightmost derivation



2/12/09

Prof. Hilfinger CS164 Lecture 8

11

Where Do Reductions Happen

Because an LR parser produces a reverse rightmost derivation:

- If $\alpha\beta\gamma$ is step of a bottom-up parse with handle $\alpha\beta$
- And the next reduction is by $A \rightarrow \beta$
- Then γ is a string of terminals!

... Because $\alpha A \gamma \rightarrow \alpha \beta \gamma$ is a step in a right-most derivation

Intuition: We make decisions about what reduction to use *after* seeing all symbols in handle, rather than *before* (as for LL(1))

2/12/09

Prof. Hilfinger CS164 Lecture 8

12

Notation

- Idea: Split the string into two substrings
 - Right substring (a string of terminals) is as yet unexamined by parser
 - Left substring has terminals and non-terminals
- The dividing point is marked by a $|$
 - The $|$ is not part of the string
 - Marks end of next potential handle
- Initially, all input is unexamined: $|x_1x_2 \dots x_n$

2/12/09

Prof. Hilfinger CS164 Lecture 8

13

Shift-Reduce Parsing

- Bottom-up parsing uses only two kinds of actions:
 - Shift:* Move $|$ one place to the right, shifting a terminal to the left string
 $E + (| int) \Rightarrow E + (int |)$
 - Reduce:* Apply an inverse production at the handle.
If $E \rightarrow E + (E)$ is a production, then
 $E + (E + (E) |) \Rightarrow E + (E |)$

2/12/09

Prof. Hilfinger CS164 Lecture 8

14

Shift-Reduce Example

$| int + (int) + (int) \$$ shift

$int + (int) + (int)$
↑

Shift-Reduce Example

$| int + (int) + (int) \$$ shift
 $int | + (int) + (int) \$$ red. $E \rightarrow int$

$int + (int) + (int)$
↑

Shift-Reduce Example

$| int + (int) + (int) \$$ shift
 $int | + (int) + (int) \$$ red. $E \rightarrow int$
 $E | + (int) + (int) \$$ shift 3 times

E
/
 $int + (int) + (int)$
↑

Shift-Reduce Example

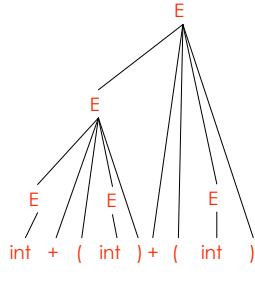
$| int + (int) + (int) \$$ shift
 $int | + (int) + (int) \$$ red. $E \rightarrow int$
 $E | + (int) + (int) \$$ shift 3 times
 $E + (int |) + (int) \$$ red. $E \rightarrow int$

E
/
 $int + (int) + (int)$
↑

Shift-Reduce Example

```

I int + (int) + (int)$ shift
int I + (int) + (int)$ red. E → int
E I + (int) + (int)$ shift 3 times
E + (int I) + (int)$ red. E → int
E + (E I) + (int)$ shift
E + (E) I + (int)$ red. E → E + (E)
E I + (int)$ shift 3 times
E + (int I)$ red. E → int
E + (E I)$ shift
E + (E) I $ red. E → E + (E)
E I $ accept
    
```



The Stack

- Left string can be implemented as a stack
 - Top of the stack is the **I**
- Shift pushes a terminal on the stack
- Reduce pops 0 or more symbols from the stack (production rhs) and pushes a non-terminal on the stack (production lhs)

2/12/09

Prof. Hilfinger CS164 Lecture 8

26

Key Issue: When to Shift or Reduce?

