

Regular Languages & Finite Automata

- Basic formal language theory result:

Regular expressions and finite automata both define the class of regular languages.

Aspect	Regular Expression	Finite Automaton
Used For	Specification (what pattern to match)	Implementation (how to match it)
Form	Symbolic pattern	State-transition diagram or table
Conversion	RE \rightarrow NFA \rightarrow DFA	DFA \rightarrow RE (via state elimination)
Tool Use	Lexer specs (Flex, Lex)	Actual automaton used for tokenizing

Thus, we are going to use:

Regular expressions for specification

Finite automata for implementation

(automatic generation of lexical analyzers)

Finite Automata

- A finite automaton is a recognizer for the strings of a regular language
- A finite automaton consists of
 - A finite input alphabet Σ
 - A set of states S
 - A start state n
 - A set of accepting states $F \subseteq S$
 - A set of transitions $\text{state}_i \xrightarrow{\text{input}} \text{state}_j$

Finite Automata

- Transition

$$s_1 \xrightarrow{a} s_2$$

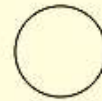
- Is read

In state s_1 on input "a" go to state s_2

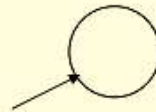
- If end of input
 - If in accepting state \Rightarrow accept, Otherwise \Rightarrow reject
- If No transition possible \Rightarrow Reject

Finite Automata State Graphs

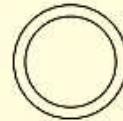
- A state



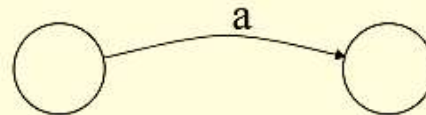
- The start state



- An accepting state



- A transition



Transition Diagrams

- Regular expressions are **declarative specifications**
- **Transition diagram is an implementation**

- An input alphabet belonging to Σ
- A set of states S
- A set of transitions $\text{state}_i \xrightarrow{\text{input}} \text{state}_j$
- A set of final states F
- A start state n

Transition $s1 \xrightarrow{a} s2$ is read:

- If end of input is reached in a final state, then accept
- Otherwise, reject

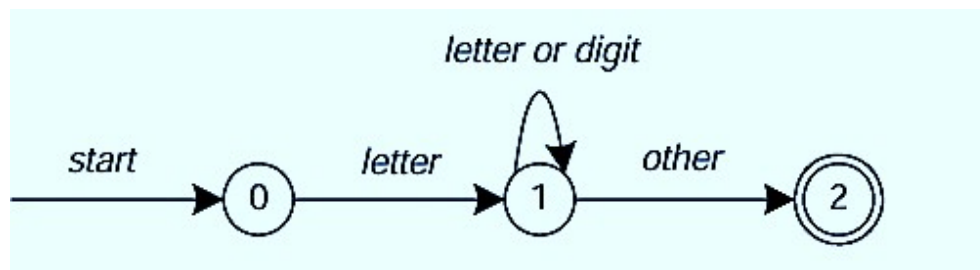
How to recognize tokens?

- **Identifier:** strings of letters or digits, starting with a letter

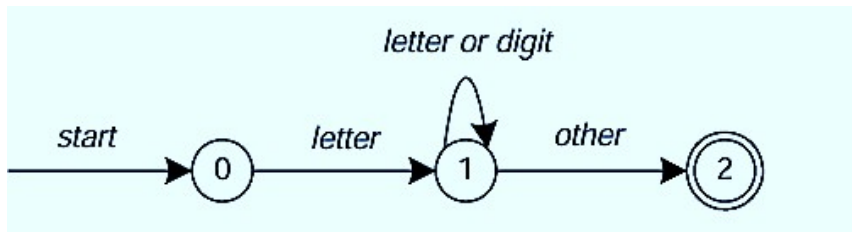
`letter(letter|digit)*`

- repetition, expressed by the “*” operator
- alternation, expressed by the “|” operator
- concatenation

- Any regular expression may be expressed as a finite state automaton (FSA).
- There is **one start state** and **one or more final or accepting states**.



Contd.,



Finite State Automaton



```
start:  goto state0

state0: read c
        if c = letter goto state1
        goto state0

state1: read c
        if c = letter goto state1
        if c = digit  goto state1
        goto state2

state2: accept string
```

