

Lexical Specification

- ① at least one: $A^+ \equiv AA^*$
- ② Union: $A \cup B \equiv A + B$
- ③ Option: $A^? \equiv A + \epsilon$
- ④ Range: 'a' + 'b' + ... + 'z' $\equiv [a-z]$
- ⑤ Excluded range:
complement of $[a-z] \equiv [^a-z]$

Based on $S \in L(LR)$, we need some extension!

1. Write a regex for the lexemes of each token class.

- Number = digit⁺
- Keyword = 'if' + 'else' + ...
- Identifier = letter(letter + digit)*
- OpenPar = '('
- ...

2. Construct R , matching all lexemes for all tokens.

$$R = \text{Keyword} + \text{Identifier} + \text{Number} + \dots \\ = R_1 + R_2 + \dots$$

3. Let input be $x_1 \dots x_n$

For $1 \leq i \leq n$ check

$$x_1 \dots x_i \in L(R)$$

4. If success, then we know that

$$x_1 \dots x_i \in L(R_j) \text{ for some } j$$

5. Remove $x_1 \dots x_i$ from input and go to (3)

Regular expressions are a concise notation for string patterns.

Use in lexical analysis requires small extensions

- To solve ambiguities - [matches as long as possible]
[highest priority match]
- To handle errors - put them last and longest priority

Good algorithms known

- Require only single pass over the input
- Few operations per character (table lookup)

Finite Automata

Regular Expression = Specification

Finite Automata = Implementation

A finite automaton consists of

- An input alphabet: Σ
- A finite set of states: S
- A start state: n
- A set of accepting states: $F \subseteq S$
- A set of transitions: state $\xrightarrow{\text{input}} \text{state}$

Transition

$$S_1 \xrightarrow{a} S_2$$

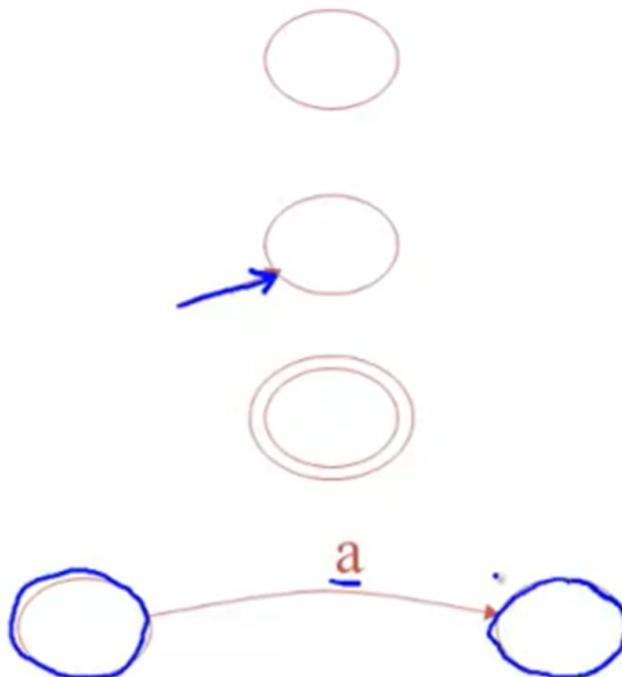
Is read In state $\underline{S_1}$ on input \underline{a} go to state
 $\underline{S_2}$

If end of input and in accepting state \Rightarrow accept

Otherwise \Rightarrow reject

Terminates in state $S \notin F$
gets stuck

- A state
- The start state
- An accepting state
- A transition

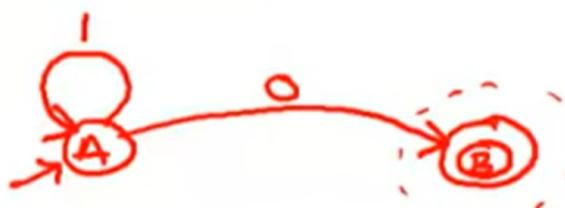


Language of a FA = set of accepted strings

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

<u>State</u>	<u>Input</u>
A	↑100
A	↓00
B	10↑0

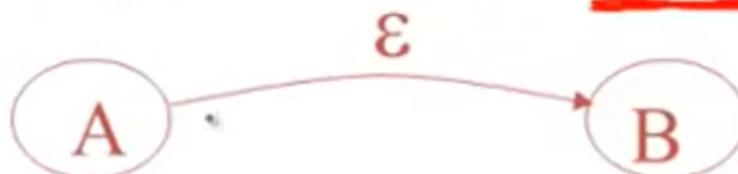
Reject



<u>State</u>	<u>Input</u>
A	↑10
A	↓10
A	11↑0
B	110↑

Accepts

- Another kind of transition: ϵ -moves



<u>State</u>	<u>Input</u>
A	$x_1 \ x_2 \ x_3$ ↑
B	$x_1 \ x_2 \ x_3$ ↑

Deterministic Finite Automata (DFA)

- One transition per input per state
- No ϵ -moves

Nondeterministic Finite Automata (NFA)

- Can have multiple transitions for one input in a given state.
- Can have ϵ -moves.

