# BU CS 332 – Theory of Computation

## Lecture 6:

- Regexes = NFAs
- Limitations of Finite Automata

Reading:

Sipser Ch 1.3

"Myhill-Nerode" note

Mark Bun

February 10, 2025

# Regular Expressions – Syntax

A regular expression $R$ is defined recursively using the following rules:

1.  $\varepsilon$, $\emptyset$, and $a$ are regular expressions for every $a \in \Sigma$

2.  If $R_1$ and $R_2$ are regular expressions, then so are
    $(R_1 \cup R_2)$, $(R_1 \circ R_2)$, and $(R_1^*)$

Examples: (over $\Sigma = \{a, b, c\}$)     (with simplified notation)
$$ab \qquad\qquad ab^*c \cup (a^*b)^* \qquad\qquad \emptyset$$

# Regular Expressions – Semantics

$L(R)$ = the language a regular expression describes

1. $L(\emptyset) = \emptyset$
2. $L(\varepsilon) = \{\varepsilon\}$
3. $L(a) = \{a\}$ for every $a \in \Sigma$
4. $L((R_1 \cup R_2)) = L(R_1) \cup L(R_2)$
5. $L((R_1 \circ R_2)) = L(R_1) \circ L(R_2)$
6. $L\big((R_1^*)\big) = (L(R_1))^*$

Example: $L(a^* b^*) = \{a^m b^n \mid m, n \geq 0\}$

# Syntactic Sugar

- For alphabet $\Sigma$, the regex $\Sigma$ represents $L(\Sigma) = \Sigma$   $= \{a, b, c\}$

  e.g.  if $\Sigma = \{a, b, c\}$
  
  then regular expression $\Sigma$ stands for $(a \cup b \cup c)$

- For regex $R$, the regex $R^+ = RR^*$

  $R^+$ is interpreted as "one or more occurrence of a string generated by $R$"

# Regexes in the Real World

`grep` = globally search for a <u>r</u>egular <u>e</u>xpression and <u>p</u>rint matching lines

```
$ grep '^xy*z' myfile
xyz
xyzde
xzz
xz
xyyz
xyyyz
xyyyyz
$ grep '^x.*z' myfile
xyz
xyzde
xxz
xzz
x\z
x*z
xz
x z
xYz
xyyz
xyyyz
xyyyyz
$ grep '^x\*z' myfile
x*z
$ grep '\\' myfile
x\z
$
```

# Regular Expressions Describe Regular Languages

Theorem: A language $A$ is regular if and only if it is described by a regular expression ← *recognized by a DFA*
⟺ *recognized by an NFA*

Theorem 1: Every regular expression has an equivalent NFA

Theorem 2: Every NFA has an equivalent regular expression

Not regular: $\{a^n b^n \mid n \geq 0\}$
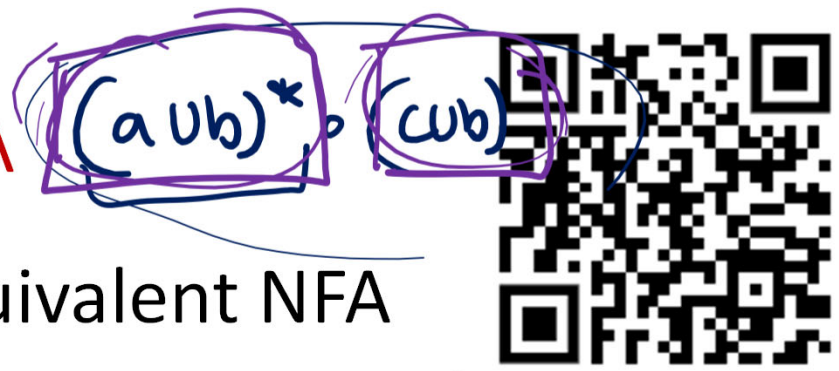
Regular: $\{a^n b^n \mid 0 \leq n \leq 2025\}$

# Regular expression -> NFA

Theorem 1: Every regex has an equivalent NFA

Proof: Induction on size of a regex

Base cases:

$L(R)$      NFA

$R = \emptyset$      $\phi$      $\rightarrow \bigcirc$

$R = \varepsilon$      $\{\varepsilon\}$      $\rightarrow \circledcirc$

$R = a$      $\{a\}$      $\rightarrow \bigcirc \xrightarrow{a} \circledcirc$

# Regular expression → NFA $(a \cup b)^* \cup (c \cup b)$

**Theorem 1:** Every regex has an equivalent NFA

Proof: Induction on size of a regex  = # of symbols, i.e.

$\phi, \varepsilon, a, (, ), \circ, \cup, *$

WTS statement is true for regexes of size $k+1$

size = length

What should the inductive hypothesis be?

a) Suppose **some** regular expression of length $k$ can be converted to an NFA

b) Suppose **every** regular expression of length $k$ can be converted to an NFA

c) Suppose **every** regular expression of length **at most** $k$ can be converted to an NFA

d) None of the above

# Regular expression → NFA

Theorem 1: Every regex has an equivalent NFA

Proof: Induction on size of a regex

IH: Assume every regex of size ≤k has an equiv. NFA
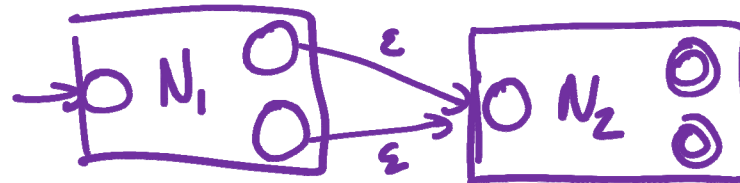WTS if R is a regex of size k+1, then R has an equiv. NFA