3-state machine

How to recognize tokens

- Consider

relop \rightarrow < | <= | = | <> | >= | >

id \rightarrow letter(letter|digit)*

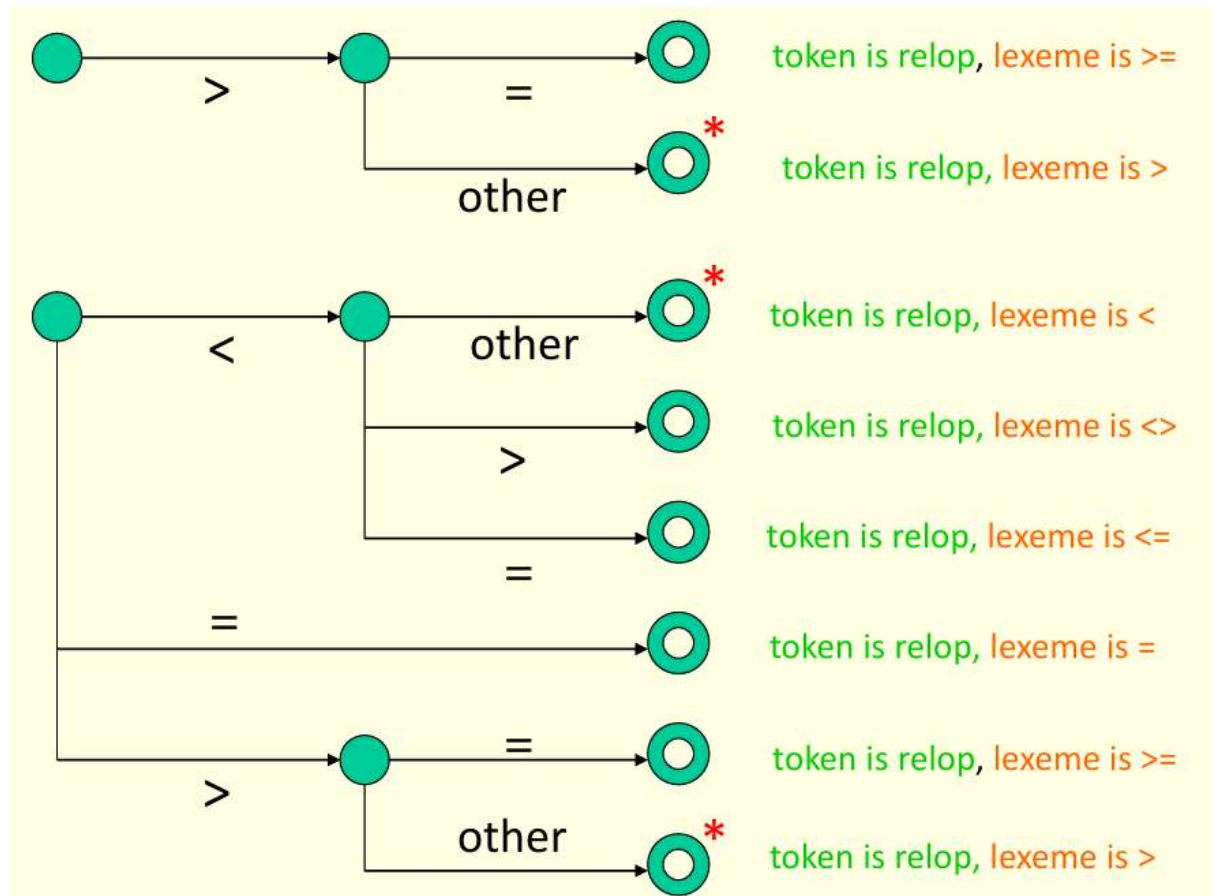
num \rightarrow digit⁺ ('.' digit⁺)? (E('+'|'-')? digit⁺)?

delim \rightarrow blank | tab | newline

ws \rightarrow delim⁺

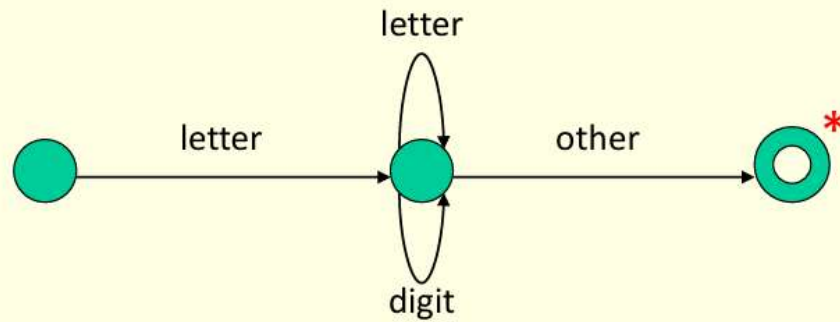
Construct an analyzer that will return <token, attribute> pairs

