

Cú pháp chương trình

Cú pháp và ngữ nghĩa

- Cú pháp của ngôn ngữ lập trình: cấu trúc hình thức của chương trình như thế nào?
 - Cú pháp được định nghĩa bằng một văn phạm hình thức
- Ngữ nghĩa của ngôn ngữ lập trình: chương trình làm gì, hoạt động như thế nào?
 - Ngữ nghĩa chương trình khó định nghĩa hơn cú pháp.

Nội dung

- Văn phạm và ví dụ về cây cú pháp
- BNF và định nghĩa cây cú pháp
- Xây dựng văn phạm
- Cấu trúc câu và từ vựng
- Các dạng văn phạm khác

Văn phạm tiếng Anh

Một câu là một cấu trúc:
danh ngữ + động từ +
danh ngữ.

$$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$$

Một danh ngữ là ...

$$\langle NP \rangle ::= \langle A \rangle \langle N \rangle$$

Một động từ là ...

$$\langle V \rangle ::= \text{loves} \mid \text{hates} \mid \text{eats}$$

Một mạo từ là ...

$$\langle A \rangle ::= \text{a} \mid \text{the}$$

Một danh từ là ...

$$\langle N \rangle ::= \text{dog} \mid \text{cat} \mid \text{rat}$$

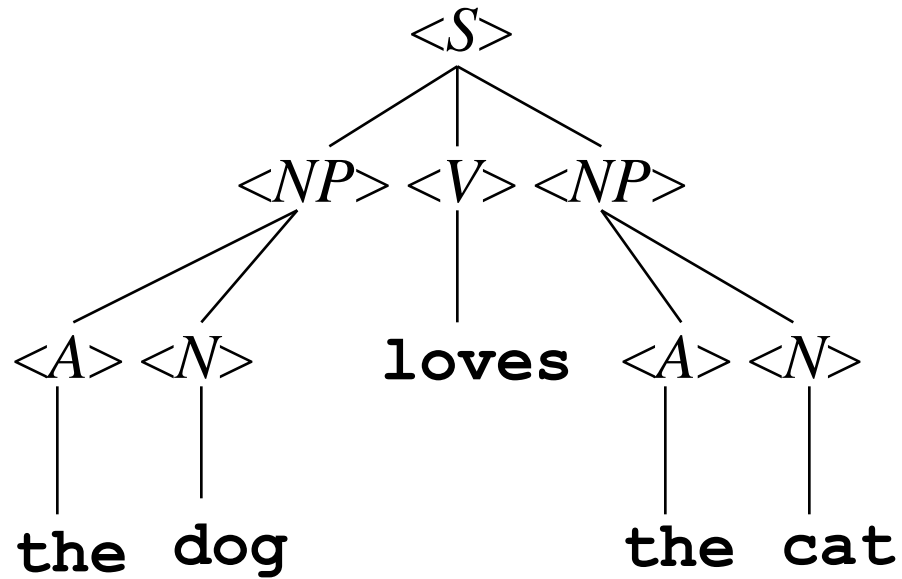
Văn phạm hoạt động như thế nào

- Văn phạm là một tập hợp các qui tắc chỉ ra cách thức xác định một *cây cú pháp*
- Cho $\langle S \rangle$ là gốc của cây
- Văn phạm xác định cách thức mà các nút con có thể được thêm vào cây
- Như vậy, qui tắc

$$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$$

cho biết có thể thêm $\langle NP \rangle$, $\langle V \rangle$, và $\langle NP \rangle$, theo trật tự trên làm các nút con của $\langle S \rangle$

Cây cú pháp



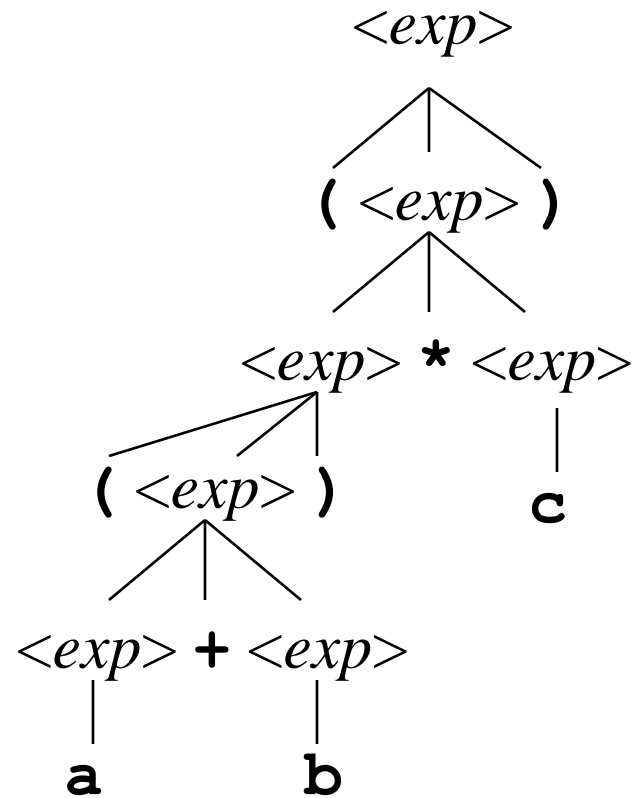
Văn phạm ngôn ngữ lập trình

$$\begin{aligned} \langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle &| \langle exp \rangle * \langle exp \rangle &| (\langle exp \rangle) \\ &| \mathbf{a} &| \mathbf{b} &| \mathbf{c} \end{aligned}$$

