



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

CS323 Lab 6

Yepang Liu

liuyp1@sustech.edu.cn

Outline

- Grammar Design Issues
- Introduction to Bison

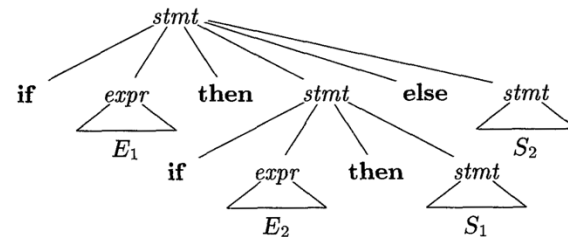
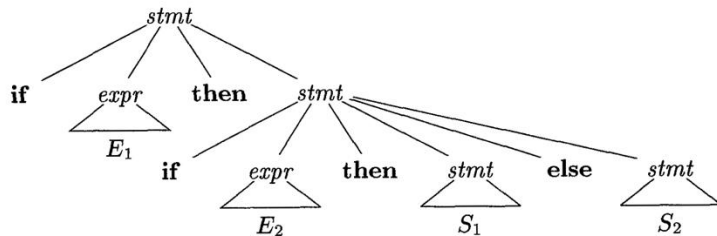
Grammar Design

- CFGs are capable of describing most, but not all, of the syntax of programming languages
 - “Identifiers should be declared before use” cannot be described by a CFG
 - Subsequent phases must analyze the output of the parser to ensure compliance with such rules
- Before parsing, we typically apply several transformations to a grammar to make it more suitable for parsing
 - Eliminating ambiguity (消除二义性)
 - Eliminating left recursion (消除左递归)
 - Left factoring (提取左公因子)

Eliminating Ambiguity (1)

stmt \rightarrow **if** *expr* **then** *stmt*
 | **if** *expr* **then** *stmt* **else** *stmt*
 | **other**

Two parse trees for **if** E_1 **then** **if** E_2 **then** S_1 **else** S_2



Which parse tree is preferred in programming?
(i.e., else matches which then?)

Eliminating Ambiguity (2)

- **Principle of proximity:** match each **else** with the closest unmatched **then**
 - **Idea of rewriting:** A statement appearing between a **then** and an **else** must be matched (must not end with an unmatched **then**)

```
stmt    →  matched_stmt  
        |  open_stmt  
matched_stmt → if expr then matched_stmt else matched_stmt  
        |  other  
open_stmt  → if expr then stmt  
        |  if expr then matched_stmt else open_stmt
```

Rewriting grammars to eliminate ambiguity is difficult.
There are no general rules to guide the process.



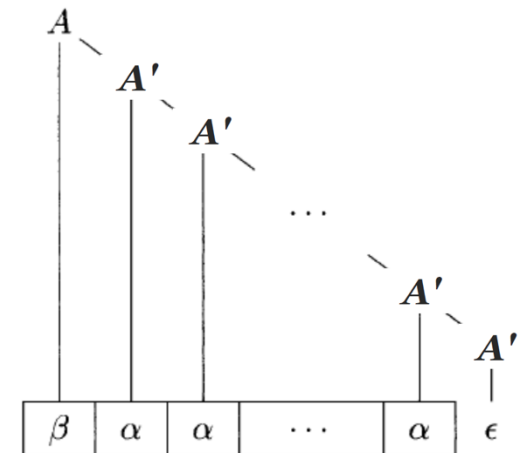
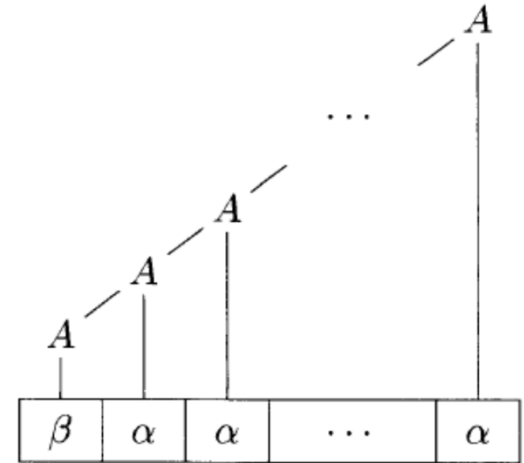
* open_stmt means the last then may not have matching else