Inductive step:

$$R = (R_1 \cup R_2)$$
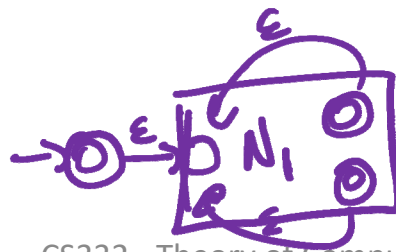
$L(R) = L(R_1) \cup L(R_2)$

By IH, ∃ NFA $N_1$ recognizing $R_1$, $N_2$ recognizing $R_2$

$$R = (R_1 R_2)$$

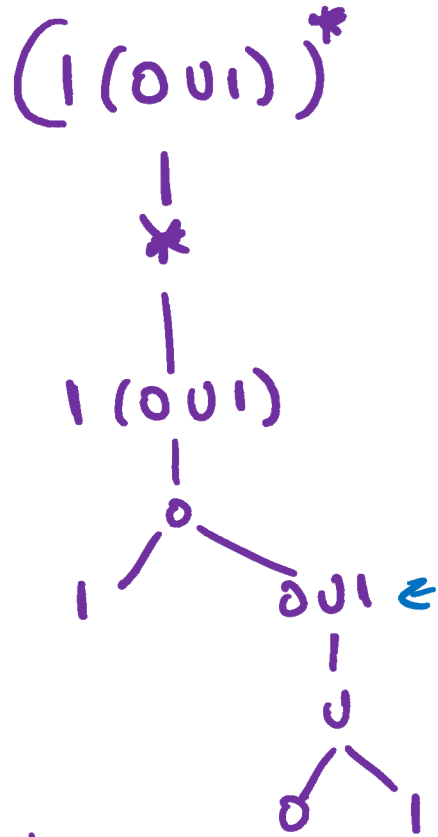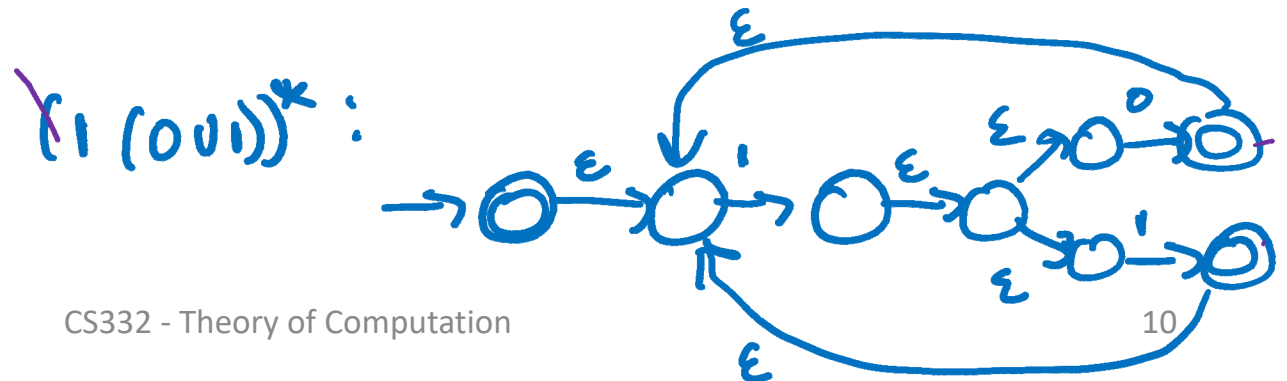$$R = (R_1^*)$$

# Example

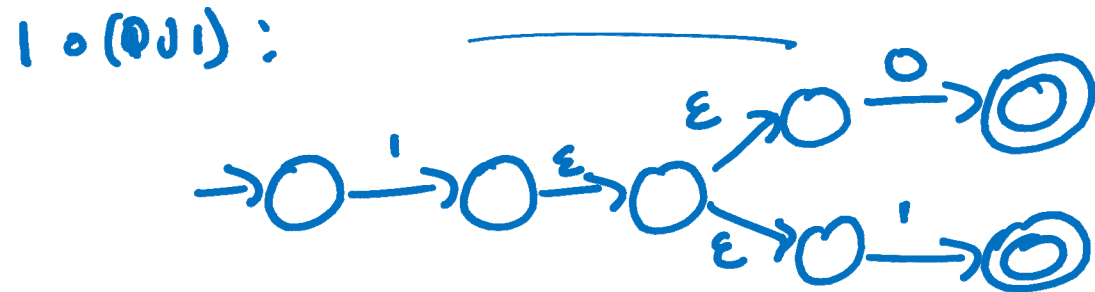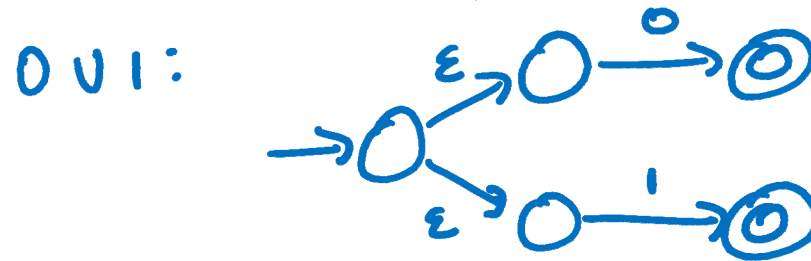$$(1(0 \cup 1))^*$$

e.g. $1\,0\,1\,1\,1\,0\,1\,1$     $1\,0\,1\,0\,1\,0$

= { even length strings w/ 1 in every odd position }

$1(0 \cup 1)$

0

$1$    $0 \cup 1 \; \varepsilon$

    $\cup$

$0$   $1$

0: →○ —0→ ◎

1: →○ —1→ ◎

$0 \cup 1$:

