

# Experiment in Compiler Construction

## Parser design

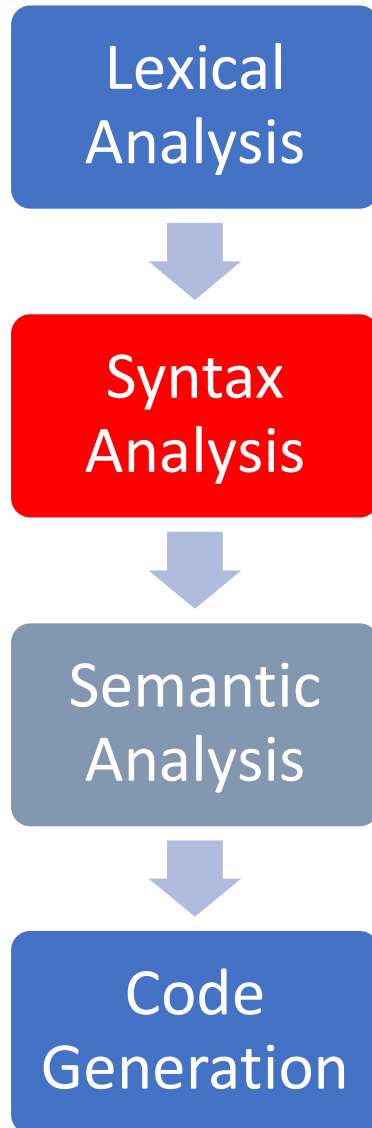
ONE LOVE. ONE FUTURE.

# Content

---

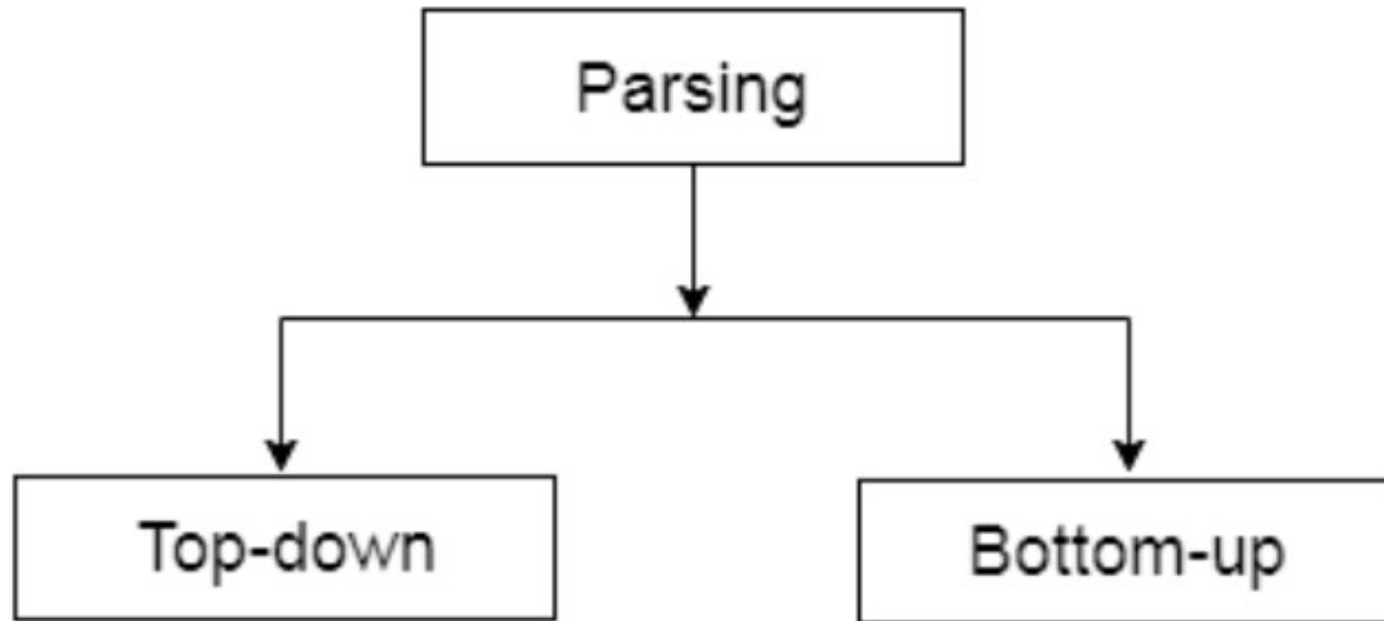
- Overview
- KPL grammar
- Parser implementation

# Tasks of a parser



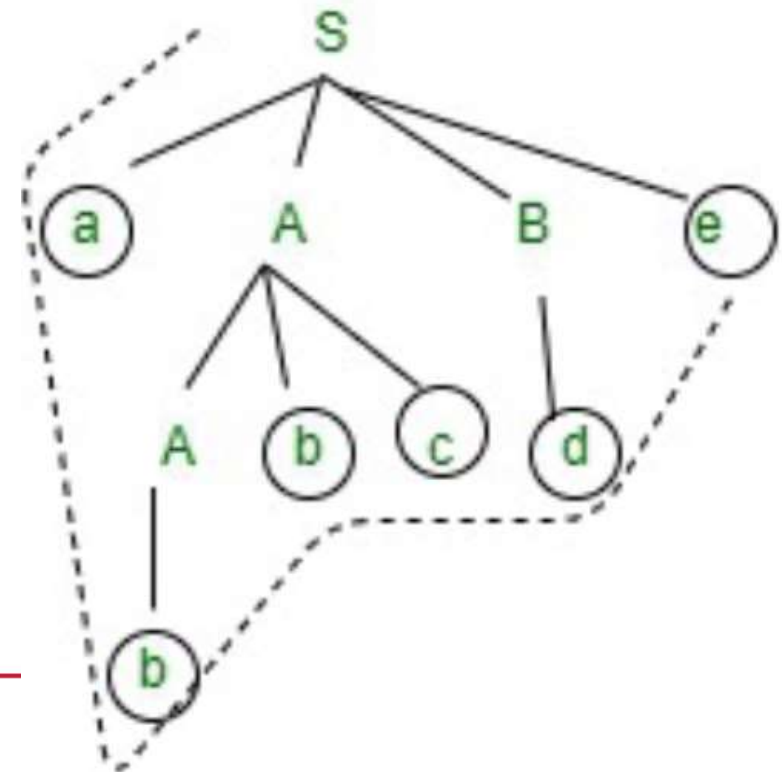
- Check the syntactic structure of a given program
  - Syntactic structure is given by Grammar
- Invoke semantic analysis and code generation
  - In an one-pass compiler, this module is very important since this forms the skeleton of the compiler