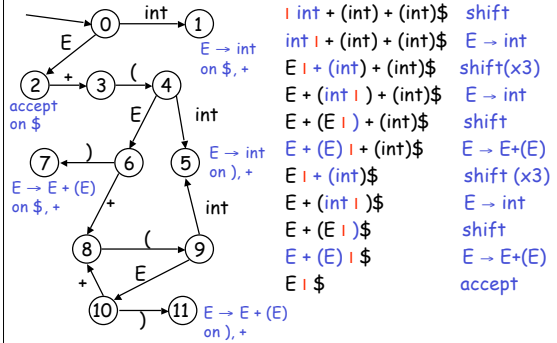
- Decide based on the left string (the stack)
- Idea: use a finite automaton (DFA) to decide when to shift or reduce
 - The DFA input is the stack up to potential handle
 - DFA alphabet consists of terminals and nonterminals
 - DFA recognizes complete handles
- We run the DFA on the stack and we examine the resulting state **X** and the token **tok** after **I**
 - If **X** has a transition labeled **tok** then *shift*
 - If **X** is labeled with "**A** → **β** on **tok**" then *reduce*

2/12/09

Prof. Hilfinger CS164 Lecture 8

27

LR(1) Parsing. An Example



Representing the DFA

- Parsers represent the DFA as a 2D table
 - As for table-driven lexical analysis
- Lines correspond to DFA states
- Columns correspond to terminals and non-terminals
- In classical treatments, columns are split into:
 - Those for terminals: action table
 - Those for non-terminals: goto table

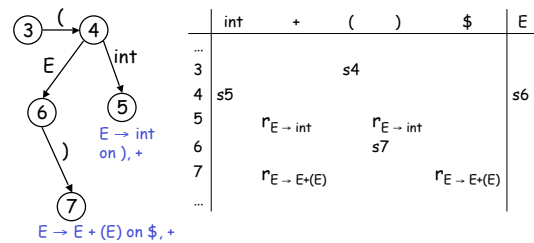
2/12/09

Prof. Hilfinger CS164 Lecture 8

29

Representing the DFA. Example

- The table for a fragment of our DFA:



2/12/09

Prof. Hilfinger CS164 Lecture 8

30

The LR Parsing Algorithm

- After a shift or reduce action we rerun the DFA on the entire stack
 - This is wasteful, since most of the work is repeated
- So record, for each stack element, state of the DFA after that state
- LR parser maintains a stack
 - $\langle \text{sym}_1, \text{state}_1 \rangle \dots \langle \text{sym}_n, \text{state}_n \rangle$
 - state_k is the final state of the DFA on $\text{sym}_1 \dots \text{sym}_k$

2/12/09

Prof. Hilfinger CS164 Lecture 8

31

The LR Parsing Algorithm

```

Let I = w1w2...wn$ be initial input
Let j = 1
Let DFA state 0 be the start state
Let stack = ⟨ dummy, 0 ⟩
repeat
  case table[top_state(stack), I[j]] of
    shift k: push ( I[j], k ); j += 1
    reduce X → α:
      pop |α| pairs,
      push (X, table[top_state(stack), X])
    accept: halt normally
    error: halt and report error
    
```

2/12/09

Prof. Hilfinger CS164 Lecture 8

32

Parsing Contexts

- Consider the state describing the situation at the $|$ in the stack $E + (\text{int}) + (\text{int})$
- Context:
 - We are looking for an $E \rightarrow E + (\bullet E)$
 - Have have seen $E + ($ (from the right-hand side)
 - We are also looking for $E \rightarrow \bullet \text{int}$ or $E \rightarrow \bullet E + (E)$
 - Have seen nothing from the right-hand side
- One DFA state describes a set of such contexts
- (Traditionally, use \bullet to show where the $|$ is.)

2/12/09

Prof. Hilfinger CS164 Lecture 8

33

LR(1) Items

- An *LR(1) item* is a pair:
 - $X \rightarrow \alpha \bullet \beta, a$
 - $X \rightarrow \alpha \beta$ is a production
 - a is a terminal (the lookahead terminal)
 - LR(1) means 1 lookahead terminal
- $[X \rightarrow \alpha \bullet \beta, a]$ describes a context of the parser
 - We are trying to find an X followed by an a , and
 - We have α already on top of the stack
 - Thus we need to see next a prefix derived from βa

