



ĐẠI HỌC BÁCH KHOA HÀ NỘI  
VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

# Experiment in Compiler Construction

## Parser design

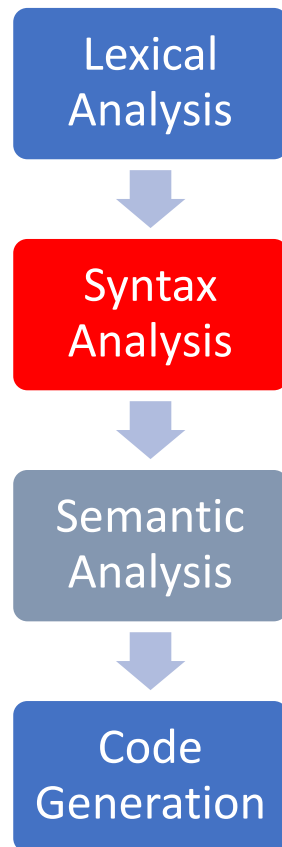
School of Information and Communication  
Technology

Hanoi university of technology

# Content

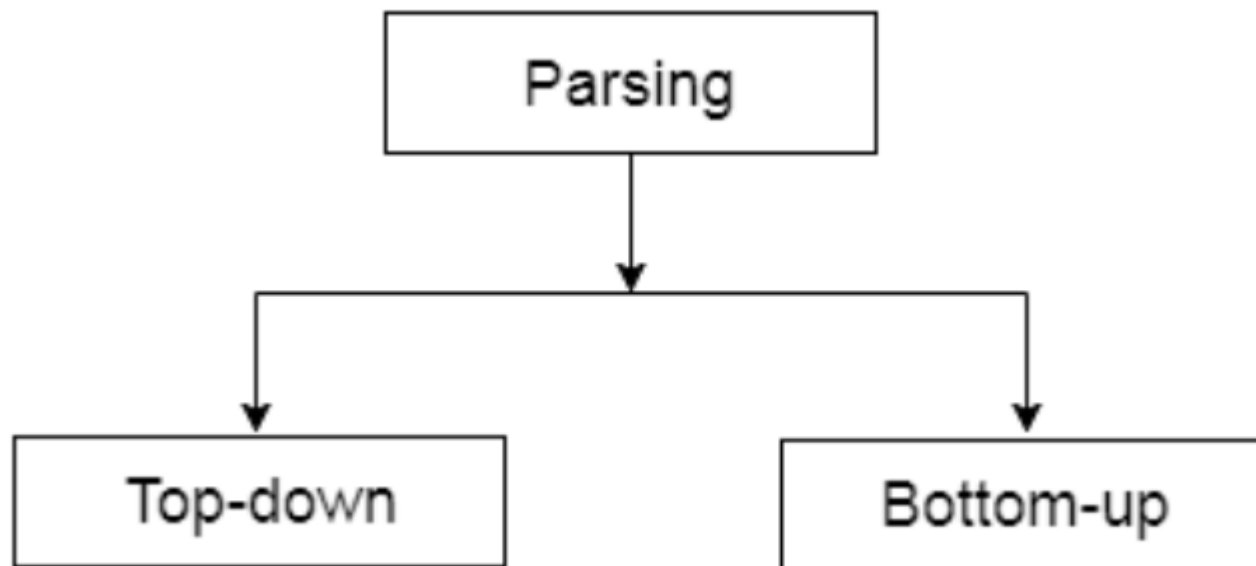
- Overview
- KPL grammar
- Parser implementation

# Tasks of a parser



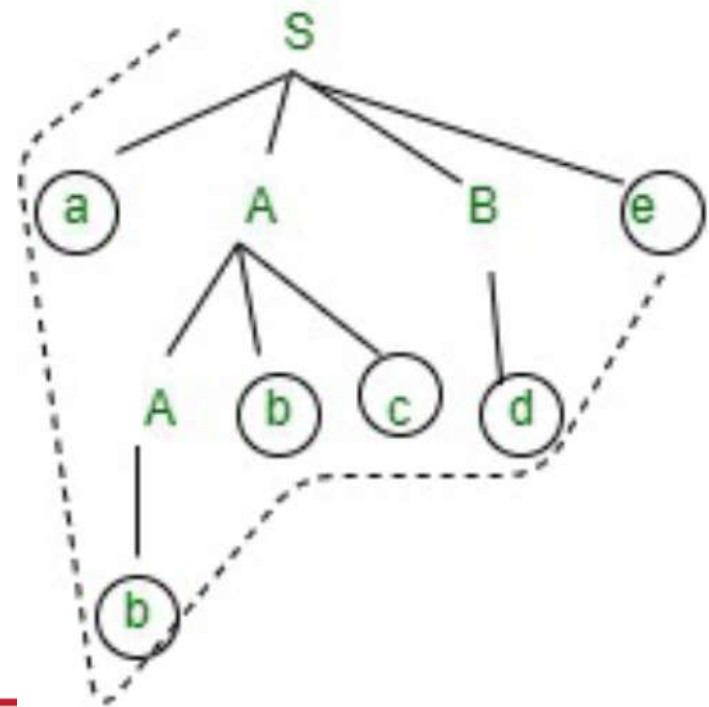
- Check the syntactic structure of a given program
  - Syntactic structure is given by Grammar
- Invoke semantic analysis and code generation
  - In an one-pass compiler, this module is very important since this forms the skeleton of the compiler

# Classification of parsing techniques



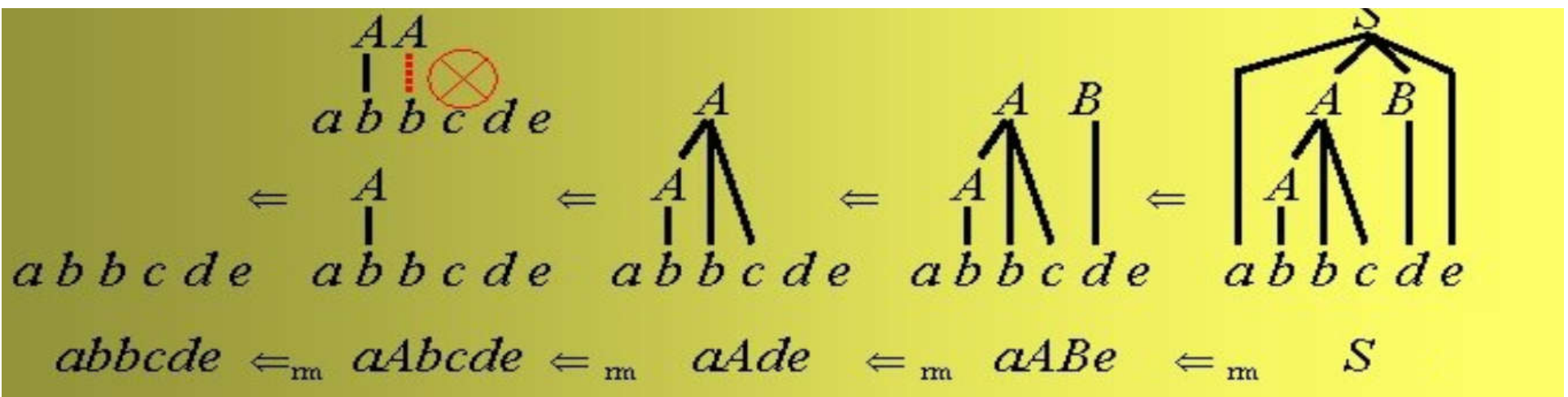
# Top down parsing

- Construct a parse tree from the root to the leaves, reading the given string from left-to-right
- It follows left most derivation.
- If a variable contains more than one possibilities, selecting 1 is difficult.
- Example: Given grammar G with a set of production rules
  - G: (1)  $S \rightarrow a A B e$
  - (2, 3)  $A \rightarrow A b c | b$
  - (4)  $B \rightarrow d$
- input: abbcde

