

CS344 - Scanning & flex

Note Title

1/27/2012

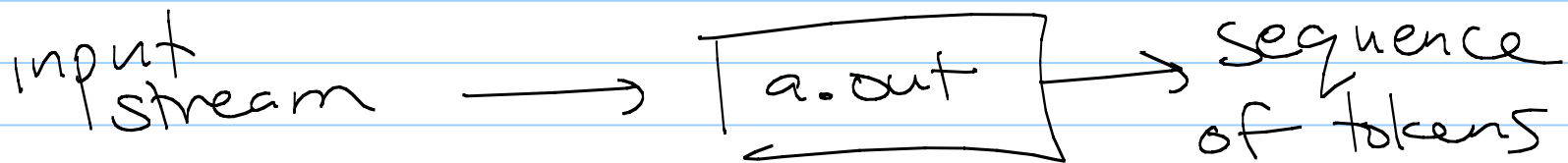
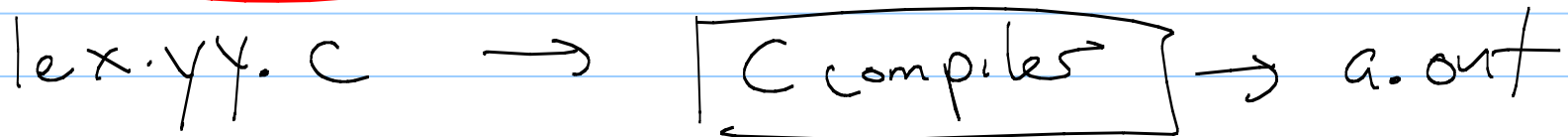
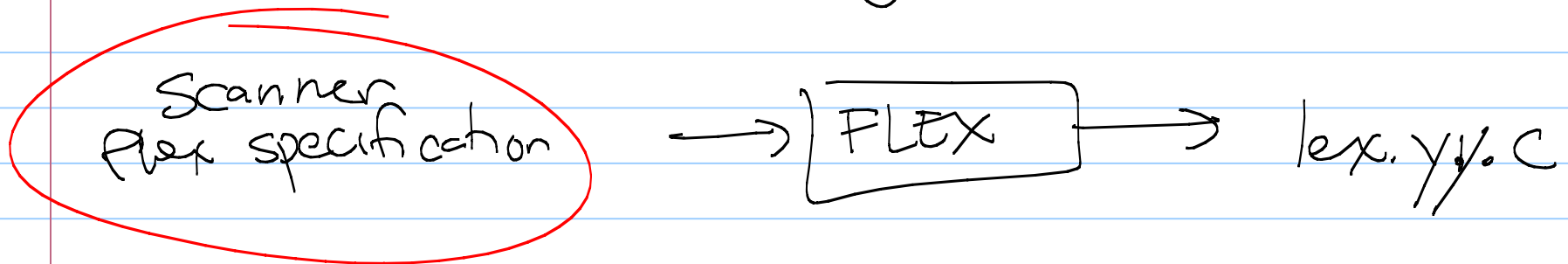
Announcements

- Turn in HW

- HW2 - due next Friday

Flex

- A C driven scanning program.



Format for .lex files:

definitions

% %

rules

% %

user code

(see examples)

Definitions

New definitions to make life easier.

Form: name definition

Ex:

digit [0-9]

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 ~~commented~~ ^{copied} verbatim (with `% { ... } %` removed)
- `% top` makes sure things are copied to top of output (for example, for `#includes`)

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

' $[\wedge A-z]$ ' - chars other than A-z
(negation)

' $[\wedge A-z \backslash n]$ ' - any char except A-z or
a newline

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

r^* , r^+

Patterns (again)

'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
variables digit

'r\$' r at end of a line

⋮
(post webpage)

Precedence:

foo | bar *
is same as (foo) | (bar)*

(since * has higher precedence than concatenation, & concatenation is higher than or)

(foo) | (bar)*

C classes

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

Ex:

`[[:alnum:]]`

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

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

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

User code

Optional, & just copied directly
to the output.

Comments

`% { // ... % }`

- C style : `/* */`

Exceptions

- No comments in the rule section when a regular expression is expected (so not beginning of line or after scanner states)
- Not on % option line of definitions

How it works

- Finds longest pattern match possible
- That match (or token) is made available to a global char pointer `yytext` w/ length in `yylen`

Then action is performed

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

Actions

(empty)

Ex.

% %
"zap me"

} remove
this string

Ex:

% %
[... [t]+
[t]+ \$

putchar(' ');
/* ignore */

* Strip tabs & spaces

→ program to strip extra whitespace

pattern action

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 return, \downarrow (when run again continues from where \downarrow left off.)

Special Actions

- ECHO

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

- REJECT tells scanner to go to second best rule

↳ Caution: slow

Ex: Word count

```
int word_count = 0;
```

```
% %
```

```
kitten
```

```
special(); REJECT;
```

```
[^ |\t|\n]+ ++word_count;
```

(default rule)

++word_count

Ex:

% %

a |
ab |
abc |
abcd |

ECHO; REJECT;

→ . / \n /* does nothing */ /
↑ anything but \n

Scans: x y z a b c d

Output? abcdabcaba

a |
ab |
abc |
abcd | ECHO; REJECT

xyzabcdabc
| abqabcd

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

%%

<STRING> [^"]* { action; }

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

After BEGIN, state is active.

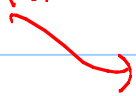
If state is inclusive, then rules with no start conditions are still active.

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

Ex: %S versus %x

%S example
%%

(INITIAL) <example> foo action();



bar other_action();
<otherex> action2();

%x example
%%

<example> foo action();

<INITIAL, example> bar other_action();
var action2();

Conditions

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

Essentially, pretend:

```
<*> ./In ECHO;
```

is a line of your file.

Ex: Scanner to ignore C comments
+ count of current input line

%X comment

%%

→ int num_line = 1; ← local var.

"/*" BEGIN(comment);

<comment> {

<comment> [^ * \n] *

<comment> "/*" + [^ * / \n]

<comment> \n ++line_num;

<comment> "/*" + "/" BEGIN(INITIAL);

Can condense

$\langle \text{comment} \rangle \{$
all rules

$\}$



$\langle \text{comment} \rangle$ rule 1
 $\langle \text{comment} \rangle$ rule 2
.
;