

CS3200 - Scanning & Parsing

Announcements

-flex HW - due next Sat.

Last week : Ch 2 of book

- Scanners and regular expressions

- DFAs & NFAs
- Flex

- Mentioned Parsing, which is the
next step
(more today)

Ex: Give the regular expression for $\{w \mid w \text{ begins with a } 1 \text{ and ends with a } 0\}$

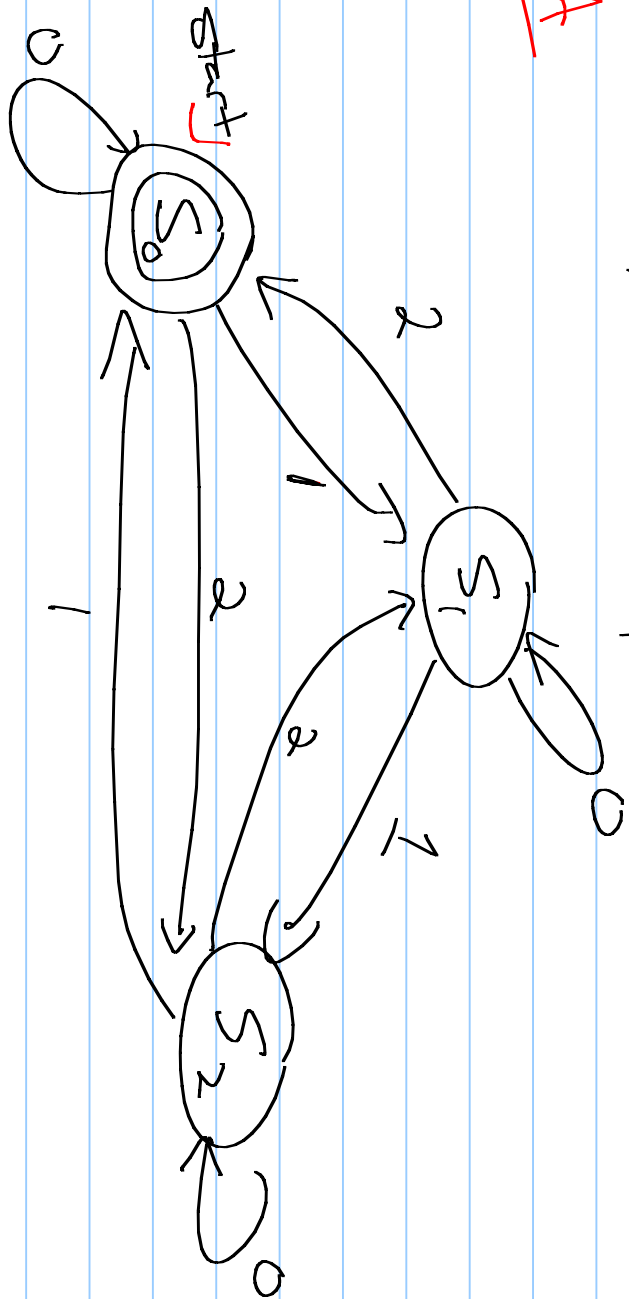
$1(0|1)^*0$

Ex: $\{w \mid w \text{ starts with } 0 \text{ and has an odd length}\}$

$0((0|1)(0|1))^*$

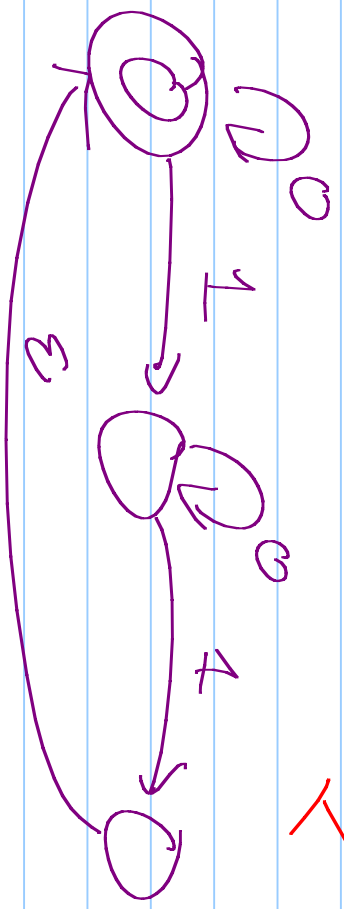
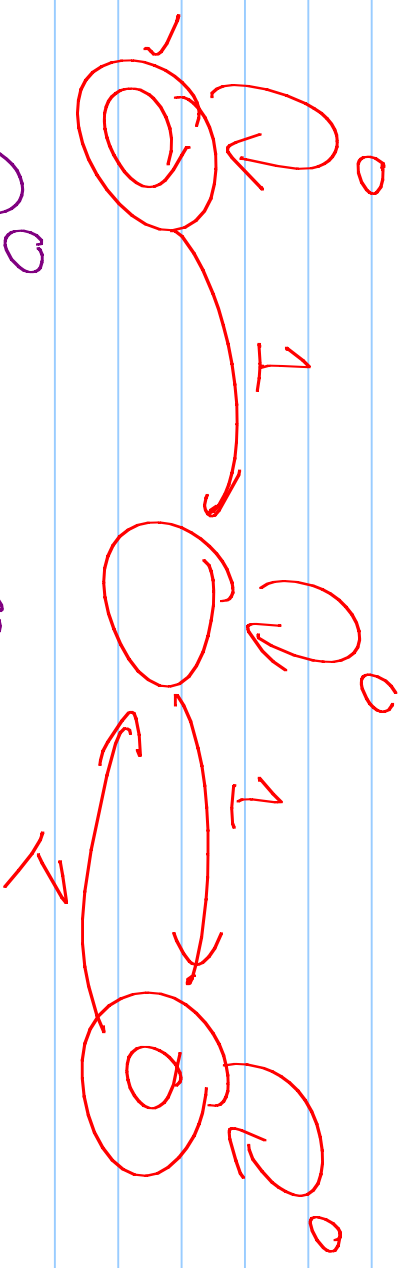
Ex: 3 symbol alphabet: $\{0,1,2\}$

