

## บทที่ 4

# วากยสัมพันธ์ (Syntax)

- 4.1 โครงสร้างเชิงศัพท์ ของภาษาโปรแกรม
- 4.2 ไวยากรณ์ไม่พืงบริบท และมีเอ็นเอฟ
- 4.3 ต้นไม้วิภัช และต้นไม้วากยสัมพันธ์แบบนามธรรม
- 4.4 ความกำกวม การสลับที่ และการทำก่อน
- 4.5 อีบีเอ็นเอฟ และ แผนภาพวากยสัมพันธ์
- 4.6 เทคนิคการวิเคราะห์กระจาย และเครื่องมือที่ใช้
- 4.7 คำศัพท์ วากยสัมพันธ์ อรรถศาสตร์  
แบบฝึกหัด

## บทที่ 4

### วากยสัมพันธ์ (Syntax)

วากยสัมพันธ์ หมายถึง โครงสร้างของภาษา ในบทที่ 1 เราได้ให้ข้อสังเกตความแตกต่างระหว่าง วากยสัมพันธ์ และอรรถศาสตร์ ของภาษาโปรแกรมแล้ว ในอดีตของภาษาโปรแกรม วากยสัมพันธ์ และอรรถศาสตร์ของภาษา ทั้งคู่ อธิบายด้วย ภาษาอังกฤษอย่างยาว และตัวอย่างจำนวนมาก ขณะที่ อรรถศาสตร์ ของภาษายังคงถูกอธิบายด้วยภาษาอังกฤษอยู่ ความรู้ขั้นสูงมากในภาษาโปรแกรม ได้พัฒนา ระบบแบบทางการ สำหรับการอธิบายวากยสัมพันธ์ ซึ่งขณะนี้ ส่วนใหญ่ใช้กันเป็นสากล ในช่วงปี ค.ศ. 1950s Noam Chomsky พัฒนา ความคิดของ ไวยากรณ์ไม่พึ่งบริบท ขึ้น ส่วน John Backus และผู้ร่วมงานของเขา Peter Naur ได้พัฒนาระบบเชิงสัญลักษณ์ (notational system) สำหรับ อธิบาย ไวยากรณ์เหล่านี้ ซึ่งครั้งแรก นำมาใช้เพื่ออธิบายวากยสัมพันธ์ของ ภาษา Algol60 สิ่งเหล่านี้เรียกว่า Backus - Naur Form หรือเรียกสั้นๆ ว่า BNFs ต่อมาได้นำมาใช้ ใน บทนิยาม ของ ภาษาโปรแกรมอื่นๆ จำนวนมาก เช่น Pascal, Modula-2, C และ Ada จริงๆ แล้ว นักเขียนโปรแกรมสมัยใหม่ และนักวิทยาศาสตร์คอมพิวเตอร์ จำเป็นต้องทราบว่าจะ อ่าน ดีความ และประยุกต์ใช้ การอธิบาย BNF ของวากยสัมพันธ์ภาษาได้อย่างไร

BNFs เหล่านี้ เกิดขึ้น ด้วย ความหลากหลายไม่มากนัก ใน รูปแบบพื้นฐานสามอย่างคือ : BNF เดิม (original BNF), ส่วนขยายของ BNF (extended BNF (EBNF)) ผู้ที่ทำให้เกิดความนิยมคือ Niklaus Wirth, และรูปแบบที่สามคือ แผนภาพวากยสัมพันธ์ (syntax diagrams)

#### 4 . 1 โครงสร้างศัพท์ ของ ภาษาโปรแกรม

(Lexical structure of programming languages)

โครงสร้างศัพท์ของภาษาโปรแกรม หมายถึง โครงสร้างของคำ หรือ โทเค็น  
ของมัน (The lexical structure of a programming language is the

structure of its words or tokens.)

โครงสร้างศัพท์ สามารถ พิจารณาแยกต่างหาก จาก โครงสร้างวากยสัมพันธ์ (syntactic structure) แต่มันเกี่ยวข้องกันมาก ในบางกรณี (ขึ้นอยู่กับ การออกแบบภาษา) ไม่สามารถ แยกส่วนต่างหากจากวากยสัมพันธ์ได้ โดยปกติ ขั้นการกวาดตรวจ (scanning phase) ของตัวแปลภาษา จัดกลุ่ม ลำดับของอักขระต่าง ๆ จาก โปรแกรม อินพุท ให้เป็น โทเค้นต่าง ๆ ซึ่งหลังจากนั้น ถูกประมวลผล โดย ขั้นการวิเคราะห์กระจาย (parsing phase) ซึ่งจะบอกความถูกต้องของ โครงสร้างวากยสัมพันธ์

โทเค้น แบ่งออกเป็นหลายชนิด ได้แก่

คำสงวน (reserved words) บางครั้ง เรียกว่า คำหลัก (keywords) เช่นคำว่า "begin", "if" และ "while"

ค่าคงที่ (constants หรือ literals) เช่น 42 เป็นค่าคงที่ตัวเลข (numeric constant) หรือ "hello" เป็นค่าคงที่สายอักขระ (string constant)

สัญลักษณ์พิเศษ (special symbols) เช่น ";", "<=" หรือ "+"

ไอดีไฟเอร์ (identifiers) เช่น x24, average, monthly-balance, หรือ write

คำสงวน มีเป็นจำนวนมาก และเนื่องจากไอดีไฟเอร์ จะมีสายอักขระ เหมือนกับ คำสงวนไม่ได้ ดังนั้น ในภาษา Pascal การประกาศตัวแปร ข้างล่างนี้ ไม่ถูกต้อง เพราะว่า "if" เป็นคำสงวน

```
var if : integer;
```

ในบางภาษา ไอดีไฟเอร์ จะมีขนาดสูงสุดคงที่ (fixed maximum size) ในขณะที่ไอดีไฟเอร์ของภาษาใหม่ ๆ ส่วนใหญ่มีความยาวเท่าไรก็ได้ (arbitrary length) บางครั้ง เมื่อ ไอดีไฟเอร์ มีขนาดอย่างไรก็ได้ แต่เฉพาะอักขระ 6 ตัว หรือ 8 ตัวแรกเท่านั้น ซึ่งมีนัยสำคัญ (significant) สิ่งนี้ ประกันความสับสนให้กับนักเขียนโปรแกรม

ปัญหาอย่างหนึ่งเกิดขึ้น ใน การบอก การจบ (end) ของ โทเค้นความยาว

แปรได้ (variable-length token) เช่น ไอเดนטיפายเออร์ และการแยก ไอเดน-  
טיפายเออร์ ออกจาก คำสงวน

ตัวอย่างเช่น ลำดับของตัวอักษร

doif

ในโปรแกรม อาจจะเป็น คำสงวนสองคำ คือ "do" และ "if" หรืออาจจะ  
เป็น ไอเดนטיפายเออร์ doif ก็ได้

ในทำนองเดียวกัน สายอักษร x12 อาจจะเป็นไอเดนטיפายเออร์ หนึ่งตัว  
คือ x12 หรือ อาจจะเป็นไอเดนטיפายเออร์ x และค่าคงที่ตัวเลข 12

เพื่อจัดความกำกวมนี้ จึงมีข้อตกลงมาตรฐาน (standard convention)  
ให้ใช้ คือ หลักของสายอักขระย่อยความยาวสูงสุด (principle of longest sub-  
string) ในการหา โทเค้น : ที่แต่ละจุด สายของอักขระที่เป็นไปได้ความยาว  
สูงสุด รวมเข้าด้วยกันเป็น หนึ่งโทเค้น สิ่งนี้หมายความว่า doif และ x12 เป็น  
ไอเดนטיפายเออร์เสมอ

ตัวอักษรแทรก (intervening characters) เช่น อักขระว่าง (blank)  
จะทำให้เกิดความแตกต่าง ดังนั้น ในภาษาส่วนใหญ่

do if

จึงไม่ใช่ ไอเดนטיפายเออร์ แต่เป็น คำสงวนสองคำ คือ "do" และ "if"

รูปแบบของโปรแกรม สามารถกระทบกับ วิธีที่จะรู้จำโทเค้น (The format  
of a program can affect the way tokens are recognized.)

ตัวอย่างเช่น หลักของสายอักขระย่อยความยาวสูงสุด ต้องการให้โทเค้นต่างๆ แยกกัน  
ด้วย ตัวคั่นโทเค้น (token delimiters) หรือ white space

การจบ ของ หนึ่งบรรทัดของข้อความ (text) อาจมีความหมายได้สองอย่าง  
คือ : มันอาจเป็น ว่างหนึ่งที่ หรือ อาจหมายถึง การจบของเอนทิตีเชิงโครงสร้าง

การย่อหน้า สามารถนำมาใช้ในภาษาโปรแกรม เพื่อบอกโครงสร้าง  
ภาษาซึ่งมีรูปแบบอิสระ หมายถึง ภาษา ซึ่งรูปแบบของมันไม่มีผลกระทบ บน โครงสร้าง  
โปรแกรม (A free-format language is one in which format has no

effect on the program structure.) และยังคงเป็นไปตามหลักของสายอักขระ  
ย่อยความสูงสุด ภาษาโปรแกรมสมัยใหม่ส่วนใหญ่ เป็นภาษารูปแบบอิสระ แต่ภาษาโปรแกรม  
ส่วนน้อย มี ข้อจำกัดรูปแบบอย่างสำคัญ ซึ่งเรียกว่า รูปแบบคงที่ (fixed format) ซึ่ง  
โทเคนทั้งหมดต้อง เกิดขึ้น ในตำแหน่งที่กำหนดไว้แล้ว ในหน้า

กระดาษ เช่นภาษา RPG, FORTRAN, และ COBOL เป็นต้น

FORTRAN เป็นตัวอย่างที่สำคัญของ ภาษาซึ่ง ผ่าฝืน (violates) รูปแบบ  
หลายอย่าง และ ข้อตกลงเรื่องโทเคน

ตัวอย่าง ภาษา FORTRAN

```
DO 99 I = 1.10 ===== (1)
```

ข้อความสั่งนี้ มีความหมายเหมือนกับ ภาษา Pascal

```
DO99I := 1.10
```

นุดอีกอย่างหนึ่งคือ เป็นการกำหนดค่า 1.1 ให้กับตัวแปรจำนวนจริง DO99I

ในทางตรงกันข้าม ข้อความสั่ง ภาษา FORTRAN ข้างล่างนี้

```
DO 99 I =1, 10 ===== (2)
```

มีความหมายเหมือนกับ ข้อความสั่ง ภาษา Pascal ดังนี้

```
for I := 1 to 10 do
```

นั่นคือ เริ่มการวนซ้ำ (loop) โดยกำหนดขอบเขตให้กับ ดรรชนี I ดังนี้

ข้อความสั่งภาษา FORTRAN ชุดแรกจึงประกอบด้วย โทเคน 3 ตัว คือ ไอเดนติไฟเออร์,  
ตัวดำเนินการกำหนดค่า ("=") และ ค่าคงที่จำนวนจริง 1.1 แต่ตรงกันข้าม ส่วนข้อ  
ความสั่ง ภาษา FORTRAN ชุดที่สอง ประกอบด้วย โทเคน 7 ตัว เหตุผลของสิ่งนี้คือ  
ภาษา FORTRAN ไม่สนใจอักขระว่างอย่างบริบูรณ์ (ignores space completely) -  
ภาษานี้จะลบอักขระว่าง (space) ออกก่อนเริ่มการประมวลผล นอกจากนี้แล้ว ภาษา  
FORTRAN ไม่มีคำสงวนใดๆ ทั้งสิ้น ดังนั้น DO, WHILE หรือคำอื่นๆ ซึ่ง อธิบายโครง-  
สร้าง อาจนำมาเป็น ไอเดนติไฟเออร์ได้

ตัวอย่าง ข้อความสั่ง ภาษา FORTRAN

```
WHILE = 2.2
```

ถูกต้องอย่างสมบูรณ์ ภาษา FORTRAN โครงสร้างโทเค็น และวากยสัมพันธ์ เป็นสิ่งคู่กัน ตัวอย่างสุดท้าย ของ โครงสร้างศัพท์ คือ บทนิยาม ของ ข้อตกลงโทเค็น ของ ภาษา C จาก คู่มือการใช้ใน Kernighan และ Richie (1988) กล่าวดังนี้:

โทเค็น แบ่งออกเป็น หก ประเภท ได้แก่ ไอ์เดนตีไฟเออร์ คำหลัก ค่าคงที่ สายอักขระ ตัวดำเนินการ และตัวคั่นอื่นๆ (: identifiers, keywords, constants, string literals, operators, and other separators.)