2/12/09

Prof. Hilfinger CS164 Lecture 8

34

Convention

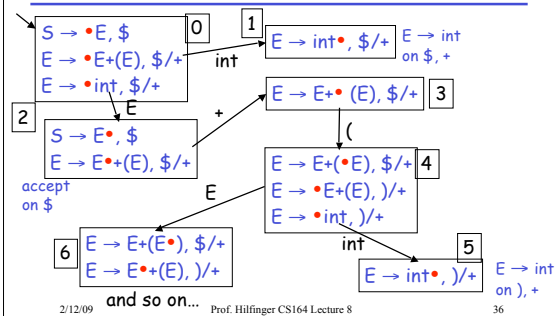
- We add to our grammar a fresh new start symbol S and a production $S \rightarrow E$
 - Where E is the old start symbol
 - No need to do this if E had only one production
- The initial parsing context contains:
 - $S \rightarrow \bullet E, \$$
 - Trying to find an S as a string derived from $E\$$
 - The stack is empty

2/12/09

Prof. Hilfinger CS164 Lecture 8

35

Constructing the Parsing DFA. Example.



2/12/09

Prof. Hilfinger CS164 Lecture 8

36

LR Parsing Tables. Notes

- Parsing tables (i.e. the DFA) can be constructed automatically for a CFG
- But we still need to understand the construction to work with parser generators
 - E.g., they report errors in terms of sets of items
- What kind of errors can we expect?

2/12/09

Prof. Hilfinger CS164 Lecture 8

37

Shift/Reduce Conflicts

- If a DFA state contains both $[X \rightarrow \alpha \cdot a \beta, b]$ and $[Y \rightarrow \gamma \cdot, a]$
- Then on input "a" we could either
 - Shift into state $[X \rightarrow \alpha a \cdot \beta, b]$, or
 - Reduce with $Y \rightarrow \gamma$
- This is called a *shift-reduce conflict*

2/12/09

Prof. Hilfinger CS164 Lecture 8

38

Shift/Reduce Conflicts

- Typically due to ambiguities in the grammar
- Classic example: the dangling else
 - $S \rightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S \mid \text{OTHER}$
- Will have DFA state containing
 - $[S \rightarrow \text{if } E \text{ then } S \cdot, \text{ else}]$
 - $[S \rightarrow \text{if } E \text{ then } S \cdot \text{ else } S, \$]$
- If *else* follows then we can shift or reduce

2/12/09

Prof. Hilfinger CS164 Lecture 8

39

More Shift/Reduce Conflicts

- Consider the ambiguous grammar
 - $E \rightarrow E + E \mid E * E \mid \text{int}$
- We will have the states containing
 - $[E \rightarrow E * \cdot E, +]$ $[E \rightarrow E * E \cdot, +]$
 - $[E \rightarrow \cdot E + E, +] \Rightarrow^E [E \rightarrow E \cdot + E, +]$
 - ...
- Again we have a shift/reduce on input +
 - We need to reduce (* binds more tightly than +)
 - Solution: declare the precedence of * and +

2/12/09

Prof. Hilfinger CS164 Lecture 8

40

More Shift/Reduce Conflicts

- In bison declare precedence and associativity of *terminal symbols*:
 - `%left +`
 - `%left *`
- Precedence of a rule = that of its last terminal
 - See bison manual for ways to override this default
- Resolve shift/reduce conflict with a *shift* if:
 - input terminal has higher precedence than the rule
 - the precedences are the same and right associative

2/12/09

Prof. Hilfinger CS164 Lecture 8

41

Using Precedence to Solve S/R Conflicts

- Back to our example:
 - $[E \rightarrow E * \cdot E, +]$ $[E \rightarrow E * E \cdot, +]$
 - $[E \rightarrow \cdot E + E, +] \Rightarrow^E [E \rightarrow E \cdot + E, +]$
 - ...
- Will choose reduce because precedence of rule $E \rightarrow E * E$ is higher than of terminal +