DEFA

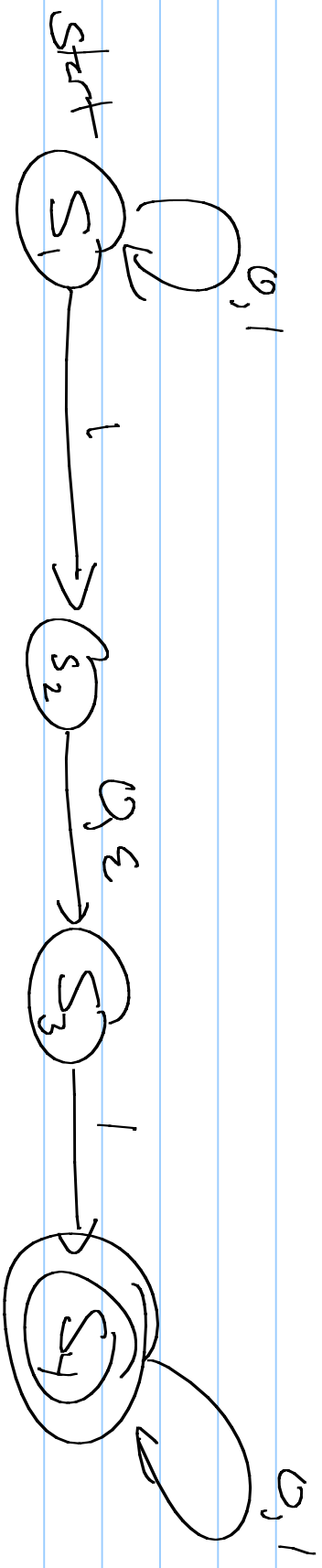


are you $\equiv 0 \pmod 3$?

