# LR Parsing LALR Parser Generators

#### Outline

- Review of bottom-up parsing
- Computing the parsing DFA
- Using parser generators

## **Bottom-up Parsing (Review)**

- A bottom-up parser rewrites the input string to the start symbol
- The state of the parser is described as

#### αΙγ

- $\alpha$  is a stack of terminals and non-terminals
- $\gamma$  is the string of terminals not yet examined
- Initially:  $|x_1x_2...x_n|$

## The Shift and Reduce Actions (Review)

- Recall the CFG:  $E \rightarrow int | E + (E)$
- A bottom-up parser uses two kinds of actions:
- <u>Shift</u> pushes a terminal from input on the stack

 $E + (int) \Rightarrow E + (inti)$ 

 <u>Reduce</u> pops 0 or more symbols off of the stack (production RHS) and pushes a nonterminal on the stack (production LHS)

 $\mathsf{E} + (\mathsf{E} + (\mathsf{E})) \Rightarrow \mathsf{E} + (\mathsf{E})) \Rightarrow \mathsf{E} + (\mathsf{E})$ 

## Key Issue: When to Shift or Reduce?

- Idea: use a deterministic finite automaton (DFA) to decide when to shift or reduce
  - The input is the stack
  - The language consists of terminals and non-terminals
- 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 <u>shift</u>
  - If X is labeled with " $A \rightarrow \beta$  on tok" then <u>reduce</u>

#### LR(1) Parsing: An Example



#### Representing the DFA

- Parsers represent the DFA as a 2D table
  - Recall table-driven lexical analysis
- Lines correspond to DFA states
- Columns correspond to terminals and nonterminals
- Typically columns are split into:
  - Those for terminals: the action table
  - Those for non-terminals: the goto table

#### Representing the DFA: Example

#### The table for a fragment of our DFA:



## 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
- Remember for each stack element on which state it brings the DFA
- LR parser maintains a stack  $\langle sym_1, state_1 \rangle \dots \langle sym_n, state_n \rangle$ state<sub>k</sub> is the final state of the DFA on  $sym_1 \dots sym_k$

## The LR Parsing Algorithm

```
let I = w$ be initial input
let j = 0
let DFA state 0 be the start state
let stack = \langle \text{ dummy, 0} \rangle
   repeat
        case action[top_state(stack), I[j]] of
                shift k: push \langle I[j++], k \rangle
                reduce X \rightarrow A:
                     pop |A| pairs,
                     push ( X, goto[top_state(stack), X] )
                accept: halt normally
                error: halt and report error
```

### Key Issue: How is the DFA Constructed?

- The stack describes the context of the parse
  - What non-terminal we are looking for
  - What production RHS we are looking for
  - What we have seen so far from the RHS
- Each DFA state describes several such contexts
  - E.g., when we are looking for non-terminal E, we might be looking either for an int or an E + (E) RHS



- An <u>LR(0) item</u> is a production with a "I" somewhere on the RHS
- The items for  $T \rightarrow (E)$  are
  - $T \rightarrow I (E)$  $T \rightarrow (IE)$  $T \rightarrow (EI)$  $T \rightarrow (EI)$
- The only item for  $X \to \varepsilon$  is  $X \to I$

# LR(0) Items: Intuition

- An item  $[X \rightarrow \alpha | \beta]$  says that
  - the parser is looking for an X
  - it has an  $\alpha$  on top of the stack
  - Expects to find a string derived from  $\boldsymbol{\beta}$  next in the input
- Notes:
  - $[X \rightarrow \alpha \mid \alpha\beta]$  means that a should follow. Then we can shift it and still have a viable prefix
  - $[X \rightarrow \alpha I]$  means that we could reduce X
    - But this is not always a good idea !