# Classification of parsing techniques



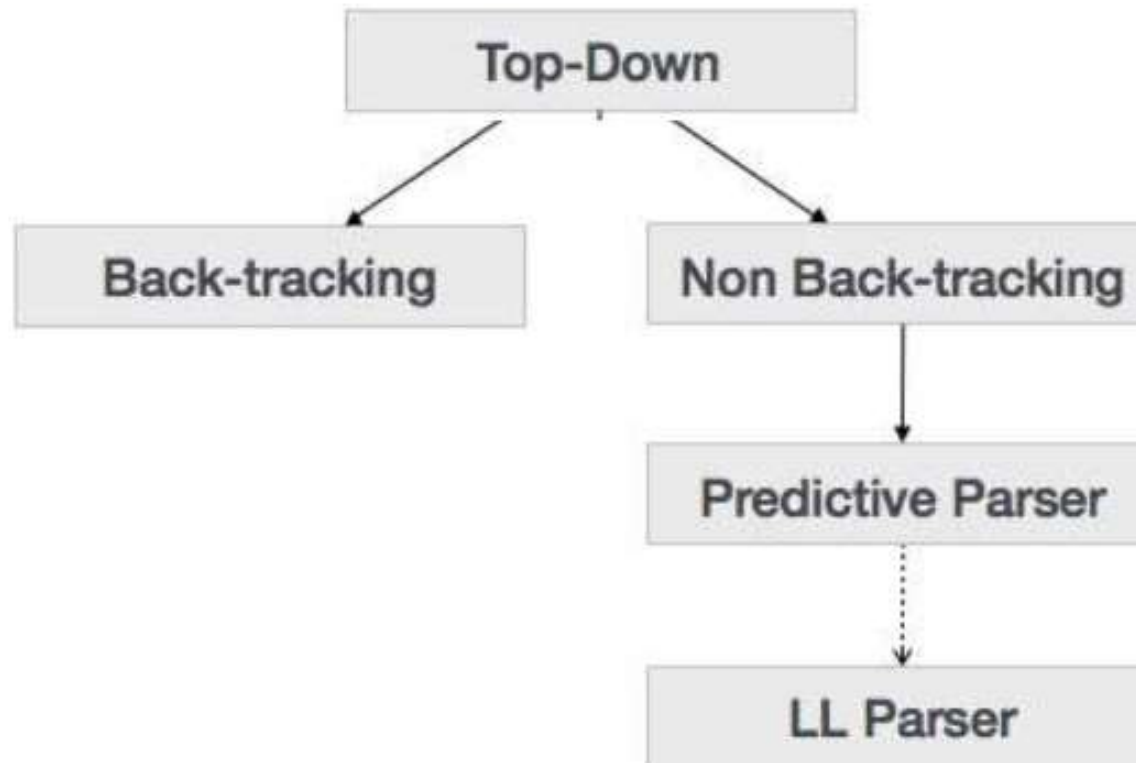
# Top down parsing

- 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





# 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.

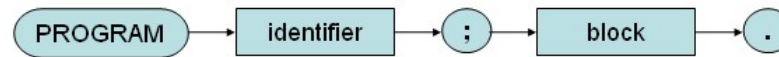
# Parsing a rule

- 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:  

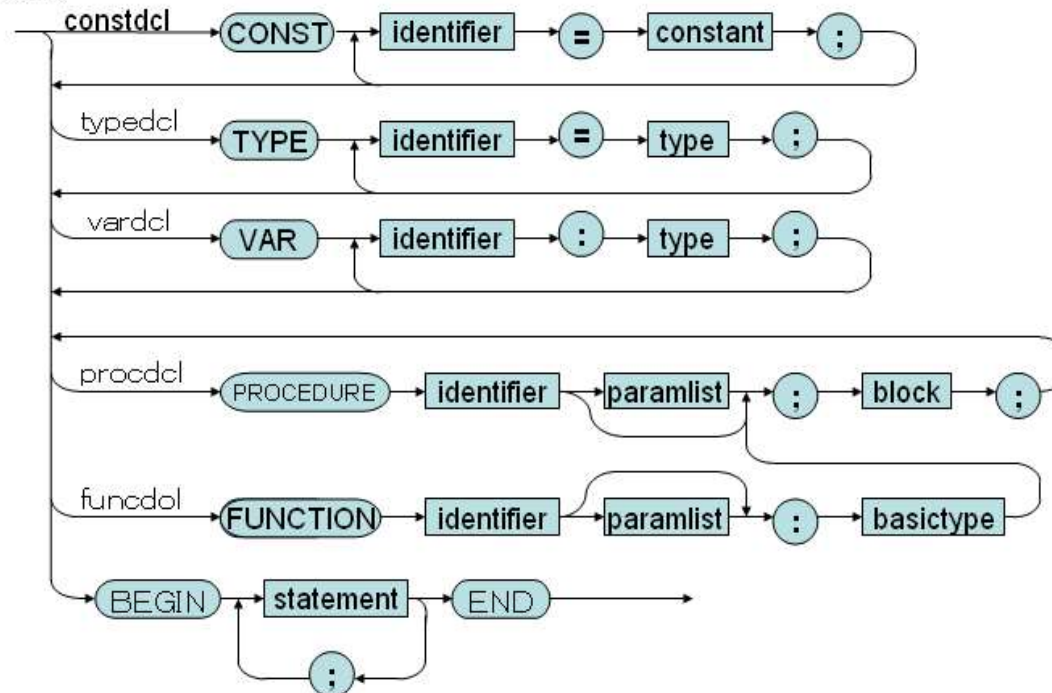
```
eat(y) :  
    if (currentsymbol == y)  
    then getNextToken()  
    else SyntaxError()
```

