	Outline
LR Parsing LALR Parser Generators	 Review of bottom-up parsing
	 Computing the parsing DFA
	 Using parser generators
	Compiler Design I (2011) 2
Bottom-up Parsing (Review)	The Shift and Reduce Actions (Review)
 A bottom-up parser rewrites the input string to the start symbol The state of the parser is described as α γ α is a stack of terminals and non-terminals γ is the string of terminals not yet examined 	 Recall the CFG: E → int E + (E) A bottom-up parser uses two kinds of actions: <u>Shift</u> pushes a terminal from input on the stack E + (1 int) ⇒ E + (int 1)
 Initially: x₁x₂x_n 	• <u>Reduce</u> pops 0 or more symbols off of the stack (production RHS) and pushes a non-terminal on the stack (production LHS) $E + (E + (E) +) \Rightarrow E + (E +)$
Compiler Design I (2011) 3	Compiler Design I (2011) 4

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 shift
 - If X is labeled with "A $\rightarrow \beta$ on tok" then \underline{reduce}

LR(1) Parsing: An Example



Representing the DFA

Compiler Design I (2011)

- 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

7



The table for a fragment of our DFA:

Compiler Design I (2011)

The LR Parsing Algorithm	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 (sym1, state1)(symn, staten) statek is the final state of the DFA on sym1 symk 	<pre>let I = w\$ be initial input let j = 0 let DFA state 0 be the start state let stack = < dummy, 0 > repeat case action[top_state(stack), I[j]] of shift k: push < I[j++], k > reduce X → A: pop A pairs, push { X, goto[top_state(stack), X] > accept: halt normally error: halt and report error</pre>
Compiler Design I (2011) 9	Compiler Design I (2011) 10
Key Issue: How is the DFA Constructed?The stack describes the context of the parse	 LR(0) Items An <u>LR(0) item</u> is a production with a "I"
 What non-terminal we are looking for What production RHS we are looking for What we have seen so far from the RHS 	somewhere on the RHS • The items for $T \rightarrow (E)$ are $T \rightarrow I (E)$
 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 	$T \rightarrow (1E)$ $T \rightarrow (E1)$ $T \rightarrow (E)$ • The only item for $X \rightarrow \varepsilon$ is $X \rightarrow 1$

LR(0) Items: Intuition

- An item $[X \rightarrow \alpha]$
 - the parser is lo
 - it has an α on to
 - Expects to find input

• Notes:

Compiler Design I (2011)

- $[X \rightarrow \alpha \mid a\beta]$ means can shift it and
- $[X \rightarrow \alpha I]$ means
 - But this is not a

LR(1) Items

 An item [X → α β] says that the parser is looking for an X it has an α on top of the stack Expects to find a string derived from β next in the input Notes: [X → α αβ] means that a should follow. Then we can shift it and still have a viable prefix [X → α] means that we could reduce X But this is not always a good idea ! 	 An LR(1) item is a pair: X → α 1 β, a X → αβ is a production a is a terminal (the lookahead terminal) LR(1) means 1 lookahead terminal [X → α 1 β, a] describes a context of the parser We are trying to find an X followed by an a, and We have (at least) α already on top of the stack Thus we need to see next a prefix derived from βa
mpiler Design I (2011) 13	Compiler Design I (2011) 14
 Note The symbol I was used before to separate the stack from the rest of input α I γ, where α is the stack and γ is the remaining string of terminals In items I is used to mark a prefix of a production RHS: X → α I β, α Here β 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 S → E Where E is the old start symbol The initial parsing context contains: S → I E , \$ Trying to find an S as a string derived from E\$ The stack is empty

LR(1) Items (Cont.) LR(1) Items (Cont.) Consider the item In context containing $E \rightarrow E + I(E)$, + $E \rightarrow E + (IE)$ + - If (follows then we can perform a shift to context • We expect a string derived from E) + containing • There are two productions for E $E \rightarrow E + (IE)$ + $E \rightarrow int$ and $E \rightarrow E + (E)$ • In context containing • We describe this by extending the context $E \rightarrow E + (E)_{I}$ + with two more items: - We can perform a reduction with $E \rightarrow E + (E)$ $E \rightarrow i \text{ int}$,) - But only if a + follows $E \rightarrow IE + (E)$ 17 Compiler Design I (2011) Compiler Design I (2011) 18 The Closure Operation Constructing the Parsing DFA (1) • The operation of extending the context with • Construct the start context: $Closure(\{S \rightarrow I \in \})$ items is called the closure operation $S \rightarrow I E$, \$ $E \rightarrow I E+(E)$, \$ Closure(Items) = $E \rightarrow I$ int , \$ $E \rightarrow I E+(E)$, + repeat $E \rightarrow I$ int . + for each $[X \rightarrow \alpha \mid Y\beta, a]$ in Items for each production $Y \rightarrow \gamma$ • We abbreviate as: for each b in First(βa) $S \rightarrow I E$, \$ $E \rightarrow I E+(E)$, \$/+ add $[Y \rightarrow I \gamma, b]$ to Items until Items is unchanged $E \rightarrow I$ int .\$/+



Compiler Design I (2011)

