Chapter 2 A Simple Compiler

Outlines

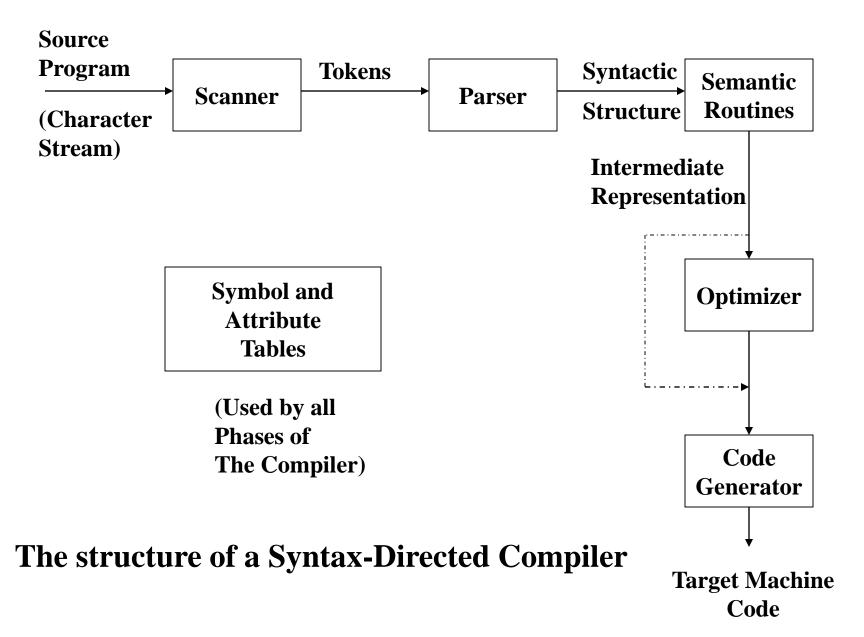
- 2.1 The Structure of a Micro Compiler
- 2.2 A Micro Scanner
- 2.3 The Syntax of Micro
- 2.4 Recursive Descent Parsing
- 2.5 Translating Micro

Micro

- Micro: a very simple language
 - Only integers
 - No declarations
 - Variables consist of A..Z, 0..9, and at most 32 characters long.
 - Comments begin with -- and end with end-of-line.
 - Three kinds of stmts:
 - assignments, e.g., a := b + c
 - read(list of ids), e.g., read(a, b)
 - write(list of exps), e.g., write(a+b)
 - Begin, end, read, and write are reserved words.
 - Tokens may not extend to the following line.

The Structure of a Micro compiler

- One-pass type, no explicit intermediate representations used
 - See P. 9, Fig. 1.3
- The interface
 - Parser is the main routine.
 - Parser calls scanner to get the next token.
 - Parser calls semantic routines at appropriate times.
 - Semantic routines produce output in assembly language.
 - A simple symbol table is used by the semantic routines.



A Micro Scanner

- The Micro Scanner will be a function of no arguments that returns *token* values
 - There are 14 tokens.

```
typedef enum token_types {
     BEGIN, END, READ, WRITE, ID, INTLITERAL,
     LPAREN, RPAREN, SEMICOLON, COMMA, ASSIGNOP,
     PLUOP, MINUSOP, SCANEOF
} token;

Extern token scanner(void);
```

A Micro Scanner (Cont'd)

 The scanner returns the longest string that constitutes a token, e.g., in

abcdef

ab, abc, abcdef are all valid tokens. The scanner will return the longest one (i.e., abcdef).

```
#include <stdio.h>
#include <ctype.h>
int in char, c;
while ((in char = getchar()) != EOF) {
    if (isspace(in char))
        continue; /* do nothing */
    else if (isalpha(in char)) {
        /*
         * ID ::= LETTER | ID LETTER
                           ID DIGIT
         *
                           ID UNDERSCORE
         *
         */
        for (c = getchar(); isalnum(c) || c == '_';
                     c = getchar())
        ungetc(c, stdin);
                                                    Continue: skip one iteration
        return ID;
                                                    getchar() , isspace(),
    } else if (isdigit(in_char)) {
        /*
                                                    isalpha(), isalnum(),
         * INTLITERAL ::= DIGIT
                           INTLITERAL DIGIT
                                                    isdigital()
         */
        while (isdigit((c = getchar())))
                                                    ungetc(): push back one character
        ungetc(c, stdin);
        return INTLITERAL;
    } else
        lexical error(in char);
}
```