ส่วน blanks, horizontal and vertical tabs, newlines, formfeeds, and comments ตามที่ได้อธิบายนั้น (รวมเรียกว่า "white space") are ignored ยกเว้น เมื่อนำมาใช้แยกโทเค็น

white space บางตัวต้องใช้แยกไอ์เดนตีไฟเออร์, คำหลัก และค่าคงที่ ซึ่ง อยู่ติดกัน ถ้า input stream ถูกแยก เป็น โทเค็นต่างๆ ตามจำนวนอักขระซึ่งกำหนดให้ โทเค็นถัดไป คือ สายของอักขระ ความยาวสูงสุด ซึ่งประกอบขึ้นเป็นหนึ่งโทเค็น ดังนั้น ภาษา C จึงติดกับ หลักของสายอักขระย่อยความยาวสูงสุด

#### 4.2 ไวยากรณ์ไม่พึ่งบริบท และรูปแบบแบกคัส-เนาร์

(Context-free grammars and BNFs)

เราจะเริ่มต้นบทนิยามของไวยากรณ์ และ BNFs ด้วยตัวอย่าง ในภาษาอังกฤษ นั้น ประโยคอย่างง่ายประกอบด้วย a noun phrase และ a verb phrase ตามด้วย a period ซึ่งสามารถแสดงออกดังนี้

1. <sentence> ::= <noun-phrase><verb-phrase>.

ต่อไปอธิบายโครงสร้างของ noun phrase และ verb phrase ดังนี้

2. <noun-phrase> ::= <article><noun>

3. <article> ::= a I the

4. <noun> ::= girl I dog

5. <verb-phrase> ::= <verb><noun-phrase>

6. <verb> ::= sees I pets

กฎไวยากรณ์เหล่านี้ แต่ละข้อประกอบด้วย string หนึ่งตัว อยู่ในเครื่องหมาย "< >" (เป็นชื่อของโครงสร้าง ซึ่งกำลังอธิบาย) ตามด้วยสัญลักษณ์ "::=" ซึ่งอ่านว่า "consists of" หรือ "is the same as" หรือ "is defined as" และจากนั้น เป็นลำดับ ของชื่อและสัญลักษณ์อื่น ๆ

เครื่องหมาย "< >" แยกความแตกต่าง ชื่อของโครงสร้าง (structure names) ออกจาก คำจริง (actual words) หรือ โทเค็น ซึ่งอาจปรากฏในภาษานั้น ตัวอย่างเช่น ในภาษา Pascal นั้น "<program>" จะหมายถึง โครงสร้างโปรแกรม ทั้งหมด (will stand for a whole program structure) ในขณะที่คำว่า "program" หมายถึง โทเค็นตัวแรก ในทุกๆ <program> เขียนดังนี้

<program> ::= program ...

สัญลักษณ์ "::=" เป็น สัญลักษณ์อภิปาษา (metasymbol) ซึ่งใช้แยก ส่วนทางซ้ายมือ ของกฎ ออกจาก ส่วนทางด้านขวามือ ของกฎ

เครื่องหมาย "< >" เป็นสัญลักษณ์อภิปาษาเช่นเดียวกับ vertical bar "|" ซึ่งหมายถึง "or" หรือ ทางเลือก (alternation) ดังนั้น กฎข้อ 6 กล่าวว่ verb อาจจะเป็นคำว่า "sees" หรือคำว่า "pets"

บางครั้ง สัญลักษณ์อภิปาษา อาจจะเป็น สัญลักษณ์จริง (actual symbol) ในภาษานั้นด้วย เช่นกัน ในกรณีเช่นนั้น เราจะใส่เครื่องหมายคำพูดให้กับ สัญลักษณ์จริง เพื่อ แยกมันออกจากสัญลักษณ์อภิปาษา สิ่งนี้เป็นความคิดที่ดี สำหรับ สัญลักษณ์พิเศษ เช่น เครื่องหมายกำกับวรรคตอน เมื่อไม่ได้ใช้เป็นสัญลักษณ์อภิปาษา

ตัวอย่างเช่น กฎข้อ 1 มี period หนึ่งตัว อยู่ท้ายสุด ขณะนั้น period ไม่ใช่ ส่วนของสัญลักษณ์อภิปาษาใดๆ ซึ่งกำลังอธิบาย ดังนั้น จึงอาจผิดพลาดได้ง่าย วิธีที่ดีกว่าคือ เขียนกฎข้อ 1 ดังนี้

<sentence> ::= <noun-phrase><verb-phrase> '.'

(เพราะฉะนั้น ขณะนี้ เครื่องหมายคำพูด เป็นสัญลักษณ์อภิปาษาด้วย)

---

หมายเหตุ ในตำราบางเล่ม ผู้แต่งอาจใช้ ตัวเอน ตัวใหญ่ เพื่อแยกความแตกต่างระหว่าง สัญลักษณ์ สัญลักษณ์อภิปาษา และชื่อ

ประโยคที่ถูกต้องใดๆ จะต้องเป็นไปตาม ไวยากรณ์ที่กล่าวข้างต้น คือสามารถสร้างชั้นได้ตามขั้นตอน ดังนี้ :

เริ่มต้นด้วย สัญลักษณ์ <sentence> และทำต่อไป เพื่อแทนที่ ส่วนทางซ้ายโดยทางเลือกต่างๆ ของส่วนทางขวามือ ในกฎที่กล่าวมาแล้ว กรรมวิธีนี้ สร้างการแปลง (derivation) ในภาษา (This process creates a derivation in the language.)

ดังนั้น เราสามารถสร้าง (construct) หรือ แปลง (derive) ประโยค "the girl sees a dog." ดังนี้

- <sentence> —> <noun-phrase><verb-phrase>. (กฎข้อ 1)
- > <article><noun><verb-phrase>. (กฎข้อ 2)
- > the <noun>< verb-phrase>. (กฎข้อ 3)
- > the girl <verb-phrase>. (กฎข้อ 4)
- > the girl <verb><noun-phrase>. (กฎข้อ 5)
- > the girl sees <noun-phrase>. (กฎข้อ 6)
- > the girl sees <article><noun>. (กฎข้อ 2)
- > the girl sees a <noun>. (กฎข้อ 3)
- > the girl sees a dog. (กฎข้อ 4)

ในทางตรงกันข้าม เราสามารถเริ่มต้น ด้วยประโยค "the girl sees a dog." และทำงานย้อนหลัง ผ่านการแปลง จนถึง <sentence> เป็นการแสดงให้เห็นว่า ประโยคที่กำหนดให้ นั้นเป็น ประโยคถูกต้องในภาษา

ไวยากรณ์อย่างง่ายนี้ แสดงให้เห็นคุณสมบัติส่วนใหญ่ ของไวยากรณ์ ของภาษาโปรแกรม แต่โปรดสังเกตว่า ไม่ใช่ ประโยคถูกต้องทั้งหมด ที่จริงๆ แล้วจะมีความหมาย (Note that not all legal sentences actually make sense.) ตัวอย่างเช่น "the dog pets the girl."

นอกจากนี้แล้ว ยังมีข้อผิดพลาดปลีกย่อยดังนี้ ในภาษาอังกฤษนั้น article ซึ่งปรากฏที่ตอนต้นของประโยค จะต้องเป็นตัวใหญ่ คุณสมบัติของ ตำแหน่ง (positional)



เช่นนี้ เป็นเรื่องยากที่จะ ทำงาน โดยใช้ ไวยากรณ์ไม่พึ่งบริบท

ต่อไปนี้ เป็นการสรุปนิยามของสิ่งต่างๆ ซึ่งเราได้กล่าวมาแล้ว ไวยากรณ์

ไม่พึ่งบริบท (A context-free grammar) ประกอบด้วย

1. ชุดของกฎไวยากรณ์ (a series of grammar rules) อธิบายดังนี้ กฎต่างๆ จะประกอบด้วย ส่วนทางซ้ายมือ ซึ่งเป็น ชื่อโครงสร้างหนึ่งชื่อ จากนั้นเป็น เครื่องหมายอภิปาษา " ::= " ตามด้วย ส่วนทางขวามือ ซึ่งประกอบด้วย ลำดับ ของ items ซึ่งอาจจะเป็น สัญลักษณ์ หรือ ชื่อโครงสร้างอื่นๆ
2. ชื่อของโครงสร้าง (the names for structures) เรียกว่า nonterminals ตัวอย่างเช่น <sentence> เพราะว่า ชื่อโครงสร้างเหล่านี้ สามารถ แบ่งย่อย ให้เป็นโครงสร้างต่อไปได้อีก
3. คำ หรือ สัญลักษณ์ โทเค้น (the words or token symbols) เรียกว่า terminals เพราะว่า คำเหล่านี้ ไม่สามารถ แบ่งย่อยได้อีก
4. กฎไวยากรณ์ (Grammar rules) เรียกอีกอย่างหนึ่งว่า productions เพราะว่า มัน "produce" strings ของภาษา โดยใช้การแปลง

productions อยู่ในรูปของ Backus - Naur form ถ้ามันถูกกำหนดให้ โดยใช้เฉพาะสัญลักษณ์อภิปาษา " ::= ", " | ", "<" และ ">" เท่านั้น (บางครั้ง เครื่องหมายวงเล็บเล็ก อนุญาตให้ใช้ได้เพื่อจัดกลุ่ม ของ items เข้าด้วยกันได้)

ในตัวอย่างข้างต้น มี terminal 7 ตัว ได้แก่ "girl", "dog", "sees", "pets", "the", "a", และ "." มี nonterminal 6 ตัว และมี 6 productions โดยทั่วไปแล้ว จำนวน production ใน ไวยากรณ์ไม่พึ่งบริบท จะมีเท่ากับ จำนวน nonterminal

ถึงแม้ว่า เราอาจจะจัด สัญลักษณ์อภิปาษา " | " ออกไป โดยเขียนทางเลือก แยกต่างหากจากกัน เช่น

<noun> ::= girl

<noun> ::= dog

ซึ่งจะทำให้ nonterminal แต่ละตัว สมัยกับ จำนวน production มากเท่ากับจำนวน ทางเลือก

ทำไมจึงเรียกว่า ไวยากรณ์ ไม่พึ่งบริบท เหตุผลง่าย ๆ คือ nonterminals ซึ่งปรากฏ ทางซ้ายมือ ของ productions มีหนึ่งตัวเท่านั้น สิ่งนี้หมายความว่า non-terminal แต่ละตัว สามารถถูกแทนที่ด้วย ทางเลือกของด้านขวามือ ตัวใดก็ได้ ไม่ว่า non-terminal นั้นจะปรากฏที่ใด พุทธิกอย่างหนึ่งคือ ไม่มี บริบท หรือ อรรถาธิบาย ภายใต้สิ่งซึ่งเกิดขึ้นเฉพาะการแทนที่เท่านั้น (In other words, there is no context under which only certain replacements can occur.)

ตัวอย่างเช่น ในไวยากรณ์ที่เราเพิ่งอภิปราย มันจะรู้เรื่องเมื่อใช้ verb คำว่า "pets" เฉพาะเมื่อ ประธานเป็น "girl" เท่านั้น สิ่งนี้เรียกว่า การพึ่งบริบท\* (a context-sensitivity)

เราสามารถเขียน ไวยากรณ์พึ่งบริบท โดยยอมให้ "context string" ปรากฏทางซ้ายมือ ของ กฎไวยากรณ์ แต่บางคน อาจถือว่า การพึ่งบริบท เป็นหัวข้อ วากยสัมพันธ์ (syntactic issue) เราจะยอมรับสิ่งนี้ อย่างไรก็ตาม สิ่งใดก็ตาม ซึ่งไม่สามารถแสดงออก ด้วย ไวยากรณ์ไม่พึ่งบริบท สิ่งนั้นคือ อรรถศาสตร์ ไม่ใช่ หัวข้อ วากยสัมพันธ์ (We shall adopt the view, however, that anything not expressible using context-free grammars is a semantic, not a syntactic issue.) มีบางสิ่ง ซึ่งสามารถแสดงออกด้วยไวยากรณ์ไม่พึ่งบริบทได้ แต่บ่อยครั้ง เป็นเรื่องดีกว่าที่จะทิ้งให้เป็น การอธิบายความหมาย เพราะว่า มันเกี่ยวข้องกับ การเขียน productions พิเศษ จำนวนมาก

ตัวอย่างของ การพึ่งบริบท (context-sensitivity) ตามข้อสังเกต ที่ว่าคือ article ซึ่งปรากฏ ตอนต้นของประโยค ใน ไวยากรณ์ข้างต้น ควรจะเป็นตัว ใหญ่ วิธีหนึ่งในการทำสิ่งนี้ คือ เขียนกฎข้อ 1 ใหม่ ดังนี้

<sentence> ::= <beginning><noun-phrase><verb-phrase> '.'

