

# Lecture 3: Lexical Analysis

---

COMP 524 Programming Language Concepts

Stephen Olivier

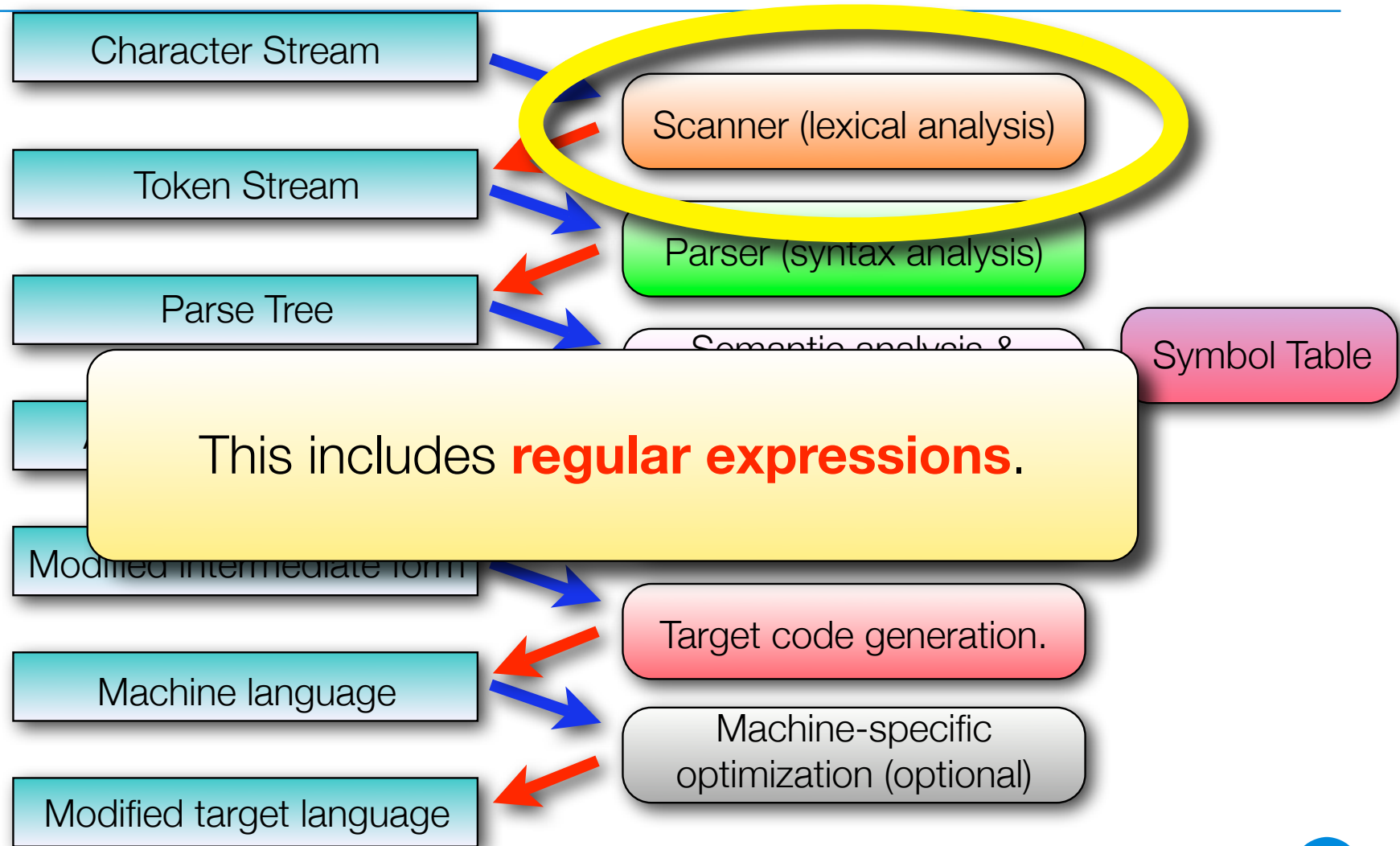
January 20, 2009

Based on notes by A. Block, N. Fisher, F. Hernandez-Campos, J. Prins and D. Stotts

The University of North Carolina at Chapel Hill



# Goal of Lecture



# Scanning

---

- The main task of scanning is to **identify tokens**.



## Pseudo-Code Scanner (Fig 2.5)

---

```
We skip any initial white spaces
we read the next character
if it is a ( we look at the next character
    if that is a * we have a comment;
        we skip forward through the terminating *)
    otherwise we return a ( and reuse the look-ahead
If it is one of the one-character tokens ([],;=+- etc.)
    we return that token
...
```



## Pseudo-Code Scanner (Fig 2.5)

---

```
We skip any initial white spaces
we read the next character
if it is
    if that
        we skip
    otherwise
If it is
    we return that token
...
```

We could just turn this into real code and use that as the scanner, and that would be fine for small programs...



## Pseudo-Code Scanner (Fig 2.5)

---

```
We skip any initial white spaces
we read the next character
if it is
    if that
        we sk
    otherwi
If it is
we return that token
...
```

... However, for larger programs that must be correct a more formal approach is more appropriate.



# Regular expressions

---

$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$non\_zero\_digit \rightarrow 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$natural\_number \rightarrow non\_zero\_digit \, digit^*$

$non\_neg\_number \rightarrow (0 \mid natural\_number) ( ( \cdot digit^* non\_zero\_digit ) \mid \epsilon )$



# Regular expressions

*digit*  $\rightarrow$  0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*non\_zero\_digit*  $\rightarrow$  1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*natural\_number*  $\rightarrow$  *non\_zero\_digit digit\**

“ $\rightarrow$ ” denotes **assignment**





# Regular expressions

---

*digit* → 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*non\_zero\_digit* → 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*natural\_number* → *non\_zero\_digit* *digit*\*

“|” denotes **or**

| ε )



# Regular expressions

---

*digit*  $\rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

*non\_zero\_digit*  $\rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

*natural\_number*  $\rightarrow non\_zero\_digit digit^*$

Thus, digit equal "0" or "1" or "2" or ....

*non*

$| \epsilon )$



# Regular expressions

“**\***” denotes **zero or more of this type**.

*non\_zero\_digit* → 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*natural\_number* → *non\_zero\_digit* *digit*\*

*non\_neg\_number* → (0 | *natural\_number*) ( ( . *digit*\* *non\_zero\_digit*) | ε )



# Regular expressions

Two REs next to each other denotes  
**concatenation.**

*non\_zero\_digit* → 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

*natural\_number* → *non\_zero\_digit* *digit*\*

*non\_neg\_number* → ( 0 | *natural\_number* ) ( ( . *digit*\* *non\_zero\_digit* ) |  $\epsilon$  )



# Regular expressions

So natural number equals at least  
**“one non-zero digit” followed by  
“zero or more digits”**.

$non\_zero\_digit \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$natural\_number \rightarrow non\_zero\_digit digit^*$

$non\_neg\_number \rightarrow (0 | natural\_number) ( ( . digit^* non\_zero\_digit ) | \epsilon )$



# Regular expressions

“ $\epsilon$ ” means **empty**.

$non\_zero\_digit \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$natural\_number \rightarrow non\_zero\_digit digit^*$

$non\_neg\_number \rightarrow (0 | natural\_number) ( ( . digit^* non\_zero\_digit ) | \epsilon )$



# Regular expressions

So, what does this mean?

$non\_zero\_digit \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$natural\_number \rightarrow non\_zero\_digit\ digit^*$

$non\_neg\_number \rightarrow (0 | natural\_number) ( ( \cdot digit^* non\_zero\_digit ) | \epsilon )$



# Regular expressions

It means “**0** or a **natural number**” followed by “**nothing**” or by “**.** and **zero or more digits** concluded by a **non\_zero number**”

$non\_zero\_digit \rightarrow 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

$natural\_number \rightarrow non\_zero\_digit\ digit^*$

$non\_neg\_number \rightarrow (0 | natural\_number) ( ( \cdot digit^* non\_zero\_digit ) | \epsilon )$





# Regular Expression Rules

---

- A RE consist of:
  - A character (e.g., “0”, “1”, ...)
  - The empty string (i.e., “ $\epsilon$ ”)
  - Two REs next to each other (e.g., “*non\_negative\_digit digit*”) to denote concatenation.
  - Two REs separated by “|” next to each other (e.g., “*non\_negative\_digit | digit*”) to denote one RE or the other.
  - An RE followed by “\*” (called the **Kleene star**) to denote zero or more iterations of the RE.
  - Parentheses (in order to avoid ambiguity).



# Regular Expression Rules

---

- A RE consist of:

- A character (e.g., “0”, “1”, ...)

- The

A RE is **NEVER** defined in terms of itself!

- Two

Thus, REs cannot **define recursive statements**.

depend

- Two REs separated by “|” next to each other (e.g., “*non\_negative\_digit* | *digit*”) to denote one RE or the other.

- An RE followed by “\*” (called the **Kleene star**) to denote zero or more iterations of the RE.

- Parentheses (in order to avoid ambiguity).



# Regular Expression Rules

---

- A RE consist of:

- A character (e.g., “0”, “1”, ...)

- The

- Two  
depend

The set of tokens that can be recognized by regular expressions is called a **regular set**.

- Two REs separated by “|” next to each other (e.g., “*non\_negative\_digit* | *digit*”) to denote one RE or the other.

- An RE followed by “\*” (called the **Kleene star**) to denote zero or more iterations of the RE.

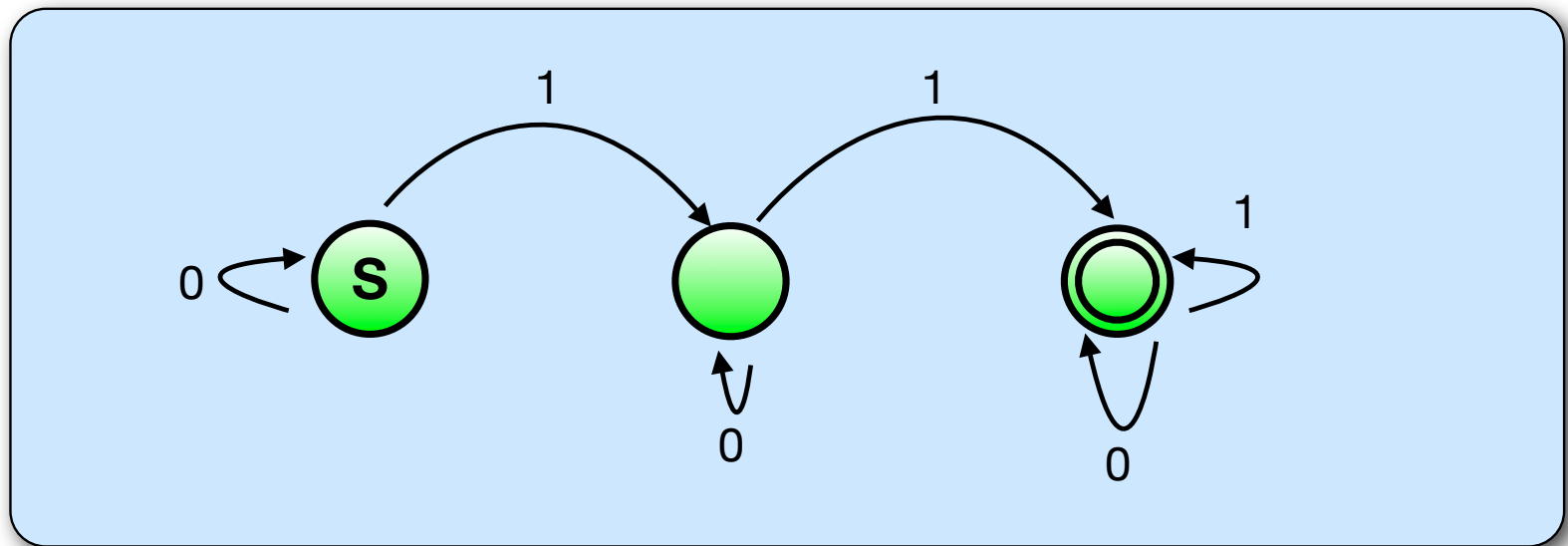
- Parentheses (in order to avoid ambiguity).



