



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY

CS323 Lab 4

Yepang Liu

liuyp1@sustech.edu.cn

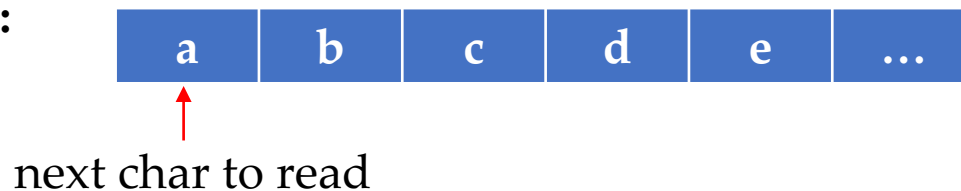
Outline

- Flex Library Functions
 - `input()`, `unput()`, `yylless()`, `yymore()`
- Grammar Design Issues

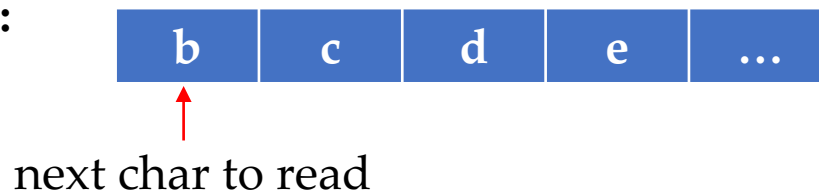
The input() Function

- The function takes no arguments
- When called, it reads a single character from the input buffer and return it to the caller

Input buffer:
(before)



Input buffer:
(after)



The unput(char c) Function

- The function puts **c** back into the input buffer

Input buffer:
(before)



↑
next char to read

Input buffer:
(after unput('a'))



↑
next char to read

Example

- End-of-line comments sanitizer

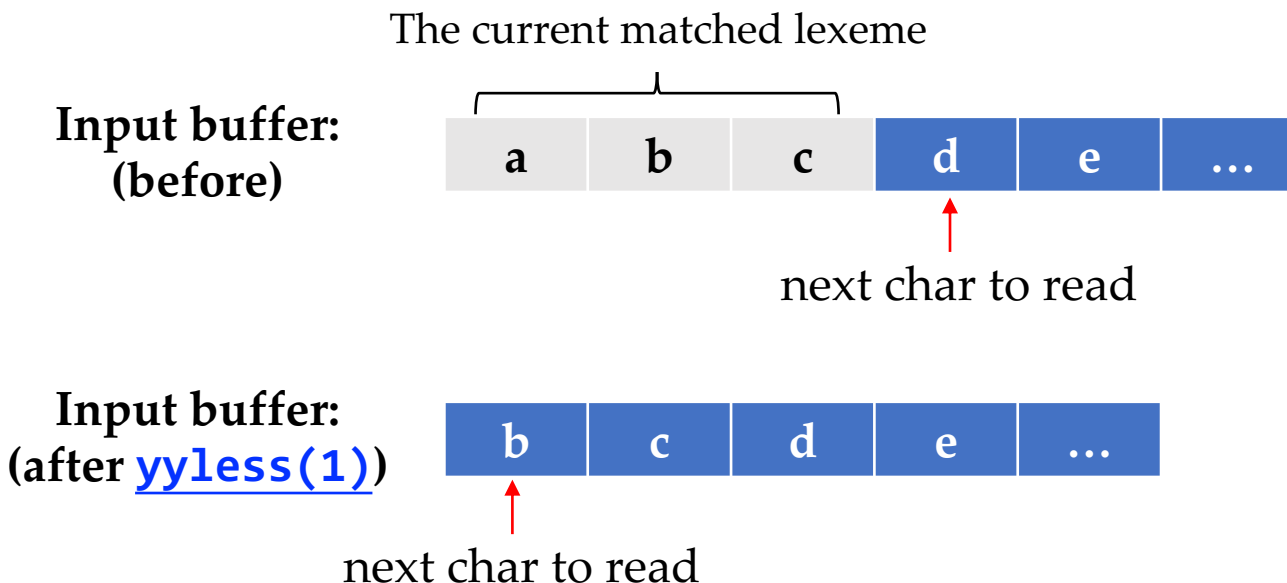
```
"//" {  
    // ignore the following chars until seeing a newline character  
    while((c = input()) != '\n');  
    // put the newline character back to the input buffer  
    unput(c);  
}
```

Steps:

- Go to the “lab4/comment_sanitizer” directory
- Run command “make sanitizer”
- Run command “./sanitizer.out test.c” and observe the effect (compare the output after running the command with the original C program in test.c)

The yyless(int n) Function

- The function returns the **yy~~leng~~-n** characters in the postfix of the current lexeme back to the input buffer



The yymore() Function

- The function causes the next matched token's `yytext` to be appended to the current `yytext`

```
abc { yymore(); }  
def { printf("%s\n", yytext); }
```

Input string:



When matching “def”, `yytext` will be “abcdef” instead of “def”.

Exercise

- Dealing with nested quotation marks when recognizing string literals

```
printf("This is a string literal without nested quotation marks");  
printf("And God said, \"Let there be light,\" and there was light.");
```

If we have the following translation rules:

```
\("[^"]*" { printf("Matched a string literal: %s\n", yytext); }  
\n {}  
. {}
```

When processing the above print statements, we will see this output:

```
Matched a string literal: "This is a string literal without nested quotation marks"  
Matched a string literal: "And God said, \"  
Matched a string literal: " and there was light."
```


Exercise

- Please modify the translation rule for string literals such that when processing the previous two print statements, we will see the correct output (hint: use `yylless` and `yymore` to manipulate the input buffer and `yytext`)

```
Matched a string literal: "This is a string literal without nested quotation marks"  
Matched a string literal: "And God said, \"Let there be light,\" and there was light."
```

Go to the “lab4/nested_quotation_marks” directory. We have provided the `lex.l` file and the input program file `test.c`. The build target “nest” can be used. You only need to modify the `lex.l` file and then try to run the command “`./nest test.c`”.