$1 \circ (0 \cup 1)$:

$(1(0 \cup 1))^*$:

Simplify to:



    CS332 - Theory of Computation    

# Regular Expressions Describe Regular Languages

**Theorem:** A language $A$ is regular if and only if it is described by a regular expression
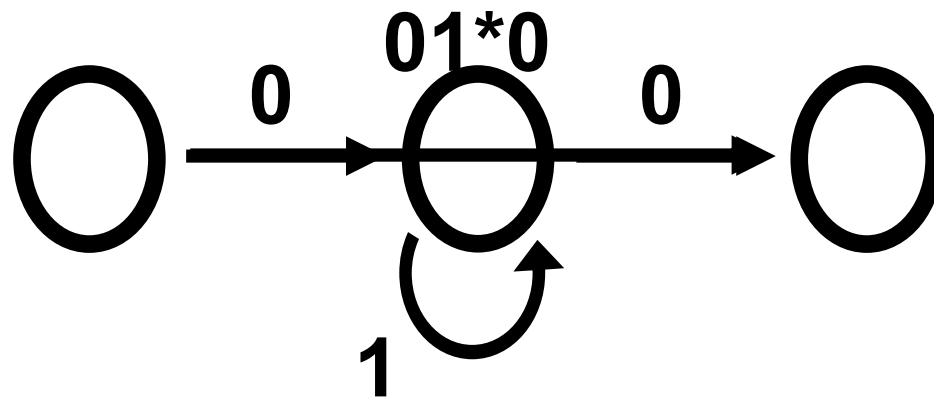
**Theorem 1:** Every regular expression has an equivalent NFA

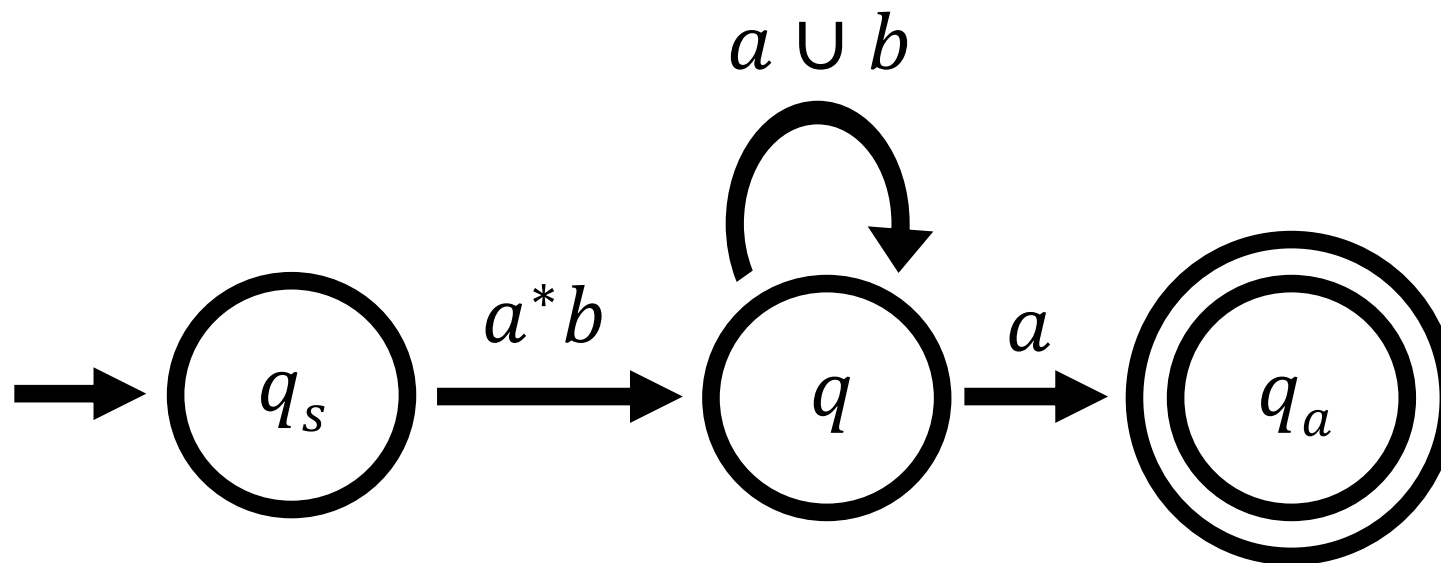**Theorem 2:** Every NFA has an equivalent regular expression

