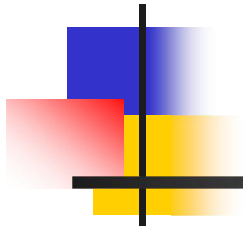


Programming Languages and Compilers (CS 421)



Choonghwan Lee

<http://www.cs.illinois.edu/class/cs421/>

Based in part on slides by Mattox Beckman, as updated by Vikram Adve, Gul Agha and Elsa Gunter



Meta-discourse

- n Language Syntax and Semantics

- n Syntax

- DFSAs and NDFSAs
- Grammars

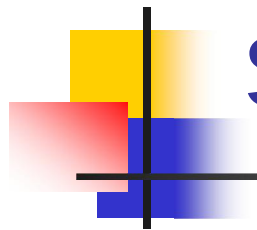
- n Semantics

- Natural Semantics
- Transition Semantics



Language Syntax

- n Syntax is the description of which strings of symbols are meaningful expressions in a language
- n It takes more than syntax to understand a language; need meaning (semantics) too
- n Syntax is the entry point



Syntax of English Language

n Pattern 1

Subject	Verb
<i>David</i>	<i>sings</i>
<i>The dog</i>	<i>barked</i>
<i>Susan</i>	<i>yawned</i>

n Pattern 2

Subject	Verb	Direct Object
<i>David</i>	<i>sings</i>	<i>ballads</i>
<i>The professor</i>	<i>wants</i>	<i>to retire</i>
<i>The jury</i>	<i>found</i>	<i>the defendant guilty</i>



Elements of Syntax

- n Character set – previously always ASCII, now often 64 character sets
- n Keywords – usually reserved
- n Special constants – cannot be assigned to
- n Identifiers – can be assigned to
- n Operator symbols
- n Delimiters (parenthesis, braces, brackets)
- n Blanks (aka white space)



Elements of Syntax

n Expressions

if ... then begin ... ; ... end else begin ... ; ... end

n Type expressions

typexpr₁ -> typexpr₂

n Declarations (in functional languages)

let *pattern₁* = *expr₁* in *expr*

n Statements (in imperative languages)

a = b + c

n Subprograms

let *pattern₁* = let rec inner = ... in *expr*



Elements of Syntax

- n Modules
- n Interfaces
- n Classes (for object-oriented languages)



Formal Language Descriptions

- n Regular expressions, regular grammars, finite state automata
- n Context-free grammars, BNF grammars, syntax diagrams
- n Whole family more of grammars and automata – covered in automata theory



Grammars

- n Grammars are formal descriptions of which strings over a given character set are in a particular language
- n Language designers write grammar
- n Language implementers use grammar to know what programs to accept
- n Language users use grammar to know how to write legitimate programs



Regular Expressions

- n Start with a given character set –
a, b, c...
- n Each character is a regular expression
 - n It represents the set of one string
containing just that character



Regular Expressions

n If **x** and **y** are regular expressions, then **xy** is a regular expression

n It represents the set of all strings made from first a string described by **x** then a string described by **y**

If **x**={a,ab} and **y**={c,d} then **xy**={ac,ad,abc,abd}.

n If **x** and **y** are regular expressions, then **x ∨ y** is a regular expression

n It represents the set of strings described by either **x** or **y**

If **x**={a,ab} and **y**={c,d} then **x ∨ y**={a,ab,c,d}



Regular Expressions

- n If **x** is a regular expression, then so is **(x)**
 - n It represents the same thing as **x**

- n If **x** is a regular expression, then so is **x***
 - n It represents strings made from concatenating zero or more strings from **x**
- If **x** = {a,ab}
 then **x*** = {"",a,ab,aa,aab,abab,aaa,aaab,...}

- n **ε**
 - n It represents {""}, set containing the empty string



Example Regular Expressions

n $(0 \vee 1)^* 1$

n The set of all strings of **0**'s and **1**'s ending in 1, $\{1, 01, 11, \dots\}$

n $a^* b(a^*)$

n The set of all strings of a's and b's with exactly one b

n $((01) \vee (10))^*$

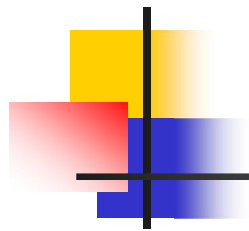
n You tell me

n Regular expressions (equivalently, regular grammars) important for lexing, breaking strings into recognized words



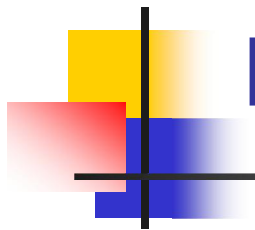
Example: Lexing

- n Regular expressions good for describing lexemes (words) in a programming language
 - n Identifier = $(a \vee b \vee \dots \vee z \vee A \vee B \vee \dots \vee Z) (a \vee b \vee \dots \vee z \vee A \vee B \vee \dots \vee Z \vee 0 \vee 1 \vee \dots \vee 9)^*$
 - n Digit = $(0 \vee 1 \vee \dots \vee 9)$
 - n Number = $0 \vee (1 \vee \dots \vee 9)(0 \vee \dots \vee 9)^* \vee \sim (1 \vee \dots \vee 9)(0 \vee \dots \vee 9)^*$
 - n Keywords: if = if, while = while,...



Implementing Regular Expressions

- n Regular expressions reasonable way to generate strings in language
- n Not so good for recognizing when a string is in language
- n Problems with Regular Expressions
 - n which option to choose,
 - n how many repetitions to make
- n Answer: finite state automata