A DFA takes only one path through the state graph.

An NFA can choose, meaning it can accept if some choices lead to an accepting state.

NFAs and DFAs recognize the same set of languages.

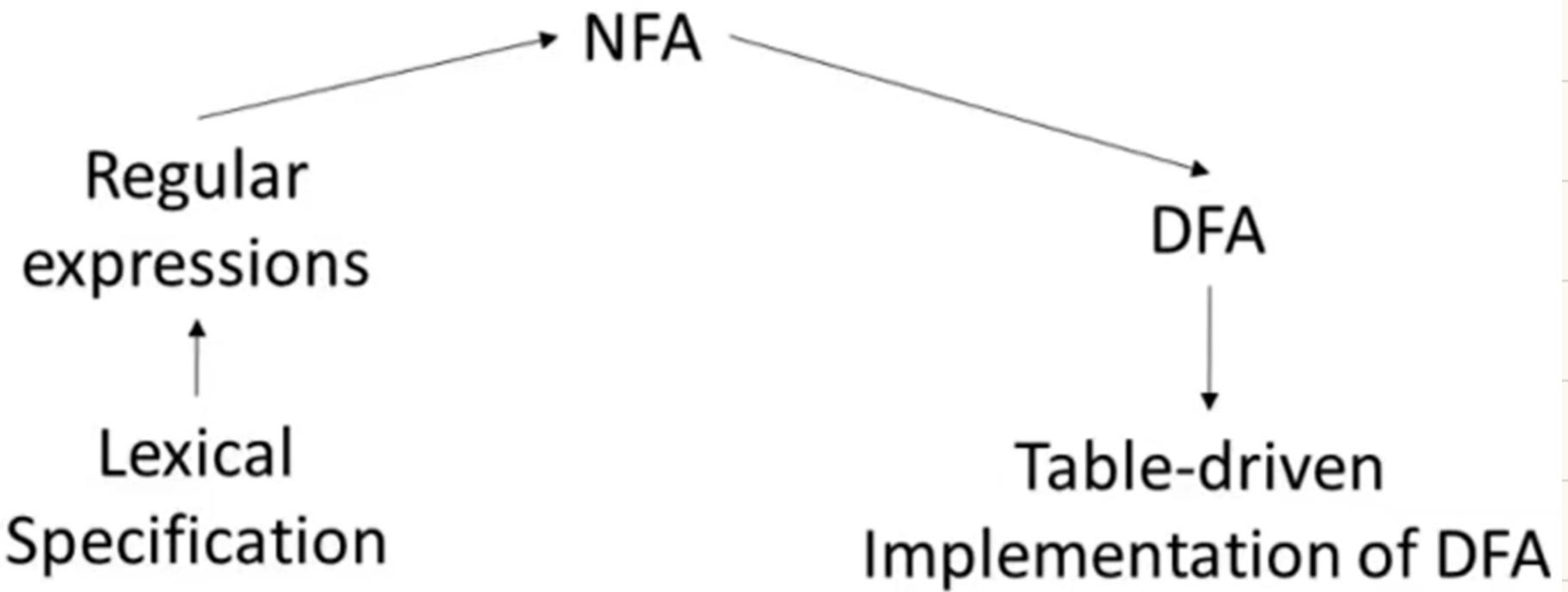
- Regular languages

DFAs are faster to execute

- There are no choices to consider.

NFAs are, in general, smaller.