# Syntax diagram of KPL

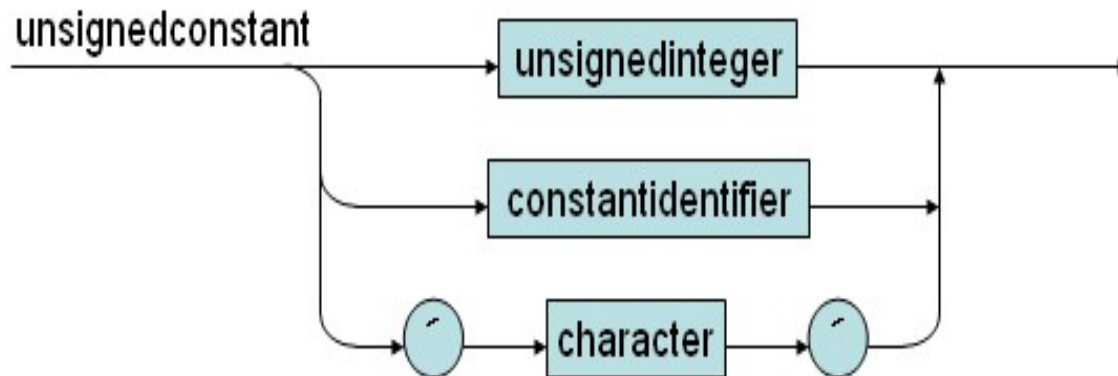
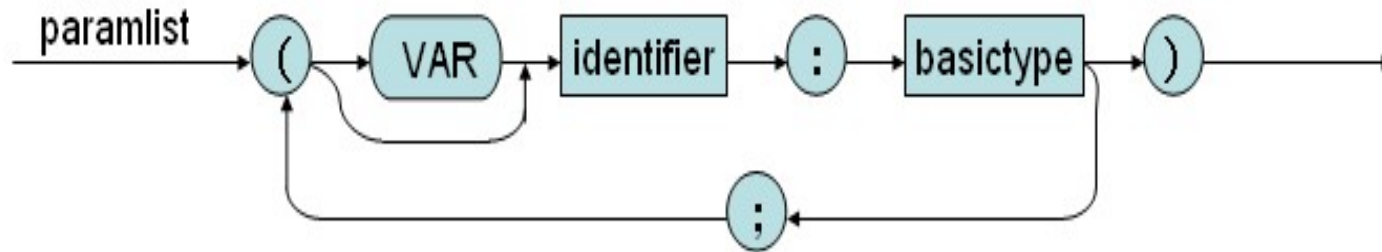
program



