

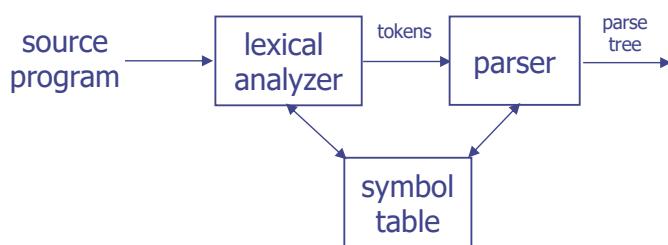
編譯器設計

Chapter 4 Syntax Analysis

Syntax Analysis

◆ The role of the syntax analyzer

- To accept a string of tokens from the scanner
- To verify the string can be generated by the grammar



◆ Three general types of parsers for grammars

- Universal parsing methods
 - e.g. Cocke-Younger-Kasami, Earley's
- Top-down 從 root 遷 $a = b + c \Rightarrow id = id + id \Rightarrow$ 都是達 tree
- Bottom-up 從 bottom 遷 $id = id + id \Rightarrow a = b + c$

too complicated
一般不用

Syntax Error Handling

◆ The goals of the error handler in a parser *report and recover*

- It should report the presence of errors clearly and accurately
- It should recover from each error quickly enough to be able to detect subsequent errors
- *It* should not significantly slow down the processing of correct programs

◆ Many errors could be classified [Ripley et al 1978]:

- 60% punctuation errors
- 20% operator and operand errors ex. Pascal $a:j = b + c;$
- 15% keyword errors
- 5% others

Error Recovery Strategies

◆ Panic mode *丢資料*

$a:j = b + c;$ *丢掉*

if *丢到句子结束或下一行的开始*

$a:j = b + c$

- The simplest method and can be used by most parsers
- On discovering an error, the parser discards input symbols one at a time until one of synchronizing tokens is found
 - ♦ The synchronizing tokens are usually delimiters, e.g. ; or end

◆ Phrase level *改正*

分界符號

- On discovering an error, a parser may perform local correction on the remaining input
 - ♦ e.g. replacing a comma by a semicolon, inserting a missing ;

◆ Error production *標 Error*

- Augment the grammar with productions that generate the erroneous constructs

把錯誤的地方改成 production rule

$a = b + c O$

◆ Global correction

將錯誤標成 error

- Make as few changes as possible in processing the input

Context-Free Grammars

$$A \rightarrow \alpha$$
$$\langle \text{stmt} \rangle \rightarrow \text{id} = \langle \text{expr} \rangle$$

- ◆ The syntax of programming language constructs can be described by *context-free grammars*
 - Or BNF (Backus-Naur Form)
- ◆ Grammars offer significant advantages
 - A grammar gives a precise syntactic specification of a programming language
 - From certain classes of grammars we can automatically construct an efficient parser
 - ◆ The parser construction process can reveal syntactic ambiguities and difficult-to-parse constructs 因為產生不出 tree
 - A properly designed grammar facilitates the translation into target code
 - New language constructs can be added easily

Context-Free Grammars

BNF 描述的就是
Context-Free Grammars

- ◆ A context-free grammar G can be denoted by
$$G = (V_N, V_T, P, S)$$
 - V_N : nonterminals
 - ◆ Syntactic variables that denotes sets of strings
 - V_T : terminals
 - ◆ The basic symbols (i.e. tokens) from which strings are formed
 - P : productions
 - ◆ The manner in which the terminals and nonterminals can be combined for form strings, e.g.
 $\text{expr} \rightarrow \text{expr op expr}$
 $\text{expr} \rightarrow \text{id}$
 $\text{op} \rightarrow + | - | * | /$
 - S : start symbol

Why Not Using Regular Languages

表達能力不足

- ◆ Regular languages are not power enough to represent certain constructs in program languages, e.g. $\{n\}^n$ usually called "balanced"
 - i.e. $\{n\}^n$ is not a regular language
- ◆ Recall the ways to prove a language is regular
 - Find a regular grammar (or regular expressions) that generates the language
 - Find an NFA that recognizes the language
- ◆ However, they can not prove $\{n\}^n$ is not regular
 - The pumping lemma of regular languages will be used
 - Prove by contradiction

regular language: $\{*\}^*$

但不能表達“幾個”
無法表達對稱

要用 Context-Free Language

只要是 Regular Language 就有此特性

Pumping Lemma

語言：句子的集合

A
 $\{i^n\}$

$\{s\}^s$

$|s| \geq n$



- ◆ If A is a regular language and string $s \in A$, where $\exists n \in \mathbb{N}$ such that $|s| \geq n$, then $s = xyz$ where

▪ $xy^iz \in A$, $i \geq 0$,

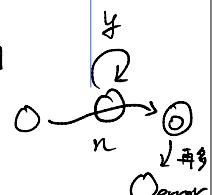
▪ $|y| > 0$, and

▪ $|xy| \leq n$

證明不是：pumping lemma \Rightarrow 不是 regular

$s = \overbrace{00\dots 0}^n \overbrace{11\dots 1}^n = xyz$ xy^z 也在語言內
限制

1. $x, y, z \neq \epsilon$
2. $|x|, |y|, |z| < n$



Example: prove $L = \{0^n 1^n \mid n \geq 0\}$ is not regular

▪ Assume L is regular. Let $s = wyz$, $s \in L$

▪ If y consists of 0's only, then

xy^z has more 0's than 1's

▪ If y consists of 1's only, then

xy^z has more 1's than 0's

▪ If y consists of 0's and 1's, then

$xy^z \neq 0^n 1^n$

$\overbrace{00\dots 0}^n \overbrace{11\dots 1}^n$

$\Rightarrow \overbrace{x\dots 0}^i \overbrace{y\dots 0}^j \overbrace{z\dots 1}^k \rightarrow 0^{n+i} 1^n$ 增加 1

$G = \{s\}, \{0, 1\}$
 $P = \{S \rightarrow s\}$
 $S \rightarrow 0s1\}$
 $0^n 1^n$

$\overbrace{0\dots 0}^n \overbrace{1\dots 1}^n$

$\Rightarrow 0^{n+i} 1^{n+j}$ 增加 0

$$L = O^\alpha | O^{\beta+1} \quad \text{False}$$

$$xyz: x = O^\alpha$$

$$y = O^\beta, \beta > 0$$

$$z = O^{\beta-\alpha-\beta} | \beta+1$$

Look at xy^iz

$$\begin{aligned} xy^iz &= O^\alpha O^{i\beta} O^{\beta-\alpha-\beta} | \beta+1 \\ &= O^{\alpha+i\beta-\beta} | \beta+1 \end{aligned}$$

$$\Rightarrow \alpha + i\beta - \beta < \beta + 1$$

$$\Rightarrow \beta(i-1) < 1$$

We know $\beta > 0 \rightarrow \beta \geq 1$

$$L = \{ O^{2n} | n : n \geq 0 \} \quad \text{False}$$

Suppose L is regular

\Rightarrow exist a p for L

$$\text{Choose } w = O^{2p} | p$$

Look at all decomp of w into

xyz :

$$x = O^\alpha$$

$$y = O^\beta, \beta \geq 1$$

$$z = O^{2p-\alpha-\beta} | p$$

Choose i s.t. $xy^iz \notin L$

xy^iz :

$$O^\alpha O^{i\beta} O^{2p-\alpha-\beta} | p$$

$$= O^{2p+\beta(i-1)-\alpha} | p$$

$$2p + \beta(i-1) = 2p$$

$$\beta(i-1) = 0$$

$$\beta \geq 1, i = 1$$

If $i = 2$, it's not regular language any more.

Pushdown Automata 算法 stack NFA + stack

◆ $M = (K, \Sigma, \Gamma, \delta, q_0, Z_0, F)$ 固定 pop stack 的頂端

- K : Finite set of states

- Σ : Input alphabet

gamma ■ Γ : Pushdown alphabet

delta ■ $\delta: K \times (\Sigma \cup \epsilon) \times \Gamma \rightarrow K \times \underline{\Gamma}^*$

stack (gamma closure)

現在還要查 stack 的資料

- q_0 : Initial state ($q_0 \in K$)

- Z_0 : Initial pushdown symbol ($Z_0 \in \Gamma$)

- F : Set of final states

◆ Example: $L = \{0^n 1^n \mid n \geq 0\}$

- $M = (\{q_0, q_1\}, \{0, 1\}, \{0, 1\}, \delta, q_0, \epsilon, \emptyset)$

- $\delta(q_0, 0, \epsilon) = (q_0, 0)$

- $\delta(q_0, 0, 0) = (q_0, 00)$

- $\delta(q_0, 1, 0) = (q_1, \epsilon)$

- $\delta(q_1, 1, 0) = (q_1, \epsilon)$

- $\delta(q_1, 0, \epsilon) = (q_1, 0)$

- $\delta(q_1, 0, 0) = (q_1, \epsilon)$

NFA

$M = (S, \Sigma, \delta, S_0, F)$

▪ S : a finite, nonempty set of states

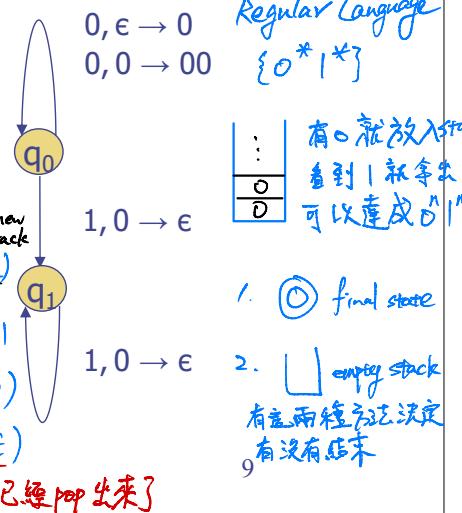
▪ Σ : an alphabet

▪ δ : mapping $S \times \Sigma \rightarrow 2^S$

▪ S_0 : start state

▪ F : the set of accepting (or final) states

Syntax Analysis



Pushdown Automata

◆ Example: $L = \{w c w^R \mid w \in \{0, 1\}^*\}$

- $M = (\{q_0, q_1\}, \{0, 1, c\}, \{R, G, B\}, \delta, q_0, R, \emptyset)$

0011 C 1100
1010 C 0101

- $\delta(q_0, 0, R) = (q_0, BR)$

- $\delta(q_0, 0, G) = (q_0, BG)$

- $\delta(q_0, 0, B) = (q_0, BB)$

- $\delta(q_0, 1, R) = (q_0, GR)$

- $\delta(q_0, 1, G) = (q_0, GG)$

- $\delta(q_0, 1, B) = (q_0, GB)$

- $\delta(q_0, c, R) = (q_1, R)$

- $\delta(q_0, c, G) = (q_1, G)$

- $\delta(q_0, c, B) = (q_1, B)$

- $\delta(q_1, 0, B) = (q_1, \epsilon)$

- $\delta(q_1, 1, G) = (q_1, \epsilon)$

- $\delta(q_1, 1, R) = (q_1, \epsilon)$

- $\delta(q_1, \epsilon, R) = (q_1, \epsilon)$

$B=0, G=1$

$\delta(q_0, 0, R) = (q_0, BR)$
 $\delta(q_0, 1, R) = (q_0, GR)$

$\delta(q_0, 0, B) = (q_0, BB)$
 $\delta(q_0, 1, B) = (q_0, GB)$

$\delta(q_0, 0, G) = (q_0, BG)$
 $\delta(q_0, 1, G) = (q_0, GG)$

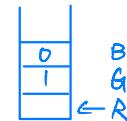
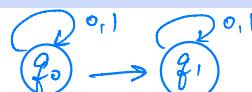
$\delta(q_0, c, R) = (q_1, R)$
 $\delta(q_0, c, G) = (q_1, G)$

$\delta(q_0, c, B) = (q_1, B)$
 $\delta(q_1, 0, B) = (q_1, \epsilon)$

$\delta(q_1, 1, G) = (q_1, \epsilon)$
 $\delta(q_1, 0, B) = (q_1, \epsilon)$

Syntax Analysis

編譯器設計



1, R → GR
1, G → GG
1, B → GB

c, R → R
c, G → G
c, B → B

0, B → ε
1, G → ε
ε, R → ε

$\delta(q_0, 0, B) = (q_1, \epsilon)$

input
top of stack
會拿出來

放入 stack 的內容
有些會直接拿掉
有些會拿掉再放回
甚至加新的東西

Reverse 後

一樣，pop 出最上面，return ε

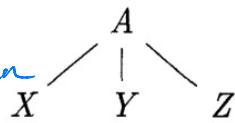
10

Parse Trees

- ◆ A parse tree pictorially shows how the start symbol of a grammar derives a string in the language

