




复旦大学

编译原理/Compiler Parsing

周雅倩
复旦大学计算机科学技术学院
zhouyaqian@fudan.edu.cn
2018/10/11




复旦大学媒体计算研究所

References

- Partial contents are copied from the slides of Alex Aiken
 - <http://www.stanford.edu/class/cs143/>
 - <https://www.coursera.org/course/compilers>
- http://en.wikipedia.org/wiki/Comparison_of_parser_generators

2




复旦大学媒体计算研究所

Outline

- [Parser overview](#)
- [Context-free grammars \(CFG\)](#)
- [LL Parsing](#) (Left to right parse, Leftmost derivation)
- [LR Parsing](#) (Left to right parse, Rightmost derivation)

3




复旦大学媒体计算研究所

Languages and Automata

- Formal languages are very important in CS
 - Especially in programming languages
- Regular languages
 - The weakest formal languages widely used
 - Many applications
- We will also study context-free languages

4



复旦大学媒体计算研究所

Limitations of Regular Languages

- Intuition: A finite automaton that runs long enough must repeat states
- Finite automaton can't remember # of times it has visited a particular state
- Finite automaton has finite memory
 - Only enough to store in which state it is
 - Cannot count, except up to a finite limit
- E.g., language of balanced parentheses is not regular: $\{ (' \beta')^i \mid i \geq 0 \}$

5



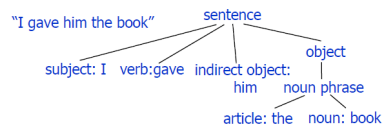
复旦大学媒体计算研究所

语法分析器简介 PARSER OVERVIEW

6

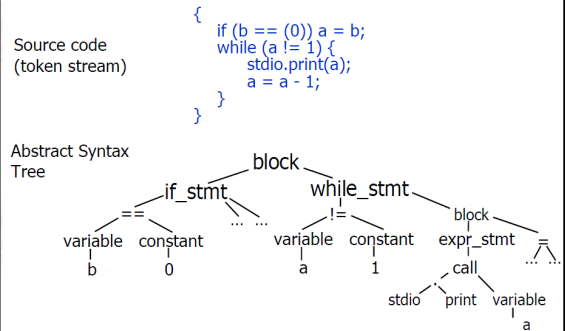
The Functionality of the Parser

- **Syntax analysis for natural languages:**
 - recognize whether a sentence is grammatically well-formed & identify the function of each component.



7

Syntax Analysis Example



Comparison with Lexical Analysis

Phase	Input	Output
Lexer	Sequence of characters	Sequence of tokens
Parser	Sequence of tokens	Parse tree

9

The Role of the Parser

- Not all sequences of tokens are programs ...
- ... Parser must distinguish between valid and invalid sequences of tokens
- We need
 - A language for describing valid sequences of tokens
 - A method for distinguishing valid from invalid sequences of tokens

10

上下文无关文法 CONTEXT-FREE GRAMMARS

11

Context-Free Grammars

- Programming language constructs **have recursive structure**
- An **exp** is
 - if (exp) {exp } else {exp } , or**
 - while (exp) {exp } , or**
 - ...
- Context-free grammars are a natural notation for this recursive structure

12

CFGs (Cont.)

- A CFG consists of
 - A set of terminals T
 - A set of non-terminals N
 - A start symbol S (a non-terminal)
 - A set of productions r
- Assuming $X \in N$
 - $X \rightarrow \varepsilon$, or
 - $X \rightarrow Y_1 Y_2 \dots Y_n$ where $Y_i \in N \cup T$

13

Terminals

- **Terminals** are called because there are no rules for replacing them
- Once generated, terminals are permanent
- Terminals ought to be tokens of the language

14

Notational Conventions (符号的使用约定)

- In these lecture notes
 - Non-terminals are written upper-case
 - Terminals are written lower-case
 - The start symbol is the left-hand side of the first production
- Productions
 - $A \rightarrow \alpha$, A left side, α right side
 - $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_k$
 - $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_k$

15

The Language of a CFG

- Read productions as replacement rules:
 - $X \rightarrow Y_1 \dots Y_n$
means X can be replaced by $Y_1 \dots Y_n$
- $X \rightarrow \varepsilon$
means X can be erased (replaced with empty string)

16

Key Ideas

1. Begin with a string consisting of the start symbol " S "
2. Replace any non-terminal X in the string by a right-hand side of some production
 - $X \rightarrow Y_1 \dots Y_n$
3. Repeat (2) until there are no non-terminals in the string

17

The Language of a CFG (Cont.)

More formally, write

$$X_1 \dots X_i \dots X_n \rightarrow X_1 \dots X_{i-1} Y_1 \dots Y_m X_{i+1} \dots X_n$$

if there is a production

$$X_i \rightarrow Y_1 \dots Y_m$$

18

Examples of CFGs

Simple arithmetic expressions:

```

E → E * E
   | E + E
   | (E)
   | id
  
```

19

The Language of a CFG (Cont.)

Write

$$X_1 \dots X_n \quad * \rightarrow \quad Y_1 \dots Y_m$$

if

$$X_1 \dots X_n \rightarrow \dots \rightarrow Y_1 \dots Y_m$$

in 0 or more steps

20

The Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language of G is:

$$\{ a_1 \dots a_n \mid S \xrightarrow{*} a_1 \dots a_n \text{ and every } a_i \text{ is a terminal} \}$$

21

Examples

- $L(G)$ is the language of CFG G
 $\{ (i)^i \mid i \geq 0 \}$
- Strings of balanced parentheses
- Grammar:
 - $S \rightarrow (S)$
 - $S \rightarrow \varepsilon$

22

Example

```

exp → if (exp ) {exp} else {exp}
      | while (exp) {exp}
      | id
  
```

23

Example (Cont.)

