

## Crafting a Compiler with C (VIII)

資科系

林偉川

### *The LL Grammar Class*

- **Elimination of left-recursion**

CFG  $G=(V,T,S,P)$  , where  $P$  has the forms

$A \rightarrow Ax_1 \mid Ax_2 \mid Ax_3 \mid \dots \mid Ax_n$

$A \rightarrow y_1 \mid y_2 \mid y_3 \mid \dots \mid y_m$

Can be changed to

$A \rightarrow y_i \mid y_i Z$  ,  $i=1,2,\dots,m$

$Z \rightarrow x_i \mid x_i Z$  ,  $i=1,2,\dots,n$

Or

$A \rightarrow y_i A'$  ,  $i=1,2,\dots,m$

$A' \rightarrow x_i A' \mid \lambda$  ,  $i=1,2,\dots,n$

## Grammar analysis algorithms

- It is necessary to **analyze the properties of a grammar** to determine if a **grammar** is readily **parsable** and to build the **tables** that can be used to **drive a parsing algorithm**
- Reinforce the **basic concepts** of **grammars** and **derivations**
- Many techniques that discussed previously need to **build parsers** are found as **components** of actual **parser generators**

3

## Grammar Representation

- **Vocabulary**
  - Terminals
  - Nonterminals
- **Productions**
  - Field is an array of structs, each of which has fields
    - Lhs, rhs, rhs\_length
- **Start symbol**

4

## Grammar data structure

```
typedef int symbol, terminal, nonterminal;
#define VOCABULARY (N_T+N_NONT)
typedef struct gram {
    int terminal[N_T], nonterminal[N_NONT];
    int start_symbol, num_productions;
    struct prod {
        int lhs, rhs_length, rhs[MAX_RHS];
    } productions[NUM_P];
    int vocabulary[VOCABULARY];
} grammar;
typedef struct prod production;
```

5

## 找出直接或間接推導成 $\lambda$

- One of the most common grammar computations is determining what **non-terminal** can **derive  $\lambda$**  ( $A \rightarrow \lambda$ )
- **Non-terminal** can **derive  $\lambda$**  may **disappear** during a parse and must be carefully handled
- $A \rightarrow BCD \rightarrow BC \rightarrow B \rightarrow \lambda$  (may take **more than one step**)
  - **Iterative marking algorithm**
    - Non-terminals derive  $\lambda$  **in one step** are marked
    - Non-terminals requiring **a parse tree height of two** or more are found

6

## 找出直接或間接推導成 $\lambda$

- 透過演算法
  - Mark\_lambda
  - Follow set

7

## Determine if a Non-terminal can **derive lambda**

```
typedef short boolean;
typedef boolean marked_vocabulary[VOCABULARY];
marked_vocabulary mark_lambda(const grammar g){
    static marked_vocabulary derives_lambda;
    boolean changes, rhs_drives_lambda;
    int v, i, j;
    production p;
    for (v=0; v<VOCABULARY; v++)
        derives_lambda[v]=FALSE;
```

8

Determine if a Nonterminal can derive lambda

```
do {
  change=FALSE;
  for (i=0; i<g.num_productions; i++) {
    p=g.productions[i];
    if (! derives_lambda[p.lhs]) {
      if (p.rhs_length == 0) { // derives  $\lambda$  directly
        changes=derives_lambda[p.lhs]=TRUE;
        continue;
      }
      rhs_derives_lambda=derives_lambda[p.rhs[0]];
    }
  }
}
```

9

Determine if a Nonterminal can derive lambda

```
for (j=1; j<p.rhs_length; j++)
  rhs_derives_lambda= rhs_derives_lambda
    && derives_lambda[p.rhs[j]];
if (rhs_derives_lambda) {
  changes-TRUE; derives_lambda[p.lhs]=TRUE;
}
}
} while (changes);
return derives_lambda;
}
```

