

## CS 3200 - Scanning &amp; Flex

Announcements

- Essay due Wed.
- No class Mon or Wed
- HW2 is posted, due Sat Jan 30

# Compilers

Front end:

- (A) Scanners
- (B) Parser
- (C) Semantic Analysis

We're focusing on Scanners first...

## Scanning (lexical analysis)

- Divide program into tokens, or small bits of meaningful units

Ex: for, recognize keywords, group operations, names, etc.

- Scanning or tokenizing makes parsing much simpler.
- While parsers can work character by character, it is slow.
- Note: Scanning is recognizing a regular language, eg via DFA

Last time: Scanners recognize:

- regular expressions  
(or reg. languages)

- can also form as DFA or NFA  
(more next week on these)

- Some limitations: can't do anything  
with memory:

Ex.:  $a^n b^n$   
exp w/ equal # of ( → )

But: good for tokenizing!  
(later will parse)

Scanners: do this in code

Find the syntax (not semantics) of code.

Output tokens.

A few types:

- Ad-hoc
- Finite automata
  - nested case statements
  - table + drivers

## Ad-hoc: Case based code

```
if current ∈ { "(", ")", "+", "-", "*", "/" }  
    return that symbol  
if current = ":"  
    read next  
    if it is = , announce "assign"  
    else announce error  
if current = "/"  
    read next  
    if it is "*" or "/"  
        read until "*" or "/" or "newline" (resp.)  
    else return divide  
etc.
```

# Ad-hoc approach

Advantage:  
code is fast & compact

Disadvantage:  
very ad-hoc!  
- hard to debug  
- no explicit representation

primary disadvantages

DFA approach:

Given a regular expression, convert to a DFA,

we'll walk through this next week - in Ch 2 of the book.

However





## Scanners Programs

In reality, this DFA is often done automatically.

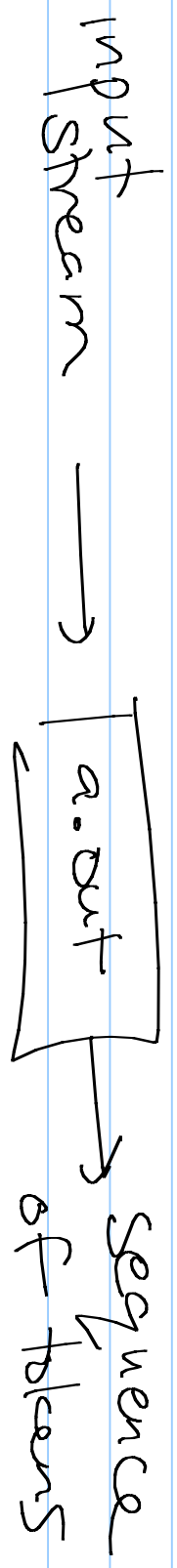
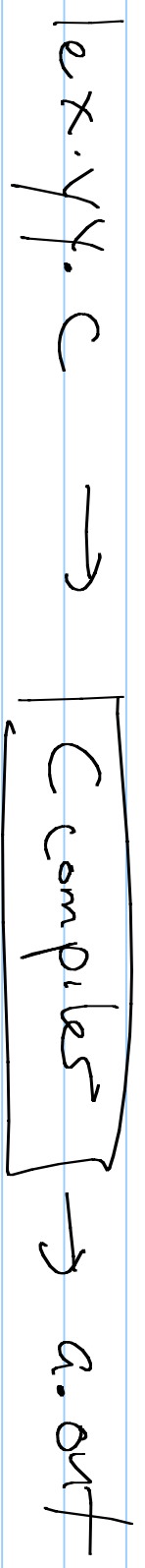
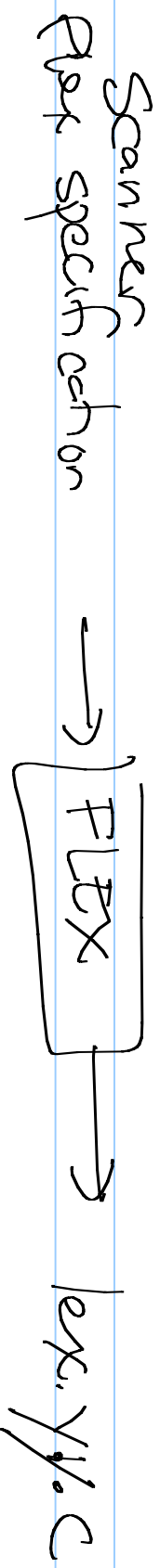
Specify the rules of regular language, as the program does this for you.

Many Such examples:

Lex (Flex), JLex / JFlex,  
Blex, REX, ...

# Flex

- A C driven scanning program.



To compile:

> flex File.lex

> gcc lex.yy.c -lfl

> ./a.out ( < otherFile.txt )  
if input

(if using stdin, might need ctrl-d)

Format for lex files:

① definitions

% %

② rules

% %

user code

(see examples)

1

## Definitions

New definitions to make life easier.

Form: name definition

Ex:

digit

$[0-9]$

$[A-Z]$  also

ID

$[a-z][a-z0-9]^*$

Note: These are regular expressions!