Ex: Strings of 0's and 1's accepted if number of 1's is even

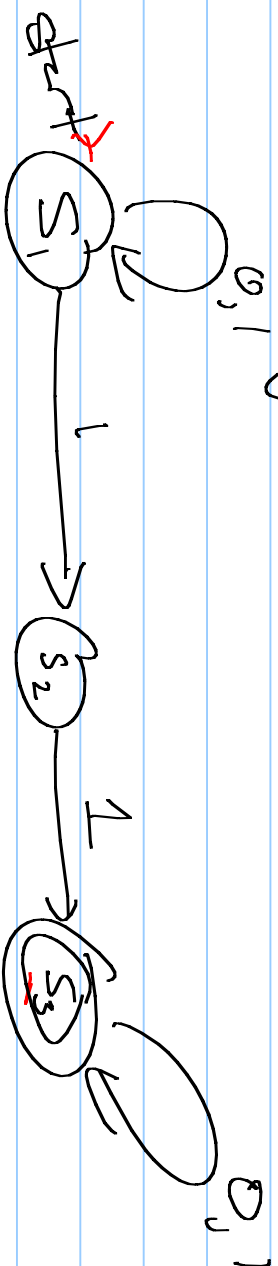


NFAs: DFAs w/ ambiguity

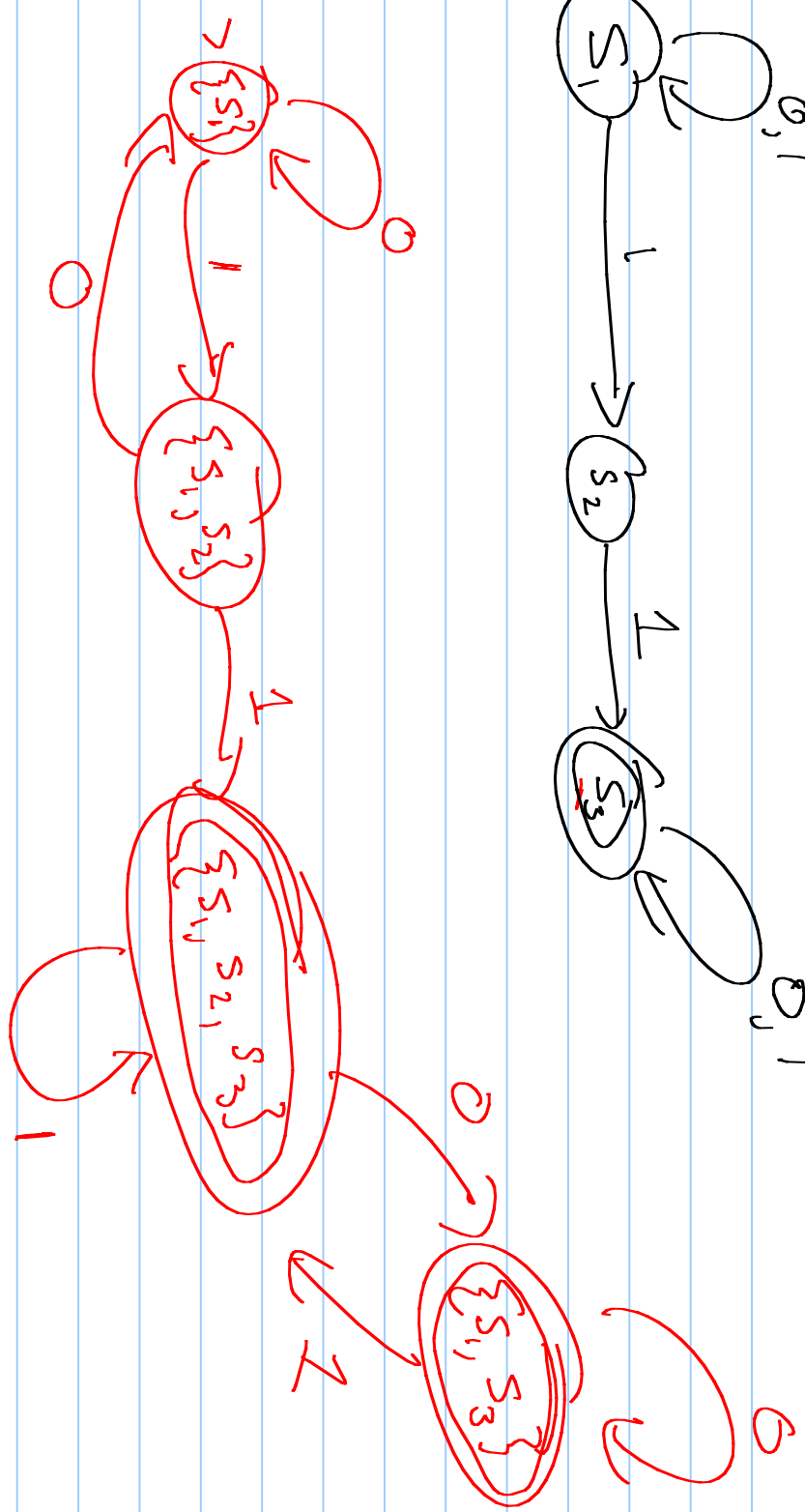


Thm: NFAs & DFAs are equivalent;
both accept reg. languages.

Converting NFAs to DFAs (p. 57)



DEFA:



Generating Scanners

Usually, easier to specify a regular expression.

But code needs to emulate a DFA - NFAs are too complex.

Unfortunately, there's not a great way to convert from reg exp \rightarrow DFA.

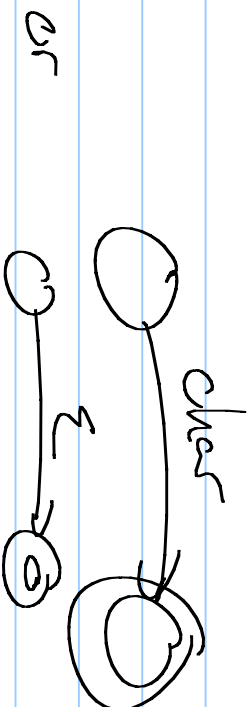
However, NFAs are easy!

\rightarrow

Constructing an NFA (p.57ish)

Given a regular expression, we can construct an NFA.