- Some elements of the language
 - id
 - if (id) {id} else {id}
 - while(id) {id}
 - if(while(id) {id}){id}else{id}
 - while(id){while(id){id}}
 - ...

24

Arithmetic Example

- Simple arithmetic expressions:

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

- Some elements of the language:

int		int+int
(int)		int*int
(int)*int		int*(int)

25

Notes

The idea of a CFG is a big step. But:

- Membership in a language is "yes" or "no"
 - we also need parse tree of the input
- Must handle errors gracefully
- Need an implementation of CFG's (e.g., bison)

26

More Notes

- Form of the grammar is important
 - Many grammars generate the same language
 - Tools are sensitive to the grammar
- Note: Tools for regular languages (e.g., flex) are also sensitive to the form of the regular expression, but this is rarely a problem in practice

27

Derivations and Parse Trees

A **derivation**(推导) is a sequence of productions

$$S \rightarrow \dots \rightarrow \dots$$

A derivation can be drawn as a tree

- Start symbol is the tree's root
- For a production $X \rightarrow Y_1 \dots Y_n$ add children Y_1, \dots, Y_n to node X

28

Derivation Example

- Grammar

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

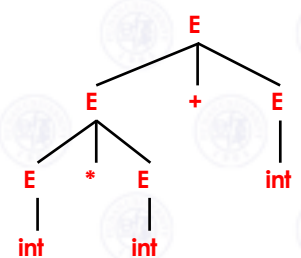
- String

int*int+int

29

Derivation Example (Cont.)

E
 $\rightarrow E+E$
 $\rightarrow E * E+E$
 $\rightarrow \text{int} * E+E$
 $\rightarrow \text{int} * \text{int} + E$
 $\rightarrow \text{int} * \text{int} + \text{int}$



30

Notes on Derivations

- A **parse tree** has
 - Terminals at the leaves
 - Non-terminals at the interior nodes
- An in-order traversal of the leaves is the original input
- The parse tree shows the association of operations, the input string does not

31

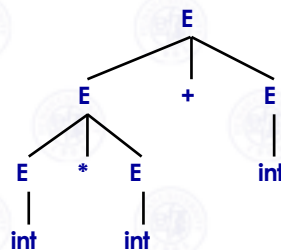
Left-most and Right-most Derivations

- Two choices to be made in each step in a derivation
 - Which nonterminal to replace
 - Which alternative to use for that nonterminal
- The example is a left-most derivation
 - At each step, replace the **left-most non-terminal**
- There is an equivalent notion of a right-most derivation.
 - At each step, replace the **right-most non-terminal**

32

Right-most Derivation in Detail

E
 $\rightarrow E+E$
 $\rightarrow E+int$
 $\rightarrow E*E+int$
 $\rightarrow E*int+int$
 $\rightarrow int*int+int$



33

Derivations and Parse Trees

- Note that right-most and left-most derivations have the same parse tree
- The difference is the order in which branches are added

34

Ambiguity

- Grammar

$$E \rightarrow E + E \mid E * E \mid (E) \mid int$$
- Strings

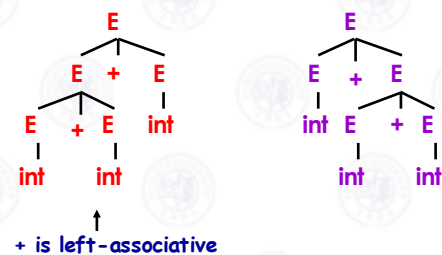
$$int + int + int$$

$$int * int + int$$

35

Ambiguity. Example

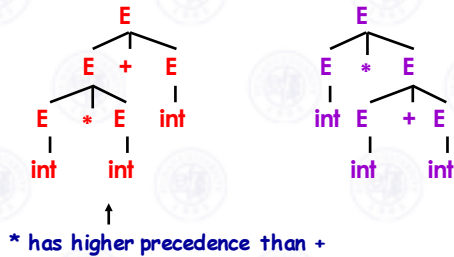
This string has two parse trees



36

Ambiguity. Example

This string has two parse trees



37

Ambiguous Grammar

- A grammar is **ambiguous** if it has more than one parse tree for some string
- Equivalently, there is more than one right-most or left-most derivation for some string

38

Ambiguity

- Ambiguity is **bad**
 - Leaves meaning of some programs ill-defined
- Ambiguity is **common** in programming languages
 - Arithmetic expressions
 - IF-THEN-ELSE

39

Dealing with Ambiguity

- There are several ways to handle ambiguity
- Most direct method is to rewrite the grammar unambiguously

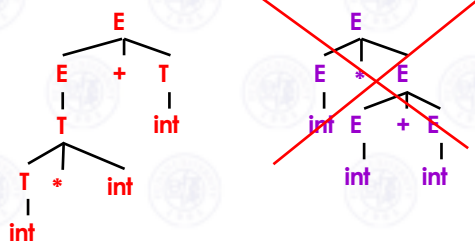
$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * \text{int} \mid \text{int} \mid (E)$$
- Enforces precedence (优先级) of ***** over **+**
- Enforces left-associativity (左结合) of **+** and *****

40

Ambiguity. Example

The `int * int + int` has only one parse tree now



41

Ambiguity: The Dangling Else

- Consider the grammar

$$E \rightarrow \text{if } E \text{ then } E$$

$$\quad \mid \text{if } E \text{ then } E \text{ else } E$$

$$\quad \mid \text{OTHER}$$
- This grammar is also ambiguous

42

The Dangling Else: Example

- The expression

if E_1 then if E_2 then E_3 else E_4

has two parse trees



- Typically we want the second form

43

The Dangling Else: A Fix

- else** matches the closest unmatched **then**
- We can describe this in the grammar (distinguish between matched and unmatched "then")