■ e.g. $A \rightarrow XYZ$

derivation generates children



- ◆ Formally, given a context-free grammar, a parse tree is a tree with the following properties

- The root is labeled by the start symbol
- Each leaf is labeled by a token or by ϵ
- Each interior node is labeled by a nonterminal
- If A is the nonterminal labeling some interior node and X_1, X_2, \dots, X_n are the labels of the children of that node from left to right, then $A \rightarrow X_1 X_2 \dots X_n$ is a production

$s^* \Rightarrow VT^*$ building parse tree

Syntax Analysis

編譯器設計

11

Parse Trees and Derivations

- ◆ There are several ways to view the process by which a grammar defines a language

- Building parse trees
- Derivations

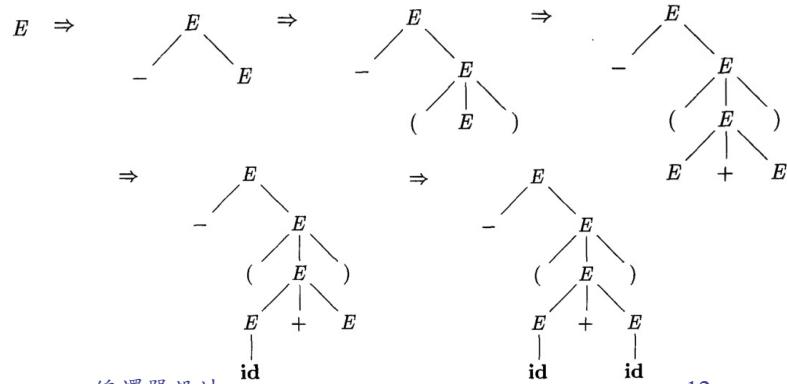
- ◆ Derivations give a precise description of the top-down construction of a parse tree

■ Example: $E \rightarrow E + E \mid E^* E \mid (E) \mid -E \mid id$ *context free grammar*

■ $w = - (id + id)$

■ Derivation:

$$\begin{aligned} E &\Rightarrow - E \\ &\Rightarrow - (E) \\ &\Rightarrow - (E + E) \\ &\Rightarrow - (id + E) \\ &\Rightarrow - (id + id) \end{aligned}$$



Syntax Analysis

編譯器設計

12

Parse Trees and Derivations

- ◆ Every parse tree has associated with it a unique leftmost and a unique rightmost derivation

- However, not every sentence has exactly one leftmost or rightmost derivation
- e.g. $\text{id} + \text{id} * \text{id}$

$$E \Rightarrow E + E$$

$$\Rightarrow \text{id} + E$$

$$\Rightarrow \text{id} + E * E$$

$$\Rightarrow \text{id} + \text{id} * E$$

$$\Rightarrow \text{id} + \text{id} * \text{id}$$

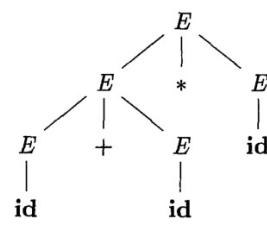
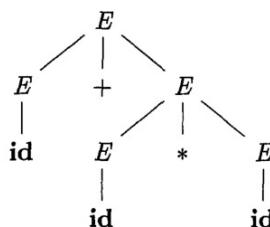
$$E \Rightarrow E * E$$

$$\Rightarrow E + E * E$$

$$\Rightarrow \text{id} + E * E$$

$$\Rightarrow \text{id} + \text{id} * E$$

$$\Rightarrow \text{id} + \text{id} * \text{id}$$



Syntax Analysis

編譯器設計

13

Ambiguity

- ◆ A grammar that produces more than one parse tree for some sentence is said to be *ambiguous*

- i.e. more than one leftmost or rightmost derivation for the same sentence

- ◆ Two options to deal with ambiguity

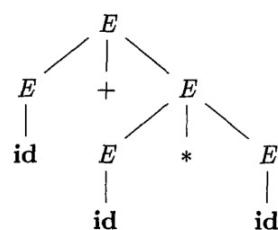
- *Disambiguating rules* are used to throw away undesirable parse trees
- To eliminate ambiguity 跟左/右計算
 - ♦ Precedence and associativity
 - ♦ Rewriting the grammar

Example: $E \rightarrow E + E \mid E * E \mid (E) \mid -E \mid \text{id}$

Rewrite : $E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$ Term

$F \rightarrow (E) \mid \text{id}$ Factor (最底層)



Syntax Analysis

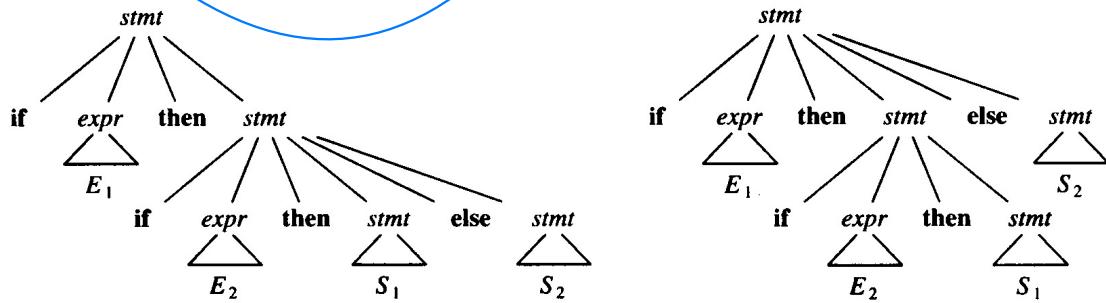
編譯器設計

14

Eliminating Ambiguity

◆ Example (Dangling Else):

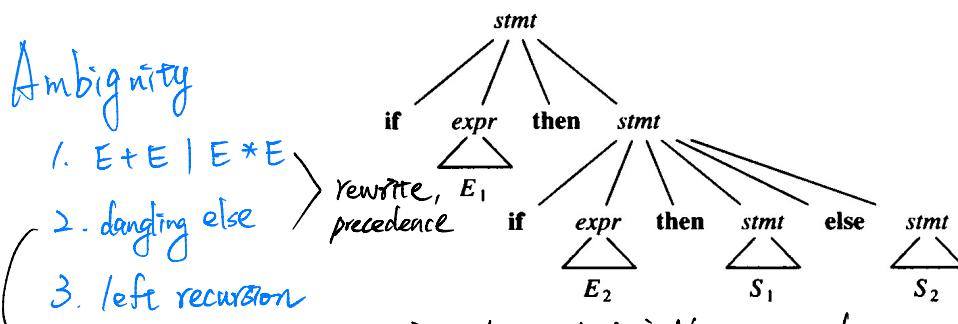
- (1) $\text{stmt} \rightarrow \text{if } \text{expr} \text{ then } \text{stmt}$
(2) $\text{stmt} \rightarrow \text{if } \text{expr} \text{ then } \text{stmt} \text{ else } \text{stmt}$
 $\text{stmt} \rightarrow \text{other}$
■ if E_1 then (if E_2 then S_1) else S_2)



Eliminating Ambiguity

◆ Rewrite the grammar:

$\text{stmt} \rightarrow \text{matched_stmt}$
| unmatched_stmt
 $\text{matched_stmt} \rightarrow \text{if } \text{expr} \text{ then } \text{matched_stmt} \text{ else } \text{matched_stmt}$
| other
 $\text{unmatched_stmt} \rightarrow \text{if } \text{expr} \text{ then } \text{stmt}$
| $\text{if } \text{expr} \text{ then } \text{matched_stmt} \text{ else } \text{unmatched_stmt}$



Left Recursion

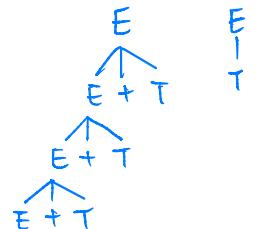
- ◆ A grammar is *left recursive* if it has a nonterminal A such that there is a derivation $A \stackrel{*}{\Rightarrow} Aa$

- e.g. $A \rightarrow Aa \mid \beta$ Compiler keeps choosing first derivation rule and causes left recursion.
It is possible that a parser takes forever if the

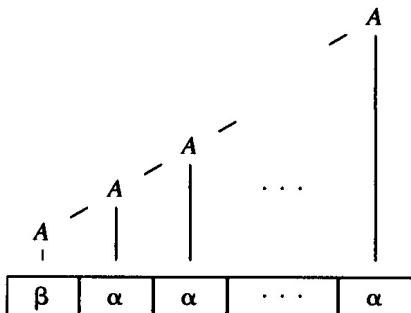
- ❖ It is possible that a parser to loop forever if the grammar is left recursive

$\text{id} + \text{id} \neq \text{id} \Rightarrow$ **parler**

$E \rightarrow E + T | I$ 如果這邊是錯的



Syntax Analysis



Greek letters represent string

編譯器設計

$\Sigma = \{\alpha, \beta\}$

$\xrightarrow{\text{字母 } \beta + n \cdot \alpha}$
 A'

$A \Rightarrow \beta A'$

$A' \Rightarrow \underline{\alpha A'} \mid \underline{\epsilon}$

n组 α 是上面的 ϵ

17

Eliminating Left Recursion

- ◆ A left-recursive pair of production $A \rightarrow A\alpha \mid \beta$ can be changed to non-left-recursive productions

$$A \rightarrow \beta A'$$

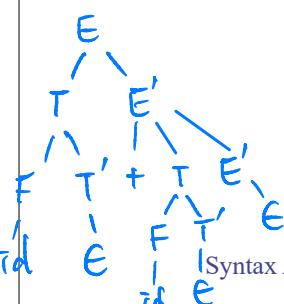
$$A' \rightarrow \alpha A' | \epsilon$$

without changing the set of strings derivable from A

- ## ◆ Example:

$$\begin{array}{l} \text{dimpie: } \\ E \rightarrow \underline{E} + T \mid \underline{T} \\ T \rightarrow \underline{T^*F} \mid F \\ F \rightarrow (E) \mid \text{id} \end{array}$$

$\text{id} + \text{id}$



$\begin{array}{l} id + id \\ E \Rightarrow TE' \\ \Rightarrow FT'E' \\ \Rightarrow idT'E' \\ \Rightarrow idEE' \\ \Rightarrow idE + TE' \end{array} \Rightarrow id + id$

編譯器設計

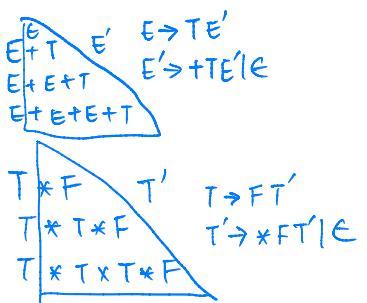
5 TEC

$$E' \rightarrow TE'|_U$$

$$T \rightarrow FT'$$

$$T' \rightarrow *FT' | \epsilon$$

$F \rightarrow (E) \mid \text{id}$



18

Eliminating Left Recursion

左邊的和右邊的第一個字一樣

◆ Algorithm

- Input. Grammar G with no cycles on ϵ -productions
- Output. An equivalent grammar with no left recursion
- Method.