# NFA → Regular expression

Theorem 2: Every NFA has an equivalent regex

Proof idea: Simplify NFA by "ripping out" states one at a time and replacing with regexes
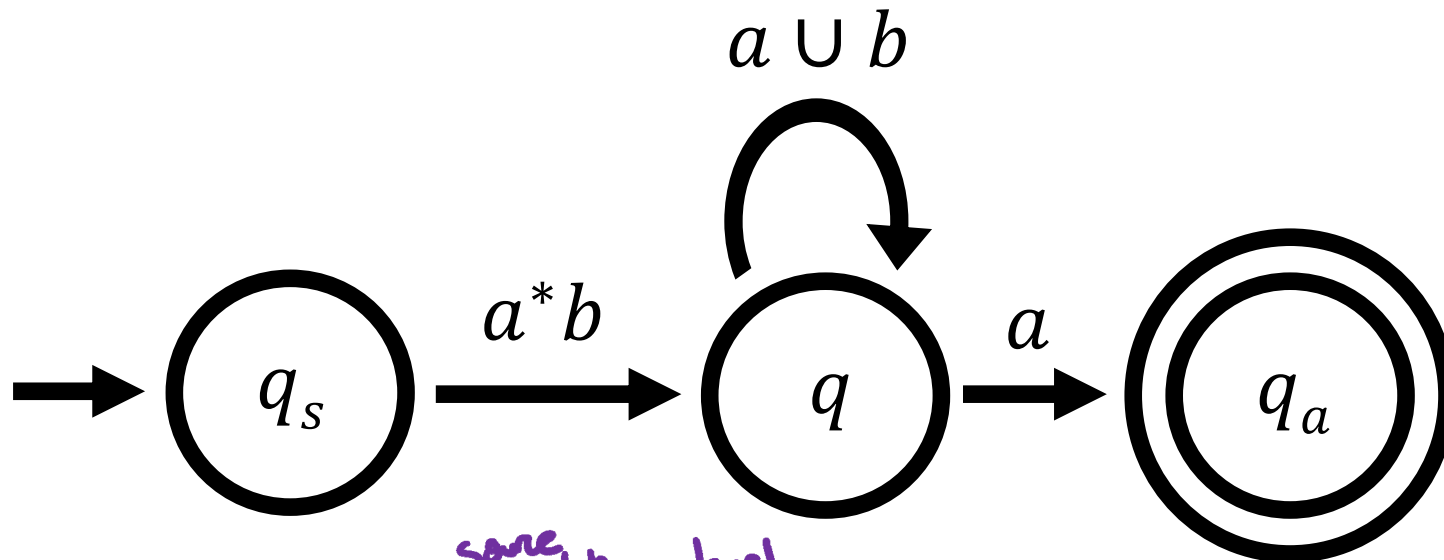
# Generalized NFAs

- **Every transition is labeled by a regex**
- One start state with only outgoing transitions
- Only one accept state with only incoming transitions
- Start state and accept state are distinct



$$a \cup b$$

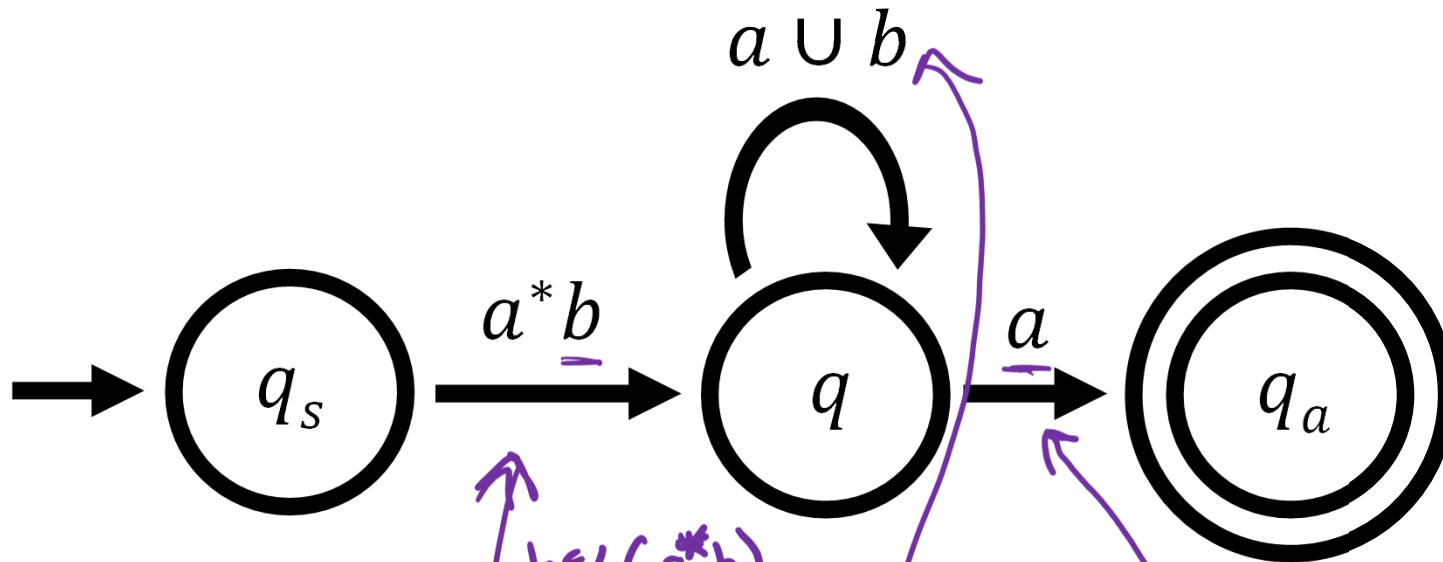$$q_s \xrightarrow{a^*b} q \xrightarrow{a} q_a$$
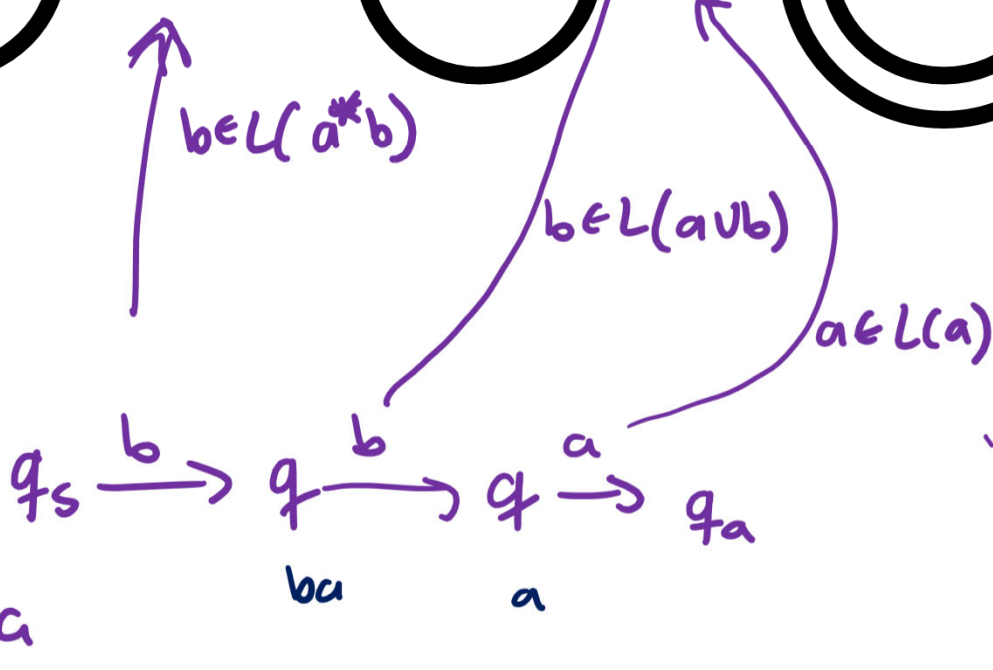
# Generalized NFA Example