# Deterministic finite automaton (DFA)

---

- Every regular set can be defined by using **deterministic finite automaton (DFA)**.
  - DFAs are turing machines that have a finite number of states and deterministically move between states.

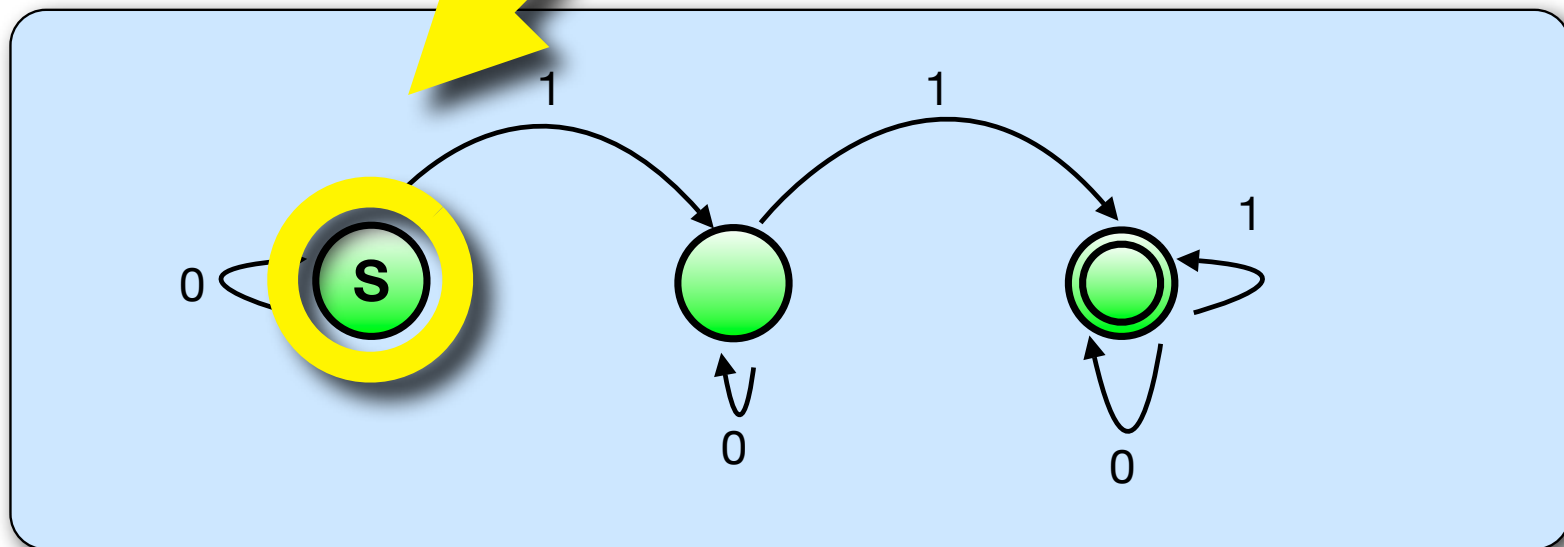


# Deterministic finite automaton (DFA)

- Every regular set can be defined by using **deterministic finite automaton**

Start State

- DFAs are turning machines that have a finite number of states and deterministically move between states.

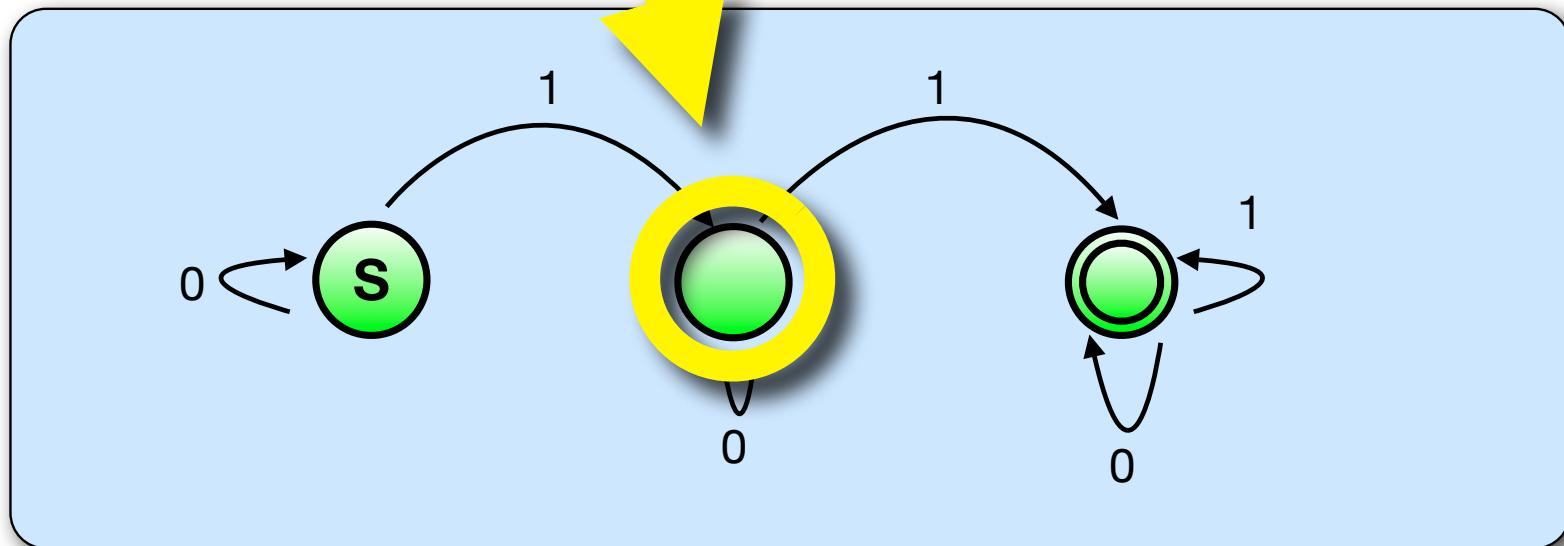


# Deterministic finite automaton (DFA)

- Every regular set can be defined by using **deterministic finite automaton**

Intermediate State

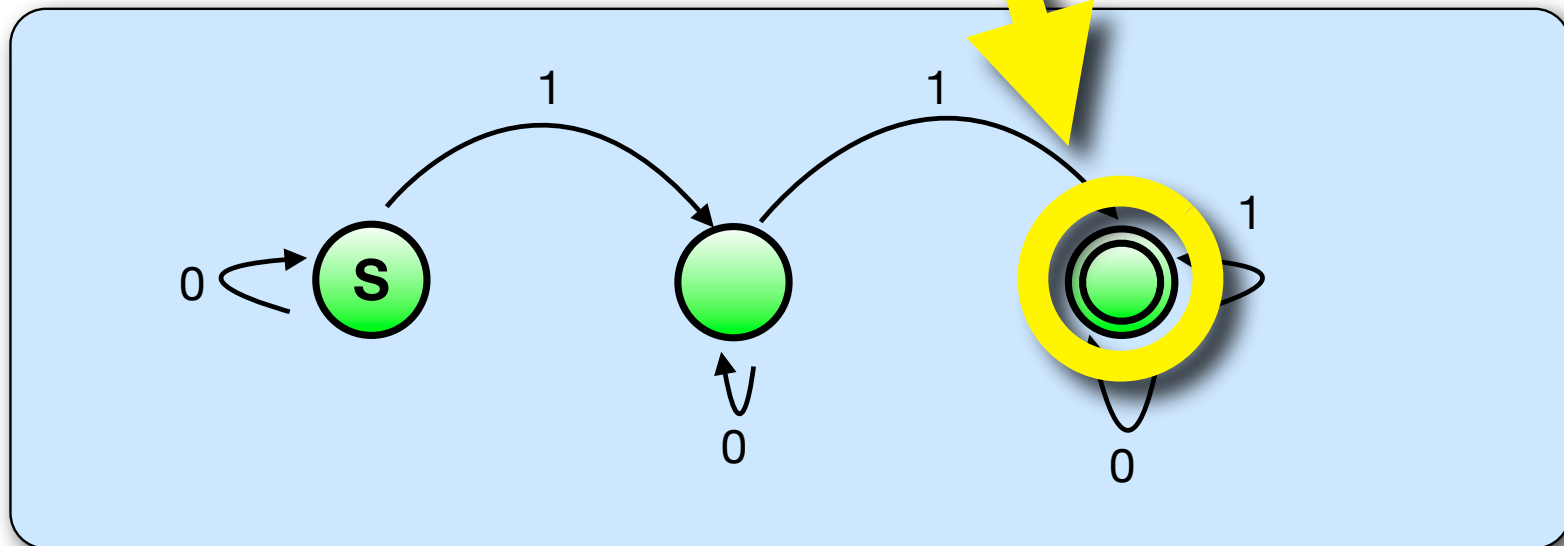
- DFAs are ~~turing machines~~ machines that have a finite number of states and deterministically move between states.



# Deterministic finite automaton (DFA)

- Every regular set can be defined by using **deterministic finite automaton**
- DFAs are turning machines that take a sequence of states and deterministically move between states.

End State (double circle)



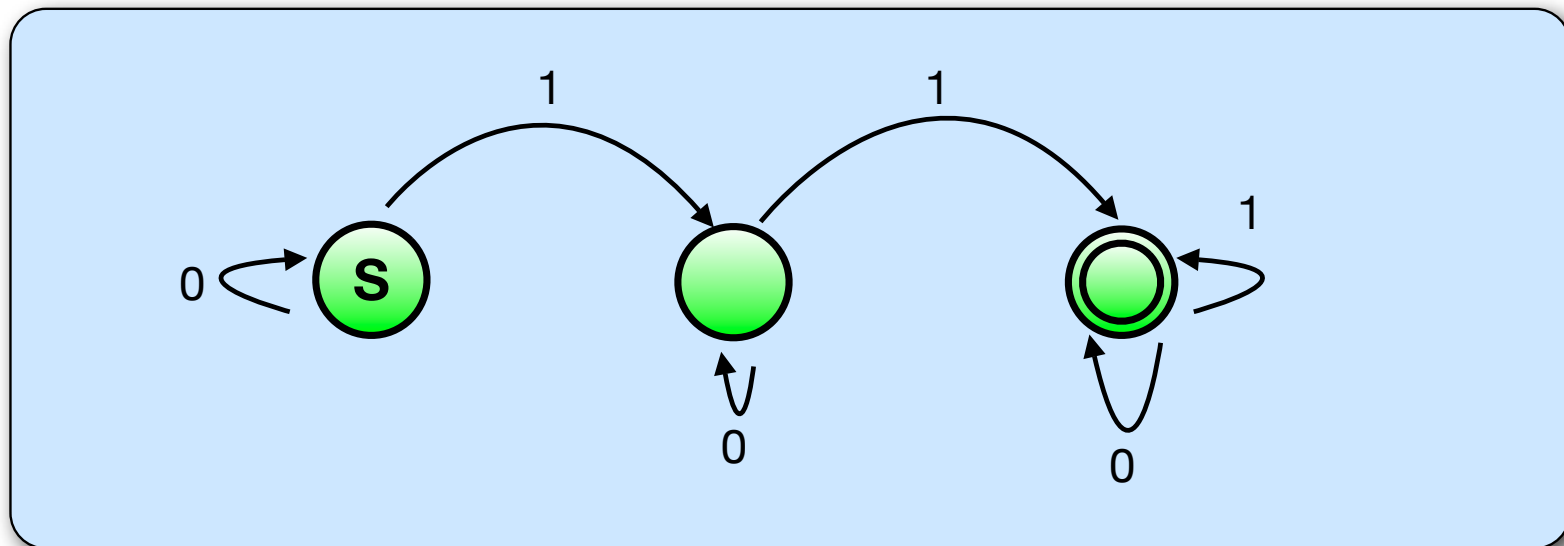
# Deterministic finite automaton (DFA)

- Every regular set can be defined by using **deterministic finite automaton**

- DFAs can be used to define regular sets. The number of states and the transitions between states define the regular set.