Simple NFA:

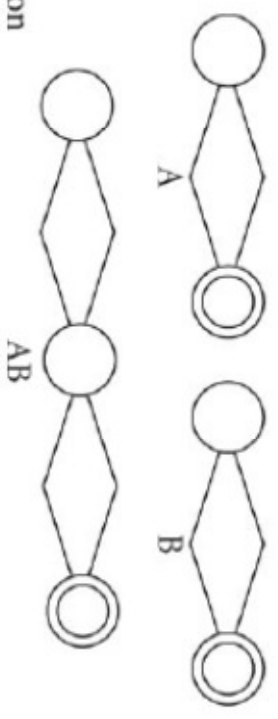


(Base case)

3 operations

Concatenation :

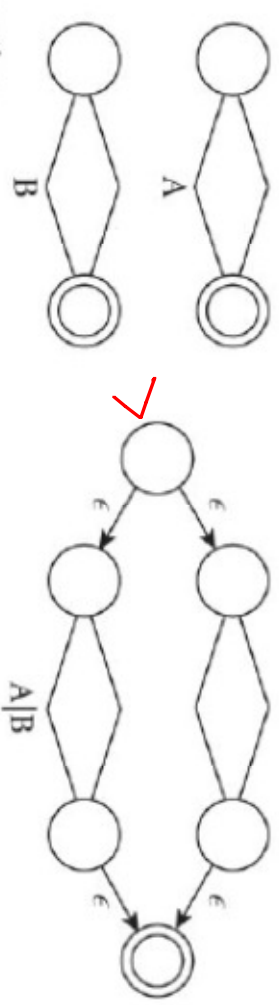
concatenation



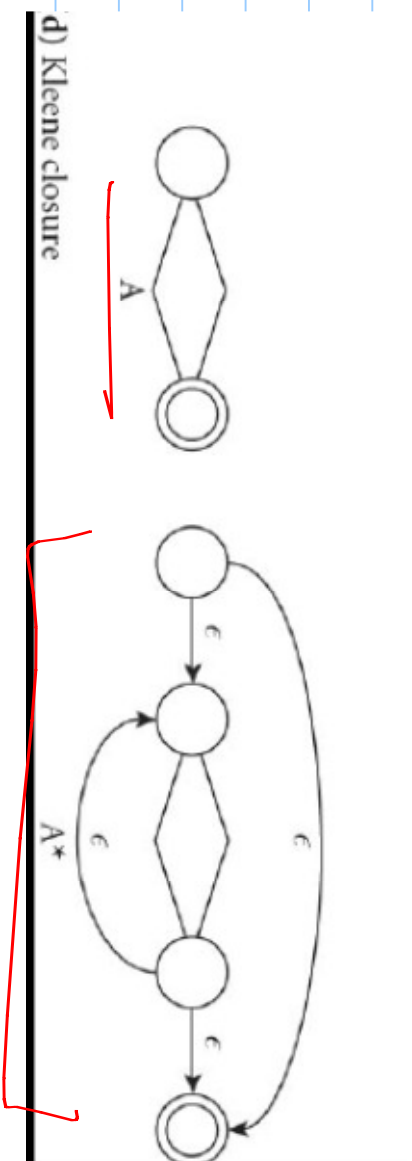
o

s

Q.1

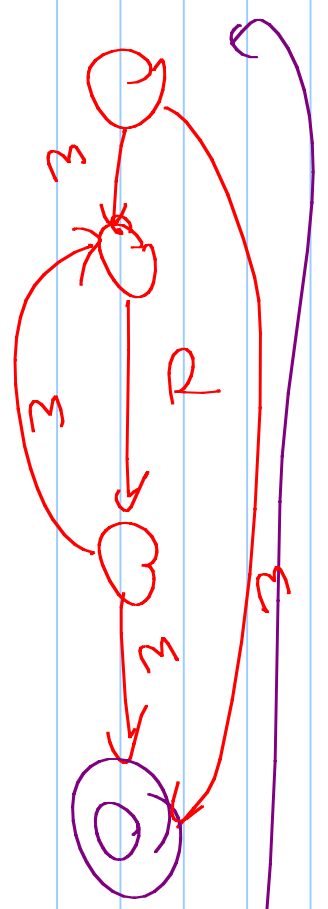
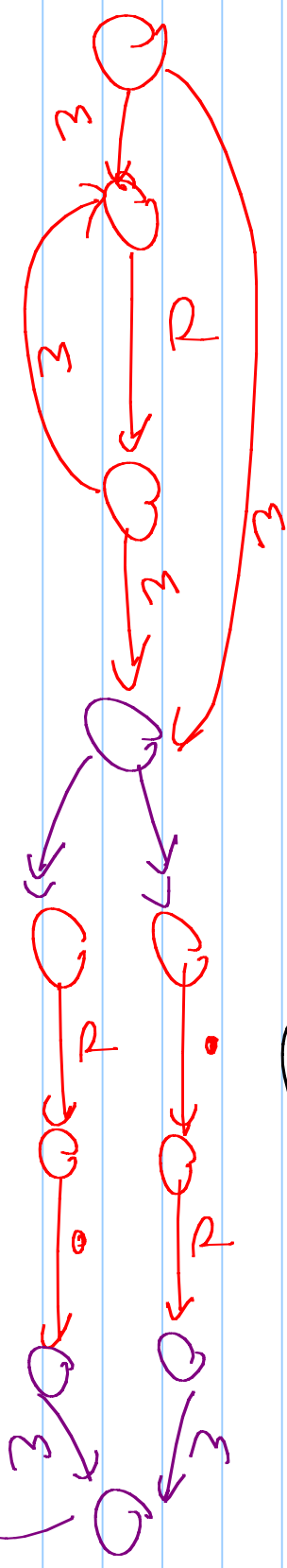