Transition diagram for relational operators

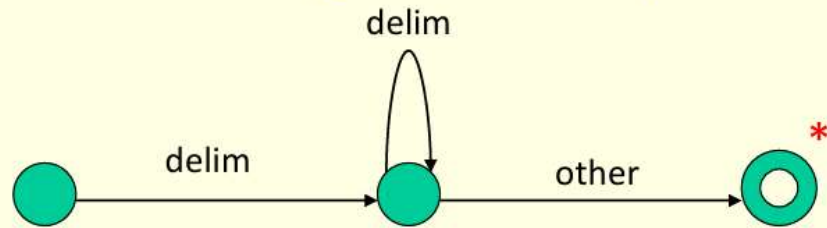


Contd.,

Transition diagram for identifier



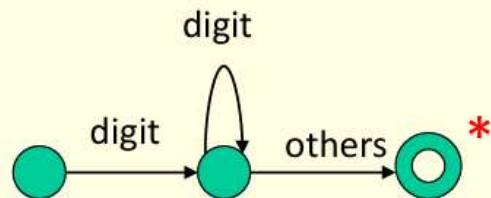
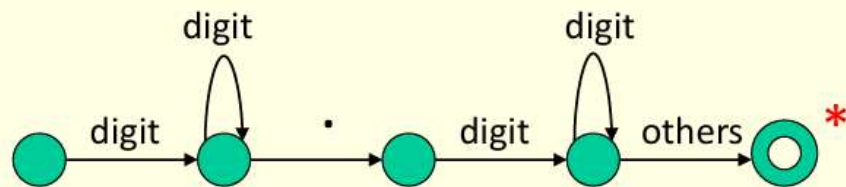
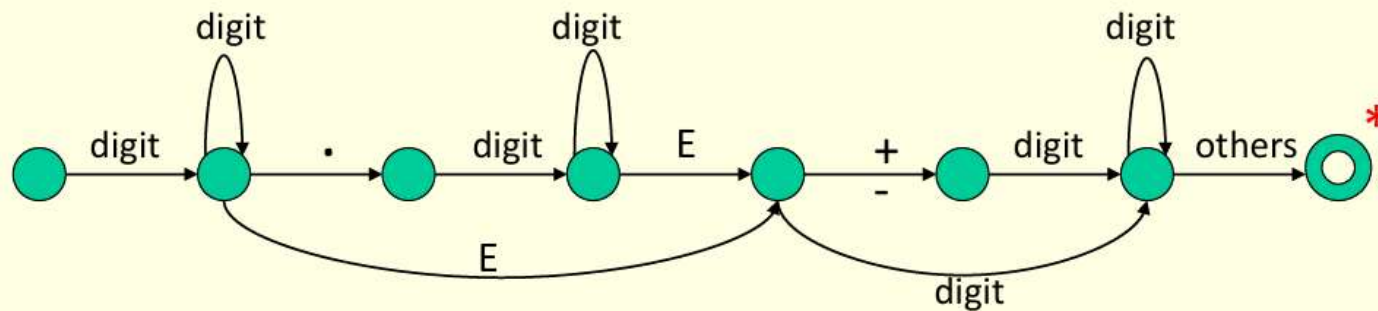
Transition diagram for white spaces



delim →→→→ blank | tab | newline

Contd.,

Transition diagram for unsigned numbers



Integer number

Real numbers

Recognizes real numbers with exponential notation.

Ex: **12.34E+56**

←←←←←←←←

Does **not** include exponential form.

Ex: **12.34**

←←←←←

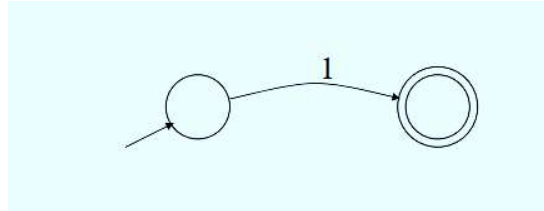
Simple DFA recognizing a sequence of digits followed by a non-digit.

Ex: **1234**

←←←←←←⁷⁵

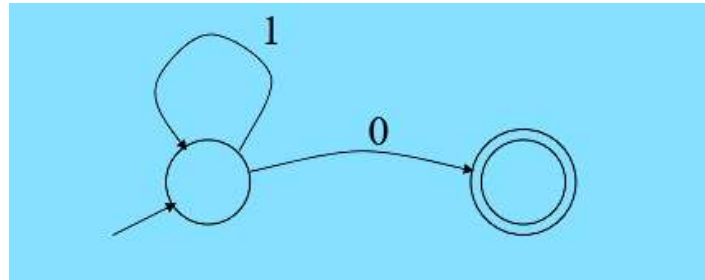
A Simple Example

- A finite automaton that accepts only “1”.



- Another Simple Example

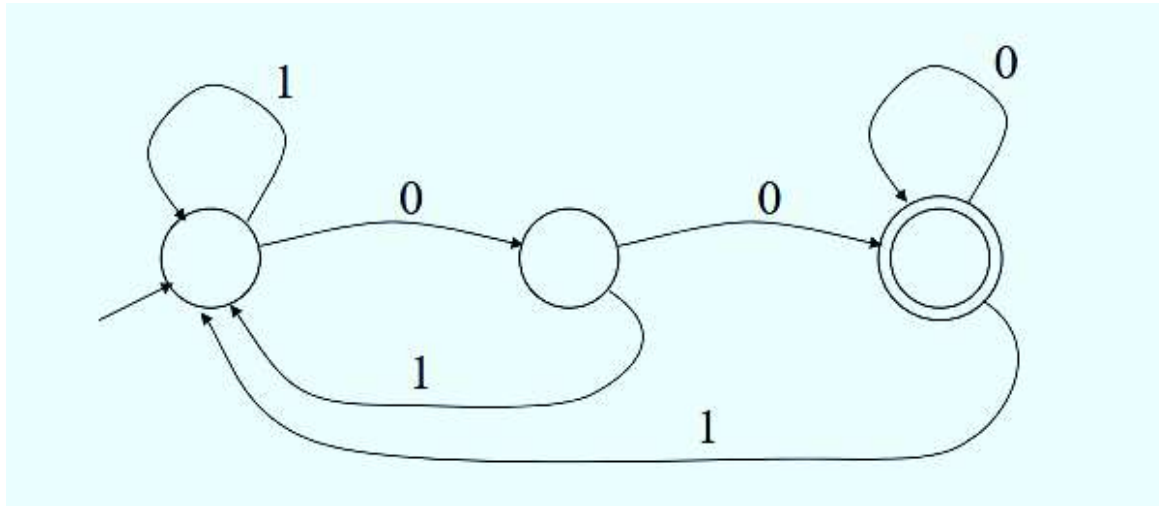
- A finite automaton accepting any number of 1's followed by a single 0
- Alphabet: {0,1}



Check that “1110” is accepted but “110...” is not

And Another Example

- Alphabet $\{0,1\}$
- What language does this recognize?

