

Chapter 2

A Simple Compiler

In this chapter, we will describe a simple programming language called:
adding calculator (ac).

Then, describe a simple compiler for ac.

Definition of ac language

Regular expression specifies Token

- The actual input characters that correspond to each terminal symbol (called token) are specified by regular expression.
- For example:
 - **assign** symbol as a terminal, which appears in the input stream as “=” character.
 - The terminal **id (identifier)** could be any alphabetic character except f, i, or p, which are reserved for special use in ac. It is specified as [a-e] | [g-h] | [j-o] | [q-z]
- Regular expression will be covered in Ch. 3.

Definition of ac language (Cont.)

Terminal	Regular Expression
floatdcl	"f"
intdcl	"i"
print	"p"
id	[a – e] [g – h] [j – o] [q – z]
assign	"="
plus	"+"
minus	"–"
inum	[0 – 9] ⁺
fnum	[0 – 9] ⁺ .[0 – 9] ⁺
blank	(" ") ⁺

Figure 2.3: Formal definition of ac tokens.

Sample AC Program

```
f b // declare a floating variable  
i a // declare an integer variable  
a = 5 // a single character identifier  
b = a + 3.2  
p b // print the variable
```

Adding calculator (ac) grammar

A **context-free grammar** (CFG) specifies the syntax of a language. It is:

**A set of productions
(or rewriting rules).**

The grammar for ac:

```
1 Prog → Dcls Stmt $  
2 Dcls → Dcl Dcls  
3           | λ  
4 Dcl → floatdcl id  
5           | intdcl id  
6 Stmt → Stmt Stmt  
7           | λ  
8 Stmt → id assign Val Expr  
9           | print id  
10 Expr → plus Val Expr  
11           | minus Val Expr  
12           | λ  
13 Val → id  
14           | inum  
15           | fnum
```

Figure 2.1: Original Material by 陳振炎教授
red material by Yung-Pin Cheng

Definition of ac language

- Two kinds of grammar symbols:
terminal (token) and **non-terminal**.
 - 1) A terminal cannot be rewritten.
 - 2) A non-terminal can be rewritten by a production rule.
- A special non-terminal is **start symbol**, which is usually the symbol on the **left-hand side (LHS)** of the grammar's first rule.
- From the start symbol, we proceed by replacing (rewriting) a symbol with the **right-hand side (RHS)** of some production of that symbol.

Definition of ac language (Cont.)

- The symbol λ denotes **empty** or **null string**, which indicates that there are no symbols on a production's RHS.
- The special symbol **\$** represents the end of the input stream.
- To show how the grammar in Fig. 2.1 defines **ac** programs, the derivation of one ac program is given in Fig. 2.2, beginning with:
The start symbol **Prog.**

Step	Sentential Form	Production Number
1	$\langle \text{Prog} \rangle$	
2	$\langle \text{Dcls} \rangle \text{Stmts } \$$	1
3	$\langle \text{Dcl} \rangle \text{Dcls} \text{Stmts } \$$	2
4	$\text{floatdcl id } \langle \text{Dcls} \rangle \text{Stmts } \$$	4
5	$\text{floatdcl id } \langle \text{Dcl} \rangle \text{Dcls} \text{Stmts } \$$	2
6	$\text{floatdcl id } \langle \text{intdcl id } \rangle \langle \text{Dcls} \rangle \text{Stmts } \$$	5
7	$\text{floatdcl id } \text{intdcl id } \langle \text{Stmts} \rangle \$$	3
8	$\text{floatdcl id } \text{intdcl id } \langle \text{Stmt} \rangle \text{Stmts } \$$	6
9	$\text{floatdcl id } \text{intdcl id } \langle \text{id assign Val Expr} \rangle \text{Stmts } \$$	8
10	$\text{floatdcl id } \text{intdcl id } \text{id assign inum } \langle \text{Expr} \rangle \text{Stmts } \$$	14
11	$\text{floatdcl id } \text{intdcl id } \text{id assign inum } \langle \text{Stmts} \rangle \$$	12
12	$\text{floatdcl id } \text{intdcl id } \text{id assign inum } \langle \text{Stmt} \rangle \text{Stmts } \$$	6
13	$\text{floatdcl id } \text{intdcl id } \text{id assign inum } \langle \text{id assign Val Expr} \rangle \text{Stmts } \$$	8
14	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id } \langle \text{Expr} \rangle \text{Stmts } \$$	13
15	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id } \langle \text{plus Val Expr} \rangle \text{Stmts } \$$	10
16	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id plus fnum } \langle \text{Expr} \rangle \text{Stmts } \$$	15
17	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id plus fnum } \langle \text{Stmts} \rangle \$$	12
18	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id plus fnum } \langle \text{Stmt} \rangle \text{Stmts } \$$	6
19	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id plus fnum } \langle \text{print id } \rangle \langle \text{Stmts} \rangle \$$	9
20	$\text{floatdcl id } \text{intdcl id } \text{id assign inum id assign id plus fnum print id } \$$	7

1 $\text{Prog} \rightarrow \text{Dcls Stmt } \$$
 2 $\text{Dcls} \rightarrow \text{Dcl Dcls}$
 3 $\text{Dcl} \rightarrow \text{floatdcl id}$
 4 $\text{Dcl} \rightarrow \text{intdcl id}$
 5 $\text{Stmts} \rightarrow \text{Stmt Stmt }$
 6 $\text{Stmt} \rightarrow \text{id assign Val Expr}$
 7 $\text{Stmt} \rightarrow \text{print id}$
 8 $\text{Expr} \rightarrow \text{plus Val Expr}$
 9 $\text{Expr} \rightarrow \text{minus Val Expr}$
 10 $\text{Val} \rightarrow \text{id}$
 11 $\text{Val} \rightarrow \text{inum}$
 12 $\text{Val} \rightarrow \text{fnum}$

Figure 2.2: Derivation of an ac program using the grammar in Figure 2.1.