This stands for:

$0^*10^*1(0|1)^*$





# Constructing DFAs

---

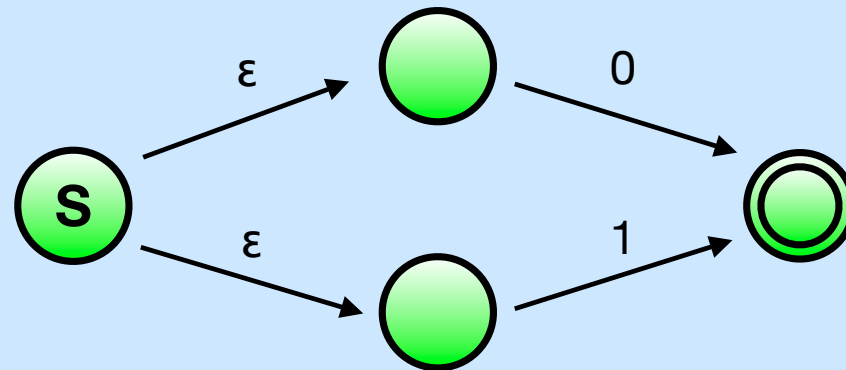
- A DFA can be constructed from a RE via two steps.
  1. Construct a **nondeterministic finite automaton (NFA)** from the RE.
  2. Construct a DFA from the NFA.
  3. Minimize the DFA



# What is an NFA?

---

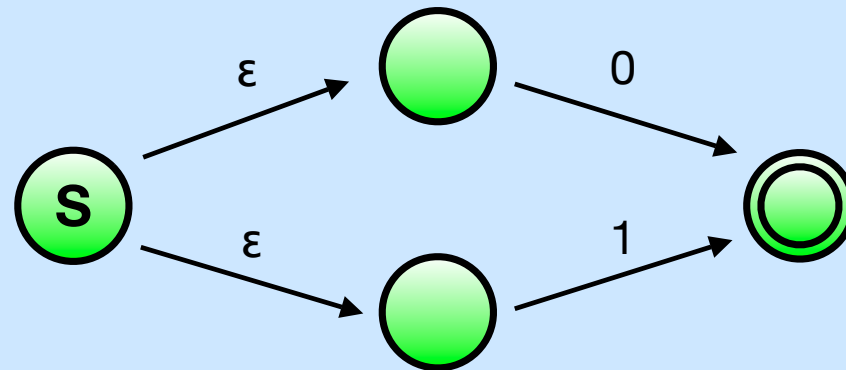
- An NFA is similar to a DFA, except that state transitions are nondeterministic.
- This nondeterminism is encapsulated via the **epsilon transition** (written as  $\epsilon$ ).



# What is an NFA?

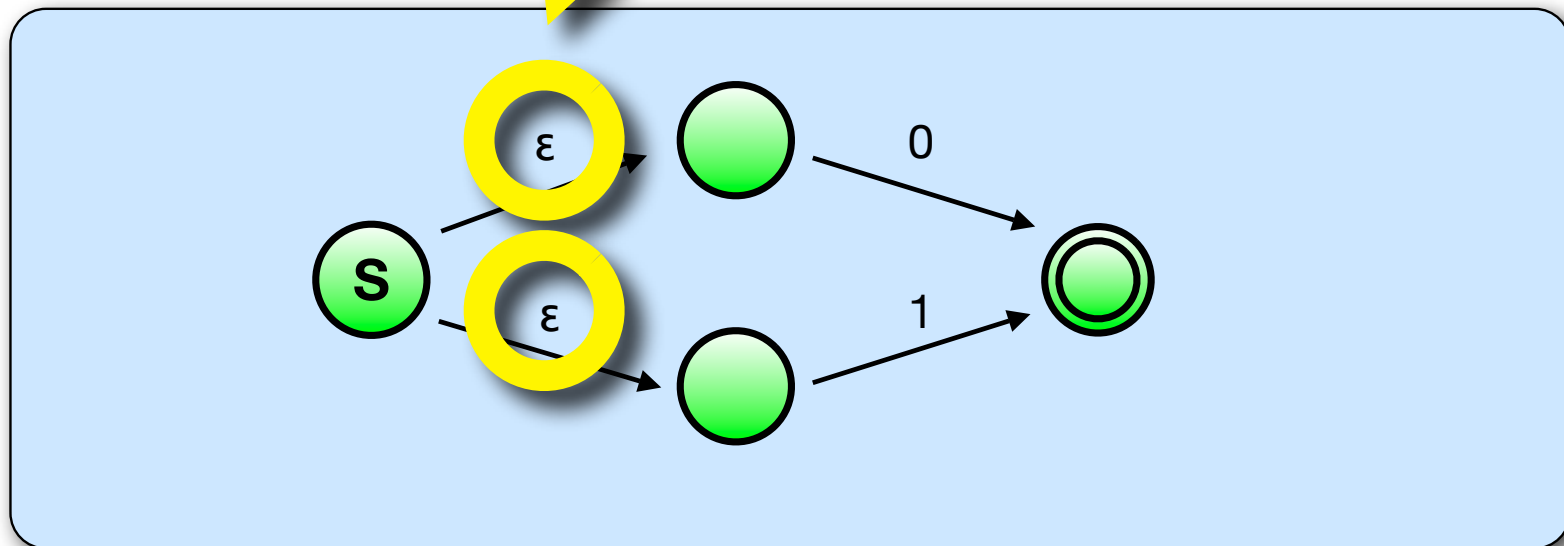
- An NFA is similar to a DFA in that state transitions are nondeterministic.
- This nondeterminism is achieved via the **epsilon transition** (written  $\epsilon$ ).

This stands for:  
 $0|1$



What is an NFA?

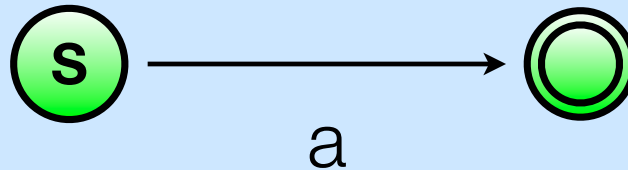
- An NFA is a finite state machine where multiple transitions are possible for the same input symbol. The  $\epsilon$  transitions imply that either transition can be taken with any (or no) input.
- This nondeterminism is encapsulated via the **epsilon transition** (written as  $\epsilon$ ).



# The four RE rules and NFA

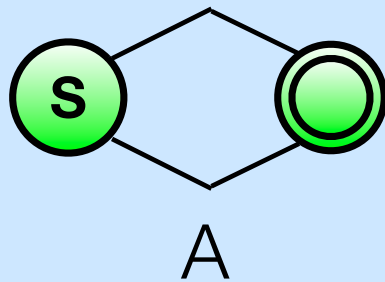
---

Rule 1--**Base case**: “a”

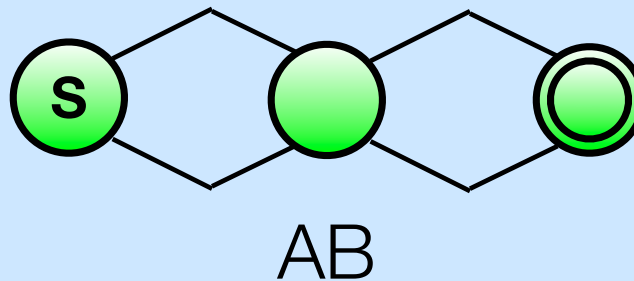
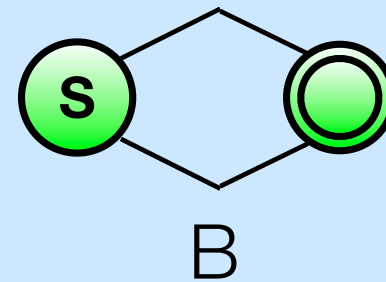


# The four RE rules and NFA

Rule 2--**Concatenation**: "AB"



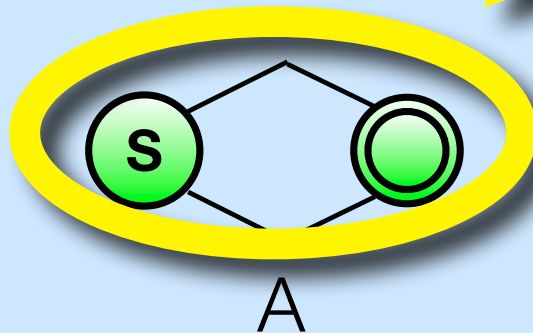
plus



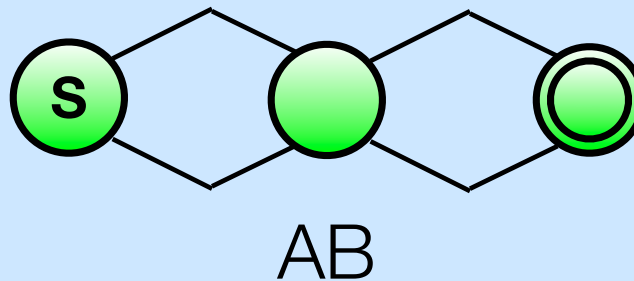
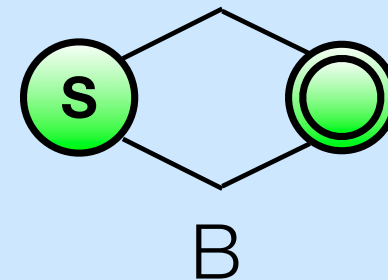
The four

Stands for Some NFA called "A"

Rule 2--**Concatenation**: "AB"