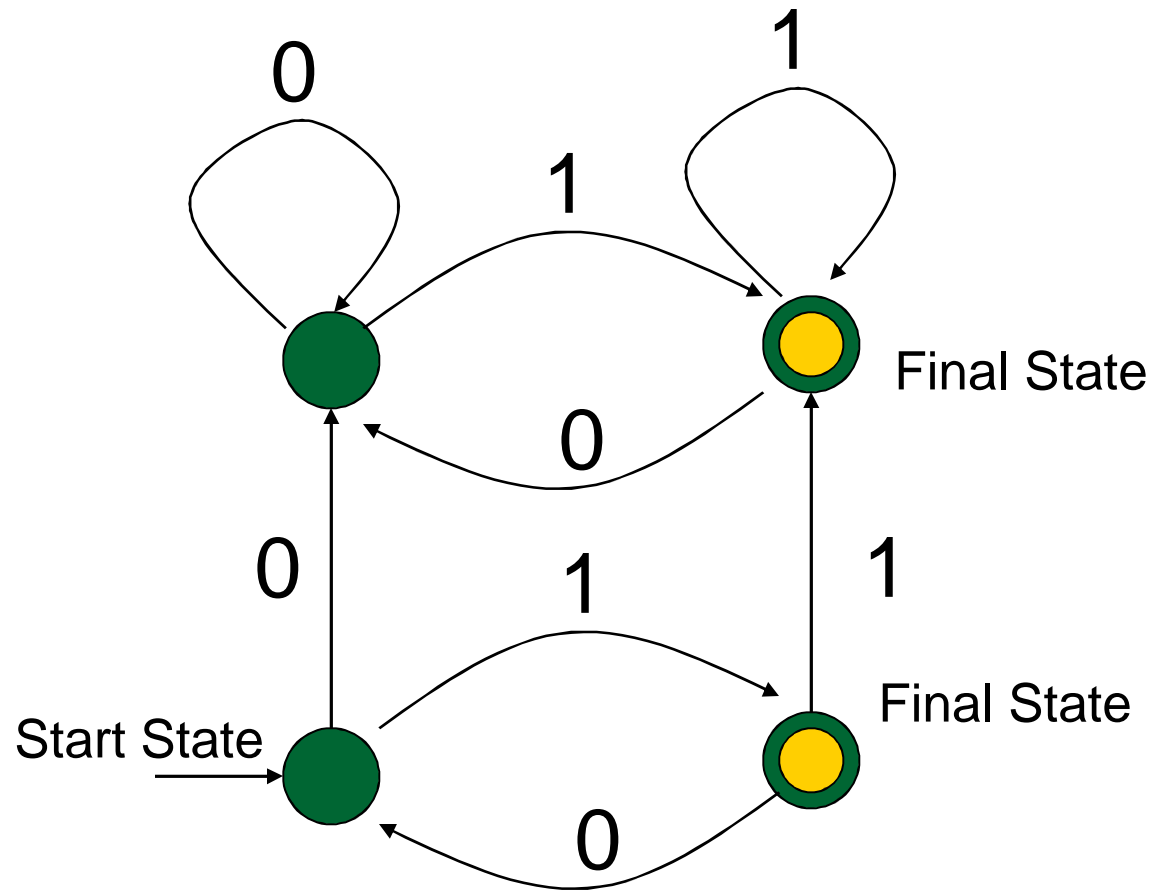


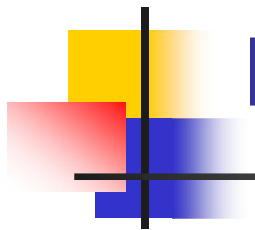
Finite State Automata

- n A finite state automata over an alphabet is:
 - n a directed graph
 - n a finite set of **states** defined by the **nodes**
 - n **edges** are labeled with elements of **alphabet**, or empty string; they define **state transition**
 - n some nodes (or *states*), marked as **final**
 - n one node marked as **start state**

n Syntax of FSA

Example FSA





Deterministic FSA's

- n If FSA has for every state *exactly one* edge for each letter in alphabet then FSA is *deterministic*
 - n No edge labeled with ϵ
- n In general FSA is *non-deterministic*.
 - n NFSA also allows edges labeled by ϵ
- n Deterministic FSA special kind of non-deterministic FSA



DFSA Language Recognition

- n Think of a DFSA as a board game; DFSA is board
- n You have string as a deck of cards; one letter on each card
- n Start by placing a disc on the start state



DFSA Language Recognition

- n Move the disc from one state to next along the edge labeled the same as top card in deck; discard top card
- n When you run out of cards,
 - n if you are in final state, you win; string is in language
 - n if you are not in a final state, you lose; string is not in language



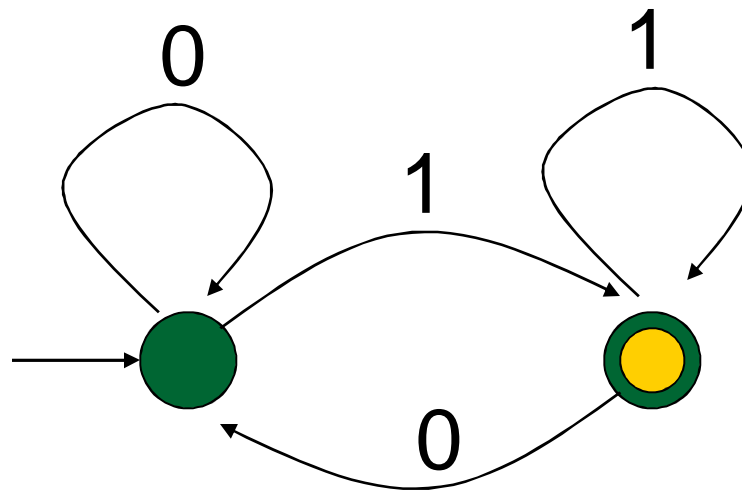
DFSA Language Recognition -Summary

- n Given a string over alphabet
- n Start at start state
- n Move over edge labeled with first letter to new state
- n Remove first letter from string
- n Repeat until string gone
- n If end in final state then string in language

- n Semantics of FSA

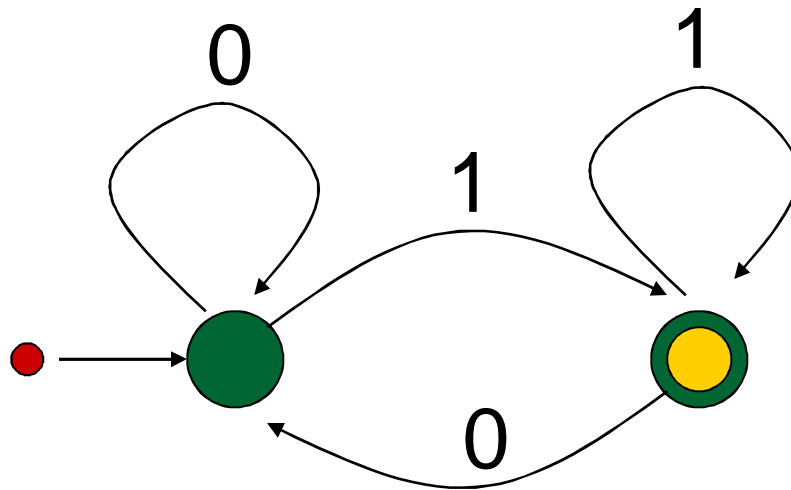
Example DFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Deterministic FSA