Definition of ac language (Cont.)

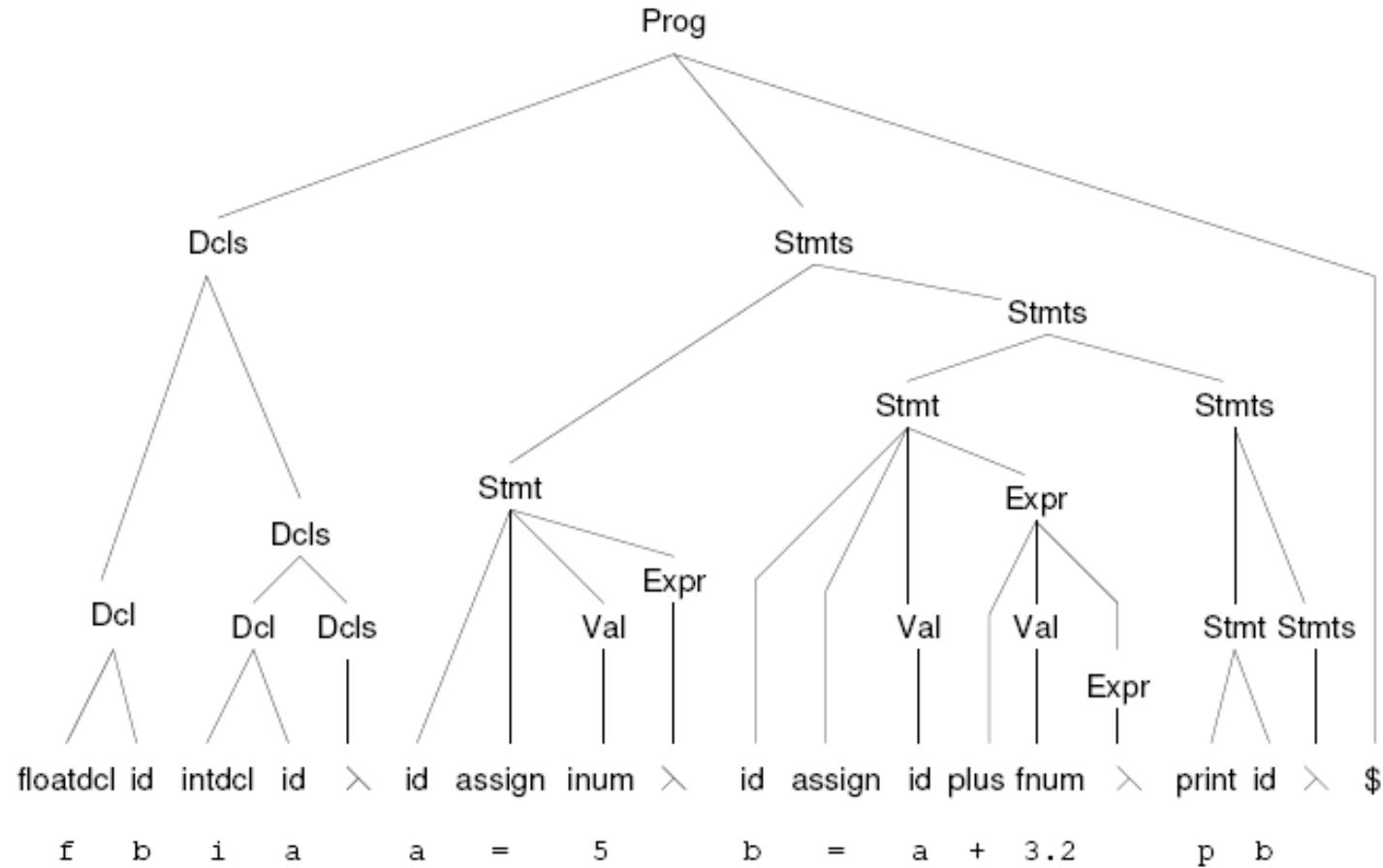


Figure 2.4: An ac program and its parse tree.

Phases of an ac compiler

- Scanning
 - The **scanner** reads a source **ac** program as a text file and produces a stream of tokens.
 - Fig. 2.5 shows a scanner that finds all tokens for ac.
 - Fig. 2.6 shows scanning a number token.
 - Each token has the two components:
 - 1)**Token type** explains the token's category. (e.g., id)
 - 2)**Token value** provides the string value of the token. (e.g., “b”)

```

function SCANNER( ) returns Token
    while s.PEEK( ) = blank do call s.ADVANCE( )
    if s.EOF( )
    then ans.type  $\leftarrow$  $
    else
        if s.PEEK( )  $\in \{0, 1, \dots, 9\}$ 
        then ans  $\leftarrow$  SCANDIGITS( )
        else
            ch  $\leftarrow$  s.ADVANCE( )
            switch (ch)
                case { a, b, ..., z } - { i, f, p }
                    ans.type  $\leftarrow$  id
                    ans.val  $\leftarrow$  ch
                case f
                    ans.type  $\leftarrow$  floatdcl
                case i
                    ans.type  $\leftarrow$  intdcl
                case p
                    ans.type  $\leftarrow$  print
                case =
                    ans.type  $\leftarrow$  assign
                case +
                    ans.type  $\leftarrow$  plus
                case -
                    ans.type  $\leftarrow$  minus
                case default
                    call LEXICALERROR( )
    return (ans)
end

```

Figure 2.5: Scanner for the ac language. The variable *s* is an input stream of characters.

```

function SCANDIGITS( ) returns token
    tok.val  $\leftarrow$  ""
    while s.PEEK( )  $\in \{0, 1, \dots, 9\}$  do
        tok.val  $\leftarrow$  tok.val + s.ADVANCE( )
        if s.PEEK( )  $\neq$  "."
            then tok.type  $\leftarrow$  inum
            else
                tok.type  $\leftarrow$  fnum
                tok.val  $\leftarrow$  tok.val + s.ADVANCE( )
                while s.PEEK( )  $\in \{0, 1, \dots, 9\}$  do
                    tok.val  $\leftarrow$  tok.val + s.ADVANCE( )
    return (tok)
end

```

Figure 2.6: Finding inum or fnum tokens for the ac language.