Figure 2.1 Scanner Loop to Recognize Identifiers and Integer Literals

```
#include <stdio.h>
#include <ctype.h>
int in_char, c;
while ((in char = getchar()) != EOF) {
    if (isspace(in char))
        /* do nothing */
        continue;
    else if (isalpha(in char))
        /* code to recognize identifiers goes here */
    else if (isdigit(in char))
        /* code to recognize int literals goes here */
    else if (in char == '(')
        return LPAREN;
    else if (in char == ')')
        return RPAREN;
    else if (in char == ';')
        return SEMICOLON;
    else if (in char == ',')
        return COMMA;
    else if (in char == '+')
        return PLUSOP;
    else if (in char == ':') {
        /* looking for ":=" */
        c = getchar();
        if (c == '=')
            return ASSIGNOP;
        else {
            ungetc(c, stdin);
            lexical_error(in_char);
        }
    } else if (in char == '-') {
        /* looking for --, comment start */
        c = getchar();
        if (c == '-') {
            while ((in char = getchar()) != '\n');
        } else {
            ungetc(c, stdin);
            return MINUSOP;
        }
    } else
        lexical error(in char);
}
```

Figure 2.2 Scanner Loop with New Code to Recognize Operators, Comments, and Delimiters

A Micro Scanner (Cont'd)

- How to handle RESERVED words?
 - Reserved words are similar to identifiers.
- Two approaches:
 - Use a separate table of reserved words
 - Put all reserved words into symbol table initially.

A Micro Scanner (Cont'd)

- Provision for saving the characters of a token as they are scanned
 - token_buffer, buffer_char(), clear_buffer(),
 check_reserved()
- Handle end of file
 - feof(stdin)

```
#include <stdio.h>
/* character classification macros */
#include <ctype.h>
extern char token buffer[];
                                   Complete Scanner Function
token scanner (void)
{
                                   for Micro
    int in char, c;
    clear buffer();
    if (feof(stdin))
        return SCANEOF;
    while ((in_char = getchar()) != EOF) {
        if (isspace(in char))
            continue; /* do nothing */
        else if (isalpha(in_char)) {
            /*
             * ID ::= LETTER | ID LETTER
                               ID DIGIT
                               ID UNDERSCORE
             */
            buffer char(in char);
            for (c = getchar(); isalnum(c) | c == '_';
                    c = getchar())
                buffer_char(c);
```

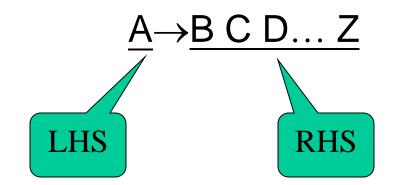
```
ungetc(c, stdin);
       return check reserved();
   } else if (isdigit(in_char)) {
        /*
         * INTLITERAL ::= DIGIT |
                          INTLITERAL DIGIT
         */
       buffer char(in char);
       for (c = getchar(); isdigit(c);
                c = getchar())
            buffer char(c);
       ungetc(c, stdin);
        return INTLITERAL;
   } else if (in char == '(')
        return LPAREN;
   else if (in char == ')')
        return RPAREN;
   else if (in char == ';')
        return SEMICOLON;
   else if (in char == ',')
        return COMMA;
   else if (in char == '+')
        return PLUSOP;
    else if (in char == ':') {
        /* looking for ":=" */
        c = getchar();
        if (c == '=')
            return ASSIGNOP;
        else {
            ungetc(c, stdin);
            lexical error(in_char);
    } else if (in char == '-') {
        /* is it --, comment start */
        c = getchar();
        if (c == '-') {
            do
                in char = getchar();
            while (in char != '\n');
        } else {
            ungetc(c, stdin);
            return MINUSOP;
        }
    } else
        lexical error(in char);
}
```

Figure 2.3 Complete Scanner Function for Micro

}

The Syntax of Micro

- Micro's syntax is defined by a context-free grammar (CFG)
 - CFG is also called BNF (Backus-Naur Form) grammar
- CFG consists of a set of production rules,