- Một biểu thức có thể là tổng của hai biểu thức, hoặc tích của hai biểu thức, hoặc một biểu thức trong cặp ngoặc đơn
- Hoặc có thể là một trong các biến **a**, **b** hay **c**

Cây cú pháp

((a+b) * c)



Tổng quan

- Văn phạm và ví dụ cây cú pháp
- BNF và định nghĩa cây cú pháp
- Xây dựng văn phạm
- Cấu trúc câu và từ vựng
- Các dạng văn phạm khác

CFG – Context Free Grammars (1950)

- Được phát triển bởi Noam Chomsky giữa thập niên 1950
- Language generators : miêu tả cú pháp của ngôn ngữ tự nhiên
- CFG là một văn phạm hình thức mà trong đó các luật là ở dạng $V \rightarrow w$

- A CFG consists of the following components:
- a set of *terminal symbols*, which are the characters of the alphabet that appear in the strings generated by the grammar.
- a set of *nonterminal symbols*, which are placeholders for patterns of terminal symbols that can be generated by the nonterminal symbols.
- a set of *productions*, which are rules for replacing (or rewriting) nonterminal symbols (on the left side of the production) in a string with other nonterminal or terminal symbols (on the right side of the production).
- a *start symbol*, which is a special nonterminal symbol that appears in the initial string generated by the grammar.

- **A CFG for Arithmetic Expressions**

- An example grammar that generates strings representing arithmetic expressions with the four operators $+$, $-$, $*$, $/$, and numbers as operands is:

- $\langle \text{expression} \rangle \rightarrow \text{number}$
- $\langle \text{expression} \rangle \rightarrow (\langle \text{expression} \rangle)$
- $\langle \text{expression} \rangle \rightarrow \langle \text{expression} \rangle + \langle \text{expression} \rangle$
- $\langle \text{expression} \rangle \rightarrow \langle \text{expression} \rangle - \langle \text{expression} \rangle$
- $\langle \text{expression} \rangle \rightarrow \langle \text{expression} \rangle * \langle \text{expression} \rangle$
- $\langle \text{expression} \rangle \rightarrow \langle \text{expression} \rangle / \langle \text{expression} \rangle$

- The only nonterminal symbol in this grammar is $\langle \text{expression} \rangle$, which is also the start symbol. The terminal symbols are $\{+, -, *, /, (,), \text{number}\}$.

BNF - Backus Normal Form (1959)

- Được tạo ra bởi John Backus để miêu tả cho Algol 58
- BNF tương đương với CFG (context-free grammars)
- Meta language là một dạng ngôn ngữ dùng để miêu tả cho một ngôn ngữ khác

Ví dụ :mô tả địa chỉ bưu điện ở mỹ :

`<postal-addr> ::= <name-part> <street-addr> <zip-part>`

(Một địa chỉ gồm 3 phần :Tên, địa chỉ và mã vùng)

`<name-part> ::=`

`<personal-part> <name-part> |`

`<personal-part> <last-name> <suffix-part>`

(Tên gồm 2 phần :phần riêng ,tên hoặc phần riêng, tên và xưng hô)

`<personal-part> ::= <first-name> | <initial> "."`

(Phần riêng :tên hoặc ký tự đầu)

`<street-addr> ::= <house-num> <street-name>`

(Địa chỉ :số nhà và tên đường)

`<zip-part> ::=`

`<town-name> "," <state-code> <ZIP-code>`

(Zip-part :tên thành phố,state-code, zipcode

`<suffix-part> ::= "Sr." | "Jr." | <roman-numeral>`

Ký hiệu bắt đầu

$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$

Một luật sinh

$\langle NP \rangle ::= \langle A \rangle \langle N \rangle$

$\langle V \rangle ::= \text{loves} \mid \text{hates} \mid \text{eats}$

$\langle A \rangle ::= \text{a} \mid \text{the}$

$\langle N \rangle ::= \text{dog} \mid \text{cat} \mid \text{rat}$

Các ký hiệu không kết thúc

tokens

Định nghĩa văn phạm BNF

- Một văn phạm BNF gồm có 4 phần:
 - Một tập hợp các *tokens*
 - Một tập hợp các *ký hiệu không kết thúc*
 - Một *ký hiệu bắt đầu*
 - Một tập các *luật sinh (productions)*

Định nghĩa văn phạm BNF

- Các *tokens* là những đơn vị nhỏ nhất của cú pháp:
 - Token có thể là một chuỗi gồm một hay nhiều ký tự
 - Token là nguyên tố: không phân tách nhỏ hơn
- Các *ký hiệu không kết thúc*:
 - Các ký hiệu không kết thúc là các chuỗi trong cặp ngoặc <>, chẳng hạn <NP>
 - Văn phạm xác định cách thức mở rộng các ký hiệu không kết thúc thành các token

Định nghĩa văn phạm BNF

- *Ký hiệu bắt đầu* là một ký hiệu không kết thúc đặc biệt, dùng làm nút gốc của cây phân tích
- *Các luật sinh* là các luật xây dựng cây cú pháp

Định nghĩa văn phạm BNF

- Mỗi luật sinh có một phần bên trái, một bộ tách $::=$, và một phần bên phải
 - Phần bên trái là một đơn ký hiệu không kết thúc
 - Phần bên phải là một dãy, mỗi thành phần của dãy này có thể là một token hay một ký hiệu không kết thúc
- Mỗi luật sinh xác định một khả năng để xây dựng cây cú pháp.