$$R(q_s, q) = a^*b$$

$$R(q_a, q) = \phi$$

$$R(q, q_s) = \phi$$

# Which of these strings is accepted?

Which of the following strings is accepted by this GNFA?



$$a \cup b$$

$$a^* b \qquad q \qquad a \qquad q_a$$

$q_s$

$b \in L(a^* b)$

$b \in L(a \cup b)$

$a \in L(a)$

a) $aaa$

b) $aabb$

c) $bbb$

d) $bba$

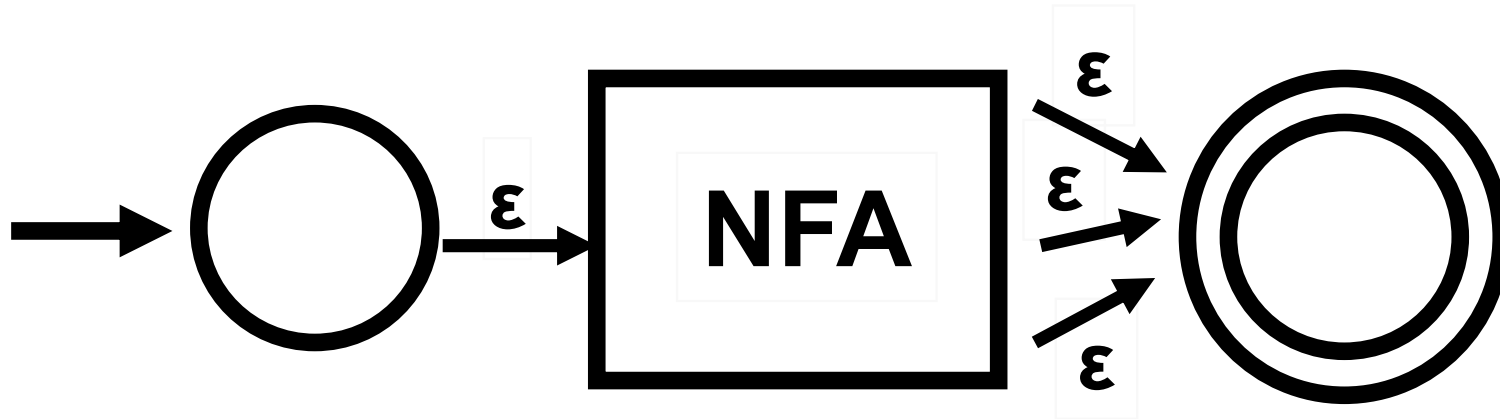$q_s \xrightarrow{b} q \xrightarrow{b} q \xrightarrow{a} q_a$
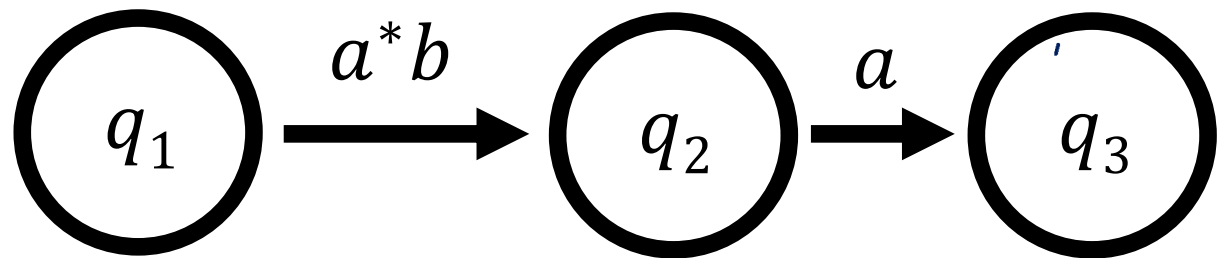
$ba$     $a$

$abba$

# NFA → Regular expression

# NFA → GNFA



- Add a new start state with no incoming arrows.
- Make a unique accept state with no outgoing arrows.

CS332 - Theory of Computation

# GNFA → Regular expression

**Idea:** While the machine has more than 2 states, rip one out and relabel the arrows with regexes to account for the missing state

# GNFA → Regular expression

Idea: While the machine has more than 2 states, rip one out and relabel the arrows with regexes to account for the missing state

a) $a^*b(a \cup b)a$
b) $a^*b(a \cup b)^*a$
c) $a^*b \cup (a \cup b) \cup a$
d) None of the above