Arrange the nonterminals in some order A_1, A_2, \dots, A_n
for i = 1 to n do

for j = 1 to i-1 do

replace each production of the form $A_i \rightarrow A_j\gamma$ by

the production $A_i \rightarrow \delta_1\gamma | \delta_2\gamma | \dots | \delta_k\gamma$,

where $A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k$ are all the current A_j -productions
eliminate the immediate left recursion among the A_i -productions

end

end

$$\begin{array}{l} A \Rightarrow B\alpha \\ \quad \uparrow \\ B \Rightarrow B\beta \\ \quad \uparrow \\ A \Rightarrow B\beta\alpha \end{array}$$

Eliminating Left Recursion

◆ Example

$$S \rightarrow Aa | b$$

$$A \rightarrow Ac | Sd | \epsilon$$

- Arrange the nonterminals in order S, A

- $i = 1$

♦ No immediate left recursion among the S -productions,

- $i = 2$

♦ Substitute the S -productions in $A \rightarrow Sd$ to obtain

$$A \rightarrow \underline{Ac} | \underline{Aad} | \underline{bd} | \epsilon$$

$$Sd \rightarrow Aad | bd$$

replacable

♦ Eliminate the left recursion and yield

$$S \rightarrow Aa | b$$

$$A \rightarrow bdA' | A'$$

$$A' \rightarrow cA' | adA' | \epsilon$$

$$A \rightarrow \underline{\beta_1} A' | \underline{\beta_2} A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \epsilon$$

$$A \rightarrow \underline{Aa} | \underline{B}$$

in replacable

$$A \rightarrow \beta A' | \alpha A'$$

$$A' \rightarrow \alpha A' | adA' | \epsilon$$

$$A \rightarrow bdA' | A' (CA' = A')$$

$$A' \rightarrow cA' | adA' | \epsilon$$

$$\begin{array}{l} A \rightarrow \underline{Ba} | \underline{Aa} | \underline{C} \\ B \rightarrow \underline{Bb} | \underline{Ab} | \underline{d} \\ A \rightarrow Bab' | CA' \\ A' \rightarrow \alpha A' | \epsilon \\ B \rightarrow Bab \\ B' \rightarrow bB' | \epsilon \end{array}$$

Left Factoring

◆ Left factoring is a grammar transformation

- Usually applied when it is not clear which of two alternative productions to use to expand a nonterminal A
- Rewrite the A -productions to defer the decision
- e.g. $stmt \rightarrow \text{if expr then stmt} \text{ else stmt}$ 等到出現 else 再選
 $stmt \rightarrow \text{if expr then stmt} \text{) 相同}$

◆ Algorithm

- For each nonterminal A find the longest prefix α common to two or more of its alternatives
- Replace $A \rightarrow \alpha\beta_1 | \alpha\beta_2 | \dots | \alpha\beta_n | \gamma$ by

$$A \rightarrow \alpha A' | \gamma$$

$$A' \rightarrow \beta_1 | \beta_2 | \dots | \beta_n$$

開頭是 terminal, 後面是 nonterminal

$$A \rightarrow \alpha\beta_1 | \alpha\beta_2$$

$$A \rightarrow \alpha\beta'$$

$$A' \rightarrow \beta_1 | \beta_2$$

Syntax Analysis

編譯器設計

21

$$\frac{\text{if } E \text{ then } S \text{ } s'}{s' \Rightarrow \beta_1 \text{ } | \text{ else } S \text{ } \beta_2}$$

Top-Down Parsing

$$L = \{cad, cabd\}$$

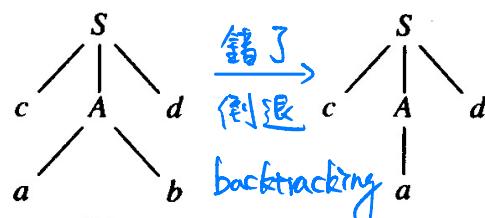
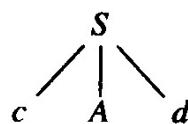
◆ Top-down parsing can be viewed as an attempt to find a leftmost derivation for an input string

- Equivalently, it can be viewed as an attempt to construct a parse tree for the input starting from the root and creating the nodes of the parse tree in preorder

◆ Example

$$S \rightarrow cAd$$

$$A \rightarrow ab | a$$



Top-Down Parsing

◆ Recursive-Descent Parsing

- Start with the root, labeled with the starting symbol, and repeatedly perform the following steps
 - At node n labeled with nonterminal A , select one of the productions for A and construct children at n for the symbols on the right side of the production
 - When a terminal is added that does not match the input string, then backtrack
 - find the next node

◆ Predictive Parsing 不會回溯的 parser

- Recursive-descent parsing without backtracking

Designing a Predictive Parser

◆ A *predictive parser* is a program consisting a procedure for every nonterminal. Each procedure does two things

$$\text{id} + \text{id} \rightarrow \boxed{A}$$

- For all A -productions it decides which production to use by looking at the lookahead symbol, say a , i.e.

- The production $A \rightarrow a$ is used if

$a \in \text{FIRST}(a)$ 第一個字母的集合

where $\text{FIRST}(a)$ is the set of terminals that begins the strings derived from a , or

- The production $A \rightarrow \epsilon$ is used if 選不到就 ϵ
 $a \notin \text{FIRST}(a)$

- The procedure uses a production by mimicking the right side:

- A nonterminal results a call to the procedure for the nonterminal
- A token matching the lookahead symbol results in the next input token being read

$$\begin{aligned} E &\rightarrow TE' \\ E &\rightarrow +TE'E \\ T &\rightarrow FT' \\ T &\rightarrow *FT'E \end{aligned}$$

if (lookahead == "+") {
 指標往下指
 match ("+");
 match (" ");
 T(); E();
}

match (" ")
match (" ");
T(); E();

遇到 terminal
→ match, point to next

Designing a Predictive Parser

◆ Example

type → simple | *id* | array [*simple*] of type
simple → integer | char | num dotdot num

```
procedure type;
begin
  if lookahead ∈ {integer, char, num} then
    simple
  else if lookahead = '^' then begin
    match('^'); match(id)
  end
  else if lookahead = array then begin
    match(array); match(['']); simple;
    match(['']); match(of); type
  end
  else error
end;
```

lookahead <= token 單位

Syntax Analysis

Pascal
var a integer
b char
C array [1..10] of integer

```
procedure match(t : token);
begin
  if lookahead = t then
    lookahead := nextoken
  else error
end;
```

```
procedure simple;
begin
  if lookahead = integer then
    match(integer)
  else if lookahead = char then
    match(char)
  else if lookahead = num then begin
    match(num); match(dotdot); match(num)
  end
  else error
end;
```

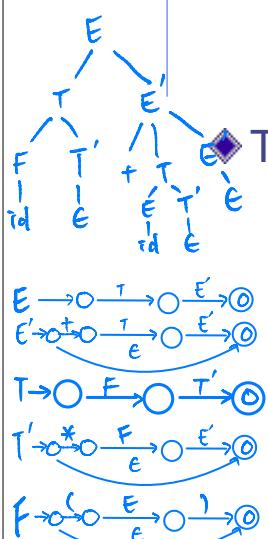
編譯器設計

25

Transition Diagrams for Predictive Parsers

◆ Features

- There is one diagram (procedure) for each nonterminal
- The labels of edges are tokens (terminals) and nonterminals
- A transition on a token is taken if that token is the next input symbol
- A transition on a nonterminal A is a call of the procedure of A



◆ To construct the transition diagram from a grammar

- Eliminate left recursion from the grammar
- Left factor the grammar
- For each nonterminal A , do the following
 - ◆ Create an initial state and a final (return) state
 - ◆ For each production $A \rightarrow X_1X_2\dots X_n$, create a path from the initial to the final state, with edges labeled X_1, X_2, \dots, X_n

Syntax Analysis

編譯器設計

26

Transition Diagrams for Predictive Parsers

◆ Example

$$E \rightarrow TE'$$

$$E' \rightarrow +TE' \mid \epsilon$$

$$T \rightarrow FT'$$

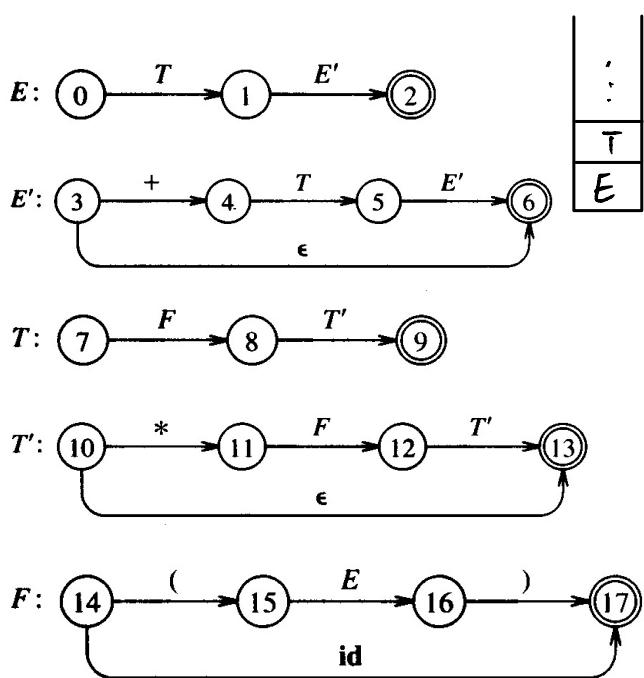
$$T' \rightarrow *FT' \mid \epsilon$$

$$F \rightarrow (E) \mid \text{id}$$

- $w = \text{id} + \text{id} * \text{id}$

```
E() {
    T(); E();
}
```

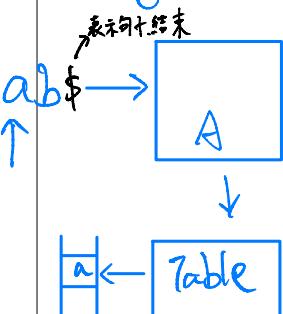
```
E'() {
    if (lookahead == '+') {
        match('+'); T(); E'();
    }
}
```



Nonrecursive Predictive Parser

- Initial
E \$ ← start symbol
- ◆ The key problem during predictive parsing is that of determining the production to be applied for a nonterminal
 - ◆ A table-drive predictive parser has an input buffer, a stack, a parsing table, and an output stream
 - ◆ The action of the parser is determined by X , the symbol on top of the stack, and a , the current input symbol
 - If $X = a = \$$, the parser halts and announces successful completion of parsing
 - If $X = a \neq \$$, the parser pops X off the stack and advances the input pointer to the next input symbol
 - If X is a nonterminal, the program consults entry $M[X, a]$ of the parsing table M .
 - If, for example, $M[X, a] = \{X \rightarrow UVW\}$, the parser replaces X on top of the stack by WVU (with U on top).
 - If $M[X, a] = \text{error}$, the parser calls an error recovery routine.

產生不同的
parsing table.



match
⇒ pop stack頂端, point to next

Nonrecursive Predictive Parser

◆ Algorithm

- Input. A string $w\$$ and a parsing table M

- Output. A leftmost derivation of w

- Method. Initially, stack = $\$S$

set ip to point to the first symbol of $w\$$

repeat

 let X be the top stack symbol and a the symbol pointed by ip

 if X is a terminal or $\$$ then

 if $X = a$ then

 pop X off the stack and advance ip

 else /* X is a nonterminal */

 if $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$ then

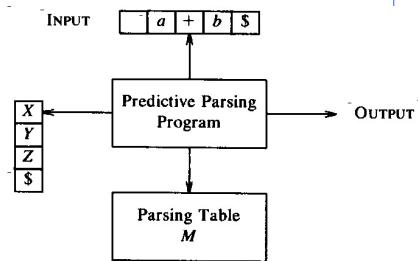
 pop X from the stack

 push Y_k, Y_{k-1}, \dots, Y_1 onto the stack, with Y_1 on the top

 Output the production $X \rightarrow Y_1 Y_2 \dots Y_k$

 else error()

 until $X = \$$ /* stack is empty */



Nonrecursive Predictive Parser

◆ Example

$E \rightarrow TE'$

$E' \rightarrow +TE' | \epsilon$

$T \rightarrow FT'$

$T' \rightarrow *FT' | \epsilon$

$F \rightarrow (E) | id$

◆ Parsing Table

- Fig. 4.17

空格是error \rightarrow syntax error

NONTERMINAL	INPUT SYMBOL					Stack	Input	Output	Leftmost Derivation
E	id	$+$	$*$	$($	$)$	$\$E$	$id + id * id \$$		$E \Rightarrow$
E'	$E \rightarrow TE'$	$E' \rightarrow +TE'$		$E \rightarrow TE'$	$E' \rightarrow \epsilon$	$\$E'T$	$id + id * id \$$	$E \rightarrow TE'$	$TE' \Rightarrow$
T	$T \rightarrow FT'$	$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$	$T \rightarrow FT'$	$T' \rightarrow \epsilon$	$\$E'T'F$	$id + id * id \$$	$T \rightarrow FT'$	$FT'E' \Rightarrow$
T'	$F \rightarrow id$			$F \rightarrow (E)$	$T' \rightarrow \epsilon$	$\$E'T'id$	$id + id * id \$$	$F \rightarrow id$	$id T'E' \Rightarrow$
F					$T' \rightarrow \epsilon$	$\$E'T'$	$+ id * id \$$		
						$\$E'T'$	$+ id * id \$$	$T' \rightarrow \epsilon$	$id E' \Rightarrow$
						$\$E'T+$	$+ id * id \$$	$E' \rightarrow +TE'$	$id + TE' \Rightarrow$
						$\$E'T$	$id * id \$$		
						$\$E'T'F$	$id * id \$$	$T \rightarrow FT'$	$id + FT'E' \Rightarrow$
						$\$E'T'id$	$id * id \$$	$F \rightarrow id$	$id + id T'E' \Rightarrow$
						$\$E'T'$	$* id \$$		
						$\$E'T'F*$	$* id \$$	$T' \rightarrow *FT'$	$id + id * FT'E' \Rightarrow$
						$\$E'T'$	$id \$$		
						$\$E'T'id$	$id \$$	$F \rightarrow id$	$id + id * id T'E' \Rightarrow$
						$\$E'T'$	$$$		
						$\$E'$	$$$	$T' \rightarrow \epsilon$	$id + id * id E' \Rightarrow$
						$\$E'$	$$$		
						$\$E'$	$$$	$E' \rightarrow \epsilon$	$id + id * id \Rightarrow$