Regular Expressions to NFAs



For each kind of rexp , define an equivalent NFA

- Notation: NFA for $\text{rexp } M$



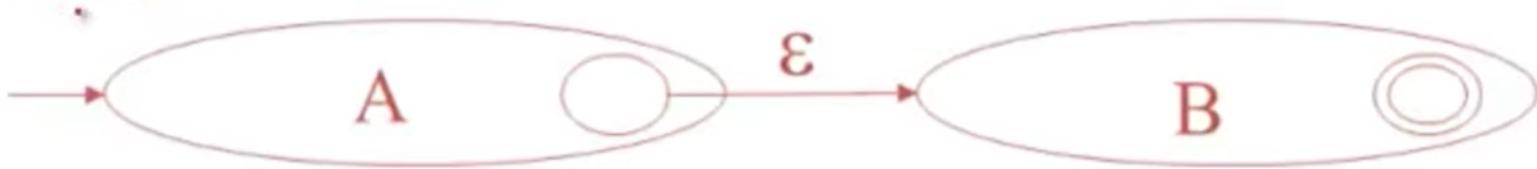
- For ϵ



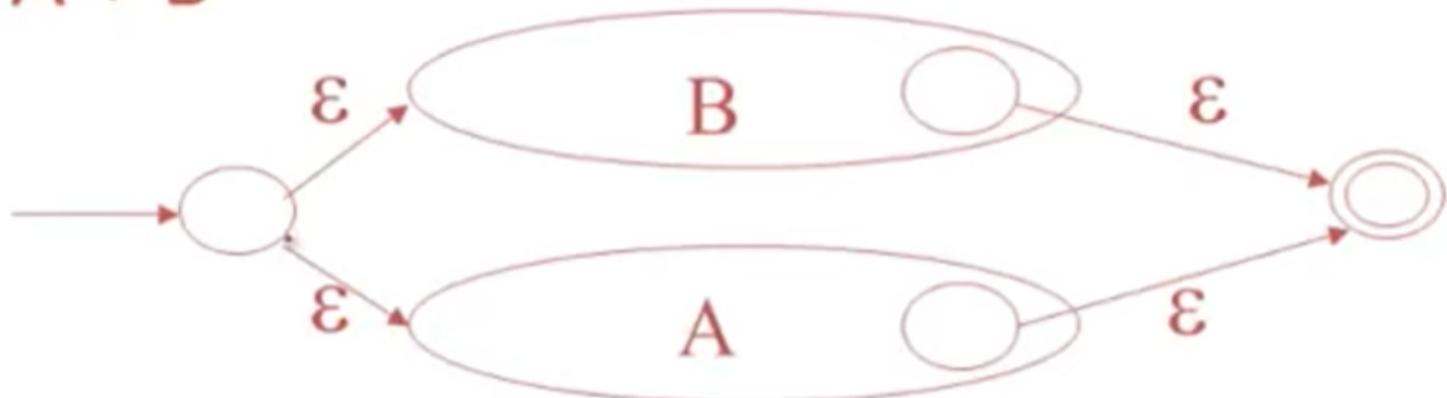
- For input a



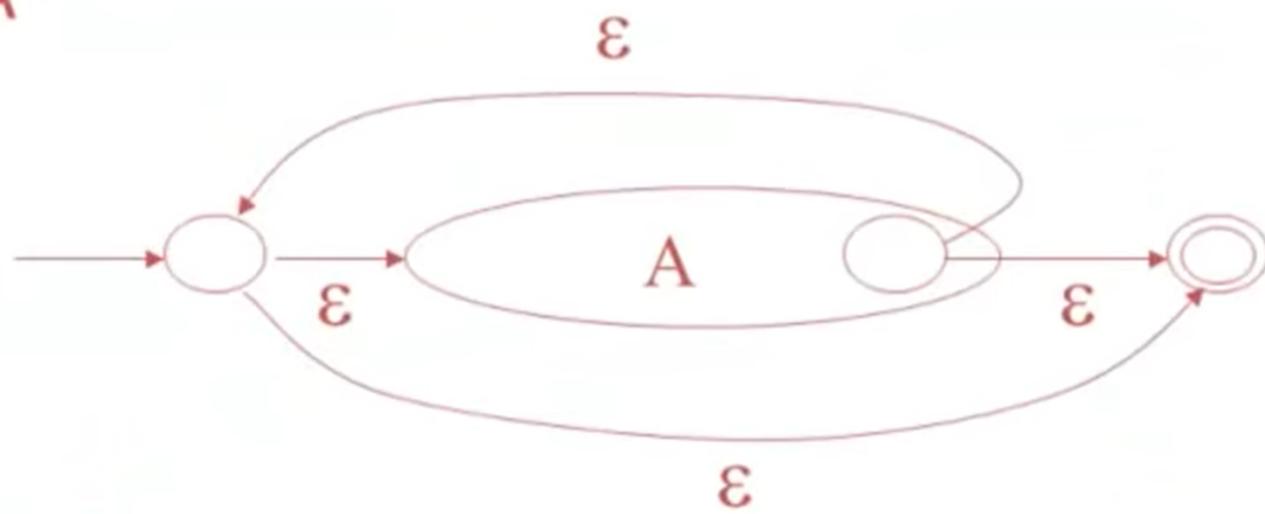
- For AB



- For $A + B$

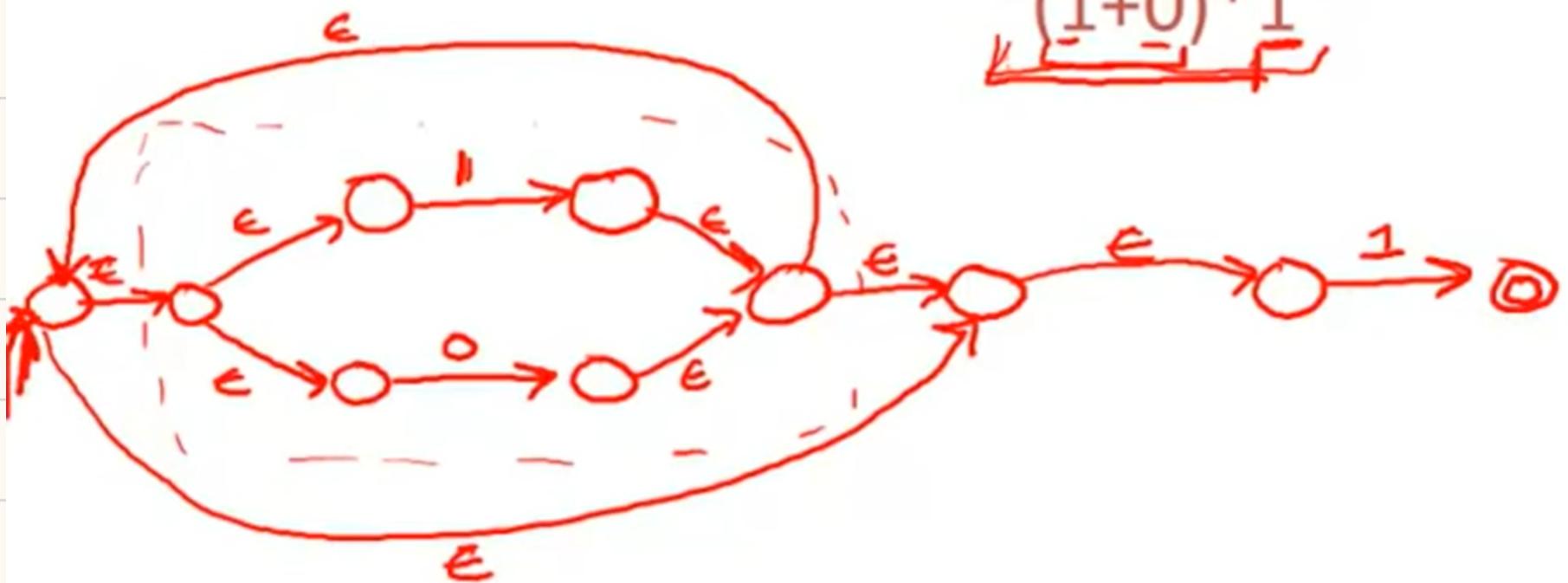


- For A^*



- Consider the regular expression

$$(1+0)^*1$$



NFA to DFA

An NFA may be in many states at any time.

NFA

states

S

start

$s \in S$

final

$F \subseteq S$

$\alpha(x) = \{y \mid s \xrightarrow{\epsilon} x_n \xrightarrow{a} y\}$

ϵ -clos

DFA

States

subsets of S (except the empty set)