\* หรือ ตัวอย่างเช่น ภาษา Pascal กล่าวว่า ตัวแปรทุกตัวในโปรแกรม ต้องมีการ ประกาศ (declared) หรือ ภาษา FORTRAN กล่าวว่า ชื่อทุกชื่อ ซึ่งขึ้นต้นด้วยตัว อักษรตัวใดตัวหนึ่งใน I ถึง N จะถือว่าเป็น ตัวแปรชนิด integer ส่วนกรณีอื่นๆ ถือ ว่าเป็น ตัวแปรชนิด real เป็นต้น

จากนั้น ใส่กฎ ฟังก์ชันบริบท ดังนี้ :

$\langle \text{beginning} \rangle \langle \text{article} \rangle ::= \text{The I A}$

ขณะนี้ การแปลง จะ เป็นดังนี้

$\langle \text{sentence} \rangle \longrightarrow \langle \text{beginning} \rangle \langle \text{noun-phrase} \rangle \langle \text{verb-phrase} \rangle.$

(กฎข้อ 1 ใหม่)

$\longrightarrow \langle \text{beginning} \rangle \langle \text{article} \rangle \langle \text{noun} \rangle \langle \text{verb-phrase} \rangle.$

(กฎข้อ 2)

$\longrightarrow \text{The } \langle \text{noun} \rangle \langle \text{verb-phrase} \rangle.$  (กฎฟังก์ชันบริบท ข้อใหม่)

$\longrightarrow . . . .$

ไวยากรณ์ไม่ฟังก์ชันบริบท ได้มีการศึกษาอย่างกว้างขวาง โดยนักทฤษฎีภาษาแบบทางการ (formal language theorists) และปัจจุบัน เป็นที่เข้าใจอย่างดีแล้วว่า มันเป็นธรรมชาติที่จะแสดงออก วากยสัมพันธ์ ของ ภาษาโปรแกรมใดๆ ในรูปแบบของ BNF ได้ จริงๆ แล้ว การทำเช่นนั้น ทำให้การเขียนตัวแปลภาษา สำหรับภาษาโปรแกรมง่ายขึ้น เพราะว่าชั้นวิเคราะห์กระจายสามารถเป็นไปอัตโนมัติ (ดูหัวข้อ 4.6)

ตัวอย่างเบื้องต้น ของการใช้ ไวยากรณ์ไม่ฟังก์ชันบริบท ในภาษาโปรแกรม ได้แก่ การอธิบาย นิพจน์คำนวณ จำนวนเต็ม อย่างง่าย ด้วย การบวกและการคูณ ซึ่งกำหนดให้ ในรูป 4.1

(A typical simple example of the use of a context-free grammar in programming languages is the description of simple integer arithmetic expressions with addition and multiplication given in Figure 4.1.)

$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \mid$

$(\langle \text{exp} \rangle) \mid \langle \text{number} \rangle$

$\langle \text{number} \rangle ::= \langle \text{number} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$

$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

รูป 4.1 A simple arithmetic expression grammar

โปรดสังเกต ธรรมชาติการเรียกซ้ำ ของกฎ : นิพจน์ หนึ่งชุดอาจจะเป็น ผลบวกของนิพจน์ สองชุด หรือ ผลคูณของนิพจน์ สองชุด ซึ่งแต่ละชุดอาจแบ่งต่อไปอีก เป็นผลบวก หรือ ผลคูณ สุดท้าย กระบวนการนี้ ต้องจบโดยการเลือก ทางเลือก <number> หรือเราอาจจะไม่ได้มาถึงสาย (string) ของ terminals

โปรดสังเกตว่า การเรียกซ้ำ ในกฎ <number > ใช้เพื่อก่อกำเนิด (generate) ลำดับของเลขโดดซ้ำกัน ตัวอย่างเช่น เลข 2 3 4 ถูกสร้าง ดังนี้

```

<number>  —> <number><digit>
           —> <number><digit><digit>
           —> <digit><digit><digit>
           —> 2<digit><digit>
           —> 23<digit>
           —> 234

```

ต่อไปคือ รูป 4.2 แสดงการเริ่มต้น ของการอธิบาย BNF สำหรับภาษา Pascal

```

<program>          ::=      <program-heading>' '<program-block>'.'
<program-heading> ::=      program <identifier> I program
                        <identifier>  '('<program-parameter>')'
<program-block>    ::=      <block>
<program-parameters> ::=  <identifier-list>
<block>            ::=      <label-declaration-part>
                        <constant-definition-part>
                        <type-definition-part>
                        <variable-declaration-part>
                        <procedure-and-function-declaration-part>
                        <statement-part>

```

```

<label-declaration-part> ::= label <label-list> ';' I <empty>
<label-list>                ::= <label-list> ',' <label> I <label>
                               . . .
                               . . .

```

## รูป 4.2 BNFs บางส่วน สำหรับภาษา Pascal

### 4.2.1 กฎ BNFs คือสมการ

#### (BNFs rules as Equations)

ทางเลือกวิธีหนึ่งของการอธิบายว่า กฎ BNF สร้างสาย (strings) ของภาษา ทำได้ไฉนนั้น เป็นดังนี้

กำหนด กฎไวยากรณ์ให้หนึ่งข้อ เช่น

$$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{number} \rangle$$

(นามธรรมจากรูป 4.1) ให้ E เป็นเซตของสาย ซึ่งแทนนิพจน์ ส่วน N หมายถึงเซตของสายซึ่งแทนเลข กฎไวยากรณ์ข้างต้น อาจมองในรูป สมการเซต (set equation) ดังนี้

$$E = E + E \cup N$$

เมื่อ  $E + E$  หมายถึง เซตสร้างโดยการต่อกัน ของสายทั้งหมดจาก E ด้วยสัญลักษณ์ "+" และจากนั้น สายทั้งหมด จาก E และ "U" หมายถึง ผลผนวกของเซต สมมติว่า เซต N ถูกสร้างเรียบร้อยแล้ว สิ่งนี้ แทน สมการการเรียกซ้ำ สำหรับเซต E เซต E เล็กที่สุด จะเป็นไปตามสมการนี้ สามารถเอามาเป็น เซตนิยามโดยกฎไวยากรณ์ คาดได้ว่า สิ่งนี้ ประกอบด้วยเซต

$$N \cup N + N \cup N + N + N \cup N + N + N \cup \dots$$

และเราสามารถพิสูจน์ได้ว่า สิ่งนี้เป็นเซตเล็กที่สุด ซึ่งเป็นไปตามสมการที่กำหนดให้ (ดูแบบฝึกหัดข้อ 40) ผลเฉลยของสมการเรียกซ้ำ เกิดขึ้นบ่อยมาก ใน การอธิบายแบบทางการ ของภาษาโปรแกรม และเป็นหัวข้อหลักของการศึกษา ในวิชาทฤษฎีของภาษาโปรแกรม ซึ่งเรียกว่า จุดคงที่น้อยที่สุด (least fixed points)

#### 4.3 ต้นไม้วิภังค์ และต้นไม้วากยสัมพันธ์แบบนามธรรม

##### (Parse trees and Abstract syntax trees)

วากยสัมพันธ์ สร้าง โครงสร้าง ไม่ใช่ความหมาย แต่ความหมายของประโยค (หรือ โปรแกรม) ต้องเกี่ยวข้องกับวากยสัมพันธ์

(Syntax establishes structure, not meaning. But the meaning of a sentence (or program) must be related to its syntax. )

ในภาษาอังกฤษ หนึ่งประโยคมี a subject และ a predicate ซึ่งเป็นความคิดเชิงความหมาย เพราะว่า subject (the "actor") และ predicate (the "action") บอกความหมายของประโยค

subject ปกติ จะอยู่ตอนต้นของประโยค และกำหนดให้ด้วย a noun phrase ดังนั้น ในวากยสัมพันธ์ ของประโยคภาษาอังกฤษ noun phrase ต้องวางที่ตำแหน่งแรก และสิ่งนั้นคือ subject

ในทำนองเดียวกัน ในไวยากรณ์ของนิพจน์ เมื่อเราเขียน

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

เราหวังว่า ให้บวกค่าของนิพจน์ ทางขวามือ สองชุด เข้าด้วยกัน เพื่อให้ได้ค่าของ นิพจน์ทางซ้ายมือ กรรมวิธีของการผูกติดความหมาย ของ ตัวสร้างหนึ่ง กับ โครงสร้างวากยสัมพันธ์ของมัน สิ่งนี้เรียกว่า วากยสัมพันธ์-ต่อตรงกับ อรรถศาสตร์

(This process of attaching the semantics of a construct to its syntactic structure is called syntax-directed semantics.)

เราต้องสร้างวากยสัมพันธ์ เพื่อให้มันสะท้อน ความหมาย สุดท้ายเราจะผูกติดกับมัน มากที่สุดเท่าที่เป็นไปได้ (วากยสัมพันธ์-ต่อตรงกับ ความหมาย ซึ่งอาจเรียกว่า ความหมาย-ต่อตรงกับ วากยสัมพันธ์)

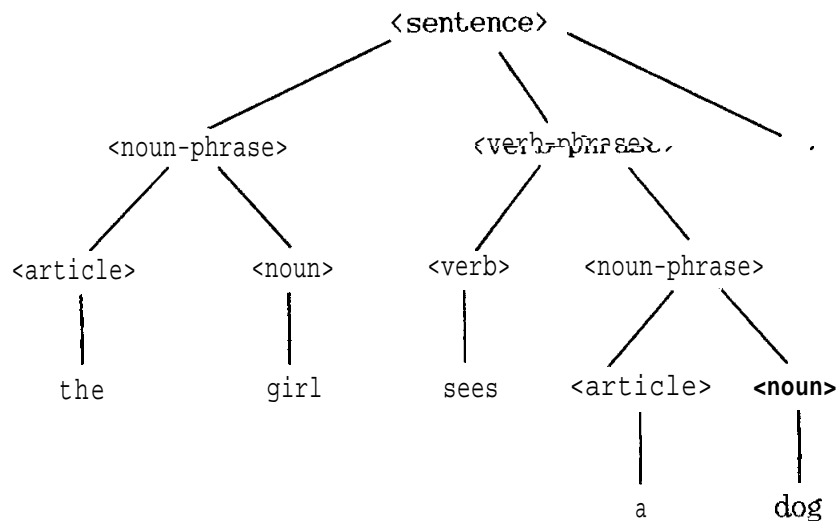
เพื่อให้ โครงสร้างวากยสัมพันธ์ ของ โปรแกรม บอกความหมายของมัน เราต้องมีวิธีแสดงออกของโครงสร้างนี้ เช่น การหาโดยการแปลง

วิธีมาตรฐาน สำหรับทำสิ่งนี้ คือ ต้นไม้วิภังค์ หรือต้นไม้วิเคราะห์กระจาย

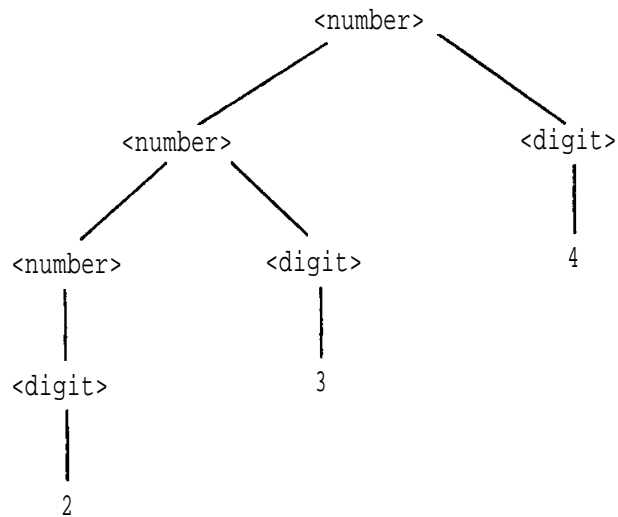
(parse tree)

ต้นไม้วิภังค์ อธิบายกรรมวิธี การแทนที่ ของ การแปลงด้วยภาพ (The parse tree describes graphically the replacement process in a derivation.)

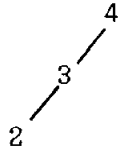
ตัวอย่าง ต้นไม้วิภังค์ สำหรับ ประโยค "the girl sees a dog."