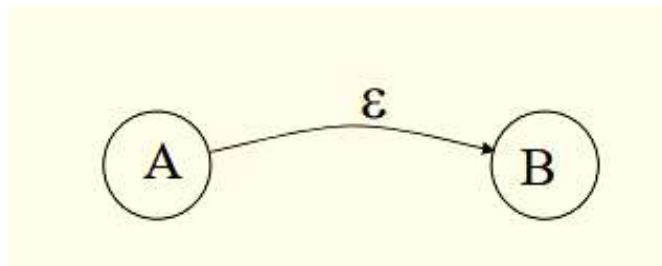


$\{0, 1\}^*$, it means it **accepts every possible binary string**, including the empty string.

ϵ , 0, 1, 00, 01, 10, 11, 000, 001, ..., 101010, etc.

Epsilon Moves

- Another kind of transition: ϵ -moves



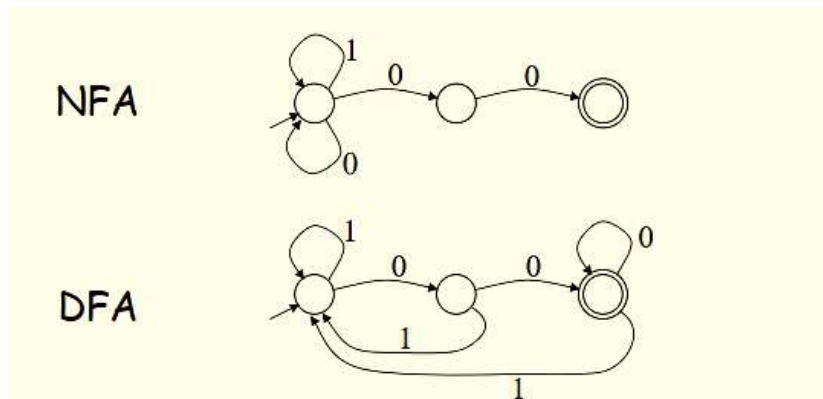
- Machine can move from state A to state B without reading input

Deterministic and Non-Deterministic Automata

- Deterministic Finite Automata (DFA)
 - One transition per input per state
 - No ϵ -moves
- Non-deterministic Finite Automata (NFA)
 - Can have multiple transitions for one input in a given state
 - Can have ϵ -moves
- Finite automata have finite memory
 - Enough to only encode the current state

NFA vs. DFA (1)

- NFAs and DFAs recognize the same set of languages (regular languages)
- DFAs are easier to implement
 - There are no choices to consider
- For a given language the NFA can be simpler than the DFA



Accepts strings like: "0 0", "1 0 0", "0 1 0", etc.

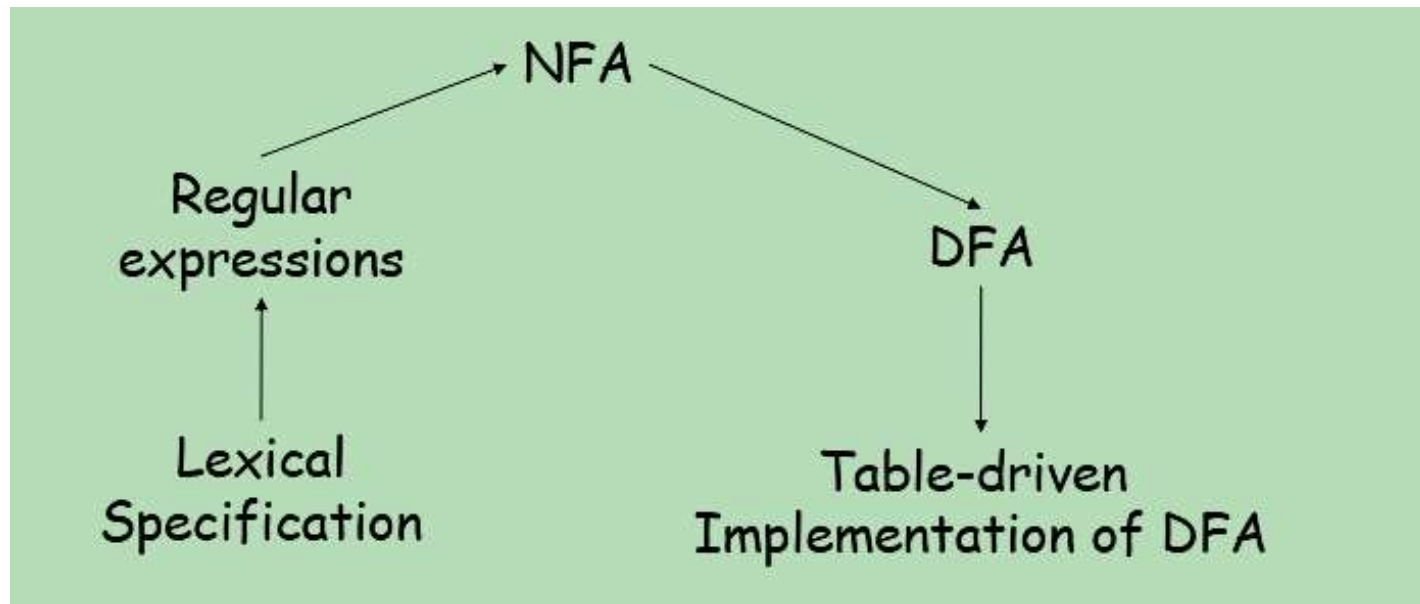
Each input at each state, Handles all possibilities deterministically.

Same language accepted, but every state has exactly one transition for each symbol (0 or 1).