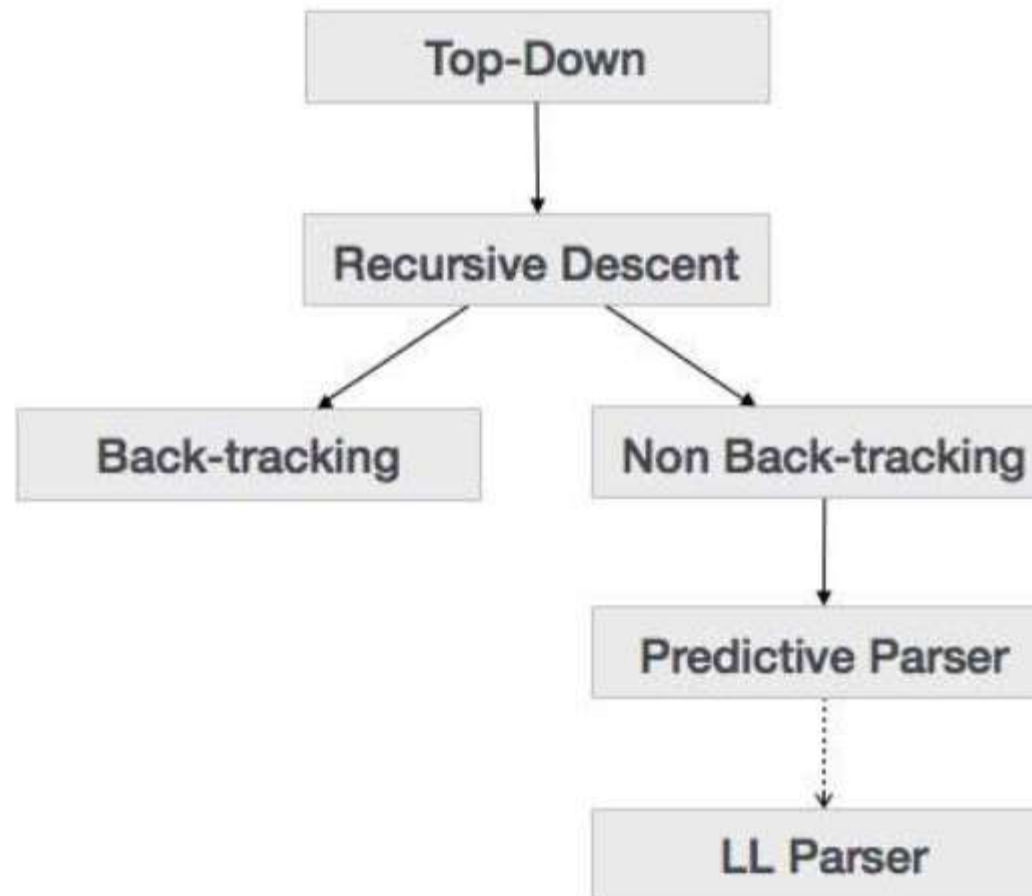


# Bottom up parsing

- Construct a parse tree from the leaves to the root: left-to-right reduction
- It follows the rightmost derivation
- Example: Given grammar  $G$  with a set of production rules
  - $G: (1) \ S \rightarrow a A B e$   
 $A \rightarrow A b c | b$   
 $B \rightarrow d$
  - input:  $abbcd e$



# Top down parsing methods



# Recursive-descent parsing

- A top-down parsing method
- *Descent*: the direction in which the parse tree is traversed (or built).
- Use a set of *mutually recursive* procedures (one procedure for each nonterminal symbol)
  - Start the parsing process by calling the procedure that corresponds to the start symbol
  - Each production becomes one branch in procedure for its LHS
- We consider a **special type of recursive-descent parsing** called **predictive parsing**
  - Use a lookahead symbol to decide which production to use



# Recursive Descent Parsing

- For every BNF rule (production) of the form

$\langle \text{phrase1} \rangle \rightarrow E$

the parser defines a function to parse phrase1 whose body is to parse the rule E

```
void compilePhrase1( )  
{ /* parse the rule E */ }
```

- Where E consists of a sequence of non-terminal and terminal symbols
- Requires **no left recursion** in the grammar.

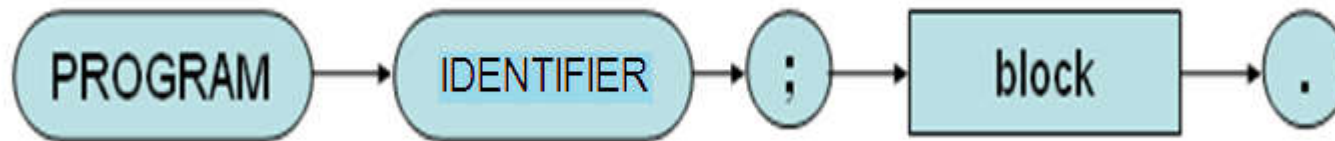
# Parsing a rule

- A sequence of non-terminal and terminal symbols,  
 $Y_1 Y_2 Y_3 \dots Y_n$   
is recognized by parsing each symbol in turn
- For each non-terminal symbol,  $Y$ , call the corresponding parse function  
**compileY**
- For each terminal symbol,  $y$ , call a function  
**eat(y)**  
that will check if  $y$  is the next symbol in the source program
  - The terminal symbols are the token types from the lexical analyzer
  - If the variable `currentsymbol` always contains the next token:  

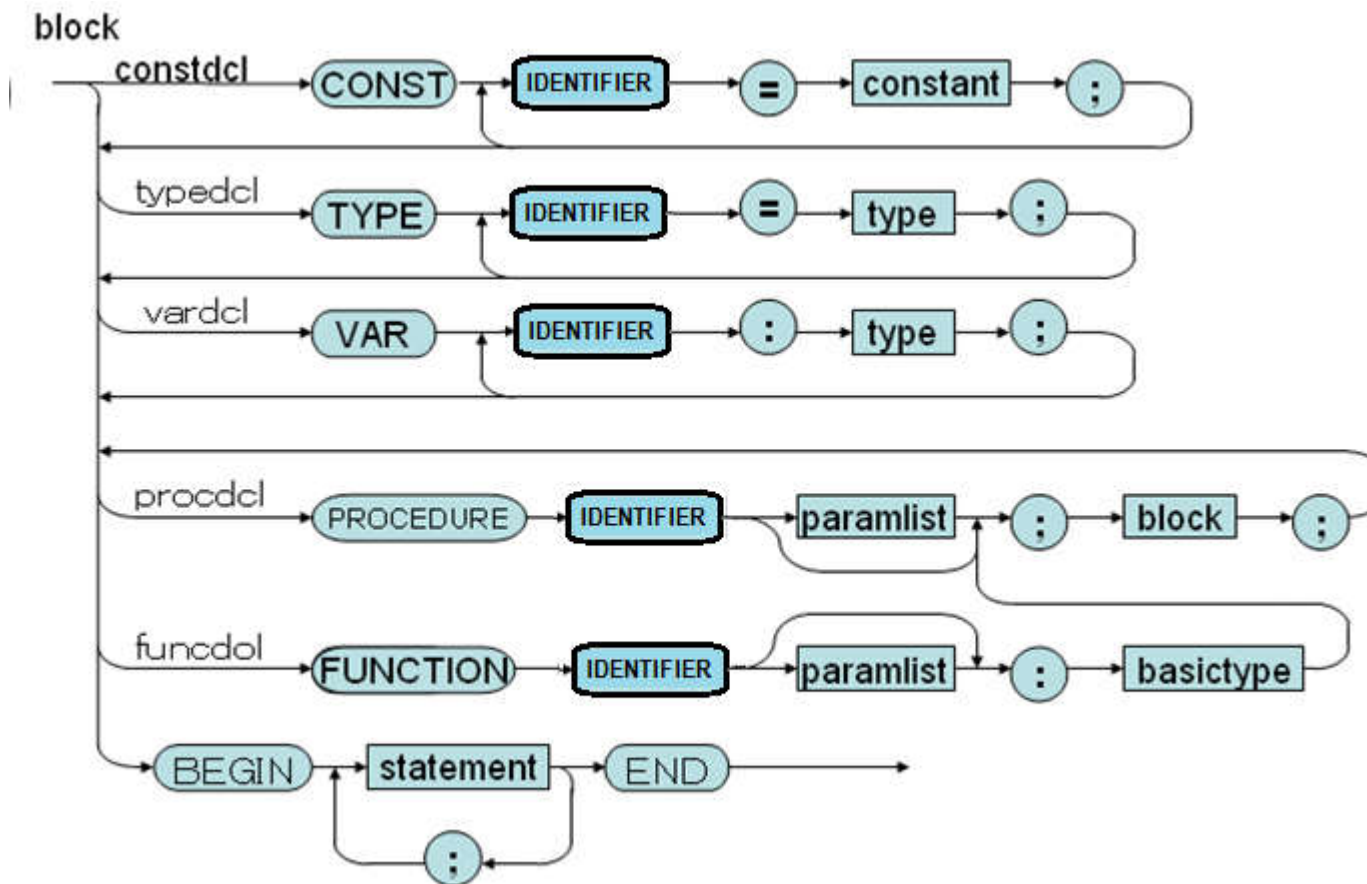
```
eat(y) :  
    if (currentsymbol == y)  
    then getNextToken()  
    else SyntaxError()
```