10

## First set

- $\text{First}(\alpha)$  is the set of all the **terminal symbols** that can **begin** a sentential form **derivable from  $\alpha$**
- $\text{First}(\alpha) = \{ a \in V_t \mid \alpha \Rightarrow^* a\beta \} \cup \{ \lambda \mid \text{if } \alpha \Rightarrow^* \lambda \text{ then } \{ \lambda \} \text{ else } \phi \}$ .
- If  $\alpha$  is the **right-hand side** of a production, then  $\text{First}(\alpha)$  contains **terminal symbols** that **begin strings derivable from  $\alpha$**
- $\text{First\_set}[X]$ , and  $\text{Follow\_set}[X]$  where  **$X$  is a nonterminal**. Elements of First and Follow set are **terminals** and  $\lambda$

11

## *The LL Grammar Class*

- **Def:**  $\text{FIRST}(\alpha) = \{ a \mid \alpha \Rightarrow^* a\beta \}$  be the set of terminals derived from  $\alpha$  (If  $\alpha \Rightarrow^* \lambda$ ,  $\lambda$  is in  $\text{FIRST}(\alpha)$ )

- ***Pairwise Disjointness Test:***

For each nonterminal,  $A$ , in the grammar that has more than one RHS, for each pair of rules,  $A \rightarrow \alpha_i$  and  $A \rightarrow \alpha_j$ , it must be true that

$$\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \phi$$

**$\rightarrow$  For each non-terminal  $A$ , the first terminal symbol of its RHS must be unique!!!**

12

## The LL Grammar Class

1.  $A \rightarrow a \mid bB \mid cAb$

2.  $A \rightarrow a \mid aB$

$\text{First}(A1) = \{ a, b, c \}$ ,  $\text{First}(A2) = \{ a \}$

$\text{FIRST}(\alpha_i) \cap \text{FIRST}(\alpha_j) = \{ a \} \neq \phi \rightarrow$  **not Pairwise Disjointness**

*Nonterminal A cannot be determined to use which of the RHS*

**Replace**

$\langle \text{variable} \rangle \rightarrow \text{id} \mid \text{id} [\langle \text{expression} \rangle]$

**with**

$\langle \text{variable} \rangle \rightarrow \text{id} \langle \text{new} \rangle$ ,  $\langle \text{new} \rangle \rightarrow \lambda \mid [\langle \text{expression} \rangle]$

**or**

$\langle \text{variable} \rangle \rightarrow \text{id} [[\langle \text{expression} \rangle]]$

(the outer brackets are **meta-symbols** of EBNF)

13

## The LL Grammar Class

### • The Left Recursion Problem

If a grammar has left recursion, either **direct** or **indirect**, it cannot be the basis for a **top-down parser**

**Ex. 1. Direct left recursive**  $A \rightarrow A + B$

**2. Indirect left recursive**  $A \rightarrow B a A$ ,  $B \rightarrow Ab$

- A grammar can be modified to **remove left recursion**
- **Left recursive** is a problem for all **top-down** parsing algorithm, but it is not a problem for **bottom-up** parsing algorithm (can eliminate the left recursive!!)

14

## *The LL Grammar Class*

- The backtracking problem  
If a sequence of **erroneous expansions** and discover a mismatch, we have to **undo** the **semantic effects** of making these erroneous expansion. The **recursive decent parser** is the type of **top-down parser** that **avoid backtracking**.

15

## *The LL Grammar Class*

- Alternates problem  
**Top-down backtracking parser** is that the **order** in which alternates are tried can affect the language accepted.  
Ex.  $S \rightarrow cAd$   
 $A \rightarrow ab \mid a \rightarrow cabd$  as the input string  
if we choose  $A \rightarrow a$  instead of  $A \rightarrow ab$ , we could fail to accepted!! (**violate pairwise-disjoint rule**)

16



## *The LL Grammar Class*

- The other characteristic of grammars that disallows top-down parsing is the **lack of pairwise disjointness**  
The inability to **determine the correct RHS** on the basis of **one token of lookahead** → **First set test**

