

Crafting a Compiler with C (VII)

資科系
林偉川

CGF Definition

Grammar and parsing !!

CFG $G=(V,T,S,P)$ is 4-tuple, where

$V (V_n)$: a finite set of **variables**(Non-terminal)

$T (V_t)$: a finite set of **terminal symbols** (token set produced by the scanner)

$S \in V_n$ is a special symbol called **start symbol**

P is a finite set of **productions** (rules) where

$A \rightarrow X_1 \dots X_m$, where $A \in V_n$, $X_i \in V_n \cup V_t$

CFG notation

- a,b,c in **lowercase** denote **symbols** in V_t
- A,B,C in **uppercase** denote symbols in V_n
- α, β, γ denote strings in V^*
- u,v,w denote strings in V_t^*
- $A \rightarrow \alpha$ or $A \rightarrow X_1 \dots X_m$, means the **left-hand side** of a production should be a **single nonterminal**, but the **right-hand side** is a string of **zero or more vocabulary symbols**
- $A \rightarrow \alpha \mid \beta \mid \dots \mid \zeta$ is an EBNF

3

CFG derivation

- $A \rightarrow \gamma$ is a production, then $\alpha A \beta \rightarrow \alpha \gamma \beta$ denote a **one-step derivation**, extend \rightarrow to \rightarrow^+ , derived in **one or more steps**, and \rightarrow^* , derived in **zero or more steps**
- $S \rightarrow^* \beta$, then β is said to be a **sentential form** of the CFG. **SF(G)** is the **set of sentential forms** of grammar G, $L(G) = \{x \in V_t^* \mid S \rightarrow^+ x\}$, is the produced language $\rightarrow L(G) = SF(G) \cap V_t^*$ the language G is simply those sentential forms of G that are terminal strings

4

CFG derivation

- If $G_1 \neq G_2$, and $L(G_1) = L(G_2)$, then G_1 and G_2 are **equivalence**
- When deriving a token sequence, if more than one non-terminal is present, we have a choice of which **non-terminal to expand** next and what **production is applied**
- **Leftmost derivation** denoted as \rightarrow_{lm} , \rightarrow_{lm}^+ and \rightarrow_{lm}^* (top-down parser)
- **Rightmost derivation** denoted as \rightarrow_{rm} , \rightarrow_{rm}^+ and \rightarrow_{rm}^* (bottom-up parser)

5

Derivation example

A grammar G as **left part** (V is variables, F is functions), the **leftmost** and **rightmost** derivation of $F(V+V)$ is denoted as **middle** and **right** part.

$E \rightarrow \text{Prefix (E)}$		$E \rightarrow_{rm} \text{Prefix(E)}$
$E \rightarrow V \text{ Tail}$		$E \rightarrow_{rm} \text{Prefix(V Tail)}$
$\text{Prefix} \rightarrow F$	$E \rightarrow_{lm} \text{Prefix(E)}$	$E \rightarrow_{rm} \text{Prefix(V+E)}$
$\text{Prefix} \rightarrow \lambda$	$E \rightarrow_{lm} F(E)$	$E \rightarrow_{rm} \text{Prefix(V+V Tail)}$
$\text{Tail} \rightarrow + E$	$E \rightarrow_{lm} F(V \text{ Tail})$	$E \rightarrow_{rm} \text{Prefix(V+V)}$
$\text{Tail} \rightarrow \lambda$	$E \rightarrow_{lm} F(V+E)$	$E \rightarrow_{rm} F(V+V)$
	$E \rightarrow_{lm} F(V+V \text{ Tail})$	
	$E \rightarrow_{lm} F(V+V)$	