Advantages

- ① only thing need to change is parsing table
- ② no recursion

Leftmost Derivation

看到 stack 顶端 \rightarrow 查表格
 放回表格内的放到 top of stack
 check if top matches Input point \rightarrow if match \$

Syntax Analysis \rightarrow pop out, pointer move to next 編譯器設計

FIRST and FOLLOW

$$\begin{array}{l} A \rightarrow \alpha \Rightarrow w \\ B \rightarrow \beta \Rightarrow Y \end{array}$$

- ◆ The construction of a predictive parser is aided by two functions associated with a grammar G

$$\text{First}(E) = \{E\}$$

◆ FIRST(α)

- The set of terminals that begins the strings derived from α
 - To compute $\text{FIRST}(X)$ for all grammar symbols X , apply the following rules until no more terminals or ϵ can be added

- If X is a terminal, then $\text{FIRST}(X) = \{X\}$
 - If $X \rightarrow \epsilon$ is a production, then add ϵ to $\text{FIRST}(X)$
 - If $X \rightarrow Y_1 Y_2 \dots Y_k$ is a production, then $\text{First}(X) = \text{First}(Y_1)$
 - Place a in $\text{FIRST}(X)$ if $\exists i (1 \leq i \leq k) a \in \text{FIRST}(Y_i)$ and $\epsilon \in \text{FIRST}(Y_j) \forall j (1 \leq j < i)$
 - Place ϵ in $\text{FIRST}(X)$ if $\text{First}(Y_1) \cup \dots \cup \text{First}(Y_k) = \{\epsilon\}$

收集句子的開頭做為set

$\forall i \ (1 \leq i \leq k) \in \text{FIRST}(Y_i)$

如果全部都會 First 要加入(ϵ)

- ① if $y_i \neq e$
- ② if $y_i = e$ (忽略)

$\text{First}(Y_1) \cup \text{First}(Y_2) \cup \dots$
 $\text{First}(Y_{n-1}) \quad Y_1 \dots Y_{n-1} \xrightarrow{*} E$

Syntax Analysis

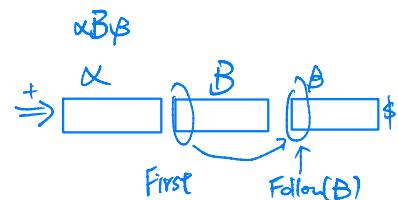
編譯器設計

31

FIRST and FOLLOW

(lowercase →
terminal)
uppercase →
non-terminal.

◆ FOLLOW(A) First 後面的第一個字母



- The set of terminals a that can appear immediately to the right of A in some sentential form, i.e.
 $a \in \text{FOLLOW}(A)$ if there exists a derivation $S \Rightarrow \alpha A a \beta$
 - To compute $\text{FOLLOW}(A)$ for all nonterminals A , apply the following rules until nothing can be added
 - Place $\$$ in $\text{FOLLOW}(S)$
 - If there is production $A \rightarrow \alpha B \beta$, then add $\text{FIRST}(\beta) - \{\epsilon\}$ to $\text{FOLLOW}(B)$
 - If there is production $A \rightarrow \alpha B$, then add $\text{FOLLOW}(A)$ to $\text{FOLLOW}(B)$
 - If there is production $A \rightarrow \alpha B \beta$ and $\epsilon \in \text{FIRST}(\beta)$, then add $\text{FOLLOW}(A)$ to $\text{FOLLOW}(B)$

FOLLOW沒有E

Syntax Analysis

編譯器設計

32

FIRST and FOLLOW

◆ Example

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

$$\begin{aligned} E &\rightarrow TE' & \alpha_A \\ E' &\rightarrow +TE' \mid \epsilon & \alpha_A \alpha_B \\ T &\rightarrow FT' & \alpha_A \\ T' &\rightarrow *FT' \mid \epsilon & \alpha_A \alpha_B \\ F &\rightarrow (E) \mid \text{id} & \alpha_A \alpha_A \end{aligned}$$

◆ FIRST

- $\text{FIRST}(E) = \text{FIRST}(T) = \text{FIRST}(F) = \{(), \text{id}\}$
- $\text{FIRST}(E') = \{+, \epsilon\}$
- $\text{FIRST}(T') = \{*, \epsilon\}$

◆ FOLLOW

- $\text{FOLLOW}(E) = \text{FOLLOW}(E') = \{(), \$\}$
- $\text{FOLLOW}(T) = \text{FOLLOW}(T') = \{+, (), \$\}$
- $\text{FOLLOW}(F) = \{+, *, (), \$\}$

$$\begin{aligned} \text{Follow}(E) &= \{\$\} \cup \{()\} \\ \text{Follow}(E') &= \text{Follow}(E) = \{\$, ()\} \\ \text{Follow}(T) &= \text{First}(E') \cup \text{Follow}(E') = \{+, \$, ()\} \\ \text{Follow}(T') &= \text{Follow}(T) = \{+, \$, ()\} \\ \text{Follow}(F) &= \text{First}(T') \cup \text{Follow}(T') = \{*, +, \$, ()\} \end{aligned}$$

Constructing Predictive Parsing Table

◆ Algorithm

- Input. Grammar G
- Output. Parsing table M
- Method.

1. For each production $A \rightarrow \alpha$, do steps 2 to 4
2. For each terminal $a \in \text{FIRST}(\alpha)$, then
 - add $A \rightarrow \alpha$ to $M[A, a]$
3. If $\epsilon \in \text{FIRST}(\alpha)$, then
 - add $A \rightarrow \alpha$ to $M[A, b]$ for each terminal $b \in \text{FOLLOW}(A)$
4. If $\epsilon \in \text{FIRST}(\alpha)$ and $\$ \in \text{FOLLOW}(A)$, then
 - add $A \rightarrow \alpha$ to $M[A, \$]$
5. Make each undefined entry of M be *error*

$X \rightarrow \alpha \mid \beta \mid \epsilon$	a	b	c	
X	$x \rightarrow \alpha$	$x \rightarrow \beta$	$\text{Follow}(x)$	

如果 α 由 $\alpha \beta$ 產生
 $\Rightarrow \text{FIRST}(\alpha), \text{FIRST}(\beta)$

如果是 ϵ

$\Rightarrow \text{FOLLOW}(x)$

Constructing Predictive Parsing Table

- Example

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid \text{id} \end{aligned}$$

	FIRST	FOLLOW
E	(id) \$
E'	+ ε) \$
T	(id	+) \$
T'	* ε	+) \$
F	(id	* +) \$

Nonterminal	Input Symbol					
	id	+	*	()	\$
E	$E \rightarrow TE'$			$E \rightarrow TE'$		
E'		$E' \rightarrow +TE'$			$E' \rightarrow \epsilon$	$E' \rightarrow \epsilon$
T	$T \rightarrow FT'$			$T \rightarrow FT'$		
T'		$T' \rightarrow \epsilon$	$T' \rightarrow *FT'$		$T' \rightarrow \epsilon$	$T' \rightarrow \epsilon$
F	$F \rightarrow \text{id}$			$F \rightarrow (E)$		

Follow(E')

Constructing Predictive Parsing Table

- Example

(if E then S)
left factoring
of example

$$\begin{aligned} S &\rightarrow iEtSS' \mid a \\ S' &\rightarrow eS \mid \epsilon \\ E &\rightarrow b \\ S &\xrightarrow{a} a \quad S &\xrightarrow{b} b \quad S &\xrightarrow{e} e \quad S &\xrightarrow{i} i \quad S &\xrightarrow{t} t \quad S &\xrightarrow{\$} \$ \\ S' &\xrightarrow{eS} eS \quad S' &\xrightarrow{\epsilon} \epsilon \\ E &\xrightarrow{b} b \end{aligned}$$

$\text{Follow}(S) = \text{First}(S') \cup \$ \cup \text{Follow}(S') = \{e, \$\}$
 $\text{Follow}(S') = \text{Follow}(S) = \{e, \$\}$
 $\text{Follow}(E) = \{t\}$

Nonterminal	Input Symbol					
	a	b	e	i	t	\$
S	$S \rightarrow a$				$S \rightarrow iEtSS'$	
S'				$S' \rightarrow eS$		$S' \rightarrow \epsilon$
E		$E \rightarrow b$				

- 格多個: ambiguous

LL(1) Grammars

*predict parser 只能處理 LL(1)
down down parsing.*

◆ A grammar whose predictive parsing table has no multiply-defined entries is said to be *LL(1)*

- 1st *L*: scanning the input from left to right
- 2nd *L*: producing a leftmost derivation
- 1: using one input symbol for lookahead at each step

$\alpha \in \text{FIRST}(\alpha)$

$\alpha \in \text{FIRST}(\beta)$

$A \rightarrow \alpha \Rightarrow^* \epsilon$

$A \rightarrow \beta$

◆ Properties

- Not ambiguous
- Not left recursive
- Where $A \rightarrow \alpha \mid \beta$ are two distinct productions in G , the following conditions hold:

*left recursion 只會出現在 Top-Down
不會出現在 Bottom-Up.*

1. ◆ For no terminal a do both α and β derive strings beginning with a
2. ◆ At most one of α and β can derive the empty string
3. ◆ If β derives the empty string then α does not derive any string beginning with a terminal in $\text{FOLLOW}(A)$

$$1. \alpha \cap \beta = \emptyset$$

$$2. \alpha \not\Rightarrow^* \epsilon, \beta \not\Rightarrow^* \epsilon$$

$$3. \text{Follow}(A) \cap \text{FIRST}(\beta) = \emptyset$$

Syntax Analysis

編譯器設計

37

Error Handling in Predicting Parsing

Recovery

$a_j = b + c$

1. panic mode

2. local correction

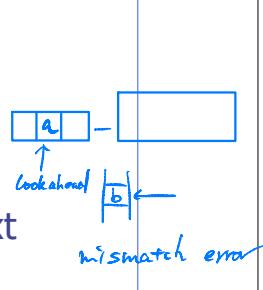
3. error production

4. global.

Table 標註 synchronize tokens

◆ An error is detected when

- The terminal on top of the stack does not match the next input symbol, or
- Nonterminal A is on top of the stack, a is the next input symbol, and $M[A, a]$ is empty



◆ Panic-mode error recovery

- Based on the idea of skipping symbols on the input until a token in a selected set of **synchronizing tokens** appears → $v, ;, \dots$
- Some heuristics of choosing the synchronizing tokens
 - ◆ Place all symbols in $\text{FOLLOW}(A)$ into the synchronizing set of A
 - ◆ Add symbols in $\text{FIRST}(A)$ to the synchronizing set of A
 - ◆ If a nonterminal can generate ϵ , then the production deriving ϵ can be used as a default
 - ◆ If a terminal on top of the stack cannot be matched, then pop the terminal

3. $\Rightarrow id = \text{expr}$
1. $\text{Follow}(s)$
2. $\text{FIRST}(\text{next line})$

Syntax Analysis

編譯器設計

38

Error Handling in Predicting Parsing

◆ Example

$E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow (E) \mid id$



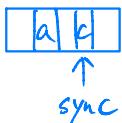
Stack	Input	Output
\$E	<u>+</u> id * <u>id\$</u>	error, skip +
\$E	id * + id\$	id \in FIRST(E)
\$ET	id * + id\$	
\$ET'F	id * + id\$	
\$ET'id	id * + id\$	
\$ET'	* + id\$	
\$ETF*	* + id\$	
\$ETF	+ id\$	error, M[F,+] = sync
\$ET'	+ id\$	F has been popped
\$E'	+ id\$	
\$E'T+	+ id\$	
\$E'T	id\$	
\$E'TF	id\$	
\$E'T'id	id\$	
\$ET'	id\$	
\$E'	id\$	
\$E'	id\$	
\$	id\$	

◆ Parsing table

- Fig. 4.22

◆ Actions

- $M[A, a] = \emptyset$ table 對應到空
input symbol is skipped
- $M[A, a] = \text{sync}$
nonterminal is popped
- Token on stack isn't matched
pop the token