$a \cup b$

$q_1 \xrightarrow{a^*b} q_2 \xrightarrow{a} q_3$

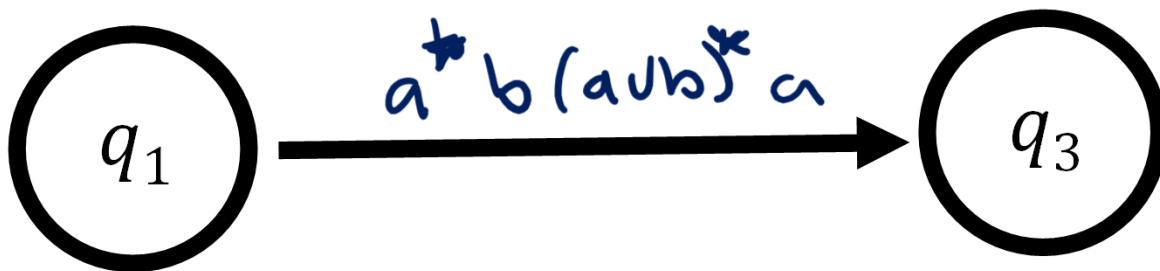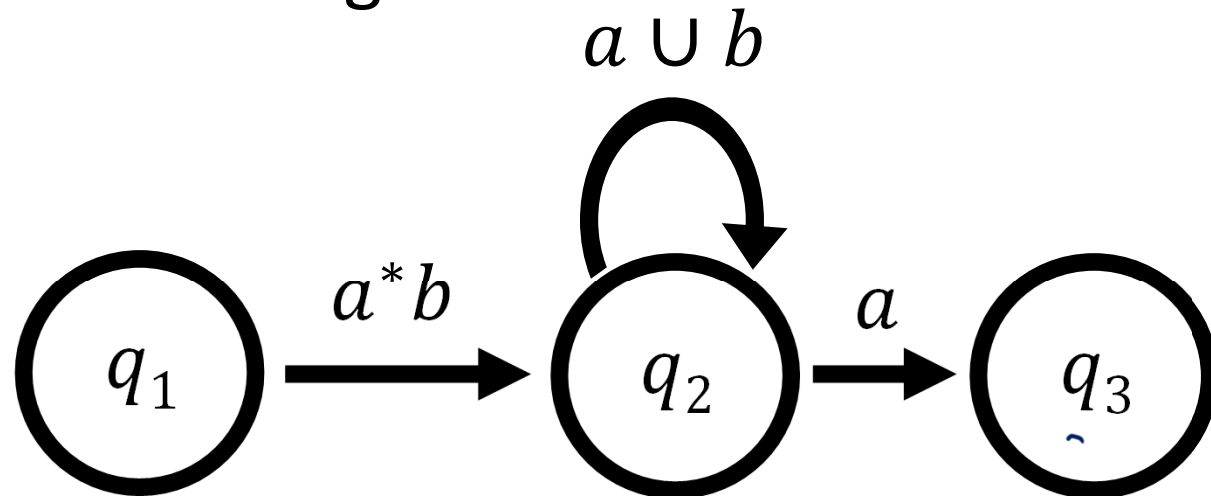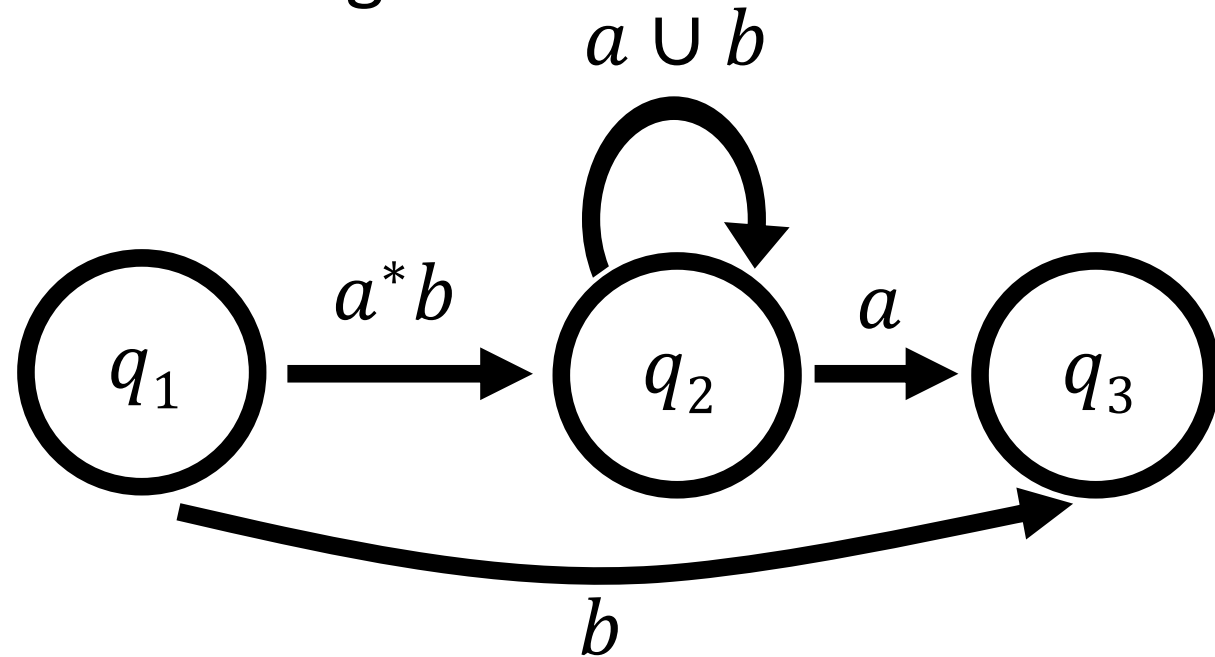$q_1 \xrightarrow{a^*b(a \cup b)^*a} q_3$