$E \rightarrow \text{MIF} \quad /* \text{all then are matched} */$
 $\quad \quad \quad \text{UIF} \quad /* \text{some then are unmatched} */$

$\text{MIF} \rightarrow \text{if } E \text{ then MIF else MIF}$

$\quad \quad \quad \text{OTHER}$

Describes the same set of strings

$\text{UIF} \rightarrow \text{if } E \text{ then } E$

$\quad \quad \quad \text{if } E \text{ then MIF else UIF}$

44

The Dangling Else: Example Revisited

- The expression if E_1 then if E_2 then E_3 else E_4



A valid parse tree
(for a UIF)

- Not valid because the then expression is not a MIF

45

Ambiguity

- No general techniques for handling ambiguity
- Impossible to convert automatically an ambiguous grammar to an unambiguous one
- Used with care, ambiguity can simplify the grammar
 - Sometimes allows more natural definitions
 - We need disambiguation mechanisms

46

LL(1) PARSING

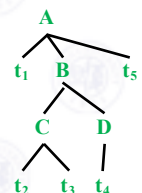
47

Intro to Top-Down Parsing

- Terminals are seen in order of appearance in the token stream:

$t_1 \ t_2 \ t_3 \ t_4 \ t_5$

- The parse tree is constructed
 - From the top
 - From left to right



48

Recursive Descent Parsing

- Consider the grammar

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

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Token stream is: $\text{int} * \text{int}$
- Start with top-level non-terminal E
- Try the rules for E in order

Different to previous grammar? Why?
Left Recursion

49

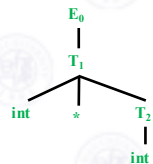
Recursive Descent Parsing. Example (Cont.)

- Try $E_0 \rightarrow T_1 + E_2$
- Then try a rule for $T_1 \rightarrow (E_3)$
But $($ does not match input token int_5
- Try $T_1 \rightarrow \text{int}$. Token matches.
But $+$ after T_1 does not match input token $*$
- Try $T_1 \rightarrow \text{int} * T_2$
This will match but $+$ after T_1 will be unmatched
- Have exhausted the choices for T_1
Backtrack to choice for E_0

50

Recursive Descent Parsing. Example (Cont.)

- Try $E_0 \rightarrow T_1$
- Follow same steps as before for T_1
And succeed with $T_1 \rightarrow \text{int} * T_2$ and $T_2 \rightarrow \text{int}$
With the following parse tree



51

Recursive Descent Parsing. Notes.

- Easy to implement by hand
- But does not always work ...

52

Recursive-Descent Parsing

- Parsing:** given a string of tokens $t_1 t_2 \dots t_n$, find its parse tree
- Recursive-descent parsing:** Try all the productions exhaustively
 - At a given moment the fringe of the parse tree is: $t_1 t_2 \dots t_k A \dots$
 - Try all the productions for A : if $A \rightarrow BC$ is a production, the new fringe is $t_1 t_2 \dots t_k B C \dots$
 - Backtrack when the fringe doesn't match the string
 - Stop when there are no more non-terminals

53

When Recursive Descent Does Not Work

- Consider a production $S \rightarrow S \alpha$:
 - In the process of parsing S we try the above rule
 - What goes wrong?
- A **left-recursive grammar** has a non-terminal S

$$S \rightarrow^+ S \alpha \text{ for some } \alpha$$
- Recursive descent does not work in such cases
 - It goes into an infinite loop

54

Elimination of Left Recursion

- Consider the left-recursive grammar

$$S \rightarrow S\alpha \mid \beta$$
- S generates all strings starting with a β and followed by a number of α
- Can rewrite using right-recursion

$$S \rightarrow \beta S'$$

$$S' \rightarrow \alpha S' \mid \epsilon$$

55

Elimination of Left-Recursion. Example

- Consider the grammar

$$S \rightarrow 1 \mid S0 \quad (\beta = 1 \text{ and } \alpha = 0)$$
- can be rewritten as

$$S \rightarrow 1S'$$

$$S' \rightarrow 0S' \mid \epsilon$$

56

More Elimination of Left-Recursion

- In general

$$S \rightarrow S\alpha_1 \mid \dots \mid S\alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$
- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$
- Rewrite as

$$S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$$

$$S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \epsilon$$

57

Summary of Recursive Descent

- Simple and general parsing strategy
 - Left-recursion must be eliminated first
 - ... but that can be done automatically
- Unpopular because of backtracking
 - Thought to be too inefficient
- In practice, backtracking is eliminated by restricting the grammar

58

Predictive Parsers

- Like recursive-descent but parser can "predict" which production to use
 - By looking at the next few tokens
 - No backtracking
- Predictive parsers accept LL(k) grammars
 - L means "left-to-right" scan of input
 - L means "leftmost derivation"
 - k means "predict based on k tokens of lookahead"
- In practice, LL(1) is used

59

LL(1) Languages

- In recursive-descent, for each non-terminal and input token there may be a choice of production
- LL(1) means that for each non-terminal and token there is only one production that could lead to success
- Can be specified as a 2D table
 - One dimension for current non-terminal to expand
 - One dimension for next token
 - A table entry contains one production

60

Predictive Parsing and Left Factoring

- Recall the grammar

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

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Impossible to predict because
 - For T two productions start with int
 - For E it is not clear how to predict
- A grammar must be **left-factored** before use for predictive parsing

61

Left-Factoring Example

- Recall the grammar

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

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Factor out common prefixes of productions

$$E \rightarrow T X$$

$$X \rightarrow + E \mid \epsilon$$

$$T \rightarrow (E) \mid \text{int} Y$$

$$Y \rightarrow * T \mid \epsilon$$

62

LL(1) Parsing Table Example

- Left-factored grammar

$$E \rightarrow T X \quad X \rightarrow + E \mid \epsilon$$

$$T \rightarrow (E) \mid \text{int} Y \quad Y \rightarrow * T \mid \epsilon$$
- The LL(1) parsing table:

	int	*	+	()	\$
T	int Y			(E)		
E	T X			T X		
X			+ E		ϵ	ϵ
Y		* T	ϵ		ϵ	ϵ

63

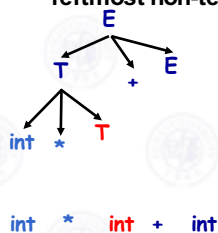
Constructing Parsing Tables

- We want to generate parsing tables from CFG

64

Top-Down Parsing. Review

- Top-down parsing expands a parse tree from the start symbol to the leaves with leftmost non-terminal.