# Syntax diagram of KPL (the whole program)

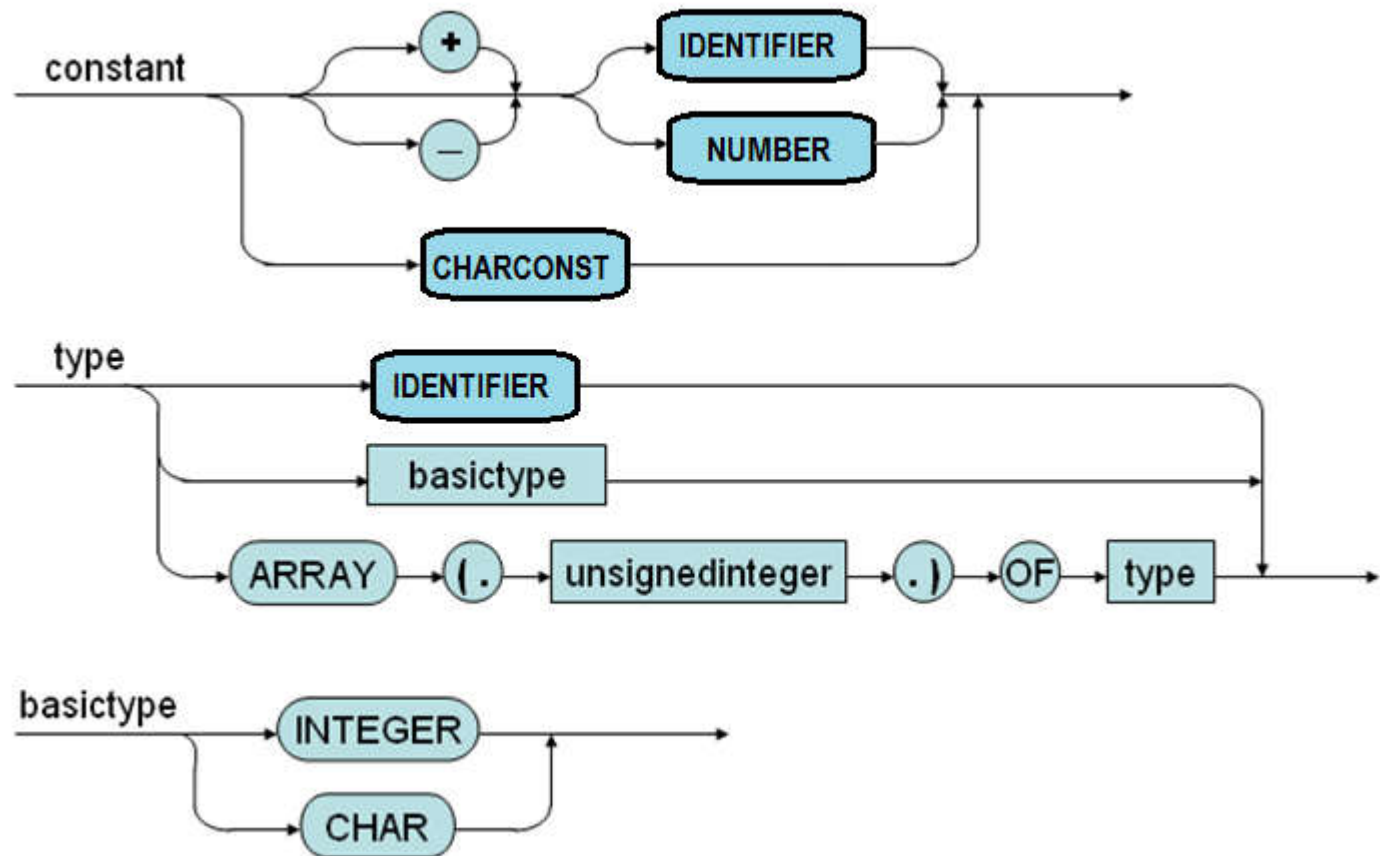
program



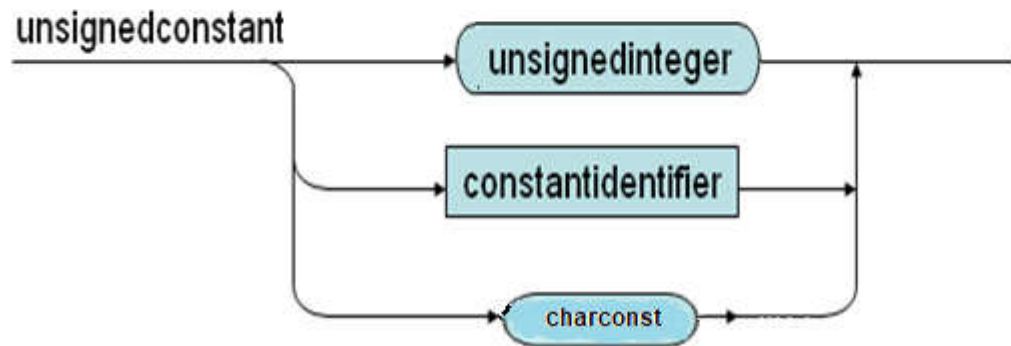
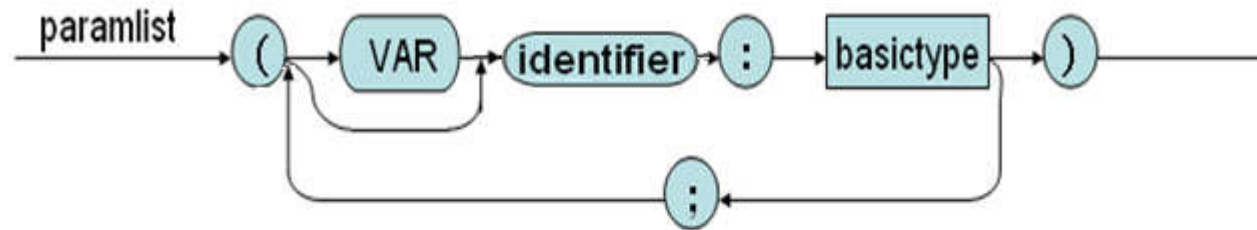
# Syntax diagrams of KPL (block)