and Kleene closure ($*$).

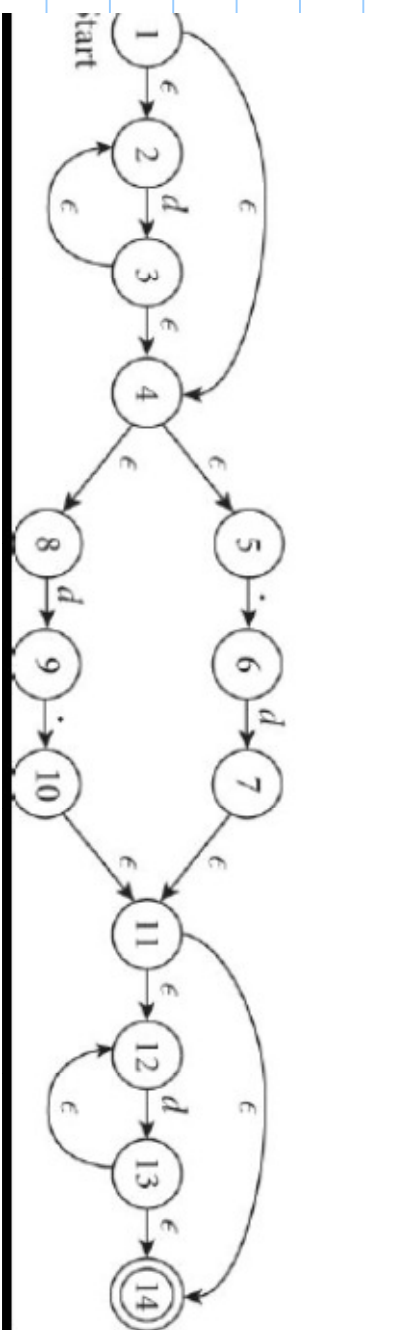


Example: decimals

$$d^* \left(\overbrace{. \overbrace{d \mid d \dots}^R}^R \right) d^*$$

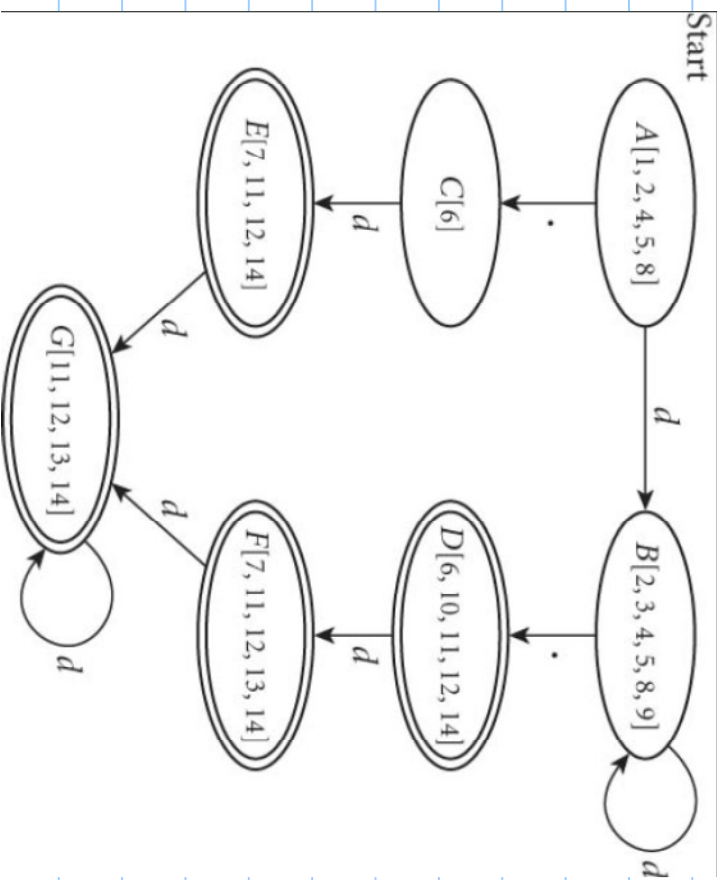


Final product:



Next: Convert to DFA. (lots of states, but same principle as we saw earlier.)

