# CS344 - Context Free Languages

**Announcements**

- HW due!

- Next HW up -
  due next Friday

## Last time : flex

A useful scanner.

Based on reg expressions as well as states.

(Examples + links should be posted now.)

Automates building DFA code.

## That buggy example:

Well, I was just wrong.

REJECT actually is also for overlapping
so was grabbing subwords:

So:    try ⎵

        ry ⎵          } +3 to wc

        y ⌐⌐

## Another example:

```
%%

pink        { npink++; REJECT; }
ink         { nink++; REJECT; }
pin         { npin++; REJECT; }
.|n ;       /* discard others */
```

So — to fix my word count program:

Avoid REJECT!

Alternate — see .lex file.

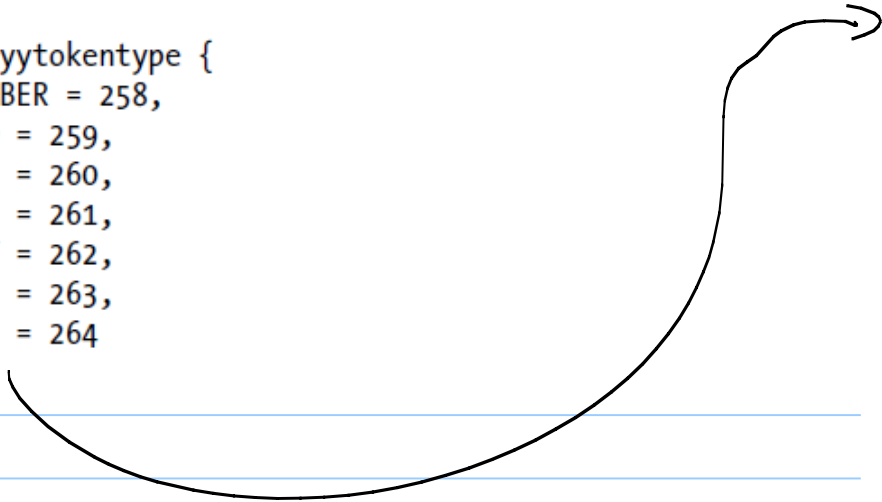key: add len(word) to charcount

(yylen)

# Flex: a tokenizer for a }; calculator:

```
/* recognize tokens for the calculator and print them out */      int yylval;
%{                                                                %}

    enum yytokentype {                                            %%
        NUMBER = 258,                                             "+"    { return ADD; }
        ADD = 259,                                                "-"    { return SUB; }
        SUB = 260,                                                "*"    { return MUL; }
        MUL = 261,                                                "/"    { return DIV; }
        DIV = 262,                                                "|"    { return ABS; }
        ABS = 263,                                                [0-9]+ { yylval = atoi(yytext); return NUMBER; }
        EOL = 264                                                 \n     { return EOL; }
                                                                  [ \t]  { /* ignore whitespace */ }
                                                                  .      { printf("Mystery character %c\n", *yytext); }
                                                                  %%
                                                                  main(int argc, char **argv)
                                                                  {
                                                                    int tok;

                                                                    while(tok = yylex()) {
                                                                      printf("%d", tok);
                                                                      if(tok == NUMBER) printf(" = %d\n", yylval);
                                                                      else printf("\n");
                                                                    }
                                                                  }
```

↰ tokens to
hand to
parser

# In Action:

```
$ flex fb1-4.l
$ cc lex.yy.c -lfl
$ ./a.out
a / 34 + |45
Mystery character a
262
258 = 34
259
263
258 = 45
264
```

# Bison accepts these tokens:

```
/* simplest version of calculator */
%{
#include <stdio.h>
%}

/* declare tokens */
%token NUMBER
%token ADD SUB MUL DIV ABS
%token EOL

%%

calclist: /* nothing */
 | calclist exp EOL { printf("= %d\n", $1); }
 ;

exp: factor          default $$ = $1
 | exp ADD factor { $$ = $1 + $3; }
 | exp SUB factor { $$ = $1 - $3; }
 ;

factor: term          default $$ = $1
 | factor MUL term { $$ = $1 * $3; }
 | factor DIV term { $$ = $1 / $3; }
 ;
```