Phases of an ac compiler (Cont.)

- Parsing
 - **Parser** processes tokens produced by the scanner, determines the syntactic validity of the token stream, and creates an **abstract syntax tree (AST)** for subsequent phases.
 - **Recursive descent (LL type parser)** is one simple parsing technique used in practical compilers.
The parsing procedure for “Stmt” is shown in Fig. 2.7.

Definition of ac language (Cont.)

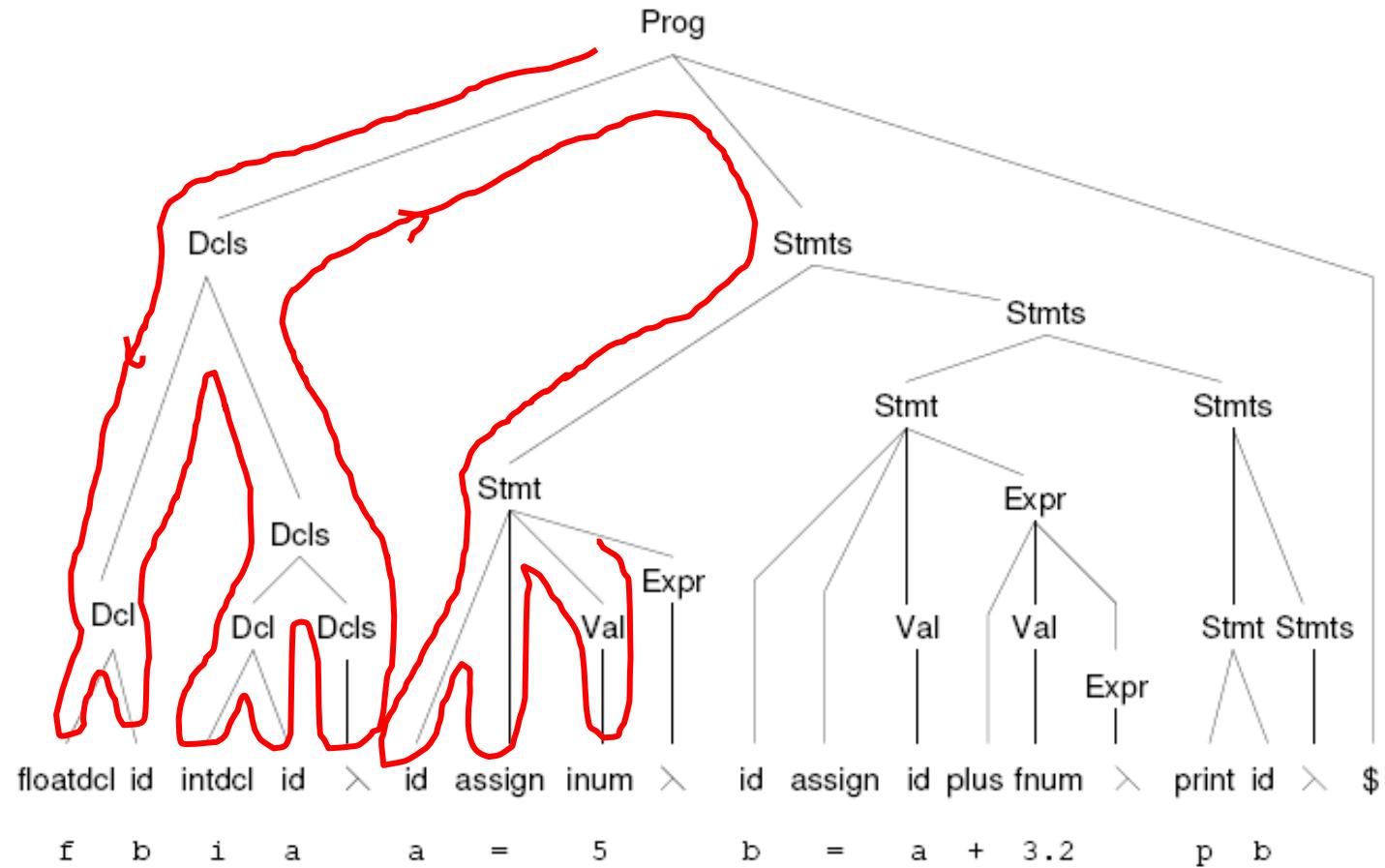


Figure 2.4: An ac program and its parse tree.

Phases of an ac compiler (Cont.)

- Each parsing procedure examines the next input token to predict which production to apply. Ex:

Stmt offers two productions:

$\text{Stmt} \rightarrow \text{id assign Val Expr}$

$\text{Stmt} \rightarrow \text{print id}$

Phases of an ac compiler (Cont.)