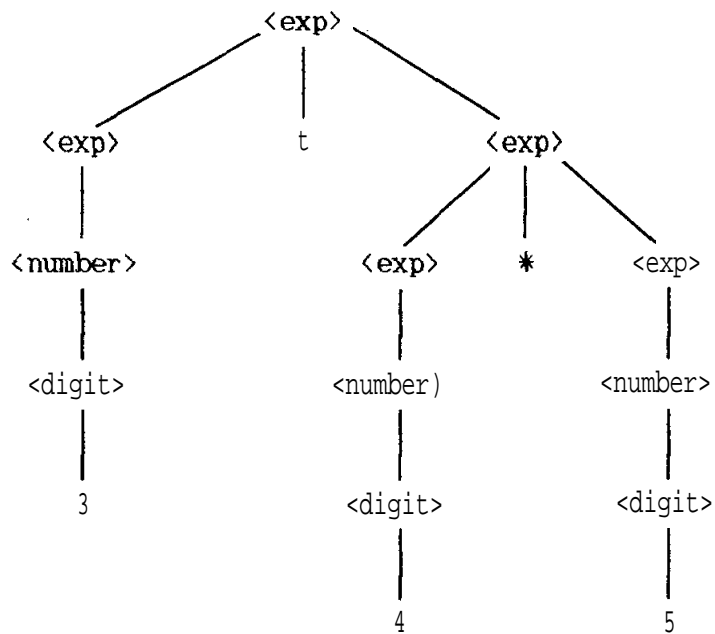
ในทำนองเดียวกัน ต้นไม้วิภังค์ สำหรับเลข 234 ในไวยากรณ์ของนิพจน์ คือ



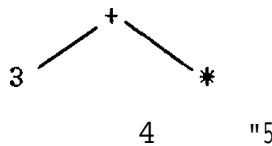
ต้นไม้วิภังค์ ถูกระบุ โดย nonterminals ที่ โหนดภายใน (interior nodes) และ terminals ที่ จุดแตกใบ (leaf) ทั้งนี้ terminals และ nonterminals ทั้งหมดในการแปลง รวมอยู่ในต้นไม้วิภังค์ แต่ terminals และ nonterminals อาจไม่จำเป็นต้องเป็นทั้งหมด เพื่อบอกความบริบูรณ์ของ โครงสร้างวากยสัมพันธ์ ของ นิพจน์ หรือ ประโยค ตัวอย่างเช่น โครงสร้างของเลข 234 บอกได้อย่างบริบูรณ์ จากต้นไม้



และต้นไม้วิภังค์ สำหรับ  $3 + 4 * 5$  เช่น



อาจทำให้เล็กลงเป็นต้นไม้





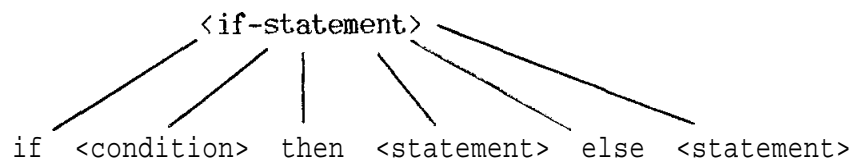
ต้นไม้เช่นนั้นเรียกว่า ต้นไม้วากยสัมพันธ์แบบนามธรรม (abstract syntax trees) หรือ ต้นไม้วากยสัมพันธ์ (syntax trees) เพราะว่ามันบอก (abstract) เฉพาะโครงสร้างสำคัญของต้นไม้วิภังค์

ต้นไม้วากยสัมพันธ์แบบนามธรรม อาจเอา terminals ซึ่งครึ่งหนึ่ง มีซ้ำกันในโครงสร้างของต้นไม้ ออกไป

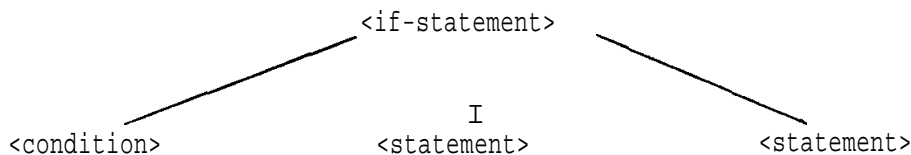
ตัวอย่าง กฎไวยากรณ์

$\langle \text{if-statement} \rangle ::= \text{if } \langle \text{condition} \rangle \text{ then } \langle \text{statement} \rangle$   
 $\qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \text{else } \langle \text{statement} \rangle$

จะให้ต้นไม้วิภังค์ ดังนี้



และต้นไม้วากยสัมพันธ์แบบนามธรรม ดังนี้



มันเป็นไปได้ ที่จะเขียนกฎ สำหรับวากยสัมพันธ์แบบนามธรรม ในวิธีที่คล้ายกับกฎ BNF สำหรับ วากยสัมพันธ์ปกติ แต่ในที่นี้ไม่ได้ทำเช่นนั้น บางครั้ง วากยสัมพันธ์ปกติ ถูกแยกต่างหาก จาก วากยสัมพันธ์นามธรรม ซึ่งเรียกว่า วากยสัมพันธ์แบบรูปธรรม (concrete syntax)

ตัวแปลภาษา บ่อยครั้งจะสร้างต้นไม้วากยสัมพันธ์ มากกว่าเป็นต้นไม้วิภังค์ เต็มรูป เพราะว่า มันกระชับ และแสดงความสำคัญ ของ โครงสร้าง ได้มากกว่า

#### 4.4 ความกำกวม การสลับที่ และ การทำก่อน

(Ambiguity, Associativity and Precedence)

การแปลงที่แตกต่างกัน สองชุด สามารถนำไปสู่ ต้นไม้วิภังค์ ต้นเดียวกันได้ (Two different derivations can lead to the same parse tree.) ตัวอย่างเช่น ในการแปลง เลข 234 ในหัวข้อที่แล้ว เราอาจเลือก เพื่อแทนที่ <digit> เป็นอันดับแรก ดังนี้

<number>  $\longrightarrow$  <number><digit>  
 $\longrightarrow$  <number> 4  
 $\longrightarrow$  <number><digit> 4  
 $\longrightarrow$  <number> 34  
...

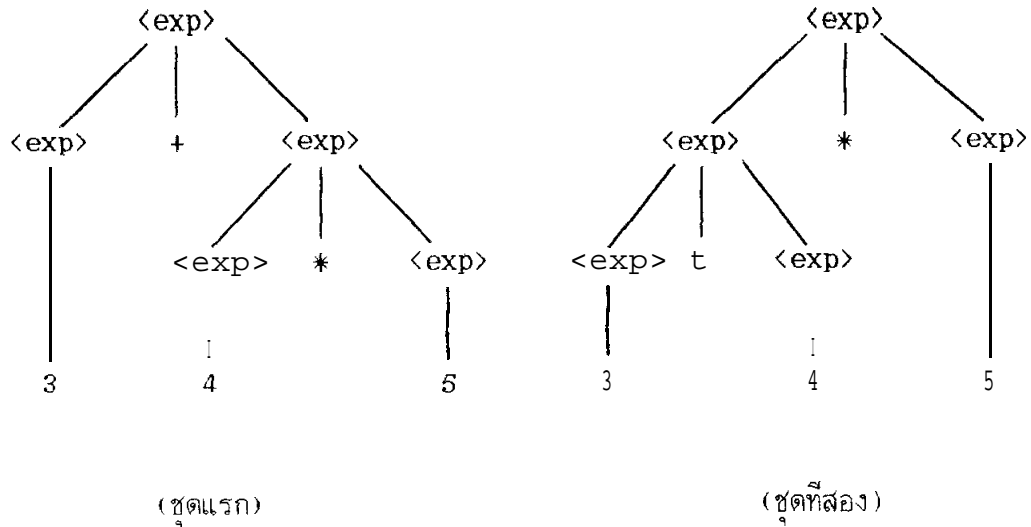
ซึ่งยังคงได้ ต้นไม้วิภังค์เหมือนเดิม อย่างไรก็ตาม การแปลงที่แตกต่างกัน สามารถนำไปสู่ ต้นไม้วิภังค์ที่แตกต่างกันได้ด้วย ตัวอย่างเช่น ถ้าเราสร้าง  $3 + 4 * 5$  จากไวยากรณ์นิพจน์ ของรูป 4.1 เราสามารถใช้การแปลง

<exp>  $\longrightarrow$  <exp> + <exp>  
 $\longrightarrow$  <exp> + <exp> \* <exp>  
(แทน <exp> ชุดที่สองด้วย <exp> \* <exp>)  
 $\longrightarrow$  <number> + <exp> \* <exp>  
 $\longrightarrow$  . . .

หรือ การแปลง

<exp>  $\longrightarrow$  <exp> \* <exp>  
 $\longrightarrow$  <exp> t <exp> \* <exp>  
(แทน <exp> ชุดแรกด้วย <exp> t <exp> )  
 $\longrightarrow$  <number> t <exp> \* <exp>  
 $\longrightarrow$  . . .

ซึ่งการแปลง สองชุดข้างต้นนี้ นำไปสู่ต้นไม้วิภังค์ สองชุด



และ ต้นไม้วากยสัมพันธ์แบบนามธรรม สองชุด



ไวยากรณ์ เช่นที่กำหนดในรูป 4.1 สำหรับ ต้นไม้วิภังค์ หรือ ต้นไม้วากยสัมพันธ์ที่แตกต่างกันสองชุด อาจเป็นไปได้สำหรับ string ชุดเดียวกัน หมายถึง **ความกำกวม** (ambiguous)

ไวยากรณ์ที่กำกวม นำเสนอสิ่งยาก เพราะว่า โครงสร้าง แสดงออกไม่ชัดเจน เพื่อให้เป็นประโยชน์ ไวยากรณ์นั้นต้องปรับปรุงแก้ไขใหม่ (revised) เพื่อความกำกวม หรือ **กฎความไม่กำกวม** ต้องถูกกล่าวขึ้นเพื่อสร้างว่า โครงสร้างชุดใดที่ต้องการ หมายถึง (a **disambiguating rule** must be stated to establish which structure is meant.)

สำหรับนิพจน์  $3 + 4 * 5$  นั้น ต้นไม้วิภังค์ชุดใดถูกต้อง ถ้าเราคิดถึงอรรถศาสตร์ซึ่งผูกติดกับนิพจน์ เราสามารถเข้าใจความหมายของการตัดสินใจนี้ว่าคืออะไร ต้นไม้วากยสัมพันธ์ต้นแรก แสดงว่า ตัวดำเนินการคูณ "\*" ใช้กับ 4 และ 5 (ผลลัพธ์คือ 20)

จากนั้นจึงนำผลลัพธ์นี้ บวกกับ 3 ได้คำตอบ 23 ในทางตรงกันข้าม ต้นไม้วากยสัมพันธ์ ชุดที่สอง สิ่งแรก บวก 3 กับ 4 (ผลลัพธ์คือ 7) จากนั้นจึงคูณด้วย 5 คำตอบคือ 35 ดังนั้น การดำเนินการ ซึ่งถูกใช้ในอันดับแตกต่างกัน และความหมายของผลลัพธ์ จะแตกต่างกันด้วย

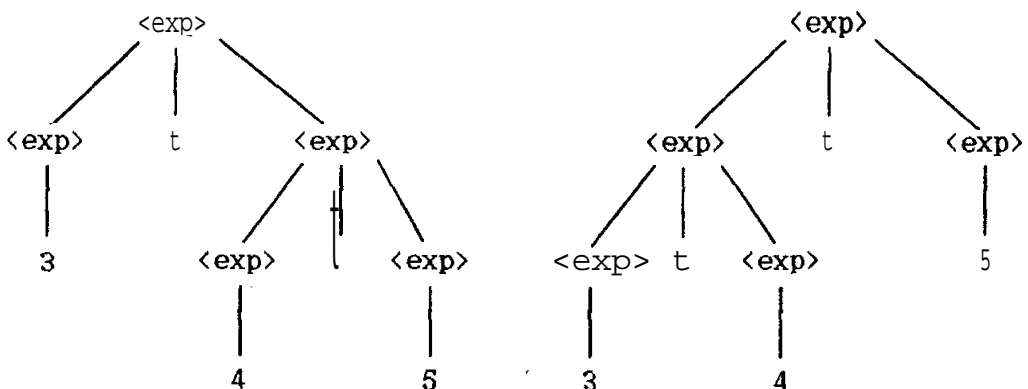
(Thus the operations are applied in a different order, and the resulting semantics are quite different.)