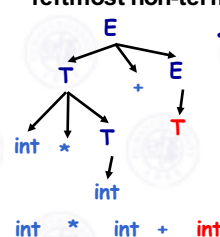


- The leaves at any point form a string $\beta A \gamma$
 - β contains only terminals
 - The input string is $\beta b \delta$
 - The prefix β matches
 - The next token is b

65

Top-Down Parsing. Review

- Top-down parsing expands a parse tree from the start symbol to the leaves with leftmost non-terminal.



- The leaves at any point form a string $\beta A \gamma$
 - β contains only terminals
 - The input string is $\beta b \delta$
 - The prefix β matches
 - The next token is b

66

Predictive Parsing. Review.

- A predictive parser is described by a table
 - For each non-terminal A and for each token b we specify a production $A \rightarrow \alpha$
 - when trying to expand A we use $A \rightarrow \alpha$ if b follows next
- Once we have the table
 - The parsing algorithm is simple and fast
 - No backtracking is necessary

67

Constructing Predictive Parsing Tables

- Consider the state $S \rightarrow^* \beta A \gamma$
 - With b the next token
 - Trying to match $\beta b \delta$
- There are two possibilities:
1. b belongs to an expansion of A
 - Any $A \rightarrow \alpha$ can be used if b can start a string derived from α
 - In this case we say that $b \in \text{First}(\alpha)$

Or...

68

Constructing Predictive Parsing Tables (Cont.)

2. b does not belong to an expansion of A
 - The expansion of A is empty and b belongs to an expansion of γ
 - Means that b can appear after A in a derivation of the form $S \rightarrow^* \beta A b \omega$
 - We say that $b \in \text{Follow}(A)$ in this case
 - What productions can we use in this case?
 - Any $A \rightarrow \alpha$ can be used if α can expand to ϵ
 - We say that $\epsilon \in \text{First}(A)$ in this case

69

Computing First Sets

Definition

$$\text{First}(X) = \{ b \mid X \rightarrow^* b \alpha \} \cup \{ \epsilon \mid X \rightarrow^* \epsilon \}$$

1. $\text{First}(b) = \{ b \}$
2. For all productions $X \rightarrow A_1 \dots A_n$
 - Add $\text{First}(A_1) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_1)$
 - Add $\text{First}(A_2) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_2)$
 - ...
 - Add $\text{First}(A_n) - \{ \epsilon \}$ to $\text{First}(X)$. Stop if $\epsilon \notin \text{First}(A_n)$
 - Add ϵ to $\text{First}(X)$

70

First Sets. Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow T X & X \rightarrow + E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow * T \mid \epsilon \end{array}$$
- First sets

$\text{First}(() = \{ (\}$	$\text{First}(T) = \{ \text{int}, (\}$
$\text{First}() = \{ \}$	$\text{First}(E) = \{ \text{int}, (\}$
$\text{First}(\text{int}) = \{ \text{int} \}$	$\text{First}(X) = \{ +, \epsilon \}$
$\text{First}(+) = \{ + \}$	$\text{First}(Y) = \{ *, \epsilon \}$
$\text{First}(*) = \{ * \}$	

71

Computing Follow Sets

Definition

$$\text{Follow}(X) = \{ b \mid S \rightarrow^* \beta X b \delta \}$$

1. Compute the First sets for all non-terminals first
2. Add $\$$ to $\text{Follow}(S)$ (if S is the start non-terminal)
3. For all productions $Y \rightarrow \dots X A_1 \dots A_n$
 - Add $\text{First}(A_1) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_1)$
 - Add $\text{First}(A_2) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_2)$
 - ...
 - Add $\text{First}(A_n) - \{ \epsilon \}$ to $\text{Follow}(X)$. Stop if $\epsilon \notin \text{First}(A_n)$
 - Add $\text{Follow}(Y)$ to $\text{Follow}(X)$

72

Follow Sets Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow TX & X \rightarrow +E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow *T \mid \epsilon \end{array}$$

- Follow sets

$$\begin{array}{ll} \text{Follow}(+) = \{ \text{int}, (\} & \text{Follow}(*) = \{ \text{int}, (\} \\ \text{Follow}() = \{ \text{int}, (\} & \text{Follow}(E) = \{ \}, \$ \} \\ \text{Follow}(X) = \{ \$,) \} & \text{Follow}(T) = \{ +,) , \$ \} \\ \text{Follow}() = \{ +,) , \$ \} & \text{Follow}(Y) = \{ +,) , \$ \} \\ \text{Follow}(\text{int}) = \{ *, +,) , \$ \} & \end{array}$$

73

Constructing LL(1) Parsing Tables

- Construct a parsing table T for CFG G

- For each production $A \rightarrow \alpha$ in G do:

- For each terminal $b \in \text{First}(\alpha)$ do
 - $T[A, b] = \alpha$
- If $\alpha \rightarrow^* \epsilon$, for each $b \in \text{Follow}(A)$ do
 - $T[A, b] = \alpha$
- If $\alpha \rightarrow^* \epsilon$ and $\$ \in \text{Follow}(A)$ do
 - $T[A, \$] = \alpha$

74

Constructing LL(1) Tables Example

- Recall the grammar

$$\begin{array}{ll} E \rightarrow TX & X \rightarrow +E \mid \epsilon \\ T \rightarrow (E) \mid \text{int } Y & Y \rightarrow *T \mid \epsilon \end{array}$$

- Where in the line of Y we put $Y \rightarrow *T$?