Định nghĩa văn phạm BNF

- Khi có nhiều hơn một luật sinh có cùng một phần bên trái giống nhau thì có thể dùng dạng thu gọn
- Văn phạm BNF có thể đặt phần bên trái, bộ tách $::=$, và một danh sách các phần bên phải – mỗi thành phần trong danh sách đó cách nhau bởi ký hiệu $|$

Ví dụ

$$\begin{aligned} \langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle &| \langle exp \rangle * \langle exp \rangle | (\langle exp \rangle) \\ &| \mathbf{a} | \mathbf{b} | \mathbf{c} \end{aligned}$$

Chú ý rằng có 6 luật sinh trong văn phạm trên
Văn phạm trên tương đương với:

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle$$

$$\langle exp \rangle ::= \langle exp \rangle * \langle exp \rangle$$

$$\langle exp \rangle ::= (\langle exp \rangle)$$

$$\langle exp \rangle ::= \mathbf{a}$$

$$\langle exp \rangle ::= \mathbf{b}$$

$$\langle exp \rangle ::= \mathbf{c}$$

Empty

- Ký hiệu $\langle empty \rangle$ được dùng để chỉ ra là văn phạm không sản sinh bất cứ gì ở vị trí nó xuất hiện
- Ví dụ, văn phạm sau định nghĩa một cấu trúc if-then với một phần else tùy chọn:

$$\begin{aligned}\langle if-stmt \rangle &::= \mathbf{if} \langle expr \rangle \mathbf{then} \langle stmt \rangle \langle else-part \rangle \\ \langle else-part \rangle &::= \mathbf{else} \langle stmt \rangle \mid \langle empty \rangle\end{aligned}$$

Cây cú pháp

- Để xây dựng một cây cú pháp, thiết lập ký hiệu bắt đầu ở nút gốc
- Thêm các ký hiệu không kết thúc vào, *dựa trên các luật sinh trong văn phạm*
- Kết thúc khi tất cả các nút lá là token
- Đọc các nút lá từ trái sang phải, đây chính là chuỗi phát sinh từ cây cú pháp

Ứng dụng

$$\begin{aligned} \langle exp \rangle ::= & \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle) \\ & \mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c} \end{aligned}$$

Vẽ cây phân tích cú pháp cho các chuỗi sau:

a+b

a*b+c

(a+b)

(a+ (b))

Định nghĩa ngôn ngữ

- Văn phạm được dùng để định nghĩa cú pháp của các ngôn ngữ lập trình
- Ngôn ngữ được định nghĩa bằng một văn phạm là một tập các chuỗi có thể phát sinh từ một cây cú pháp nào đó của văn phạm

Nội dung

- Văn phạm và ví dụ cây cú pháp
- BNF và định nghĩa cây cú pháp
- Xây dựng văn phạm
- Cấu trúc câu và từ vựng
- Các dạng văn phạm khác

Xây dựng văn phạm

- Ví dụ: các khai báo trong ngôn ngữ Java: tên kiểu, danh sách các biến cách nhau bởi dấu phẩy, và một dấu chấm phẩy
- Mỗi biến có thể được theo sau bởi một bộ khởi tạo:

```
float a;  
boolean a,b,c;  
int a=1, b, c=1+2;
```

Ví dụ (tiếp theo)

- Nếu chưa kể đến bộ khởi tạo:

$\langle var\text{-}dec \rangle ::= \langle type\text{-}name \rangle \langle declarator\text{-}list \rangle ;$

- Các kiểu dữ liệu cơ bản:

$\langle type\text{-}name \rangle ::= \text{boolean} \mid \text{byte} \mid \text{short} \mid \text{int}$
 $\qquad \qquad \qquad \mid \text{long} \mid \text{char} \mid \text{float} \mid \text{double}$

$\langle declarator\text{-}list \rangle ::= \langle declarator \rangle$
 $\qquad \qquad \qquad \mid \langle declarator \rangle , \langle declarator\text{-}list \rangle$

$\langle declarator \rangle ::= \langle variable\text{-}name \rangle$
 $\qquad \qquad \qquad \mid \langle variable\text{-}name \rangle = \langle expr \rangle$

Nội dung

- Văn phạm và ví dụ cây cú pháp
- BNF và định nghĩa cây cú pháp
- Xây dựng văn phạm
- **Cấu trúc câu và từ vựng**
- Các dạng văn phạm khác

Token

- Các token là những mẫu văn bản trong chương trình mà chúng không thể được chia nhỏ hơn nữa
- Định danh (**count**), từ khóa (**if**), toán tử (**==**), hằng số (**123.4**), etc.
- Chương trình máy tính là một dãy các ký tự
- Làm thế nào một dãy các ký tự có thể phân tích thành các token?

Cấu trúc từ vựng và cấu trúc câu

- Văn phạm định nghĩa *cấu trúc từ vựng*: một tập tin văn bản được phân tách thành các token như thế nào?
- Văn phạm định nghĩa *cấu trúc câu*: chương trình được xây dựng như thế nào từ một chuỗi các token ?