ถ้าเราเอาความหมายปกติ ของนิพจน์  $3 + 4 * 5$  จากวิชาคณิตศาสตร์ เราควรเลือก ต้นไม้ต้นแรก มากกว่าที่จะเป็นต้นที่สอง เพราะว่าการคูณ มี การทำก่อน สูงกว่าการบวก (since multiplication has precedence over addition) สิ่งนี้ เป็นการเลือกปกติ ในภาษาโปรแกรม ถึงแม้ว่าบางภาษา (เช่น APL) มีการเลือกที่แตกต่างออกไป เราสามารถแสดงออกความจริงที่ว่า การคูณ มี การทำก่อน สูงกว่าการบวกได้อย่างไร? เราควร กล่าวกฎไม่กำกวม แยกต่างหากจากไวยากรณ์ หรือเราควรปรับปรุงแก้ไข (revise) ไวยากรณ์ วิธีปกติ ของการปรับปรุงแก้ไขไวยากรณ์ คือ เขียนกฎไวยากรณ์ใหม่ (เรียกว่า "<term>") ซึ่งสร้าง การต่อเรียงการทําก่อน (precedence cascade) บังคับให้ การจับคู่ ของ "\*" ที่จุดต่ำกว่า ใน ต้นไม้วักย

$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{term} \rangle$

$\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{term} \rangle \mid (\langle \text{exp} \rangle) \mid \langle \text{number} \rangle$

แต่การแก้ปัญหาความกำกวม ยังไม่สมบูรณ์ เพราะว่า กฎ สำหรับ <exp> ยังคงวิเคราะห์กระจาย นิพจน์  $3 + 4 + 5$  เป็น  $(3 + 4) + 5$  หรืออาจเป็น  $3 + (4 + 5)$ . พุดอีกอย่างหนึ่งคือ เราสามารถทำการบวก ให้เป็น การสลับที่แบบขวา (right-associative) หรือ การสลับที่แบบซ้าย (left-associative) ได้ :



กรณีของการบวก สิ่งนี้ไม่มีผลกระทบทับผลลัพธ์ แต่ถ้าเรารวมการลบ (subtraction) ไว้ด้วย มีการกระทบแน่นอน ตัวอย่างเช่น นิพจน์  $8 - 4 - 2$

กรณี ถ้า "-" เป็น การสลับที่แบบซ้าย ผลลัพธ์คือ  $8 - 4 - 2 = 2$

แต่ถ้า "-" เป็น การสลับที่แบบขวา ผลลัพธ์คือ  $8 - 4 + 2 = 6$

ดังนั้นจึงจำเป็นต้องแทนกฎ

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

ด้วยกฎ

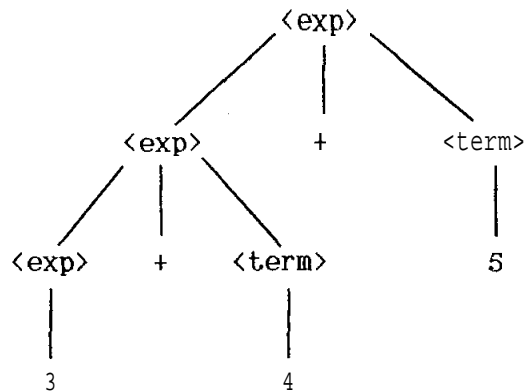
$$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{term} \rangle \quad \text{===== (1)}$$

หรือกฎ

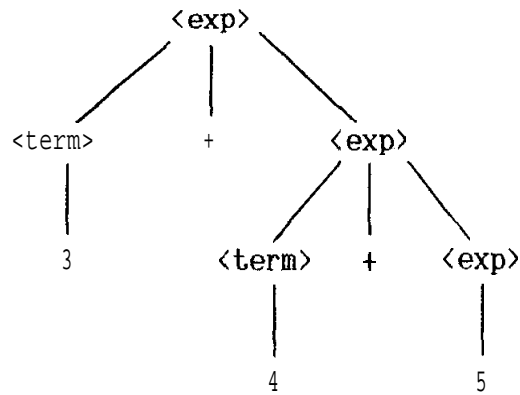
$$\langle \text{exp} \rangle ::= \langle \text{term} \rangle + \langle \text{exp} \rangle \quad \text{===== (2)}$$

จุดไหน ใน กฎสองข้อนี้ เป็นจุดที่ถูกต้อง กฎข้อแรกนั้น เป็น การเรียกซ้ำแบบซ้าย (left-recursive) ในขณะที่กฎข้อสอง เป็น การเรียกซ้ำแบบขวา (right-recursive)

กฎการเรียกซ้ำแบบซ้าย สำหรับการดำเนินการ ทำให้มันเป็นการสลับที่แบบซ้าย เช่น ในต้นไม้วิภังค์ ข้างล่างนี้



ในขณะที่ กฎการเรียกซ้ำแบบขวา ทำให้ การดำเนินการ เป็นการสลับที่แบบขวา



ไวยากรณ์ ชุดแก้ไขปรับปรุงใหม่ สำหรับ นิพจน์คำนวณอย่างง่าย ซึ่งแสดงออกทั้ง การทำก่อน และการสลับที่ กำหนดให้แล้วในรูป 4.3

ขณะนี้ BNF สำหรับ นิพจน์คำนวณอย่างง่าย ไม่กำกวม (การนิสัจของกฎเหล่านี้ ต้องใช้ เทคนิคขั้นสูงมากขึ้น จากทฤษฎีของการวิเคราะห์กระจาย)

นอกจากนี้แล้ว ต้นไม้วิภังค์ จะสมนัยกับความหมายของการดำเนินการคำนวณ เช่น ที่นิยามใช้ปกติ บางครั้ง กรรมวิธีของ การเขียน ไวยากรณ์ชั้นใหม่ เพื่อขจัดความกำกวม เป็นสาเหตุที่ทำให้ ไวยากรณ์ มีความซับซ้อนอย่างมาก แต่กรณีต่างๆ เช่นนั้น เราชอบที่จะกล่าวไวยากรณ์ที่ไม่กำกวม มากกว่า (ดู การอภิปราย if-statement ในบทที่ 7)

```

<exp> ::= <exp> + <term> | <term>
<term> ::= <term> * <factor> | <factor>
<factor> ::= (<exp>) | <number>
<number> ::= <number><digit> | <digit>
<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  
```

รูป 4.3 ไวยากรณ์ชุดแก้ไขปรับปรุง สำหรับนิพจน์คำนวณอย่างง่าย

#### 4.5 อีบีเอ็นเอฟ และแผนภาพวากยสัมพันธ์

(EBNFs and syntax diagrams)

จากกฎไวยากรณ์ ในหัวข้อที่แล้ว

$\langle \text{number} \rangle ::= \langle \text{number} \rangle \langle \text{digit} \rangle \mid \langle \text{digit} \rangle$

ก่อนำเนิต (generate) เลขหนึ่งจำนวน เป็นลำดับของเลขโดด ดังนี้ :

$\langle \text{number} \rangle \longrightarrow \langle \text{number} \rangle \langle \text{digit} \rangle$

$\longrightarrow \langle \text{number} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle$

$\longrightarrow \langle \text{number} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle \langle \text{digit} \rangle$

. . .

$\longrightarrow \langle \text{digit} \rangle . . . \langle \text{digit} \rangle$

(arbitrary repetitions of  
digits)

ในทำนองเดียวกัน กฎ

$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{term} \rangle \mid \langle \text{term} \rangle$

ก่อนำเนิต นิพจน์ หนึ่งชุด เป็น ลำดับ ของ term แยกกันด้วยเครื่องหมาย "+" ดังนี้

$\langle \text{exp} \rangle \longrightarrow \langle \text{exp} \rangle + \langle \text{term} \rangle$

$\longrightarrow \langle \text{exp} \rangle + \langle \text{term} \rangle t \langle \text{term} \rangle$

$\longrightarrow \langle \text{exp} \rangle + \langle \text{term} \rangle t \langle \text{term} \rangle t \langle \text{term} \rangle$

. . .

$\longrightarrow \langle \text{term} \rangle + . . . t \langle \text{term} \rangle$

สถานะการณืเช่นนี้ เกิดขึ้นบ่อยมาก ดังนั้น สัญกรณ์พิเศษหนึ่งตัว สำหรับกฎไวยากรณ์เช่นนี้ ถูกยอมรับว่า แสดงธรรมชาติการทำซ้ำของโครงสร้างของมันได้ชัดเจนมากกว่า

$\langle \text{number} \rangle ::= \langle \text{digit} \rangle \{ \langle \text{digit} \rangle \}$

และ

$\langle \text{exp} \rangle ::= \langle \text{term} \rangle \{ + \langle \text{term} \rangle \}$

สัญกรณ์นี้ (this notation) วงเล็บปีกกา ( $\{ \}$ ) หมายถึง "zero or more repetitions of" ดังนั้น กฎนี้ แสดงว่า เลข หมายถึง ลำดับของเลขโดดหนึ่งตัว หรือ เลขโดด มากกว่าหนึ่งตัว และกฎที่แสดงว่า นิพจน์ หมายถึง term ตามด้วยการกระทำซ้ำของ "+" และอีกหนึ่งเทอม อาจจะไม่มี หรือมีหลายชุดได้ (zero or more repetition of a "+" and another term.) สัญกรณ์นี้ วงเล็บปีกกา จึงเป็น สัญลักษณ์อภิภาษา ตัวใหม่ (new metasymbols) และเรียกลักษณะนี้ว่า **extended Backus-Naur form** หรือ ตัวย่อว่า EBNF

สัญกรณ์ใหม่นี้ การสลับที่แบบซ้ายของตัวดำเนินการ "+" ซึ่งแสดงออกโดยการเรียกซ้ำแบบซ้าย ของกฎเดิม ใน BNF จะไม่เห็นเลย (obscures) สิ่งนี้เราต้องสมมติว่า ตัวดำเนินการใดๆ ซึ่งเกี่ยวกับการทำซ้ำ ของวงเล็บปีกกา เป็นการสลับที่แบบซ้าย จริงๆ แล้ว ถ้าตัวดำเนินการ เป็นการสลับที่แบบขวา กฎไวยากรณ์ซึ่งตรงกันควรจะเป็น การเรียกซ้ำแบบขวา และกฎการเรียกซ้ำแบบขวา ปกติ จะไม่เขียน โดยใช้วงเล็บปีกกา ปัญหานี้เป็นช่องโหว่ของ ไวยากรณ์ EBNF ต้นไม้วิเศษ และต้นไม้วากยสัมพันธ์ ไม่สามารถเขียนได้โดยตรงจากไวยากรณ์ แต่ ข้อสมมติ (assumptions) ต้องกระทำเกี่ยวกับโครงสร้างของมัน เพราะฉะนั้น ปกติเราจะใช้ BNF เพื่อเขียนต้นไม้วิเศษ

สถานะการที่สอง คือ โครงสร้าง ซึ่งมีส่วนละเว้นได้ (optional part) เช่น ส่วนที่เป็น else ของ ข้อความสั่ง if ในภาษา Pascal แสดงออกด้วย BNF ดังนี้

```
<if-statement> ::= if <condition> then <statement> |
                 if <condition> then <statement> else <statement>
```

สิ่งนี้เขียนง่ายกว่า เมื่อ แสดงออกด้วย EBNFs ดังนี้

```
<if-statement> ::= if <condition> then <statement>
                  Celse <statement>]
```



เมื่อ วงเล็บใหญ่ ( E J ) หมายถึง สัญลักษณ์อภิปาษาตัวใหม่ แสดงถึง ส่วนละ  
เว้นได้ ของ โครงสร้าง

สถานะการณือีกอย่างหนึ่งซึ่ง โครงสร้าง อาจละเว้นได้ เกิดขึ้นเมื่อ กฎ  
ไวยากรณ์ แสดงถึง สายว่าง (empty string) ซึ่งเขียนด้วย สัญลักษณ์อภิปาษาตัวใหม่  
ε (ดูแบบฝึกหัด ข้อ 39) หรือ เป็น nonterminal <empty>

ตัวอย่าง กฎไวยากรณ์ในรูป 4.2

```
<label-declaration-part> ::= label <label-list> ';' |  
                               <empty>
```

สร้าง label-declaration-part ให้เป็นส่วนละเว้นได้ สิ่งนี้แสดงด้วย  
EBNF ดังนี้

```
<label-declaration-part> ::= [label <label-list> ';' ]
```

ตัวดำเนินการ (ทวิภาค) สลับที่แบบขวา (Right-associative (binary)  
operators) สามารถเขียนโดยใช้ สัญลักษณ์อภิปาษาตัวใหม่ เหล่านี้ ตัวอย่างเช่น ถ้า  
"@" เป็น ตัวดำเนินการสลับที่แบบขวา ดังนี้ เขียนด้วย BNF ดังนี้

```
<exp> ::= <term> @ <exp> | <term>
```

ดังนั้น กฎข้อนี้ เขียนใหม่ด้วย EBNF เป็นดังนี้

```
<exp> ::= <term> [ @ <exp> ]
```

สำหรับความบริบูรณ์ ของ ไวยากรณ์ในรูป 4.3 สำหรับนิพจน์คำนวณอย่างง่าย เขียนด้วย  
EBNF ได้ดังนี้

```
<exp> ::= <term> {+ <term> %}
```

```
<term> ::= <factor> { * <factor> }
```

```
<factor> ::= ( <exp> ) | <number>
```

```
<number> ::= <digit> { <digit> }
```