17

## Explanation of first\_set algorithm

- For an arbitrary string  $\alpha$ , **compute\_first( $\alpha$ )** can return the **set of terminals** defined by **FIRST( $\alpha$ )**
- If  $\alpha$  happens to be exactly **one symbol long**, **compute\_first( $\alpha$ )** will simply return **first\_set[ $\alpha$ ]**
- **Fill\_first\_set** initializes first\_set. The algorithm operates iteratively, first considering **single productions**, then considering **chains of productions**

18

## Compute First(alpha)

```
typedef set_of_terminal_or_lambda termset;
termset follow_set[NUM_NONTERMINAL];
termset first_set[SYMBOL];
marked_vocabulary derives_lambda=mark_lambda(g);
termset compute_first(string_of_symbols alpha) {
    int i,k;
    termset result;
    k=length(alpha);
```

19

## Compute First(alpha)

```
if (k == 0) result=SET_OF( $\lambda$ );
else {
    result=first_set[alpha[0]];
    for (i=1; i<k &&  $\lambda \in$  first_set[alpha[i-1]]; i++)
        result=result  $\cup$  (first_set[ alpha[i] ] - SET_OF( $\lambda$ ));
    if (i == k &&  $\lambda \in$  first_set[ alpha[k-1] ])
        result=result  $\cup$  SET_OF( $\lambda$ );
}
return result;
}
```

20

## Compute the First set

To compute **First(X)** for all grammar symbols **X**, we can have the following algorithm:

1. If **X** is **terminal**, then **FIRST(X)** is **{ X }**
2. If **X** is **non-terminal** and  $X \rightarrow a\alpha$  is a production, add **a** to **FIRST(X)**. If  $X \rightarrow \lambda$  is a production, then add  **$\lambda$**  to **FIRST(X)**.
3. If  $X \rightarrow Y_1 Y_2 \dots Y_k$  is a production, then for all  $i \in \{1, 2, \dots, k\}$  such that  $Y_1 Y_2 \dots Y_{i-1} \Rightarrow^* \lambda$ , add every non- **$\lambda$**  symbol in **FIRST( $Y_i$ )** to **FIRST(X)**. If  **$\lambda$**  is in **FIRST( $Y_j$ )** for  $j=1, 2, \dots, k$ , then add  **$\lambda$**  to **FIRST(X)**.

21

## Compute First set for V

```
extern grammar g;
void fill_first_set(void) {
    nonterminal A;
    terminal a;
    production p;
    boolean changes;
    int i,j;
    for (i=0; i<NUM_NONTERMINAL; i++) { //1st loop
        A=g.nonterminals[i];
        if (derives_lambda[A]) first_set[A]=SET_OF( $\lambda$ );
        else first_set[A]= $\phi$ ;
    }
}
```

22

## Compute First set for V

```
for (i=0; i<NUM_TERMINAL; i++) { // 2nd loop
  a=g.terminals[i];
  first_set[a]=SET_OF(a);
  for (j=0; j<NUM_NONTERMINAL; j++) {
    A=g.nonterminals[j];
    if (there exists a production  $A \rightarrow a\beta$ )
      first_set[A]=first_set[A]  $\cup$  SET_OF(a);
  }
}
```

23

## Compute First set for V

```
do { // 3rd loop
  change=FALSE;
  for (i=0; i < g.num_productions; i++) {
    p=g.productions[i];
    first_set[p.lhs]=first_set[p.lhs]  $\cup$ 
      compute_first[p.rhs];
    if (first_set changed) changes=TRUE;
  }
} while (changes);
}
```

24

## Example of grammar

$E \rightarrow \text{Prefix } ( E ) \mid V \text{ Tail}$

$\text{Prefix} \rightarrow F \mid \lambda$

$\text{Tail} \rightarrow + E \mid \lambda$

Step	First_set							
	E	Prefix	Tail	(	)	V	F	+
(1) First loop	$\phi$	{ $\lambda$ }	{ $\lambda$ }					
(2) Second (nested) loop	{V}	{F, $\lambda$ }	{+, $\lambda$ }	{(}	{)}	{V}	{F}	{+}
(3) Third loop, production 1	{V,F,(}	{F, $\lambda$ }	{+, $\lambda$ }	{(}	{)}	{V}	{F}	{+}