Outline

- Flex Library Functions
 - `input()`, `unput()`, `yylless()`, `yymore()`
- Grammar Design Issues

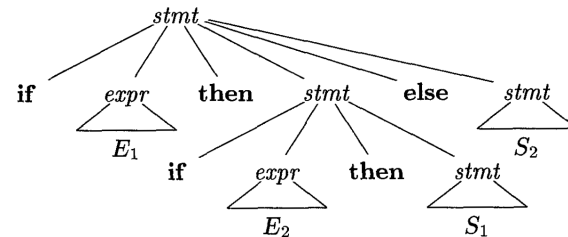
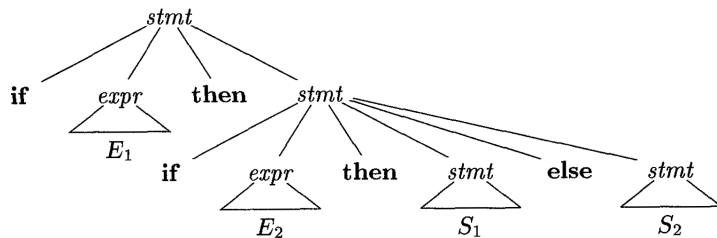
Grammar Design

- 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 → **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.



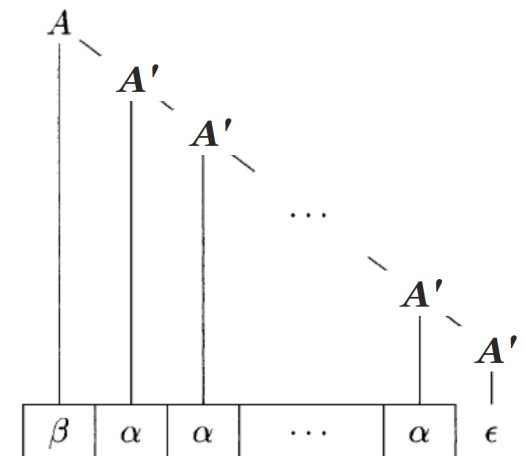
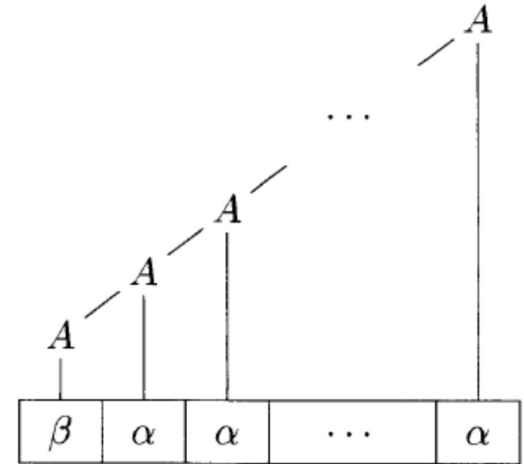
* open_stmt means the last then may not have matching else

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

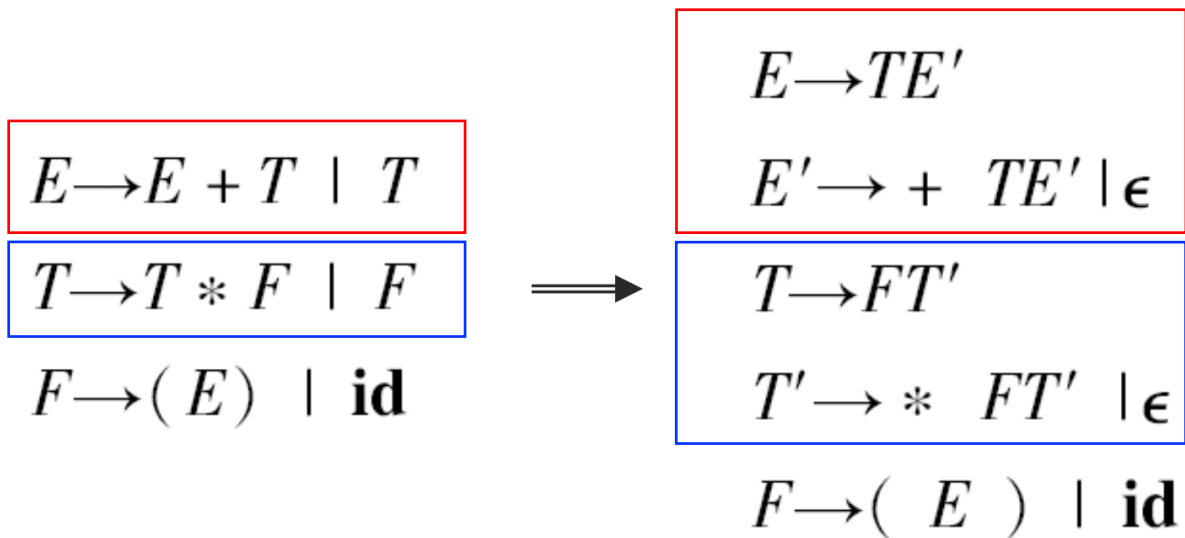
- Simple grammar: $A \rightarrow A\alpha \mid \beta$
 - It generates sentences starting with the symbol β followed by zero or more α 's
- Replace the grammar by:
 - $A \rightarrow \beta A'$
 - $A' \rightarrow \alpha A' \mid \epsilon$
 - It is right recursive now



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

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