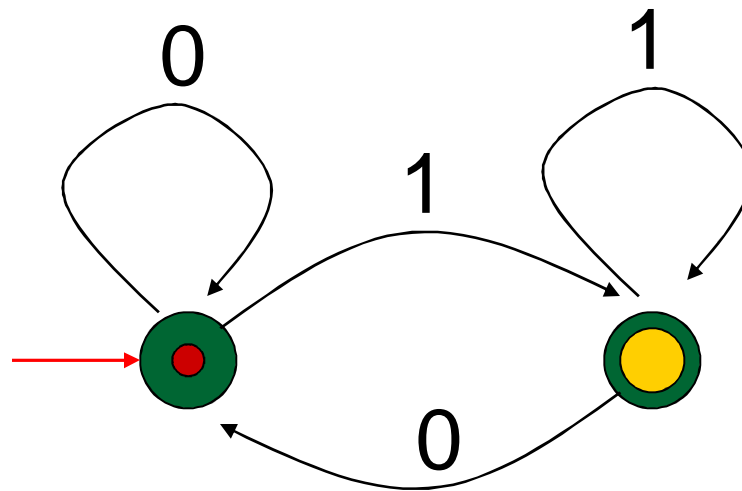
Example DFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string 0 1 1 0 1



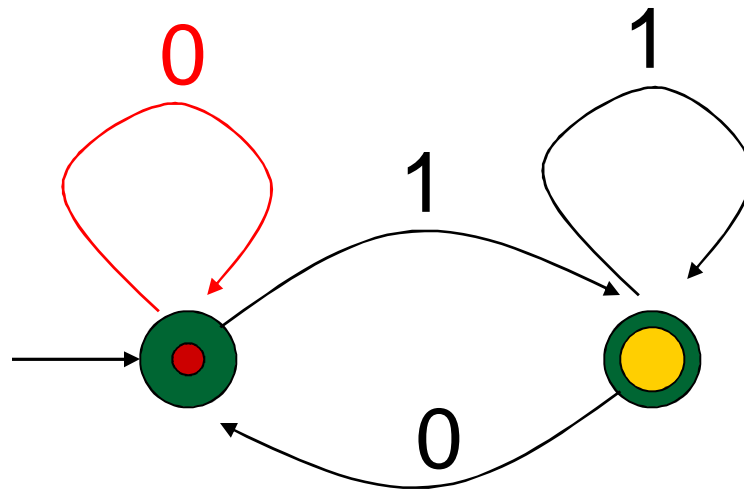
Example DFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string 0 1 1 0 1



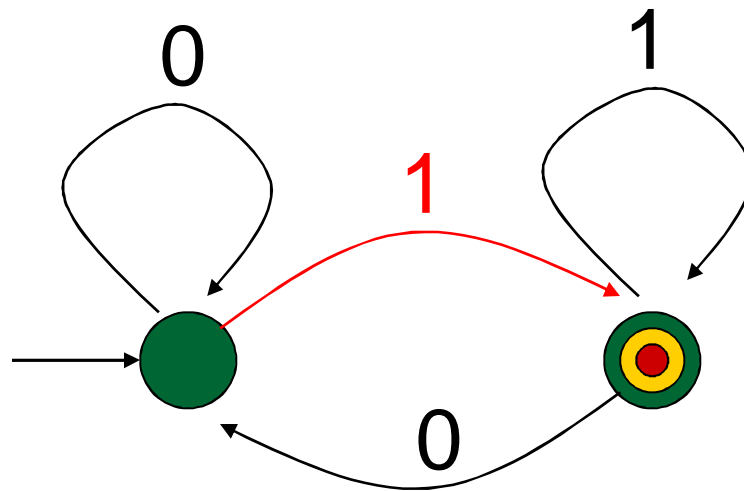
Example DFSA

- Regular expression: $(0 \vee 1)^* 1$
- Accepts string ~~0~~ 1 1 0 1



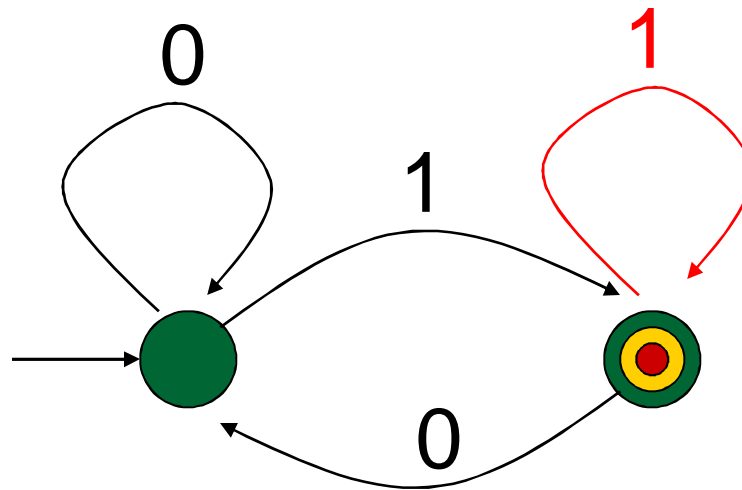
Example DFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ 1 1 0 1



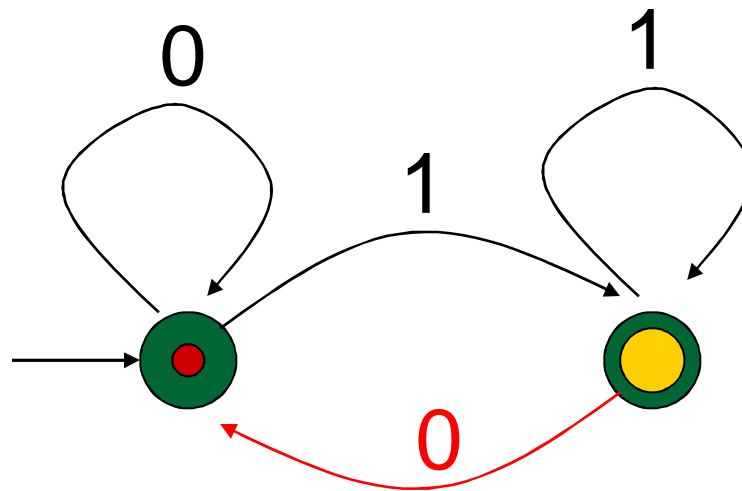
Example DFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ ~~1~~ ~~1~~ 0 1



Example DFSA

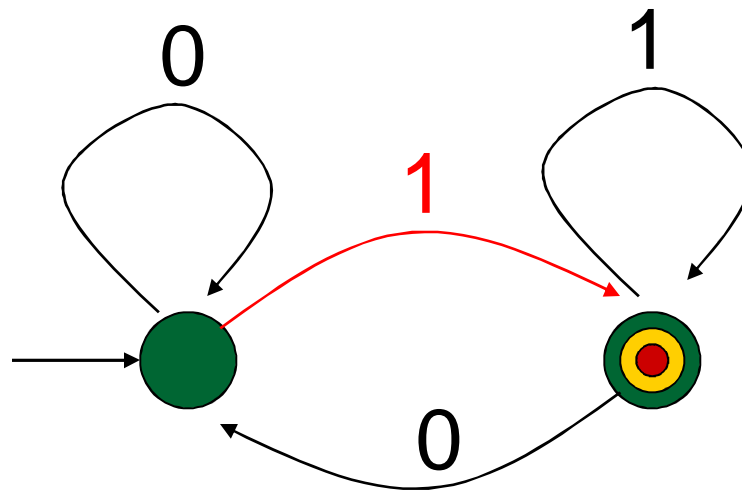
- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ ~~1~~ ~~1~~ ~~0~~ 1