25

## Follow set

- When constructing parser, we often analyze a grammar to compute a set  $\text{Follow}(A)$ , where A is any non-terminal.
- $\text{Follow}(A)$  is the set of terminals that may follow A in some sentential form
- If A appears as the rightmost symbol in a sentential form,  $\lambda$  is include in  $\text{Follow}(A)$
- $\text{Follow}(A) = \{ a \in V_t \mid S \rightarrow + \dots A a \} \cup \{ \text{if } S \rightarrow + \alpha A \text{ then } \{ \lambda \} \text{ else } \phi \}$

26

## Follow set

- Follow(A) provides the **lookahead** that might signal the **recognition of a production with A** as the left-hand side
- **Def: FOLLOW(A) = { a | S =>\*  $\alpha A a \beta$  for some  $\alpha, \beta$  }** be the set of terminals appear **immediately to the right of A** for some  $\alpha, \beta$ . **If A can be the rightmost symbol in some sentential form, then add  $\lambda$  to FOLLOW(A)**

27

## Compute the Follow set

To compute Follow(A) for all non-terminals A, we can have the following algorithm:

1.  **$\lambda$  is in Follow(A), where A is the start symbol**
2. **If there is a production  $A \rightarrow \alpha B \beta$ ,  $\beta \neq \lambda$ , then everything in **First( $\beta$ )** but  $\lambda$  is in **Follow(B)****
3. **If there is a production  $A \rightarrow \alpha B$ , or a production  $A \rightarrow \alpha B \beta$  where **First( $\beta$ ) contains  $\lambda$**  (i.e.  $\beta \Rightarrow^* \lambda$ ), then everything in **Follow(A)** is in **Follow(B)****

28

## Compute Follow set for all Nonterminal

```
extern grammar g;
void fill_follow_set(void) {
    nonterminal A, B;
    boolean changes;
    int i;
    for (i=0; i<NUM_NONTERMINAL; i++) {
        A=g.nonterminals[i]; follow_set[A]= $\phi$ ;
    } //initialization
    follow_set[g.start_symbol]=SET_OF( $\lambda$ );
```

29

## Compute Follow set for all Nonterminal

```
do {
    change=FALSE;
    for (each production  $A \rightarrow \alpha B \beta$ ; i++) {
        follow_set[B]=follow_set[B]  $\cup$ 
            (compute_first[ $\beta$ ]-SET_OF( $\lambda$ ));
        if ( $\lambda \in$  compute_first( $\beta$ ))
            follow_set[B]=follow_set[B]  $\cup$  follow_set[A];
        if (follow_set[B] changed) changes=TRUE;
    }
} while (changes);
}
```

30

## Example of grammar

$E \rightarrow \text{Prefix } ( E ) \mid V \text{ Tail}$

$\text{Prefix} \rightarrow F \mid \lambda$

$\text{Tail} \rightarrow + E \mid \lambda$

Step	Follow_set		
	E	Prefix	Tail
(1) Initialization	$\{\lambda\}$	$\phi$	$\phi$
(2) Process Prefix in production 1	$\{\lambda\}$	$\{( )\}$	$\phi$
(3) Process E in production 1	$\{\lambda, )\}$	$\{( )\}$	$\phi$
(4) Process Tail in production 2	$\{\lambda, )\}$	$\{( )\}$	$\{\lambda, )\}$