- In the lines of $\text{First}(*T) = \{ * \}$

- Where in the line of Y we put $Y \rightarrow \epsilon$?

- In the lines of $\text{Follow}(Y) = \{ \$, +,) \}$

75

Using Parsing Tables

- Method similar to recursive descent, except

- For each non-terminal S
- We look at the next token a
- And choose the production shown at $[S, a]$

- We use a stack to keep track of pending non-terminals

- We reject when we encounter an error state
- We accept when we encounter end-of-input

76

LL(1) Parsing Algorithm

initialize $\text{stack} = \langle S \$ \rangle$ and next (pointer to tokens)

repeat

case stack of

$\langle X, \text{rest} \rangle$: if $T[X, \text{next}] = Y_1 \dots Y_n$
 then $\text{stack} = \langle Y_1 \dots Y_n \text{rest} \rangle$;
 else error ();

$\langle t, \text{rest} \rangle$: if $t == \text{next} ++$
 then $\text{stack} = \langle \text{rest} \rangle$;
 else error ();

until $\text{stack} == \langle \rangle$

77

LL(1) Parsing Example (Table)

Stack	Input	Action
E \$	int * int \$	TX
TX \$	int * int \$	int Y
int YX \$	int * int \$	terminal
YX \$	* int \$	* T
*TX \$	* int \$	terminal
TX \$	int \$	int Y
int YX \$	int \$	terminal
YX \$	\$	ϵ
X \$	\$	ϵ
\$	\$	ACCEPT

78

Notes on LL(1) Parsing Tables

- LL(1) languages are those defined by a parsing table for the LL(1) algorithm
- If any entry is multiply defined then G is not LL(1)
 - If G is ambiguous
 - If G is left recursive
 - If G is not left-factored
 - And in other cases as well
- Most programming language grammars are not LL(1)
- There are tools that build LL(1) tables

79

LR PARSING

80

Bottom-up parsing

- Bottom-up parsing is more general than top-down parsing
 - Builds on ideas in top-down parsing
 - Preferred method in practice
- also called as LR Parsing
 - L means that read tokens from **left to right**
 - R means that it constructs a **rightmost derivation**
 - LR(0), LR(1), SLR(1), LALR(1)

81

The Idea

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

str ← input string of terminals

repeat

Identify β in **str** such that $A \rightarrow \beta$ is a production
(i.e., $\text{str} = \alpha \beta \gamma$)

Replace β by **A** in **str** (i.e., **str** becomes $\alpha A \gamma$)

until str = S

82

An Introductory Example

- LR parsers don't need left-factored grammars and can also handle left-recursive grammars
- Consider the following grammar:

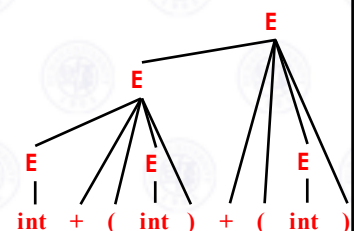
$$E \rightarrow E + (E) \mid \text{int}$$
 - Why is this not LL(1)?
- Consider the string: **int + (int) + (int)**

83

A Bottom-up Parse in Detail

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

A rightmost derivation in reverse



84

Where Do Reductions Happen

- Some characters:
 - Let $\alpha\beta\gamma$ be a step of a bottom-up parse
 - Assume the next reduction is by $A \rightarrow \beta$
 - Then γ is a string of terminals!
- Why?
 - Because $\alpha A \gamma \rightarrow \alpha\beta\gamma$ is a step in a right-most derivation

85

Top-down vs. Bottom-up

- Bottom-up: Don't need to figure out as much of the parse tree for a given

Top-down Bottom-up

86

Notation

- Idea: Split 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 **I**
 - The **I** is not part of the string
- Initially, all input is unexamined:

$$I x_1 x_2 \dots x_n$$

87

Shift-Reduce Parsing

- Bottom-up parsing uses only three kinds of actions:

Shift
Reduce
Accept

88

Shift

Shift: Move **I** one place to the right

- Shifts a terminal to the left string

$$E + (I \text{ int}) \Rightarrow E + (\text{int } I)$$

89

Reduce

Reduce: Apply an inverse production at the right end of the left string

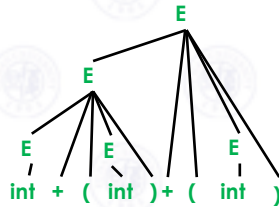
- If $E \rightarrow E + (E)$ is a production, then

$$E + (\underline{E + (E)} I) \Rightarrow E + (E I)$$

90

Shift-Reduce Example

i int + (int) + (int)\$ shift
 int i + (int) + (int)\$ red. $E \rightarrow \text{int}$
 E i + (int) + (int)\$ shift 3 times
 E + (int i) + (int)\$ red. $E \rightarrow \text{int}$
 E + (E i) + (int)\$ shift
 E + (E i) + (int)\$ red. $E \rightarrow E + (E)$
 E i + (int)\$ shift 3 times
 E + (int i)\$ red. $E \rightarrow \text{int}$
 E + (E i)\$ shift
 E + (E i)\$ red. $E \rightarrow E + (E)$
 E i\$ accept



91

Key Issue: When to Shift or Reduce?