Example DFSA

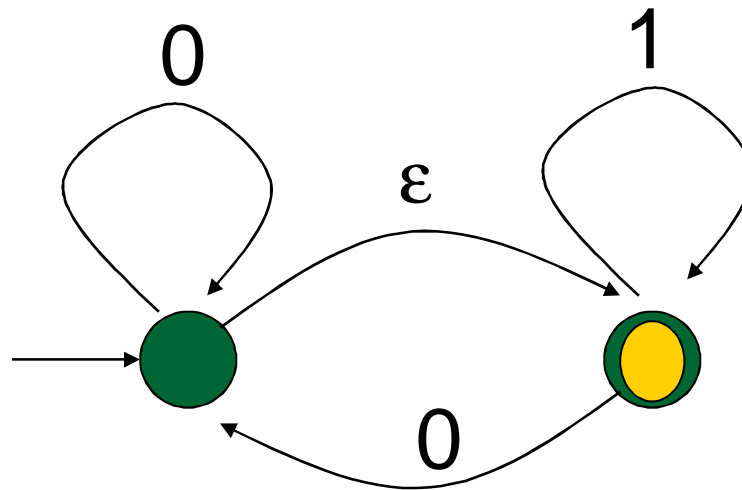
n Regular expression: $(0 \vee 1)^* 1$

n Accepts string ~~0~~~~1~~~~1~~~~0~~~~1~~



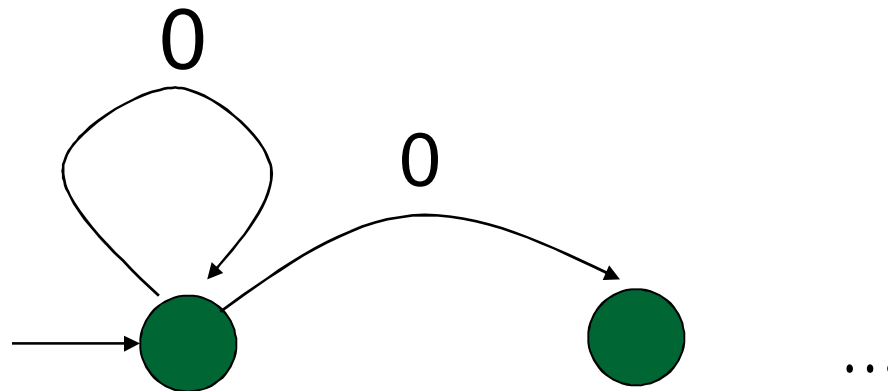
Non-deterministic FSA's

- n NFSA generalize DFSA in two ways:
- n Include edges labeled by ϵ
 - n Allows process to non-deterministically change state



Non-deterministic FSA's

- n Each state can have zero, one or more edges labeled by each letter
 - n Given a letter, non-deterministically choose an edge to use



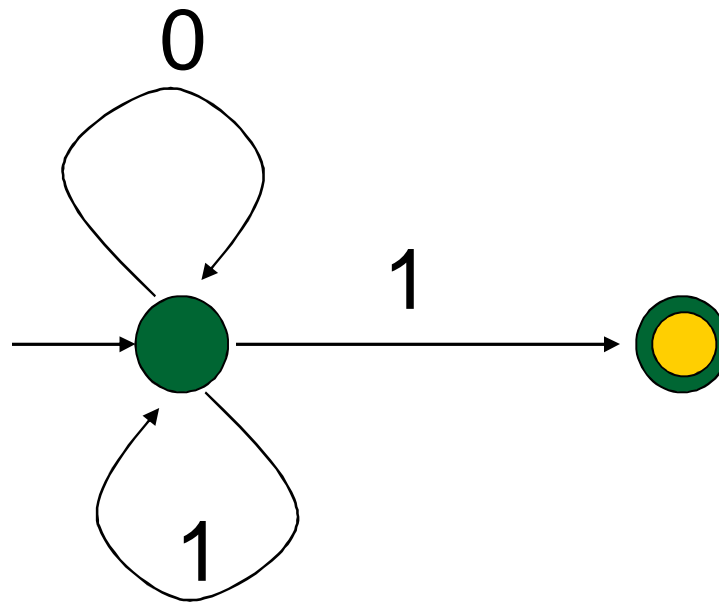


NFSA Language Recognition

- n Play the same game as with DFSA
- n Free move: move across an edge with empty string label without discarding card
- n When you run out of letters, if you are in final state, you win; string is in language
- n You can take one or more moves back and try again
- n If have tried all possible paths without success, then you lose; string not in language

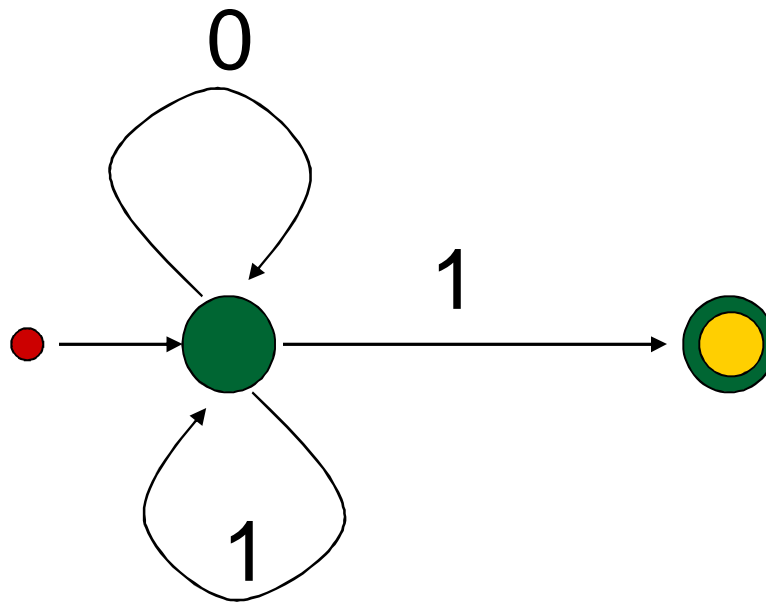
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Non-deterministic FSA