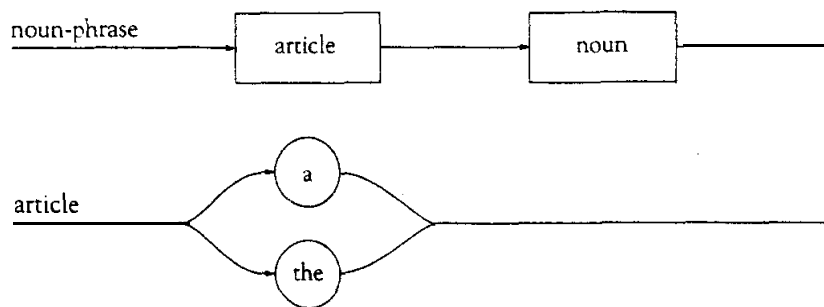
```
<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

#### รูป 4.4 กฎ EBNFs สำหรับนิพจน์คำนวณอย่างง่าย

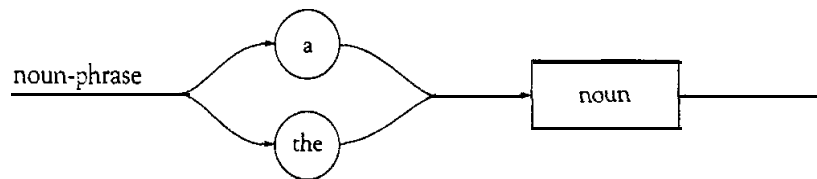
แผนภาพวากยสัมพันธ์ เป็นการแทนที่เชิงภาพที่เป็นประโยชน์สำหรับกฎไวยากรณ์ ซึ่งแสดงลำดับของ terminals และ nonterminals ที่พบทางด้านขวามือ ของกฎ

(Syntax diagram is a useful **graphical** representation for a grammar rule, which indicates the sequence of terminals and nonterminals encountered in the right-hand side of the rule.)

ตัวอย่างเช่น แผนภาพวากยสัมพันธ์ของ <noun-phrase> และ (article) ของไวยากรณ์ภาษาอังกฤษ ในหัวข้อ 4.2 วาดภาพดังนี้



หรือรวม ทั้งสองรูป ให้เป็น แผนภาพเดียว ดังนี้



ในแผนภาพวากยสัมพันธ์ เราใช้ วงกลม หรือ วงรี สำหรับ terminals และ ใช้ สี่เหลี่ยมจัตุรัส หรือ สี่เหลี่ยมผืนผ้า สำหรับ nonterminals การต่อสิ่งเหล่านี้เข้าด้วยกันใช้ เส้น (line) และลูกศร (arrows) เพื่อแสดงลำดับที่เหมาะสม

แผนภาพวากยสัมพันธ์ อาจรวมหลาย productions เข้าด้วยกัน ให้เป็นหนึ่ง

แผนภาพได้

แผนภาพวากยสัมพันธ์ สำหรับไวยากรณ์ ในรูป 4.4 กำหนดให้แล้วในรูป 4.5  
โปรดสังเกตว่า การใช้ลูป (loops) ในแผนภาพ เพื่อแสดงการทำซ้ำ ซึ่งกำหนดโดย  
วงเล็บปีกกาใน EBNFs

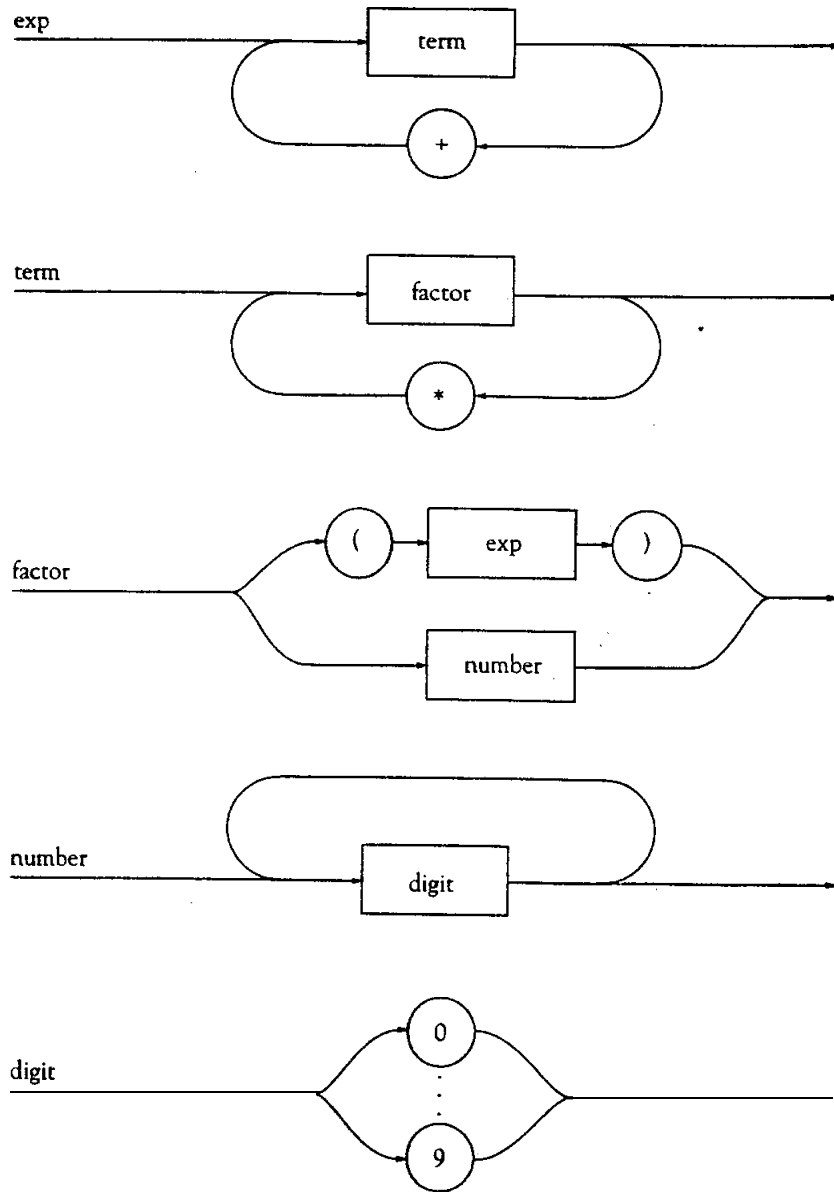
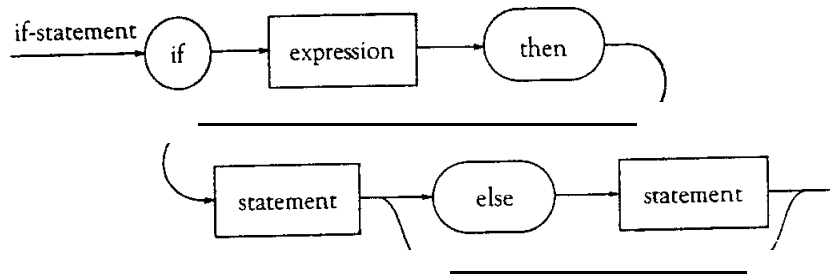
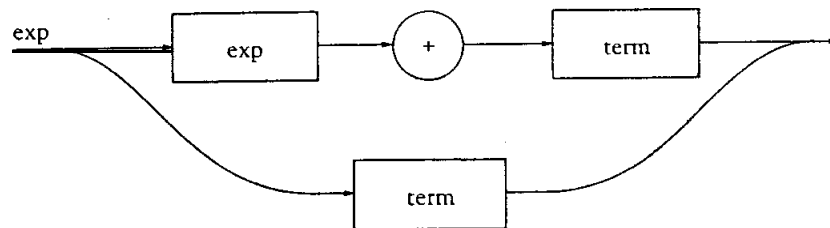


Figure 4-5 Syntax Diagrams for a Simple-Expression Grammar

ส่วนการแสดงออกของวงเล็บใหญ่ ( [ ] ) ในแผนภาพวากยสัมพันธ์ จงดูตัวอย่าง ข้อความสั่ง-if ของภาษา Pascal ข้างล่างนี้



แผนภาพวากยสัมพันธ์ปกติจะเขียนจาก EBNF ไม่ใช่เขียนจาก BNF ดังนั้นแผนภาพของ <exp> ข้างล่างนี้ จึง ผิด (incorrect)



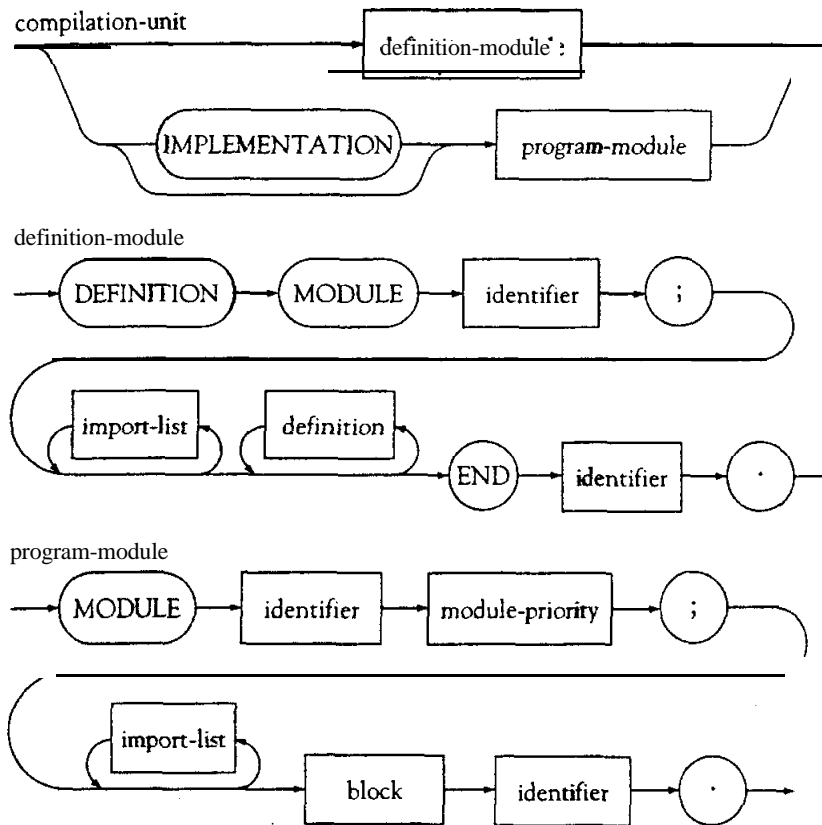
ตัวอย่างของ EBNFs และ แผนภาพวากยสัมพันธ์ของกฎไวยากรณ์ ของภาษา Modula-2 กำหนดให้แล้วในรูป 4.6