LHS must be a single nonterminal

RHS consists 0 or more terminals or nonterminals

- Two kinds of symbols
 - Nonterminals
 - Delimited by < and >
 - Represent syntactic structures
 - Terminals
 - Represent tokens
- E.g.
- Start or goal symbol
- λ : empty or null string

• E.g.

```
<statement list> \rightarrow <statement><statement tail> <statement tail> \rightarrow \lambda <statement tail> \rightarrow <statement tail>
```

Extended BNF: some abbreviations

```
1. optional: [] 0 or 1
        \langle stmt \rangle \rightarrow if \langle exp \rangle then \langle stmt \rangle
         <stmt> → if <exp> then <stmt> else <stmt>
        can be written as
        \langle \text{stmt} \rangle \rightarrow \text{if } \langle \text{exp} \rangle \text{ then } \langle \text{stmt} \rangle \text{ [else } \langle \text{stmt} \rangle \text{]}
2. repetition: {} 0 or more
        <stmt list> \rightarrow <stmt> <tail>
        <tail> \rightarrow \lambda
        \langle tail \rangle \rightarrow \langle stmt \rangle \langle tail \rangle
        can be written as
        \langle \text{stmt list} \rangle \rightarrow \langle \text{stmt} \rangle
```

Extended BNF: some abbreviations

- Extended BNF == BNF
 - Either can be transformed to the other.
 - Extended BNF is more compact and readable

```
→ begin <statement list> end
     cprogram>
                      → <statement> {<statement>}
     <statement list>
                      \rightarrow ID := <expression> ;
3.
   <statement>
                      \rightarrow read ( <id list> );
4.
   <statement>
                      → write ( <expr list> );
5. <statement>
                      \rightarrow ID \{, ID\}
6. <id list>
7. <expr list>
                      → <expression> {, <expression>}
                      → <primary> {<add op> <primary>}
8.
     <expression>
                      → ( <expression> )
     cprimary>
9.
     primary>
10.
                      \rightarrow ID
     cprimary>
                      → INTLITERAL
11.
     <add op>
                      → PLUSOP
12.
                      → MINUSOP
     <add op>
13.
                      <system goal>
14.
```

Figure 2.4 Extended CFG Defining Micro

• The derivation of

begin ID:= ID + (INTLITERAL – ID); end

```
cprogram>
begin <statement list> end
                                                            (Apply rule 1)
                                                            (Apply rule 2)
begin <statement> {<statement>} end
                                                            (Choose 0 repetitions)
begin <statement> end
                                                            (Apply rule 3)
begin ID := <expression> ; end
                                                            (Apply rule 8)
begin ID := <primary> {<add op> <primary>} ; end
begin ID := <primary> <add op> <primary> ; end
                                                            (Choose 1 repetition)
begin ID := <primary> + <primary> ; end
                                                            (Apply rule 12)
                                                            (Apply rule 10)
begin ID := ID + <pri>mary> ; end
                                                            (Apply rule 9)
begin ID := ID + ( <expression> ); end
                                                            (Apply rule 8)
begin ID := ID + ( <primary> {<add op> <primary>} ) ; end
begin ID := ID + ( <primary> <add op> <primary> ) ; end
                                                            (Choose 1 repetition)
begin ID := ID + ( <primary> - <primary> ); end
                                                            (Apply rule 13)
begin ID := ID + ( INTLITERAL - <primary> ) ; end
                                                            (Apply rule 11)
                                                            (Apply rule 10)
begin ID := ID + ( INTLITERAL - ID ); end
```

- A CFG defines a *language*, which is a set of sequences of **tokens**
- Syntax errors & semantic errors

Associativity

Operator precedence

$$A+B*C$$

• A grammar fragment defines such a precedence relationship

```
<expression> → <factor> {<add op> <factor>}
<factor> → <primary> {<mult op> <primary>}
<primary> → (<expression>)
<primary> → ID
<primary> → INTLITERAL
```