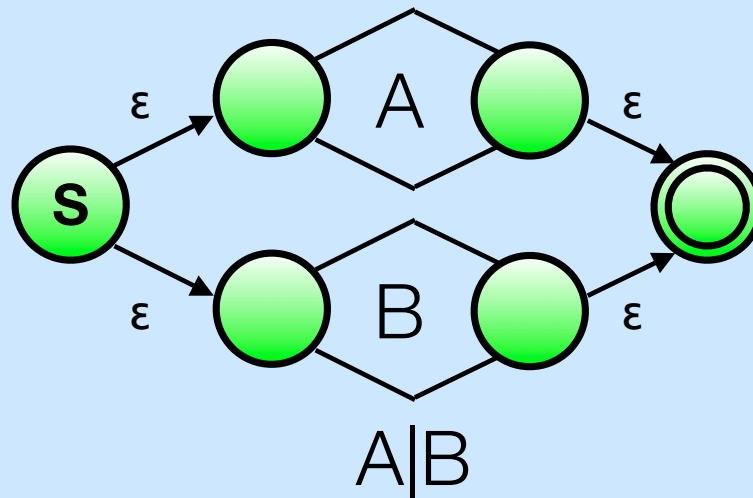
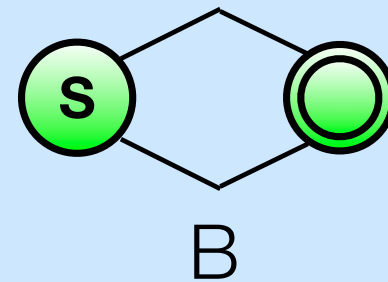
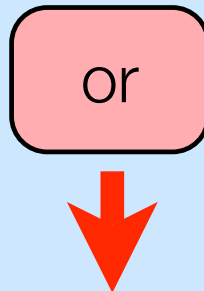
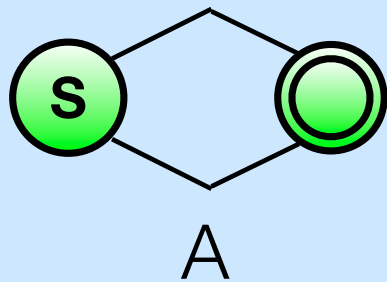


plus



## The four RE rules and NFA

Rule 3--**Alternation**: “A|B”

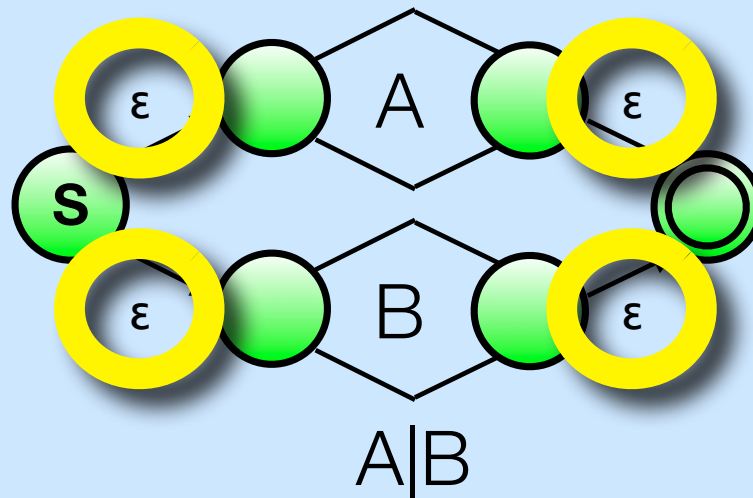
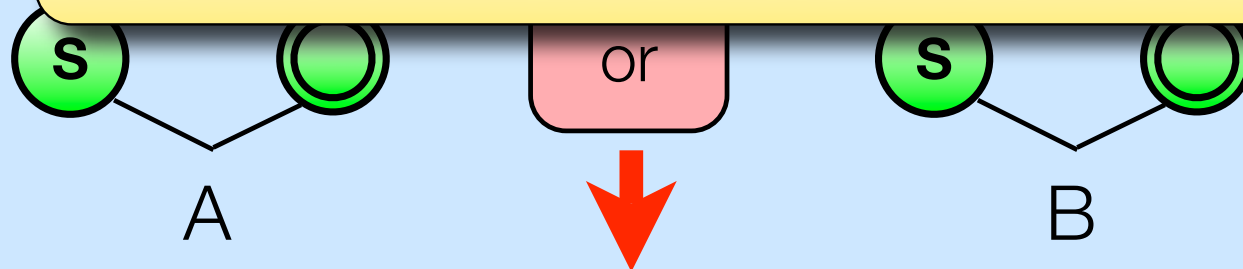




# The four RE rules and NFA

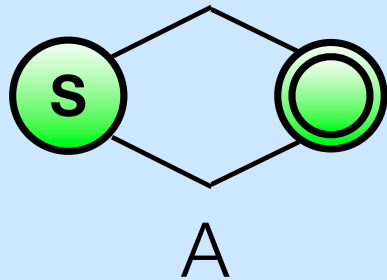
Rule 2: **Alternation**: "A|B"

Notice the epsilon transitions.

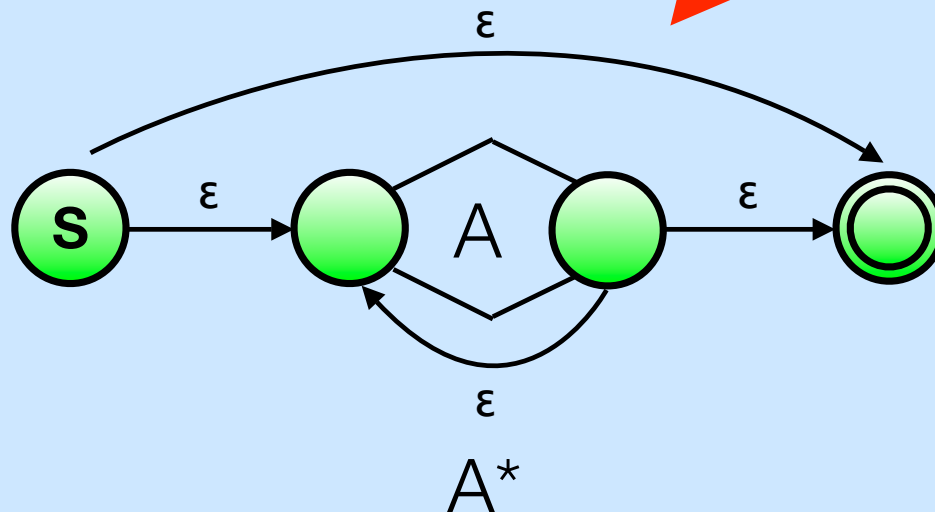


# The four RE rules and NFA

Rule 4--**Kleene Closure**: " $A^*$ "



empty or repeated

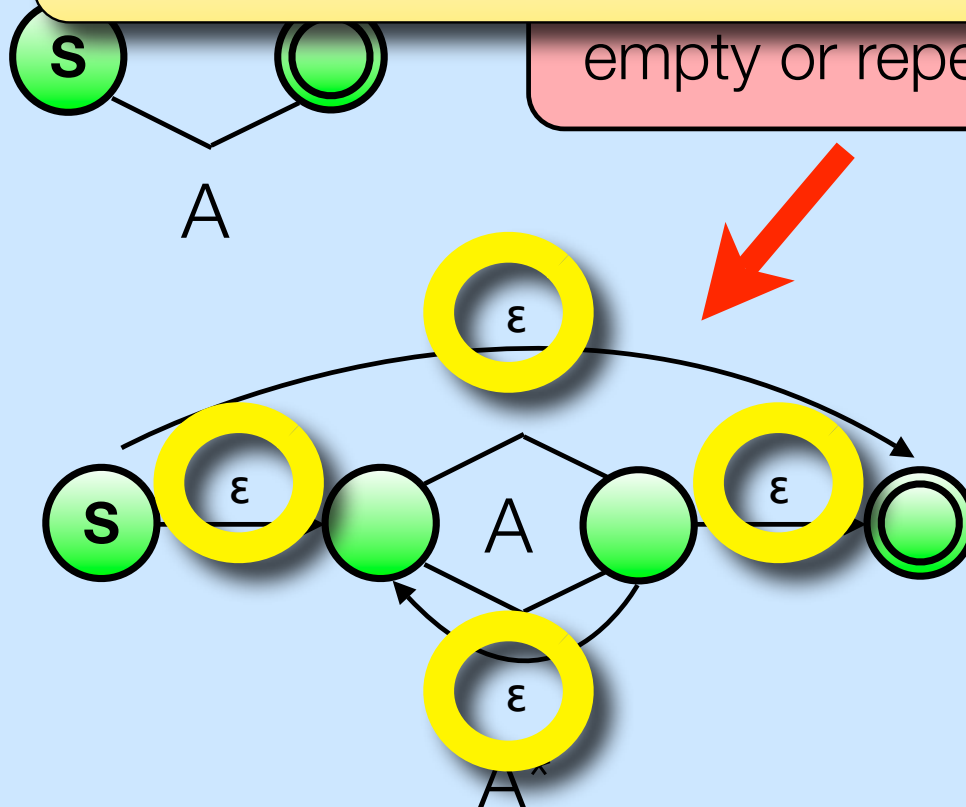


# The four RE rules and NFA

Rule 4: **Kleene Closure**: " $A^*$ "

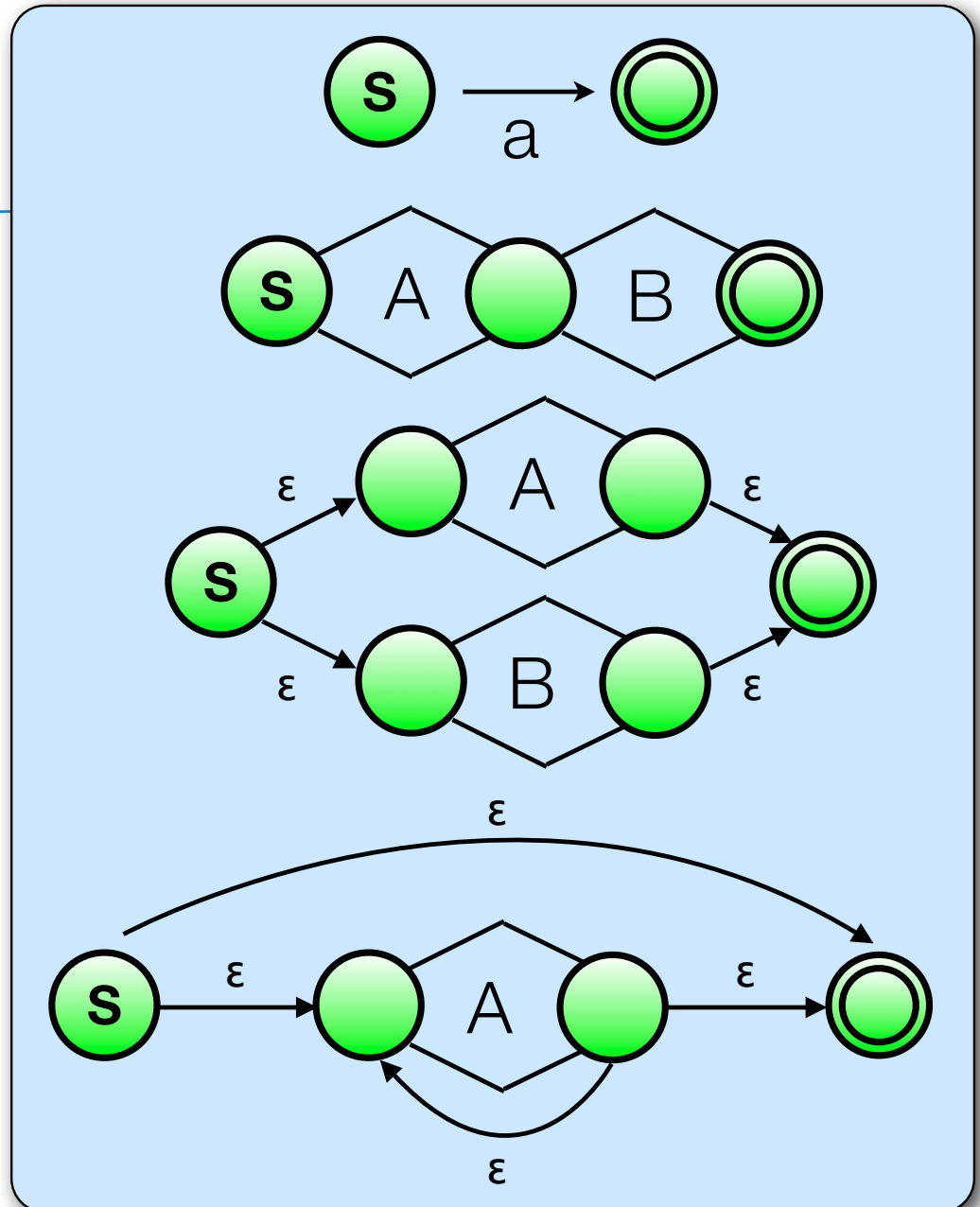
Notice the epsilon transitions.