Eliminating Left Recursion

- A grammar is **left recursive** if it has a nonterminal A such that there is a derivation $A \Rightarrow^+ A\alpha$ for some string α
 - $S \rightarrow Aa \mid b$
 - $A \rightarrow Ac \mid Sd \mid \epsilon$
 - Because $S \Rightarrow Aa \Rightarrow Sda$
- **Immediate left recursion (立即左递归)**: the grammar has a production of the form $A \rightarrow A\alpha$
- Top-down parsing methods cannot handle left-recursive grammars (bottom-up parsing methods can handle...)

Eliminating Immediate Left Recursion

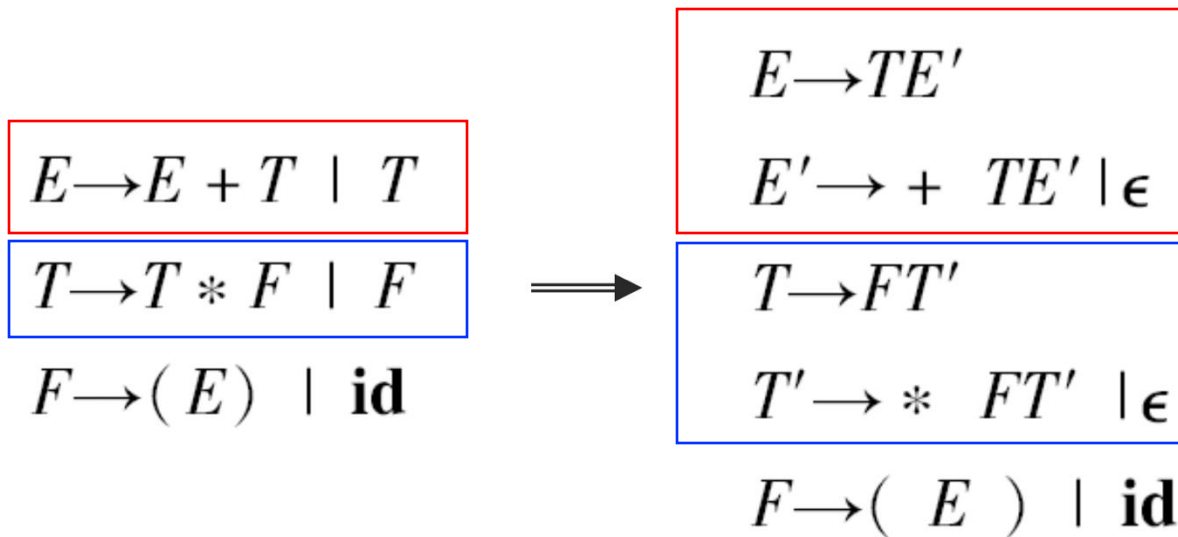
- Simple grammar: $A \rightarrow A\alpha \mid \beta$
 - It generates sentences starting with the symbol β followed by zero or more α 's
- Replace the grammar by:
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' \mid \epsilon$
 - It is right recursive now



Eliminating Immediate Left Recursion

- The general case: $A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$
- Replace the grammar by:
 - $A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$
 - $A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \epsilon$

Example



Left Factoring (提取左公因子)

- If we have the following two productions

$$\begin{aligned} stmt &\rightarrow \text{if } expr \text{ then } stmt \text{ else } stmt \\ &| \text{ if } expr \text{ then } stmt \end{aligned}$$

- On seeing input **if**, we cannot immediately decide which production to choose
- In general, if $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$ are two productions, and the input begins with a nonempty string derived from α . We may defer choosing productions by expanding A to $\alpha A'$ first

$$\begin{aligned} A &\rightarrow \alpha A' \\ A' &\rightarrow \beta_1 \mid \beta_2 \end{aligned}$$