# GNFA → Regular expression

**Idea:** While the machine has more than 2 states, rip one out and relabel the arrows with regexes to account for the missing state

$$a \cup b$$

$q_1$ $\xrightarrow{a^*b}$ $q_2$ $\xrightarrow{a}$ $q_3$

$b$

$q_1$ $\xrightarrow{a^*b(a \cup b)^*a \cup b}$ $q_3$

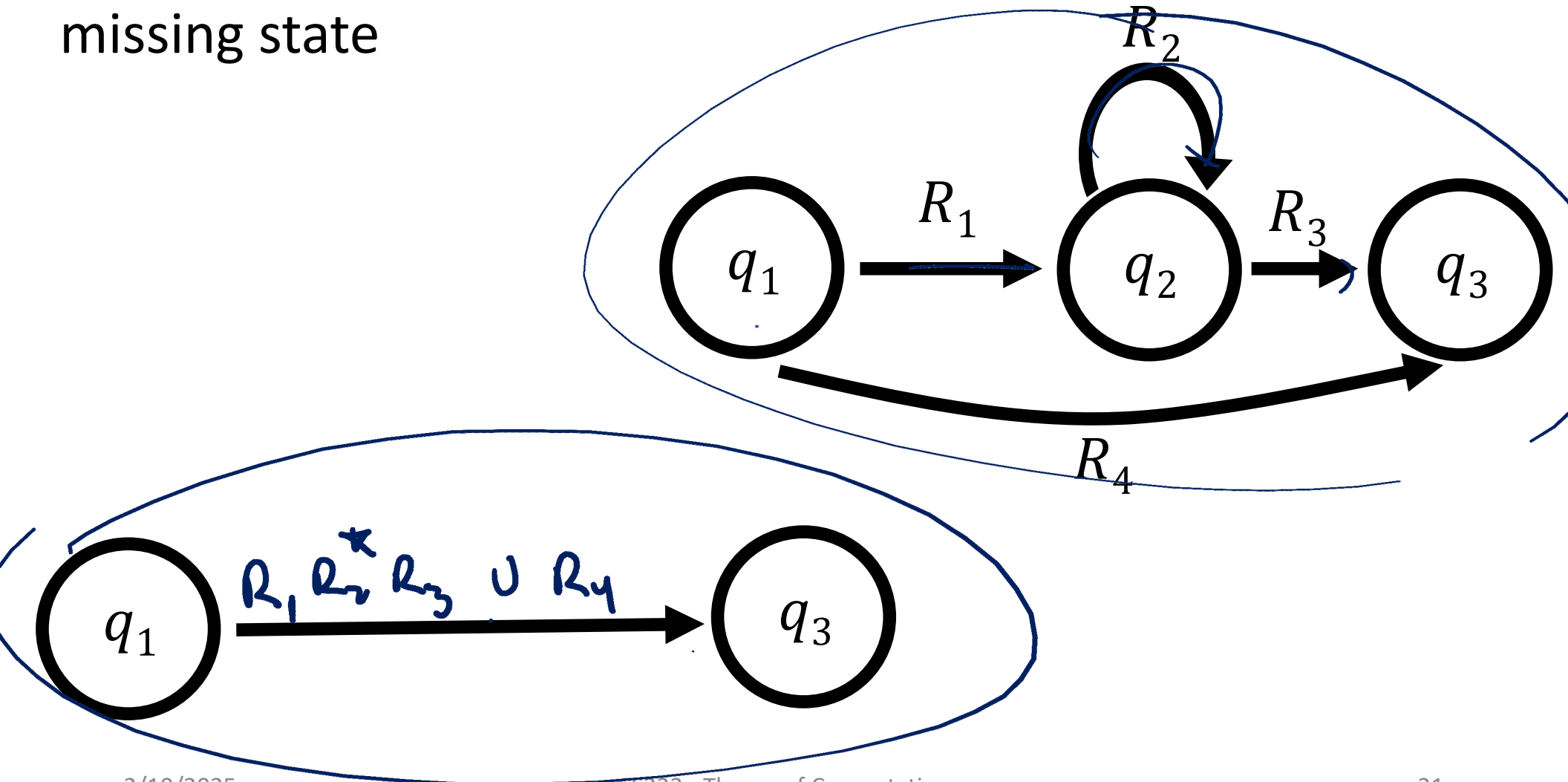# GNFA -> Regular expression

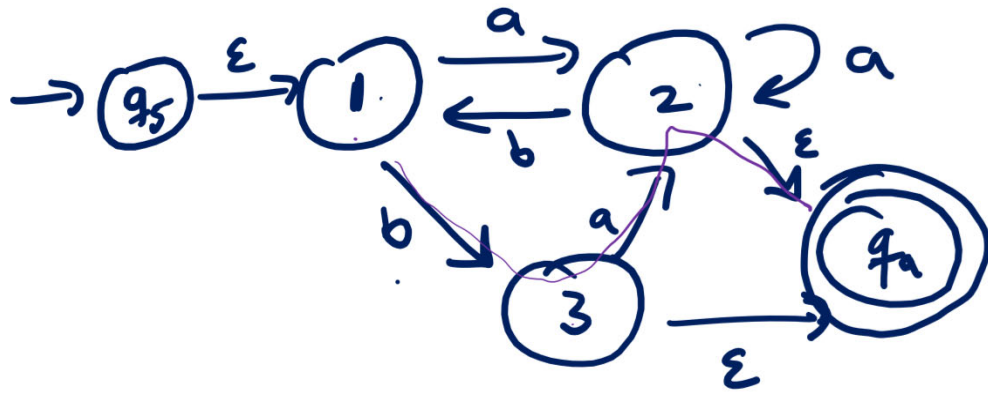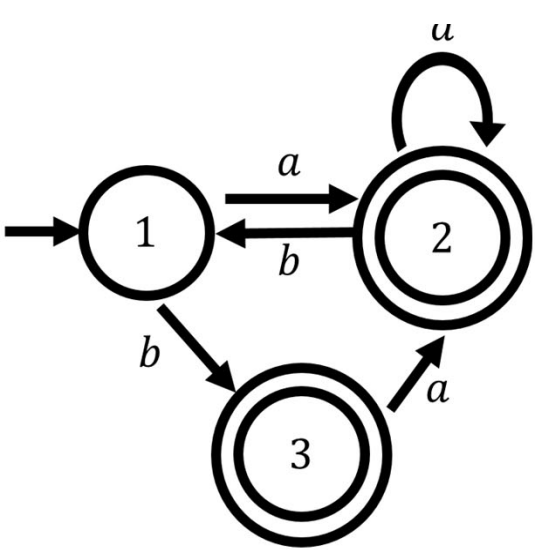Idea: While the machine has more than 2 states, rip one out and relabel the arrows with regexes to account for the missing state
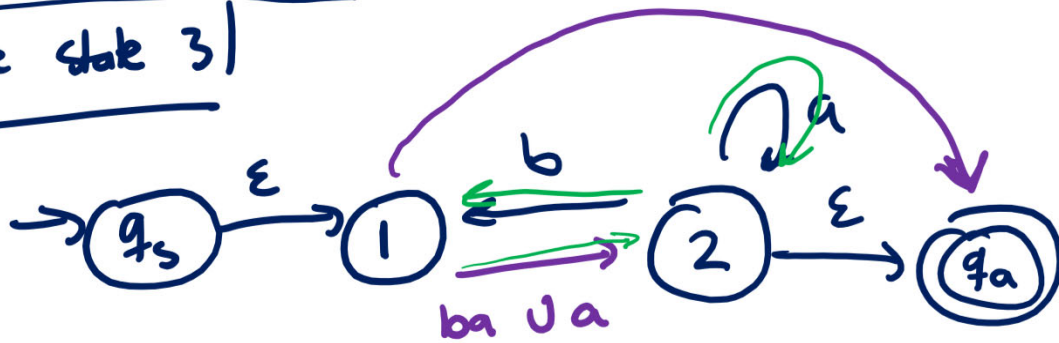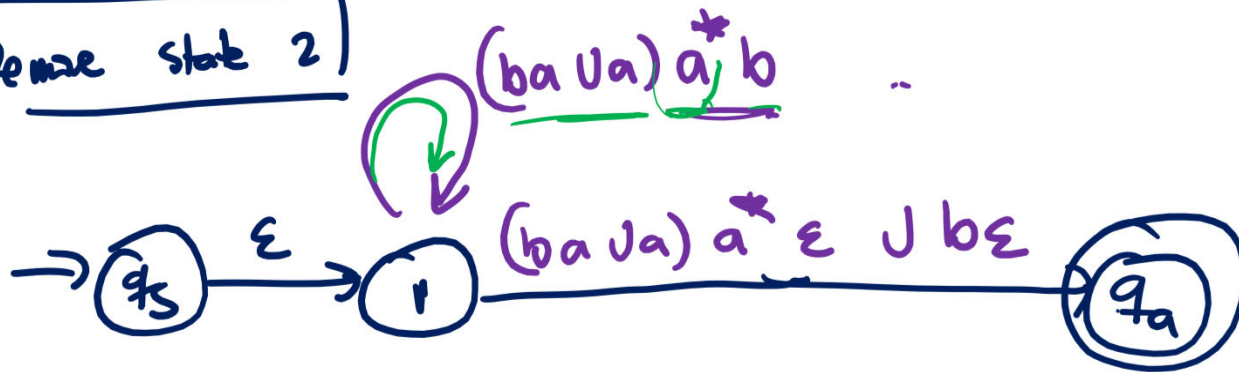
$R_2$

$q_1$  $R_1$  $q_2$  $R_3$  $q_3$

$R_4$

$q_1$  $R_1 R_2^* R_3 \cup R_4$  $q_3$

Top-left automaton: states 1, 2, 3 with start arrow into 1. Transitions: 1 →a→ 2, 2 →b→ 1, 1 →b→ 3, 3 →a→ 2, self-loop a on 2. States 2 and 3 are accepting.

Top-right: GNFA construction

$q_s \xrightarrow{\varepsilon} 1 \xrightarrow{a} 2$, $2 \xleftarrow{b} 1$, self-loop $a$ on 2, $2 \xrightarrow{\varepsilon} q_a$, $1 \xrightarrow{b} 3$, $3 \xrightarrow{a} 2$, $3 \xrightarrow{\varepsilon} q_a$