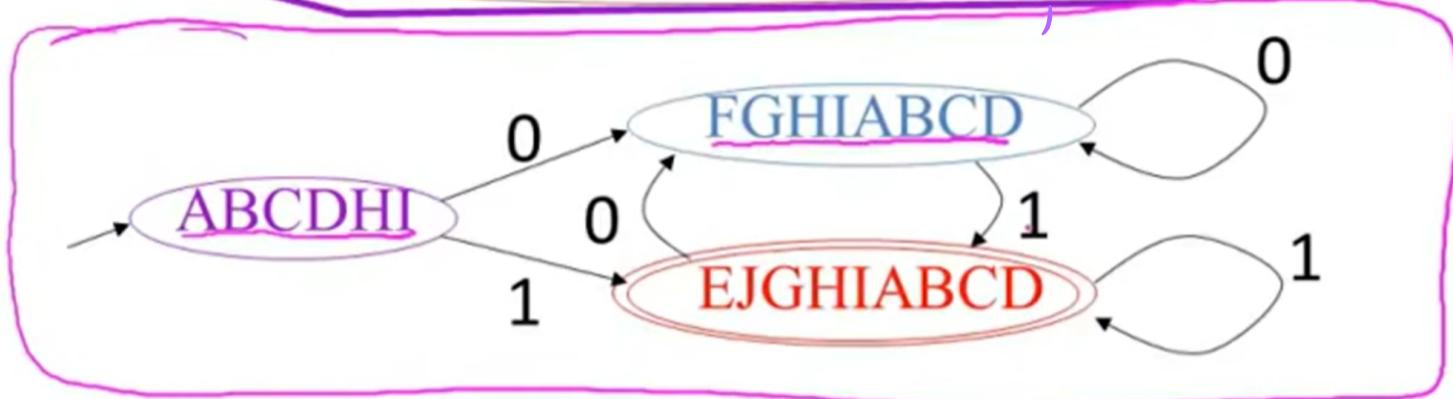
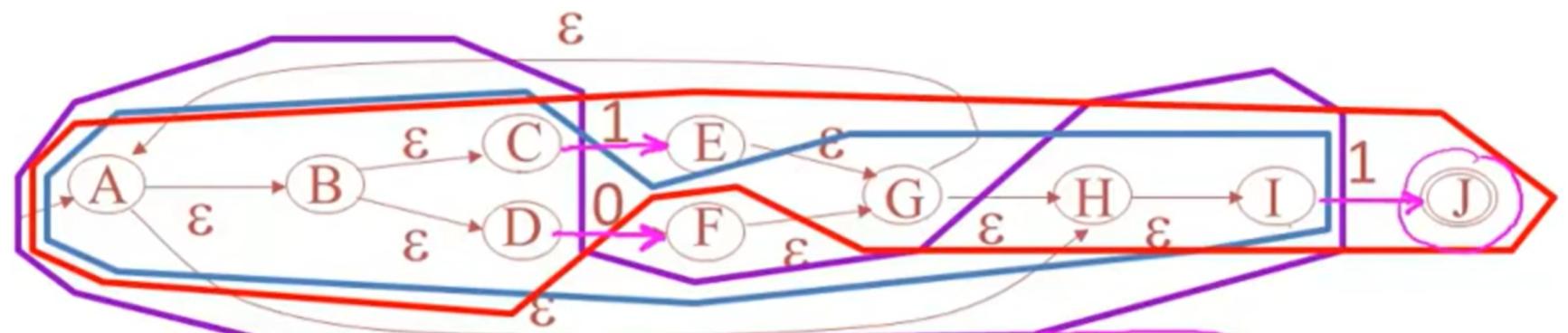
start

ϵ -clos(s)