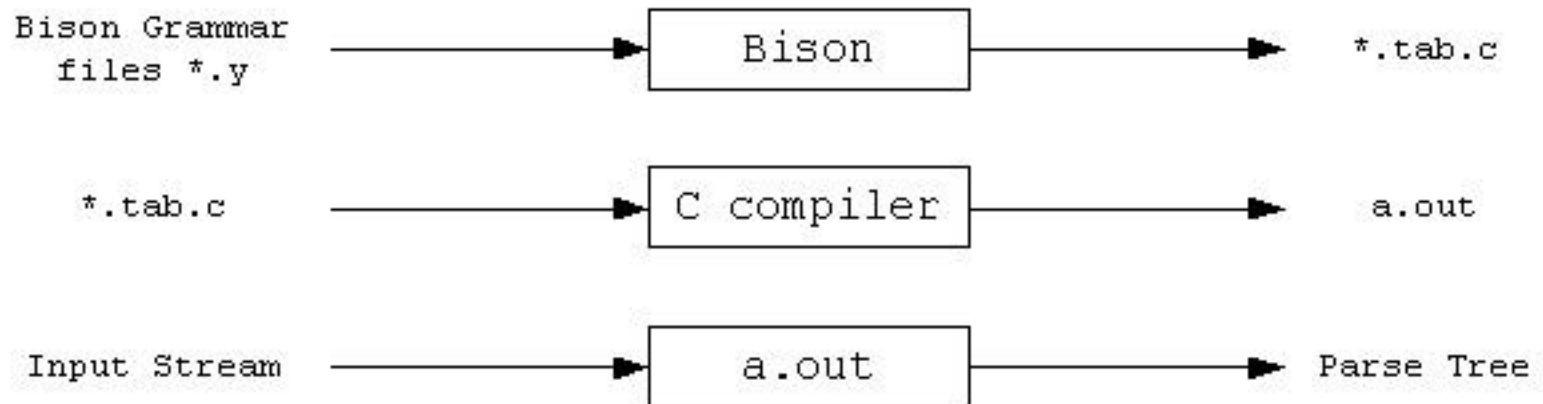
Outline

- Grammar Design Issues
- Introduction to Bison

The Parser Generator Bison

- Bison generates a parser, which accepts the input token stream from Flex, to do syntax analysis, according to the specified CFG.
- Bison的前身为基于Unix的Yacc (Yet another compiler compiler)。Yacc的发布时间比Lex还要早，其采用的语法分析技术的理论基础早在20世纪50年代就已经由Knuth逐步建立了起来，而Yacc本身则是贝尔实验室的S.C. Johnson基于这些理论在1975年到1978年写成的。
- 到了1985年，当时在UC Berkeley的一个研究生Bob Corbett在BSD下重写了Yacc，取名为Bison (美洲野牛，yak牦牛的近亲)，后来GNU Project接管了这个项目，为其增加了许多新的特性，于是就有了我们今天所用的GNU Bison。

Input/Output of Bison (Yacc)



Structure of YACC Source Programs

- **Declarations (声明)**

- Ordinary C declarations
- Grammar tokens

- **Translation rules (翻译规则)**

- Rule = a production + semantic action

- **Supporting C routines (辅助性C语言例程)**

- Directly copied to `y.tab.c`
- Can be invoked in the semantic actions
- Other procedures such as error recovery routines may be provided

```
declarations
%%
translation rules
%%
supporting C routines
```

Translation Rules

```
⟨head⟩    :  ⟨body⟩1  { ⟨semantic action⟩1 }  
           |  ⟨body⟩2  { ⟨semantic action⟩2 }  
           ...  
           |  ⟨body⟩n  { ⟨semantic action⟩n }  
           ;
```

DO NOT miss it

- The first head in the list of rules is taken as the start symbol of the grammar
- A semantic action is a sequence of C statements
 - `$$` is a Bison's internal variable that holds the semantic value of the left-hand side of a production rule (i.e., the whole construct)
 - `$i` holds the semantic value of *i*th grammar symbol of the body
- A semantic action is performed when we apply the associated production for reduction (归约, the reverse of rewrite)
 - Most of the time, the purpose of an action is to compute the semantic value of the whole construct from the semantic values of its parts, i.e., compute `$$` using `$i`'s