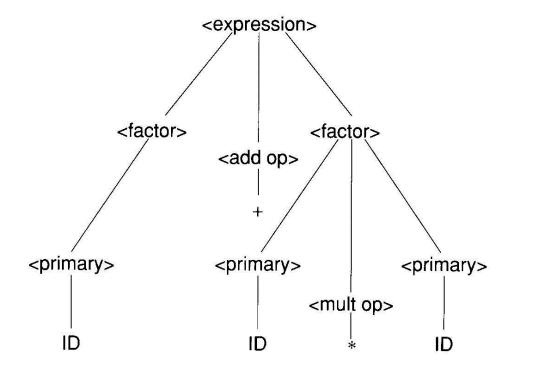


Figure 2.5 Derivation Tree for A+B*C

• With parentheses, the desired grouping can be forced

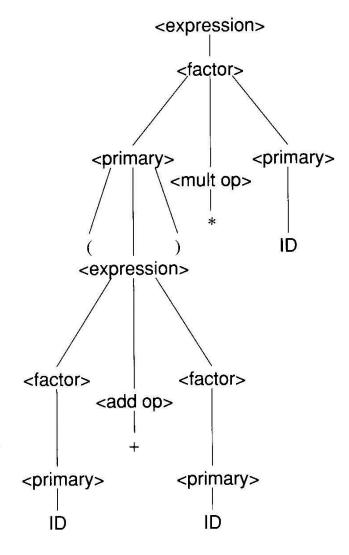


Figure 2.6 Derivation Tree for (A+B)*C

Recursive Descent Parsing

- There are many parsing techniques.
 - Recursive descent is one of the simplest parsing techniques
- Basic idea
 - Each nonterminal has a parsing procedure
 - For symbol on the RHS: a sequence of matching
 - To Match a nonterminal A
 - Call the parsing procedure of A
 - To match a terminal symbol t
 - Call match(t)
 - » match(t) calls the scanner to get the next token. If it is, everything is correct. If it is not t, we have found a syntax error.

Recursive Descent Parsing

- If a nonterminal has several productions, choose an appropriate one based on the next input token.
- Parser is started by invoking system_goal().

```
void system goal(void)
    /* <system goal> ::= cprogram> SCANEOF */
    program();
    match (SCANEOF);
void program(void)
    /* cprogram> ::= BEGIN <statement list> END */
    match (BEGIN);
    statement list();
   match (END);
```

```
void statement_list(void)
    /*
     * <statement list> ::= <statement>
                            { <statement> }
    statement();
   while (TRUE)
        switch (next_token()) {
        case ID:
        case READ:
        case WRITE:
                                  處理{<statement>}不出
            statement();
                                         現的情形
            break:
        default:
            return;
```

next_token(): a function that returns the next token.
 It does not call scanner(void).

```
void statement(void)
    token tok = next_token();
    switch (tok) {
    case ID:
        /* <statement> ::= ID := <expression> ; */
        match(ID); match(ASSIGNOP);
        expression(); match(SEMICOLON);
        break;
    case READ:
        /* <statement> ::= READ ( <id list> ) ; */
        match (READ); match (LPAREN);
        id list(); match(RPAREN);
        match (SEMICOLON);
        break;
    case WRITE:
        /* <statement> ::= WRITE ( <expr list> ) ; */
        match(WRITE); match(LPAREN);
        expr_list(); match(RPAREN);
        match (SEMICOLON);
        break;
                                              <statement>
     default:
         syntax error (tok);-
                                              必出現一次
         break;
```

```
void id list(void)
    /* <id list> ::= ID { , ID } */
    match (ID);
    while (next token() == COMMA) {
        match (COMMA);
        match (ID);
void expression(void)
    token t;
    /*
     * <expression> ::= <primary>
     *
                         { <add op> <primary> }
     */
    primary();
    for (t = next token(); t == PLUSOP || t == MINUSOP;
                  t = next token()) {
        add op();
        primary();
```

```
void expr list(void)
 /* <expr list> ::= <expression> { , <expression> } */
    expression();
    while (next token() == COMMA) {
        match (COMMA);
        expression();
void add op(void)
    token tok = next token();
    /* <addop> ::= PLUSOP | MINUSOP */
    if (tok == PLUSOP | tok == MINUSOP)
        match (tok);
    else
        syntax error(tok);
```

```
void primary (void)
   token tok = next_token();
    switch (tok) {
    case LPAREN:
        /* /* con> ) */
        match(LPAREN); expression();
        match (RPAREN);
        break;
    case ID:
        /* <primary> ::= ID */
        match (ID);
        break;
    case INTLITERAL:
        /* /* primary> ::= INTLITERAL */
        match(INTLITERAL);
        break;
    default:
        syntax_error(tok);
        break;
```