Một văn phạm cho cấu trúc từ vựng và cấu trúc câu

- Theo quan điểm này, sử dụng một văn phạm chung cho cấu trúc câu và cấu trúc từ
- Không hiệu quả trong thực tế: các khoảng trắng hay các ghi chú làm cho văn phạm khó đọc và rườm rà

$\langle \text{if-stmt} \rangle ::= \mathbf{if} \langle \text{white-space} \rangle \langle \text{expr} \rangle \langle \text{white-space} \rangle$
 $\quad \mathbf{then} \langle \text{white-space} \rangle$
 $\quad \langle \text{stmt} \rangle \langle \text{white-space} \rangle \langle \text{else-part} \rangle$
 $\langle \text{else-part} \rangle ::= \mathbf{else} \langle \text{white-space} \rangle \langle \text{stmt} \rangle \mid \langle \text{empty} \rangle$

Hai văn phạm tách biệt

- Thông thường, các ngôn ngữ sử dụng hai văn phạm tách biệt
 - Một văn phạm chỉ ra cách thức cấu trúc các token từ một văn bản các ký tự
 - Một văn phạm chỉ ra cách xây dựng cây cú pháp từ một dãy các token

<program-file> ::= <end-of-file> | <element> <program-file>

<element> ::= <token> | <one-white-space> | <comment>

<one-white-space> ::= <space> | <tab> | <end-of-line>

<token> ::= <identifier> | <operator> | <constant> | ...

Lưu ý

- Các ngôn ngữ lập trình trước đây thường không tách bạch cấu trúc từ vựng với cấu trúc câu
 - Trong các thể hệ Fortran và Algol đầu tiên, khoảng trắng có thể ở bất kỳ vị trí nào, ngay cả ở giữa một từ khóa
 - Các ngôn ngữ khác như PL/I cho phép các từ khóa được dùng làm các định danh

Lưu ý

- Một số ngôn ngữ có định dạng cố định cho cấu trúc từ vựng – theo các vị trí cột
 - Mỗi dòng một phát biểu
 - Một vài cột đầu tiên dành cho nhãn của phát biểu
- Ví dụ: một vài thế hệ đầu của Fortran, Cobol, và Basic
- Hầu hết các ngôn ngữ lập trình hiện nay có định dạng tự do cho cấu trúc từ vựng

Nội dung

- Văn phạm và ví dụ cây cú pháp
- BNF và định nghĩa cây cú pháp
- Xây dựng văn phạm
- Cấu trúc câu và từ vựng
- Các dạng văn phạm khác

Các dạng văn phạm khác

- Các biến thể của BNF
- Các biến thể của EBNF
- Các sơ đồ cú pháp

Các biến thể của BNF

- Một số dùng \rightarrow hay $=$ thay vì $::=$
- Một số cho phép dùng cặp dấu nháy đơn bao các tokens, ví dụ phân biệt ‘|’ như là một token với | như là meta-symbol

Các biến thể của BNF

- Cú pháp bổ sung cho BNF:
 - $\{x\}$ nghĩa là không hay lặp lại x nhiều lần
 - $[x]$ có nghĩa x là tùy chọn (i.e. $x \mid \langle empty \rangle$)
 - $()$ cho nhóm
 - $|$ chọn một trong số các chọn lựa
 - Cặp dấu ngoặc kép bao ký hiệu để phân biệt với các meta-symbols

EBNF Examples

$\langle \text{if-stmt} \rangle ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle [\text{else } \langle \text{stmt} \rangle]$

$\langle \text{stmt-list} \rangle ::= \{ \langle \text{stmt} \rangle ; \}$

$\langle \text{thing-list} \rangle ::= \{ (\langle \text{stmt} \rangle \mid \langle \text{declaration} \rangle) ; \}$

- Các mở rộng của BNF được gọi là Extended BNF: EBNF
- Có nhiều biến thể của BNF

Các sơ đồ cú pháp

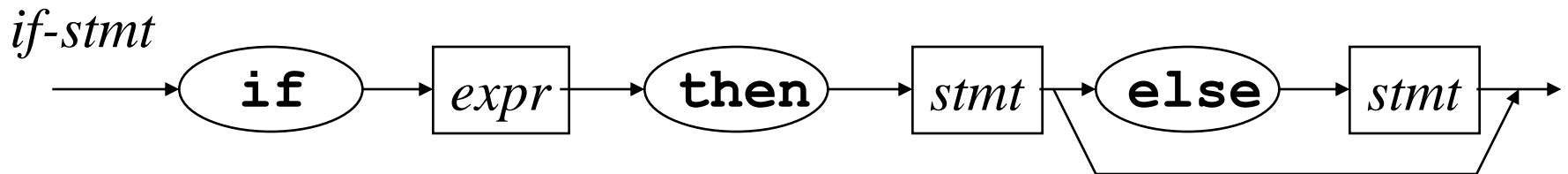
- Văn phạm EBNF
- Một luật sinh đơn giản là một chuỗi các ký hiệu kết thúc và không kết thúc:

$\langle if-stmt \rangle ::= \mathbf{if} \langle expr \rangle \mathbf{then} \langle stmt \rangle \mathbf{else} \langle stmt \rangle$