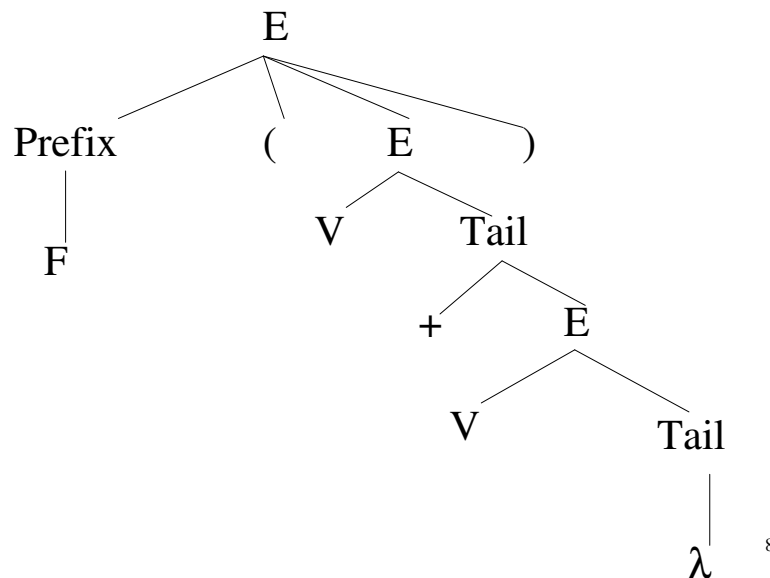
6

Leftmost or rightmost derivation

- A sentential form produced via a **leftmost derivation** sequence is called **left sentential form**
- An alternative to a leftmost derivation is a **rightmost derivation** (**canonical** derivation)
- A **bottom-up parser** discovers the production used to derive a token sequence, it discovers a **rightmost derivation** in the **reverse direction** of a leftmost derivation

7

Parse tree of grammar



8

Handle and phrase

- A **simple or prime phrase** is a phrase that contains **no smaller phrase** and is a sequence of **symbols directly derived from a non-terminal**
- The **handle** of a sentential form is the **leftmost simple phrase**
- Simple phrases **cannot overlap** → leftmost simple phrase is **unambiguous**
- The sentential form **F(V Tail)**, F and V Tail are simple phrases and **F** is the **handle**

9

Handle and phrase

- **Handles** are important because they represent **individual derivation step**, which can be recognized by various **parsing techniques**
- $\{ [^i]^i \mid i \geq 1 \}$ is not regular and can be generated by the CFG as followed:
S → [T]
T → [T] | λ

10

Handle and phrase

- Although CFGs are widely used to define the **syntax of programming languages**, not all **syntactic rules** are **expressible** using CFGs → **variables must be declared before they are used** cannot be expressed in a CFG and there is no way to **transmit the exact set of variables that has been declared** to the body of a program (**static semantics** and **semantic routine with scope and type rules**)

11

Linguistics

Type 0: Natural Language

Type 1: Context Sensitive Language ($\alpha A \beta \rightarrow \alpha \delta \beta$)

Type 2: Context Free Language ($A \rightarrow \alpha \delta \beta$)

Type 3: Regular Language (No memory)

- **Context sensitive** and **type 0 grammars** are more powerful than CFG but they are far less useful → the **efficient parsers** for these extended grammar classes does not exist just **CFG has (Tradeoff)**

12

Useless part of CFG production

CFG are a **definitional mechanism** and may have error

$S \rightarrow A \mid B$

$A \rightarrow a$

$B \rightarrow Bb$ \rightarrow **derived no terminal string** (nonreduced)

$C \rightarrow c$ \rightarrow **unreachable** so is useless (cannot from S)