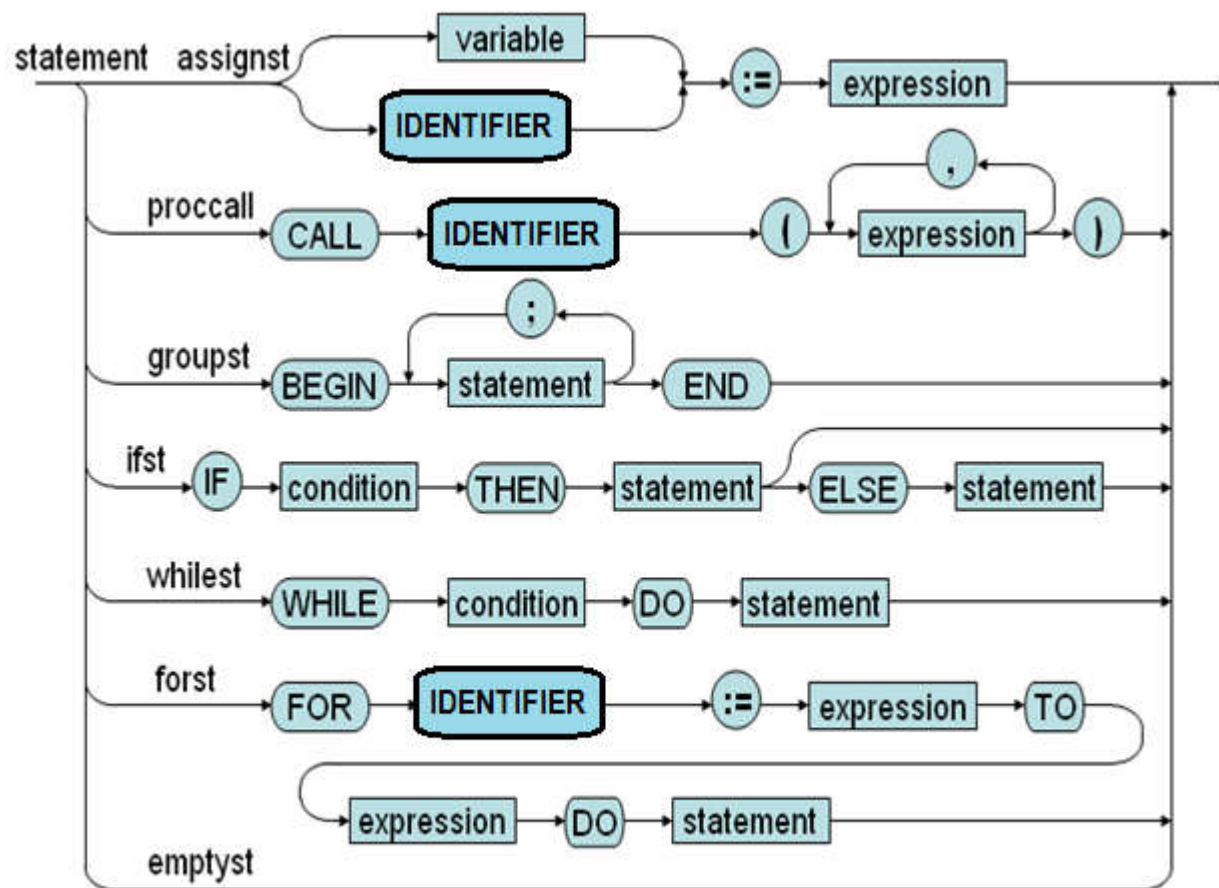
# Syntax diagrams of KPL(declarations)



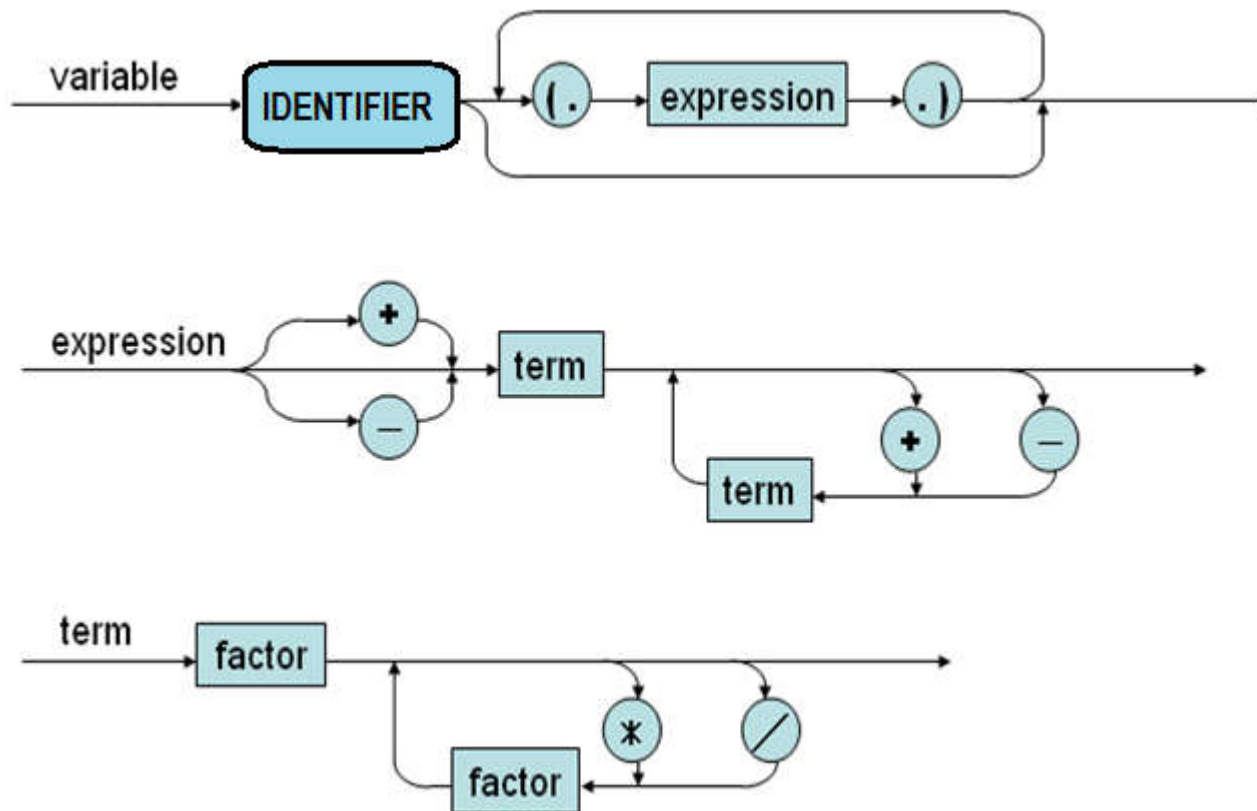
# Syntax diagrams of KPL (list of parameters, unsigned constants)