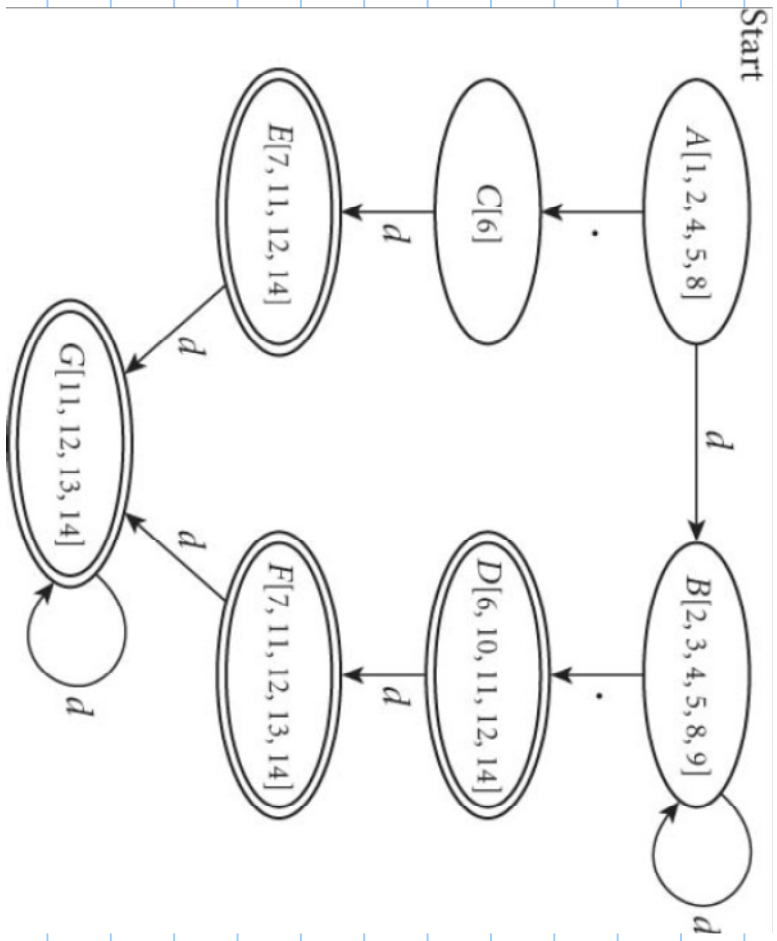
Result:
(see p. 57-58)



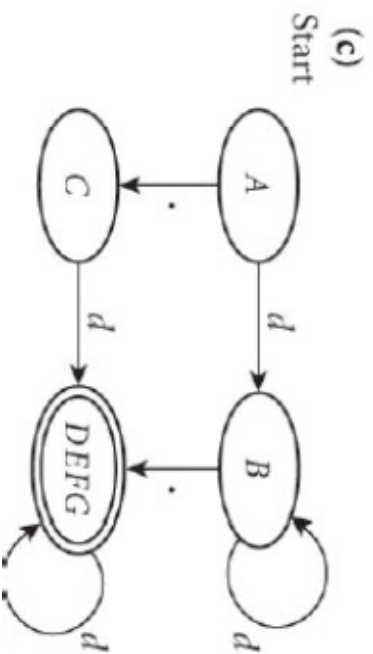
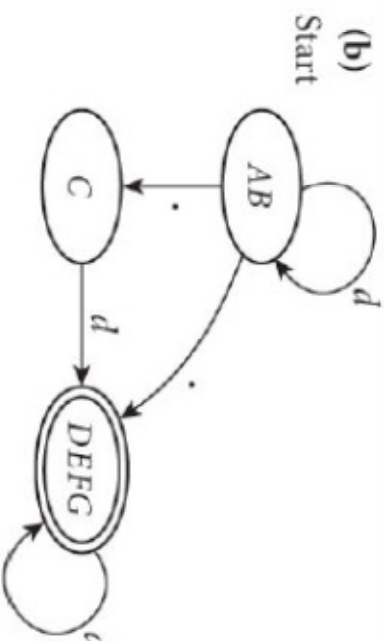
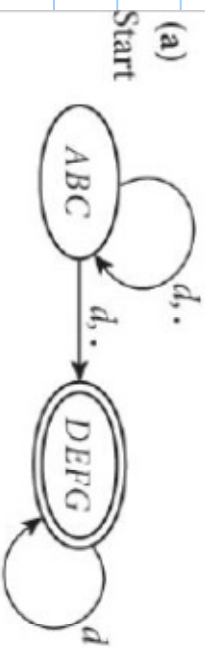
Note: This DFA is a bit redundant.

Not minimal.

Can easily find their equivalence classes and minimize.