- If **id** is the next input token, the parse proceeds with the production:
 $\text{Stmt} \rightarrow \text{id assign Val Expr}$
and the **predict set** for the production is $\{\text{id}\}$.
- If **print** is the next, the parse proceeds with:
 $\text{Stmt} \rightarrow \text{print id}$
the **predict set** for the production is $\{\text{print}\}$.
- If the next input token is neither **id** nor **print**,
neither rule can be applied. (aka, **syntax error**)

Notes

- OK, we can **peek (lookahead)** to choose a production to expand (rewrite) if two or more productions are available
- If a lookahead cannot determine a single production to rewrite, what happens? (See future chapters)
 1. Rewrite grammar
 2. Peek more lookaheads

Phases of an ac compiler (Cont.)

Computing the predict sets used in Stmt is relatively easy, because each production for Stmt begins with a **distinct terminal symbol**:

id or print

In fig 2.7, markers 1 and 6 pick which of the two productions to apply by examining the next input token.

1	Prog	\rightarrow	Dcls	Stmts	\$
2	Dcls	\rightarrow	Dcl	Dcls	
3			λ		
4	Dcl	\rightarrow	floatdcl	id	
5			intdcl	id	
6	Stmts	\rightarrow	Stmt	Stmts	
7			λ		
8	Stmt	\rightarrow	id	assign	Val Expr
9			print	id	
10	Expr	\rightarrow	plus	Val Expr	
11			minus	Val Expr	
12			λ		
13	Val	\rightarrow	id		
14			inum		
15			fnum		

```
procedure STMT()
```

```
    if ts.PEEK( ) = id
```

①

```
        then
```

```
            call MATCH(ts, id)
```

②

```
            call MATCH(ts, assign)
```

③

```
            call VAL( )
```

④

```
            call EXPR( )
```

⑤

```
        else
```

```
            if ts.PEEK( ) = print
```

⑥

```
                then
```

```
                    call MATCH(ts, print)
```

```
                    call MATCH(ts, id)
```

```
                else
```

```
                    call ERROR( )
```

⑦

```
end
```

1 Prog → Dcls Stmt \$
2 Dcls → Dcl Dcls
3 | λ
4 Dcl → floatdcl id
5 | intdcl id
6 Stmt → Stmt Stmt
7 | λ
8 Stmt → id assign Val Expr
9 | print id
10 Expr → plus Val Expr
11 | minus Val Expr
12 | λ
13 Val → id
14 | inum
15 | fnum

For each production, write the LHS code

```
procedure Stmt()
  if ts.PEEK( ) = id
    then
      call MATCH(ts, id) ①
      call MATCH(ts, assign) ②
      call VAL() ③
      call EXPR() ④
    else
      if ts.PEEK( ) = print
        then
          call MATCH(ts, print) ⑤
          call MATCH(ts, id) ⑥
      else
        call ERROR()
    end
  end
```

1	Prog	\rightarrow	Dcls	Stmts	\$
2	Dcls	\rightarrow	Dcl	Dcls	
3			λ		
4	Dcl	\rightarrow	floatdcl	id	
5			intdcl	id	
6	Stmts	\rightarrow	Stmt	Stmts	
7			λ		
8	Stmt	\rightarrow	id	assign	Val Expr
9			print	id	
10	Expr	\rightarrow	plus	Val Expr	
11			minus	Val Expr	
12			λ		
13	Val	\rightarrow	id		
14			inum		
15			fnum		

Figure 2.7: Recursive-descent parsing procedure for Stmt. The variable ts is an input stream of tokens.

Phases of an ac compiler (Cont.)

- Consider the productions for Stmt:

 - $\text{Stmts} \rightarrow \text{Stmt Stmt}$
 - $\text{Stmts} \rightarrow \lambda$

- The predict sets for Stmt can be computed by inspecting the following:
 - $\text{Stmts} \rightarrow \text{Stmt Stmt}$
begins with the non-terminal Stmt.
 - Sol: Find those symbols that predict **any** rule for Stmt.
 - $\text{Stmts} \rightarrow \lambda$ derives no symbols.
 - Sol: Look for what symbol(s) could occur **following** such a production (that is \$).

For production rule that begins with a nonterminal you still need to tell how to keep going based on peeking lookahead

procedure STMTS()

if $ts.\text{PEEK}() = \text{id}$ or $ts.\text{PEEK}() = \text{print}$ ⑧

then

call STMT() ⑨

call STMTS() ⑩

else

if $ts.\text{PEEK}() = \$$ ⑪

then

★ do nothing for λ -product/ ⑫

else call ERROR()

end

1	Prog	$\rightarrow Dcls\ Stmt\ $$
2	Dcls	$\rightarrow Dcl\ Dcls$
3		$ \lambda$
4	Dcl	$\rightarrow \text{floatdcl}\ id$
5		$ \text{intdcl}\ id$
6	Stmts	$\rightarrow Stmt,\ Stmts$
7		$ \lambda$
8	Stmt	$\rightarrow id\ assign\ Val\ Expr$
9		$ print\ id$
10	Expr	$\rightarrow plus\ Val\ Expr$
11		$ minus\ Val\ Expr$
12		$ \lambda$
13	Val	$\rightarrow id$
14		$ inum$
15		$ fnum$

Figure 2.8: Recursive-descent parsing procedure for Stmt.

