



# Module 10

## Compiler Basics



# Module Ten

- This week, we are going to talk about :

- The Compiler

Language Definition - Grammar

Lexical Analysis

Syntactical Analysis

Semantic Analysis

- Material adapted from L. Beck, System Software, 1997  
A. Aho, Compilers, 1986



# Basic Compiler Function

- Translate

Text written in one language  
converted into a second language
- System Software
- High Level language

FORTRAN, ALGOL, Pascal, C, C++  
Write solution to problem as engineer knows
- Machine code

Solution as executed by the hardware.

# Basic Compiler Tasks

- Lexical Analysis - Scanner
  - Read the input text and find the token objects.
- Syntactical Analysis - Parser
  - From the tokens, discover the structure of the program
- Semantic Analysis - Code Generator
  - From the structure, create the machine language.



# What is this ?

- IF THEN THEN THEN = ELSE ELSE ELSE = THEN

IF (THEN) THEN

THEN = ELSE;

ELSE

ELSE = THEN;



# A program

```
PROGRAM Stats
```

```
VAR
```

```
    sum, sumSQ, I, value, mean, variance : INTEGER
```

```
BEGIN
```

```
    sum      := 0;
```

```
    sumSQ    := 0;
```

```
    FOR I := 1 TO 100 DO
```

```
        BEGIN
```

```
            READ ( value );
```

```
            sum      := sum + value;
```

```
            sumSQ    := sumSQ + value * value
```

```
        END;
```

```
    mean      := sum DIV 100;
```

```
    variance := sumSQ DIV 100 - mean * mean;
```

```
    WRITE ( mean, variance )
```

```
END.
```



# Grammar

- Basic Language definition
- Formal description of the Programming Language
  - SYNTAX - form of legal statements
  - Structure - appropriate collection of the statements
  - Rules - 'how to' for constructing statements
- Example:     ADD
  - $I := J + K$      ;  $I, J, K,$      integers
  - $Z := X + Y$      ;  $X, Y, Z,$      real
- The meanings of the statements is called Semantics.



# Grammar

- BNF - Backus - Naur Form
  - Notation for writing a grammar.       $\langle \text{symbol} \rangle ::= \text{tokens}$
  - Early 1960's
  - Used to define ALGOL
  - Top - Down set of rules for defining the objects within the programming language
  - Simple and widely used
  - Generally sufficient





# Backus - Naur Form

- Set of Rules that define the syntax of a language object
- $\langle \text{symbol} \rangle ::= \text{tokens}$
- This  $\langle \text{non-terminal symbol} \rangle$  is defined to be, this set of tokens
- Tokens - terminal symbol
- $\langle \text{symbol} \rangle ::= \langle \text{set of symbols or tokens} \rangle$



# A Simple Pascal Grammar

1.  $\langle \text{prog} \rangle ::= \text{PROGRAM } \langle \text{prog-name} \rangle \text{ VAR } \langle \text{dec-list} \rangle \text{ BEGIN } \langle \text{stmt-list} \rangle \text{ END.}$
2.  $\langle \text{prog-name} \rangle ::= \text{id}$
3.  $\langle \text{dec-list} \rangle ::= \langle \text{dec} \rangle \mid \langle \text{dec-list} \rangle , \langle \text{dec} \rangle$
4.  $\langle \text{dec} \rangle ::= \langle \text{id-list} \rangle : \langle \text{type} \rangle$
5.  $\langle \text{id-list} \rangle ::= \text{id} \mid \langle \text{id-list} \rangle , \text{id}$
6.  $\langle \text{type} \rangle ::= \text{INTEGER}$
7.  $\langle \text{stmt-list} \rangle ::= \langle \text{stmt} \rangle \mid \langle \text{stmt-list} \rangle ; \langle \text{stmt} \rangle$
8.  $\langle \text{stmt} \rangle ::= \langle \text{assign} \rangle \mid \langle \text{read} \rangle \mid \langle \text{write} \rangle \mid \langle \text{for} \rangle$

# A Simple Pascal Grammar

- 9. `<assign>` ::= `id := <exp>`
- 10. `<exp>` ::= `<term> | <exp> + <term> | <exp> - <term>`
- 11. `<term>` ::= `<factor> | <term> * <factor> | <term> DIV <factor>`
- 12. `<factor>` ::= `id | int | ( <exp> )`
- 13. `<read>` ::= `READ ( <id-list> )`
- 14. `<write>` ::= `WRITE ( <id-list> )`
- 15. `<for>` ::= `FOR <index-exp> DO <body>`
- 16. `<index-exp>` ::= `id := <exp> TO <exp>`
- 17. `<body>` ::= `<stmt> | BEGIN <stmt-list> END`



# BNF Examples

- $\langle \text{stmt} \rangle ::= \langle \text{assign} \rangle \mid \langle \text{read} \rangle \mid \langle \text{write} \rangle \mid \langle \text{for} \rangle$
- $\langle \text{read} \rangle ::= \text{READ} ( \langle \text{id-list} \rangle )$
- $\langle \text{id-list} \rangle ::= \text{id} \mid \langle \text{id-list} \rangle , \text{id}$



# Rule order

- Note the order of the rules.
- $\langle \text{exp} \rangle ::= \langle \text{term} \rangle \mid \langle \text{exp} \rangle + \langle \text{term} \rangle \mid \langle \text{exp} \rangle - \langle \text{term} \rangle$
- $\langle \text{term} \rangle ::= \langle \text{factor} \rangle \mid \langle \text{term} \rangle * \langle \text{factor} \rangle \mid \langle \text{term} \rangle \mathbf{DIV} \langle \text{factor} \rangle$
- $\langle \text{factor} \rangle ::= \text{id} \mid \text{int} \mid ( \langle \text{exp} \rangle )$

- Rule order assures that

$$Q + 7 * A - 1$$

- Parses as if parentheses existed:  $Q + ( 7 * A ) - 1$



# Grammar Review

- Formal definition of a programming Language
- Rules for forming the structures of the language.
- BNF notation
- Meanings



# Summary

- Compiler Basics

Language Definition - Grammar

Next: Lexical Analysis