



CS 412 Introduction to Compilers

Andrew Myers
Cornell University

Lecture 8: Parser generators and abstract syntax trees

LR(1) parsing

- As much power as possible out of 1 look-ahead symbol parsing table
- LR(1) grammar = recognizable by a shift/reduce parser with 1 look-ahead.
- LR(1) item = LR(0) item + look-ahead symbols possibly following production

LR(0): $S \rightarrow \cdot S + E$

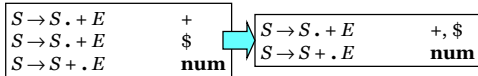
LR(1): $S \rightarrow \cdot S + E \quad +$

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

2

LR(1) state

- LR(1) state = set of LR(1) items
- LR(1) item = LR(0) item + set of look-ahead symbols
- No two items in state have same production + dot configuration



CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

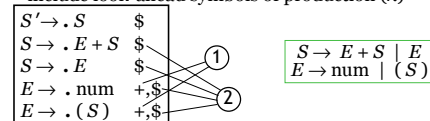
3

LR(1) closure

Consider $A \rightarrow \beta \cdot C \delta \lambda$

Closure formed just as for LR(0) *except*

1. Look-ahead symbols include characters following the non-terminal symbol to the right of dot: $\text{FIRST}(\delta)$
2. If non-terminal symbol may produce last symbol of production (δ is nullable), look-ahead symbols include look-ahead symbols of production (λ)

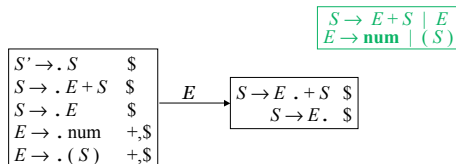


CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

4

LR(1) DFA construction

- Given LR(1) state, for each symbol (terminal or non-terminal) following a dot, construct a state with dot shifted across symbol, perform closure



CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

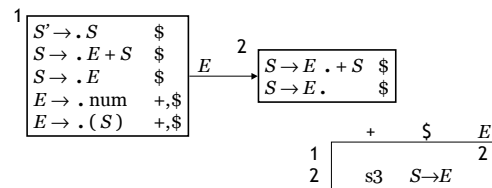
5

LR(1) example

Reductions unambiguous if:
look-aheads are disjoint,

not to right of any dot in state

$S \rightarrow E + S \mid E$
 $E \rightarrow \text{num} \mid (S)$



CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

6

LALR grammars

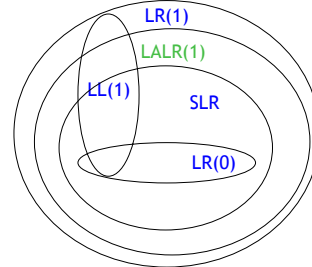
- Problem with LR(1): too many states
- LALR(1) (Look-Ahead LR)
 - Merge any two LR(1) states whose items are identical except look-ahead
 - Results in smaller parser tables—works extremely well in practice
 - Usual technology for automatic parser generators

$$\begin{array}{|l} S \rightarrow id \cdot + \\ S \rightarrow E \cdot \$ \end{array} + \begin{array}{|l} S \rightarrow id \cdot \$ \\ S \rightarrow E \cdot + \end{array} = ?$$

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

7

Classification of Grammars



CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

8

How are parsers written?

- Automatic parser generators: yacc, bison, CUP
- Accept LALR(1) grammar specification
 - *plus*: declarations of precedence, associativity
 - output: LR parser code (inc. parsing table)
- Some parser generators accept LL(1), e.g. **javacc** – less powerful
- Rest of this lecture: how to use parser generators
- Can we use parsers for programs other than compilers?

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

9

Associativity

$$S \rightarrow S + E \mid E$$

$$E \rightarrow \text{num} \mid (S)$$



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

What happens if we run this grammar through LALR construction?

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

10

Conflict!