Let's do semantics !!

```

procedure STMT()
    if ts.PEEK( ) = id
        then
            call MATCH(ts, id)
            call MATCH(ts, assign)
            call VAL( )
            call EXPR( )
    else
        if ts.PEEK( ) = print
            then
                call MATCH(ts, print)
                call MATCH(ts, id)
        else
            call ERROR( )
end

```

- ①
- ②
- ③
- ④
- ⑤
- ⑥
- ⑦

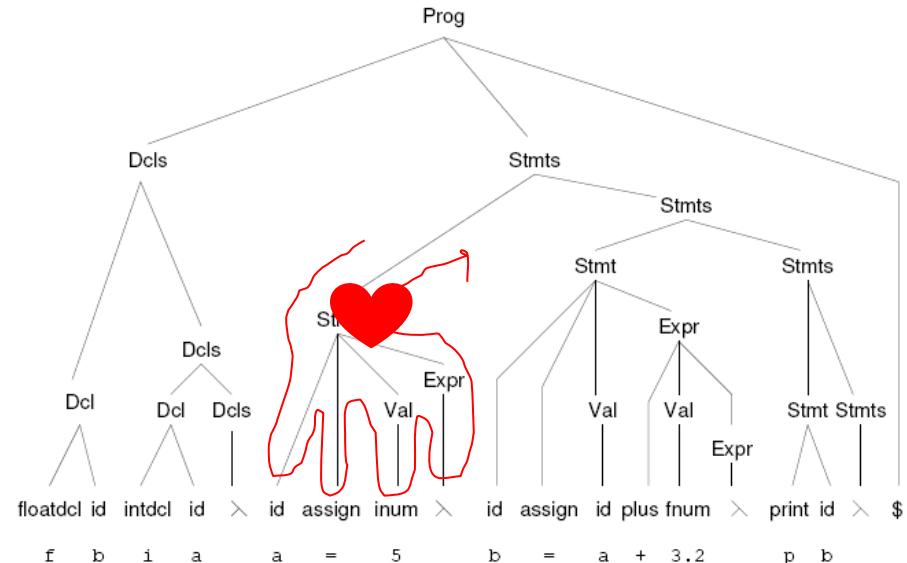


Figure 2.4: An ac program and its parse tree.

OK, when parser does parsing, there are some points that reach a point where you can do semantics. We should do $a=5$ here

Semantics Actions

- You can execute semantics immediately in an interpreter
- In a Compiler,
 - we don't execute semantics immediately
 - We do not produce assembly code yet.

Phases of an ac compiler (Cont.)

- If all of the tokens are processed, an **abstract syntax tree (AST)** will be generated.
 - An example is shown in fig 2.9.
- AST serves as a representation of a program for all phases after **syntax analysis**.

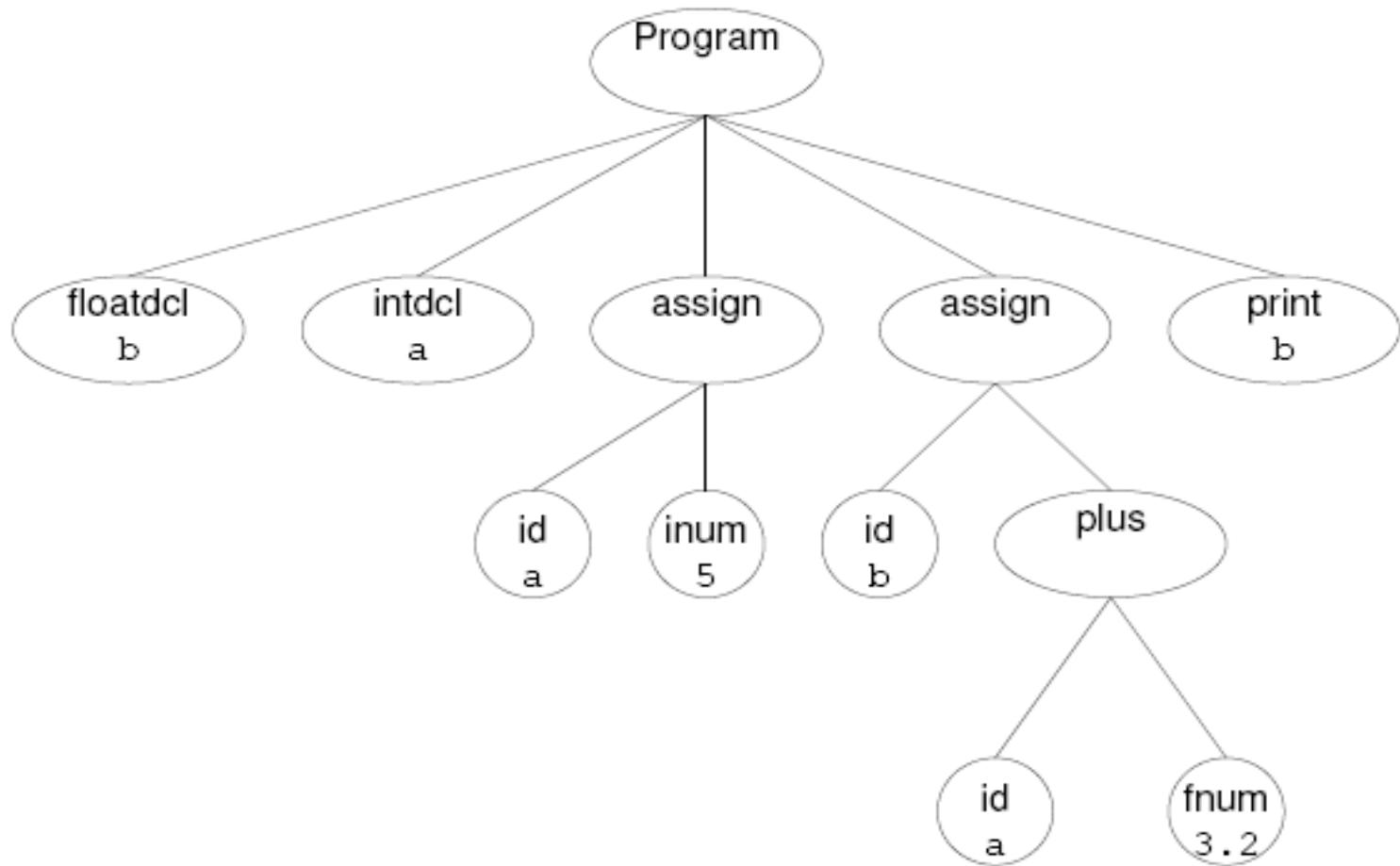


Figure 2.9: An abstract syntax tree for the ac program shown in Figure 2.4.