```

<compilation-unit> ::= <definition-module> I
                        [IMPLEMENTATION]<program-module>
<definition-module> ::= DEFINITION MODULE <identifier>';'
                        {<import-list>}{<definition>}
                        END <identifier>'.'
<program-module> ::= MODULE <identifier>
                        [<module-priority>]';' {<input-list>}
                        <block><identifier>'.'
                        . . .
                        . . .

```

รูป 4.6a ตัวอย่าง EBNFs สำหรับภาษา Modula-2



รูป 4.6b ตัวอย่างแผนภาพวากยสัมพันธ์ สำหรับ ภาษา Modula-2

#### 4.6 เทคนิคการวิเคราะห์กระจายและเครื่องมือที่ใช้

(Parsing techniques and tools)

ไวยากรณ์ ซึ่งเขียนด้วยรูปแบบ BNF , EBNF หรือ แผนภาพวากยสัมพันธ์ อธิบาย สายของ โทเค้น ซึ่งถูกต้อง เชิงวากยสัมพันธ์ ในภาษาโปรแกรม

(A grammar written in BNF, EBNF, or a syntax diagrams describes the strings of tokens that are syntactically legal in a programming language.)

ดังนั้น ไวยากรณ์ จึงอธิบายการกระทำโดยนัย ซึ่ง ตัววิเคราะห์กระจาย (parser) ต้องเอา สายของ โทเค้น มาวิเคราะห์กระจาย อย่างถูกต้อง นั่นคือ การสร้าง (construct) ต้นไม้การแปลง หรือ ต้นไม้วิภีช สำหรับ string ไม่ว่าจะ เป็น โดยนัย หรือ ชัดแจ้ง รูปแบบอย่างง่ายที่สุด ของตัววิเคราะห์กระจาย คือ ตัวรู้จำ ซึ่ง หมายถึง โปรแกรม ยอมรับ หรือ ปฏิเสธ สายของโทเค้น ขึ้นอยู่กับว่า มันเป็น สาย โทเค้น ถูกต้องหรือไม่ ในภาษานั้น (The simplest form of a parser

is a recognizer a program that accepts or rejects strings, based on whether they are legal strings in the language.) ตัววิเคราะห์

กระจาย โดยทั่วไป สร้างต้นไม้วิเคราะห์กระจาย (หรือ ต้นไม้วากยสัมพันธ์ แบบนามธรรม) และทำการดำเนินการอื่นๆ ให้เป็นผลสำเร็จ เช่น คำนวณค่าของนิพจน์

กำหนดไวยากรณ์ ให้ ในหนึ่งรูปแบบ จากสามรูปแบบที่อภิปรายมาแล้ว มัน สัมพันธ์กับ การกระทำของตัววิเคราะห์กระจายอย่างไร? วิธีหนึ่ง ของ การวิเคราะห์ กระจาย คือพยายามจับคู่ (match) อินพุท กับ ส่วนทางขวามือของกฎไวยากรณ์ เมื่อการ จับคู่เกิดขึ้น ทางด้านขวามือ ถูกแทนที่ หรือ ลดทอน (reduced) โดย nonterminal ทางซ้ายมือ ตัววิเคราะห์กระจายเช่นนี้ เรียกว่า ตัววิเคราะห์กระจายจากล่างขึ้นบน (bottom-up parsers) เพราะว่า มันสร้าง การแปลง และ ต้นไม้วิเคราะห์วิภีช จาก จุดแตกใบ ไปสู่ ราก บางครั้งเรียกว่า ตัววิเคราะห์กระจายแบบเลื่อน-ลดทอน (shift-reduce parser) เพราะว่ามันเลื่อนโทเค้น ไปไว้บนกองซ้อนก่อน จึงลดทอน strings ให้เป็น nonterminals



วิธีวิเคราะห์กระจายอีกวิธีหนึ่ง ที่สำคัญ ได้แก่ แบบบนลงล่าง (top-down) วิธีนี้ nonterminals ถูกขยายให้จับคู่กับโทเค็นที่เข้ามา และสร้างการแปลงโดยตรง เทคนิคการวิเคราะห์กระจายทั้งสองวิธีนี้ สามารถทำให้เป็นไปอัตโนมัติ กล่าวคือ เขียนโปรแกรมขึ้นมาเพื่อให้ แปลการอธิบาย BNF ให้เป็นตัววิเคราะห์กระจายอย่างอัตโนมัติ

เนื่องจากการวิเคราะห์กระจายจากล่างขึ้นบน ค่อนข้างดีกว่า (more powerful) การวิเคราะห์กระจายแบบจากบนลงล่าง ซึ่งปกติ เป็นวิธีที่ช้อมมากกว่า เช่น ตัวก่อกำเนิดของตัววิเคราะห์กระจาย (parser generators) (ในอดีต เรียกว่า compiler compilers)

ตัวก่อกำเนิดของตัววิเคราะห์กระจาย ชนิดหนึ่ง ซึ่งใช้กันอย่างแพร่หลาย คือ YACC (ย่อมาจาก yet another compiler compiler) ซึ่งเราจะได้ศึกษาต่อไปในหัวข้อนี้

อย่างไรก็ตาม มีอีกวิธีหนึ่งซึ่งเป็นวิธีเก่า สำหรับ การสร้างตัววิเคราะห์กระจาย ด้วยมือ จากไวยากรณ์ ซึ่งเป็นวิธีที่มีประสิทธิผลมากและยังคงใช้กันอยู่บ่อยๆ สิ่งที่สำคัญคือ มันดำเนินการโดยการหมุน (turning) nonterminals ให้เป็นกลุ่มของกระบวนการเรียกซ้ำร่วมกัน ซึ่ง การกระทำ ขึ้นอยู่กับ ทางด้านขวามือ ของ BNFs ดังนั้น จึงมีชื่อว่า การวิเคราะห์กระจายแบบเรียกซ้ำลงล่าง (recursive-descent parsing)

ส่วนทางขวามือ ถูก ตีความในกระบวนการงาน ดังนี้ โทเค็นถูกจับคู่โดยตรงกับ อินพุท โทเค็น ขณะที่สร้างโดย ตัวกราดตรวจ (scanner)

Nonterminals ถูกตีความ เช่น เรียกกระบวนการงาน ซึ่งสมนัยกับ nonterminals

ตัวอย่าง ในไวยากรณ์ภาษาอังกฤษอย่างง่าย กระบวนการ สำหรับ <sentence>, <noun-phrase> และ <article> จะเขียนดังนี้ (ด้วยรหัสเทียม คล้ายภาษา Pascal)

```

procedure sentence;
begin
    nounphrase;
    verbphrase;
end;

procedure nounphrase;
begin
    article;
    noun ;
end;

procedure article;
begin
    if token = 'a' then
        Gettoken
    else if token = 'the' then
        Gettoken
    else Error;
end;

```

ในรหัส ข้างต้นนี้ เราใช้ตัวแปรส่วนกลางชื่อ token เพื่อเก็บโทเค้นปัจจุบัน (current token) ซึ่งถูกสร้างโดยตัวกราดตรวจ (scanner)

กระบวนการ Gettoken ของตัวกราดตรวจ จะถูกเรียก เมื่อใดก็ตาม ที่ต้องการ โทเค้นตัวใหม่ การจับคู่ของโทเค้น สมัยกับ การทดสอบที่ประสบผลสำเร็จ สำหรับโทเค้นนั้น ตามด้วย การเรียก Gettoken

การวิเคราะห์กระจาย เริ่มต้นด้วย ตัวแปร token เก็บโทเค็นตัวแรกเรียบร้อยแล้ว ดังนั้น การเรียก Gettoken ต้องอยู่ก่อน การเรียกครั้งแรก ของ กระบวนการ sentence

สมมติว่า มีกระบวนการ Error อยู่ด้วย ซึ่งยกเลิก (aborts) การวิเคราะห์กระจาย (โปรดสังเกตว่า error จำเป็นเฉพาะถูกต้องพบ เมื่อ โทเค็นจริง ถูกคาดหวัง ขณะอยู่ใน กระบวนการ สำหรับ <article> เราควร มีการตรวจสำหรับ "a" หรือ "the" ในกระบวนการ สำหรับ <sentence> แต่ใน โครงร่างนี้ มันไม่จำเป็น)

ถ้าเราประยุกต์ใช้ กรรมวิธีนี้ กับ การอธิบาย BNF ของรูป 4.3 จะพบปัญหาของกฎ เรียกซ้ำแบบซ้าย สำหรับนิพจน์เช่น

$$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{term} \rangle$$

ถ้าเราพยายาม เขียนกฎนี้ เป็น กระบวนการ เรียกซ้ำแบบลงล่าง (recursive-descent) จะได้

```
procedure exp;
begin
  exp;
  if token = '+' then begin
    gettoken;
    term;
  end;
end;
```

โชคไม่ดี เมื่อกระบวนการนี้ถูกเรียก มันทำให้เกิด ลูของการเรียกซ้ำไม่รู้จบ (infinite recursive loop)ทันที (สิ่งนี้เรียกว่า "head recursion" ซึ่งไม่เหมือนกับ tail recursion) สำหรับเหตุผลนี้ นักทฤษฎีภาษา จึงได้มีการศึกษา เทคนิค left recursion removal แต่เทคนิคทั่วไป ไม่จำเป็น ณ ที่นี้ ถ้าเราเขียน กฎใน EBNF หรือเป็น syntax diagrams การเรียกซ้ำแบบซ้าย จะถูกแทนที่อย่างง่าย โดย การวนซ้ำ (สังเกตกับวงเล็บปีกกา ของ EBNF)

## ตัวอย่างเช่น กฎ

`<exp> ::= <term> It <term>}`

ใน EBNF สมนัยกับ recursive-descent code ดังนี้

```
procedure exp;
begin
    term;
    while token = '+' do begin
        gettoken;
        term;
    end;
end;
```

ดังนั้น วงเล็บปีกกา ใน EBNF แทน left recursion removal โดยใช้ การวนซ้ำ (loop)

ในทางตรงกันข้าม กฎการเรียกซ้ำแบบขวา ไม่มีปัญหาเช่นนั้น ในการวิเคราะห์  
กระจายแบบการเรียกซ้ำแบบลงล่าง และกฎ

`<exp> ::= <term> @ <exp>`

สมนัยโดยตรง กับ recursive-descent code

```
procedure exp;
begin
    term;
    if token <> '@' then Error
    else begin
        gettoken
        exp;
    end;
end;
```

มีปัญหาค้างกัน ในกฎ BNF ซึ่งแสดง โครงสร้างส่วนละเว้นได้ เช่น BNF สำหรับข้อความสั่ง - if ดังนี้

```
<if-statement> ::= if <condition> then <statement> I
                    if <condition> then <statement>
                    else <statement>
```

สิ่งนี้ไม่สามารถ แปลให้เป็นรหัส (code) ได้โดยตรง เพราะว่าทางเลือกทั้งสอง มี prefix เหมือนกัน อย่างไรก็ตาม ใน EBNF กฎข้อนี้ เขียนด้วย วงเล็บใหญ่ และ เอาส่วนของทางเลือก ที่ร่วมกันออก ตัวอย่างเช่น

```
<if-statement> ::= if <condition> then <statement>
                    [else <statement>]
```

กฎข้างต้นนี้ สมัยโดยตรงกับ recursive-descent code เมื่อวิเคราะห์ กระจาย ข้อความสั่ง - if ข้างล่างนี้

```
procedure ifstatement;
begin
    if token <> 'if' then Error
    else begin
        gettoken;
        condition;
        if token <> 'then' then Error
        else begin
            get token;
            statement;
            if token = 'else' then begin
                gettoken;
                statement;
            end;
        end;
    end;
end;
```

```
end;  
end;  
end;
```

กรรมวิธีของการเขียน ส่วนละเว้นได้ (optional parts) ของ BNF rules ใน EBNFs โดย เอา prefix ร่วมออก และการใช้วงเล็บ [ ] เรียกว่า left factoring และนี่เป็นเรื่องจำเป็นสำหรับ recursive-descent parsing ดังนั้น กฎ EBNF หรือ syntax diagrams สมัยอย่างธรรมชาติ กับ recursive-descent parser และสิ่งนี้เป็น หนึ่งในเหตุผลหลักที่เลือกใช้ ในรูป 4.7 จะเป็นโครงร่างของ recursive descent parser ที่สมบูรณ์ สำหรับ ไวยากรณ์ของนิพจน์ ในรูป 4.4

```

procedure exp;
begin
  term;
  while Token = '+' do begin
    GetToken;
    term;
  end;
end;

procedure term;
begin
  factor;
  while Token = '*' do begin
    GetToken;
    factor;
  end;
end;

procedure factor;
begin
  if Token = '(' then begin
    GetToken;
    exp;
    if Token = ')' then GetToken
    else Error
  end else number;
end;

procedure number;
begin
  digit;
  while Token in ['0'..'9'] do digit;
end;

procedure digit;
begin
  if Token in ['0'..'9'] then GetToken
  else Error;
end;

procedure parse:
begin
  GetToken;
  exp;
end;

```

**Figure 4-7** Sketch of a Recursive-Descent Parser for Simple Arithmetic Expressions

ตัวก่อกำเนิดของตัววิเคราะห์กระจาย ซึ่งเป็นที่นิยมกันมากตัวหนึ่ง คือ YACC ซึ่งมีอยู่บนระบบ UNIX ส่วนใหญ่ ผู้เขียนคือ Steve Johnson ในช่วงกลางปี ค.ศ. 1970s มัน generates ก่อกำเนิด โปรแกรม ภาษา C ซึ่งใช้อัลกอริทึม แบบล่างขึ้นบน เพื่อวิเคราะห์กระจาย ไวยากรณ์ ไวยากรณ์ที่กำหนดให้ อยู่ในรูปแบบคล้าย BNF และเกี่ยวข้องกับ การกระทำด้วย ภาษา C

YACC มีประโยชน์ ไม่ใช่เฉพาะผู้เขียนตัวแปลภาษาเท่านั้น แต่ยังเป็นประโยชน์ สำหรับนักออกแบบภาษาด้วย (YACC is useful not only to the translator writer, it is also useful to the langauge designer) : กำหนด ไวยากรณ์ให้ มันจัดหา การอธิบายของปัญหาและความกำกวมที่เป็นไปได้ อย่างไรก็ตาม เพื่อให้สามารถอ่านสารสนเทศนี้ได้ เราต้องมีความรู้ เฉพาะด้าน ของเทคนิค การ วิเคราะห์กระจาย แบบ จากล่างขึ้นบน มากกว่า