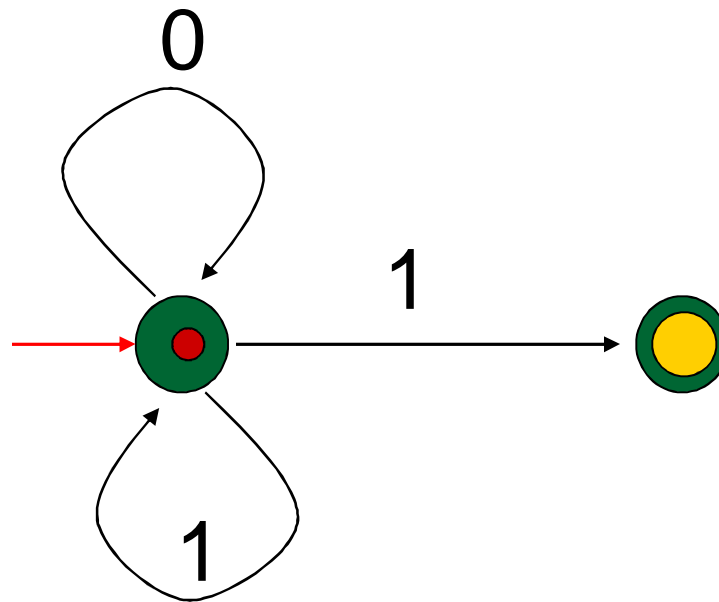
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string 0 1 1 0 1



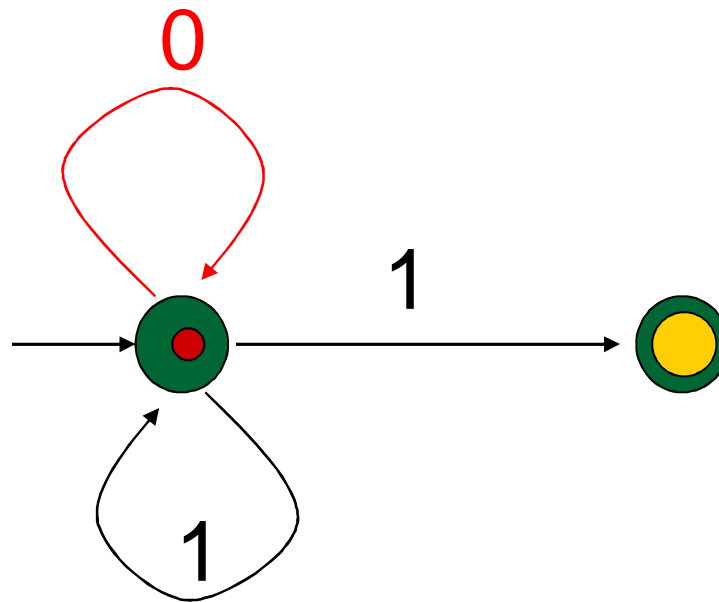
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string 0 1 1 0 1



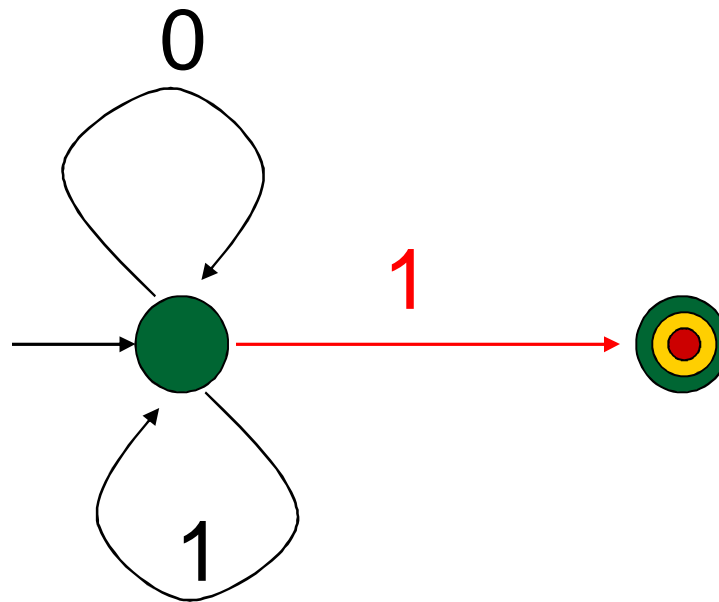
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ 1 1 0 1



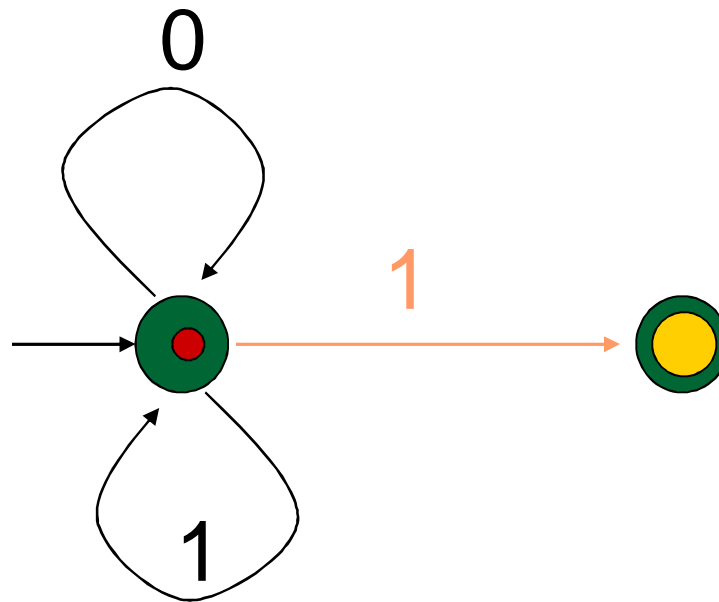
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ 1 1 0 1
- n Guess



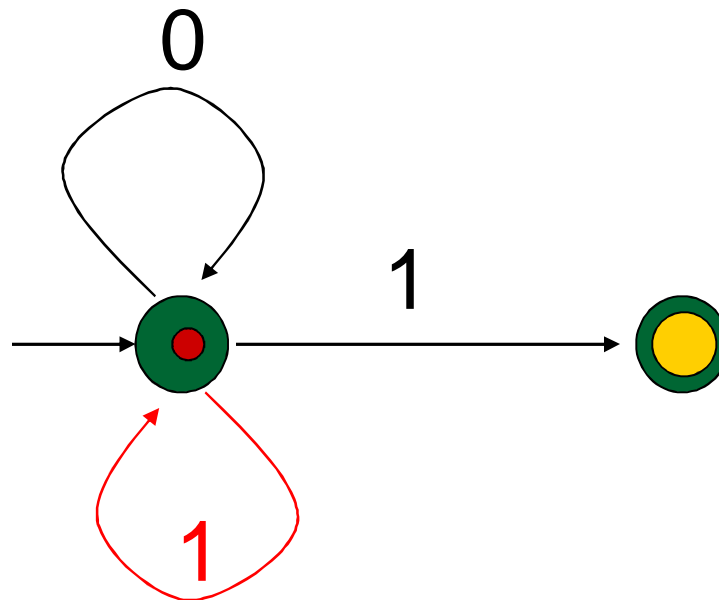
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string 0 1 1 0 1
- n Backtrack