Phases of an ac compiler (Cont.)

The symbol-table constructor traverses
the AST to record
all *identifiers* and
their *types*
in a **symbol table**.

The constructor is shown in fig 2.10.

Phases of an ac compiler (Cont.)

```
/* Visitor methods */  
procedure VISIT( SymDeclaring n )  
    if n.GETTYPE( ) = floatdcl  
    then call ENTERSYMBOL( n.GETID( ), float )  
    else call ENTERSYMBOL( n.GETID( ), integer )  
end  
  
/* Symbol table management */  
procedure ENTERSYMBOL( name, type )  
    if SymbolTable[name] = null  
    then SymbolTable[name] ← type  
    else call ERROR( "duplicate declaration" )  
end  
  
function LOOKUPSYMBOL( name ) returns type  
    return (SymbolTable[name])  
end
```

Figure 2.10: Symbol table construction for ac.

Original Material by 陳振炎教授,
red material by Yung-Pin Cheng

Phases of an ac compiler (Cont.)

An **ac** symbol table has 23 identifier entries with their types:
integer, float, or unused (null).

Fig 2.11 shows the symbol table constructed.

Phases of an ac compiler (Cont.)

Symbol	Type	Symbol	Type	Symbol	Type
a	integer	k	null	t	null
b	float	l	null	u	null
c	null	m	null	v	null
d	null	n	null	w	null
e	null	o	null	x	null
g	null	q	null	y	null
h	null	r	null	z	null
j	null	s	null		

Figure 2.11: Symbol table for the ac program from Figure 2.4.

Phases of an ac compiler (Cont.)

- Type checking
 - The ac language offers two types:
integer and *float*,
and all identifiers must be type-declared in a
program before they can be used.
 - This process walks the AST **bottom-up** from
its leaves toward its root.

Phases of an ac compiler (Cont.)

- At each node, appropriate analysis is applied:
 - For constants and symbol references, the visitor methods simply set the supplied node's type based on the node's contents.
 - For nodes that compute value, such as **plus** and **minus**, the appropriate type is computed by calling the utility methods.
 - For an assignment operation, the visitor makes certain that the value computed by the second child is of the same type as the assigned identifier (the first child).

The results of applying semantic analysis to the AST of fig 2.9 are shown in fig 2.13.

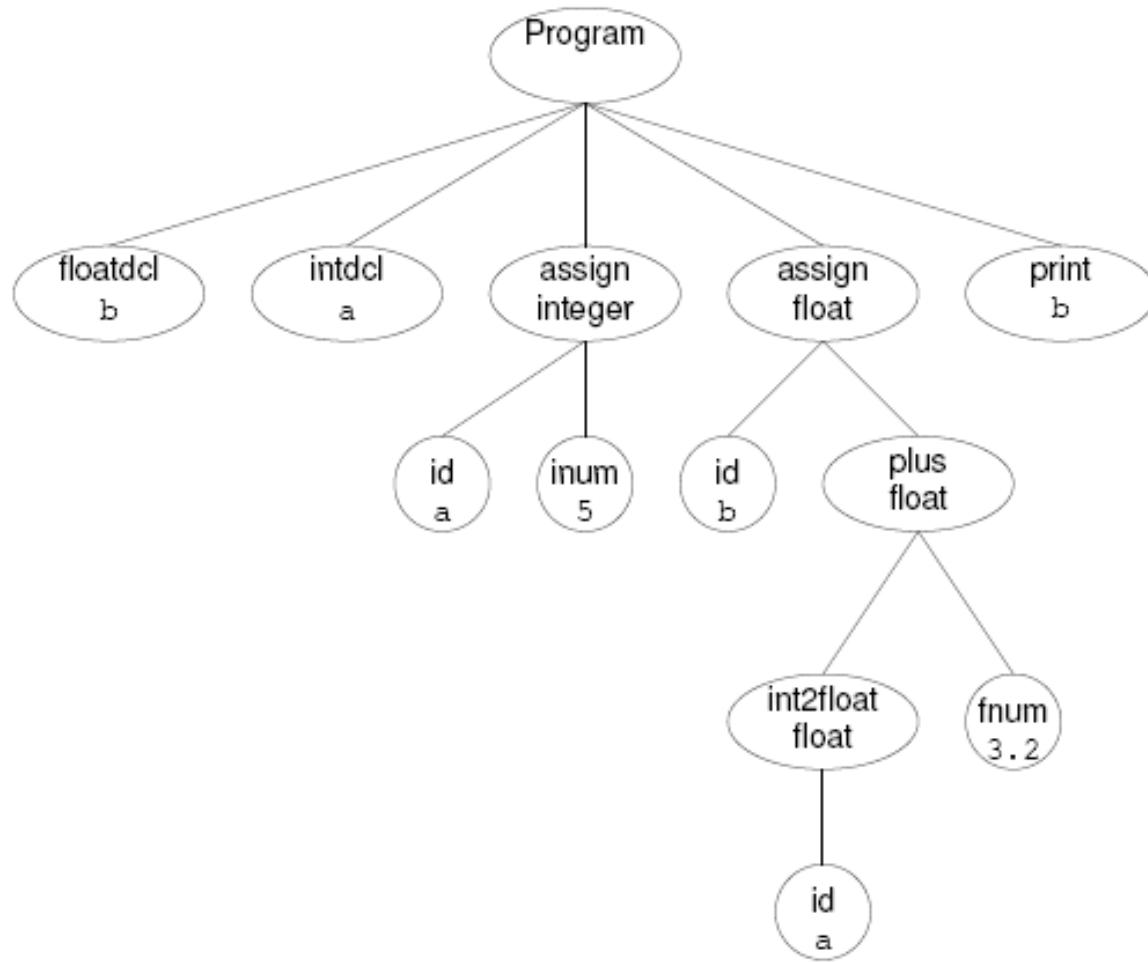


Figure 2.13: AST after semantic analysis.