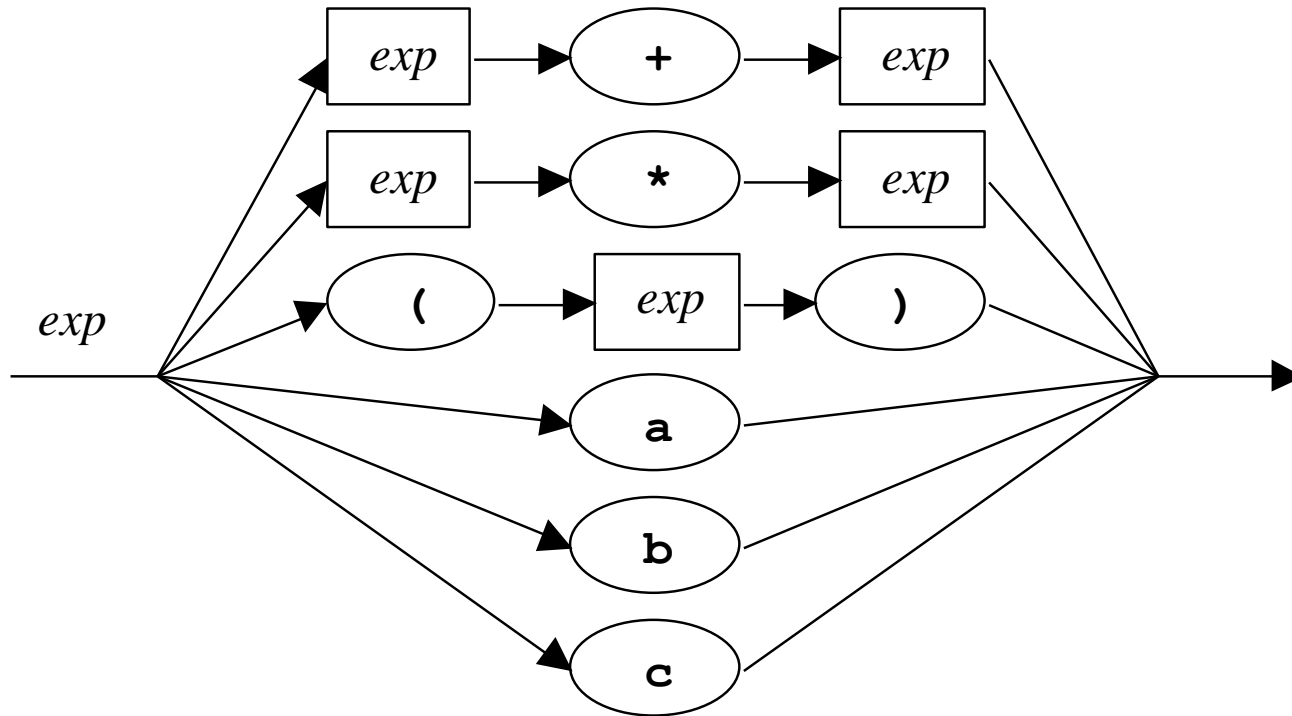
Đường vòng

$\langle if\text{-}stmt \rangle ::= \mathbf{if} \ \langle expr \rangle \ \mathbf{then} \ \langle stmt \rangle \ [\mathbf{else} \ \langle stmt \rangle]$



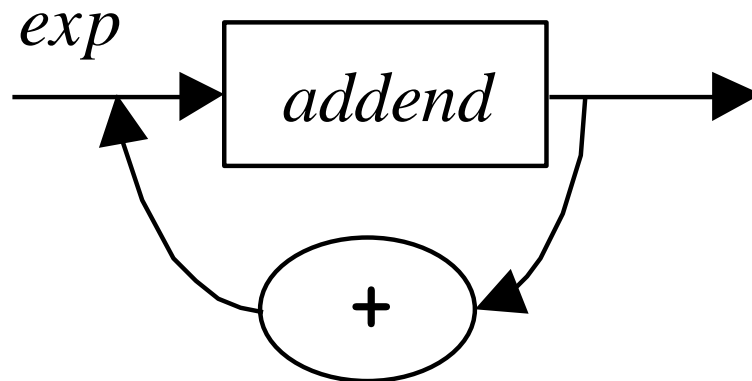
Rẽ nhánh

$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle)$
 $\mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c}$



Vòng lặp

$\langle exp \rangle ::= \langle addend \rangle \{ + \langle addend \rangle \}$



Ví dụ

WhileStatement:

while (*Expression*) *Statement*

DoStatement:

do *Statement* while (*Expression*) ;

ForStatement:

for (*ForInit*_{opt} ; *Expression*_{opt} ; *ForUpdate*_{opt})
 Statement

Cho luật sinh sau

$\langle \text{assign} \rangle \rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\langle \text{id} \rangle \rightarrow A \mid B \mid C$

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

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

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

Example: Derivation and parser tree of sentence $A = B + C * A$

$\langle \text{assign} \rangle$

$\Rightarrow \langle \text{id} \rangle = \langle \text{expr} \rangle$

$\Rightarrow A = \langle \text{expr} \rangle$

$\Rightarrow A = \langle \text{expr} \rangle + \langle \text{term} \rangle$

$\Rightarrow A = \langle \text{term} \rangle + \langle \text{term} \rangle$

$\Rightarrow A = \langle \text{factor} \rangle + \langle \text{term} \rangle$

$\Rightarrow A = \langle \text{id} \rangle + \langle \text{term} \rangle$

$\Rightarrow A = B + \langle \text{term} \rangle$

$\Rightarrow A = B + \langle \text{term} \rangle * \langle \text{factor} \rangle$