- DFA can be exponentially larger than NFA

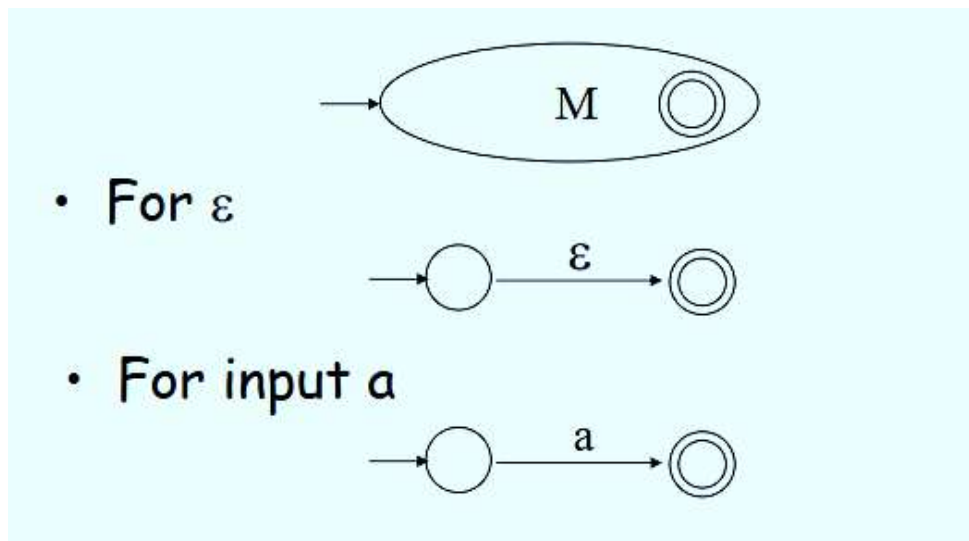
Regular Expressions to Finite Automata

- High-level sketch



Regular Expressions to NFA (1)

- For each kind of reg. expr, define an NFA: **Basic NFA fragments used in Thompson's construction** of regular expressions



- M represents a **general NFA fragment**

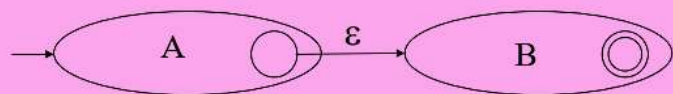
ϵ -NFA Fragment (For ϵ)

- Two states: a **start state** and a **final state** with an **ϵ -transition** (epsilon move) between them.
- This NFA accepts the **empty string (ϵ)** — i.e., without consuming any input character.

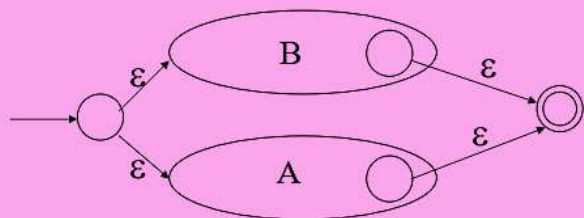
Basic Symbol NFA (For input a): This NFA accepts **only the string "a"**.

Regular Expressions to NFA (2)

- For AB



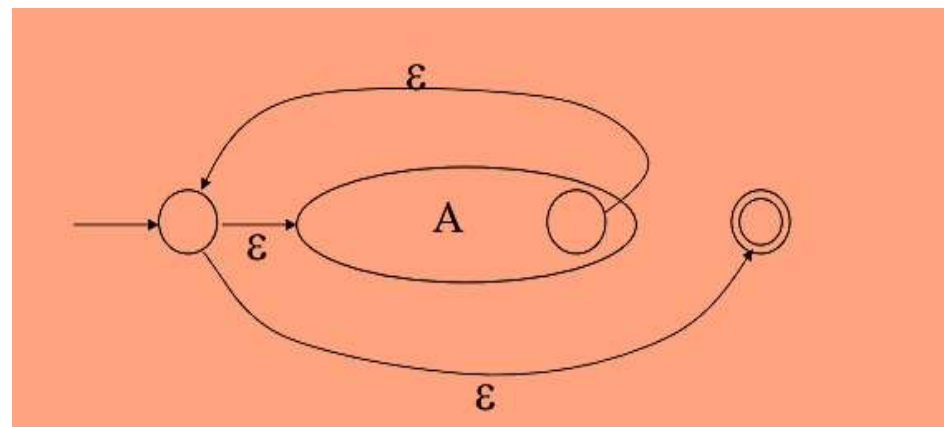
- For A + B



Regex	Meaning	NFA Mechanism
AB	Concatenation	$A \rightarrow \epsilon \rightarrow B$
A + B	Union	$\epsilon \rightarrow A$ and $\epsilon \rightarrow B$; $A \rightarrow \epsilon \rightarrow \text{Final}$, $B \rightarrow \epsilon \rightarrow \text{Final}$
A*	Zero/more of A	Loop: $\text{Start} \rightarrow \epsilon \rightarrow A \rightarrow \epsilon \rightarrow \text{Start}$; $\epsilon \rightarrow \text{Final}$ for empty

In **Thompson's construction**, each regular expression is broken into:

- Basic parts like a, b, ϵ ,
- And combined using operations: **Concatenation** (ab), **Union** (a | b), **Kleene star** (a*)