Finally forms the reduced grammar

$S \rightarrow A$

$A \rightarrow a$

13

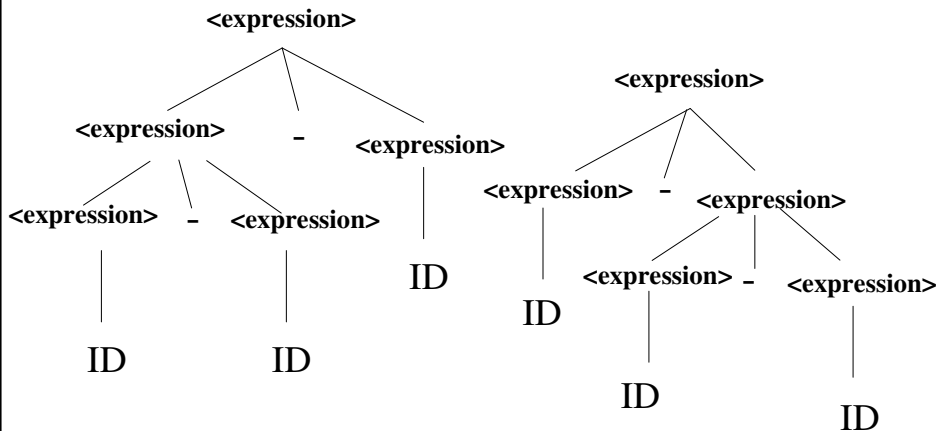
Ambiguous grammar

- A serious grammar flaw is that **a grammar allows a program** to have **two or more** different **parse trees**
- A grammar allows **different parse trees** for the same **terminal strings**
- Example
 $\langle \text{exp} \rangle \rightarrow \langle \text{exp} \rangle - \langle \text{exp} \rangle$ (**left recursive**)
 $\langle \text{exp} \rangle \rightarrow \text{ID}$

14

Parse tree

There are 2 different parse trees for ID-ID-ID



15

Grammar notation

- A grammar can produce two parse trees is called ambiguous.
- Ambiguous grammar are rare used because a **unique structure** cannot be guaranteed for **all inputs**, and a **unique translation**, guided by the parse tree structure, may not be obtained
- Normally restrict to **unambiguous grammars** to guarantee **unique structure**
- **Operator Precedence** is a problem.
- **Operator** in an arithmetic expression is generated **lower** in the parse tree can be used to indicate that it **has precedence over** an operator produced **higher up** in the tree.

16

Grammar notation

$\langle \text{Assign} \rangle \rightarrow \langle \text{id} \rangle := \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \langle \text{id_tail} \rangle \mid B \langle \text{id_tail} \rangle \mid \dots \mid Z \langle \text{id_tail} \rangle$

$\langle \text{id_tail} \rangle \rightarrow A \langle \text{id_tail} \rangle \mid B \langle \text{id_tail} \rangle \mid \dots \mid Z \langle \text{id_tail} \rangle \mid$
 $0 \langle \text{id_tail} \rangle \mid \dots \mid 9 \langle \text{id_tail} \rangle \mid \lambda$

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle \langle \text{addop} \rangle \langle \text{expr} \rangle \mid (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

$\langle \text{addop} \rangle \rightarrow + \mid -$

Syntactic ambiguity of language structures is a problem.

17

Grammar notation

A grammar can be written to separate the + and * so they are consistently in a **higher to lower ordering** in the parse tree.

$\langle \text{Assign} \rangle \rightarrow \langle \text{id} \rangle := \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \langle \text{id_tail} \rangle \mid B \langle \text{id_tail} \rangle \mid \dots \mid Z \langle \text{id_tail} \rangle$

$\langle \text{id_tail} \rangle \rightarrow A \langle \text{id_tail} \rangle \mid B \langle \text{id_tail} \rangle \mid \dots \mid Z \langle \text{id_tail} \rangle \mid$
 $0 \langle \text{id_tail} \rangle \mid \dots \mid 9 \langle \text{id_tail} \rangle \mid \lambda$

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid$

$\langle \text{expr} \rangle - \langle \text{term} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid$

$\langle \text{term} \rangle / \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

$\langle \text{factor} \rangle \rightarrow (\langle \text{expr} \rangle) \mid \langle \text{id} \rangle$

18

Transforming EBNF to standard form

For (each production $P=A \rightarrow \alpha [x_1 \dots x_n] \beta$) {

Create a new non-terminal, N . **Option**

Replace production P with $P'=A \rightarrow \alpha N \beta$

Add the productions: $N \rightarrow x_1 \dots x_n$ and $N \rightarrow \lambda$

}

For (each production $Q=B \rightarrow \gamma \{y_1 \dots y_m\} \delta$) {

Create a new non-terminal, M . **Option List**

Replace production Q with $Q'=B \rightarrow \gamma M \delta$

Add the productions: $M \rightarrow y_1 \dots y_m M$ and $M \rightarrow \lambda$

}

19

CFG Important Rules

Rule 2. Eliminate the 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

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

Or

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

20

Parser and Recognizer

- Given an **input string** as a **sequence of tokens**, is this input **syntactically valid**? (can it be generated from the grammar?) → a **recognizer**
- Is this input valid, what is its **structure** (**parse tree**)? An **algorithm** that answers this general question is a **parser**