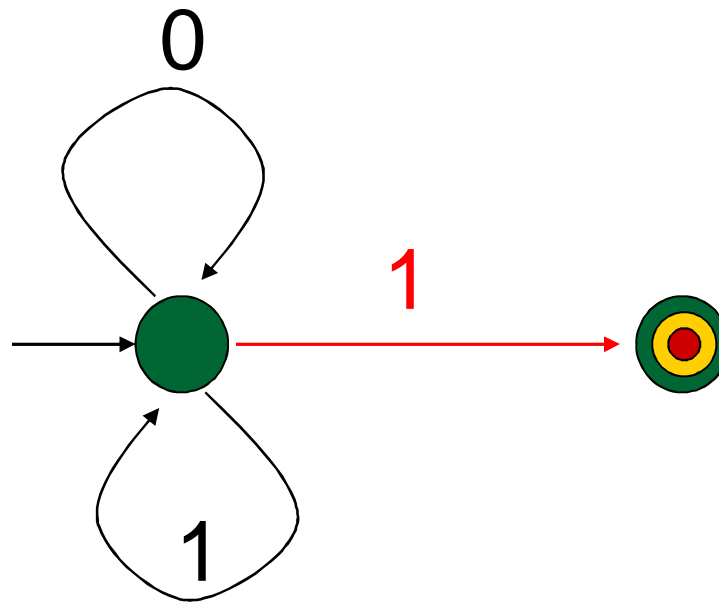
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ 1 1 0 1
- n Guess again



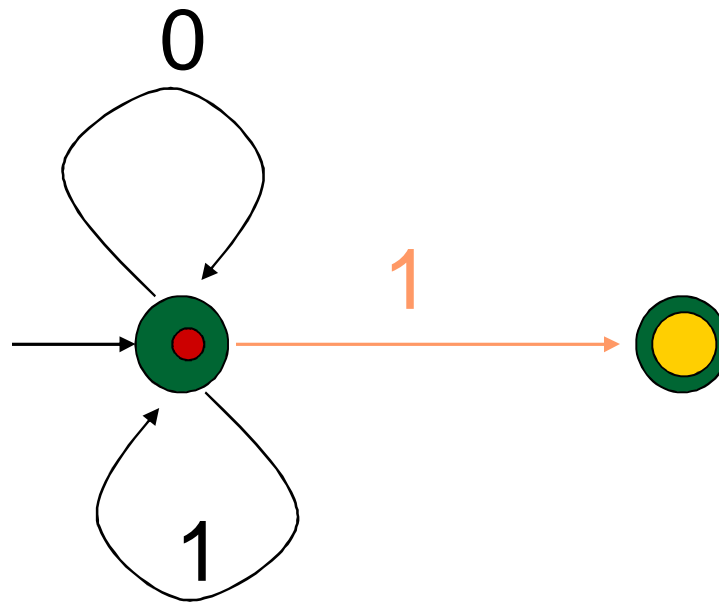
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ ~~1~~ ~~1~~ 0 1
- n Guess



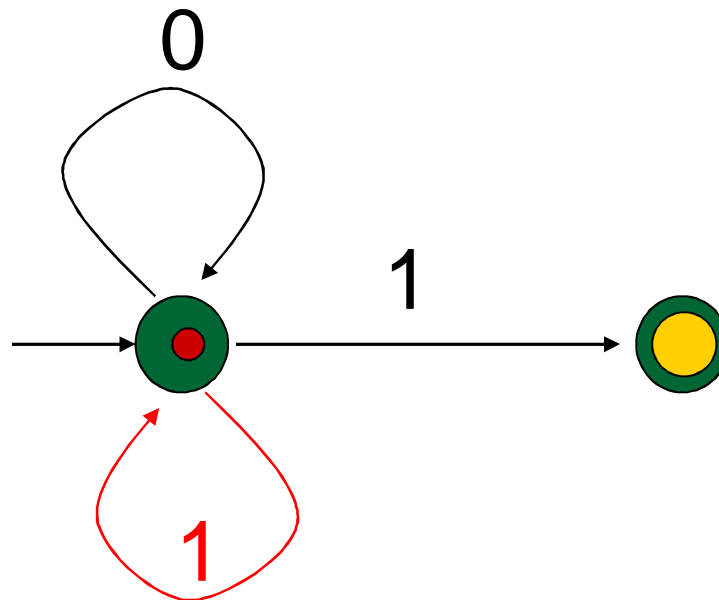
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ 1 1 0 1
- n Backtrack



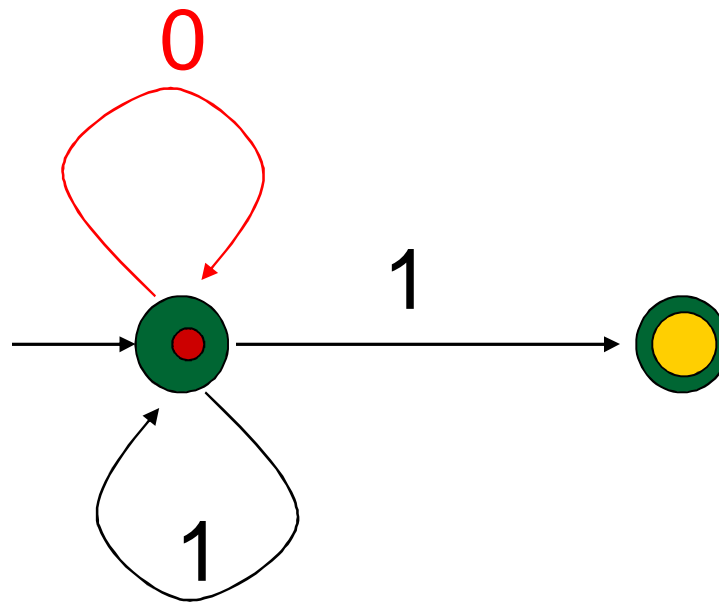
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ ~~1~~ ~~1~~ 0 1
- n Guess again