# Syntax diagram of KPL(lệnh)

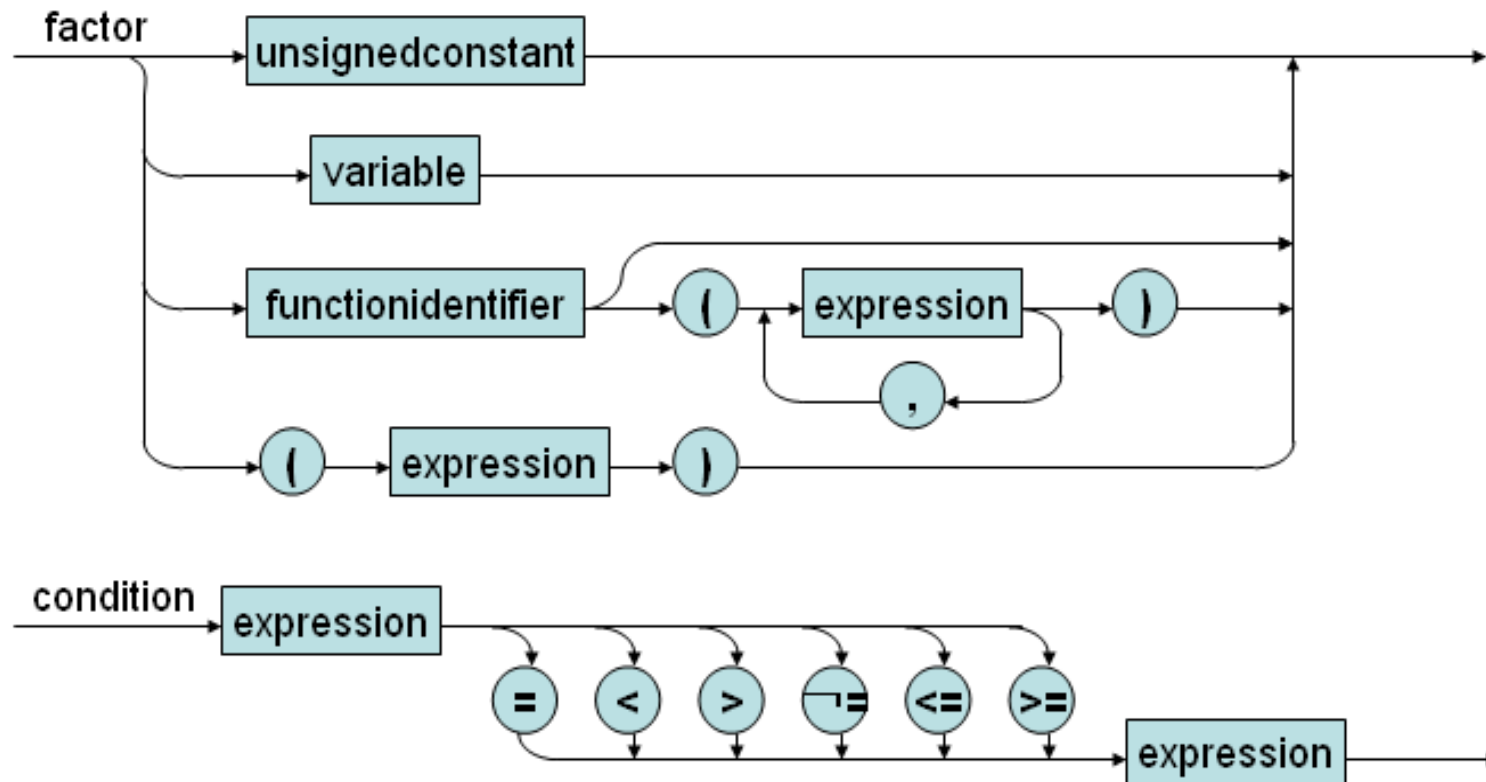


# Syntax diagrams of KPL(expression)





# Syntax diagram of KPL (factor, condition)



# KPL Grammar in BNF

- Construct a grammar  $G$  based on syntax diagram
- Perform left recursive elimination (already)
- Perform left factoring

# KPL Grammar in BNF

- 01) `<Prog> ::= KW_PROGRAM TK_IDENT SB_SEMICOLON <Block> SB_PERIOD`
- 02) `<Block> ::= KW_CONST <ConstDecl> <ConstDecls> <Block2>`
- 03) `<Block> ::= <Block2>`
- 04) `<Block2> ::= KW_TYPE <TypeDecl> <TypeDecls> <Block3>`
- 05) `<Block2> ::= <Block3>`
- 06) `<Block3> ::= KW_VAR <VarDecl> <VarDecls><Block4>`
- 07) `<Block3> ::= <Block4>`
- 08) `<Block4> ::= <SubDecls><Block5>`
- 09) `<Block4> ::= |<Block5>`
- 10) `<Block5> ::= KW_BEGIN <Statements> KW_END`
- 11) `<ConstDecls> ::= <ConstDecl> <ConstDecls>`
- 12) `<ConstDecls> ::= ε`
- 13) `<ConstDecl> ::= TK_IDENT SB_EQUAL <Constant> SB_SEMICOLON`
- 14) `<TypeDecls> ::= <TypeDecl> <TypeDecls>`
- 15) `<TypeDecls> ::= ε`
- 16) `<TypeDecl> ::= TK_IDENT SB_EQUAL <Type> SB_SEMICOLON`
- 17) `<VarDecls> ::= <VarDecl> <VarDecls>`
- 18) `<VarDecls> ::= ε`
- 19) `<VarDecl> ::= TK_IDENT SB_COLON <Type> SB_SEMICOLON`

# KPL Grammar in BNF

- 20) `<SubDecls> ::= <FunDecl> <SubDecls>`
- 21) `<SubDecls> ::= <ProcDecl> <SubDecls>`
- 22) `<SubDecls> ::= ε`
  
- 23) `<FunDecl> ::= KW_FUNCTION TK_IDENT <Params> SB_COLON <BasicType>  
SB_SEMICOLON  
    <Block> SB_SEMICOLON`
  
- 24) `<ProcDecl> ::= KW_PROCEDURE TK_IDENT <Params> SB_SEMICOLON <Block>  
SB_SEMICOLON`
  
- 25) `<Params> ::= SB_LPAR <Param> <Params2> SB_RPAR`
- 26) `<Params> ::= ε`
  
- 27) `<Params2> ::= SB_SEMICOLON <Param> <Params2>`
- 28) `<Params2> ::= ε`
  
- 29) `<Param> ::= TK_IDENT SB_COLON <BasicType>`
- 30) `<Param> ::= KW_VAR TK_IDENT SB_COLON <BasicType>`
  
- 31) `<Type> ::= KW_INTEGER`
- 32) `<Type> ::= KW_CHAR`
- 33) `<Type> ::= TK_IDENT`
- 34) `<Type> ::= KW_ARRAY SB_LSEL TK_NUMBER SB_RSEL KW_OF <Type>`

# KPL Grammar in BNF

- 35) `<BasicType> ::= KW_INTEGER`
- 36) `<BasicType> ::= KW_CHAR`
  
- 37) `<UnsignedConstant> ::= TK_NUMBER`
- 38) `<UnsignedConstant> ::= TK_IDENT`
- 39) `<UnsignedConstant> ::= TK_CHAR`
  
- 40) `<Constant> ::= SB_PLUS <Constant2>`
- 41) `<Constant> ::= SB_MINUS <Constant2>`
- 42) `<Constant> ::= <Constant2>`
- 43) `<Constant> ::= TK_CHAR`
  
- 44) `<Constant2> ::= TK_IDENT`
- 45) `<Constant2> ::= TK_NUMBER`
  
- 46) `<Statements> ::= <Statement> <Statements2>`
  
- 47) `<Statements2> ::= SB_SEMICOLON <Statement> <Statements2>`
- 48) `<Statements2> ::= ε`