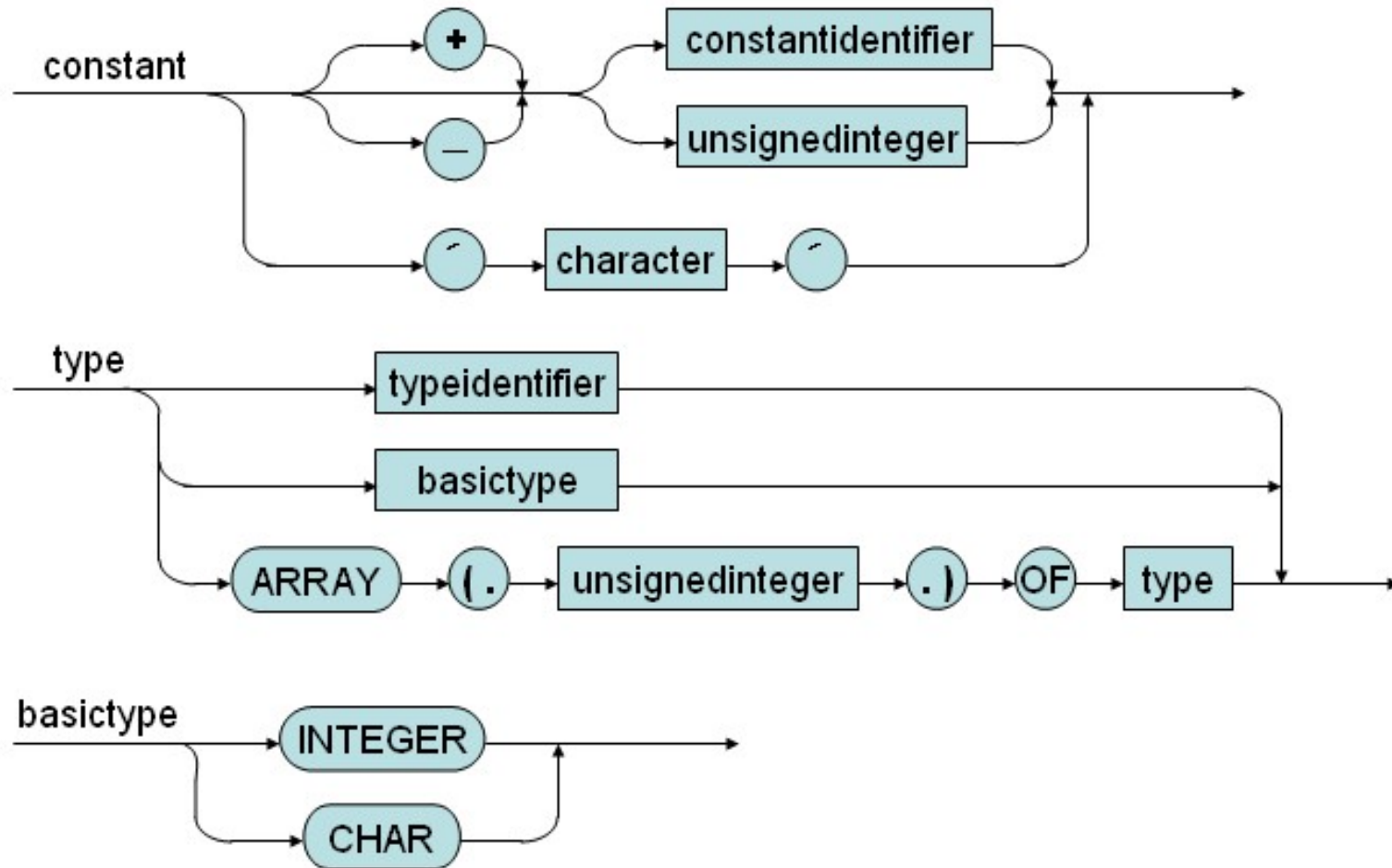
block



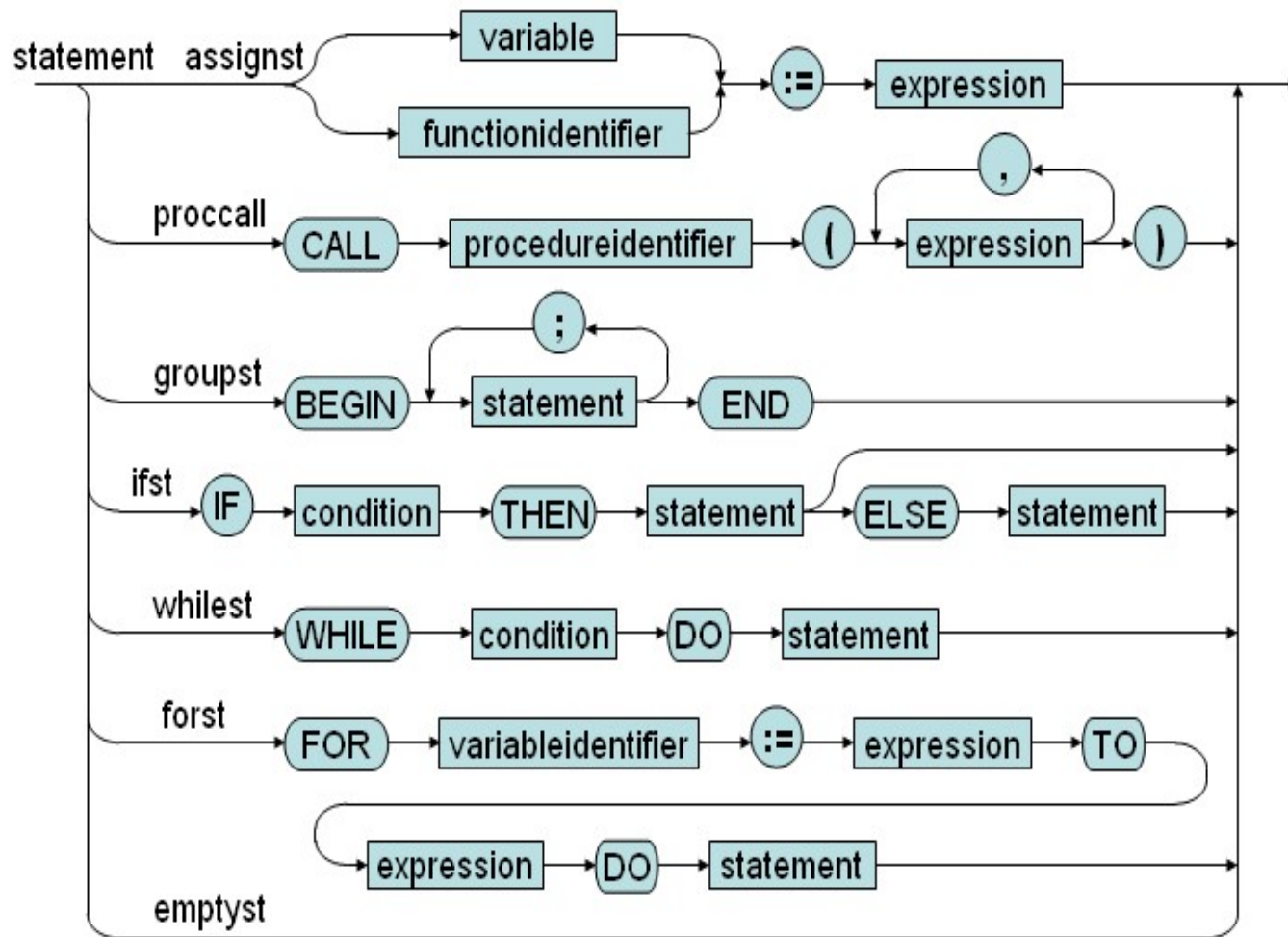
# Syntax diagram of KPL



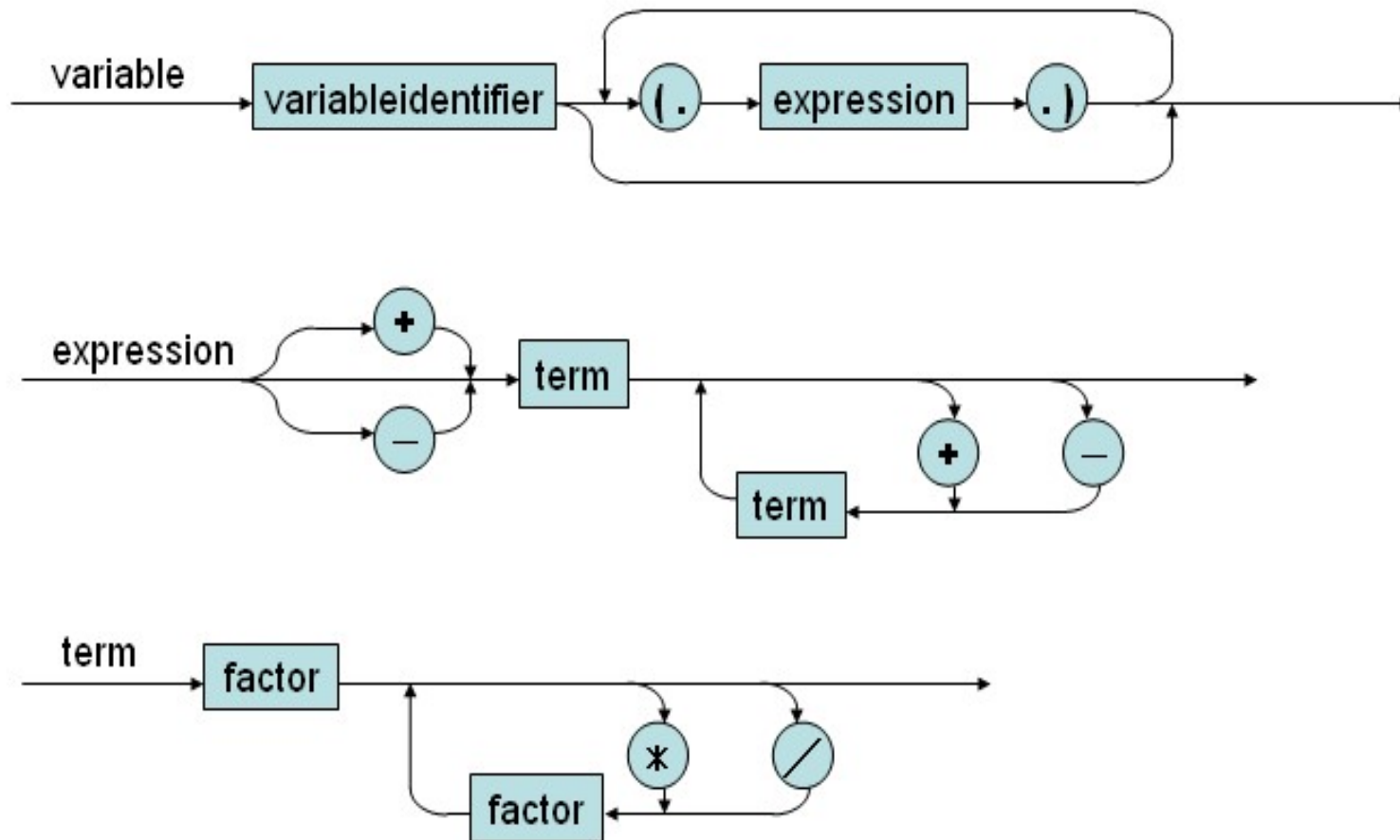
# Syntax diagram of KPL



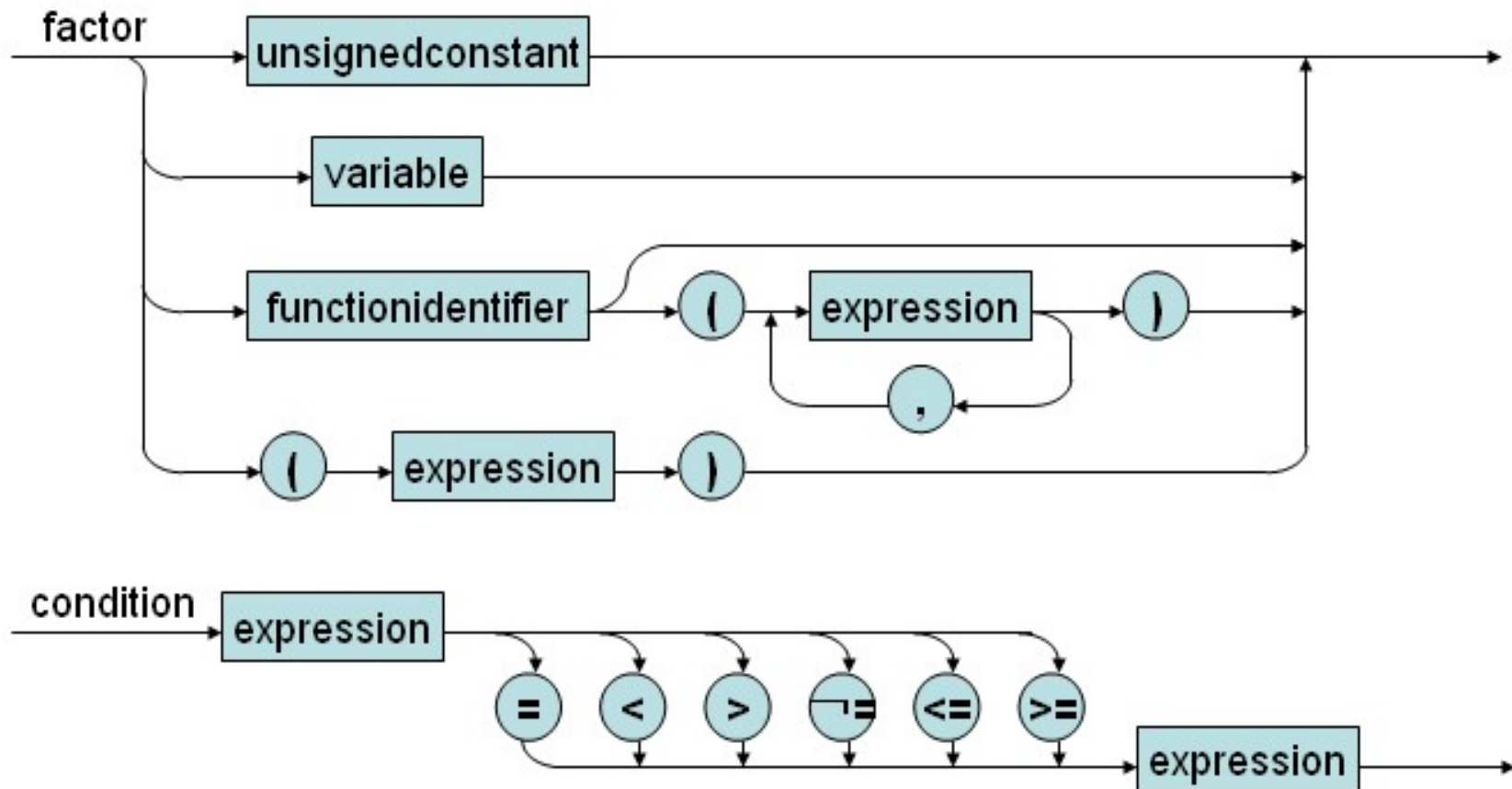
# Syntax diagram of KPL



# Syntax diagram of KPL

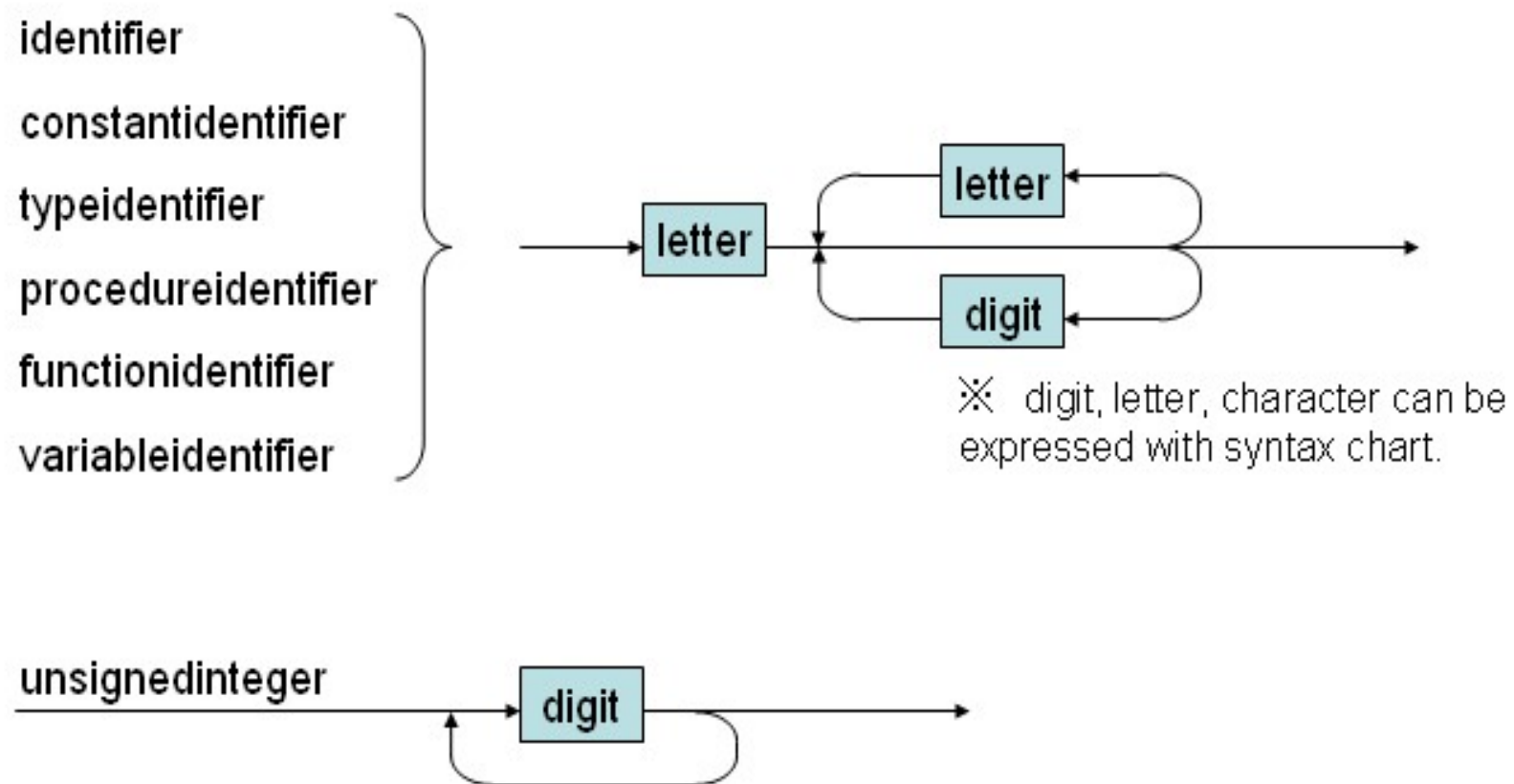


# Syntax diagram of KPL





# Syntax diagram of KPL



# 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`

# KPL Grammar in BNF

- 11)  $\langle \text{ConstDecls} \rangle ::= \langle \text{ConstDecl} \rangle \langle \text{ConstDecls} \rangle$
- 12)  $\langle \text{ConstDecls} \rangle ::= \varepsilon$
- 13)  $\langle \text{ConstDecl} \rangle ::= \text{TK\_IDENT SB\_EQUAL} \langle \text{Constant} \rangle \text{SB\_SEMICOLON}$
- 14)  $\langle \text{TypeDecls} \rangle ::= \langle \text{TypeDecl} \rangle \langle \text{TypeDecls} \rangle$
- 15)  $\langle \text{TypeDecls} \rangle ::= \varepsilon$
- 16)  $\langle \text{TypeDecl} \rangle ::= \text{TK\_IDENT SB\_EQUAL} \langle \text{Type} \rangle \text{SB\_SEMICOLON}$
- 17)  $\langle \text{VarDecls} \rangle ::= \langle \text{VarDecl} \rangle \langle \text{VarDecls} \rangle$
- 18)  $\langle \text{VarDecls} \rangle ::= \varepsilon$
- 19)  $\langle \text{VarDecl} \rangle ::= \text{TK\_IDENT SB\_COLON} \langle \text{Type} \rangle \text{SB\_SEMICOLON}$
- 20)  $\langle \text{SubDecls} \rangle ::= \langle \text{FunDecl} \rangle \langle \text{SubDecls} \rangle$
- 21)  $\langle \text{SubDecls} \rangle ::= \langle \text{ProcDecl} \rangle \langle \text{SubDecls} \rangle$
- 22)  $\langle \text{SubDecls} \rangle ::= \varepsilon$

# KPL Grammar in BNF

```
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)  $\langle \text{Params} \rangle ::= \text{SB\_LPAR } \langle \text{Param} \rangle \langle \text{Params2} \rangle \text{SB\_RPAR}$

26)  $\langle \text{Params} \rangle ::= \varepsilon$

27)  $\langle \text{Params2} \rangle ::= \text{SB SEMICOLON } \langle \text{Param} \rangle \langle \text{Params2} \rangle$

28)  $\langle \text{Params2} \rangle ::= \varepsilon$

29) **<Param>** ::= TK IDENT SB COLON **<BasicType>**