31

## Example of grammar 1

$S \rightarrow aSe \mid B$

$B \rightarrow bBe \mid C$

$C \rightarrow cCe \mid d$

Step	First_set								
	S	B	C	a	b	c	d	e	
(1) First loop	$\phi$	$\phi$	$\phi$						
(2) Second (nested) loop	$\{a\}$	$\{b\}$	$\{c,d\}$	$\{a\}$	$\{b\}$	$\{c\}$	$\{d\}$	$\{e\}$	
(3) Third loop, production 2	$\{a,b\}$	$\{b\}$	$\{c,d\}$	$\{a\}$	$\{b\}$	$\{c\}$	$\{d\}$	$\{e\}$	
(4) Third loop, production 4	$\{a,b\}$	$\{b,c,d\}$	$\{c,d\}$	$\{a\}$	$\{b\}$	$\{c\}$	$\{d\}$	$\{e\}$	
(5) Third loop, production 2	$\{a,b,c,d\}$	$\{b,c,d\}$	$\{c,d\}$	$\{a\}$	$\{b\}$	$\{c\}$	$\{d\}$	$\{e\}$	

32



## Example of grammar 1

$S \rightarrow aSe \mid B$

$B \rightarrow bBe \mid C$

$C \rightarrow cCe \mid d$

Step	Follow_set		
	S	B	C
(1) Initialization	{ $\lambda$ }	$\phi$	$\phi$
(2) Process S in production 1	{e, $\lambda$ }	$\phi$	$\phi$
(3) Process B in production 2	{e, $\lambda$ }	{e, $\lambda$ }	$\phi$
(4) Process B in production 3	{e, $\lambda$ }	{e, $\lambda$ }	$\phi$
(5) Process C in production 4	{e, $\lambda$ }	{e, $\lambda$ }	{e, $\lambda$ }
(6) Process C in production 5	{e, $\lambda$ }	{e, $\lambda$ }	{e, $\lambda$ }

33

## Example of grammar 2

$S \rightarrow ABc$

$A \rightarrow a \mid \lambda$

$B \rightarrow b \mid \lambda$

Step	First_set					
	S	A	B	a	b	c
(1) First loop	$\phi$	{ $\lambda$ }	{ $\lambda$ }			
(2) Second (nested) loop	$\phi$	{a, $\lambda$ }	{b, $\lambda$ }	{a}	{b}	{c}
(3) Third loop, production 1	{a,b,c}	{a, $\lambda$ }	{b, $\lambda$ }	{a}	{b}	{c}

Step	Follow_set		
	S	A	B
(1) Initialization	{ $\lambda$ }	$\phi$	$\phi$
(2) Process A in production 1	{ $\lambda$ }	{b,c}	$\phi$
(3) Process B in production 1	{ $\lambda$ }	{b,c}	{c}

34

## Homework

$S \rightarrow AbB \mid d$        $\langle \text{program} \rangle \rightarrow \text{begin } \langle \text{stmts} \rangle \text{ end } \$$   
 $A \rightarrow CAB \mid B$        $\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle ; \langle \text{stmts} \rangle$   
 $B \rightarrow cSd \mid \lambda$        $\langle \text{stmts} \rangle \rightarrow \lambda$   
 $C \rightarrow a \mid ed$        $\langle \text{stmt} \rangle \rightarrow \text{simplestmt}$   
                                  $\langle \text{stmt} \rangle \rightarrow \text{begin } \langle \text{stmts} \rangle \text{ end}$   
 $\langle \text{Assign} \rangle \rightarrow \langle \text{id} \rangle := \langle \text{expr} \rangle$   
 $\langle \text{id} \rangle \rightarrow A \langle \text{id\_tail} \rangle \mid B \langle \text{id\_tail} \rangle \mid \dots \mid Z \langle \text{id\_tail} \rangle$   
 $\langle \text{id\_tail} \rangle A \langle \text{id\_tail} \rangle \mid B \langle \text{id\_tail} \rangle \mid \dots \mid Z \langle \text{id\_tail} \rangle \mid$   
                                  $0 \langle \text{id\_tail} \rangle \mid \dots \mid 9 \langle \text{id\_tail} \rangle \mid \lambda$   
 $\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid$   
                                  $\langle \text{expr} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$   
 $\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$   
                                  $\langle \text{term} \rangle / \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$   
 $\langle \text{factor} \rangle \rightarrow ( \langle \text{expr} \rangle ) \mid \langle \text{id} \rangle$

35