# KPL Grammar in BNF

- 49) `<Statement> ::= <AssignSt>`
- 50) `<Statement> ::= <CallSt>`
- 51) `<Statement> ::= <GroupSt>`
- 52) `<Statement> ::= <IfSt>`
- 53) `<Statement> ::= <WhileSt>`
- 54) `<Statement> ::= <ForSt>`
- 55) `<Statement> ::= ε`
  
- 56) `<AssignSt> ::= <Variable> SB_ASSIGN <Expression>`
- 57) `<AssignSt> ::= TK_IDENT SB_ASSIGN <Expression>`
  
- 58) `<CallSt> ::= KW_CALL TK_IDENT <Arguments>`
  
- 59) `<GroupSt> ::= KW_BEGIN <Statements> KW_END`
  
- 60) `<IfSt> ::= KW_IF <Condition> KW_THEN <Statement> <ElseSt>`
  
- 61) `<ElseSt> ::= KW_ELSE <Statement>`
- 62) `<ElseSt> ::= ε`
  
- 63) `<WhileSt> ::= KW_WHILE <Condition> KW_DO <Statement>`
- 64) `<ForSt> ::= KW_FOR TK_IDENT SB_ASSIGN <Expression> KW_TO  
                  <Expression> KW_DO <Statement>`

# KPL Grammar in BNF

65)  $\langle \text{Arguments} \rangle ::= \text{SB\_LPAR } \langle \text{Expression} \rangle \langle \text{Arguments2} \rangle \text{SB\_RPAR}$

66)  $\langle \text{Arguments} \rangle ::= \varepsilon$

67)  $\langle \text{Arguments2} \rangle ::= \text{SB\_COMMA } \langle \text{Expression} \rangle \langle \text{Arguments2} \rangle$

68)  $\langle \text{Arguments2} \rangle ::= \varepsilon$

69)  $\langle \text{Condition} \rangle ::= \langle \text{Expression} \rangle \langle \text{Condition2} \rangle$

70)  $\langle \text{Condition2} \rangle ::= \text{SB\_EQ } \langle \text{Expression} \rangle$

71)  $\langle \text{Condition2} \rangle ::= \text{SB\_NEQ } \langle \text{Expression} \rangle$

72)  $\langle \text{Condition2} \rangle ::= \text{SB\_LE } \langle \text{Expression} \rangle$

73)  $\langle \text{Condition2} \rangle ::= \text{SB\_LT } \langle \text{Expression} \rangle$

74)  $\langle \text{Condition2} \rangle ::= \text{SB\_GE } \langle \text{Expression} \rangle$

75)  $\langle \text{Condition2} \rangle ::= \text{SB\_GT } \langle \text{Expression} \rangle$

# KPL Grammar in BNF

- 76) `<Expression> ::= SB_PLUS <Expression2>`
- 77) `<Expression> ::= SB_MINUS <Expression2>`
- 78) `<Expression> ::= <Expression2>`
  
- 79) `<Expression2> ::= <Term> <Expression3>`
  
- 80) `<Expression3> ::= SB_PLUS <Term> <Expression3>`
- 81) `<Expression3> ::= SB_MINUS <Term> <Expression3>`
- 82) `<Expression3> ::=  $\epsilon$`
  
- 83) `<Term> ::= <Factor> <Term2>`
  
- 84) `<Term2> ::= SB_TIMES <Factor> <Term2>`
- 85) `<Term2> ::= SB_SLASH <Factor> <Term2>`
- 86) `<Term2> ::=  $\epsilon$`
- 87) `<Factor> ::= <UnsignedConstant>`
- 88) `<Factor> ::= <Variable>`
- 89) `<Factor> ::= <FunctionApptication>`
- 90) `<Factor> ::= SB_LPAR <Expression> SB_RPAR`
  
- 91) `<Variable> ::= TK_IDENT <Indexes>`
- 92) `<FunctionApplication> ::= TK_IDENT <Arguments>`
  
- 93) `<Indexes> ::= SB_LSEL <Expression> SB_RSEL <Indexes>`
- 94) `<Indexes> ::=  $\epsilon$`



# Input – output in KPL

- Input: Use functions
  - ReadI: Read an integer. No parameter
  - ReadC: Read a character. No parameter

Example

```
var a: integer;  
a:= ReadI;
```

- Output: Use procedures
  - WriteI: Print an integer. 1 parameter
  - WriteC: Print a character. 1 parameter
  - WriteLn: Print the newline character.

Ví dụ

```
call WriteI(a);  
call WriteLn;
```

# KPL program

- Write a function that calculates the square of an integer
- Write a program to calculate the sum of the squares of the first n natural numbers. n is read from the keyboard

# Solution

```
program example5;  
  (* sum of the squares of the first n natural  
  numbers *)  
  var n : integer; i: integer; sum: integer;  
  
  function f(k : integer) : integer;  
    begin  
      f := k * k;  
    end;  
  
  BEGIN  
    n := readI;  
    sum := 0;  
    for i:=1 to n do  
      sum:= sum + f(i);  
    call writeln;  
    call writeI(f(n));  
  END. (* example*)
```

# Implementation

- In general, KPL is a LL(1) grammar
- design a top-down parser
  - *lookAhead* token
  - Parsing terminals
  - Parsing non-terminals
    - Constructing a parsing table
      - Computing FIRST() and FOLLOW()

- Example

02) Block ::= KW\_CONST ConstDecl ConstDecls Block2 =>RHS1

03) Block ::= Block2 =>RHS2

FIRST(RHS1)={KW\_CONST}

FIRST(RHS2)={KW\_TYPE, KW\_VAR, KW\_FUNCTION, KW\_PROCEDURE, KW\_BEGIN}

FIRST(RHS1)  $\cap$  FIRST(RHS2) =  $\emptyset$

LookAhead =KW\_BEGIN =>RHS2 is chosen =>LL(1)

# Recursive-descent parsing

- A top-down parsing method
- The term *descent* refers to the direction in which the parse tree is traversed (or built).
- Use a set of *mutually recursive* procedures (one procedure for each nonterminal symbol)
  - Start the parsing process by calling the procedure that corresponds to the start symbol
  - Each production becomes one branch in procedure for its LHS
- We consider a special type of recursive-descent parsing called predictive parsing
  - Use a lookahead symbol to decide which production to use

# Recursive Descent Parsing

- For every BNF rule (production) of the form

$\langle \text{phrase1} \rangle \rightarrow E$

the parser defines a function to parse phrase1 whose body is to parse the rule E

```
void compilePhrase1( )  
{ /* parse the rule E */ }
```

- Where E consists of a sequence of non-terminal and terminal symbols
- Requires **no left recursion** in the grammar.