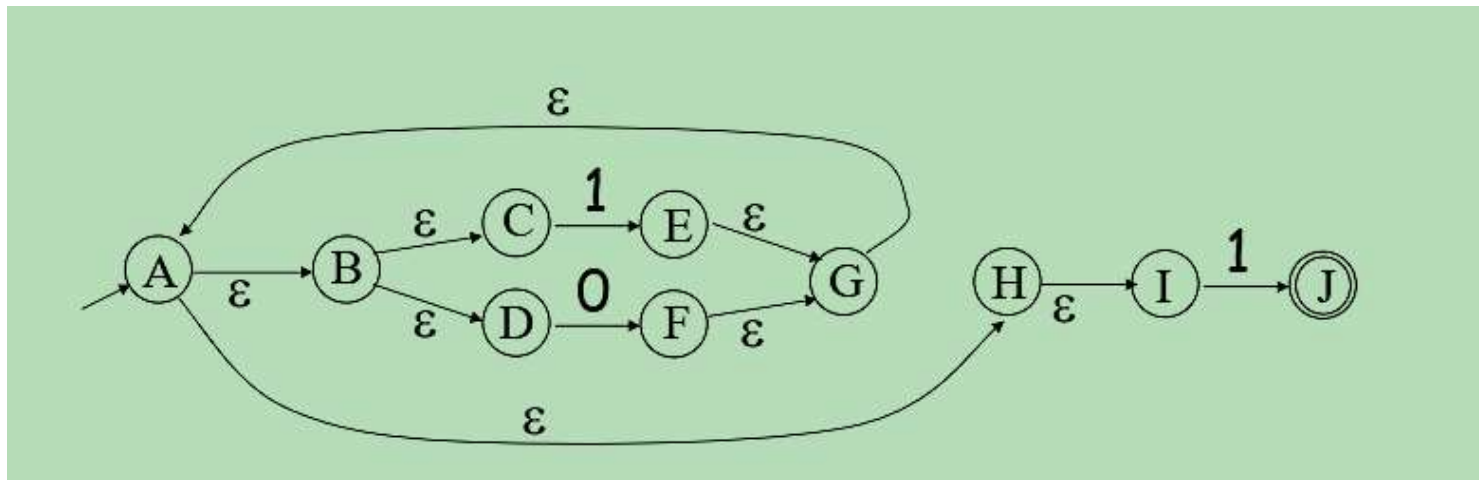
For A*

→→→→→

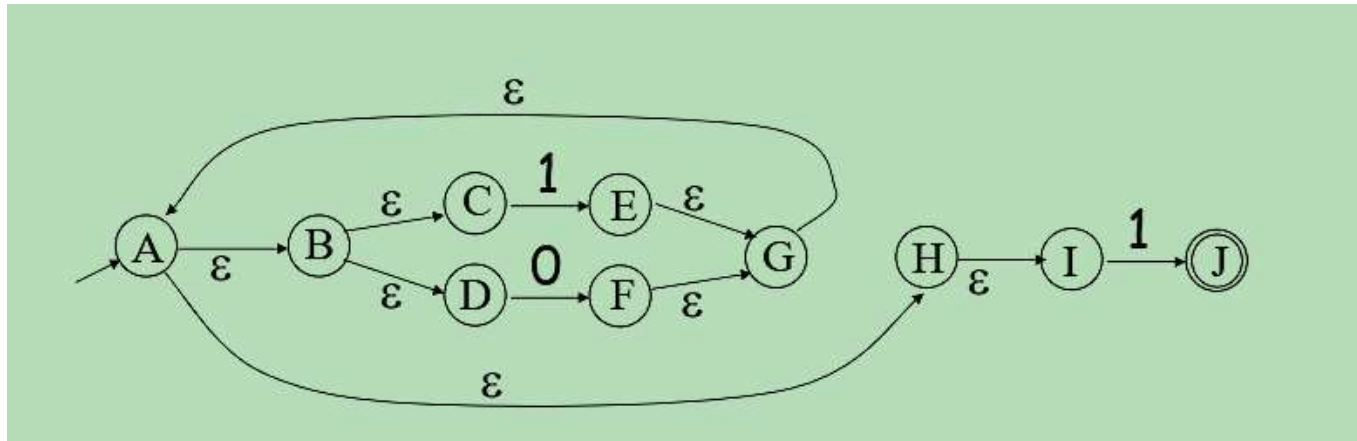
Example of Regular Expression \rightarrow NFA conversion

- Consider the regular expression: $(1+0)^*1$

The NFA is



Thompson's Construction \rightarrow ϵ -NFA



- A is the initial state
- Loop for $(1 + 0)^*$:
 - Through ϵ -transitions, it branches to:
 - $B \rightarrow \epsilon \rightarrow C \rightarrow 1 \rightarrow E \rightarrow \epsilon \rightarrow G$
 - $B \rightarrow \epsilon \rightarrow D \rightarrow 0 \rightarrow F \rightarrow \epsilon \rightarrow G$
 - Then $G \rightarrow \epsilon \rightarrow A$ to allow repetition.
- Final part $A \rightarrow \epsilon \rightarrow H \rightarrow \epsilon \rightarrow I \rightarrow 1 \rightarrow J$ (final)
- So it accepts any combination of 0s and 1s followed by a 1.
- This is the ϵ -NFA formed using **Thompson's Construction**.