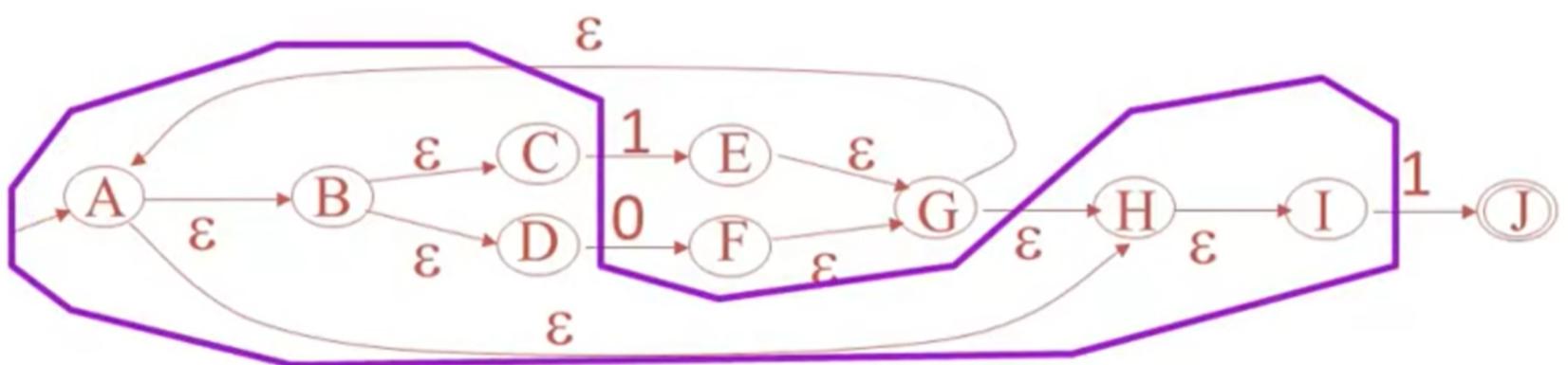
final

$\{x \mid x \cap F \neq \emptyset\}$

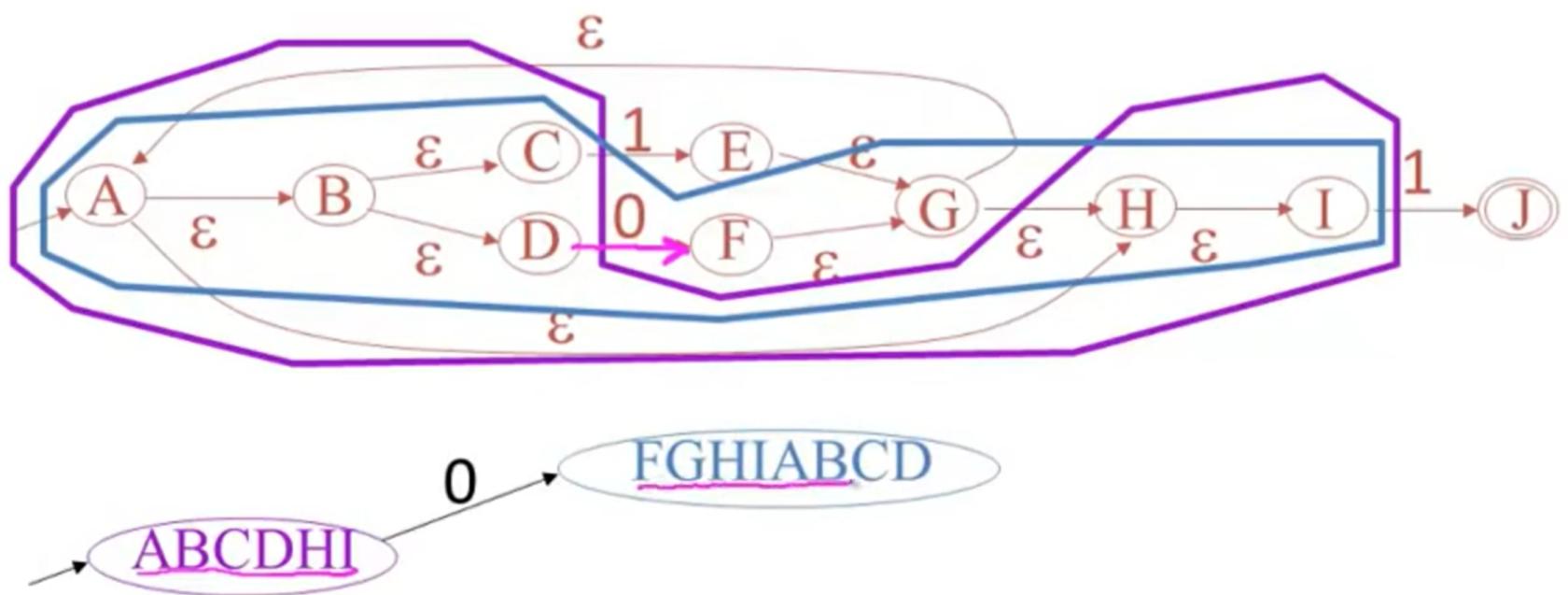
$X \xrightarrow{a} Y$ if $Y = \epsilon\text{-clos}(\alpha(X))$