Example: A Parser for Calculation

- When processing arithmetic expressions, lexical analyzer will recognize the following types of tokens: INT, ADD, SUB, MUL, DIV
 - When analyzing the input string “3 + 6 / 2”, the lexer will output this string of tokens: **INT ADD INT DIV INT**
- After lexical analysis, the parser will check if the string of tokens produced by lexer has a valid structure or not according to the syntactic specification described by the following context-free grammar

```
Calc -> Exp
Exp -> Factor | Exp ADD Factor | Exp SUB Factor
Factor -> Term | Factor MUL Term | Factor DIV Term
Term -> INT
```

*Red for non-terminals, Blue for terminals, Calc is the start symbol

Valid and Invalid Inputs

Valid input expressions:

- 3
- $3 + 5 * 4$
- $3 + 6 / 2$
- $3 + 2 - 1$
- ...

Invalid input expressions:

- -3
- $3 + 5 *$
- ...

Flex/Bison Code for Calculator

syntax.y

lex.l

```
%{
    #include "syntax.tab.h"
    #include "stdlib.h"
}%
%%
[0-9]+ { yylval = atoi(yytext); return INT; }
"+" { return ADD; }
"-" { return SUB; }
"*" { return MUL; }
"/" { return DIV; }
[ \n\r\t] {}
. { fprintf(stderr, "unknown symbol: %s\n", yytext);
  exit(1); }
```

yyval:

- Flex internal variable that is used to store the attribute of a recognized token
- Its data type is YYSTYPE (int by default)*
- After storing values to `yyval` in Flex code, the values will be propagated to Bison (i.e., the syntax analyzer part) and can be retrieved using `$n`

```
%{
    #include "lex.yy.c"
    void yyerror(const char*);
}%
%token INT
%token ADD SUB MUL DIV
%%
Calc: /* to allow empty input */
    | Exp { printf("= %d\n", $1); }
    ;
Exp: Factor
    | Exp ADD Factor { $$ = $1 + $3; }
    | Exp SUB Factor { $$ = $1 - $3; }
    ;
Factor: Term
    | Factor MUL Term { $$ = $1 * $3; }
    | Factor DIV Term { $$ = $1 / $3; }
Term: INT
    ;
%%
void yyerror(const char *s) {
    fprintf(stderr, "%s\n", s);
}
int main() {
    yyparse(); // will invoke yylex()
}
```

Can be customized by putting command like `#define YYSTYPE char` at the beginning of .l and .y files.

A Further Look at Semantic Actions

```
%{  
    #include "lex.yy.c" syntax.y  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

When applying this rule for reduction, the semantic analysis is successful (i.e., the start symbol calc can generate the input expression), we output the calculation result.

A Further Look at Semantic Actions

```
%{
    #include "lex.yy.c"
    void yyerror(const char*);
}%
%token INT
%token ADD SUB MUL DIV
%%
Calc: /* to allow empty input */
    | Exp { printf("= %d\n", $1); }
    ;
Exp: Factor
    | Exp ADD Factor { $$ = $1 + $3; }
    | Exp SUB Factor { $$ = $1 - $3; }
    ;
Factor: Term
    | Factor MUL Term { $$ = $1 * $3; }
    | Factor DIV Term { $$ = $1 / $3; }
    ;
Term: INT
    ;
%%
void yyerror(const char *s) {
    fprintf(stderr, "%s\n", s);
}
int main() {
    yyparse();
}
```

syntax.y

Intermediate calculation of semantic values of a language construct represented by the head symbol

A Further Look at Semantic Actions