30) **<Param>** ::= KW VAR TK IDENT SB COLON **<BasicType>**

# KPL Grammar in BNF

- 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>`
  
- 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`

# KPL Grammar in BNF

46)  $\langle \text{Statements} \rangle ::= \langle \text{Statement} \rangle \langle \text{Statements2} \rangle$

47)  $\langle \text{Statements2} \rangle ::= \text{SB\_SEMICOLON} \langle \text{Statement} \rangle \langle \text{Statements2} \rangle$

48)  $\langle \text{Statements2} \rangle ::= \varepsilon$

49)  $\langle \text{Statement} \rangle ::= \langle \text{AssignSt} \rangle$

50)  $\langle \text{Statement} \rangle ::= \langle \text{CallSt} \rangle$

51)  $\langle \text{Statement} \rangle ::= \langle \text{GroupSt} \rangle$

52)  $\langle \text{Statement} \rangle ::= \langle \text{IfSt} \rangle$

53)  $\langle \text{Statement} \rangle ::= \langle \text{WhileSt} \rangle$

54)  $\langle \text{Statement} \rangle ::= \langle \text{ForSt} \rangle$

55)  $\langle \text{Statement} \rangle ::= \varepsilon$

# KPL Grammar in BNF

- 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> ::=  $\epsilon$`
  
- 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)  $\langle \text{Arguments2} \rangle ::= \varepsilon$
  
- 68)  $\langle \text{Condition} \rangle ::= \langle \text{Expression} \rangle \langle \text{Condition2} \rangle$
  
- 69)  $\langle \text{Condition2} \rangle ::= \text{SB\_EQ } \langle \text{Expression} \rangle$
- 70)  $\langle \text{Condition2} \rangle ::= \text{SB\_NEQ } \langle \text{Expression} \rangle$
- 71)  $\langle \text{Condition2} \rangle ::= \text{SB\_LE } \langle \text{Expression} \rangle$
- 72)  $\langle \text{Condition2} \rangle ::= \text{SB\_LT } \langle \text{Expression} \rangle$
- 73)  $\langle \text{Condition2} \rangle ::= \text{SB\_GE } \langle \text{Expression} \rangle$
- 74)  $\langle \text{Condition2} \rangle ::= \text{SB\_GT } \langle \text{Expression} \rangle$

# KPL Grammar in BNF

- 75)  $\langle \text{Expression} \rangle ::= \text{SB\_PLUS } \langle \text{Expression2} \rangle$