Subset Construction \rightarrow DFA: NFA to DFA \rightarrow The Trick

- Simulate the NFA
- Each state of DFA = a non-empty subset of states of the NFA
- Start state = the set of NFA states reachable through ϵ -moves from NFA start state
- Add a transition $S \xrightarrow{a} S'$ to DFA iff
 - S' is the set of NFA states reachable from any state in S after seeing the input a
 - considering ϵ -moves as well

Example

Given ϵ -NFA \rightarrow States

- q_0 (start state)
- q_1
- q_2

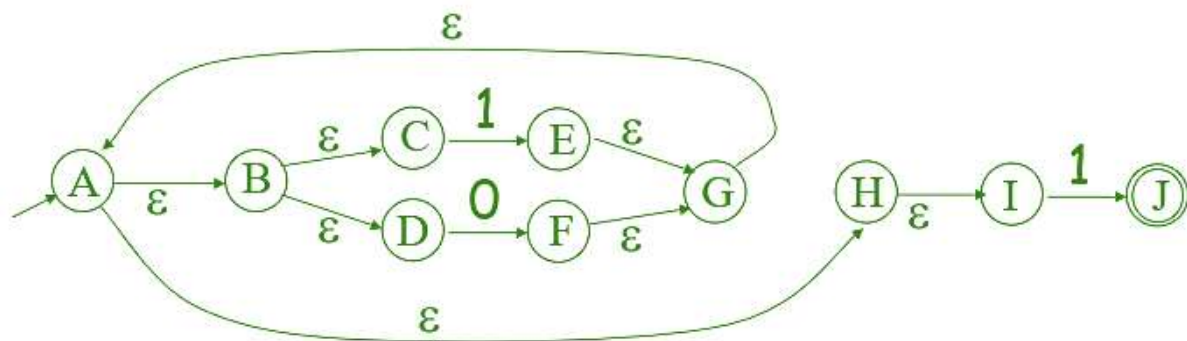
Transitions:

- $q_0 \xrightarrow{\epsilon} q_1$
- $q_1 \xrightarrow{\epsilon} q_2$
- $q_0 \xrightarrow{a} q_0$
- $q_1 \xrightarrow{b} q_1$

ϵ -Closures: The **ϵ -closure** of a state is the set of states reachable from it using **only ϵ -transitions**, including the state itself.

1. $\epsilon\text{-closure}(q_0) = \{q_0, q_1, q_2\}$
($q_0 \rightarrow q_1 \rightarrow q_2$ through ϵ -transitions)
2. $\epsilon\text{-closure}(q_1) = \{q_1, q_2\}$
($q_1 \rightarrow q_2$ through ϵ)
3. $\epsilon\text{-closure}(q_2) = \{q_2\}$
(no ϵ -transitions from q_2)

NFA to DFA Example

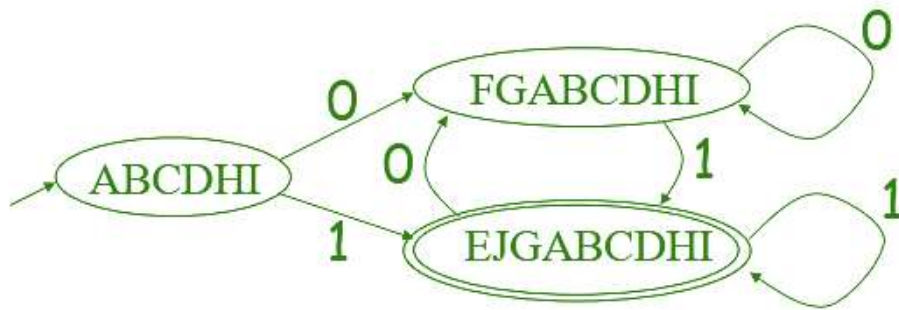


Step 1: ϵ -closure(A)

From A:

- $A \rightarrow \epsilon \rightarrow B, H$
- $B \rightarrow \epsilon \rightarrow C, D$
- $C \rightarrow 1 \rightarrow E$
- $D \rightarrow 0 \rightarrow F$
- $H \rightarrow \epsilon \rightarrow I$
- $(E, F \rightarrow \epsilon \rightarrow G), G \rightarrow \epsilon \rightarrow A$
- $I \rightarrow 1 \rightarrow J$

ϵ -closure(A) = {A, B, C, D, H, I}



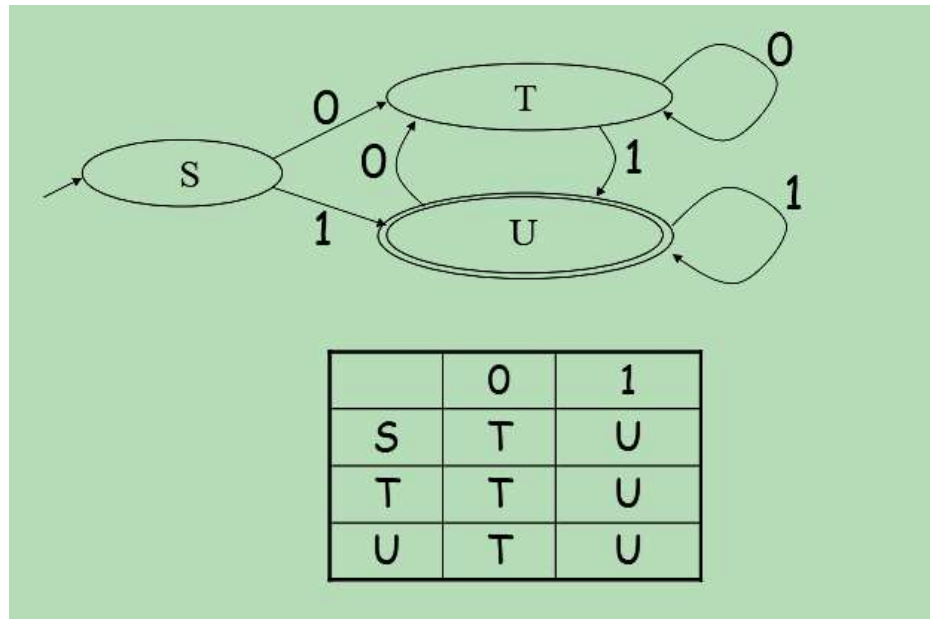
- On input 0 from $S_0 = \{A, B, C, D, H, I\}$:
- Only D has a 0-transition: $D \rightarrow 0 \rightarrow F$
- Then:
 - $F \rightarrow \epsilon \rightarrow G$
 - $G \rightarrow \epsilon \rightarrow A$
 - From $A \rightarrow \epsilon \rightarrow \{B, C, D, H, I\} \rightarrow$ and it loops