**CFG**

```
term: NUMBER    default $$ = $1
 | ABS term     { $$ = $2 >= 0? $2 : - $2; }
 ;
%%
main(int argc, char **argv)
{
  yyparse();
}

yyerror(char *s)
{
  fprintf(stderr, "error: %s\n", s);
}
```

# Building:

```
# part of the makefile
fb1-5:  fb1-5.l fb1-5.y
        bison -d fb1-5.y
        flex fb1-5.l
        cc -o $@ fb1-5.tab.c lex.yy.c -lfl
```

# Running:

```
$ ./fb1-5
2 + 3 * 4
= 14
2 * 3 + 4
= 10
20 / 4 - 2
= 3
20 - 4 / 2
= 18
```

# Back to what Bison is:

```
exp: factor            default $$ = $1
  | exp ADD factor { $$ = $1 + $3; }
  | exp SUB factor { $$ = $1 - $3; }
  ;

factor: term            default $$ = $1
  | factor MUL term { $$ = $1 * $3; }
  | factor DIV term { $$ = $1 / $3; }
  ;
```

nice

Essentially, this is a CFG!
But only works on a particular
type of grammar.

# Context-Free Languages

Recall that for any context free languages, there are an infinite # of grammars that can produce it.

We wish to somehow give a definition of a "good" set of productions.

Goal: Parsing (well) —
given a language, detect if a string is in that language.

Ex: (BAD)

$S_0 \rightarrow S \mid X \mid Z$

$S \rightarrow A$

$A \rightarrow B$ → useless - stuck

$C \rightarrow Aa$

$X \rightarrow C$

unreachable
→ useless rule

$(Y) \rightarrow aY \mid a$

$Z \rightarrow \varepsilon$

capitols → non-terminals
lowercase terminals

chain
$S \rightarrow X \rightarrow C \rightarrow Aa$
$\quad\quad\quad\quad\quad\quad\quad\downarrow$
$\quad\quad\quad\quad\quad\quad B \ldots$

# Chomsky Normal Forms (CNF)

Each rule in the grammar is either:

- $A \longrightarrow BC$
  where neither B or C is the start variable, & both are nonterminals

- $A \longrightarrow a$
  where a is a terminal

- $S \longrightarrow \varepsilon$
  where S is the start symbol

**Thm**: All grammars can be converted to CNF.

**Procedure**:

① Create a new start symbol, $S_0$, & send $S_0 \rightarrow S$

(might need $S_0 \rightarrow \varepsilon$)

(1a) Eliminate useless rules
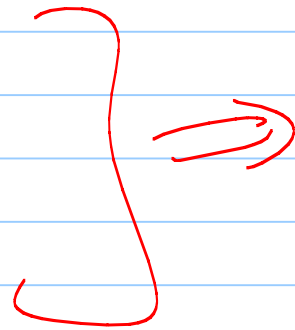
(just delete ones that can't be reached)

② Remove nullable variables.
$$A \to \varepsilon$$

How?
Remove all $\varepsilon$ productions.
Then fix.

Ex: $A \to CBC$
$\vdots$
$B \to \varepsilon \mid b$

$\Rightarrow$

$A \to CBC \mid CC$

$B \to b$

③ Remove unit rules:

$$S \to A$$

How? Must have:
$$X \to Z_1, \; Z_1 \to Z_2, \; \dots, \; Z_k \to Y$$
<span style="color:red">$(X, Y)$ is a unit pair</span>
(since we removed $\varepsilon$-transitions in ②)

Then:

<span style="color:red">add $X \to Y$</span>

<span style="color:red">but if $Y \to$ non-term!</span>

For each unit pair $(A, B)$
and rule $B \to w$,

add $A \to w$ to a <u>new</u>
grammar.