- 76)  $\langle \text{Expression} \rangle ::= \text{SB\_MINUS } \langle \text{Expression2} \rangle$
- 77)  $\langle \text{Expression} \rangle ::= \langle \text{Expression2} \rangle$
  
- 78)  $\langle \text{Expression2} \rangle ::= \langle \text{Term} \rangle \langle \text{Expression3} \rangle$
  
- 79)  $\langle \text{Expression3} \rangle ::= \text{SB\_PLUS } \langle \text{Term} \rangle \langle \text{Expression3} \rangle$
- 80)  $\langle \text{Expression3} \rangle ::= \text{SB\_MINUS } \langle \text{Term} \rangle \langle \text{Expression3} \rangle$
- 81)  $\langle \text{Expression3} \rangle ::= \varepsilon$
  
- 82)  $\langle \text{Term} \rangle ::= \langle \text{Factor} \rangle \langle \text{Term2} \rangle$
  
- 83)  $\langle \text{Term2} \rangle ::= \text{SB\_TIMES } \langle \text{Factor} \rangle \langle \text{Term2} \rangle$
- 84)  $\langle \text{Term2} \rangle ::= \text{SB\_SLASH } \langle \text{Factor} \rangle \langle \text{Term2} \rangle$
- 85)  $\langle \text{Term2} \rangle ::= \varepsilon$
- 86)  $\langle \text{Factor} \rangle ::= \langle \text{UnsignedConstant} \rangle$
- 87)  $\langle \text{Factor} \rangle ::= \langle \text{Variable} \rangle$
- 88)  $\langle \text{Factor} \rangle ::= \langle \text{FunctionApptication} \rangle$
- 89)  $\langle \text{Factor} \rangle ::= \text{SB\_LPAR } \langle \text{Expression} \rangle \text{SB\_RPAR}$

# KPL Grammar in BNF

90)  $\langle \text{Variable} \rangle ::= \text{TK\_IDENT } \langle \text{Indexes} \rangle$

91)  $\langle \text{FunctionApplication} \rangle ::= \text{TK\_IDENT } \langle \text{Arguments} \rangle$

92)  $\langle \text{Indexes} \rangle ::= \text{SB\_LSEL } \langle \text{Expression} \rangle \text{ SB\_RSEL } \langle \text{Indexes} \rangle$

93)  $\langle \text{Indexes} \rangle ::= \varepsilon$

- **Input: Use functions**

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

Example

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

- **Output: Use procedures**

- Writel: Print an integer. 1 parameter
- WriteC: Print a character. 1 parameter
- WriteLn: Print the newline character.

Example