不要再丟了

後面要處理 Syntax Analysis

編譯器設計

39

Bottom-Up Parsing

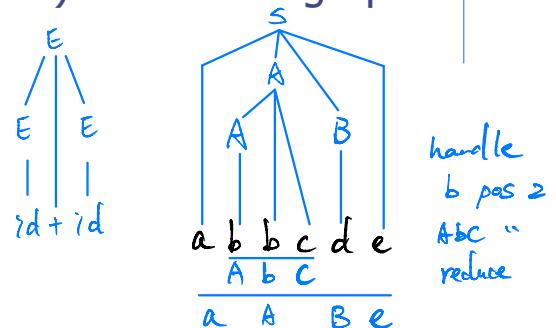
◆ Attempt to construct a parse tree for an input string beginning at the leaves (the bottom) and working up towards to the start symbol

- Example. Consider the grammar

$S \rightarrow aABe$
 $A \rightarrow Abc \mid b$
 $B \rightarrow d$

The sentence abbcde can be reduced to S:

reduce ↓
abbcde
aAbcde
aAde
aABe
S



Rightmost Derivation

$S \Rightarrow$
 $aABe \Rightarrow$
 $aAde \Rightarrow$
 $aAbcde \Rightarrow$
abbcde

Bottom-Up Parsing

◆ A general style of bottom-up syntax analysis, known as *shift-reduce parsing*, will be introduced

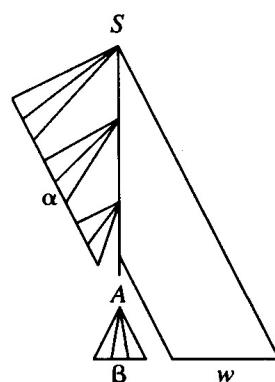
- The process of constructing a parse tree can be thought of as “reducing” a string w to the start symbol
- It is equivalent to performing a rightmost derivation in reverse
- At each reduction step, a particular substring matching the right side of a production is replaced by the symbol on the left
 - ◆ The matched substring can be called a *handle*
 - ◆ The reduction of a handle represents one step along the rightmost derivation in reverse
- Formally, a handle of γ is
 - ◆ A production $A \rightarrow \beta$, and
 - ◆ A position of γ where the substring β may be found
 - ◆ e.g. $A \rightarrow b$ at position 2 is a handle of $abbcde$
 - ◆ $A \rightarrow Abc$ at position 2 is a handle of $aAbcde$

找 handle (位置)
替换

Handle Pruning 剪枝

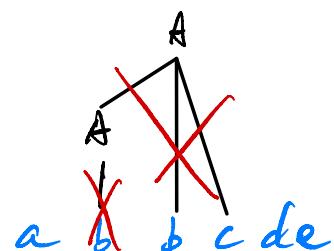
◆ Suppose $A \rightarrow \beta$ is a handle

- It represents a leftmost complete subtree consisting of a node and all its children
- Reducing β to A in $\alpha\beta w$ can be thought of “pruning the handle”
 - ◆ i.e. removing the children of A



◆ A rightmost derivation in reverse can be obtained by handle pruning

最後剩下 S 就對



Stack Implementation of Shift-Reduce Parsing

對應 parser

◆ Two problems for parsing by handle pruning

- To locate the substring to be reduced
- To determine what production to choose in case there is more than one production with that substring

不用清去 left recursion

$E \rightarrow E + E \mid E * E \mid (E) \mid id$ ◆ A convenient way to implement a shift-reduce parser is to use a stack to hold grammar symbols

- Actions
 - ♦ Shift
 - the next input symbol is shifted onto the top of the stack
 - ♦ Reduce
 - The handle is at the top of the stack
 - Remove the handle and place the left side nonterminal on the stack
 - ♦ Accept
 - Announce successful completion of parsing
 - ♦ error

Stack Implementation of Shift-Reduce Parsing

- ◆ Example.
- (1) $E \rightarrow E + E$
 - (2) $E \rightarrow E * E$
 - (3) $E \rightarrow (E)$
 - (4) $E \rightarrow id$

The precedence of * and + is not defined.
Grammar is ambiguous.

Stack	Input	Action	Stack	Input	Action
\$	id + id * id\$	shift	\$	id + id * id\$	shift
\$id	+ id * id\$	reduce by (4)	\$id	+ id * id\$	reduce by (4)
\$E	I + id * id\$	shift	\$E	+ id * id\$	shift
\$E+	id * id\$	shift	\$E+	id * id\$	shift
\$E+ id	* id\$	reduce by (4)	\$E+ id	* id\$	reduce by (4)
\$E+ E	* id\$	reduce by (1)	\$E+ E	* id\$	shift
\$E	* id\$	shift	\$E+ E *	id\$	shift
\$E*	id\$	shift	\$E+ E * id	\$	reduce by (4)
\$E* id	\$	reduce by (4)	\$E+ E * E	\$	reduce by (2)
\$E* E	\$	reduce by (2)	\$E+ E	\$	reduce by (1)
\$E	\$	accept	\$E	\$	accept

Conflicts During Shift-Reduce Parsing

- ◆ For some grammars, a shift-reduce parser can't decide
 - whether to shift or to reduce (a *shift/reduce conflict*), or
 - which of several reductions to make (a *reduce/reduce conflict*)
- Example 1

Stack	Input	Action	Stack	Input	Action
\$E + E	* id\$	shift	\$E + E	* id\$	reduce by $E \rightarrow E + E$
\$E + E *	id\$		\$E	* id\$	

- Example 2
- $stmt \rightarrow \text{if } expr \text{ then } stmt$
| $\text{if } expr \text{ then } stmt \text{ else } stmt$
| other

Stack	Input
... if $expr$ then $stmt$	else ... \$

shift-reduce

LR Parsers

◆ LR(k) parsing $LR(1)$ 就可以了

$LALR(1)$ 在 $LR(1)$ 中
大部分 language 都可用

- an efficient, bottom-up parsing technique
- L: scanning the input from left to right
- R: producing a rightmost derivation in reverse
- k: the number of input symbols of lookahead

◆ Advantages

- Can recognize virtually all programming-language constructs for which context-free grammars can be written
- The most general nonbacktracking shift-reduce parsing → 因為有 parsing table
- The class of grammars that can be parsed using LR methods is a proper superset of the class of grammars that can be parsed with predictive parsers
- Can detect a syntactic error as soon as possible

◆ Drawback

LR parser 的 symbol table 複雜 \rightarrow yacc

- Too much work to construct an LR parser by hand

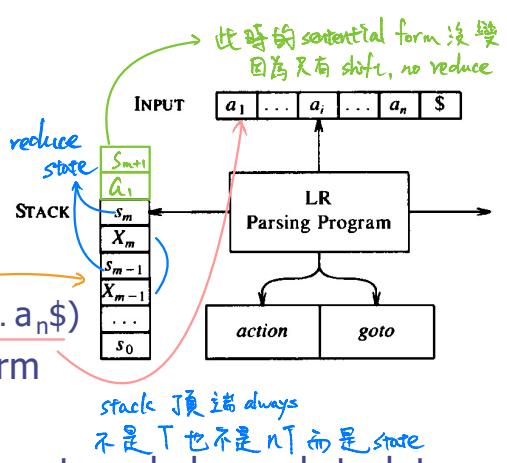


CPP is not $LR(1)$
Because the class is too
complicated, can't be represented
in nonterminal.

LR Parsing Algorithm

◆ Configuration

- The stack contents
- The unexpected input
 $(s_0 X_1 s_1 X_2 s_2 \dots X_m s_m, a_i a_{i+1} \dots a_n \$)$
- It represents the sentential form
 $X_1 X_2 \dots X_m a_i a_{i+1} \dots a_n$

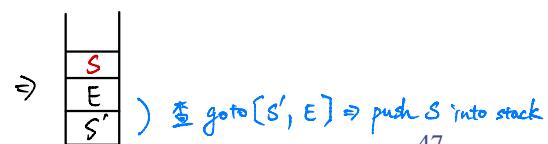


◆ Next move is determined by the input symbol a_i and stack top s_m

- 查 parsing table / action: $\text{action}[s_m, a_i] = \text{shift } s$ (文法符號和 state 是交替的)
- $\text{action}[s_m, a_i] = \text{reduce } A \rightarrow \beta$ where $s = \text{goto}[s_{m-r}, A]$ and $|\beta| = r$ $\Rightarrow \text{pop } 2*|\beta| \text{ (including state)}$
- $\text{action}[s_m, a_i] = \text{accept}$
- $\text{action}[s_m, a_i] = \text{error}$

Syntax Analysis

編譯器設計



LR Parsing Algorithm

- Input. String w and an LR parsing table
- Output. A bottom-up parsing for w
- Method. Initially, stack = s_0 and input = $w\$$

```

set ip to point to the first symbol of  $w\$$ 
repeat forever
    let  $s$  be the state on stack top and  $a$  the symbol pointed by  $ip$ 
    if  $\text{action}[s, a] = \text{shift } s'$  then
        push  $a$  and  $s'$  onto the stack and advance the  $ip$ 
    else if  $\text{action}[s, a] = \text{reduce } A \rightarrow \beta$  then
        push  $2*|\beta|$  symbols off the stack and now  $s'$  is the top state
        push  $A$  and then  $\text{goto}[s', A]$  on top of the stack
        output  $A \rightarrow \beta$ 
    else if  $\text{action}[s, a] = \text{accept}$  then
        return
    else error()
end

```

Syntax Analysis

編譯器設計

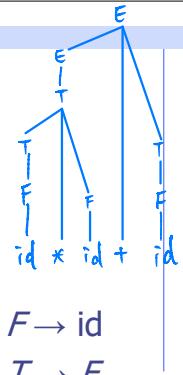
$$E \rightarrow E+E \mid E \times E \mid (E) \mid id$$

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

$$T \rightarrow T * T \mid F$$

$$F \rightarrow (E) \mid id$$

LR Parsing Algorithm



◆ Example.

$$(1) E \rightarrow E + T$$

$$(2) E \rightarrow T$$

$$(3) T \rightarrow T * F$$

$$(4) T \rightarrow F$$

$$(5) F \rightarrow (E)$$

$$(6) F \rightarrow id$$

Parsing table

■ Fig. 4.37

STATE	action					goto
	id	+	*	()	
0	s5		s4			1 2 3
1		s6				
2	r2	s7		r2	r2	
3	r4	r4		r4	r4	
4	s5					8 2 3
5	r6	r6		r6	r6	
6	s5		s4			9 3
7	s5		s4			10
8	s6		s11			
9	r1	s7		r1	r1	
10	r3	r3		r3	r3	
11	r5	r5		r5	r5	

Syntax Analysis

↓ 編譯器設計

Table is the most difficult part

49

Constructing LR Parsing Tables

◆ Three techniques for constructing an LR parsing table

simple ■ Simple LR (SLR)

complicated ■ Canonical LR (LR) 標準 (少用因為太大)

medium ■ Lookahead LR (LALR) yacc

◆ LR(0) item 不管 input 只看文法

■ A production of G with a dot at some position of the right side,

■ e.g. production $A \rightarrow XYZ$ yields

Production $A \rightarrow \epsilon$ yields

LR(0) item $\begin{cases} A \rightarrow \cdot XYZ \\ A \rightarrow X \cdot YZ \\ A \rightarrow XY \cdot Z \\ A \rightarrow XYZ \cdot \end{cases}$

$A \rightarrow \cdot$

$$L(G) = \{w \mid S \xrightarrow{*} w, w \in V_T\}$$

$$G(V_n, V_T, P, S)$$

◆ Augmented grammar $G' (V_n \cup S', V_T, P \cup \{S' \rightarrow S\}, S')$

■ G with a new start symbol S' and production $S' \rightarrow S$

■ Acceptance occurs when the parser is to reduce by $S' \rightarrow S$

Syntax Analysis

編譯器設計

也需要看 input (\$)
判定結束

Closure Operation

- If I is a set of items, then $\text{closure}(I)$ is the set of items constructed from I :
 - Initially, every item in I is added to $\text{closure}(I)$
 - If $A \rightarrow \alpha \cdot B\beta$ is in $\text{closure}(I)$ and $B \rightarrow \gamma$ is a production, then add the item $B \rightarrow \cdot \gamma$ to $\text{closure}(I)$
 - Repeat until no more new items can be added

```

function closure(I)
  J = I
  repeat
    for each item  $A \rightarrow \alpha \cdot B\beta \in J$ 
      and each  $B \rightarrow \gamma$ 
      such that  $B \rightarrow \cdot \gamma \notin J$  do
        add  $B \rightarrow \cdot \gamma$  to  $J$ 
    until no more items can be added to  $J$ 
    return  $J$ 
  end