# LR(1) Items

• An <u>LR(1) item</u> is a pair:

 $X \rightarrow \alpha I \beta$ , a

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

#### Note

- The symbol I was used before to separate the stack from the rest of input
  - $\alpha$  I  $\gamma,$  where  $\alpha$  is the stack and  $\gamma$  is the remaining string of terminals
- In items I is used to mark a prefix of a production RHS:

 $X \rightarrow \alpha \, \iota \, \beta$ , a

- Here  $\beta$  might contain terminals as well
- In both case the stack is on the left of I

#### Convention

- We add to our grammar a fresh new start symbol S and a production  $\mathsf{S} \to \mathsf{E}$ 
  - Where E is the old start symbol
- The initial parsing context contains:  $S \rightarrow IE$ , \$
  - Trying to find an S as a string derived from E\$
  - The stack is empty

# LR(1) Items (Cont.)

- In context containing  $E \rightarrow E + I(E)$ , +
  - If (follows then we can perform a shift to context containing

 $E \rightarrow E + (IE)$ , +

• In context containing

 $E \rightarrow E + (E)I$ , +

- We can perform a reduction with  $\mathsf{E} \to \mathsf{E}$  + (  $\mathsf{E}$  )
- But only if a + follows

# LR(1) Items (Cont.)

# • Consider the item $E \rightarrow E + (IE)$ , +

- We expect a string derived from E) +
- There are two productions for E  $E \rightarrow int$  and  $E \rightarrow E + (E)$
- We describe this by extending the context with two more items:

 $E \rightarrow i \text{ int } , )$  $E \rightarrow i E + (E) , )$ 

## The Closure Operation

• The operation of extending the context with items is called the closure operation

```
Closure(Items) =
repeat
for each [X \rightarrow \alpha \mid Y\beta, a] in Items
for each production Y \rightarrow \gamma
for each b in First(\beta a)
add [Y \rightarrow \mid \gamma, b] to Items
until Items is unchanged
```

### Constructing the Parsing DFA (1)

• Construct the start context:  $Closure(\{S \rightarrow I E, \$\})$ 

$$S \rightarrow I E , \$$$
  

$$E \rightarrow I E+(E), \$$$
  

$$E \rightarrow I int , \$$$
  

$$E \rightarrow I E+(E), +$$
  

$$E \rightarrow I int , +$$

• We abbreviate as:

$$S \rightarrow I E$$
 , \$  
E  $\rightarrow I E+(E)$  , \$/+  
E  $\rightarrow I int$  , \$/+

#### Constructing the Parsing DFA (2)

- A DFA state is a closed set of LR(1) items
- The start state contains  $[S \rightarrow IE, \$]$

- A state that contains  $[X \rightarrow \alpha I, b]$  is labelled with "reduce with  $X \rightarrow \alpha$  on b"
- And now the transitions ...

## The DFA Transitions

- A state "State" that contains  $[X \rightarrow \alpha | y\beta, b]$ has a transition labeled y to a state that contains the items "Transition(State, y)"
  - y can be a terminal or a non-terminal

```
Transition(State, y)

Items = \emptyset

for each [X \rightarrow \alpha | y\beta, b] in State

add [X \rightarrow \alpha y | \beta, b] to Items

return Closure(Items)
```

#### Constructing the Parsing DFA: Example



## 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?

#### Shift/Reduce Conflicts

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

## Shift/Reduce Conflicts

- Typically due to ambiguities in the grammar
- Classic example: the dangling else  $S \rightarrow if E$  then S | if E then S else S | OTHER
- Will have DFA state containing  $[S \rightarrow if E \text{ then } S \text{ I}, else]$  $[S \rightarrow if E \text{ then } S \text{ I} \text{ else } S, x]$
- If else follows then we can shift or reduce
- Default (yacc, ML-yacc, etc.) is to shift
  - Default behavior is as needed in this case

## More Shift/Reduce Conflicts

...

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

...