(Note that $(A, A)$ is a unit pair,
so all rules $A \to w$
will stick around.)

$$X \rightarrow ABCDy$$

④ Get rid of "long" righthand
      sides.

4a: Create $V_c \rightarrow c$ for every
      character.

    Replace $c$ with $V_c$ everywhere.

    Now either

$$A \rightarrow CDEF$$

    or

$$V_c \rightarrow c.$$

To demo:

$>A \rightarrow ABx \mid \varepsilon$

$B \rightarrow By \mid \varepsilon$

add start $\quad S_0 \rightarrow A \mid \varepsilon \qquad$ (step 1)

$A \rightarrow \underline{ABx} \mid Bx \mid x \mid Ax$

$B \rightarrow By \mid y$

add dummy non-terms:

$$S_0 \rightarrow A \mid \varepsilon$$

$$\begin{bmatrix} A \rightarrow \underline{AB\times} \mid B\times \mid \times \mid A\times \\ B \rightarrow By \mid y \end{bmatrix}$$

New: $S_0 \rightarrow A \mid \varepsilon$

$V_x \rightarrow x$

$V_y \rightarrow y$

$$\begin{bmatrix} A \rightarrow CV_x \\ C \rightarrow AB \end{bmatrix}$$

$A \rightarrow ABV_x \mid BV_x \mid V_x \mid AV_x$

$B \rightarrow BV_y \mid V_y$

4b:     $A \longrightarrow B_1 B_2 B_3 \ldots B_k$

How to replace with only 2 nonterminals on the right?

$A \rightarrow B_1 X_1$

$X_1 \rightarrow B_2 X_2$

$X_2 \rightarrow B_3 X_3$

$\vdots$

~~Unit Pars : $S_0 \rightarrow A$~~    CNF !

$S_0 \rightarrow \varepsilon \mid CV_x \mid BV_x \mid x \mid AU_x$

$V_x \rightarrow x$

$V_y \rightarrow y$

$A \rightarrow CV_x \mid BV_x \mid x \mid AU_x$

$B \rightarrow BU_y \mid y$

$C \rightarrow$

$\qquad CD \mid BD \mid V_x B \mid AD$

$D \rightarrow V_x B$

**Ex:** Convert:

$$S \to ASA \mid aB$$
$$A \to B \mid S$$
$$B \to b \mid \varepsilon$$

① $S_0 \to S$

$$S \to ASA \mid aB \mid a \mid AS \mid SA \mid S$$
$$A \to \cancel{\varepsilon} \mid S \mid b$$
$$B \to b$$

$$(S_0, S) \quad (A, S)$$

~~$S_0$~~

$$bSA$$
$$ASb$$
$$bSb$$

$$S \rightarrow AY_- \mid V_a B \mid a \mid AS \mid SA \mid \cancel{X} \mid XA)$$

$$A \rightarrow \cancel{X} \mid b \quad AX \mid XS \mid BY \mid$$

$$B \rightarrow b \quad\quad ZB) \; WB) \; SS)$$

$$BS \mid SB$$

$$X \rightarrow SS$$

$$Y \rightarrow SA \quad\quad V_a \rightarrow a$$

$$Z \rightarrow AS \quad\quad W \rightarrow BS$$

Now— why do we care??

Parsing: building those parse trees we saw

In general, there are an exponential number of parse trees for a given input.

So how to check quickly?

Even in CNF, might be $2^b$ possible parse trees.

# Cocke-Younger-Kasami (CYK) algorithm <span style="color:red">70's</span>

Uses a table & dynamic programming to give a parse tree in $O(n^3)$ time.

Grammar <u>must</u> be in CNF!

## Other options

- $n^3$ is still pretty slow.

$5 + 2$
num $\rightarrow$ 5

← size of my program

In general, can't really do better. However, certain <u>classes</u> could be done faster.

- LL(1), LR(1) } $O(n)$ algorithm

↑
lookahead