2/12/09

Prof. Hilfinger CS164 Lecture 8

42

Using Precedence to Solve S/R Conflicts

- Same grammar as before
 $E \rightarrow E + E \mid E * E \mid \text{int}$
- We will also have the states
 $[E \rightarrow E + \bullet E, +]$ $[E \rightarrow E + E \bullet, +]$
 $[E \rightarrow \bullet E + E, +] \Rightarrow^E [E \rightarrow E \bullet + E, +]$
 ...
- Now we also have a shift/reduce on input +
 - We choose reduce because $E \rightarrow E + E$ and + have the same precedence and + is left-associative

2/12/09

Prof. Hilfinger CS164 Lecture 8

43

Using Precedence to Solve S/R Conflicts

- Back to our dangling else example
 $[S \rightarrow \text{if } E \text{ then } S \bullet, \text{ else}]$
 $[S \rightarrow \text{if } E \text{ then } S \bullet \text{ else } S, x]$
- Can eliminate conflict by declaring *else* with higher precedence than *then*
- However, best to avoid overuse of precedence declarations or you'll end with unexpected parse trees

2/12/09

Prof. Hilfinger CS164 Lecture 8

44

Reduce/Reduce Conflicts

- If a DFA state contains both
 $[X \rightarrow \alpha \bullet, a]$ and $[Y \rightarrow \beta \bullet, a]$
 - Then on input "a" we don't know which production to reduce
- This is called a *reduce/reduce conflict*

2/12/09

Prof. Hilfinger CS164 Lecture 8

45

Reduce/Reduce Conflicts

- Usually due to gross ambiguity in the grammar
- Example: a sequence of identifiers
 $S \rightarrow \epsilon \mid \text{id} \mid \text{id } S$
- There are two parse trees for the string *id*
 $S \rightarrow \text{id}$
 $S \rightarrow \text{id } S \rightarrow \text{id}$
- How does this confuse the parser?

2/12/09

Prof. Hilfinger CS164 Lecture 8

46

More on Reduce/Reduce Conflicts

- Consider the states
 $[S' \rightarrow \bullet S, \$]$ $[S \rightarrow \text{id} \bullet, \$]$
 $[S \rightarrow \bullet, \$] \Rightarrow^{\text{id}} [S \rightarrow \bullet, \$]$
 $[S \rightarrow \bullet \text{id}, \$]$ $[S \rightarrow \bullet \text{id}, \$]$
 $[S \rightarrow \bullet \text{id } S, \$]$ $[S \rightarrow \bullet \text{id } S, \$]$
- Reduce/reduce conflict on input \$
 $S' \rightarrow S \rightarrow \text{id}$
 $S' \rightarrow S \rightarrow \text{id } S \rightarrow \text{id}$
- Better rewrite the grammar: $S \rightarrow \epsilon \mid \text{id } S$

2/12/09

Prof. Hilfinger CS164 Lecture 8

47

Relation to Bison

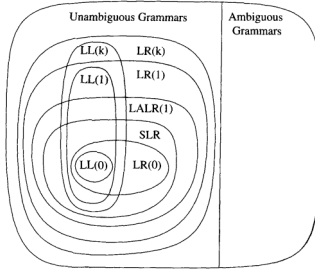
- Bison builds this kind of machine.
- However, for efficiency concerns, collapses many of the states together.
- Causes some additional conflicts, but not many.
- The machines discussed here are LR(1) engines. Bison's optimized versions are LALR(1) engines.

2/12/09

Prof. Hilfinger CS164 Lecture 8

48

A Hierarchy of Grammar Classes



From Andrew Appel,
"Modern Compiler
Implementation in Java"

2/12/09

Prof. Hilfinger CS164 Lecture 8

49

Notes on Parsing

• Parsing

- A simple parser: LL(1), recursive descent
- A more powerful parser: LR(1)
- An efficiency hack: LALR(1)
- We use LALR(1) parser generators
- Earley's algorithm provides a complete algorithm for parsing context-free languages.

2/12/09

Prof. Hilfinger CS164 Lecture 8

50