## More Shift/Reduce Conflicts

- In yacc declare precedence and associativity:
   %left +
   %left \*
- Precedence of a rule = that of its last terminal
   See yacc manual for ways to override this default
- Resolve shift/reduce conflict with a <u>shift</u> if:
  - no precedence declared for either rule or terminal
  - input terminal has higher precedence than the rule
  - the precedences are the same and right associative

#### Using Precedence to Solve S/R Conflicts

- Back to our example:  $\begin{bmatrix} E \rightarrow E * I E, +] & [E \rightarrow E * E I, +] \\ [E \rightarrow I E + E, +] \Rightarrow^{E} & [E \rightarrow E I + E, +] \\ \vdots & \vdots & \vdots & \vdots \\ \end{bmatrix}$
- Will choose reduce because precedence of rule  $E \rightarrow E^* E$  is higher than of terminal +

### Using Precedence to Solve S/R Conflicts

• Same grammar as before  $E \rightarrow E + E \mid E^*E \mid int$ 

...

- We will also have the states  $\begin{bmatrix} E \rightarrow E + I E, +] & \begin{bmatrix} E \rightarrow E + E I, +] \\ E \rightarrow I E + E, +\end{bmatrix} \Rightarrow^{E} & \begin{bmatrix} E \rightarrow E I + E, +] \end{bmatrix}$
- Now we also have a shift/reduce on input +
  - We choose reduce because  $\mathsf{E}\to\mathsf{E}+\mathsf{E}$  and + have the same precedence and + is left-associative

...

## Using Precedence to Solve S/R Conflicts

- Back to our dangling else example  $[S \rightarrow if E \text{ then } S \text{ I}, else]$  $[S \rightarrow if E \text{ then } S \text{ I} else S, x]$
- Can eliminate conflict by declaring else having higher precedence than then
- But this starts to look like "hacking the tables"
- Best to avoid overuse of precedence declarations or we will end with unexpected parse trees

The term "precedence declaration" is misleading!

These declarations do not define precedence: they define conflict resolutions

I.e., they instruct shift-reduce parsers to resolve conflicts in certain ways

The two are not quite the same thing!

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

#### **Reduce/Reduce Conflicts**

- Usually due to gross ambiguity in the grammar
- Example: a sequence of identifiers

 $S \rightarrow \epsilon$  | id | id S

- There are two parse trees for the string id  $S \rightarrow id$  $S \rightarrow id S \rightarrow id$
- How does this confuse the parser?

#### More on Reduce/Reduce Conflicts

- Consider the states  $[S \rightarrow id I, \$]$   $[S' \rightarrow I S, \$]$   $[S \rightarrow id I S, \$]$
- Reduce/reduce conflict on input \$
  - $S' \rightarrow S \rightarrow id$  $S' \rightarrow S \rightarrow id S \rightarrow id$
- Better rewrite the grammar:  $S \rightarrow \epsilon \mid id S$

#### Using Parser Generators

- Parser generators automatically construct the parsing DFA given a CFG
  - Use precedence declarations and default conventions to resolve conflicts
  - The parser algorithm is the same for all grammars (and is provided as a library function)
- But most parser generators do not construct the DFA as described before
  - Because the LR(1) parsing DFA has 1000s of states even for a simple language

# LR(1) Parsing Tables are Big

- But many states are similar, e.g. 1
  5  $E \rightarrow int I, $/+$   $e \rightarrow int on $, +$ and  $E \rightarrow int I, )/+$   $e \rightarrow int on ), +$
- <u>Idea</u>: merge the DFA states whose items differ only in the lookahead tokens
  - We say that such states have the same core
- We obtain  $\begin{bmatrix} 1' \\ E \rightarrow int I, \$/+/ \end{bmatrix} \xrightarrow{E \rightarrow int}_{on \$, +, }$

#### The Core of a Set of LR Items

# <u>Definition</u>: The core of a set of LR items is the set of first components

- Without the lookahead terminals

- Example: the core of  $\{ [X \to \alpha \, | \, \beta, \, b], \, [Y \to \gamma \, | \, \delta, \, d] \}$  is

 $\{X \to \alpha \mid \beta, Y \to \gamma \mid \delta\}$ 

#### LALR States

• Consider for example the LR(1) states  $\{[X \rightarrow \alpha I, a], [Y \rightarrow \beta I, c]\}$ 

 $\{[X \rightarrow \alpha I, b], [Y \rightarrow \beta I, d]\}$ 

- They have the same core and can be merged
- And the merged state contains:  $\{[X \rightarrow \alpha I, a/b], [Y \rightarrow \beta I, c/d]\}$
- These are called LALR(1) states
  - Stands for LookAhead LR
  - Typically 10 times fewer LALR(1) states than LR(1)