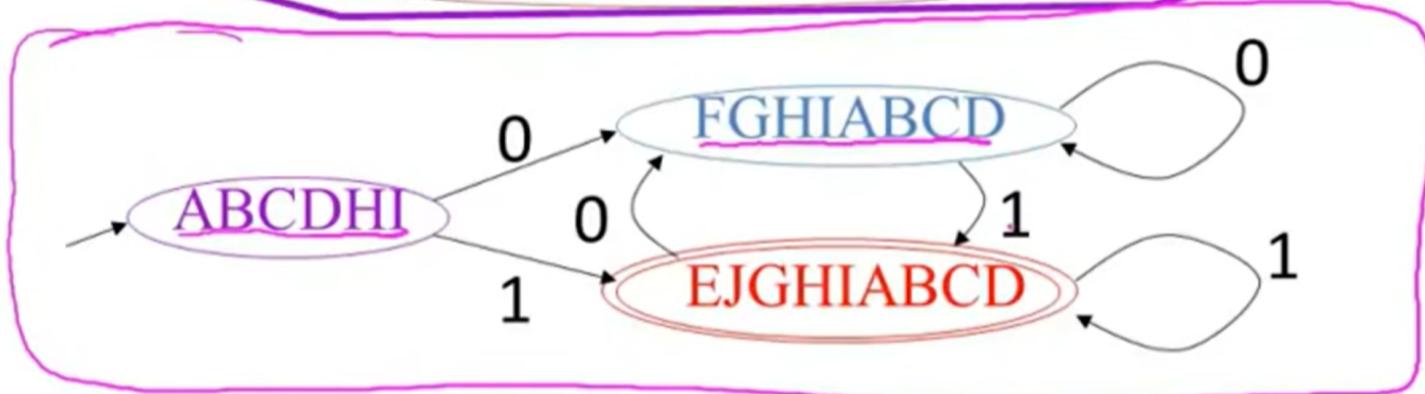
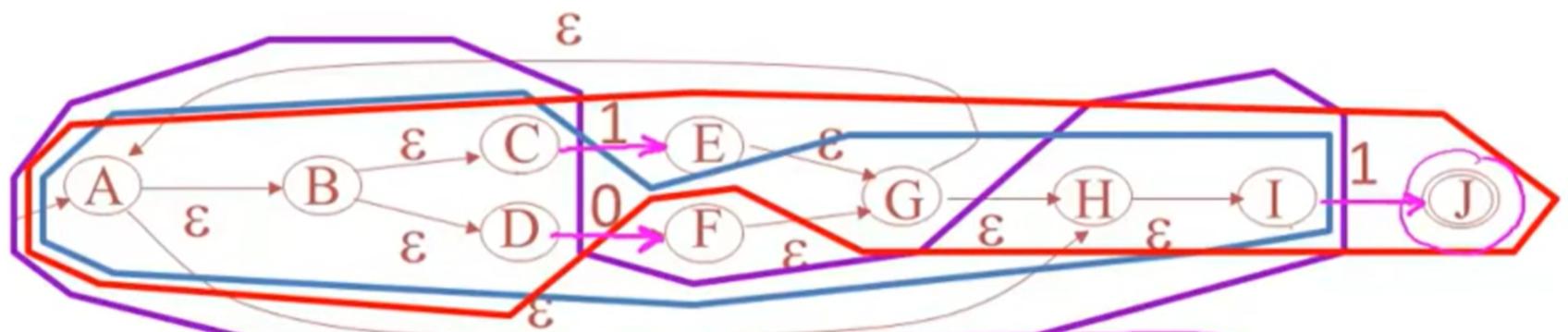
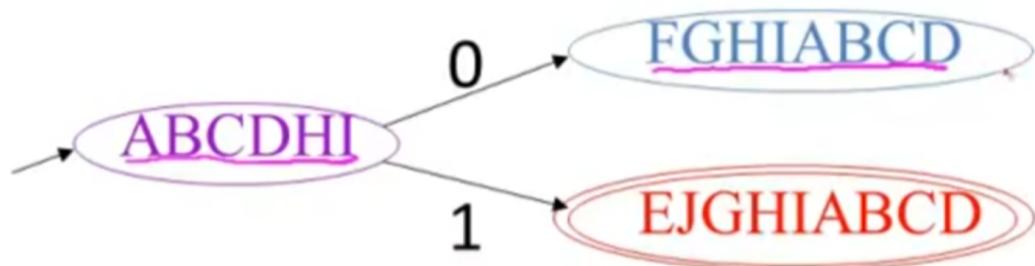
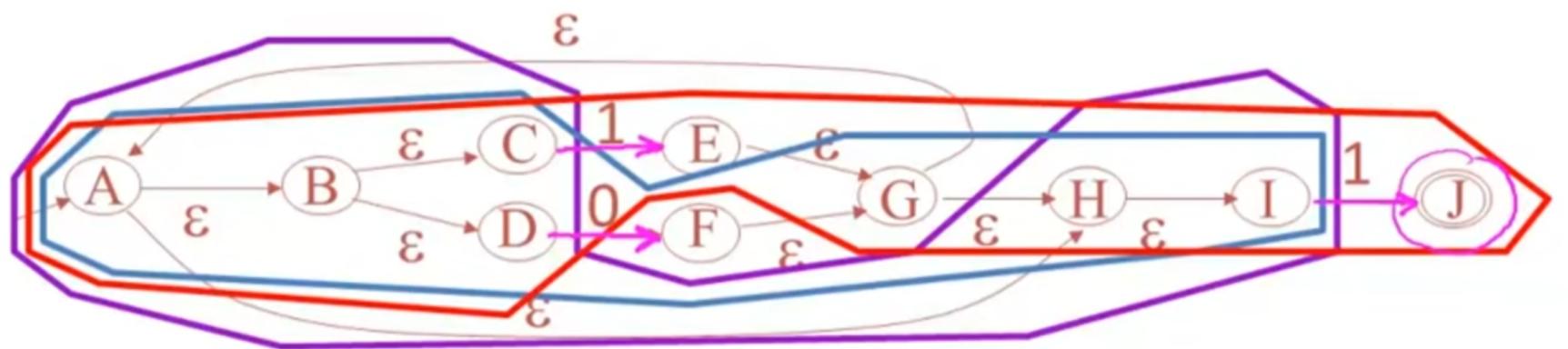


