

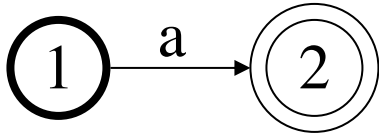
# Lexical Analysis

CMPT 379: Compilers

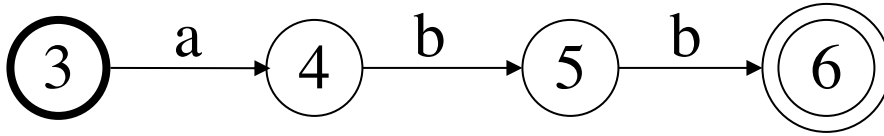
Instructor: Anoop Sarkar

[anoopsarkar.github.io/compilers-class](https://anoopsarkar.github.io/compilers-class)

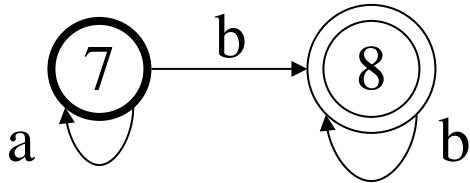
# Lexical Analysis using NFAs



TOKEN\_A = a

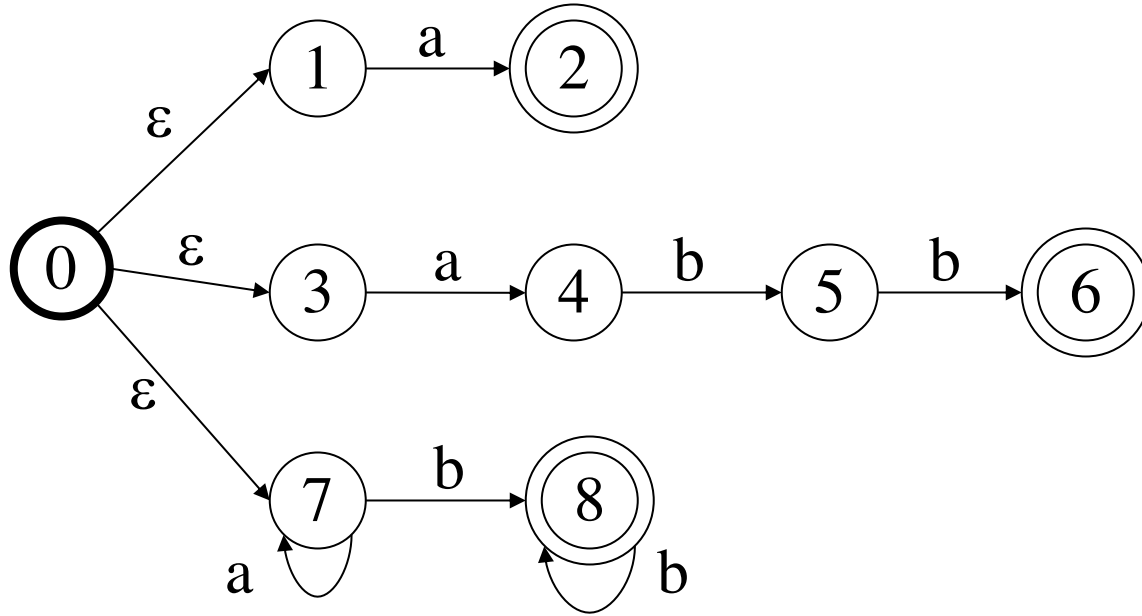


TOKEN\_B = abb



TOKEN\_C = a\*b+

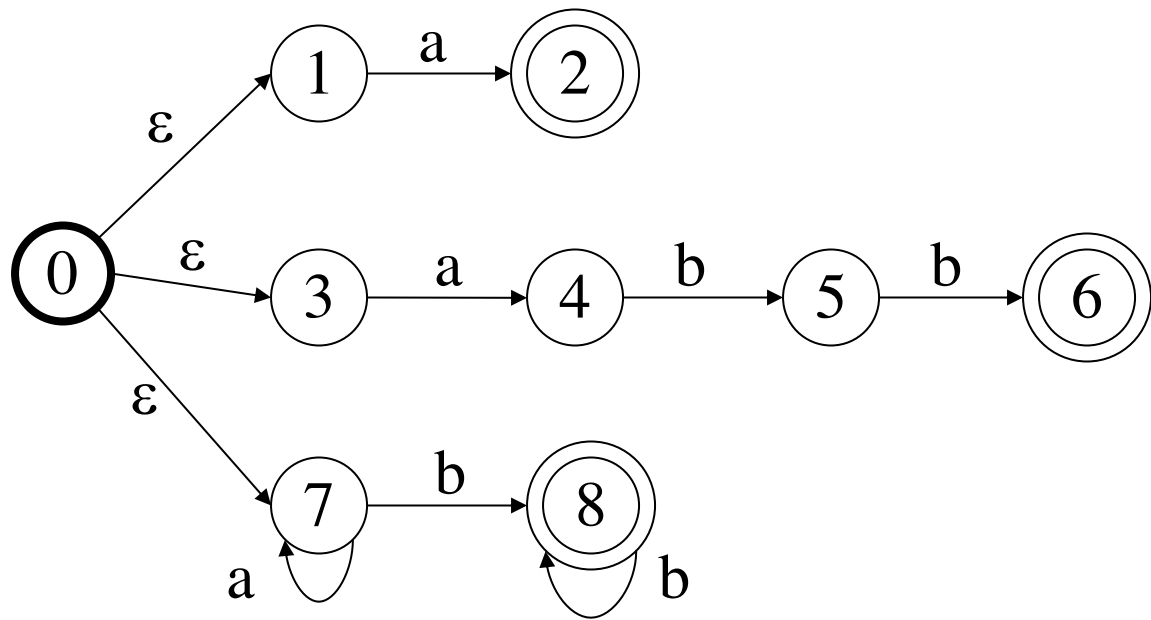
# Lexical Analysis using NFAs



TOKEN\_A = a

TOKEN\_B = abb

TOKEN\_C = a\*b+



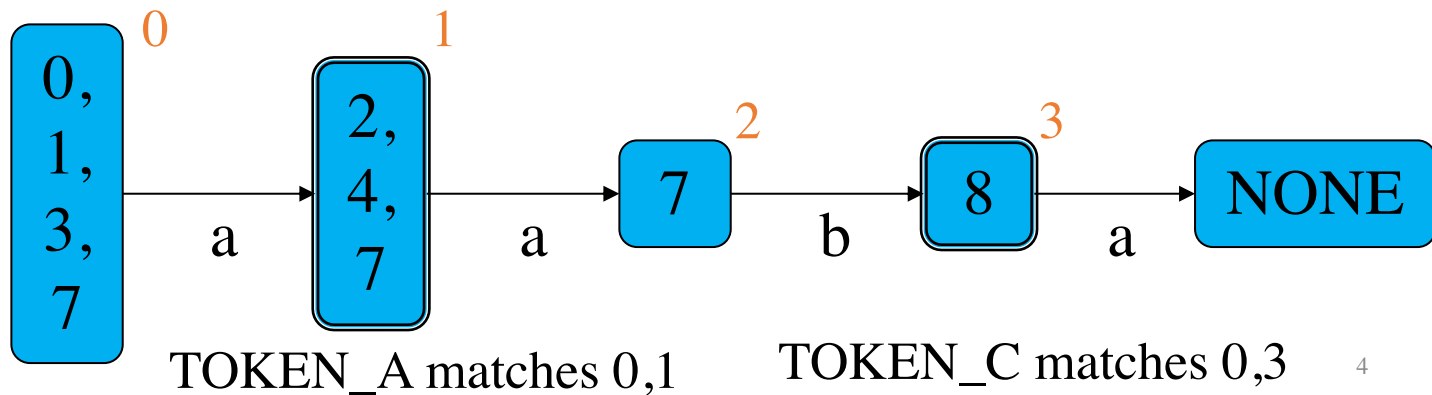
TOKEN\_A = a

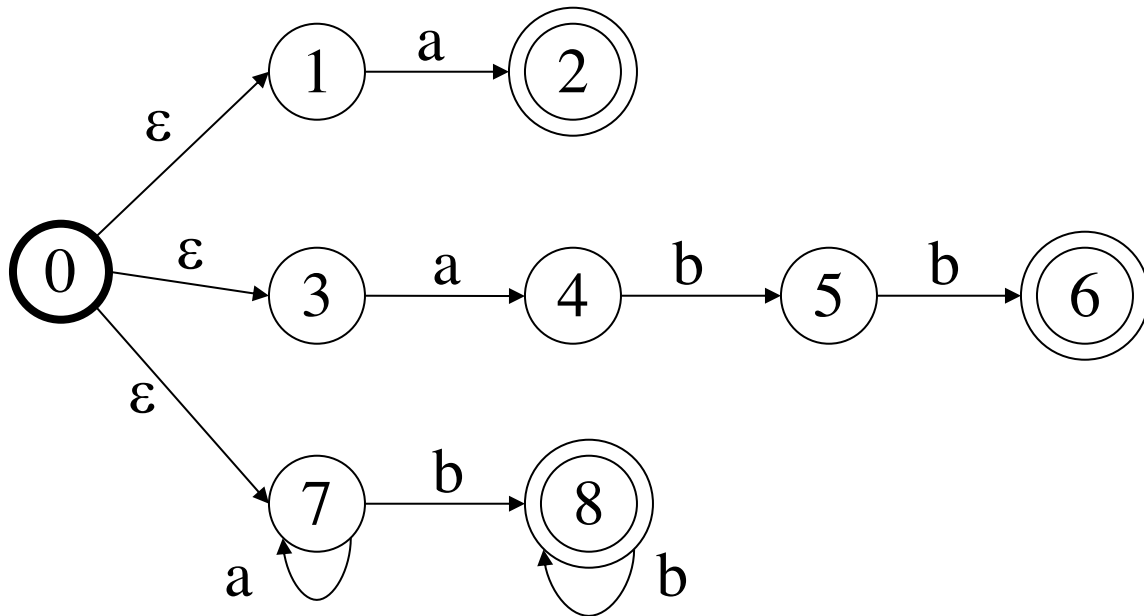
TOKEN\_B = abb

TOKEN\_C = a\*b+

Input: aaba

<sub>0</sub>a<sub>1</sub>a<sub>2</sub><sub>3</sub>b<sub>4</sub>a<sub>5</sub>





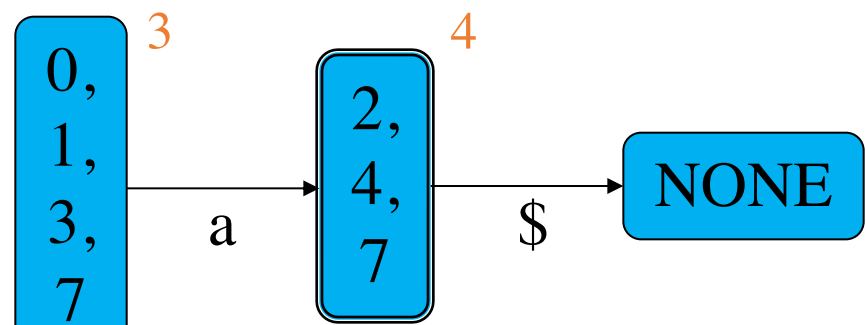
TOKEN\_A = a

TOKEN\_B = abb

TOKEN\_C = a\*b+

Input: aaba

<sub>0</sub>a<sub>1</sub>a<sub>2</sub><sub>3</sub>b<sub>4</sub>a



TOKEN\_A matches 3,4

Output:  
 TOKEN\_C aab [0,3]  
 TOKEN\_A a [3,4]