**Remove state 3**

$b\varepsilon \quad (= b)$

$q_s \xrightarrow{\varepsilon} 1 \xrightarrow{ba \cup a} 2 \xrightarrow{\varepsilon} q_a$, $2 \xleftarrow{b} 1$, self-loop $a$ on 2, $1 \xrightarrow{b} q_a$ (arc)

**Remove State 3**

$q_s \to q_a$

$$\varepsilon\big((ba \cup a)a^* b\big)^* \circ$$
$$\big((ba \cup a)a^* \varepsilon \cup b\varepsilon\big)$$

$$\equiv \big((ba \cup a)a^* b\big)^* \circ$$
$$\big((ba \cup a)a^* \cup b\big)$$

**Remove State 2**

$q_s \xrightarrow{\varepsilon} 1 \to q_a$

self-loop on 1: $(ba \cup a)a^* b$

$1 \xrightarrow{(ba \cup a)a^* \varepsilon \cup b\varepsilon} q_a$

# Limitations of Finite Automata

# Motivating Questions

- We've seen techniques for showing that languages are regular

  - Construct a DFA
  - Construct an NFA
  - Construct a regex
  - Use closure properties

- How can we tell if we've found the smallest DFA recognizing a language?

- Are all languages regular? How can we prove that a language is not regular?