Phases of an ac compiler (Cont.)

- Code generation
 - Code generation proceeds by traversing (visiting) the AST nodes, starting at its root and working toward its leaves.
 - A visitor generates code based on the node's type (fig 2.14).
 - Fig 2.15 shows the target code generated from the AST in fig 2.9.

```
procedure VISIT( Assigning n )
    call CODEGEN( n.child2 )
    call EMIT( "s" )
    call EMIT( n.child1.id )
    call EMIT( "0 k" )
end
```

(14)

```
procedure VISIT( Computing n )
    call CODEGEN( n.child1 )
    call CODEGEN( n.child2 )
    call EMIT( n.operation )
end
```

(15)

```
procedure VISIT( SymReferencing n )
    call EMIT( "1" )
    call EMIT( n.id )
end
```

```
procedure VISIT( Printing n )
    call EMIT( "1" )
    call EMIT( n.id )
    call EMIT( "p" )
    call EMIT( "si" )
end
```

(16)

```
procedure VISIT( Converting n )
    call CODEGEN( n.child )
    call EMIT( "5 k" )
end
```

(17)

```
procedure VISIT( Consting n )
    call EMIT( n.val )
end
```

Figure 2.14: Code generation for ac
Original Material by 陳振炎教授,
red material by Yung-Pin Cheng

Code	Source	Comments
5 sa 0 k	a = 5	Push 5 on stack Pop the stack, storing (<u>s</u>) the popped value in register <u>a</u> Reset precision to integer
1a 5 k 3.2 + sb 0 k	b = a + 3.2	Load (<u>1</u>) register <u>a</u> , pushing its value on stack Set precision to float Push 3.2 on stack Add: 5 and 3.2 are popped from the stack and their sum is pushed Pop the stack, storing the result in register <u>b</u> Reset precision to integer
1b p si	p b	Push the value of the <u>b</u> register Print the top-of-stack value Pop the stack by storing into the <u>i</u> register

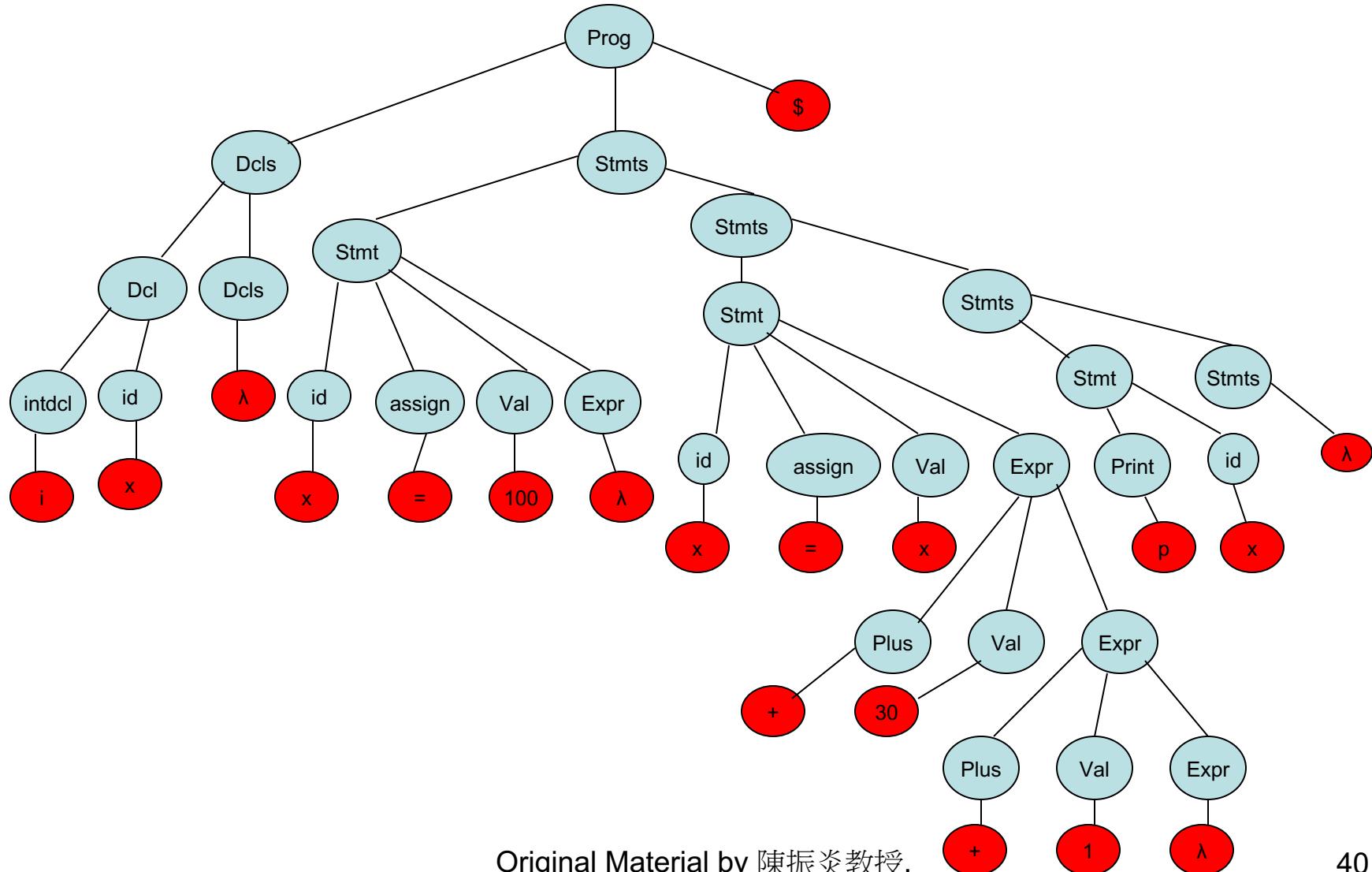
Figure 2.15: Code generated for the AST shown in Figure 2.9.