```

operation	
1. closure	$E' \rightarrow \cdot E$
2. goto	$(A \rightarrow \alpha \cdot B\beta)$
	$E \rightarrow E + T$
	$E \rightarrow T$
	$T \rightarrow T * F$
	$T \rightarrow F$
	$F \rightarrow (E)$
	$F \rightarrow \text{id}$
	$E' \rightarrow \cdot E$
	$E \rightarrow \underline{\cdot E + T} \cdot r$ (gamma)
	$E \rightarrow \cdot T$
	$T \rightarrow \cdot T * F$
	$T \rightarrow \cdot F$
	$F \rightarrow \cdot (E)$
	$F \rightarrow \cdot \text{id}$

Goto Operation

move dot from left to right

- If I is a set of items and X is a grammar symbol, then $\text{goto}(I, X)$ is defined to be the closure of the set of all items $A \rightarrow \alpha X \cdot \beta$ such that $A \rightarrow \alpha \cdot X\beta$ is in I

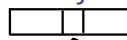
可以是 terminal
也可以是 NT

Example

$E' \rightarrow E$
 $E \rightarrow E + T$
 $E \rightarrow T$
 $T \rightarrow T * F$
 $T \rightarrow F$
 $F \rightarrow (E)$
 $F \rightarrow \text{id}$

no "+", drop

$\text{goto}(\{E' \rightarrow E \cdot, E \rightarrow E \cdot + T\}, +) =$
 $\text{closure}(\{E \rightarrow E + \cdot T\}) = \{$
 $E \rightarrow E + \cdot T$
 $T \rightarrow \cdot T * F$
 $T \rightarrow \cdot F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot \text{id}\}$



- Viable prefixes

- The set of prefixes of right sentential forms that can appear on the stack of a shift-reduce parser

合起來叫sentential form
stack 部分是 viable prefix

Set-of-Items Construction

```

procedure items(G')
    C = {closure({S' → · S})} I₀
    repeat
        for each set of items I in C and each grammar symbol X
            such that goto(I, X) is not empty and not in C do
                add goto(I, X) to C
        until no more set of items can be added to C
    end

```

Example $G =$

$E' \rightarrow E$	$I_0: E' \rightarrow \cdot E$
$E \rightarrow E + T$	$E \rightarrow \cdot E + T$
$E \rightarrow T$	$E \rightarrow \cdot T$
$T \rightarrow T * F$	$T \rightarrow \cdot T * F$
$T \rightarrow F$	$T \rightarrow \cdot F$
$F \rightarrow (E)$	$F \rightarrow \cdot (E)$
$F \rightarrow \text{id}$	$F \rightarrow \cdot \text{id}$

$\text{goto}(I_0, E) = I_1$
 $\text{goto}(I_0, T) = I_2$
 $\text{goto}(I_0, F) = I_3$
 $\text{goto}(I_0, ') = I_4$
 $\text{goto}(I_0, \text{id}) = I_5$

最新的就加入 closure

Set-of-Items Construction $C = \{I_0, I_1, \dots, I_{10}, I_{11}\}$

$I_0: E' \rightarrow \cdot E$	$\text{goto}(I_0, "(") =$	$\text{goto}(I_2, "*") =$	$\text{goto}(I_6, F) = I_3$
$E \rightarrow \cdot E + T$	$I_4: F \rightarrow (\cdot E)$	$I_7: T \rightarrow T * \cdot F$	$\text{goto}(I_6, ")") = I_4$
$E \rightarrow \cdot T$	$E \rightarrow \cdot E + T$	$F \rightarrow \cdot (E)$	$\text{goto}(I_6, \text{id}) = I_5$
$T \rightarrow \cdot T * F$	$E \rightarrow \cdot T$	$F \rightarrow \cdot \text{id}$	$\text{goto}(I_7, F) =$
$T \rightarrow \cdot F$	$T \rightarrow \cdot T * F$	$I_8: F \rightarrow (E \cdot)$	$I_{10}: T \rightarrow T * F \cdot$
$F \rightarrow \cdot (E)$	$T \rightarrow \cdot F$	$E \rightarrow E \cdot + T$	$\text{goto}(I_7, ")") = I_4$
$F \rightarrow \cdot \text{id}$	$F \rightarrow \cdot (E)$	$\text{goto}(I_4, T) = I_2$	$\text{goto}(I_7, \text{id}) = I_5$
	$F \rightarrow \cdot \text{id}$	$\text{goto}(I_4, F) = I_3$	$\text{goto}(I_8, "+") =$
<i>遍到 NT 才停</i>		$\text{goto}(I_4, "(") = I_4$	$I_{11}: F \rightarrow (E) \cdot$
$\text{goto}(I_0, E) =$	$\text{goto}(I_0, \text{id}) =$	$\text{goto}(I_6, T) =$	$\text{goto}(I_8, "+") = I_6$
$I_1: E' \rightarrow E \cdot$	$I_5: F \rightarrow \text{id} \cdot$	$I_9: E \rightarrow E + T \cdot$	$\text{goto}(I_9, "*") = I_7$
$E \rightarrow E \cdot + T$		$T \rightarrow T \cdot * F$	
		$F \rightarrow \cdot (E)$	
		$F \rightarrow \cdot \text{id}$	

ϵ -closure

goto

$\text{goto}(I_0, T) = I_2: E \rightarrow T \cdot$
 $T \rightarrow T \cdot * F$
 $\text{goto}(I_0, F) = I_3: T \rightarrow F \cdot$

$\text{goto}(I_1, "+") = I_6: E \rightarrow E + T \cdot$
 $T \rightarrow T \cdot * F$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot \text{id}$

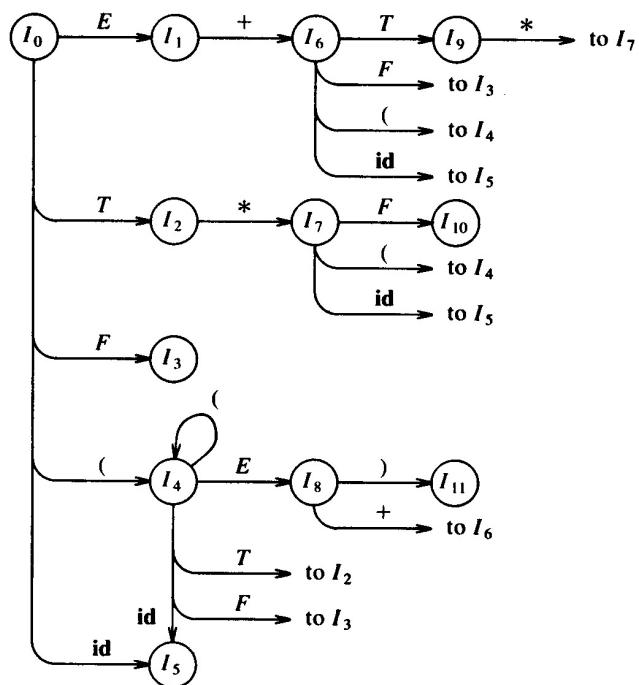
$\text{goto}(I_2, "+") = I_4: T \rightarrow T * F \cdot$
 $F \rightarrow \cdot (E)$
 $F \rightarrow \cdot \text{id}$

Transition Diagram of DFA for Viable Prefixes

- ◆ The *goto* function for the sets of items can be represented as the transition diagram of a DFA D
 - Each set of items I_i is a state of D

Rightmost Derivation

$E \Rightarrow$
 $E + T \Rightarrow$
 $E + F \Rightarrow$
 $E + id \Rightarrow$
 $T + id \Rightarrow$
 $T^* F + id \Rightarrow$
 $T^* id + id \Rightarrow$
 $F * id + id \Rightarrow$
 $id + id * id$



Transition Diagram of DFA for Viable Prefixes

- ◆ D recognizes exactly the viable prefixes of G if
 - Each state is a final state, and
 - I_0 is the initial state
- ◆ For every grammar G , the *goto* function of the collection of sets of items defines a DFA that recognizes the viable prefix of G
 - If each item is treated as a state, an NFA N can be formed:
 - There is a transition from $A \rightarrow \alpha \cdot X \beta$ to $A \rightarrow \alpha X \cdot \beta$ label X , and
 - There is a transition from $A \rightarrow \alpha \cdot B \beta$ to $B \rightarrow \cdot \gamma$ labeled ϵ
 - Then *closure*(I) for the set of items I is exactly the ϵ -closure of a set of NFA states
 - Thus $\text{goto}(I, X)$ gives the transition from I on symbol X in the DFA constructed from N by the subset construction
 - The procedure *items*(G') is just the subset construction itself applied to the NFA N

$\epsilon \downarrow E' \Rightarrow \cdot E \xrightarrow{E} E' \Rightarrow E \cdot$
 $\epsilon \downarrow E \Rightarrow \cdot E + T \xrightarrow{E} E \Rightarrow E \cdot + T \dots$
 $\epsilon \downarrow E \Rightarrow \cdot T \xrightarrow{T} E \Rightarrow T \cdot$
 $\epsilon \downarrow T \Rightarrow \cdot T * F \xrightarrow{T} T \Rightarrow T \cdot * F \xrightarrow{*} T * \cdot F \dots$

Constructing an SLR Parsing Table

◆ Algorithm

- Input. An augmented grammar G'
- Output. The SLR parsing table for G'
- Method.
 1. Construct $C = \{I_0, I_1, I_2, \dots, I_n\}$, the sets of LR(0) items
 2. State i is constructed from I_i . Actions for state i
 - If $A \rightarrow \alpha \cdot a\beta \in I_i$ and $\text{goto}(I_i, a) = I_j$, then
 $\text{action}[i, a] = \text{shift } j$
 - If $A \rightarrow \alpha \cdot \in I_i$, then
 $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha \forall a \in \text{FOLLOW}(A)$
 - If $S' \rightarrow S \cdot \in I_i$, then
 $\text{action}[i, \$] = \text{accept}$
 3. If $\text{goto}(I_i, A) = I_j$, then
 $\text{goto}[i, A] = j$
 4. The set of items containing $S' \rightarrow \cdot S$ is the initial state

Constructing an SLR Parsing Table

■ Example

- (0) $E' \rightarrow E$
- (1) $E \rightarrow E + T$
- (2) $E \rightarrow T$
- (3) $T \rightarrow T * F$
- (4) $T \rightarrow F$
- (5) $F \rightarrow (E)$
- (6) $F \rightarrow \text{id}$

- $\text{FOLLOW}(E) = \{+,), \$\}$
- $\text{FOLLOW}(T) = \{+, *,), \$\}$
- $\text{FOLLOW}(F) = \{+, *,), \$\}$

$$C = \{I_0, I_1, I_2, \dots, I_{11}\}$$

1. $\text{goto}(I_i, a) = I_j$
 $\text{action}[i, a] = \text{shift } j$
2. $A \rightarrow \alpha \cdot \in I_i$
 $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
3. $\text{goto}(I_i, A) = I_j$
 $\text{goto}[i, A] = j$
4. $S' \rightarrow S \cdot$
 $\text{action}[i, \$] = \text{accept}$

	action V_F					goto V_N		
	id	+	*	()	\$	E	T	F
0	s5			s4		1	2	3
1		s6						
2		r2	s7		r2	r2		
3		r4	r4		r4	r4		
4	s5			s4		8	2	3
5		r6	r6		r6	r6		
6	s5			s4			9	3
7	s5			s4				10
8		s6			s11			
9		r1	s7		r1	r1		
10		r3	r3		r3	r3		
11		r5	r5		r5	r5		

Constructing an SLR Parsing Table

- Some grammars are not SLR, e.g.

(0) $S' \rightarrow S$
 (1) $S \rightarrow L = R$
 (2) $S \rightarrow R$
 (3) $L \rightarrow *R$
 (4) $L \rightarrow id$
 (5) $R \rightarrow L$

$I_0: S' \rightarrow \cdot S$
 $S \rightarrow \cdot L = R$
 $S \rightarrow \cdot R$
 $L \rightarrow \cdot *R$
 $L \rightarrow \cdot id$
 $R \rightarrow \cdot L$

$goto(I_0, S) =$
 $I_1: S' \rightarrow S \cdot$

$goto(I_0, L) =$
 $I_2: S \rightarrow L \cdot = R$
 $R \rightarrow L \cdot$

$goto(I_0, R) =$
 $I_3: S \rightarrow R \cdot$

$goto(I_0, "=") =$
 $I_4: L \rightarrow * \cdot R$
 $R \rightarrow \cdot L$
 $L \rightarrow \cdot * R$
 $L \rightarrow \cdot id$

$goto(I_0, id) =$
 $I_5: L \rightarrow id \cdot$

$goto(I_2, "=") =$
 $I_6: S \rightarrow L = \cdot R$
 $R \rightarrow \cdot L$

$goto(I_4, R) =$
 $I_7: L \rightarrow * R \cdot$

$goto(I_4, L) =$
 $I_8: R \rightarrow L \cdot$

$goto(I_4, "=") = I_4$
 $goto(I_4, id) = I_5$

$goto(I_6, R) =$
 $I_9: S \rightarrow L = R \cdot$

$goto(I_6, L) = I_8$
 $goto(I_6, "=") = I_4$
 $goto(I_6, id) = I_5$

Constructing an SLR Parsing Table

- Some grammars are not SLR, e.g.