```
call Writel(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(sum);  
  END. (* example5*)
```

# Parser 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}

$\text{FIRST(RHS1)} \cap \text{FIRST(RHS2)} = \emptyset$

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

# lookAhead token

- Look ahead the next token

```
Token *currentToken;  
Token *lookAhead;  
  
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 the 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**

1)  $\langle \text{Prog} \rangle ::= \text{KW\_PROGRAM TK\_IDENT SB\_SEMICOLON } \langle \text{Block} \rangle \text{ SB\_PERIOD}$

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

# Parsing statements

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> ::= <Whilst?
KW_FOR	54) <Statement> ::= <ForSt?
-----	
SB_SEMICOLON	55) $\epsilon$
KW_END	55) $\epsilon$
KW_ELSE	55) $\epsilon$
-----	
Others	Error

# Parsing a 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

35) `<BasicType> ::= KW_INTEGER`

36) `<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: Recursion OK, you should process the FOLLOW SET

10)  $\langle \text{ConstDecls} \rangle ::= \langle \text{ConstDecl} \rangle \langle \text{ConstDecls} \rangle$

11)  $\langle \text{ConstDecls} \rangle ::= \varepsilon$

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

# Sometimes you should refer to syntax diagrams

## Syntax of Term (using BNF)

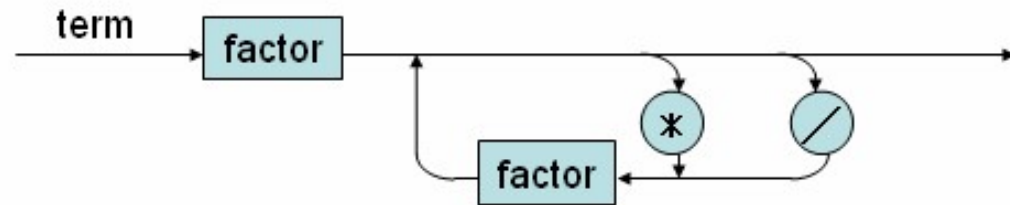
82)  $\langle \text{Term} \rangle ::= \langle \text{Factor} \rangle \langle \text{Term2} \rangle$

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

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

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

## Syntax of Term (using Syntax Diagram)





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

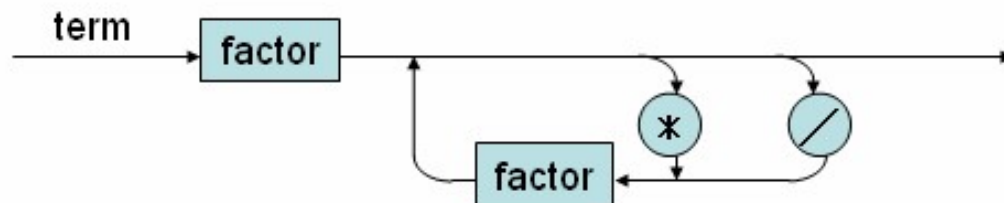
```
void compileTerm(void)
{ compileFactor();
  compileTerm2();
}

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:
```