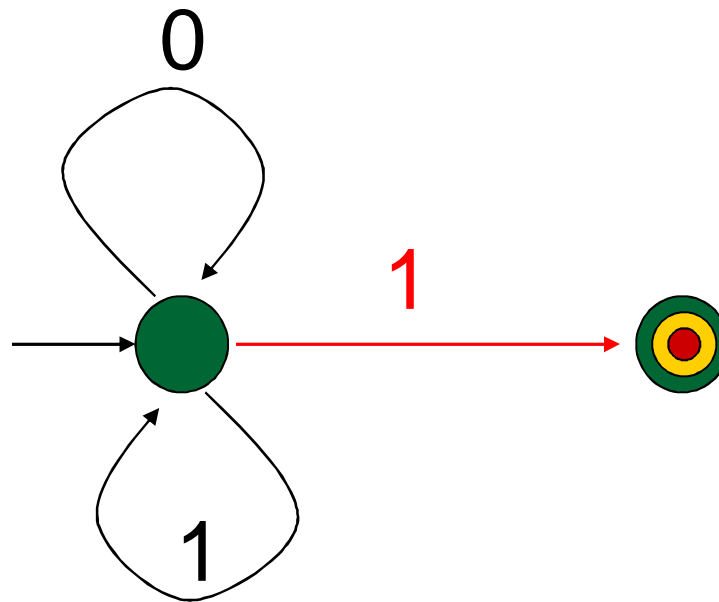
Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~ ~~1~~ ~~1~~ ~~0~~ 1



Example NFSA

- n Regular expression: $(0 \vee 1)^* 1$
- n Accepts string ~~0~~~~1~~~~1~~~~0~~~~1~~
- n Guess (Hurray!!)





Rule Based Execution

- n Search
- n When stuck backtrack to last point with choices remaining
- n Executing the NFSA in last example was example of rule based execution
- n FSA's are rule-based programs; transitions between states (labeled edges) are rules; set of all FSA's is programming language



Rule Based Execution

- n Search
- n When stuck backtrack to last point with choices remaining
- n FSA's are rule-based programs; transitions between states (labeled edges) are rules; set of all FSA's is programming language



Where We Are Going

- n We want to turn strings (code) into computer instructions
- n Done in phases
- n Turn strings into abstract syntax trees (parse)
- n Translate abstract syntax trees into executable instructions (interpret or compile)



Lexing and Parsing

- n Converting strings to abstract syntax trees done in two phases
 - n **Lexing:** Converting string (or streams of characters) into lists (or streams) of tokens (the “words” of the language)
 - n **Parsing:** Convert a list of tokens into an abstract syntax tree



Lexing

- n Different syntactic categories of "words": tokens

Example:

- n Convert sequence of characters into sequence of strings, integers, and floating point numbers.
- n "asd 123 jkl 3.14" will become:
[String "asd"; Int 123; String "jkl"; Float 3.14]



Lexing

- n Each category described by regular expression (with extended syntax)
- n Words recognized by (encoding of) corresponding finite state automaton
- n Problem: we want to pull words out of a string; not just recognize a single word

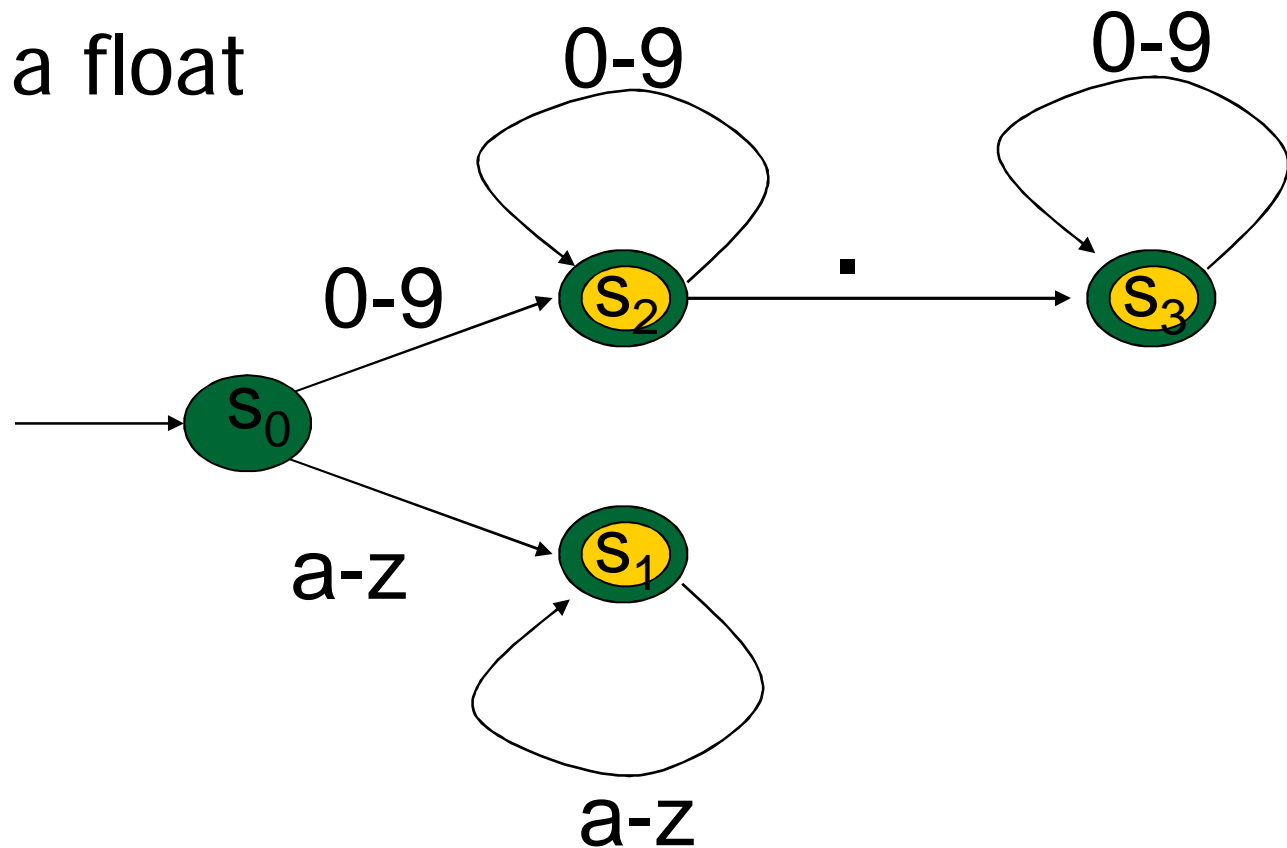


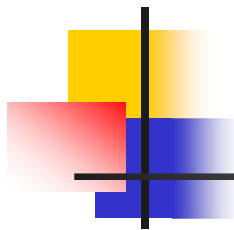
Lexing

- n Modify behavior of DFA
- n When we encounter a character in a state for which there is no transaction
 - n Stop processing the string
 - n If in an accepting state, return the token that corresponds to the state, and the remainder of the string
 - n If not, fail
- n Add recursive layer to get sequence

Example

- s_1 return a string
- s_2 return an integer
- s_3 return a float





Lex, ocamllex

- n Could write the reg exp, then translate to DFA by hand
 - n A lot of work
- n Better: Write program to take reg exp as input and automatically generates automata
- n Lex is such a program
- n ocamllex version for ocaml



How to do it

- n To use regular expressions to parse our input we need:
 - n Some way to identify the input string — call it a lexing buffer
 - n Set of regular expressions,
 - n Corresponding set of actions to take when they are matched.