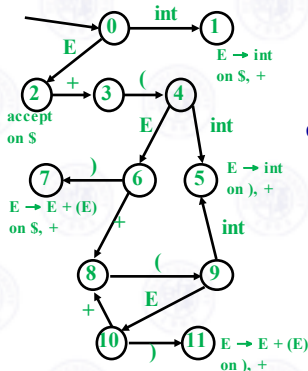
- Decide based on the left string (the stack)
- Idea: use a finite automaton (DFA) to decide when to shift or reduce
 - The DFA input is the stack
 - The language consists of terminals and non-terminals

92

LR Parsing Engine

- Basic mechanism:
 - Use a set of parser states
 - Use a stack of states
 - Use a parsing table to:
 - Determine what action to apply (shift/reduce)
 - Determine the next state
- The parser actions can be precisely determined from the table

93



i int + (int) + (int)\$ shift
 int i + (int) + (int)\$ $E \rightarrow \text{int}$
 E i + (int) + (int)\$ shift(x3)
 E + (int i) + (int)\$ $E \rightarrow \text{int}$
 E + (E i) + (int)\$ shift
 E + (E i) + (int)\$ $E \rightarrow E + (E)$
 E i + (int)\$ shift(x3)
 E + (int i)\$ $E \rightarrow \text{int}$
 E + (E i)\$ shift
 E + (E i)\$ $E \rightarrow E + (E)$
 E i\$ accept

94

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 non-terminals
- Typically columns are split into:
 - Those for terminals: action table
 - Those for non-terminals: goto table

95

The LR Parsing Table

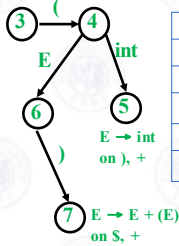
	Terminals	Non-terminals
State	Next action and next state	Next state

Algorithm:
 look at entry for current state S and input terminal c
 if $\text{Table}[S, c] = s(S')$ then **shift**:
 push(S')
 if $\text{Table}[S, c] = A \rightarrow \alpha$ then **reduce**:
 pop($|\alpha|$); $S' = \text{top}()$; push(A); push($\text{Table}[S', A]$)

96

Representing the DFA. Example

- The table for a fragment of our DFA:



	int	+	()	\$	E
3				s4		
4	s5					g6
5		rE→ int		rE→ int		
6	s8			s7		
7		rE→ E+(E)			rE→ E+(E)	

97

The LR Parsing Algorithm

- After a shift or reduce action we return 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
 - $\langle \text{sym}_1, \text{state}_1 \rangle \dots \langle \text{sym}_n, \text{state}_n \rangle$
 - state_k is the final state of the DFA on $\text{sym}_1 \dots \text{sym}_k$

98

LR Parsing Notes

- Can be used to parse more grammars than LL
- Most programming languages grammars are LR
- Can be described as a simple table
- There are many tools for building the table
- How is the table constructed?

99

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 a **E + (E)** rhs

100

LR(0) Items

- The LR(0) item of a grammar G
 - For example: $A \rightarrow XYZ$
 - $A \rightarrow \bullet XYZ$
 - $A \rightarrow X \bullet YZ$
 - $A \rightarrow XY \bullet Z$
 - $A \rightarrow XYZ \bullet$
- The parser begins with $\bullet SS$ (state 0).

101

The Closure Operation

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

Closure(I) =
 repeat
 for each $[X \rightarrow a \bullet Yb]$ in I
 for each production $Y \rightarrow g$
 add $[Y \rightarrow \bullet g]$ to Items
 until I is unchanged
 return I

102

Goto Operation

- $\text{Goto}(I, X)$ is defined to be the closure of the set of all items $[A \rightarrow aX\beta]$ such that $[A \rightarrow a\bullet X\beta]$ is in I .

I is the set of items
 X is a grammar symbol
 $\text{Goto}(I, X) =$
 set J to the empty set
 for any item $A \rightarrow a\bullet X\beta$ in I
 add $A \rightarrow aX\bullet\beta$ to J
 return $\text{Closure}(J)$

103

The Sets-of-items Construction

Initialize $T = \{\text{closure}(\{[S' \rightarrow \bullet S\$]\})\};$
 Initialize $E = \{\};$
 repeat
 for each state (a set of items) I in T
 let J be $\text{Goto}(I, X)$
 $T = T \cup \{J\};$
 $E = E \cup \{I \rightarrow X \rightarrow J\};$
 until E and T did not change

For the symbol $\$$ we do not compute $\text{Goto}(I, \$)$.

104

LR Parser

- To be add.....

2018/10/11

Example(1)

- Grammar G_1 :

- $S' \rightarrow S\$$
- $S \rightarrow (L)$
- $S \rightarrow x$
- $L \rightarrow S$
- $L \rightarrow L, S$

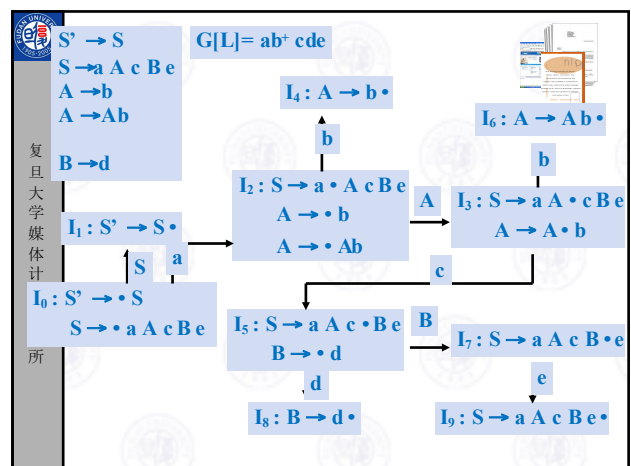
106

Example(2)

- Grammar G_2 :

$S \rightarrow a A c B e$
 $A \rightarrow b$
 $A \rightarrow Ab$
 $B \rightarrow d$