empty or repeated



Some examples:

- $0|1^*$
- $AB^*$
- $F|(GH^*)$
- $Z^*|\epsilon|Y^*X$



# Constructing DFAs

---

- A DFA can be constructed from a RE via two steps.
  1. Construct a **nondeterministic finite automaton (NFA)** from the RE.
  2. Construct a DFA from the NFA.
  3. Minimize the DFA



# Constructing DFAs

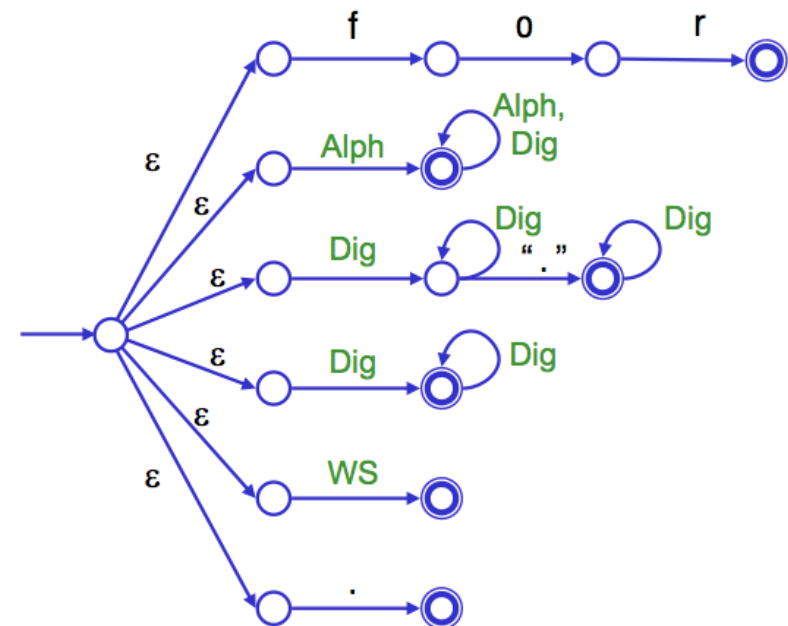
- A DFA can be constructed from a RE via two steps.
  1. Construct a **nondeterministic finite automaton (NFA)** from the RE.

*FOR\_KEYWORD* → for

*IDENTIFIER* → Alph ( Alph | Dig )<sup>\*</sup>

*REAL* → Dig Dig<sup>\*</sup> . Dig<sup>\*</sup>

*INT* → Dig Dig<sup>\*</sup>



# Constructing DFAs

---

- A DFA can be constructed from a RE via two steps.
  1. Construct a **nondeterministic finite automaton (NFA)** from the RE.
  2. Construct a DFA from the NFA.
  3. Minimize the DFA



## Constructing a DFA from an NFA.

---

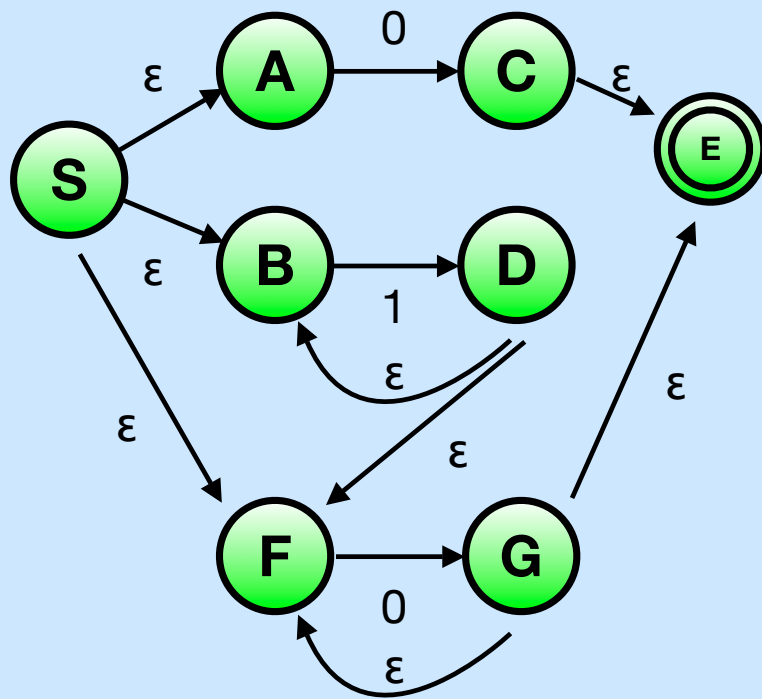
- Construct the DFA by “collapsing” the states of an NFA.
- Three steps
  1. Identify set of states that can be reached from the **start state** via **epsilon-transitions** and make this **one state**.
  2. For a given DFA state (which is a set of NFA states) consider each possible input and combine the **resulting NFA states into one DFA state**.
  3. Repeat Step 2 until all states have been added.





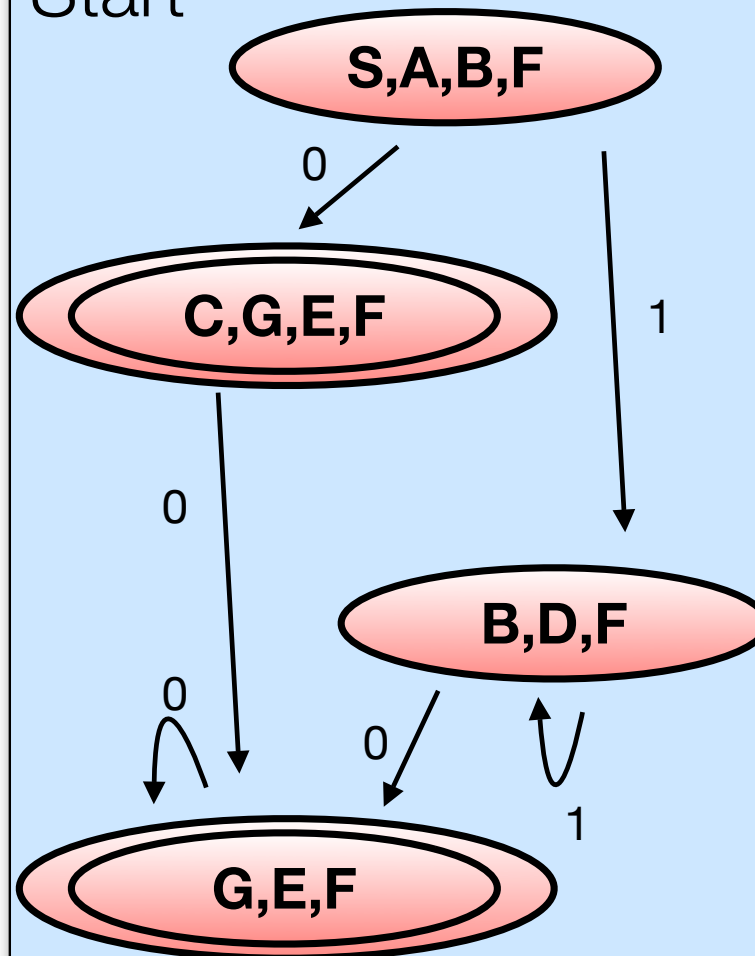
# An example

$0|1^*00^*$



NFA

Start



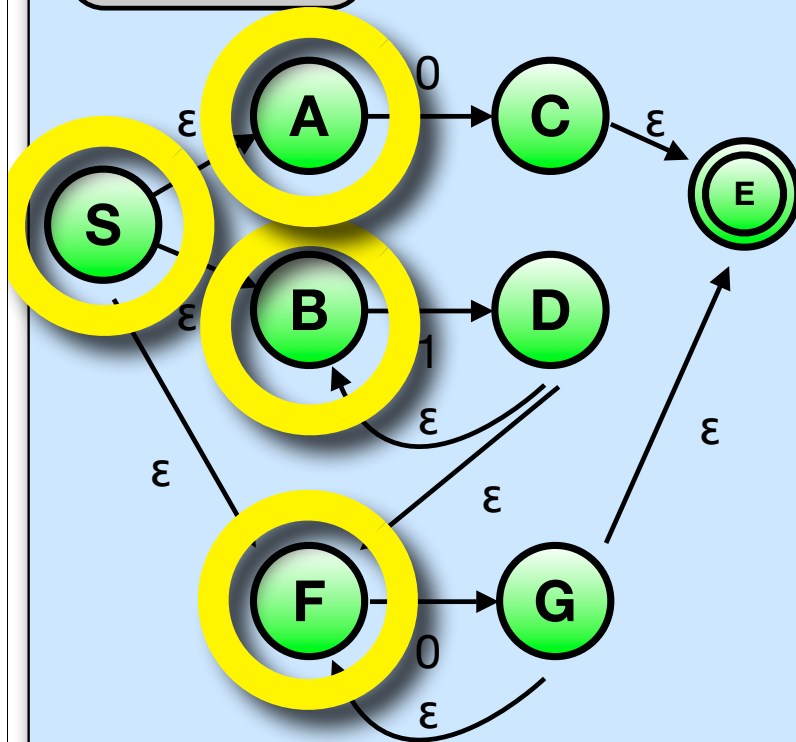
DFA



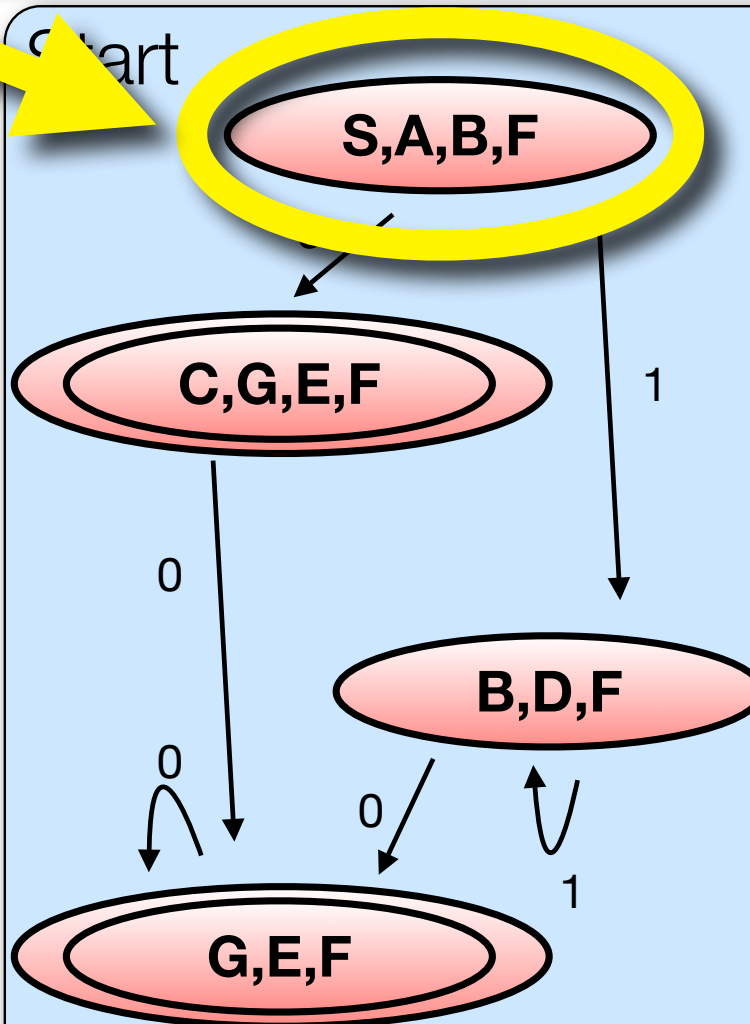
All the states that we can reach via  $\epsilon$  are in this state