How to do it

- n The lexer will take the regular expressions and generate a state machine.
- n The state machine will take our lexing buffer and apply the transitions...
- n If we reach an accept state from which we can go no further, the machine will perform the appropriate action.



Mechanics

- n Put table of reg exp and corresponding actions (written in ocaml) into a file *<filename>.ml*
- n Call

`ocamllex <filename>.ml`
- n Produces Ocaml code for a lexical analyzer in file *<filename>.ml*



Sample Input

```
rule main = parse
  ['0'-'9']+ { print_string "Int\n"}
  | ['0'-'9']+ '.' ['0'-'9']+ { print_string "Float\n"}
  | ['a'-'z']+ { print_string "String\n"}
  | _ { main lexbuf }
{
  let newlexbuf = (Lexing.from_channel stdin) in
  print_string "Ready to lex.\n";
  main newlexbuf
}
```



General Input

{ *header* }

let *ident* = *regexp* ...

rule *entrypoint* [*arg1*... *argn*] = parse
 regexp { *action* }

| ...

| *regexp* { *action* }

and *entrypoint* [*arg1*... *argn*] = parse ...and

...

{ *trailer* }



Ocamllex Input

- n *header* and *trailer* contain arbitrary ocaml code put at top and bottom of *<filename>.ml*
- n *let ident = regexp ...* Introduces *ident* for use in later regular expressions



Ocamlex Input

- n *<filename>.ml* contains one lexing function per *entrypoint*
 - n Name of function is name given for *entrypoint*
 - n Each entry point becomes an Ocaml function that takes $n+1$ arguments, the extra implicit last argument being of type `Lexing.lexbuf`
- n *arg1... argn* are for use in *action*



Ocamllex Regular Expression

- n Single quoted characters for letters:
'a'
- n *_*: (underscore) matches any letter
- n *Eof*: special "end_of_file" marker
- n Concatenation same as usual
- n *"string"*: concatenation of sequence of characters
- n *e_1 / e_2* : choice - what was *$e_1 \vee e_2$*



Ocamllex Regular Expression

- n $[c_1 - c_2]$: choice of any character between first and second inclusive, as determined by character codes
- n $[^c_1 - c_2]$: choice of any character NOT in set
- n e^* : same as before
- n $e+$: same as $e e^*$
- n $e?$: option - was $e_1 \vee \varepsilon$



Ocamllex Regular Expression

- n $e_1 \# e_2$: the characters in e_1 but not in e_2 ; e_1 and e_2 must describe just sets of characters
- n *ident*: abbreviation for earlier reg exp in $\text{let } \textit{ident} = \textit{regex}$
- n $e_1 \text{ as } \textit{id}$: binds the result of e_1 to *id* to be used in the associated *action*



Ocamllex Manual

n More details can be found at

<http://caml.inria.fr/pub/docs/manual-ocaml/manual026.html>



Example : test.mll

```
{ type result = Int of int | Float of float |  
  String of string }
```

```
let digit = ['0'-'9']
```

```
let digits = digit +
```

```
let lower_case = ['a'-'z']
```

```
let upper_case = ['A'-'Z']
```

```
let letter = upper_case | lower_case
```

```
let letters = letter +
```



Example : test.ml

```
rule main = parse
  (digits)'.'digits as f { Float (float_of_string f) }
| digits as n           { Int (int_of_string n) }
| letters as s           { String s}
| _ { main lexbuf }
{ let newlexbuf = (Lexing.from_channel stdin) in
  print_string "Ready to lex.";
  print_newline ();
  main newlexbuf }
```



Example

```
# #use "test.ml";;
```

```
...
```

```
val main : Lexing.lexbuf -> result = <fun>
```

```
val __ocaml_lex_main_rec : Lexing.lexbuf -> int ->  
  result = <fun>
```

Ready to lex.

hi there 234 5.2

```
- : result = String "hi"
```

What happened to the rest?!?



Example

```
# let b = Lexing.from_channel stdin;;
```

```
# main b;;
```

```
hi 673 there
```

```
- : result = String "hi"
```

```
# main b;;
```

```
- : result = Int 673
```

```
# main b;;
```

```
- : result = String "there"
```



Problem

- n How to get lexer to look at more than the first token at one time?
- n Answer: *action* has to tell it to -- recursive calls
- n Side Benefit: can add "state" into lexing
- n Note: already used this with the _ case



Example

```
rule main = parse
  (digits) '.' digits as f { Float
    (float_of_string f) :: main lexbuf }
| digits as n              { Int (int_of_string n) ::
  main lexbuf }
| letters as s             { String s :: main
  lexbuf }
| eof                     { [] }
| _                       { main lexbuf }
```



Example Results

Ready to lex.

hi there 234 5.2

- : result list = [String "hi"; String "there"; Int
234; Float 5.2]

#

Used Ctrl-d to send the end-of-file signal



Dealing with comments

First Attempt

```
let open_comment = "("*"  
let close_comment = "*")"  
rule main = parse  
  (digits) '.' digits as f { Float (float_of_string  
    f) :: main lexbuf}  
| digits as n          { Int (int_of_string n) ::  
  main lexbuf }  
| letters as s         { String s :: main lexbuf}
```




Dealing with comments

| open_comment { comment lexbuf }

| eof { [] }

| _ { main lexbuf }

and comment = parse

close_comment { main lexbuf }

| _ { comment lexbuf }



Dealing with nested comments

```
rule main = parse ...
| open_comment      { comment 1 lexbuf }
| eof               { [] }
| _ { main lexbuf }
and comment depth = parse
  open_comment      { comment (depth+1) lexbuf }
| close_comment     { if depth = 1
                      then main lexbuf
                      else comment (depth - 1) lexbuf }
| _                 { comment depth lexbuf }
```



Dealing with nested comments

rule main = parse

(digits) '.' digits as f { Float (float_of_string f) ::
main lexbuf }

| digits as n { Int (int_of_string n) :: main
lexbuf }

| letters as s { String s :: main lexbuf }

| open_comment { (comment 1 lexbuf }

| eof { [] }

| _ { main lexbuf }



Dealing with nested comments

and comment depth = parse

```
open_comment      { comment (depth+1) lexbuf  
}
```

```
| close_comment   { if depth = 1  
                    then main lexbuf  
                    else comment (depth - 1) lexbuf }
```

```
| _               { comment depth lexbuf }
```