## Definitions cont.:

- An unindented comment (~~/~~\*) is copied verbatim to output, up to the next ~~\*~~

- Any indented text or text enclosed in {} ... {} is also copied verbatim (with %{ ... %} removed)

- %top makes sure things are copied to top of output (for example, for #includes)

## 2 Rules Section

Format: `pattern action`

where `pattern` is unindented, `action` is on the same line

Any indented or `%{ %}` can be used to declare variables. local to the scanning routine.

(Other things may cause compile issues)

## Allowed Patterns

'X' - matches the character X  
'.' - any char except newline

'[xyz]' - matches x, y, or z

'[a-j-o-z]' - matches a, b, j, k, l, m, n, o, z



## More patterns

$[^A-Z]'$  - chars other than A-Z  
(negation)

$[^A-Z\n ]'$  - any char except A-Z or a newline

$[a-z]\{-\}[aeiou]$  - any lower case consonant

$\{0\}$  or more  $\{1\}$  or more  
 $\{#\}$ ,  $\{+\}$

## Patterns (again)

$\frac{r?}{r?}$  0 or 1 r's

$r\{2-5\}$  Between 2 + 5 r's

$r\{2,\}$  2 or more r's

$r\{4\}$  exactly 4 r's

$\{name\}$  expansion of name definition

$r\$$  r at end of a line

(post webpage)

Precedence:

$fb^*$  |  $ba^*$   
is same as  $(fb)^* | (ba)^*$

(since  $*$  has higher precedence than concatenation,  $fb^*$  concatenation is higher than  $or$ )

## C classes

`[[:alpha:]]` matches anything that satisfies `.isalpha()`

Ex:

`[[:alnum:]]`

`[[:alpha:]][[:digit:]]`

`[[:alpha:]][0-9]`

`[a-zA-Z][0-9]`

3

User code

Optional, if just copied directly  
to the output.

(If empty, leave off  $\%s\%$ )

## Comments

- C style : ~~/\*~~ ~~\*/~~

Exceptions Do: pattern rule ~~No comment~~ ~~\*/~~

- No comments in the rule section when a regular expression is expected (so not beginning of line or after spaces)
- Not an option line of `defn` headers

## How it works

- Finds longest pattern match possible

- That match (or token) is made available for a global char pointer  
yytext w/ length yylen

Then action is performed

- If no match, next char goes to std  
(So % %  
is valid.)

# Actions

Ex.      % %  
"zap me"

Ex.      %      %  
[ [ + +  
          + +  
          \$      put char ( ' ' );  
                  / # ignore # /



## Actions (cont)

- If action contains a  $\{$ , then action: spans until next  $\}$  (and may go over many lines)
- Action | means "Same action as the next rule"
- Can be arbitrary C code, including a  $\{$  when run again continues from where it left off.

## Special Actions

- ECHO

- BEGIN followed by name of a  
start condition places  
Scanners in that condition  
(more on this later...)

- REJECT tells scanner to go  
to second best rule

↳ Caution: slow

Ex: Count the # of words  
pattern to look for:

Characters [a-z A-Z ']

Characters\* [ ' / n | ]

Ex:

% %

a |

ab |

- abc

abcd | ECHO; REFLECT;

Scans: xyzabcd

Output? xyzabcdabcbaabcd

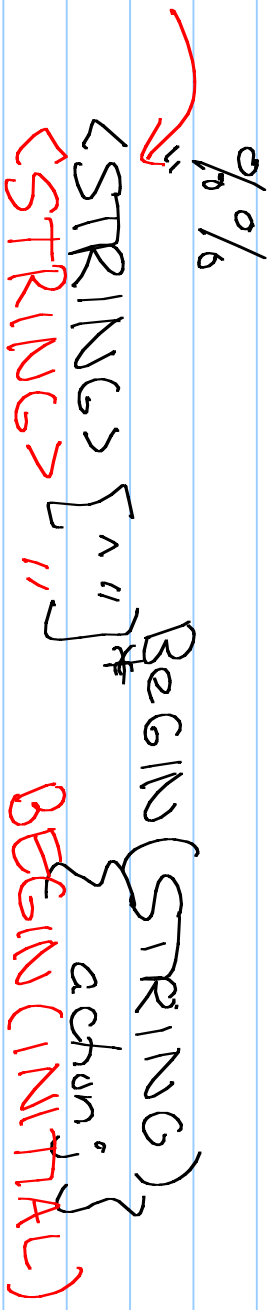
# Conditional Rules

- State based! activated using BEGIN

Define a set of states

- INITIAL is there by default
- Rest defined in %s or %X in first section

Ex: %s STRING



%s are inclusive start conditions  
%x are exclusive start conditions

After BEGIN, state is active.

If state is inclusive, then rules with no start conditions are still active.

If state is exclusive, then rules with no start conditions are inactive.

Ex: %s versus %x

%s example  
%%

<example> foo action();

bar other\_action();

vs:

%x example  
%%

<example> foo action();

<INITIAL, example> bar other\_action();

## Conditions

- `<*>` matches all states
- default rule is in all states.

Essentially, pretend:

```
<*>.\n ECHO;
```

is a line of your file.



Ex: Scanner to ignore comments  
%X comment  
%X comment  
%X comment

```
int num_line = 1;
```

```
" /*" BEGIN (comment);
```

```
<comment> [\n*\n]*  
<comment> "*" + [\n*\n]
```

```
<comment> \n* + line_num;  
<comment> "*" + "/"  
BEGIN (INITIAL);
```

Can condense

← comment > {

all rules

}