An example

$0|1^*00^*$



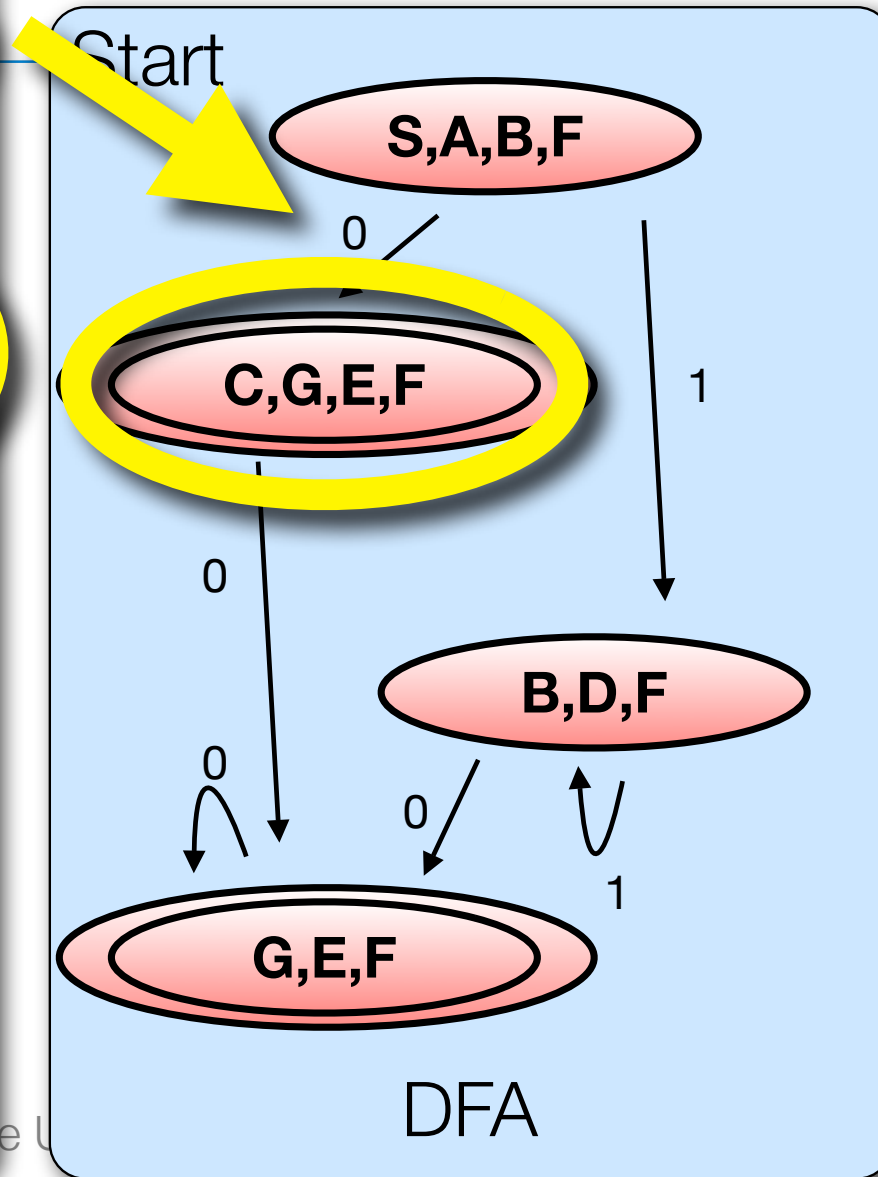
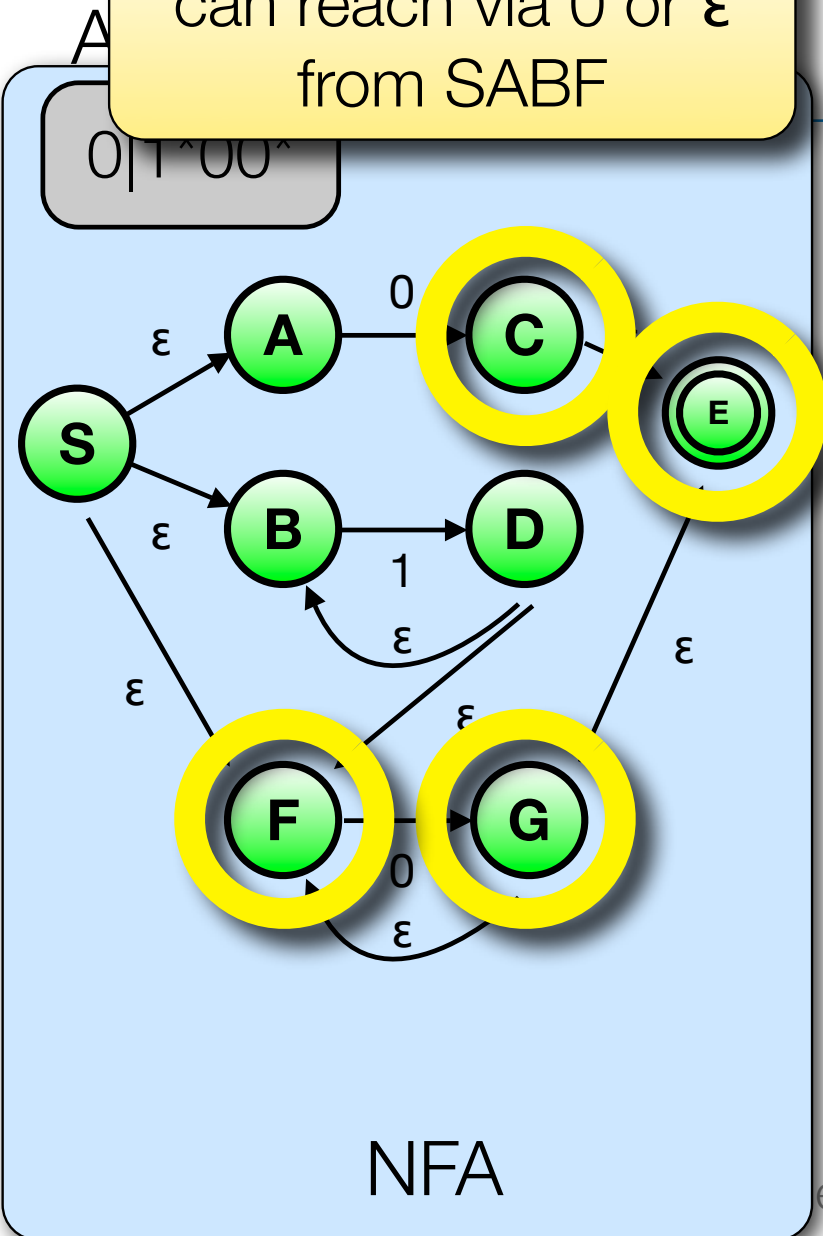
NFA



DFA



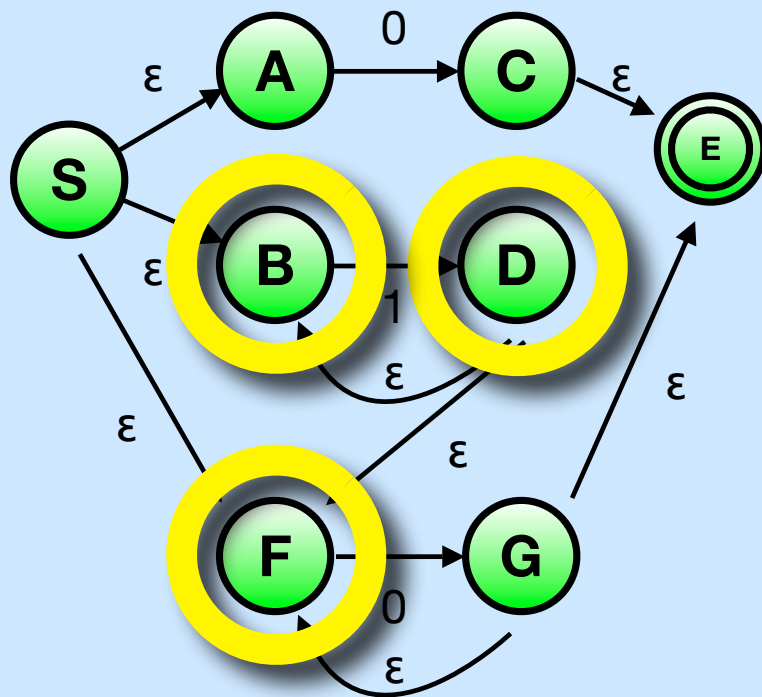
All the states that we can reach via 0 or  $\epsilon$  from SABF