21

Parser classification

- **Top-down parser** – it discovers the parse tree corresponding to a **token sequence** by starting at the **top of the tree** (**start symbol**) and then expanding it (via **predictions**) in a **depth-first** manner (LL → leftmost)
- A top-down parse corresponds to a **preorder** traversal of the parse tree (中左右)
- Top-down parsing is **predictive** because they always **predict the production** that is to be **matched** before matching begins

22

Parser classification

- **Bottom-up parser** discovers the structure of a parse tree by **beginning at its bottom** (the **leaves** of the tree, **terminal symbols**) and **determines the productions used to generate the leaves** (LR \rightarrow rightmost)
- Then the productions used to generate the **immediate parents** of the leaves are discovered
- The parser continues until it **reaches the production** used to **expand the start symbol**
- A bottom-up parser corresponds to a **postorder** traversal of the parse tree (左右中)

23

Example grammar

$\langle \text{program} \rangle \rightarrow \mathbf{begin} \langle \text{stmts} \rangle \mathbf{end} \$$

$\langle \text{stmts} \rangle \rightarrow \langle \text{stmt} \rangle ; \langle \text{stmt} \rangle$

$\langle \text{stmts} \rangle \rightarrow \lambda$

$\langle \text{stmt} \rangle \rightarrow \mathbf{simplestmt}$

$\langle \text{stmt} \rangle \rightarrow \mathbf{begin} \langle \text{stmts} \rangle \mathbf{end}$

For parsing the example

begin simplestmt; simplestmt; end\$

24

Explanation of parse tree

- A **bottom-up parser** proceeds by **discovering subtrees** and linking them into **increasingly larger trees**
- A **top-down parser** represents a **leftmost derivation** and produces a **leftmost parse**
- A **bottom-up parser** represents a **rightmost derivation** and produces a **rightmost parse**
- The best known **top-down** and **bottom-up parsing** strategies are called **LL** and **LR**.

27

Explanation of parse tree

- The first **L** states that the token sequence will **parse from left to right**. The second letter (**L** or **R**) states whether a **leftmost** or **rightmost parse** will be produced
- We may characterize the parsing technique by including **the number of lookahead symbols** (symbols beyond the current token) the parser may use to make parsing. **One-symbol lookahead** is most common → **LL(1)** or **LR(1)**

28

CFG Important Rules

Rule 2. Eliminate the 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

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

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

Or

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

$$Z \rightarrow x_i Z \mid \varepsilon, \quad i=1,2,\dots,m$$

29

Example

$$E \rightarrow E + T$$

$$E \rightarrow T$$

Can be changed to

$$E \rightarrow T \mid TZ$$

$$Z \rightarrow +T \mid +TZ$$

Or

$$A \rightarrow TZ$$

$$Z \rightarrow +TZ \mid \lambda$$

30

Example

BNF: (left-recursive grammar)

$\langle \text{expr} \rangle \rightarrow \langle \text{expr} \rangle + \langle \text{term} \rangle \mid \langle \text{expr} \rangle - \langle \text{term} \rangle \mid$
 $\langle \text{term} \rangle$

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{term} \rangle /$
 $\langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

EBNF: (non left-recursive grammar)

$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \{ (+ \mid -) \langle \text{term} \rangle \}$

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \}$

31

Example

$\langle \text{term} \rangle \rightarrow \langle \text{term} \rangle (* \mid /) \langle \text{factor} \rangle \mid \langle \text{factor} \rangle$

By rule 2. Eliminate left-recursive

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \mid \langle \text{factor} \rangle \langle Z \rangle$

$\langle Z \rangle \rightarrow (* \mid /) \langle \text{factor} \rangle \mid (* \mid /) \langle \text{factor} \rangle \langle Z \rangle$

By EBNF definition

$\langle Z \rangle \rightarrow (* \mid /) \langle \text{factor} \rangle \{ \langle Z \rangle \}$

$\langle Z \rangle \rightarrow (* \mid /) \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \}$

$\langle Z \rangle \rightarrow \{ (* \mid /) \langle \text{factor} \rangle \}$

By rule 5. $\langle Z \rangle$ is unit production

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ \langle Z \rangle \}$

$\langle \text{term} \rangle \rightarrow \langle \text{factor} \rangle \{ (* \mid /) \langle \text{factor} \rangle \}$

32