# Lexical Analyzer using DFAs

- Each token is defined using a regexp  $r_i$
- Merge all regexps into one big regexp
  - $R = (r_1 \mid r_2 \mid \dots \mid r_n)$
- Convert  $R$  to an NFA, then DFA, then minimize
  - remember original NFA final states with each DFA state

# Lexical Analyzer using DFAs

- The DFA recognizer must find the *longest leftmost match* for a token
  - continue matching and report the last final state reached once DFA simulation cannot continue
  - e.g. longest match: `<print>` and not `<pr>`, `<int>`
  - e.g. leftmost match: for input string `aabaaaaab` the regexp `a+b` will match `aab` and not `aaaaab`
- If two patterns match the same token, pick the one that was listed earlier in R
  - e.g. prefer final state (in the original NFA) of  $r_2$  over  $r_3$

# Lookahead operator

- Implementing  $r_1/r_2$  : match  $r_1$  when followed by  $r_2$
- e.g.  $a^*b+/a^*c$  accepts a string  $bac$  but not  $abd$
- The lexical analyzer matches  $r_1 \epsilon r_2$  up to position  $q$  in the input
- But remembers the position  $p$  in the input where  $r_1$  matched but not  $r_2$
- Reset to start state and start from position  $p$



TOKEN\_A = (ab)\*a

TOKEN\_B = (ab)\*a(ca)\*

TOKEN\_C = bab(bab)\*

TOKEN\_D = a\*ba(ba)\*

Q: Use the ordered token definitions shown here and provide the tokenized output for the input string *abacabababa* using the greedy longest match lexical analysis method.

# Summary

- Token  $\Rightarrow$  Pattern
  - Pattern  $\Rightarrow$  Regular Expression
  - Regular Expression  $\Rightarrow$  NFA
    - Thompson's Rules
  - NFA  $\Rightarrow$  DFA
    - Subset construction
  - DFA  $\Rightarrow$  minimal DFA
    - Minimization
- $\Rightarrow$  Lexical Analyzer (multiple patterns)**

Extra Slides

TOKEN\_A = (ab)\*a

TOKEN\_B = (ab)\*a(ca)\*

TOKEN\_C = bab(bab)\*

TOKEN\_D = a\*ba(ba)\*

Q: Use the ordered token definitions shown here and provide the tokenized output for the input string *abacabababa* using the greedy longest match lexical analysis method.

A:

TOKEN\_B (abaca)

TOKEN\_D (bababa)