# A LALR(1) DFA

- Repeat until all states have distinct core
  - Choose two distinct states with same core
  - Merge the states by creating a new one with the union of all the items
  - Point edges from predecessors to new state
  - New state points to all the previous successors



#### Conversion LR(1) to LALR(1): Example.



#### The LALR Parser Can Have Conflicts

Consider for example the LR(1) states
 {[X → α I, a], [Y → β I, b]}
 {[X → α I, b], [Y → β I, a]}
 {IX → α I, b], [Y → β I, a]}
 And the merged LALR(1) state

 $\{[X \rightarrow \alpha \text{ I, a/b}], [Y \rightarrow \beta \text{ I, a/b}]\}$ 

- Has a <u>new</u> reduce/reduce conflict
- In practice such cases are rare

# LALR vs. LR Parsing: Things to keep in mind

- LALR languages are not natural
  - They are an efficiency hack on LR languages
- Any reasonable programming language has a LALR(1) grammar
- LALR(1) parsing has become a standard for programming languages and for parser generators

#### A Hierarchy of Grammar Classes



#### Semantic Actions in LR Parsing

- We can now illustrate how semantic actions are implemented for LR parsing
- Keep attributes on the stack
- On shifting a, push attribute for a on stack
- On reduce  $X \to \alpha$ 
  - pop attributes for  $\alpha$
  - compute attribute for X
  - and push it on the stack

# Performing Semantic Actions: Example

- Recall the example
  - $\begin{array}{ll} E \rightarrow T + E_1 & \{ \text{ E.val} = \text{T.val} + E_1.val \} \\ & | & T & \{ \text{ E.val} = \text{T.val} \} \\ & T \rightarrow \text{int} * T_1 & \{ \text{ T.val} = \text{int.val} * T_1.val \} \\ & | & \text{int} & \{ \text{ T.val} = \text{int.val} \} \end{array}$
- Consider the parsing of the string  $3 \times 5 + 8$

#### Performing Semantic Actions: Example

```
int * int + int
int_3 | * int + int
int_3 * int + int
int_3 * int_5 + int_7
int_3 * T_5 + int
T<sub>15</sub> + int
T_{15} + | int
T_{15} + int_8
T_{15} + T_8
T_{15} + E_8
E_{23}
```

shift shift shift reduce  $T \rightarrow int$ reduce  $T \rightarrow int * T$ shift shift reduce  $T \rightarrow int$ reduce  $E \rightarrow T$ reduce  $E \rightarrow T + E$ accept

3 \* 5 + 8

#### Notes

- The previous example shows how synthesized attributes are computed by LR parsers
- It is also possible to compute inherited attributes in an LR parser

## Notes on Parsing

- Parsing
  - A solid foundation: context-free grammars
  - A simple parser: LL(1)
  - A more powerful parser: LR(1)
  - An efficiency hack: LALR(1)
  - LALR(1) parser generators
- Next time we move on to semantic analysis

#### Supplement to LR Parsing

## Strange Reduce/Reduce Conflicts due to LALR Conversion (and how to handle them)

### Strange Reduce/Reduce Conflicts

- P parameters specification
- R result specification
- N a parameter or result name
- T a type name
- NL a list of names

#### Strange Reduce/Reduce Conflicts

- In P an id is a
  - N when followed by , or :
  - T when followed by id
- In R an id is a
  - N when followed by :
  - T when followed by,
- This is an LR(1) grammar
- But it is not LALR(1). Why?
  - For obscure reasons

### A Few LR(1) States



Compiler Design I (2011)

# What Happened?

- Two distinct states were confused because they have the same core
- Fix: add dummy productions to distinguish the two confused states
- E.g., add

# $R \rightarrow id \ bogus$

- bogus is a terminal not used by the lexer
- This production will never be used during parsing
- But it distinguishes R from P

#### A Few LR(1) States After Fix



Compiler Design I (2011)