Figure 2.7 Remaining Parsing Procedures for Micro

Translating Micro

- Target language: 3-addr code (quadruple)
 - OP A, B, C
- Note that we did not worry about registers at this time.
 - temporaries: Sometimes we need to hold temporary values.
 - E.g. A+B+C
 ADD A,B,TEMP&1
 ADD TEMP&1,C,TEMP&2

Translating Micro (Cont'd)

- Action Symbols
 - The bulk of a translation is done by semantic routine
 - Action symbols can be added to a grammar to specify when semantic processing should take place
 - Be placed anywhere in the RHS of a production
 - translated into procedure call in the parsing procedures
 - #add corresponds to a semantic routine named add()
 - No impact on the languages recognized by a parser driven by a CFG

```
→ #start begin <statement list> end
cprogram>
                  → <statement> {<statement>}
<statement list>
                  → <ident> := <expression> #assign ;
<statement>
                  \rightarrow read ( <id list> );
<statement>
                  \rightarrow write ( <expr list> );
<statement>
                  → <ident> #read id {, <ident> #read id }
<id list>
                  → <expression> #write_expr
<expr list>
                       {, <expression> #write expr}
                  → <primary>
<expression>
                       {<add op> <primary> #gen_infix}
primary>
                  \rightarrow ( <expression> )
                  → <ident>
cprimary>
<primary>
                  → INTLITERAL #process_literal
<add op>
                  → PLUSOP #process op
                 → MINUSOP #process_op
<add op>
                  → ID #process_id
<ident>
                  → program> SCANEOF #finish
<system goal>
```

Figure 2.9 Grammar for Micro with Action Symbols

Semantic Information

- Semantic routines need certain information to do their work.
 - These information is stored in *semantic records*.
 - Each kind of grammar symbol has a semantic record

```
#define MAXIDLEN 33
typedef char string[MAXIDLEN];

typedef struct operator { /* for operators */
    enum op { PLUS, MINUS } operator;
} op_rec;

/* expression types */
enum expr { IDEXPR, LITERALEXPR, TEMPEXPR };

/* for <primary> and <expression> */
typedef struct expression {
    enum expr kind;
    union {
        string name; /* for IDEXPR, TEMPEXPR */
        int val; /* for LITERALEXPR */
    };
} expr_rec;
```

Figure 2.8 Semantic Records for Micro Grammar Symbols

```
void expression (void)
    token t;
     * <expression> ::= <primary>
                         { <add op> <primary> }
     */
   primary();
    for (t = next token(); t == PLUSOP || t == MINUSOP;
                  t = next token()) {
        add op();
        primary();
}
```

Old parsing procedure P. 37

```
void expression(expr rec *result)
                                            <expression> → <primary>
   expr rec left operand, right operand;
                                                {<add op> <primary> #gen_infix
   op rec op;
   primary(& left operand);
   while (next token() == PLUSOP ||
          next token() == MINUSOP) {
       add op(& op);
                                                    New parsing procedure
       primary(& right operand);
       left operand = gen infix(left operand, op,
                                right operand);
                                                    routines
   *result = left operand;
```

A Parsing Procedure Including Semantic Figure 2.11 Processing

which involved semantic

Semantic Information (Cont'd)

Subroutines for symbol table and temporaries

```
extern void enter(string s);
void check id(string s)
     if (! lookup(s)) {
         enter(s);
         generate("Declare", s, "Integer", "");
char *get temp(void)
   /* max temporary allocated so far */
    static int max temp = 0;
    static char tempname[MAXIDLEN];
   max temp++;
    sprintf(tempname, "Temp&%d", max temp);
    check id(tempname);
    return tempname;
```

/* Put s unconditionally into symbol table. */

/* Is s in the symbol table? */

extern int lookup(string s);

Semantic Information (Cont'd)