(0) $S' \rightarrow S$
 (1) $S \rightarrow L = R$
 (2) $S \rightarrow R$
 (3) $L \rightarrow * R$
 (4) $L \rightarrow id$
 (5) $R \rightarrow L$

■ FOLLOW(S) = {\$}

■ FOLLOW(L) = {=, \$}

■ FOLLOW(R) = {=, \$}
 {=, \$}

$id = id$
 reduce
 or shift?

$I_2: S \rightarrow L \cdot = R$
 $R \rightarrow L \cdot$

$LR(0)$ item lookahead
 $R \Rightarrow L$
 $R \Rightarrow L \cdot$
 \nmid
 reduce
 not reduce

⇒ 此語言不屬於SLR(1)

	action				goto		
	id	=	*	\$	S	L	R
0							
1							
2		s6 r5			confice		
3							
4							
5							
6							
7							
8							
9							

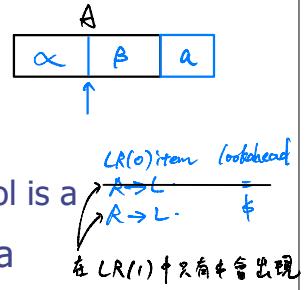
Constructing Canonical LR Parsing Tables

◆ Why does SLR sometimes fail $\{=, \$\}$

- $R \rightarrow L \cdot \in I_2 \Rightarrow \text{action}[2, '='] = \text{reduce } R \rightarrow L$
 $\because = \in \text{FOLLOW}(R)$
- However, there is no sentential form begins $R = \dots$
- Extra information will be incorporated by redefining items

◆ LR(1) items 看左也看右

- The general form of an items becomes $[A \rightarrow \alpha \cdot \beta, a]$
 - The lookahead a has no effect when $\beta \neq \epsilon$
 - $[A \rightarrow \alpha \cdot, a]$ calls for a reduction by $A \rightarrow \alpha$ if next symbol is a
- Formally, LR(1) is valid for a viable prefix γ if there is a derivation $S \xrightarrow{*} \delta A w \Rightarrow \delta \alpha \beta w$, where
 - $\gamma = \delta \alpha$, and
 - Either a is the first symbol of w , or w is ϵ and a is $\$$



Sets of LR(1) Items

```

function closure(I)
  J = I
  repeat
    for each item  $[A \rightarrow \alpha \cdot B \beta, a] \in J$ 
      each  $B \rightarrow \gamma$  in  $G'$  and
      each terminal  $b \in \text{FIRST}(\beta a)$ 
      such that  $[B \rightarrow \cdot \gamma, b] \notin J$  do
        add  $[B \rightarrow \cdot \gamma, b]$  to  $J$ 
    until no more items can be added to  $J$ 
    return  $J$ 
  end

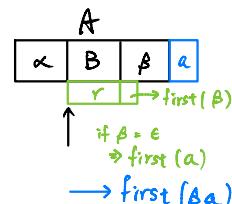
procedure items( $G'$ )
  C = {closure({ $S' \rightarrow \cdot S, \$$ }})  $\leftarrow \{S' \rightarrow \cdot S\}$ 
  repeat
    for each set of item  $I \in C$  and each grammar symbol  $X$ 
      such that  $goto(I, X) \notin C$  do
        add  $goto(I, X)$  to  $C$ 
    until no more sets of items can be added to  $C$ 
  end

```

```

function goto(I, X)
  J = set of items  $[A \rightarrow \alpha X \cdot \beta, a]$ 
  such that  $[A \rightarrow \alpha \cdot X \beta, a] \in I$ 
  return closure(J)
end

```



Sets of LR(1) Items

Example $G' =$
 $S' \rightarrow S$
 $S \rightarrow CC$
 $C \rightarrow cC$
 $C \rightarrow d$
 $\rightarrow \boxed{c^*d \quad c^*d}$
 $I_0: S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot CC, \$$ FIRST($\$$)
 $C \rightarrow \cdot cC, c/d$ FIRST($c\$$)
 $C \rightarrow \cdot d, c/d$
 $A \rightarrow \alpha \cdot B \beta, \alpha$
 $[A \rightarrow \alpha \cdot B \beta, \alpha]$

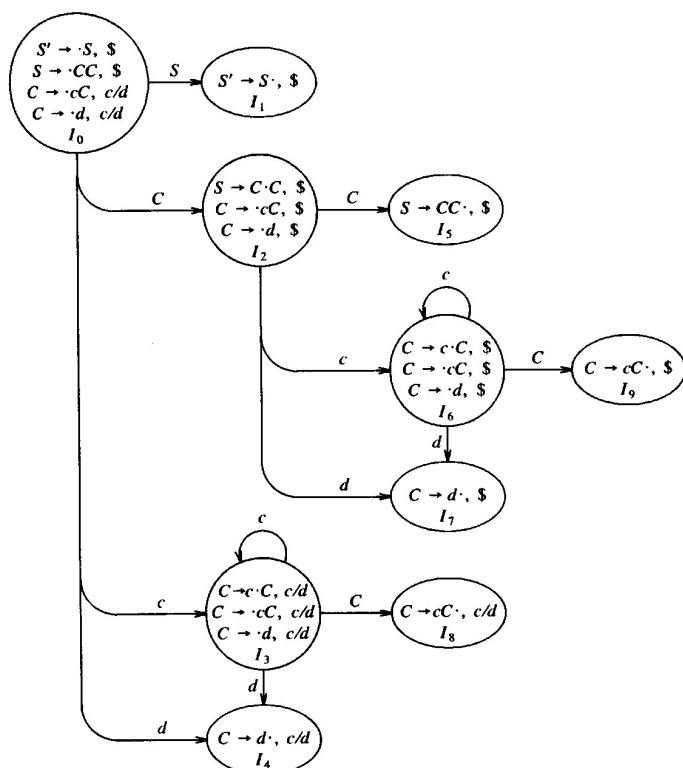
$goto(I_0, S) = I_1: S' \rightarrow S \cdot, \$$
 $goto(I_0, C) = I_2: S \rightarrow C \cdot, C, \$$
 $goto(I_0, c) = I_3: C \rightarrow c \cdot, C, c/d$
 $goto(I_0, d) = I_4: C \rightarrow d \cdot, c/d$

$goto(I_0, d) = I_4: C \rightarrow d \cdot, c/d$
 $goto(I_2, C) = I_5: S \rightarrow CC \cdot, \$$
 $goto(I_2, c) = I_6: C \rightarrow c \cdot, C, \$$
 $goto(I_2, d) = I_7: C \rightarrow d \cdot, \$$

$goto(I_2, d) = I_7: C \rightarrow d \cdot, \$$
 $goto(I_3, C) = I_8: C \rightarrow cC \cdot, c/d$
 $goto(I_3, c) = I_3$
 $goto(I_3, d) = I_4$
 $goto(I_6, C) = I_9: C \rightarrow cC \cdot, \$$
 $goto(I_6, c) = I_6$
 $goto(I_6, d) = I_7$

Sets of LR(1) Items

The *goto* graph



Construction of the LR Parsing Table

◆ Algorithm

- Input. An augmented grammar G'
- Output. The canonical LR parsing table for G'
- Method.
 1. Construct $C = \{I_0, I_1, I_2, \dots, I_n\}$, the sets of LR(1) items
 2. State i is constructed from I_i . Actions for state i
 - If $[A \rightarrow \alpha \cdot a\beta, b] \in I_i$ and $\text{goto}(I_i, a) = I_j$, then $\text{action}[i, a] = \text{shift } j$
 - If $[A \rightarrow \alpha \cdot, a] \in I_i$ and $A \neq S'$ then $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
 - If $[S' \rightarrow S \cdot, \$] \in I_i$, then $\text{action}[i, \$] = \text{accept}$
 3. If $\text{goto}(I_i, A) = I_j$, then $\text{goto}[i, A] = j$
 4. The set of items containing $[S' \rightarrow \cdot S, \$]$ is the initial state

Constructing an LR Parsing Table

Example G' =

$S' \rightarrow S$

$S \rightarrow CC$

$C \rightarrow cC$

$C \rightarrow d$

$\text{goto}(I_6, d) = I_7$

SLR

和 LR

state 是 $S \sim 10$ 倍

	action			goto	
	c	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
8	r2	r2			
9			r2		

Constructing LALR Parsing Tables

◆ Lookahead LR (LALR) technique

- Often used in practice because the tables are much smaller than the LR tables
- Yet most common syntactic constructs of programming languages can be expressed by an LALR grammar
- Consider the following sets of LR(1) items

$$I_3: C \rightarrow c \cdot C, c/d$$

$$C \rightarrow \cdot cC, c/d$$

$$C \rightarrow \cdot d, c/d$$

$$I_4: \frac{C \rightarrow d \cdot}{\text{core}}, c/d$$

$$I_8: C \rightarrow cC \cdot, c/d$$

• The only difference is the lookahead symbols

• Note the G' generates the language c^*dc^*d

- An LALR parsing table can be obtained by merging the sets of items of a LR table with the same core

$$\begin{array}{|c|c|} \hline C & C \\ \hline c^*d & c^*d \\ \hline \end{array} \$$$

SLR單純根據 follow

$$I_6: C \rightarrow c \cdot C, \$$$

$$C \rightarrow \cdot cC, \$$$

$$C \rightarrow \cdot d, \$$$

$$I_7: C \rightarrow d \cdot, \$$$

$$I_9: C \rightarrow cC \cdot, \$$$

$$I_3, I_6 \Rightarrow I_{36}$$

$$I_4, I_7 \Rightarrow I_{47} = \{(c \rightarrow d, c/d/\$)\}$$

Construction of the LALR Parsing Table

◆ Algorithm

- Input. An augmented grammar G'
- Output. The LALR parsing table for G'
- Method.

- Construct $C = \{I_0, I_1, I_2, \dots, I_n\}$, the sets of LR(1) items
- For each core in C , find all sets with that core and replace these sets by their union. Let result be $C' = \{J_0, J_1, J_2, \dots, J_m\}$
- State i is constructed from J_i . Actions for state i
 - If $[A \rightarrow \alpha \cdot a\beta, b] \in J$ and $\text{goto}(J_i, a) = J_j$, then
 $\text{action}[i, a] = \text{shift } j$
 - If $[A \rightarrow \alpha \cdot, a] \in J_i$ and $A \neq S'$ then
 $\text{action}[i, a] = \text{reduce } A \rightarrow \alpha$
 - If $[S' \rightarrow S \cdot, \$] \in J_i$, then
 $\text{action}[i, \$] = \text{accept}$
- If $J_i = I_1 \cup I_2 \cup \dots \cup I_k$, then
 $\text{goto}(I_1, X) = \text{goto}(I_2, X) = \dots = \text{goto}(I_k, X) = J_j$
 $\Rightarrow \text{goto}[i, X] = j$

Construction of the LALR Parsing Table

$$\begin{array}{ll} I_3: C \rightarrow c \cdot C, c/d & I_6: C \rightarrow c \cdot C, \$ \\ C \rightarrow \cdot cC, c/d & C \rightarrow \cdot cC, \$ \\ C \rightarrow \cdot d, c/d & C \rightarrow \cdot d, \$ \end{array}$$

$$I_4: C \rightarrow d \cdot, c/d \quad I_7: C \rightarrow d \cdot, \$$$

$$I_8: C \rightarrow cC \cdot, c/d \quad I_9: C \rightarrow cC \cdot, \$$$

$$C' = \{I_0, I_1, I_2, I_{3b}, I_{47}, I_5, I_{89}\}$$

	action			goto	
	c	d	\$	S	C
0	s3	s4		1	2
1			acc		
2	s6	s7			5
3	s3	s4			8
4	r3	r3			
5			r1		
6	s6	s7			9
7			r3		
89	r2	r2	r2		
9			r2		

	action			goto	
	c	d	\$	S	C
0	s36	s47		1	2
1			acc		
2	s36	s47			5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
6	s6	s7			9
7			r3		
89	r2	r2	r2		
9			r2		

Error Detection by LALR

◆ LALR parser may proceed to do some reductions after the LR parser has declared an error

- LALR parser will never shift another symbol after the LR parser declares an error
- Example $w = ccd\$$

	action			goto	
	c	d	\$	S	C
0	s36	s47		1	2
1			acc		
2	s36	s47	/11		5
36	s36	s47			89
47	r3	r3	r3		
5			r1		
6	s6	s7			9
7			r3		
89	r2	r2	r2		
9			r2		

LALR

Stack	Input	Action
\$0	ccd\$	shift 36
\$0 c 36	cd\$	shift 36
\$0 c 36 c 36	d\$	shift 36 47
\$0 c 36 c 36 d 47	\$	reduce by $C \rightarrow d$
\$0 c 36 c 36 C 89	\$	reduce by $C \rightarrow cC$
\$0 c 36 C 89	\$	reduce by $C \rightarrow cC$
\$0 C 2	\$	error

LR

Stack	Input	Action
\$0	ccd\$	shift 3
\$0 c 3	cd\$	shift 3
\$0 c 3 c 3	d\$	shift 4
\$0 c 3 c 3 d 4	\$	error

always find error at the same point.
but LALR need more reduction steps.