All the states that we  
can reach via 1 or  $\epsilon$   
from SABF

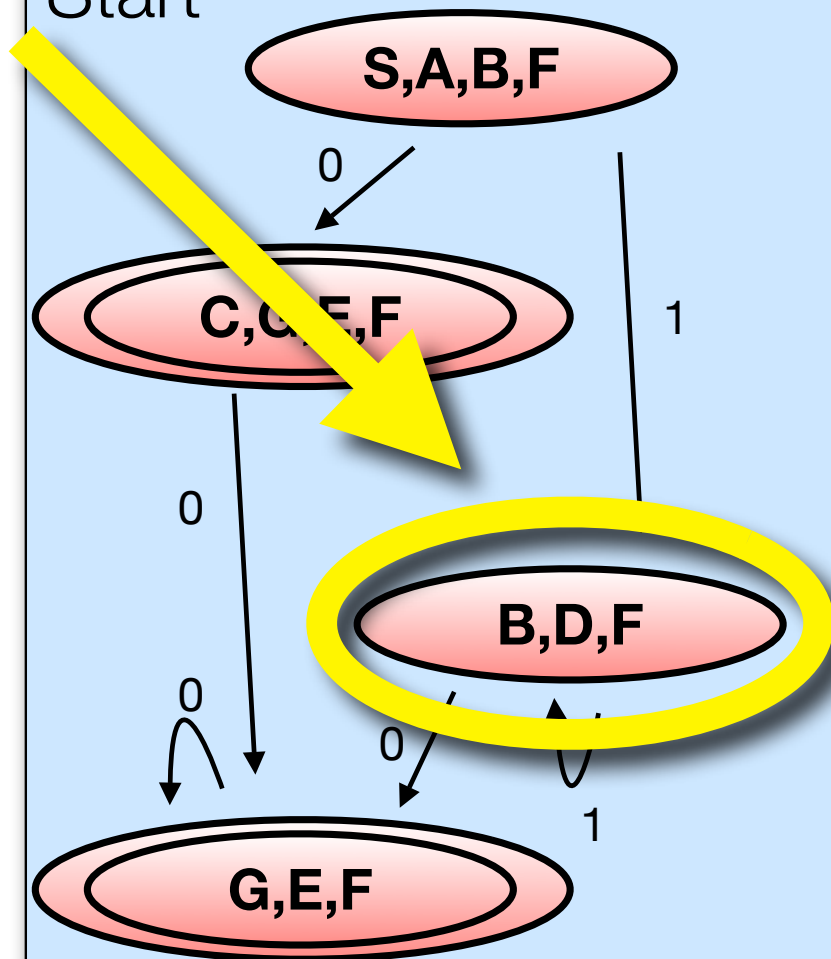
A

$0|1^*00^*$



NFA

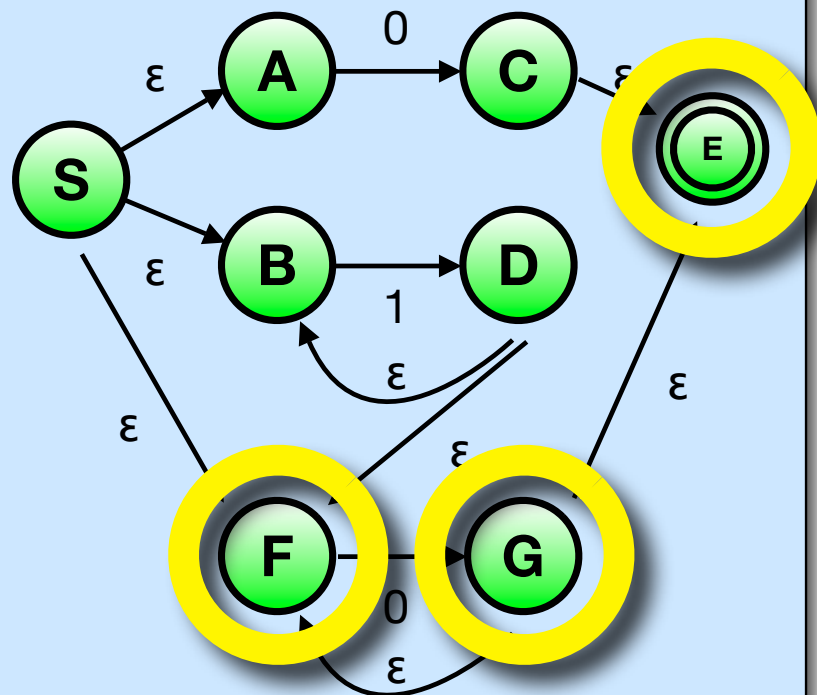
Start



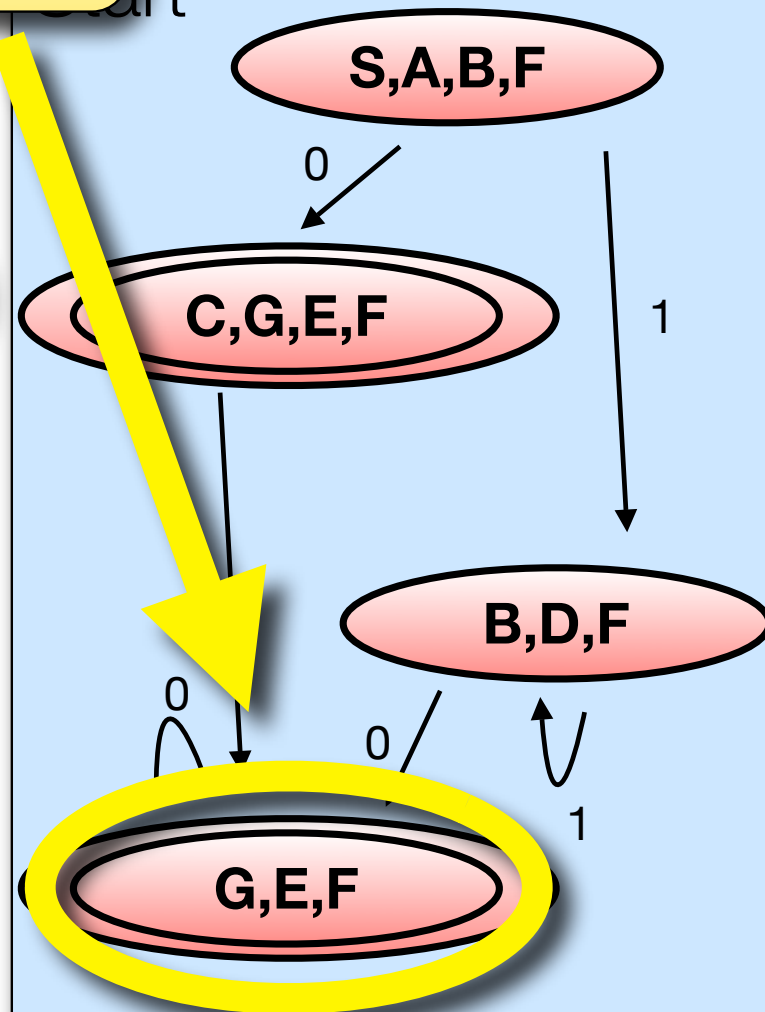
DFA



All the states that we can reach via 0 or  $\epsilon$  from B,D,F or C,G,F



NFA



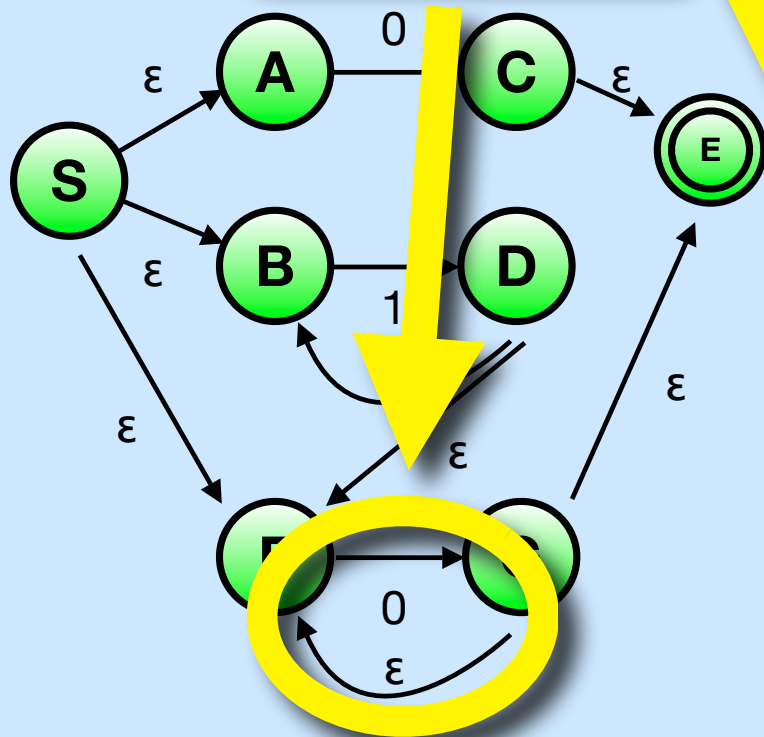
DFA



# An example

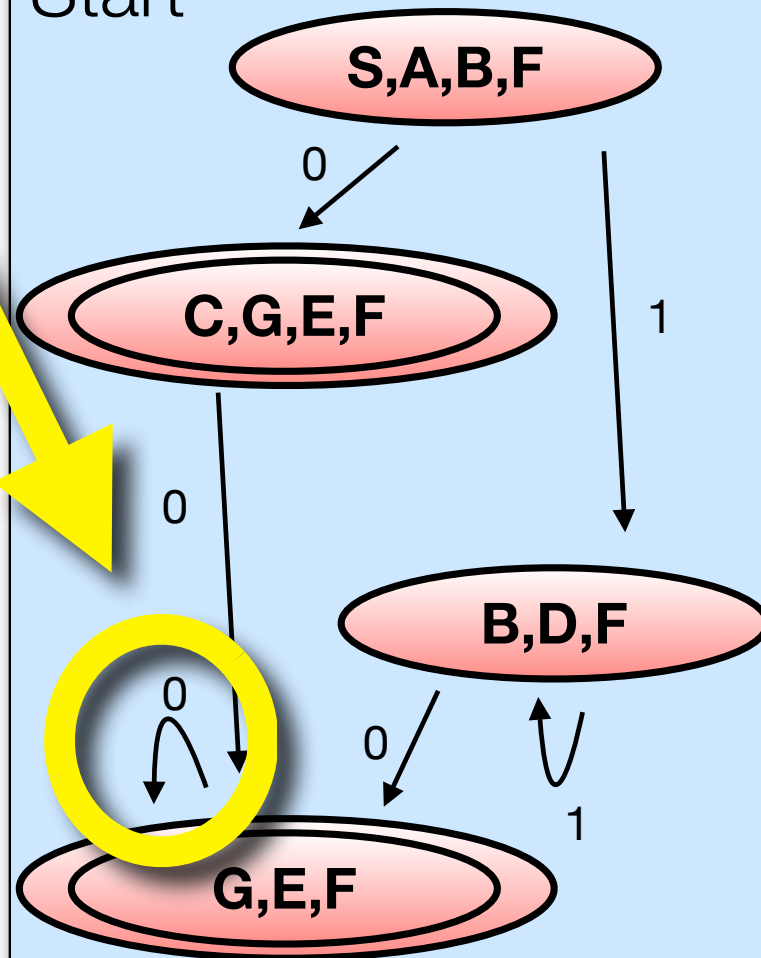
0|1\*0

Self Loop.