ϵ -closure(F) = {F, G, A, B, C, D, H, I}

Implementation

- A DFA can be implemented by a 2D table T
 - One dimension is “states”
 - Other dimension is “input symbols”
 - For every transition **state**_i $\xrightarrow{\text{input}(a)}$ **state**_k define $T[i,a] = k$
- DFA “execution”
 - If in **state**_i and input a, read $T[i,a] = k$ and skip to state **state**_k
 - Very efficient

Table Implementation of a DFA



- NFA → DFA conversion is at the heart of tools such as [lex](#), [ML-Lex](#) or [flex](#)
- But, DFAs can be huge
- In practice, [lex/ML-Lex/flex](#)-like tools trade off speed for space in the choice of NFA and DFA representations

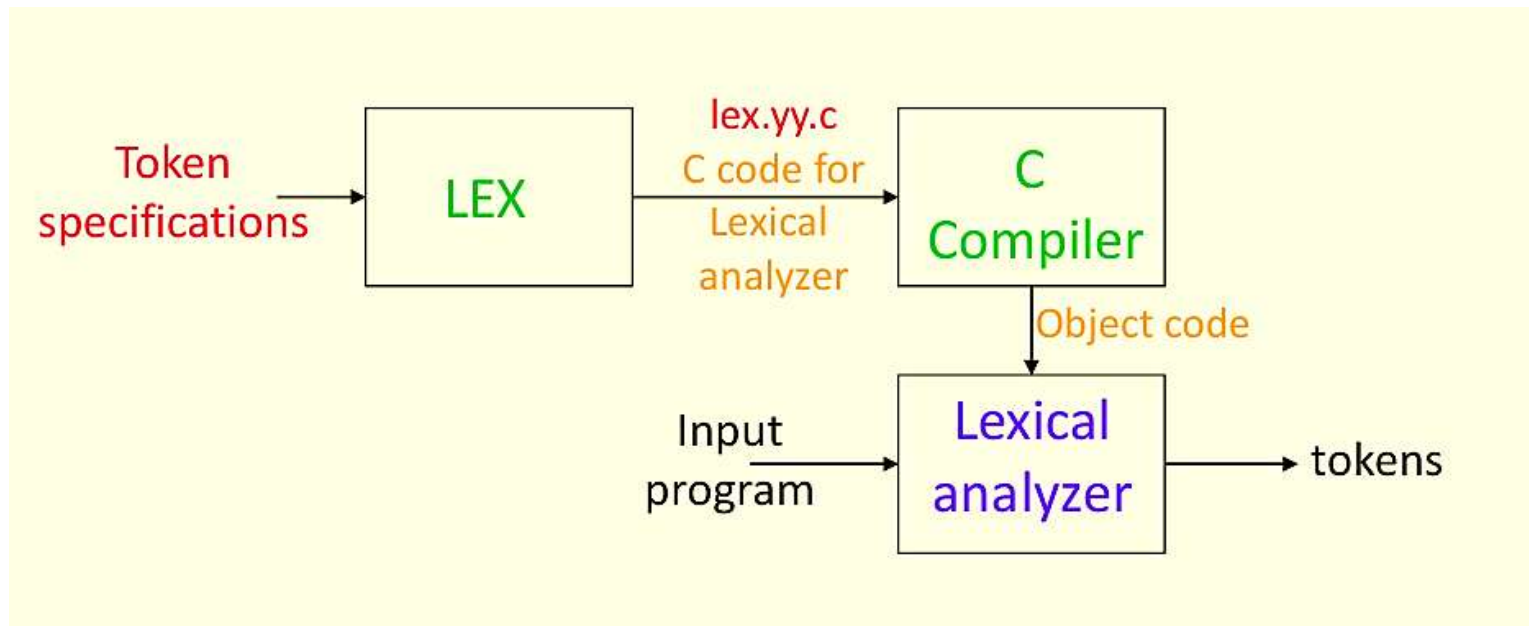
Theory vs Practice

- Two differences:
 - DFAs recognize lexemes.
 - A lexer must return a type of acceptance (token type) rather than simply an accept/reject indication.
 - DFAs consume the complete string and accept or reject it.
 - A lexer must find the end of the lexeme in the input stream and then find the next one, etc.

Lexical analyser generator

- Input to the generator
 - List of regular expressions in priority order
 - Associated actions for each of regular expression
(generates kind of token and other keeping information)
- Output of the generator
 - Program that reads input character stream and breaks that into tokens
 - Reports lexical errors (unexpected characters), if any

LEX: A lexical analyzer generator



How does LEX work?

- Regular expressions describe the languages that can be recognized by finite automata
- Translate each token regular expression into a non deterministic finite automaton (NFA)
- Convert the NFA into an equivalent DFA
- Minimize the DFA to reduce number of states
- Emit code driven by the DFA tables