# Parsing a rule

- A sequence of non-terminal and terminal symbols,  
 $Y_1 Y_2 Y_3 \dots Y_n$   
is recognized by parsing each symbol in turn
  - For each non-terminal symbol,  $Y$ , call the corresponding parse function `compileY`
  - For each terminal symbol,  $y$ , call a function
- that will check if  $y$  is the next symbol in the source program
- The terminal symbols are the token types from the lexical analyzer
  - If the variable `currentsymbol` always contains the next token:

```
eat(y) :  
    if (LookAhead == y)  
    then getNextToken()  
    else SyntaxError()
```

# lookAhead token

- Look ahead the next token

```
Token *currentToken;    // Token vừa đọc  
Token *lookAhead;      // Token xem trước
```

```
void scan(void) {  
    Token* tmp = currentToken;  
    currentToken = lookAhead;  
    lookAhead = getValidToken();  
    free(tmp);  
}
```



# Parsing terminal symbol

```
void eat(TokenType tokenType) {  
    if (lookAhead->tokenType == tokenType) {  
        printToken(lookAhead);  
        scan();  
    } else  
        missingToken(tokenType, lookAhead->lineNo, lookAhead->colNo);  
}
```

# Invoking parser

```
int compile(char *fileName) {  
    if (openInputStream(fileName) == IO_ERROR)  
        return IO_ERROR;  
  
    currentToken = NULL;  
    lookAhead = getValidToken();  
  
    compileProgram();  
  
    free(currentToken);  
    free(lookAhead);  
    closeInputStream();  
    return IO_SUCCESS;  
}
```

# Parsing non-terminal symbol

Example: **Program**

**Prog ::= KW\_PROGRAM TK\_IDENT SB\_SEMICOLON Block SB\_PERIOD**

```
void compileProgram(void) {  
    assert("Parsing a Program ....");  
    eat(KW_PROGRAM);  
    eat(TK_IDENT);  
    eat(SB_SEMICOLON);  
    compileBlock();  
    eat(SB_PERIOD);  
    assert("Program parsed!");  
}
```

# Parsing statement

Example: **Statement**

FIRST(Statement) = {TK\_IDENT, KW\_CALL, KW\_BEGIN, KW\_IF, KW\_WHILE,  
KW\_FOR,  $\epsilon$ }

FOLLOW(Statement) = {SB\_SEMICOLON, KW\_END, KW\_ELSE}

/\* Predict parse table for Expression \*/

Input

Production

TK_IDENT	49) Statement ::= AssignSt
KW_CALL	50) Statement ::= CallSt
KW_BEGIN	51) Statement ::= GroupSt
KW_IF	52) Statement ::= IfSt
KW_WHILE	53) Statement ::= WhileSt
KW_FOR	54) Statement ::= ForSt

SB_SEMICOLON	55) $\epsilon$
KW_END	55) $\epsilon$
KW_ELSE	55) $\epsilon$

Others

Error

# Parsing statement

Example: **Statement**

```
void compileStatement(void) {
    switch (lookAhead->tokenType)
    {
        case TK_IDENT:
            compileAssignSt();
            break;
        case KW_CALL:
            compileCallSt();
            break;
        case KW_BEGIN:
            compileGroupSt();
            break;
        case KW_IF:
            compileIfSt();
            break;
        case KW_WHILE:
            compileWhileSt();
            break;
```

```
        case KW_FOR:
            compileForSt();
            break;
            // check FOLLOW tokens
        case SB_SEMICOLON:
        case KW_END:
        case KW_ELSE:
            break;
            // Error occurs
        default:
            error(ERR_INVALIDSTATEMENT,
lookAhead->lineNo, lookAhead-
>colNo);
            break;
    }
}
```

# LHS with more than 1 RHS

## Two alternatives for Basic Type

34) `BasicType ::= KW_INTEGER`

35) `BasicType ::= KW_CHAR`

```
void compileBasicType(void) {  
    switch (lookAhead->tokenType) {  
        case KW_INTEGER:  
            eat(KW_INTEGER);  
            break;  
        case KW_CHAR:  
            eat(KW_CHAR);  
            break;  
        default:  
            error(ERR_INVALIDBASICTYPE, lookAhead->lineNo,  
lookAhead->colNo);  
            break;  
    }  
}
```

# Loop processing

Loop for sequence of constant declarations

10) `ConstDecls ::= ConstDecl ConstDecls`

11) `ConstDecls ::=  $\epsilon$`

```
void compileConstDecls(void) {  
    while (lookAhead->tokenType == TK_IDENT)  
        compileConstDecl();  
}
```

# Sometimes you should refer to syntax diagrams

## Syntax of Term (using BNF)

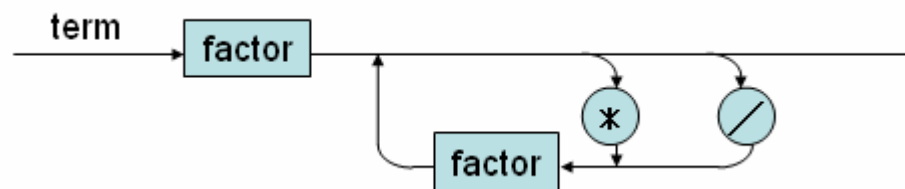
82)  $\text{Term} ::= \text{Factor Term2}$

83)  $\text{Term2} ::= \text{SB\_TIMES Factor Term2}$

84)  $\text{Term2} ::= \text{SB\_SLASH Factor Term2}$

85)  $\text{Term2} ::= \varepsilon$

## Syntax of Term (using Syntax Diagram)





# Process rules for Term : 2 functions with Follow set checking

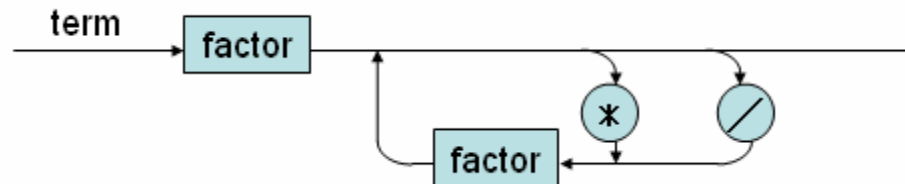
```
void compileTerm(void)
{ compileFactor();
  compileTerm2();
}
```

```
void compileTerm2(void) {
  switch (lookAhead->tokenType) {
  case SB_TIMES:
    eat(SB_TIMES);
    compileFactor();
    compileTerm2();
    break;
  case SB_SLASH:
    eat(SB_SLASH);
    compileFactor();
    compileTerm2();
    break;
  // check the FOLLOW set
  case SB_PLUS:
  case SB_MINUS:
  case KW_TO:
  case KW_DO:
```

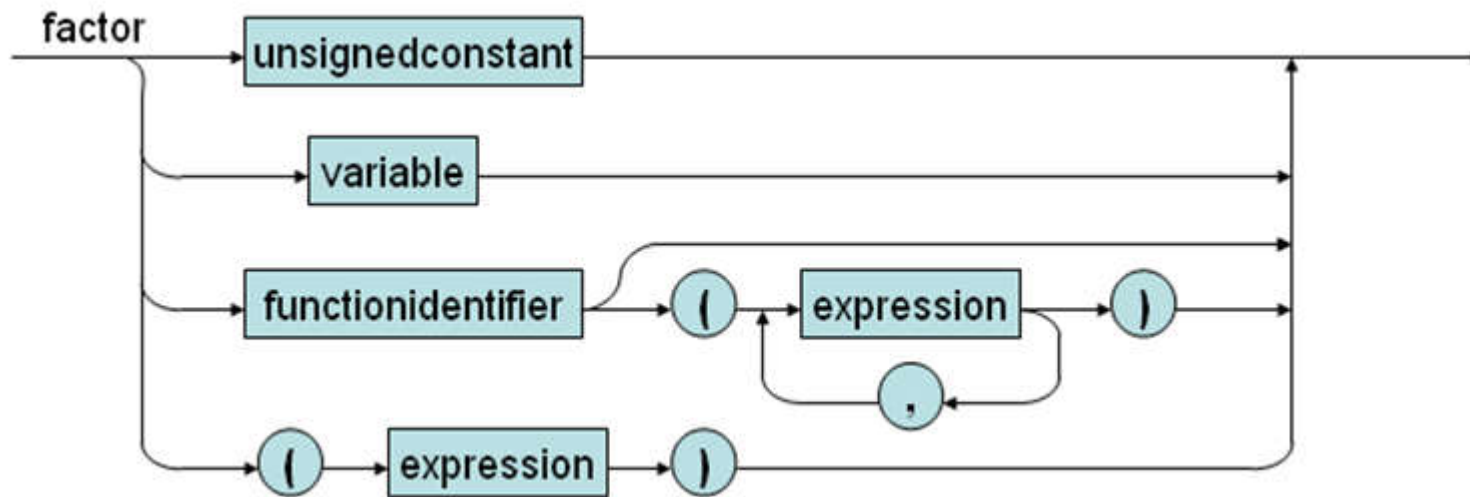
```
  case SB_RPAR:
    case SB_COMMA:
    case SB_EQ:
    case SB_NEQ:
    case SB_LE:
    case SB_LT:
    case SB_GE:
    case SB_GT:
    case SB_RSEL:
    case SB_SEMICOLON:
    case KW_END:
    case KW_ELSE:
    case KW_THEN:
      break;
    default:
      error(ERR_INVALIDTERM, lookAhead->lineNo,
        lookAhead->colNo);
  }
}
```