NFA

Start

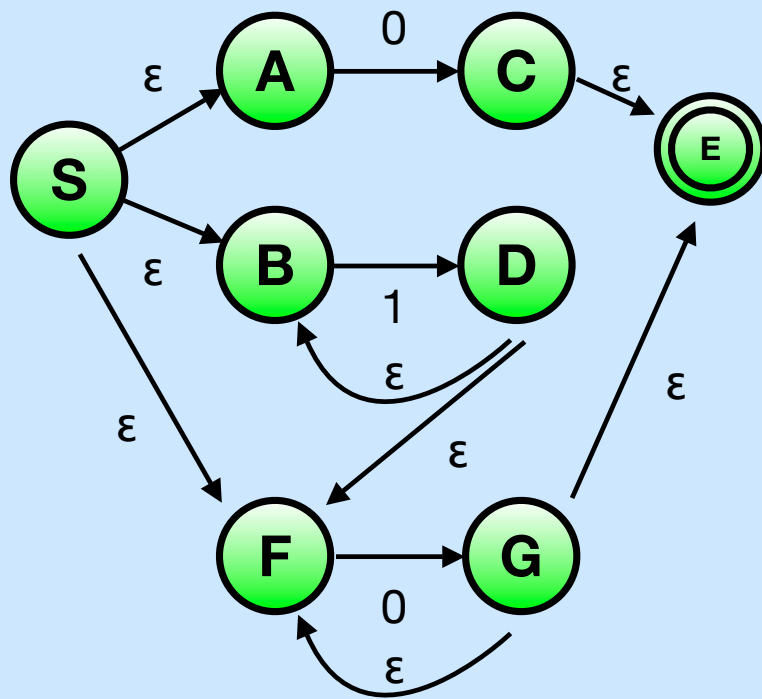


DFA



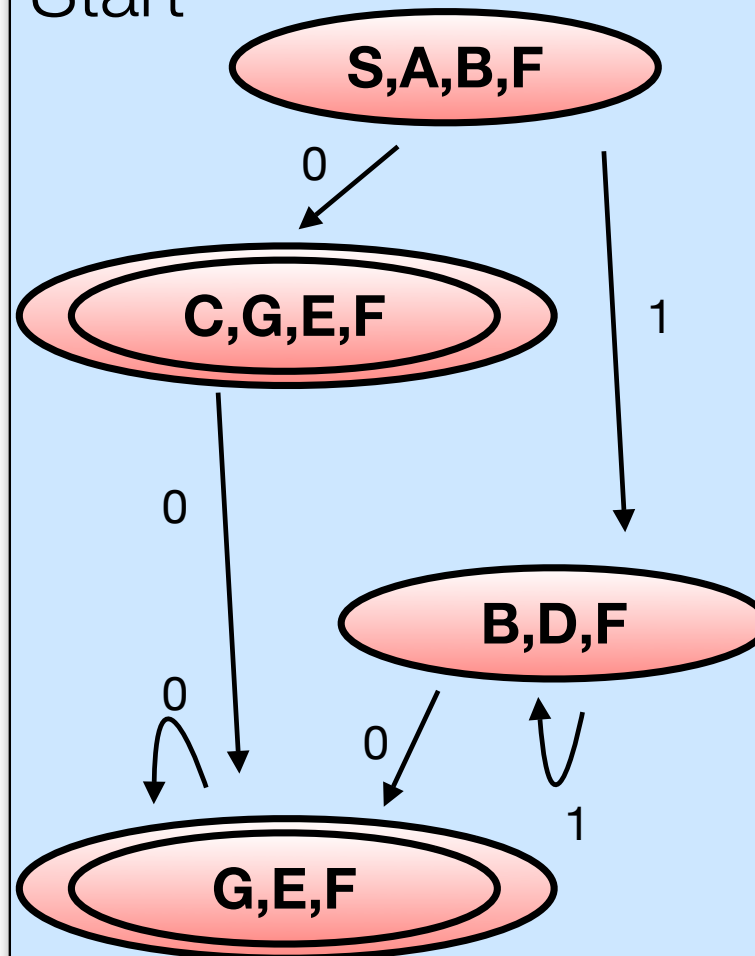
# An example

$0|1^*00^*$



NFA

Start



DFA



## Minimize via partitioning

---

- First, partition states into final and non-final
- Second, determine the effect of the state transition based on what partition the transition goes to.
- Third, Create new partition for those states that have different transitions.
- Fourth, repeat.



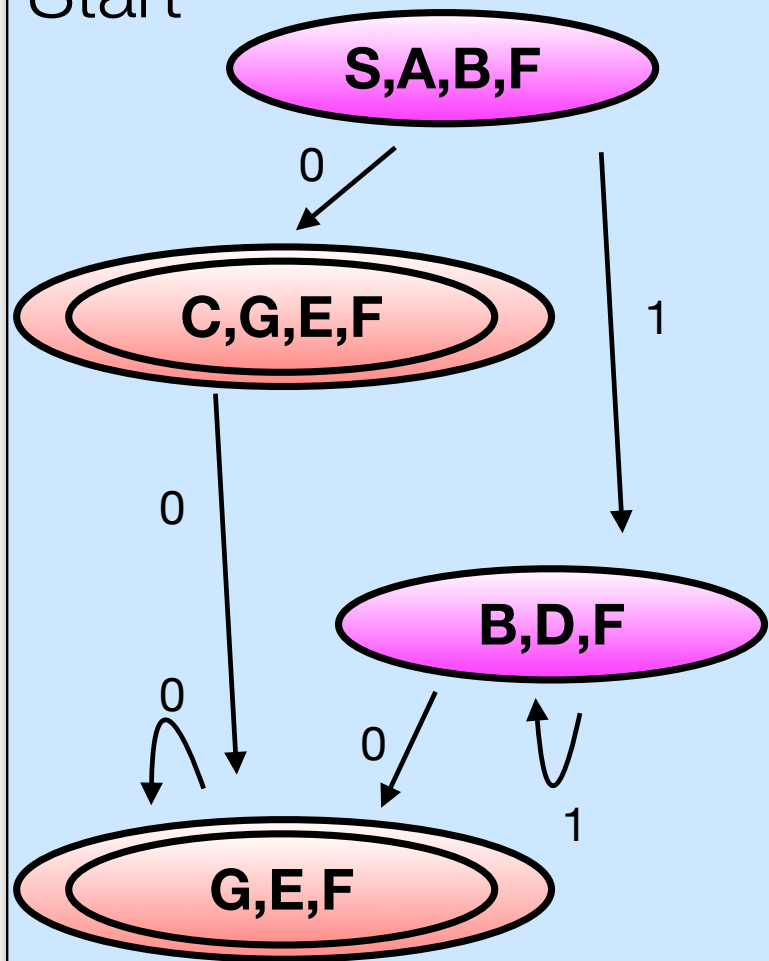


X-1

X-2

An example

Start



DFA

State	0	1
<b>SABF</b>	X-2	X-1
<b>BDF</b>	X-2	X-1
<b>CGEF</b>	X-2	N/A
<b>GEF</b>	X-2	N/A

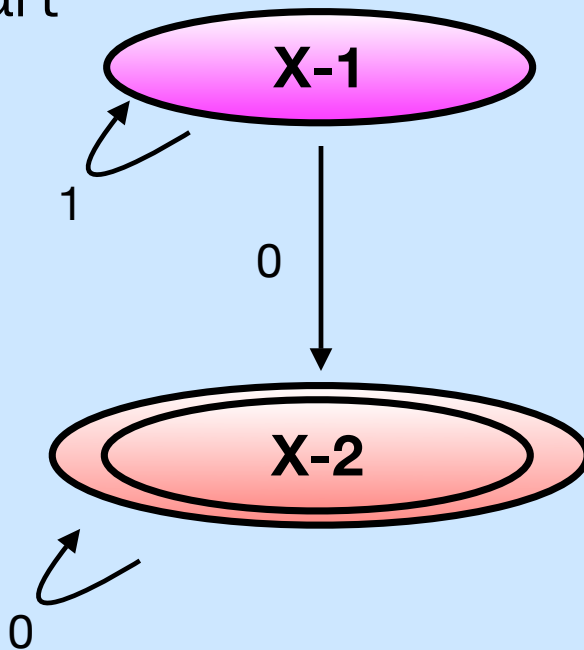


X-1

X-2

An example

Start



State	0	1
SABF	X-2	X-1
BDF	X-2	X-1
CGEF	X-2	N/A
GEF	X-2	N/A



# Scanner Code

---

- Can create Scanner from the DFA one of two ways:
  - Nested case statements (Handwritten)
  - Tables (easy to generate from code, hard to write by hand)



```

state := start
loop
  case state of
    start :
      erase text of current token
      case input_char of
        '\t', '\n', '\r' : no_op
        '[' : state := got_lbrac
        ']' : state := got_rbrac
        ',' : state := got_comma
        ...
        '(' : state := saw_lparen
        '.' : state := saw_dot
        '<' : state := saw_lthan
        ...
        'a'..'z', 'A'..'Z' :
          state := in_ident
        '0'..'9' : state := in_int
        ...
      else error
    ...
    saw_lparen: case input_char of
      '*' : state := in_comment
      else return lparen
    in_comment: case input_char of
      '*' : state := leaving_comment
      else no_op
    leaving_comment: case input_char of
      ')' : state := start
      else state := in_comment
    ...
    saw_dot : case input_char of
      '.' : state := got_dotdot
      else return dot
    ...

```

```

else return dot
', ' : state := got_comma
saw_dot : case input_char of
...
else state := in_comment
')' : state := start

```

```

saw_lthan : case input_char of
  '=' : state := got_le
  else return lt
...
in_ident : case input_char of
  'a'..'z', 'A'..'Z', '0'..'9', '_' : no_op
  else
    look up accumulated token
    in keyword table
    if found, return keyword
    else return identifier
...
in_int : case input_char of
  '0'..'9' : no_op
  '_' :
    peek at character beyond input_char;
    if '0'..'9', state := saw_real_dot
    else
      unread peeked-at character
      return intconst
  'a'..'z', 'A'..'Z', '_' : error
  else return intconst
...
saw_real_dot : ...
...
got_lbrac : return lbrac
got_rbrac : return rbrac
got_comma : return comma
got_dotdot : return dotdot
got_le : return le
...
append input_char to text of current token
read new input_char

```

```

read new input_char
append input_char to text of current token
...
got_le : return le
got_dotdot : return dotdot
got_comma : return comma

```



## Two complications--Nested Case

---

- **Keywords**

- It is possible to maintain a DFA for keywords, but the number of states would be even larger! So, they are handled as exceptions to the rule.

- **“Dot-Dot”**

- Pascal uses “..” to denote a range of numbers; however, to determine the meaning of the “..” we need to “look ahead” after reading the first “.” to determine if “.” denotes the end of a token or a beginning of a new token.
  - “3.14” one token
  - “2 .. 5” three tokens



```

state = 1..number of states
action_rec = record
  action : (move, recognize, error)
  new_state : state

```

This code specifies a two-dimensional transition table, which tells “whether to move, return token, or announce error”

```

state := start_state
image := null
repeat
  loop
    read cur_char
    case scan_tab[cur_char, cur_state].action
      move:
        cur_state := scan_tab[cur_char, cur_state].new_state
      recognize:
        tok := scan_tab[cur_char, cur_state].token_found
        exit inner loop
      error:
        -- print error message and recover; probably start over
        append cur_char to image
    -- end inner loop
  until tok not in [white_space, comment]
  look image up in keyword_tab and replace tok with appropriate keyword if found
  return (tok, image)

```

```

return (tok, image)
look image up in keyword_tab and replace tok with appropriate keyword if found
until tok not in [white_space, comment]
-- end inner loop
append cur_char to image

```



```
state = 1..number of states
action_rec = record
  action : (move, recognize, error)
  new_state : state
```

A second table tells when we might have hit the end of a token (for backing up)

```
state := start_state
image := null
repeat
  loop
    read cur_char
    case scan_tab[cur_char, cur_state].action
      move:
        cur_state := scan_tab[cur_char, cur_state].new_state
      recognize:
        tok := scan_tab[cur_char, cur_state].token_found
        exit inner loop
      error:
        -- print error message and recover; probably start over
    append cur_char to image
  -- end inner loop
until tok not in [white_space, comment]
look image up in keyword_tab and replace tok with appropriate keyword if found
return (tok, image)
```

```
return (tok, image)
look image up in keyword_tab and replace tok with appropriate keyword if found
until tok not in [white_space, comment]
-- end inner loop
append cur_char to image
```



# Pragmas

---

- **Pragmas** are “comments” that provide direction for the compiler.
  - For example, “Variable x is used a lot, keep it in memory if possible.”
- These are often handled by the parser since this makes the grammar much simpler.