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

$$\begin{array}{|l} E \rightarrow E + E \cdot + \\ E \rightarrow E \cdot + E +, \$ \end{array}$$

shift/reduce conflict

$$1+2+3$$

^

shift: 1+(2+3)
reduce: (1+2)+3

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

11

Grammar in CUP

non terminal E; terminal PLUS, LPAREN...
precedence left PLUS;

“When shifting + conflicts with reducing a production containing +, choose reduce”

E ::= E PLUS E
| LPAREN E RPAREN
| NUMBER ;

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

12

Precedence

- Also can handle operator precedence

$$E \rightarrow E + E \mid T$$

$$T \rightarrow T \times T \mid \text{num} \mid (E)$$


$$E \rightarrow E + E \mid E \times E$$

$$\mid \text{num} \mid (E)$$

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

13

Conflicts w/o precedence

$$E \rightarrow E + E \mid E \times E$$

$$\mid \text{num} \mid (E)$$

$$E \rightarrow E \cdot + E \quad \dots$$

$$E \rightarrow E \times E \cdot \quad +$$

$$E \rightarrow E + E \cdot \quad \times$$

$$E \rightarrow E \cdot \times E \quad \dots$$

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

14

Precedence in CUP

precedence left PLUS;
 precedence left TIMES; // TIMES > PLUS
 $E ::= E \text{ PLUS } E \mid E \text{ TIMES } E \mid \dots$

$$E \rightarrow E \cdot + E \quad \dots$$

$$E \rightarrow E \times E \cdot \quad +$$

$$E \rightarrow E + E \cdot \quad \times$$

$$E \rightarrow E \cdot \times E \quad \dots$$

Rule: in conflict, choose **reduce** if production symbol higher precedence than shifted symbol; choose **shift** if vice-versa

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

15

Summary

- Look-ahead information makes SLR(1), LALR(1), LR(1) grammars expressive
- Automatic parser generators support LALR(1)
- Precedence, associativity declarations simplify grammar writing
- Easiest and best way to read structured human-readable input

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

16

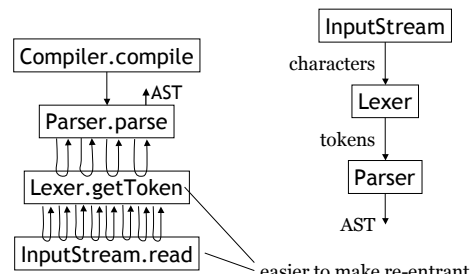
Compiler 'main program'

```
class Compiler {
    void compile() throws CompileError {
        Lexer l = new Lexer(input);
        Parser p = new Parser(l);
        AST tree = p.parse();
        // calls l.getToken() to read tokens
        if (typeCheck(tree))
            IR = genIntermediateCode(tree);
        IR.emitCode();
    }
}
```

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

17

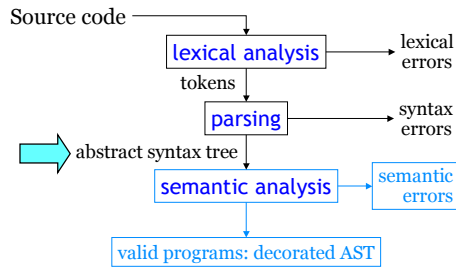
Thread of Control



CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

18

Semantic Analysis

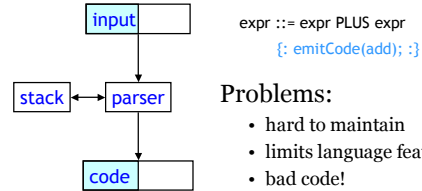


CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

19

Do we need an AST?

- Old-style compilers: semantic actions generate code during parsing!
- Especially for stack machine:



Problems:

- hard to maintain
- limits language features
- bad code!

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

20

AST

- **Abstract Syntax Tree** is a tree representation of the program. Used for
 - semantic analysis (type checking)
 - some optimization (*e.g.* constant folding)
 - intermediate code generation (sometimes intermediate code = AST with somewhat different set of nodes)
- Compiler phases = recursive tree traversals
- Object-oriented languages convenient for defining AST nodes

CS 412/413 Spring '01 Introduction to Compilers -- Andrew Myers

21