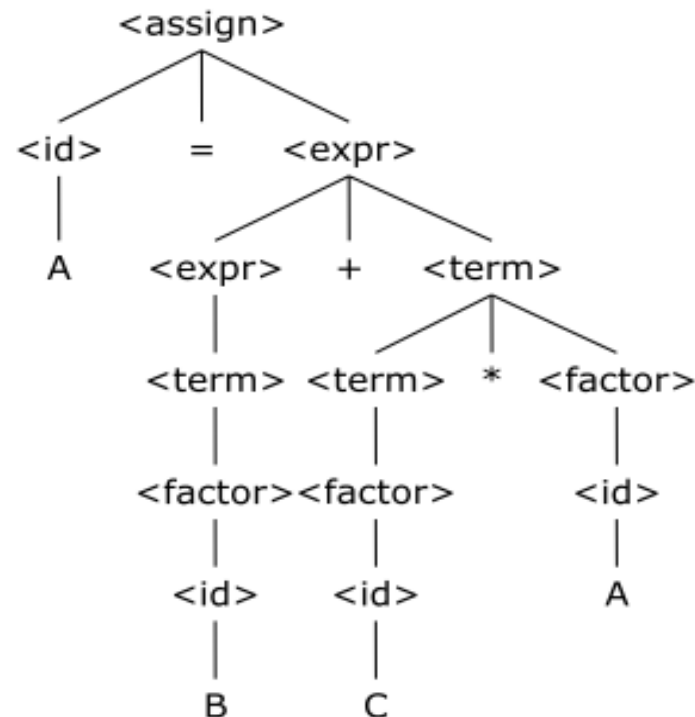
$\Rightarrow A = B + \langle \text{factor} \rangle * \langle \text{factor} \rangle$

$\Rightarrow A = B + \langle \text{id} \rangle * \langle \text{factor} \rangle$

$\Rightarrow A = B + C * \langle \text{factor} \rangle$

$\Rightarrow A = B + C * \langle \text{id} \rangle$

$\Rightarrow A = B + C * A$



- **Ba cách xây dựng bộ phân tích từ vựng (Scanner)**

- Dùng phần mềm phát sinh tự động(đầu vào là một văn phạm chính quy)
- Xây dựng sơ đồ trạng thái DFA
- Xây dựng bảng thực thi dòng, cột

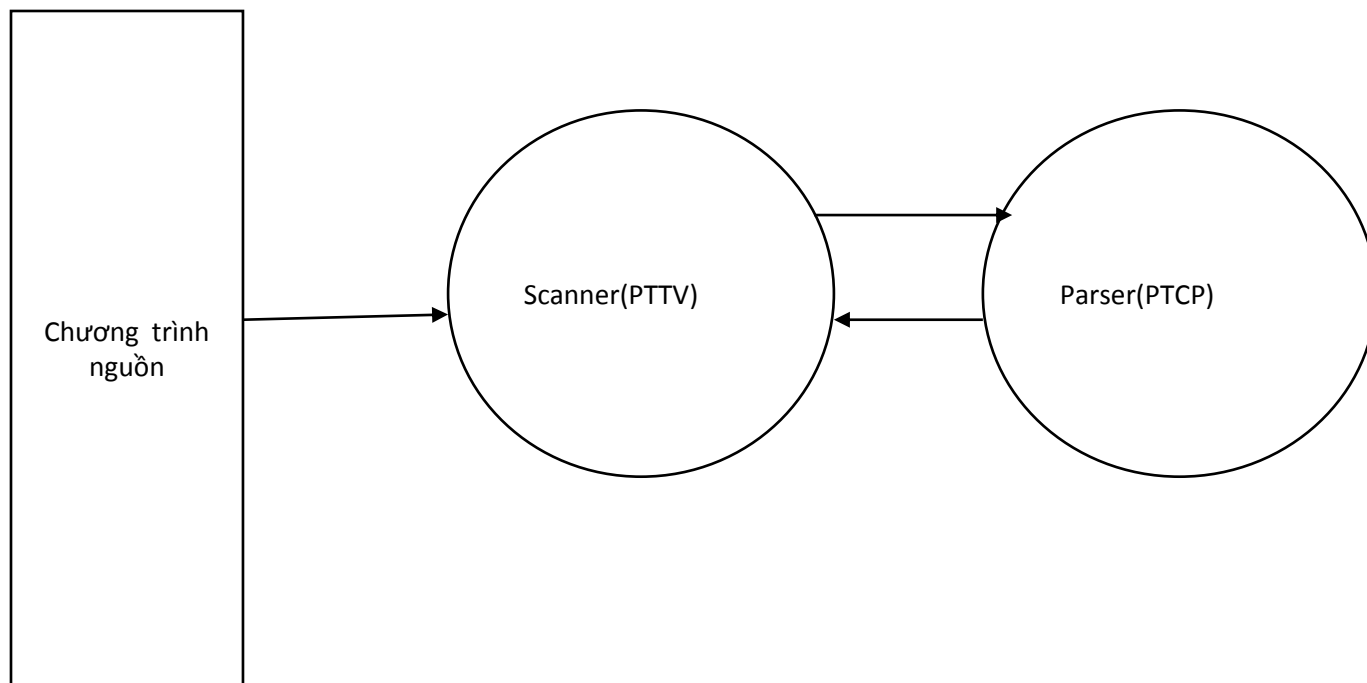
- Luật tìm lexeme dài nhất

- Scanner chỉ trả về token khi đọc đến ký tự kết thúc token(thường là khoảng trắng)
- Trong một số trường hợp Scanner phải đọc nhiều hơn một ký tự để nhận biết đã kết thúc token