Semantic routines

```
void start(void)
    /* Semantic initializations, none needed. */
void finish (void)
    /* Generate code to finish program. */
    generate("Halt", "", "", "");
void assign(expr rec target, expr rec source)
    /* Generate code for assignment. */
    generate("Store", extract(source),
              target.name, "");
```

```
op_rec process_op(void)
    /* Produce operator descriptor. */
    op rec o;
    if (current_token == PLUSOP)
        o.operator = PLUS;
    else
        o.operator = MINUS;
    return o;
expr_rec gen_infix(expr_rec el, op_rec op,
                    expr rec e2)
{
    expr rec e rec;
    /* An expr rec with temp variant set. */
    e rec.kind = TEMPEXPR;
    /*
     * Generate code for infix operation.
     * Get result temp and set up semantic record
     * for result.
     */
    strcpy(erec.name, get_temp());
    generate (extract (op), extract (e1),
              extract(e2), erec.name);
     return erec;
```

```
void read_id(expr_rec in_var)
{
    /* Generate code for read. */
    generate("Read", in_var.name,
             "Integer", "");
}
expr_rec process_id(void)
    expr rec t;
    /*
     * Declare ID and build a
     * corresponding semantic record.
     */
    check id(token buffer);
    t.kind = IDEXPR;
    strcpy(t.name, token buffer);
    return t;
```

```
expr rec process literal(void)
    expr rec t;
    /*
     * Convert literal to a numeric representation
     * and build semantic record.
     */
    t.kind = LITERALEXPR;
    (void) sscanf(token buffer, "%d", & t.val);
    return t;
void write_expr(expr_rec out_expr)
    generate("Write", extract(out expr),
             "Integer", "");
}
```

Figure 2.10 Action Routines for Micro

Step	Parser Action	Remaining Input	Generated Code
(1)	Call system_goal()	begin A:=BB-314+A; end SCANEOF	
(2)	Call program()	begin A:=BB-314+A; end SCANEOF	
(3)	Semantic Action: start()	begin A:=BB-314+A; end SCANEOF	
(4)	match (BEGIN)	begin A:=BB-314+A; end SCANEOF	
(5)	Call statement_list()	A:=BB-314+A; end SCANEOF	
(6)	Call statement()	A:=BB-314+A; end SCANEOF	
(7)	Call ident()	A:=BB-314+A; end SCANEOF	
(8)	match(ID)	A:=BB-314+A; end SCANEOF	
(9)	Semantic Action:	:=BB-314+A; end SCANEOF	Declare A,Integer
	process_id()		
(10)	match (ASSIGNOP)	:=BB-314+A; end SCANEOF	
(11)	Call expression()	BB-314+A; end SCANEOF	
(12)	Call primary()	BB-314+A; end SCANEOF	
(13)	Call ident()	BB-314+A; end SCANEOF	
(14)	match (ID)	BB-314+A; end SCANEOF	
(15)	Semantic Action:	-314+A; end SCANEOF	Declare BB,Integer
	process_id()		
(16)	Call add_op()	-314+A; end SCANEOF	
(17)	match (MINUSOP)	-314+A; end SCANEOF	
(18)	Semantic Action:	314+A; end SCANEOF	
	process_op()		
(19)	Call primary()	314+A; end SCANEOF	
(20)	match (INTLITERAL)	314+A; end SCANEOF	
(21)	Semantic Action:	+A; end SCANEOF	
	<pre>process_literal()</pre>		
(22)	Semantic Action:	+A; end SCANEOF	Declare Temp&1,Integer
	gen_infix()		Sub BB,314,Temp&1
(23)	Call add_op()	+A; end SCANEOF	

(24)	match (PLUSOP)	+A; end SCANEOF	
(25)	Semantic Action:	A; end SCANEOF	
	process_op()		
(26)	Call primary()	A; end SCANEOF	
(27)	Call ident ()	A; end SCANEOF	
(28)	match(ID)	A; end SCANEOF	
(29)	Semantic Action:	; end SCANEOF	
	<pre>process_id()</pre>		
(30)	Semantic Action:	; end SCANEOF	Declare Temp&2,Integer
	gen_infix()		Add Temp&1,A,Temp&2
(31)	Semantic Action: assign()	; end SCANEOF	Store Temp&2,A
(32)	match (SEMICOLON)	; end SCANEOF	
(33)	match (END)	end SCANEOF	
(34)	match (SCANEOF)	SCANEOF	
(35)	Semantic Action: finish()		Halt