# Process term with syntax diagram

```
void compileTerm(void)
{
    compileFactor();
    while(lookAhead->tokenType == SB_TIMES ||
        lookAhead->tokenType == SB_SLASH)
    {
        switch (lookAhead->tokenType)
        {
            case SB_TIMES:
                eat(SB_TIMES);
                compileFactor();
                break;
            case SB_SLASH:
                eat(SB_SLASH);
                compileFactor();
                break;
        }
    }
}
```



# Syntax diagram of factor in KPL



$\text{FIRST}(\text{unsignedconstant}) = \{\text{TK\_NUMBER}, \text{TK\_IDENT}, \text{TK\_CHAR}\}$

$\text{FIRST}(\text{variable}) = \{\text{TK\_IDENT}\}$

$\text{FIRST}(\text{functioncall}) = \{\text{TK\_IDENT}\}$

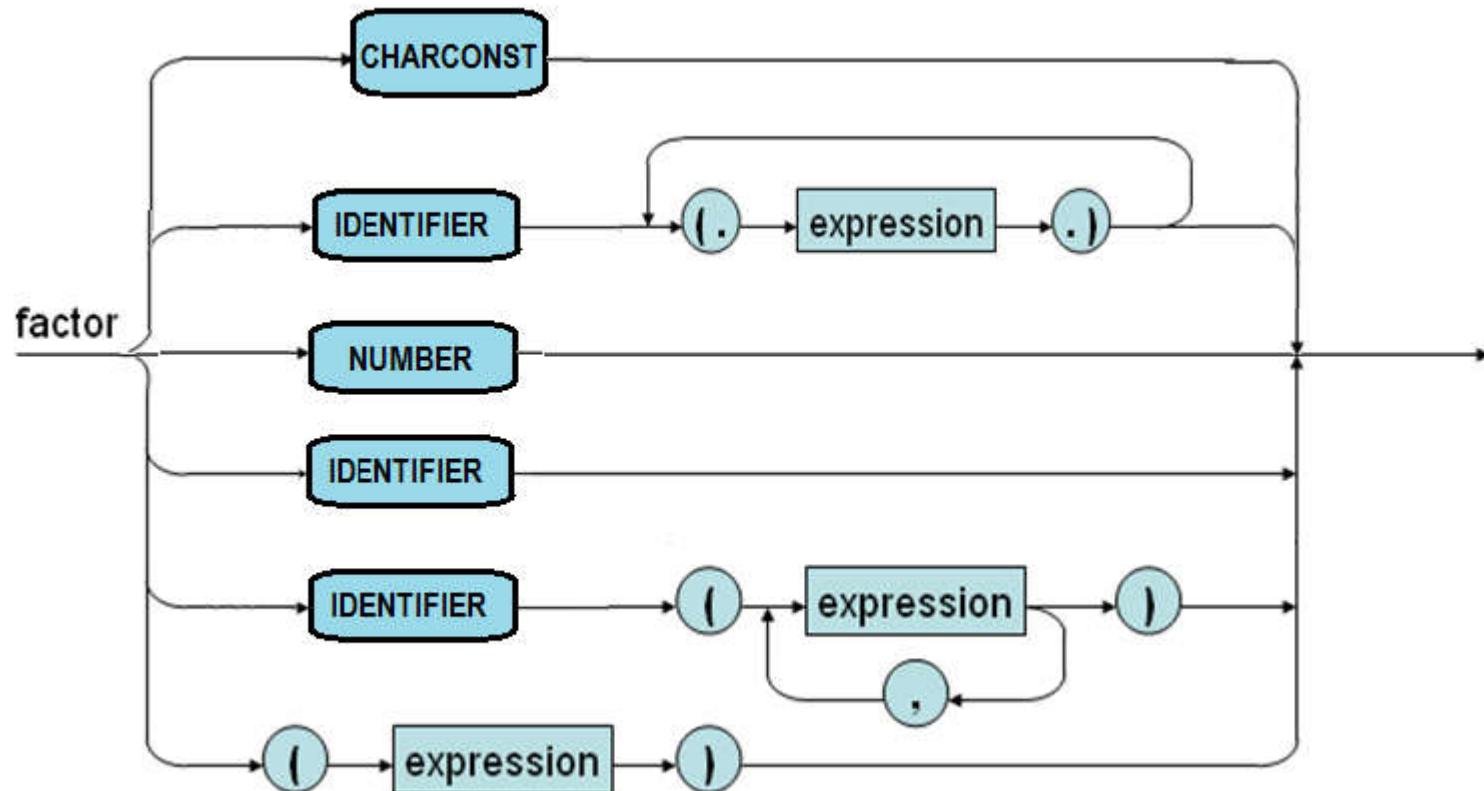
$\text{FIRST}(\text{unsignedconstant}) \cap \text{FIRST}(\text{functioncall}) = \{\text{TK\_IDENT}\}$

$\text{FIRST}(\text{unsignedconstant}) \cap \text{FIRST}(\text{variable}) = \{\text{TK\_IDENT}\}$

$\text{FIRST}(\text{variable}) \cap \text{FIRST}(\text{functioncall}) = \{\text{TK\_IDENT}\}$

=>violation of LL(1) condition

# After separating and merging



```

void compileFactor(void) {
    switch (lookAhead->tokenType) {
    case TK_NUMBER:
        eat(TK_NUMBER);
        break;

    case TK_CHAR:
        eat(TK_CHAR);
        break;

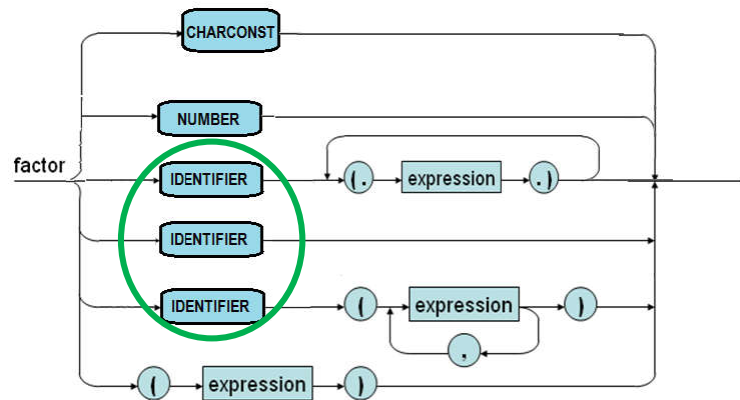
    case TK_IDENT:
        eat(TK_IDENT);
        switch (lookAhead->tokenType) {
            case SB_LSEL:
                compileIndexes();
                break;

            case SB_LPAR:
                compileArguments();
                break;

            default: break;
        }
        break;
    }
}

```

# Compile a factor



```

case SB_LPAR:
    eat(SB_LPAR);
    compileExpression();
    eat(SB_RPAR);
    break;

default:
    error(ERR_INVALIDFACTOR,
    lookAhead->lineNo, lookAhead->colNo);
}
}

```