- abbcde
- abbce



LR(0) Parsing Table

		ACTION						GOTO		
		a	c	e	b	d	#	S	A	B
0	s ₂							1		
1							acc			
2					s ₄				3	
3		s ₅		s ₆						
4	r ₂	r ₂	r ₂	r ₂	r ₂	r ₂				
5					s ₈					7
6	r ₃	r ₃	r ₃	r ₃	r ₃	r ₃				
7			s ₉							
8	r ₄	r ₄	r ₄	r ₄	r ₄	r ₄				
9	r ₁	r ₁	r ₁	r ₁	r ₁	r ₁				

Parsing the String abcde##

Step	states.	Syms.	The rest of input	action	goto
1	0	#	abcde\$	s ₂	
2	02	#a	bbce\$	s ₄	
3	024	#ab	bce\$	r ₂	3
4	023	#aA	bcde\$	s ₆	
5	0236	#aAb	cde\$	r ₃	3
6	023	#aA	cde\$	s ₅	
7	0235	#aAc	de\$	s ₈	
8	02358	#aAcd	e\$	r ₄	7
9	02357	#aAcB	e\$	s ₉	
10	023579	#aAcBe	\$	r ₁	1
11	01	#S	\$	acc	

Parsing the String abbce#

Step	states.	Syms.	The rest of input	action
goto				
1	0	#	abbce#	s ₂
2	02	#a	bbce#	s ₄
3	024	#ab	bce#	r ₂
4	023	#aA	bce#	s ₆
5	0236	#aAb	ce#	r ₃
6	023	#aA	ce#	s ₅
7	0235	#aAc	e#	error

Example(3)

■ Grammar G3:

1. $S \rightarrow E\$$
2. $E \rightarrow T + E$
3. $E \rightarrow T$
4. $T \rightarrow x$

■ G3 is not LR(0).

- We can extend LR(0) in a simple way.
- **SLR(1)**: simple LR
- Reduce actions are indicated by **Follow set**.

SLR(1)

$R = \{ \}$;
 for each state I in T
 for each item $A \rightarrow \alpha \cdot$ in I
 for each token a in $\text{Follow}(A)$
 $R = R \cup \{ (I, a, A \rightarrow \alpha) \}$

$(I, a, A \rightarrow \alpha)$ indicates that in state I , on lookahead symbol a , the parser will reduce by rule $A \rightarrow \alpha$.

LR(0) Limitations

- An LR(0) machine only works if states with reduce actions have a single reduce action
- With more complex grammar, construction gives states with shift/reduce or reduce/reduce conflicts
- Need to use look-ahead to choose

LR(1) Items

- An **LR(1) item** is a pair:
 - $X \rightarrow \alpha \cdot \beta, a$
 - $X \rightarrow \alpha \beta$ is a production
 - a is a terminal (the lookahead terminal)
 - LR(1) means 1 lookahead terminal
- $[X \rightarrow \alpha \cdot \beta, a]$ 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

115

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
- For grammar:
 - $E \rightarrow \text{int}$ and $E \rightarrow E + (E)$
- The initial parsing context contains:
 - $S \rightarrow \cdot E \$, ?$
 - Trying to find an S as a string derived from $E \$$
 - The stack is empty

116

LR(1) Items (Cont.)

- In context containing
 - $E \rightarrow E + \cdot (E), +$
 - If $($ follows then we can perform a shift to context containing
 - $E \rightarrow E + (\cdot E), +$
- In context containing
 - $E \rightarrow E + (E) \cdot , +$
 - We can perform a reduction with $E \rightarrow E + (E)$
 - But only if a $+$ follows

117

LR(1) Items (Cont.)

- Consider the item
 - $E \rightarrow E + (\cdot E), +$
- We expect a string derived from E
- We describe this by extending the context with two more items:
 - $E \rightarrow \cdot \text{int},)$
 - $E \rightarrow \cdot E + (E),)$

118

The Closure Operation

The operation of extending the context with items

```

Closure(I) =
repeat
  for each  $[A \rightarrow \alpha \cdot X \beta, z]$  in  $I$ 
    for each production  $X \rightarrow \gamma$ 
      for each  $w \in \text{First}(\beta z)$ 
        add  $[X \rightarrow \cdot \gamma, w]$  to  $I$ 
until  $I$  is unchanged
return  $I$ 

```

119

The Goto Operation

```

Goto(I, X) =
J = {};
for each item  $(A \rightarrow \alpha \cdot X \beta, z)$  in  $I$ 
  add  $(A \rightarrow \alpha X \cdot \beta, z)$  to  $J$ ;
return Closure(J);

```

120

Example(1)

■ Grammar:

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

■ Construct the start context: Closure($\{S \rightarrow \bullet E \$, ?\}$)

$$S \rightarrow \bullet E \$, ?$$

$$E \rightarrow \bullet E + (E) \$$$

$$E \rightarrow \bullet \text{int}, \$$$

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

$$E \rightarrow \bullet \text{int}, +$$



$$S \rightarrow \bullet E \$, ?$$

$$E \rightarrow \bullet E + (E), \$/+$$

$$E \rightarrow \bullet \text{int}, \$/+$$

121

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?¹²²

Shift/Reduce Conflicts

- If a DFA state contains both
 $[X \rightarrow \alpha \bullet a \beta, b]$ and $[Y \rightarrow \gamma \bullet, a]$

- Then on input "a" we could either
 - Shift into state $[X \rightarrow \alpha a \bullet \beta, b]$, or
 - Reduce with $Y \rightarrow \gamma$

- This is called a shift-reduce conflict

123

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
 $[S \rightarrow \text{if } E \text{ then } S \bullet, \text{ else}]$
 $[S \rightarrow \text{if } E \text{ then } S \bullet \text{ else } S, x]$
- If **else** follows then we can shift or reduce
- Default (bison) is to shift
 - Default behavior is as needed in this case

124

More Shift/Reduce Conflicts

- Consider the ambiguous grammar

$$E \rightarrow E + E \mid E * E \mid \text{int}$$

- We will have the states containing

$$[E \rightarrow E * \bullet E, +] \quad [E \rightarrow E * E \bullet, +]$$

$$[E \rightarrow \bullet E + E, +] \Rightarrow^E [E \rightarrow E \bullet + E, +]$$

...