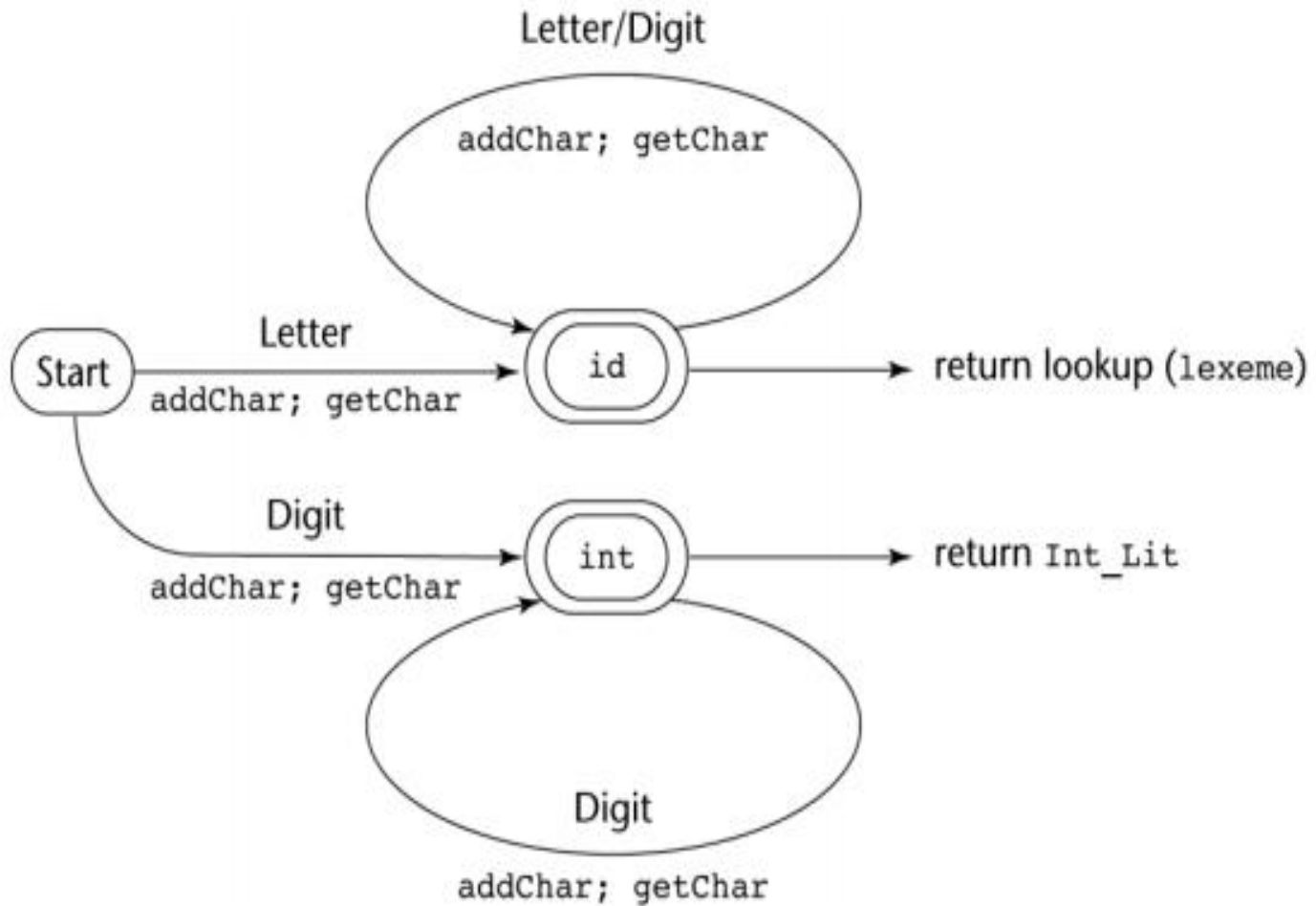
Example

`sum = oldsum - value / 100;`

<i>Token</i>	<i>Lexeme</i>
IDENT	sum
ASSIGN_OP	=
IDENT	oldsum
SUBTRACT_OP	-
IDENT	value
DIVISION_OP	/
INT_LIT	100
SEMICOLON	;



State Diagram



Copyright © 2006 Addison-Wesley. All rights reserved.

1-15

Lexical Analysis – Implementation

```
int lex() {  
    getChar();  
    switch (charClass) {  
        case LETTER:  
            addChar();  
            getChar();  
            while (charClass == LETTER || charClass == DIGIT) {  
                addChar();  
                getChar();  
            }  
            return lookup(lexeme);  
            break;  
        ...  
    }
```