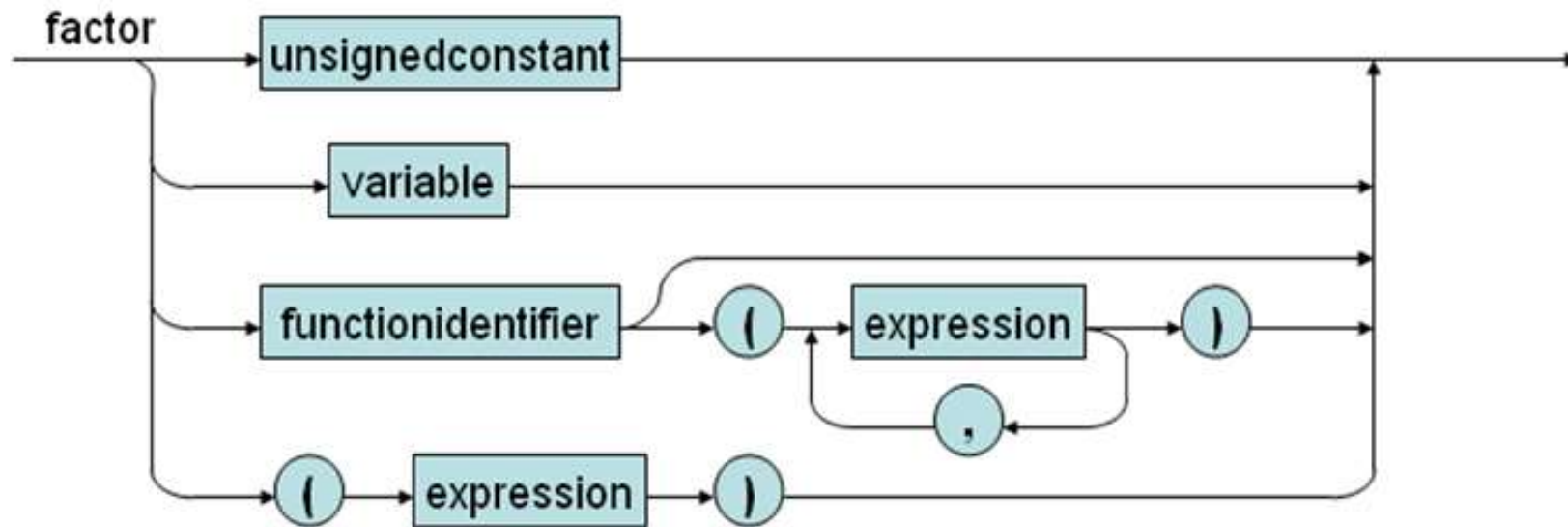
```
  case SB_RPAR:
    case SB_COMMA:
    case SB_EQ:
    case SB_NEQ:
    case SB_LE:
    case SB_LT:
    case SB_GE:
    case SB_GT:
    case SB_RSEL:
    case SB_SEMICOLON:
    case KW_END:
    case KW_ELSE:
    case KW_THEN:
      break;
    default:
      error(ERR_INVALIDTERM,
        lookAhead->lineNo,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



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

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

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

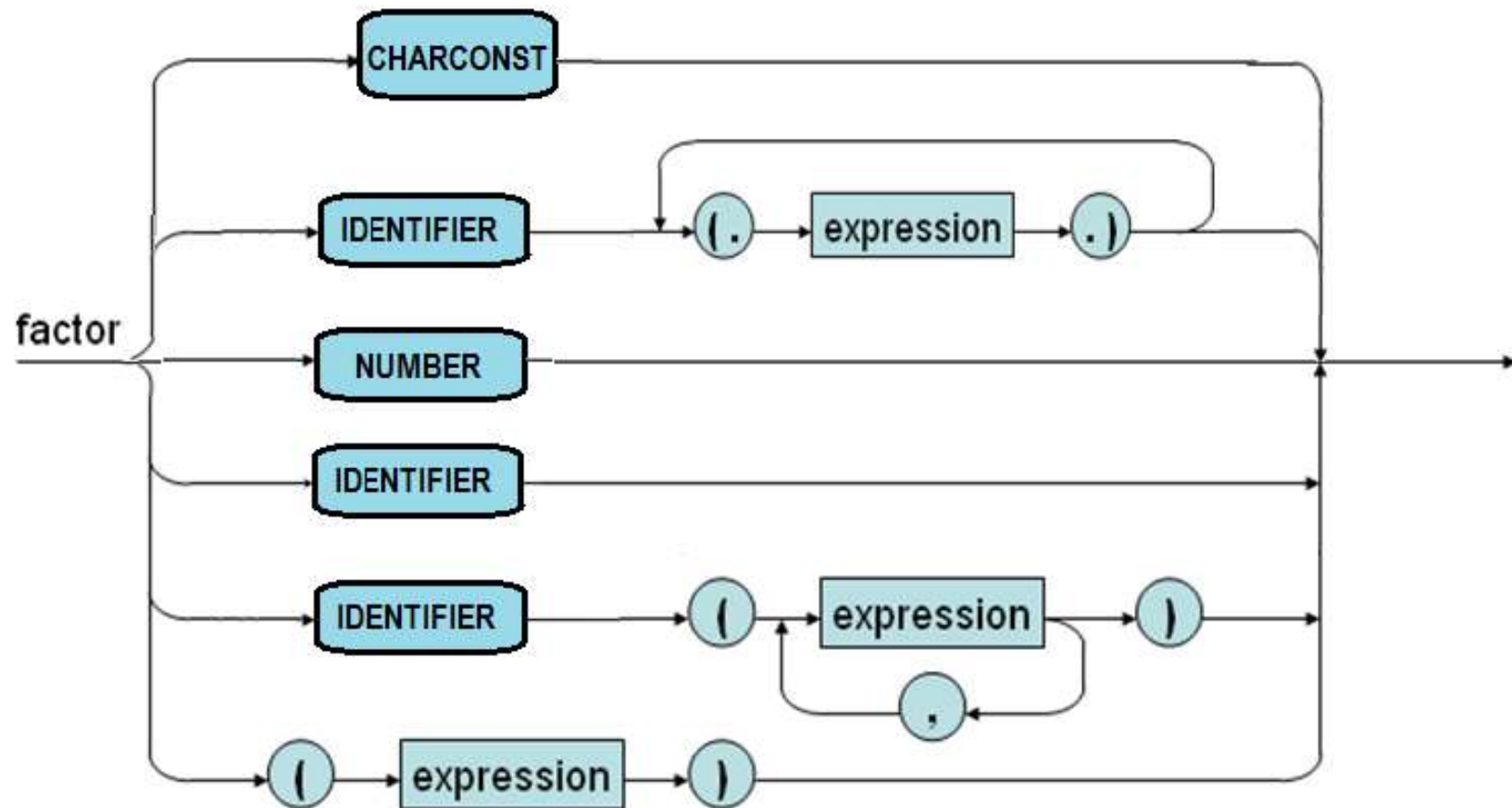
$\text{FIRST}(\text{unsignedconstant}) \cap \text{FIRST}(\text{functioncall}) = \{\text{TK\_IDENT}\}$

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

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

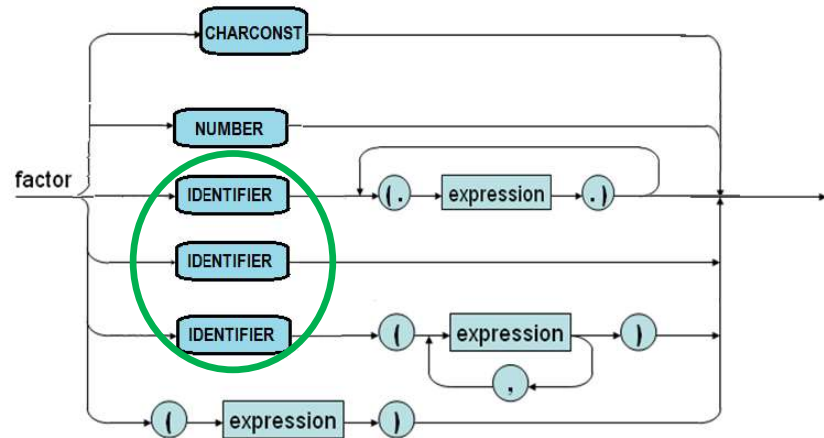
=> violation of LL(1) condition

# After separating and merging



# Compile a factor

```
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;  
}
```



```
case SB_LPAR:  
    eat(SB_LPAR);  
    compileExpression();  
    eat(SB_RPAR);  
    break;  
default:  
    error(ERR_INVALIDFACTOR,  
        lookAhead->lineNo, lookAhead->colNo);  
}
```