...

- Again we have a shift/reduce on input +
 - We need to reduce ($*$ binds more tightly than $+$)
 - Recall solution: declare the precedence of $*$ and $+$

125

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

126

Reduce/Reduce Conflicts

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

127

More on Reduce/Reduce Conflicts

- Consider the states:

$[S' \rightarrow \cdot S, \$]$	\xrightarrow{id}	$[S \rightarrow id \cdot, \$]$
$[S' \rightarrow \cdot, \$]$		$[S \rightarrow id \cdot S, \$]$
$[S' \rightarrow \cdot id, \$]$		$[S \rightarrow \cdot, \$]$
$[S' \rightarrow \cdot id S, \$]$		$[S \rightarrow \cdot id, \$]$
		$[S \rightarrow \cdot id S, \$]$
- Reduce conflict on input $\$$
 $S' \rightarrow S \rightarrow id$
 $S' \rightarrow S \rightarrow id S \rightarrow id$
- Better rewrite the grammar: $S \rightarrow \varepsilon \mid id S$

128

LR(1) but not SLR(1)

- Let G have productions
 $S \rightarrow aAb \mid Ac$
 $A \rightarrow a \mid \varepsilon$
 - $V(a) = \{$
 $[S \rightarrow a.Ab]$
 $[A \rightarrow a.]$
 $[A \rightarrow \cdot a]$
 $[A \rightarrow \cdot]$
 $\}$
- $FOLLOW(A) = \{b, c\}$
- reduce-reduce conflict

129

LR(1) Parsing Tables are Big

- But many states are similar, e.g.

1 $E \rightarrow int \cdot, \$ / +$	$E \rightarrow int$ on $\$, +$	and	5 $E \rightarrow int \cdot,) / +$	$E \rightarrow int$ on $), +$
--	-----------------------------------	-----	---------------------------------------	----------------------------------
- Idea: merge the DFA states whose items differ only in the lookahead tokens
 – We say that such states have the same core
- We obtain

1' $E \rightarrow int \cdot, \$ / + /)$	$E \rightarrow int$ on $\$, +,)$
---	--------------------------------------

130

The Core of a Set of LR Items

- Definition: The core of a set of LR items is the set of first components
 – Without the lookahead terminals
- Example: the core of
 $\{ [X \rightarrow \alpha \cdot \beta, b], [Y \rightarrow \gamma \cdot \delta, d] \}$
 is
 $\{ X \rightarrow \alpha \cdot \beta, Y \rightarrow \gamma \cdot \delta \}$

131

LALR States

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

132

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

Diagram illustrating the merging process:

```

  graph LR
    A((A)) --> B((B))
    C((C)) --> B
    D((D)) --> E((E))
    E --> F((F))
    B --> BE((BE))
    BE --> F
    F --> BE
  
```

Point edges from predecessors to new state
New state points to all the previous successors

133

Conversion LR(1) to LALR(1) Example

Diagram illustrating the conversion of LR(1) to LALR(1) by merging states with the same core.

LR(1) states (left): 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11. Transitions include 'int', '+', '(', ')', and 'E'.

LALR(1) states (right): 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11. Transitions are simplified by merging states with the same core (e.g., 3 and 8, 4 and 9, 6 and 10, 7 and 11).

134

The LALR Parser Can Have Conflicts

- Consider for example the LR(1) states
 - $\{[X \rightarrow \alpha \cdot, a], [Y \rightarrow \beta \cdot, b]\}$
 - $\{[X \rightarrow \alpha \cdot, b], [Y \rightarrow \beta \cdot, a]\}$
- And the merged LALR(1) state
 - $\{[X \rightarrow \alpha \cdot, a/b], [Y \rightarrow \beta \cdot, a/b]\}$
- Has a new reduce-reduce conflict
- In practice such cases are rare

135

LALR vs. LR Parsing

- LALR languages are not natural
 - They are an efficiency hack on LR languages
- Any reasonable programming language has a LALR(1) grammar
- LALR(1) has become a standard for programming languages and for parser generators
- LALR parsers are automatically generated by compiler compilers such as Yacc and GNU Bison.

136

A Hierarchy of Grammar Classes

Diagram illustrating the hierarchy of grammar classes:

```

  graph TD
    subgraph Unambiguous_Grammars [Unambiguous Grammars]
      LLk[LL(k)]
      LRk[LR(k)]
      LL1[LL(1)]
      LR1[LR(1)]
      LALR1[LALR(1)]
      SLR[SLR]
      LL0[LL(0)]
      LR0[LR(0)]
    end
    subgraph Ambiguous_Grammars [Ambiguous Grammars]
      LR0
    end
    LLk --> LL1
    LRk --> LR1
    LL1 --> LL0
    LR1 --> LR0
    LALR1 --> LR1
    SLR --> LR1
  
```

137

Implementation Issues

LL	LR(LALR)
Left derivation	Right derivation
More intuitive	Understanding how it works is quite difficult
Easy for hand-written codes	Error reporting can be quite hard
LL(*) are also very powerful	Handle a larger range of languages
Functional languages <ul style="list-style-type: none"> Clean, Haskell, Lisp or ML 	GCC 4.1: <ul style="list-style-type: none"> replaced by a new, faster hand-written recursive-descent parser GLR parser

138



复旦大学
媒体计算
研究所

作业

- 3.4
- 3.14



139