```
%{
    #include "lex.yy.c"
    void yyerror(const char*);
}%
%token INT
%token ADD SUB MUL DIV
%%
Calc: /* to allow empty input */
    | Exp { printf("= %d\n", $1); }
    ;
Exp: Factor
    | Exp ADD Factor { $$ = $1 + $3; }
    | Exp SUB Factor { $$ = $1 - $3; }
    ;
Factor: Term
    | Factor MUL Term { $$ = $1 * $3; }
    | Factor DIV Term { $$ = $1 / $3; }
    ;
Term: INT
    ;
%%
void yyerror(const char *s) {
    fprintf(stderr, "%s\n", s);
}
int main() {
    yyparse();
}
```

syntax.y

Although no action is specified, Bison will still propagate the attribute of the only symbol INT to the head Term

Same as writing: `$$ = $1;`

Illustrative Example

```
%{  
    #include "lex.yy.c" syntax.y  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

Input: 3 + 5

Token string: INT ADD INT

3 INT

Illustrative Example

```
%{  
    #include "lex.yy.c" syntax.y  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

Input: 3 + 5

Token string: INT ADD INT

3 Term

3 INT

Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT

3 Factor

3 Term

3 INT

Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT

3 Exp
|
3 Factor
|
3 Term
|
3 INT

Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT

3 Exp ADD
|
3 Factor
|
3 Term
|
3 INT

Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT

3 Exp ADD
|
3 Factor
|
3 Term
|
3 INT

INT 5

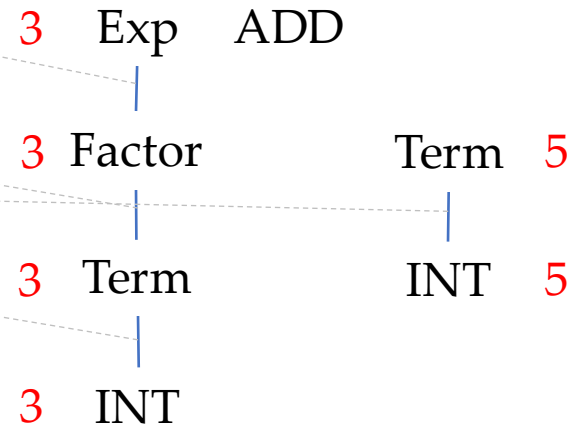
Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT



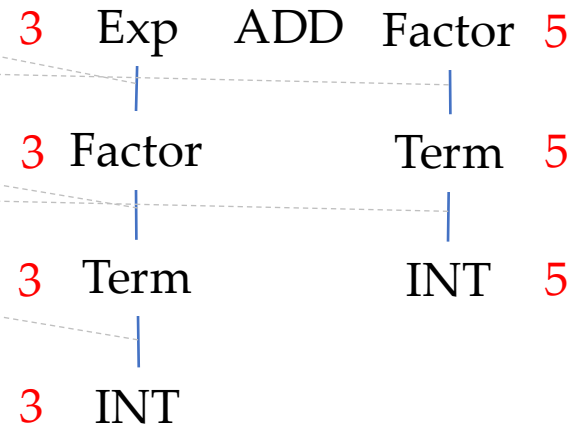
Illustrative Example

```
%{  
    #include "lex.yy.c"  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT

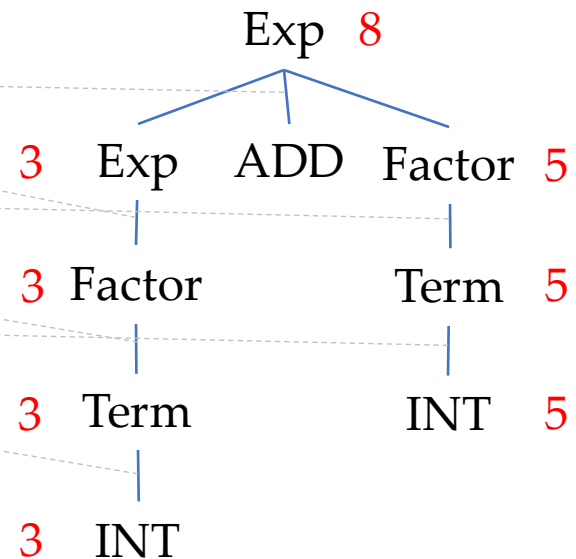


Illustrative Example

```
%{  
    #include "lex.yy.c" syntax.y  
    void yyerror(const char*);  
%}  
%token INT  
%token ADD SUB MUL DIV  
%%  
Calc: /* to allow empty input */  
    | Exp { printf("= %d\n", $1); }  
    ;  
Exp: Factor  
    | Exp ADD Factor { $$ = $1 + $3; }  
    | Exp SUB Factor { $$ = $1 - $3; }  
    ;  
Factor: Term  
    | Factor MUL Term { $$ = $1 * $3; }  
    | Factor DIV Term { $$ = $1 / $3; }  
    ;  
Term: INT  
    ;  
%%  
void yyerror(const char *s) {  
    fprintf(stderr, "%s\n", s);  
}  
int main() {  
    yyparse();  
}
```

Input: 3 + 5

Token string: INT ADD INT



Illustrative Example

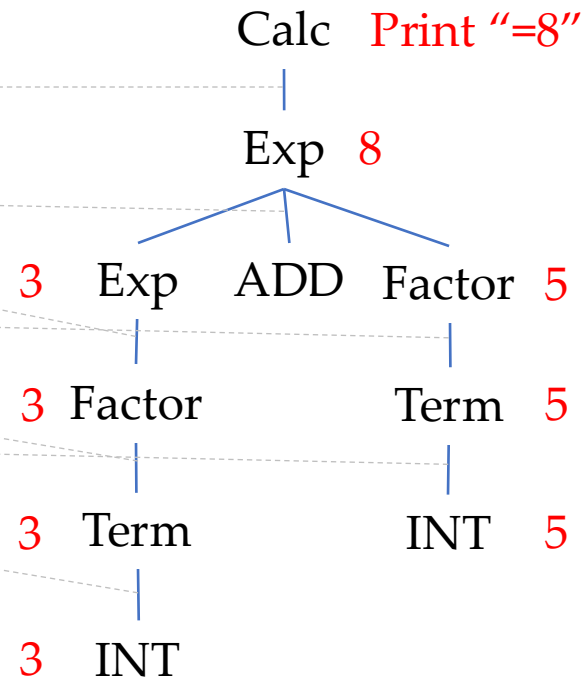
```

%{
    #include "lex.yy.c"
    void yyerror(const char*);
}%
%token INT
%token ADD SUB MUL DIV
%%
Calc: /* to allow empty input */
    | Exp { printf("= %d\n", $1); }
    ;
Exp: Factor
    | Exp ADD Factor { $$ = $1 + $3; }
    | Exp SUB Factor { $$ = $1 - $3; }
    ;
Factor: Term
    | Factor MUL Term { $$ = $1 * $3; }
    | Factor DIV Term { $$ = $1 / $3; }
    ;
Term: INT
    ;
%%
void yyerror(const char *s) {
    fprintf(stderr, "%s\n", s);
}
int main() {
    yyparse();
}
    
```

syntax.y

Input: 3 + 5

Token string: INT ADD INT



Compile the Example

- Compile the Bison code into C code: a header file and an implementation file
 - `bison -d syntax.y` (why it is compilable without `lex.yy.c`?)
 - The command produces `syntax.tab.h` and `syntax.tab.c`
- Compile the flex source code
 - `flex lex.l`
 - The command produce `lex.yy.c`
- Putting things together
 - `gcc syntax.tab.c -lfl -ly -o calc.out`
 - The options `-lfl` and `-ly` tell gcc to include Flex and Bison libraries (in some environments, the options may not be needed)

Run the Calculator

- In the runs below, we use pipes* to pass the output of the first command to be the input of the second command

```
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "3" | ./calc.out
= 3
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "3+5*4" | ./calc.out
= 23
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "3+6/2" | ./calc.out
= 6
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "3+ 2-1" | ./calc.out
= 4
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "-3" | ./calc.out
syntax error
cs323@deb-cs323-compilers:~/Desktop/lab6$ echo "3+5*" | ./calc.out
syntax error
```

* <https://linuxhint.com/what-is-pipe-in-linux/>