```

%{
#include <stdio.h>
#include <ctype.h>
%}

Ktoksn NUMBER

%%
command : exp {printf("%d\n", $1);}
        ; /* allows printing of the result */

exp : exp '+' term {$$ = $ + $3;}
    | term {$$ = $1;}

term : term '*' factor {$$ = $1 * $3;}
     | factor {$$ = $1;}

factor : NUMBER {$$ = $1;}
       | '(' exp ')' {$$ = $2;}

%%

void main(void)
{yyparse();}

int yylex(void)
{int c ;
 while((c = getchar()) != '\n');
 /* eliminates blanks */
 if (isdigit(c))
 {yyval = 0 ;
  while (isdigit(c))
   {yyval = 10*yyval + c - '0';
    c = getchar();}
  ungetc(c, stdin);
  return(NUMBER);}
 if (c == '\n') return 0;
 /* makes the parse stop */
 return(c);
}

void yyerror(char *s) /* allows for printing
error message */
{printf("%s\n", s);}

```

Figure 4-8 YACC Input for Simple Arithmetic Expressions

#### 4.7 โครงสร้างศัพท์ วากยสัมพันธ์ อรรถศาสตร์

(Lexics versus Syntax versus Semantics)

ไวยากรณ์ ไม่เพียงบริบท รวม การอธิบาย คำศัพท์ต่างๆ ของภาษา โดยรวมสายของตัวอักษร ซึ่งประกอบ เป็นโทเค้น ในกฎไวยากรณ์

(A context-free grammar includes a description of the tokens of a language by including the strings of characters that form the tokens in the grammar rules.)

ตัวอย่างเช่น ในไวยากรณ์ภาษาอังกฤษ ของหัวข้อ 4 . 2 โทเค้น คือ คำในภาษาอังกฤษ ได้แก่ "a", "the", "girl", "dog", "sees", "pet" และตัวอักษร ". " ส่วนใน ไวยากรณ์ นิพจน์ จำนวนเต็ม อย่างง่าย ใน หัวข้อ 4 . 3 โทเค้น หมายถึง สัญลักษณ์คำนวณ "+", "\*", วงเล็บ "(" และ ")" และเลขโดด 0 ถึง 9 รายละเอียดเฉพาะ ของ การจัดรูปแบบ เช่น ข้อตกลงของ white-space ที่กล่าวถึง ใน หัวข้อ 4.1 ถูกทิ้งให้กับ ตัวกราดตรวจ (scanner) และจำเป็นต้องกล่าวถึง ข้อตกลง แยกต่างหากจากไวยากรณ์

ประเภทของโทเค้น บางอย่าง เช่น ค่าคงที่ ไอเดนติไฟเออร์ หมายถึง ไม่ใช่ ลำดับคงที่ ของ ตัวอักษร ในตัวมันเอง แต่ถูกสร้าง ให้เป็นเซตคงที่ ของ ตัวอักษร เช่น เลขโดด 0 . . 9

(Some typical token categories, such as constants and identifiers, are not fixed sequences of characters in themselves, but are built up out of a fixed set of characters, such as the digits 0 . . 9)

ประเภทของโทเค้นเหล่านี้ บ่อยครั้ง มี โครงสร้างของมัน นิยามโดยไวยากรณ์ เช่น ตัวอย่าง กฎไวยากรณ์ สำหรับ <number> และ <digit> ใน ไวยากรณ์ของนิพจน์ อย่างไรก็ตาม มันเป็นไปได้ ที่ จะใช้ ตัวกราดตรวจ (scanner) เพื่อรู้จำ (to recognise) โครงสร้างเหล่านี้ เนื่องจากกำลังเรียกซ้ำแบบเต็ม ของ ตัววิเคราะห์กระจาย ไม่มีความจำเป็น และ ตัวกราดตรวจสามารถ รู้จำได้ โดยการดำเนินการ ทำซ้ำอย่างง่าย

สิ่งนี้มีประสิทธิภาพมากกว่า - มันทำการรู้จำ เลข และ ไอเดนติไฟเออร์ เร็วกว่า และ  
ง่ายกว่า และ มันลดขนาด ของ ตัวกราดตรวจ และ จำนวน กฎ BNFs

เพื่อแสดงให้เห็นความจริงที่ว่า เลข (a number) ใน ไวยากรณ์ของนิพจน์  
ควรจะเป็น โทเค้น ไม่ใช่แทนด้วย nonterminal เราจะเขียน ไวยากรณ์ ของ  
รูป 4-4 ใหม่ ในรูป 4-9 ดังนี้

```
<exp> ::= <term>{+ <term>}  
<term> ::= <factor>{* <factor>}  
<factor> ::= (<exp>) | NUMBER
```

รูป 4-9 เลข เป็น โทเค้น ในการคำนวณอย่างง่าย

การทำให้ สาย NUMBER เป็น อักษรตัวใหญ่ ใน ไวยากรณ์ใหม่นี้ เรากำลัง  
พูดว่า มัน ไม่ใช่ โทเค้น ซึ่ง รู้จำได้ โดยภาษา แต่ มันเป็น โครงสร้างอย่างหนึ่ง ซึ่ง  
บอกได้ โดยตัวกราดตรวจ โครงสร้างของ โทเค้น ประเภทนี้ ต้องระบุให้เป็นส่วน  
ของ ข้อตกลงคำศัพท์ ของภาษา

สัญกรณ์ ซึ่ง ใช้บ่อยมาก สำหรับ ข้อกำหนดนี้ คือ นิพจน์ปกติ (regular  
expressions) อย่างไรก็ตาม มีนักออกแบบภาษา จำนวนมาก ซึ่งเลือก ที่จะ รวม  
โครงสร้างของ โทเค้น เช่นนี้ ให้เป็นส่วนหนึ่งของ ไวยากรณ์ ด้วย ความเข้าใจว่า นัก  
ทำภาษาให้เกิดผล (implementor) อาจจะ รวมสิ่งเหล่านี้ ใน ตัวกราดตรวจ แทนที่  
จะเป็น ตัววิเคราะห์กระจาย ดังนั้น การอธิบาย ภาษา โดยใช้ BNF, EBNF หรือ  
แผนภาพวากยสัมพันธ์ ไม่เพียงแต่ รวม วากยสัมพันธ์ไว้ แต่ยังรวม โครงสร้างคำศัพท์ (หรือ  
lexics) ส่วนใหญ่ ของ ภาษาโปรแกรมไว้ด้วย ถึงแม้ว่า ขอบเขต (boundary)  
ระหว่าง โครงสร้างวากยสัมพันธ์ และ โครงสร้างคำศัพท์ ปกติ ไม่เขียนชัดเจน แต่ขึ้นอยู่กับ  
จุดของการมอง ของ นักออกแบบภาษา และ นักทำภาษาให้เกิดผล

เช่นเดียวกัน มันเป็นจริงด้วย สำหรับวากยสัมพันธ์ และอรรถศาสตร์ เราได้

มีการเข้าถึง มาแล้วว่า วากยสัมพันธ์ คือ สิ่งใดก็ตาม ซึ่งถูกนิยาม ได้ด้วย ไวยากรณ์ ไม่พึ่งบริบท และอรรถศาสตร์ คือ สิ่งใดก็ตาม ซึ่งนิยามไม่ได้ ด้วย ไวยากรณ์ ไม่พึ่งบริบท อย่างไรก็ตาม มีผู้เขียนหนังสือ จำนวนมาก ได้รวมคุณสมบัติ ซึ่งเราเรียกว่า อรรถศาสตร์ เป็น คุณสมบัติเชิงวากยสัมพันธ์ ของภาษา ตัวอย่างเช่น กฎต่างๆ ของ การประกาศ ก่อน ใช้ตัวแปร และการไม่มีการประกาศใหม่ ของ ไอเดนติไฟเออร์ ภายใน โปรซีเดอร์ สิ่งเหล่านี้เป็นกฎ ซึ่งเรียกว่า ไวยากรณ์ไม่พึ่งบริบท (context-sensitive grammar) และไม่สามารถเขียนด้วยกฎ พึ่งบริบท ดังนั้น เราจึงชอบมากกว่า ที่จะคิดว่า มันเป็น ความหมาย ไม่ใช่กฎ วากยสัมพันธ์

ข้อขัดแย้งอีกประการหนึ่งระหว่างวากยสัมพันธ์ และอรรถศาสตร์ เกิดขึ้น เมื่อภาษานั้น ต้องการ สายอักขระบางอย่าง ให้เป็น ไอเดนติไฟเออร์นิยามมาแล้ว (predefined identifiers) ไม่ใช่ คำสงวน (reserved words) จากหัวข้อ 4.1 ซึ่งกล่าวว่า คำสงวน หมายถึง สายคงที่ของ อักขระ ซึ่งเป็นโทเค็นโดยตัวมันเอง และนำไปใช้ เป็น ไอเดนติไฟเออร์ ไม่ได้

(Recall that reserved words are fixed strings of characters that are tokens themselves and that cannot be used as identifiers.)

ตัวอย่างเช่น "begin", "end", "while", "do", และ "procedure" strings ทั้งหมดนี้ เป็นคำสงวน ในภาษา Pascal

อย่างไรก็ตาม string "true" และ "false" ไม่ใช่คำสงวน แต่เป็น predefined identifiers นั่นคือ มันเป็น ไอเดนติไฟเออร์ ซึ่ง มีความหมายคงที่ ในภาษา แต่ ความหมายนี้ สามารถเปลี่ยนแปลงได้ โดยการ ประกาศให้ใหม่ ภายใน โปรแกรม

(Predefined identifiers are identifiers that have a fixed meaning in the language, but this meaning can be changed by redeclaring them within a program.)

อย่างไรก็ตาม การทำเช่นนี้ เป็นแนวปฏิบัติของการเขียนโปรแกรมที่แย่มาก เพราะว่ามันจะ ไม่มีความหมายปกติ เช่น ค่าคงที่แบบบูล ด้วย ค่าปกติ วิธีที่ดีกว่า คือ ทำให้ "true" และ "false" เป็น syntactic entities คือ เป็นคำสงวน ซึ่งสัมพันธ์ กับ อรรถศาสตร์ คงที่

## แบบฝึกหัด

1. Add subtraction and division to the
  - (a) BNF
  - (b) EBNFand (c) syntax diagrams of simple arithmetic expression (Figures 4-3, 4-4, and 4-5) Be sure to give them the appropriate precedence.
  
2. Add the mod and power operation to
  - (a) the arithmetic BNF or
  - (b) EBNF.Use % for the mod operation and ^ for the power operation. Recall that mod is left-associative, like division, but that power is right-associative. (Thus  $2^{2^3} = 256$ , not 64)
  
3. Unary minuses can be added in several ways to the arithmetic expression grammar of Figure 4-3 or the grammar from Exercise 1. Revise the BNF and EBNF for each of the cases that follow so that it satisfies the stated rule:
  - (a) At most one unary minus is allowed in each expression, and it must come at the beginning of an expression, so  $-2 - 3$  is legal (and equals -5) and  $-2 - (-3)$  is legal, but  $-2 - -3$  is not.
  - (b) At most one unary minus is allowed before a number or left parenthesis, so  $-2 - -3$  is legal but  $- -2$  and  $-2 - - -3$  are not.

(c) Arbitrarily many **unary** minuses are allowed before numbers and left parentheses, so everything above is legal, but, for example,  $2 = t3$  is not.

4. จงวาดรูปต้นไม้วิภังค์ และต้นไม้วากยสัมพันธ์แบบนามธรรม ของนิพจน์คำนวณข้างล่างนี้  
(Draw parse trees and abstract syntax trees for the arithmetic expressions) :

(a) ((2))

(b)  $3 + 4 * 5 + 6 * 7$

(c)  $3 * 4 + 5 * 6 + 7$

(d)  $3 * (4 + 5) * (6 + 7)$

(e)  $(2 + (3 + (4 + 5)))$

5. เป็นไปได้หรือไม่ ที่จะมียาษา ซึ่งไม่มีคำสงวนใดๆ ให้อภิปราย (Is it possible to have a language without any reserved words? Discuss.)

6. A number is defined in the grammar of Figure 4-3 using a left-recursive rule. However, it could also be defined using a right-recursive rule :

$\langle \text{number} \rangle ::= \langle \text{digit} \rangle \langle \text{number} \rangle | \langle \text{digit} \rangle$

which is better, or does it matter? Why?

7. Given the following BNF :

$\langle \text{exp} \rangle ::= (\langle \text{list} \rangle) | a$

$\langle \text{list} \rangle ::= \langle \text{list} \rangle, \langle \text{exp} \rangle | \langle \text{exp} \rangle$

- (a) Write **EBNF** rules and syntax diagrams for the language.

(b) Draw the parse tree for  $((a, a), a, (a))$ .

(c) Write a recursive-descent **recognizer** for the language.