Shift/Reduce Conflicts

 If a DFA state contains both [X → α aβ, b] and [Y → γ , a] Then on input "a" we could either Shift into state [X → α β, b], or Reduce with Y → γ This is called a <i>shift-reduce conflict</i> 	 Typically due to ambiguities in the grammar Classic example: the dangling else → if E then S if E then S else S OTHER Will have DFA state containing [S → if E then S I, else] [S → if E then S I 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 	
Compiler Design I (2011) 25	Compiler Design I (2011)	
More Shift/Reduce Conflicts	More Shift/Reduce Conflicts	
 Consider the ambiguous grammar E→E+E E*E int We will have the states containing [E→E*IE, +] [E→E*EI, +] [E→IE+E, +] ⇒^E [E→EI+E, +] 	 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 	

Shift/Reduce Conflicts

Using Precedence to Solve S/R Conflicts
 Same grammar as before E→E+E E*E int We will also have the states [E→E+1E, +] [E→E+E1, +] [E→1E+E, +] □ Now we also have a shift/reduce on input + We choose reduce because E→E+E and + have the same precedence and + is left-associative
Compiler Design I (2011) 30
Precedence Declarations Revisited
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!

Reduce/Reduce Conflicts	Reduce/Reduce Conflicts	
 If a DFA state contains both [X → α ι, a] and [Y → β ι, a] - Then on input "a" we don't know which production to reduce 	 Usually due to gross ambiguity in the grammar Example: a sequence of identifiers S → ε id id S 	
 This is called a <i>reduce/reduce conflict</i> 	 There are two parse trees for the string id S → id S → id S → id How does this confuse the parser? 	
Compiler Design I (2011) 33	Compiler Design I (2011) 34	
More on Reduce/Reduce Conflicts	Using Parser Generators	
 Consider the states [S → id I, \$] [S' → I S, \$] [S → id I S, \$] [S → I, \$] ⇒^{id} [S → id I S, \$] [S → I id, \$] ⇒^{id} [S → I, \$] [S → I id, \$] [S → I id, \$] [S → I id S, \$] [S → I id S, \$] Reduce/reduce conflict on input \$ S' → S → id S' → S → id S → id Better rewrite the grammar: S → ε id S 	 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 	

35 Compiler Design I (2011)

Compiler Design I (2011)





The Core of a Set of LR Items

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



The LALR Parser Can Have Conflicts

•	Consider for example the LR(1) states
	$\{[X \rightarrow \alpha I, a], [Y \rightarrow \beta I, b]\}$
	$\{[X \rightarrow \alpha I, b], [Y \rightarrow \beta I, a]\}$
•	And the merged LALR(1) state
	{[X $\rightarrow \alpha$ I, a/b], [Y $\rightarrow \beta$ 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



Compiler Design I (2011)

Semantic Actions	in LR Parsing	Performing Semantic Actions: Example
• We can now illus are implemented	trate how semantic actions for LR parsing	 Recall the example
 Keep attributes 	on the stack	$F \rightarrow T + F_{1} = \{F \text{ val} = T \text{ val} + F_{1} \text{ val} \}$
		$\int T = \{F_{i} \in T_{i} \in T_{i} \in T_{i} \}$
• On shifting a pu	sh attribute for a on stack	$T \leq int * T \qquad \{T, val = int val * T val \}$
On shijing u, pu		$\downarrow int \qquad (Type = int yel)$
• On reduce $X \rightarrow 0$		INT { 1.val = INT.val }
 pop attributes to compute attributes to and push it on th 	te for X le stack	\cdot Consider the parsing of the string 3 * 5 + 8
Compiler Design I (2011)	tie Actioner Exemple	45 Compiler Design I (2011) 4
Pertorming Seman	tic Actions: Example	Notes
int * int + int int ₃ $ $ * int + int int ₃ * $ $ int + int	shift <mark>3 * 5 + 8</mark> shift shift	 The previous example shows how synthesized attributes are computed by LR parsers
$int_3 * int_5 + int$	reduce $T \rightarrow int$	This class prescible the computer intervited
$\operatorname{int}_3 * \operatorname{T}_5 + \operatorname{int}$	reduce $T \rightarrow int * T$	• It is also possible to compute innerited
T_{15} + Int T + I int	Shift shift	diffibules in an ER paiser
$T_{15} + int_{0}$	reduce $T \rightarrow int$	
$T_{15} + T_8$	reduce $E \rightarrow T$	
$T_{15} + E_8$	reduce $E \rightarrow T + E$	
E ₂₃	accept	
Compiler Design L(2011)		47 Compiler Design I (2011) 4

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)
$\begin{array}{llllllllllllllllllllllllllllllllllll$	 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

Compiler Design I (2011)

A Few LR(1) States



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

54

- But it distinguishes R from P

Compiler Design I (2011)

Compiler Design I (2011)

 $T \rightarrow I id$

 $R \rightarrow . T$

 $T \rightarrow . id$ $N \rightarrow . id$

 $R \rightarrow N : T$

 $R \rightarrow . id bogus$

id

2

 $T \rightarrow id I$

 $N \rightarrow id I$

 $R \rightarrow id I bogus$

4