Conflicts Caused by Merging States

- ◆ Merging states with common cores might cause conflicts

- No shift-reduce conflicts will be introduced
 - Because shift actions depend only on the core, not the lookahead
- Reduce-reduce conflicts might be produced by merging

Example $G' =$

$$S' \rightarrow S$$

$$S \rightarrow aAd \mid bBd \mid aBe \mid bAe$$

$$A \rightarrow c$$

$$B \rightarrow c$$

$$goto(I_0, a) =$$

$$I_2: S \rightarrow a \cdot Ad \mid a \cdot Be, \$$$

$$A \rightarrow \cdot c, d$$

$$B \rightarrow \cdot c, e$$

$$goto(I_2, c) =$$

$$I_4: A \rightarrow c \cdot, d$$

$$B \rightarrow c \cdot, e$$

Syntax Analysis

$$I_0: S' \rightarrow \cdot S, \$$$

$$S \rightarrow \cdot aAd \mid \cdot bBd \mid \cdot aBe \mid \cdot bAe, \$$$

$$goto(I_0, S) = I_1: S' \rightarrow S \cdot, \$$$

$$goto(I_0, b) =$$

$$I_3: S \rightarrow b \cdot Bd \mid b \cdot Ae, \$$$

$$A \rightarrow \cdot c, e$$

$$B \rightarrow \cdot c, d$$

$$goto(I_3, c) =$$

$$I_5: A \rightarrow c \cdot, e$$

$$B \rightarrow c \cdot, d$$

編譯器設計



	d	e
45	rA→c rB→c	rA→c rB→c

merge I_4 and $I_5 =$

$$I_{45}: A \rightarrow c \cdot, d/e$$

$$B \rightarrow c \cdot, d/e$$

71

Using Ambiguous Grammars

- ◆ Every ambiguous grammar fails to be LR

- Certain types of ambiguous grammars are useful in the specification and implementation of languages
 - An ambiguous grammar provides a shorter, more natural specification than any equivalent unambiguous grammar
e.g. $E \rightarrow E + E \mid E^* E \mid (E) \mid id$ 加規則, 優先順序
 - Ambiguous grammars can be used to isolate commonly occurring syntactic constructs for special case optimization
- Disambiguating rules must be specified to allow only one parse tree for each sentence
 - e.g. using precedence and associativity to resolve conflicts
 - In this way, the overall language specification still remains unambiguous

Top-Down
 $LL(1)$

Bottom-Up
 $SLR(1)$
 $LALR(1)$
 $LR(1)$

Using Precedence and Associativity

Example G'

$$(0) E' \rightarrow E$$

$$(1) E \rightarrow E + E$$

$$(2) E \rightarrow E * E$$

$$(3) E \rightarrow (E)$$

$$(4) E \rightarrow \text{id}$$

$$I_0: E' \rightarrow \cdot E$$

$$E \rightarrow \cdot E + E$$

$$E \rightarrow \cdot E * E$$

$$E \rightarrow \cdot (E)$$

$$E \rightarrow \cdot \text{id}$$

$$\text{goto}(I_0, E) =$$

$$I_1: E' \rightarrow E \cdot$$

$$E \rightarrow E \cdot + E$$

$$E \rightarrow E \cdot * E$$

$$\text{goto}(I_0, '(') =$$

$$I_2: E \rightarrow (\cdot E)$$

$$E \rightarrow \cdot E + E$$

$$E \rightarrow \cdot E * E$$

$$E \rightarrow \cdot (E)$$

$$E \rightarrow \cdot \text{id}$$

$$\text{goto}(I_0, \text{id}) =$$

$$I_3: E \rightarrow \text{id} \cdot$$

$$\text{goto}(I_1, '+') =$$

$$I_4: E \rightarrow E + \cdot E$$

$$E \rightarrow \cdot E + E$$

$$E \rightarrow \cdot E * E$$

$$E \rightarrow \cdot (E)$$

$$E \rightarrow \cdot \text{id}$$

$$\text{goto}(I_1, '*') =$$

$$I_5: E \rightarrow E * \cdot E$$

$$E \rightarrow \cdot E + E$$

$$E \rightarrow \cdot E * E$$

$$E \rightarrow \cdot (E)$$

$$E \rightarrow \cdot \text{id}$$

$$\text{goto}(I_2, E) =$$

$$I_6: E \rightarrow (E \cdot)$$

$$E \rightarrow E \cdot + E$$

$$E \rightarrow E \cdot * E$$

$$\text{goto}(I_4, E) =$$

$$I_7: E \rightarrow E + E \cdot$$

$$E \rightarrow E \cdot + E$$

$$E \rightarrow E \cdot * E$$

$$\% \text{left } '*'$$

$$\% \text{left } '+'$$

$$\text{left-association}$$

$E + E * E$

$\text{goto}(I_5, E) =$

$I_8: E \rightarrow E * E \cdot$

$$E \rightarrow E \cdot + E$$

$$E \rightarrow E \cdot * E$$

$\text{goto}(I_6, ') =$

$I_9: E \rightarrow (E) \cdot$

	action			goto
	+	*	...	
...				
7	r1	s5		
8	r2	r2		
9				

The “Dangling-else” Ambiguity

Example G' if else

$$(0) S' \rightarrow S$$

$$(1) S \rightarrow iSeS$$

$$(2) S \rightarrow iS$$

$$(3) S \rightarrow a$$

$$I_0: S' \rightarrow \cdot S$$

$$S \rightarrow \cdot iSeS$$

$$S \rightarrow \cdot iS$$

$$S \rightarrow \cdot a$$

$$\text{goto}(I_0, S) =$$

$$I_1: S' \rightarrow S \cdot$$

$$\text{goto}(I_0, i) =$$

$$I_2: S \rightarrow i \cdot SeS$$

$$S \rightarrow i \cdot S$$

$$S \rightarrow \cdot iSeS$$

$$S \rightarrow \cdot iS$$

$$S \rightarrow \cdot a$$

$$\text{goto}(I_0, a) =$$

$$I_3: S \rightarrow a \cdot$$

$$\text{goto}(I_2, S) =$$

$$I_4: S \rightarrow iS \cdot eS \text{ 應該是 shift}$$

$$S \rightarrow iS \cdot$$

$$\text{goto}(I_4, e) =$$

$$I_5: S \rightarrow iSe \cdot S$$

$$S \rightarrow \cdot iSeS$$

$$S \rightarrow \cdot iS$$

$$S \rightarrow \cdot a$$

$$\text{action}$$

$$i \quad e \quad \dots$$

$$\dots$$

$$4 \quad \quad \quad s5$$

$$8 \quad \quad \quad$$

$$9 \quad \quad \quad$$

yacc default:

shift-reduce conflict

都是 shift, 符合我們的 Project, 因此不需更改

Ambiguities from Special-Case Productions

Example G'

- (0) $E' \rightarrow E$
 - (1) $E \rightarrow E \text{ sub } E \text{ sup } E$
 - (2) $E \rightarrow E \text{ sub } E$
 - (3) $E \rightarrow E \text{ sup } E$
 - (4) $E \rightarrow \{E\}$
 - (5) $E \rightarrow c$
- % right sub sup

- $I_0: E' \rightarrow \cdot E$
- $E \rightarrow \cdot E \text{ sub } E \text{ sup } E$
 - $E \rightarrow \cdot E \text{ sub } E$
 - $E \rightarrow \cdot E \text{ sup } E$
 - $E \rightarrow \cdot \{E\}$
 - $E \rightarrow \cdot c$

$a_{\text{sub}}^{\text{sup}}$

$$I_7: E' \rightarrow E \cdot$$

$$E \rightarrow E \cdot \text{ sub } E \text{ sup } E$$

right associated \Rightarrow shift

$$E \rightarrow E \cdot \text{ sub } E \cdot \text{ sup } E$$

$$E \rightarrow E \cdot \text{ sub } E$$

$$E \rightarrow E \cdot \text{ sub } E \cdot$$

$$E \rightarrow E \cdot \text{ sup } E$$

$$I_8: E' \rightarrow E \cdot$$

$$E \rightarrow E \cdot \text{ sub } E \text{ sup } E$$

$$E \rightarrow E \cdot \text{ sub } E$$

$$E \rightarrow E \cdot \text{ sup } E$$

$$E \rightarrow E \cdot \text{ sup } E \cdot$$

$$I_{11}: E' \rightarrow E \cdot$$

$$E \rightarrow E \cdot \text{ sub } E \text{ sup } E$$

$$E \rightarrow E \cdot \text{ sub } E \text{ sup } E \cdot$$

$$E \rightarrow E \cdot \text{ sub } E$$

$$E \rightarrow E \cdot \text{ sup } E$$

$$E \rightarrow E \cdot \text{ sup } E \cdot$$

	action				goto
	sub	sup	}	\$	
...					
7	s4 r2	s10 r2	r2	r2	
8	s4 r3	s5 r3	r3	r3	
9					
11	s4 r1 r3	s5 r1 r3	r1 r3	r1 r3	

↓
reduce - reduce conflict
保證長的，因為長的是special case

Syntax Analysis

編譯器設計

Operator-Precedence Parsing

◆ Works on a class of grammars: *operator grammars*

■ No production right side is ϵ or has two adjacent nonterminals

■ e.g. $E \rightarrow EAE | (E) | -E | id$

$A \rightarrow + | - | * | /$

is not an operator grammar

■ $E \rightarrow E+E | E-E | E^* E | E/E | E | (E) | -E | id$

is an operator grammar

$E + E$

$E \text{ op } E \Rightarrow$ 這不是 operator grammar

* $\bullet > +$
precedence
的次序

◆ Advantages

■ Easy to implement

◆ Disadvantages

■ Hard to handle tokens like the minus sign (with 2 precedences)

■ Can't be sure the parser accepts exactly the desired language

■ Only a small class of grammars can be parsed

Syntax Analysis

編譯器設計

Operator-Precedence Parsing

◆ Precedence relations

Relation	Meaning
$a < \cdot b$	a "yields precedence to" b
$a \doteq b$	a "has the same precedence as" b
$a \cdot > b$	a "takes precedence over" b

- Example $E \rightarrow E + E \mid E - E \mid E^* E \mid E / E \mid (E) \mid -E \mid id$

	id	+	*	\$
id	不會碰到	$\cdot >$	$\cdot >$	$\cdot >$
+	$< \cdot$	$\cdot >$	$< \cdot$	$\cdot >$
*	$< \cdot$	$\cdot >$	$\cdot >$	$\cdot >$
\$	$< \cdot$	$< \cdot$	$< \cdot$	

left

right

Operator-Precedence Parsing

◆ The handle can be found by the steps

- Scan the string from left until the first $\cdot >$ is encountered
- Then scan back backwards over any \doteq until a $< \cdot$ is encountered
- The handle contains everything within the $< \cdot$ and $\cdot >$ above, including any intervening or surrounding nonterminals
- e.g. $w = id + id * id$

$\$ < \cdot id \cdot > + < \cdot id \cdot > * < \cdot id \cdot > \$$ 1. skip $\doteq \cdot < \cdot$, stop by $\cdot >$
 $\$ < \cdot E + < \cdot id \cdot > * < \cdot id \cdot > \$$ 2. look backward to find $< \cdot$
 $\$ < \cdot E + < \cdot E * < \cdot id \cdot > \$$ \Rightarrow handle
 $\$ < \cdot E + < \cdot E * E \cdot > \$$
 $\$ < \cdot E + E \cdot > \$$
 $\$ E \$$