step-by-step analyse



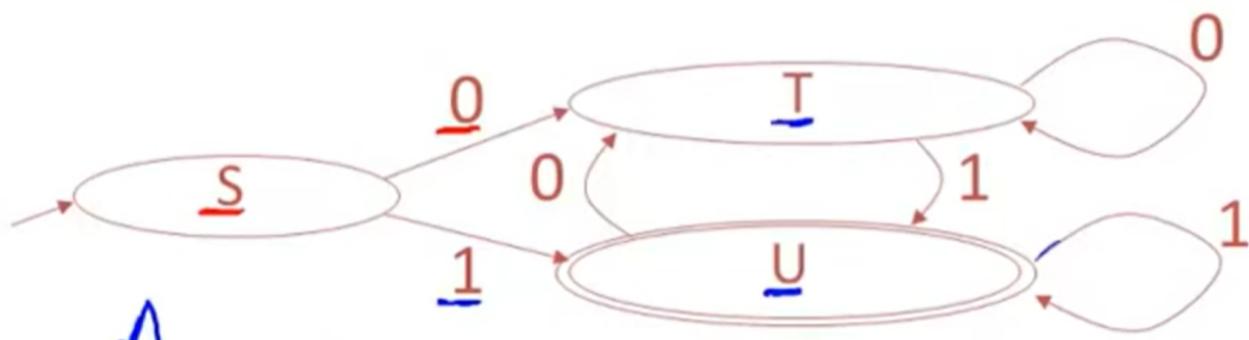
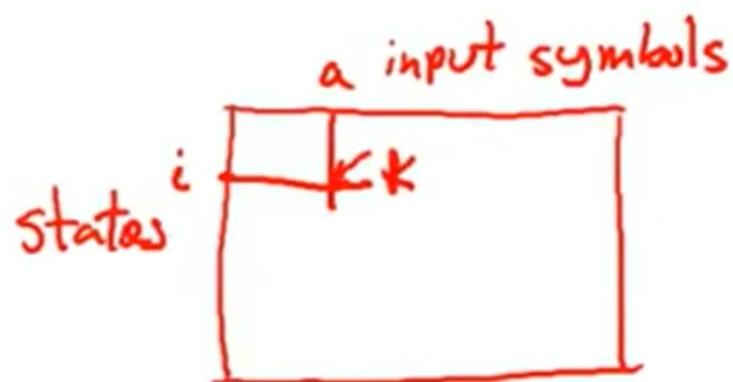
DFA start state
 $\epsilon\text{-clos}(A)$





Implementing Finite Automata

- A DFA can be implemented by a 2D table T
 - One dimension is states
 - Other dimension is input symbol
 - For every transition $S_i \xrightarrow{a} S_k$ define $T[i,a] = k$



A	B	C
S	T	U
T	T	U
U	T	U

```
i = 0;  
state = 0;  
while (input[i]) {  
    state = A[state, input[i++]];  
}
```

NFA \rightarrow DFA conversion is key.

Tools trade between speed and space

DFA's faster, less compact

NFA's slower, concise