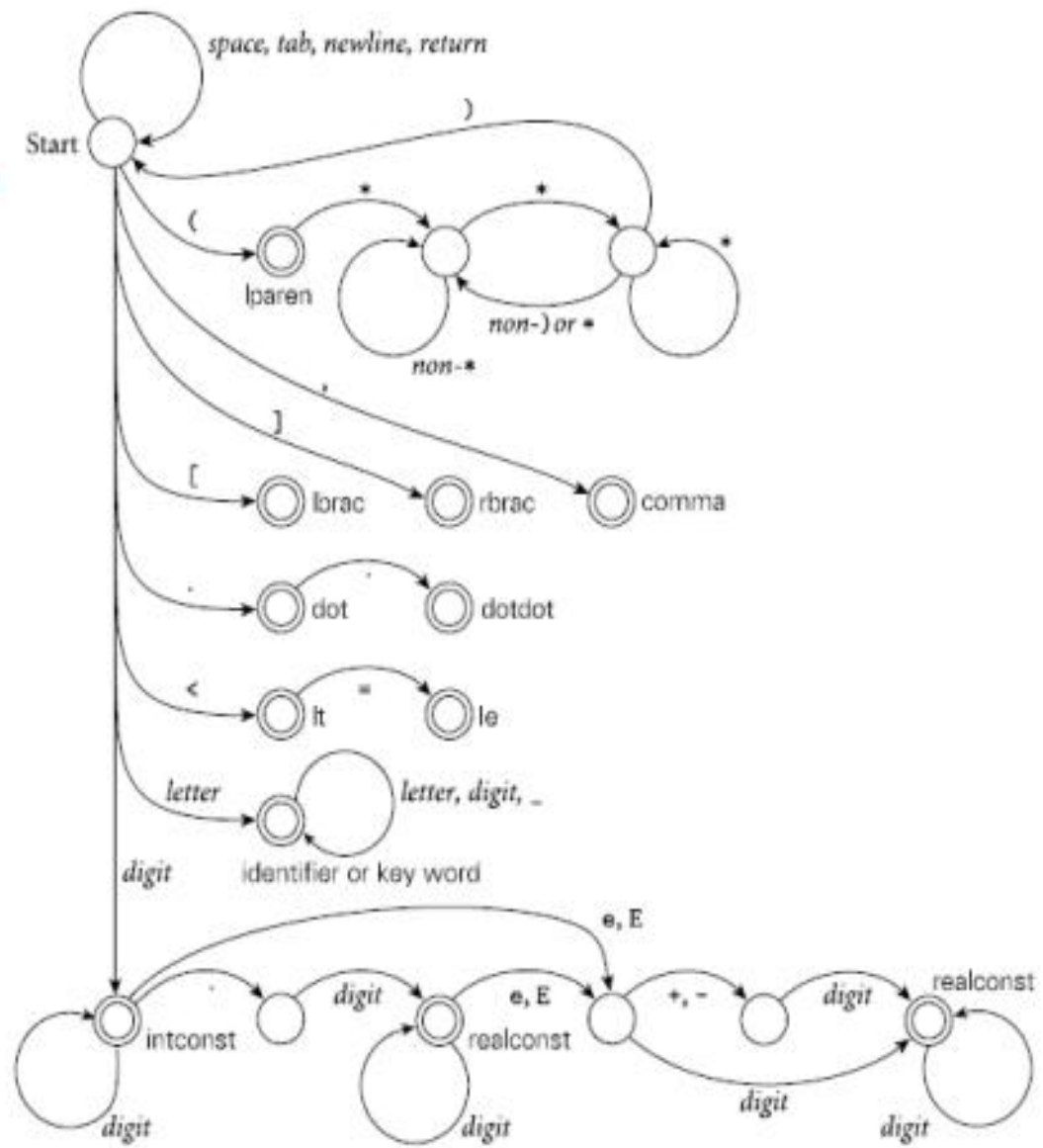
Copyright © 2006 Addison-Wesley. All rights reserved.

1-16

Lexical Analysis – Implementation

```
    case DIGIT:
        addChar();
        getChar();
        while (charClass == DIGIT) {
            addChar();
            getChar();
        }
        return INT_LIT;
    } /* End of switch */
} /* End of function lex() */
```

State Diagram (Pascal)



•Phân tích cú pháp

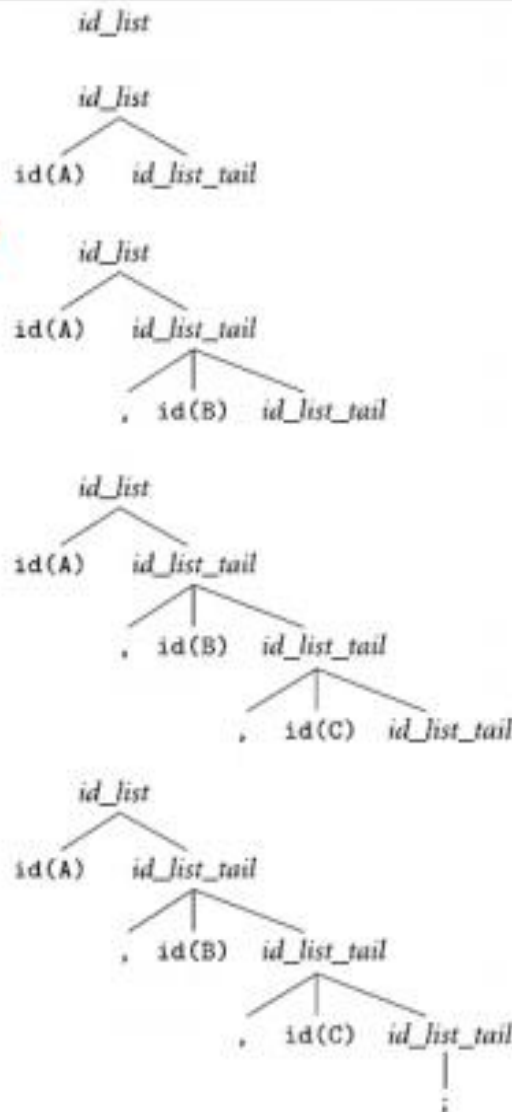
I. Nhiệm vụ

- Tìm tất cả các lỗi cú pháp , đưa ra các thông báo lỗi
- Xây dựng cây phân tích cú pháp (lưu vết) cho chương trình

II. Có hai kỹ thuật phân tích

- Top-down : xây dựng cây phân tích cú pháp bắt đầu từ gốc
- Bottom up : xây dựng cây phân tích cú pháp bắt đầu từ lá

Top-down and bottom-up parsing of the input string A, B, C;



$id_list \rightarrow id\ id_list_tail$
 $id_list_tail \rightarrow ,\ id\ id_list_tail$
 $id_list_tail \rightarrow ;$