Process to minimize (p. 59)



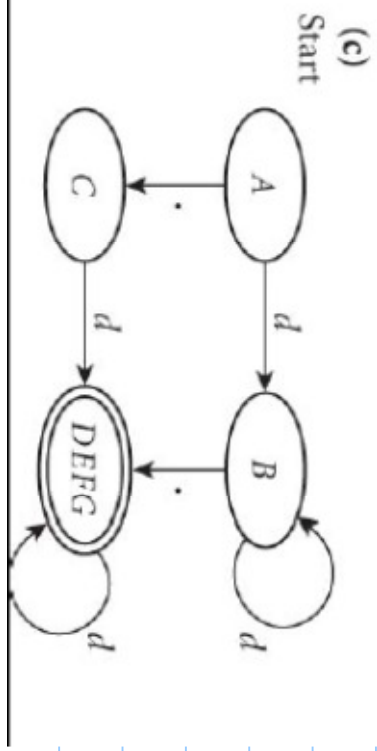
Now:

Given DFA, generate case statements to simulate it.

State = 1
repeat;
read curr_char
Case state is:

A: case curr_char = d
State = B
case curr_char = .
State = C

B:



So: Find result:

Automatic generation of a DFA
in some language given reg.
expression input.

(This is what flex is doing for you!)

Limitations of Regular Expressions

Certain languages are not regular.

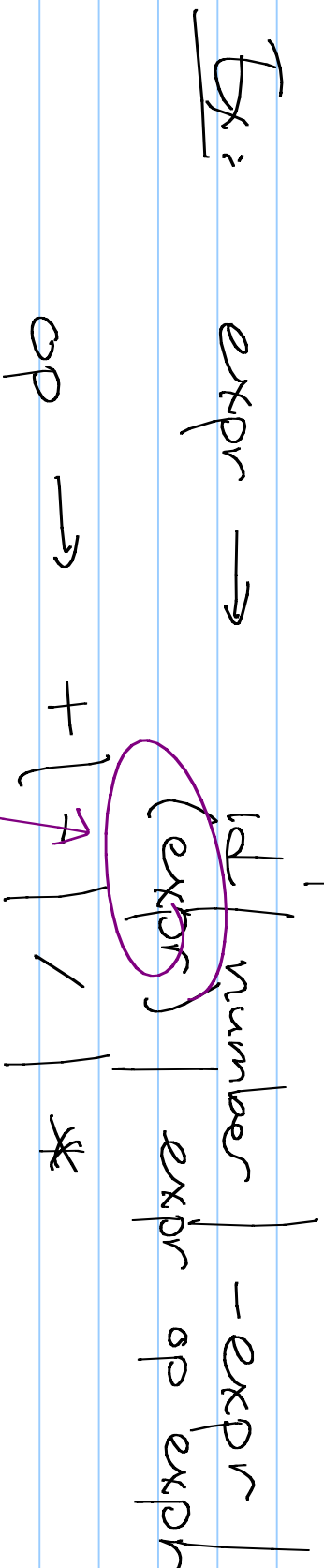
Ex: $\{w \mid w \text{ has an equal number of } 0\text{'s and } 1\text{'s}\}$

Somewhat, this needs a type of memory, which regular expressions do not have. \cup

Ex: $0^n 1^n$

Why do we need this?

Need to "nest" expressions.



Regular expressions can't quite do this.
non-regular

Context Free Languages

Described in terms of productions
(called Backus-Naur Form, or BNF)

- A set of terminals T
- A set of non-terminals N
- A start symbol S (a non-terminal)
- A set of productions

$$\text{Ex: } \{0^n 1^n \mid n \geq 0\} = L$$

non-terminals \rightarrow (capital letters)
 $S \rightarrow 0S1 \leftarrow$ terminals

$S \rightarrow 01$

Q: Is $0011 \in L$?

yes: $S \rightarrow 0S1 \rightarrow 0011 \checkmark$

Is $00011 \in L$? No

word + palindrome =

Ex: $\sum |w|$ has an equal number of 0's & 1's

S \rightarrow 0011

1010

S \rightarrow 1100

S \rightarrow 1100

S \rightarrow 1001

\rightarrow 100110

S \rightarrow 0110

\rightarrow 1010

Q: is this enough?