Homework

1. Based on the grammar in Figure 2.1, construct the parse tree for the following input:

(a) i x x=100 x=x+30+1 p x

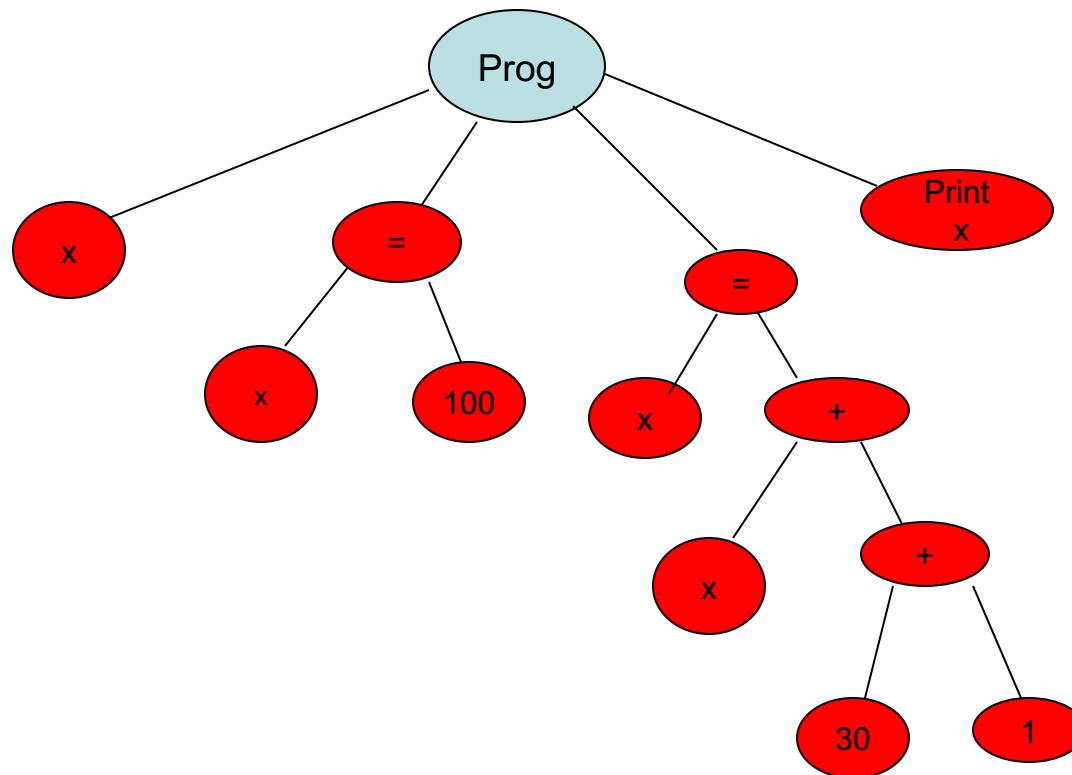
HW Solution 1.(a)



Homework (Cont.)

- 2. For the input shown in Exercise 1:
 - (a) Construct the input's AST
 - (b) According to the definition of **ac**, is the input semantically correct? If not, what changes to the **ac** language would render the input semantically correct?
 - (c) With **ac** changes in place that make the input semantically correct, show the code that would be generated on behalf of the input.

HW Solution 2.(a)



HW Solution 2.(b)

(b) No.

we can do semantic analysis below.

In the AST, the three nodes 30, +, and 1 can be changed to one node 31.

This analysis simplifies the AST.

HW Solution 2.(c)

(c) 100 x=100 push 100 to stack
sx pop top of stack, store it to reg. x
0 k set precision 0

lx x=x+31 load x (push x to stack)
31 push 31 to stack
+ pop 31, pop x, add them, push the result

sx
0 k

lx p x load x (push x to stack)
p print top of stack
si pop top of stack, store it to reg. i

Homework (Cont.)

- 3. Extend the ac scanner (Figure 2.5) in the following ways:
 - (a) A floatdcl can be represented as either f or float, allowing a more Java-like syntax for declarations.
 - (b) An intdcl can be represented as either i or int
 - (c) A num may be entered in exponential (scientific) form. That is, an **ac** num may be suffixed with an optionally signed exponent (1.0e10, 123e-22 or 0.31415926535e1)

HW Solution 3

(a) (b)

Terminal	RegularExpression
floatdcl	“f” (“f” “l” “o” “a” “t”)
intdcl	“i” (“i” “n” “t”)

(c) inum	$[0-9]^+ \text{ e } -? [0-9]^+$
fnum	$[0-9] . [0-9]^+ \text{ e } -? [0-9]^+$