S \rightarrow 101

S \rightarrow 10

Expression grammars: Simple calculator

expr \rightarrow term / expr add_op term

term \rightarrow factor | term mult_op factor

factor \rightarrow id | number | -factor | (expr)

add_op \rightarrow + | -

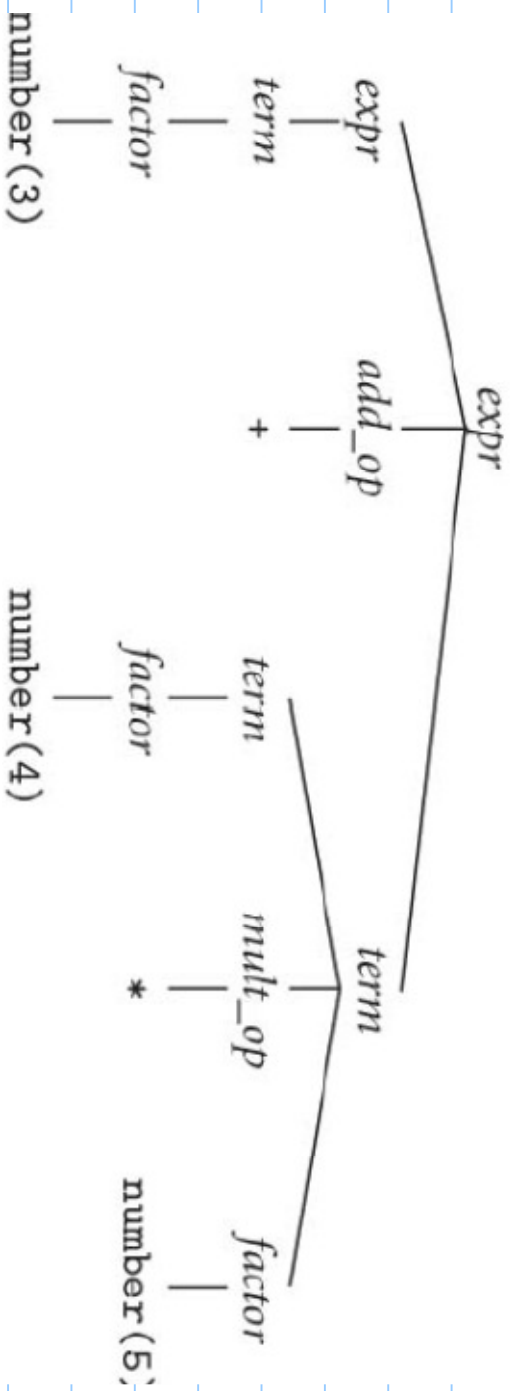
mult_op \rightarrow * | /

Example: Show how 3 rules can generate $3 + 2 + 5$

Parse Tree

Ex: $3 + 4 * 5$

derivation:



Parsers:

These are also tools to do this for you.

For flex, the complementary tool is lexon.

Example:

Considers a simple calculator.
Needs to accept valid numerical expressions, & then compute results.

Flex: a tokenizer for a calculator:

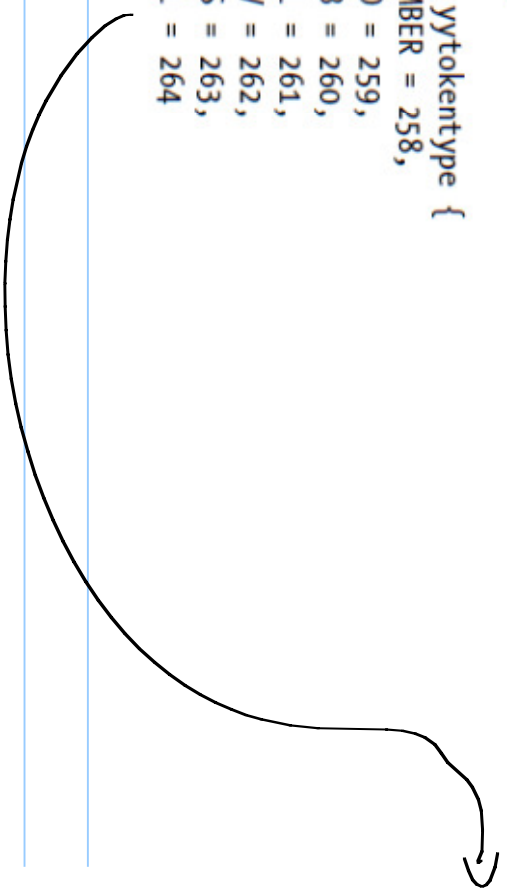
```
/* recognize tokens for the calculator and print them out */
%{
enum yytokentype {
    NUMBER = 258,
    ADD = 259,
    SUB = 260,
    MUL = 261,
    DIV = 262,
    ABS = 263,
    EOL = 264
};

int yy[1val];

%%
"+" { return ADD; }
"_" { return SUB; }
"*" { return MUL; }
"/" { return DIV; }
"|" { return ABS; }
[0-9]+ { yy[1val] = atoi(yytext); return NUMBER; }
\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", yy[1val]);
        else printf("\n");
    }
}
```



In Action:

```
$ flex fb1-4.1
$ 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; }
;
```

CFL

```
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: fb1-5.1 fb1-5.y
      bison -d fb1-5.y
      flex fb1-5.1
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; }
;
```

Essentially, this is a CFG ^{nice}!

But only works on a particular
type of grammar.

