




UMB CS622

# Regular Expressions

Wednesday September 22, 2021

Expressions		
Small Expression	Regular Expression	Large Expression
		
\$4.23	<code>^\\$d+\.?\d{2}/</code>	\$6.23

# *Announcements*

- HW1 graded
  - Use gradescope for grade questions / disputes
- HW2 due Sun 9/26 11:59pm EST

# *HW1 Review:* Inductive Proofs

## Must state:

- Induction on what
  - Often, “length of input string”
  - But not always!
- Base Case
- Inductive Case
  - with inductive hypothesis

Every statement and logical step must have justification

Usually taken from:

- Other theorems
- Definitions
- Given assumptions

# HW1 Review: Problem 4

Q: Prove that if some DFA  $M = (Q, \Sigma, \delta, q_0, F)$  has a state  $q$  such that  $\delta(q, a) = q$ , for all  $a \in \Sigma$ , then  $\hat{\delta}(q, w) = q$  for all possible strings  $w \in \Sigma^*$ . Use induction on the length of  $w$ .

A: *Claim.* If a DFA has a state  $q$  such that  $\forall a \in \Sigma \delta(q, a) = q$ , then  $\forall w \in \Sigma^* \hat{\delta}(q, w) = q$ .

*Proof.* By induction on  $w$ .

*Basis:* Trivially,  $\hat{\delta}(q, \epsilon) = q$  by the definition of  $\hat{\delta}$ .

*Induction step:* Let  $w = w'x$  where  $x \in \Sigma$ , assume the inductive hypothesis  $\hat{\delta}(q, w') = q$ . The objective is to show  $\hat{\delta}(q, w) = q$  using the claim's precondition  $\forall a \delta(q, a) = q$ .

$$\begin{aligned} \hat{\delta}(q, w) &= \hat{\delta}(q, w'x) && \text{by substitution of } w = w'x \\ &= \delta(\hat{\delta}(q, w'), x) && \text{by the definition of } \hat{\delta} \\ &= \delta(q, x) && \text{by the inductive hypothesis} \\ &= q && \text{by the precondition} \end{aligned}$$

□

Clearly stated base case and inductive step, with IH

Every logical step has justification

# HW1 Review: Problem 3 (part 2)

Prove that the following language is regular:

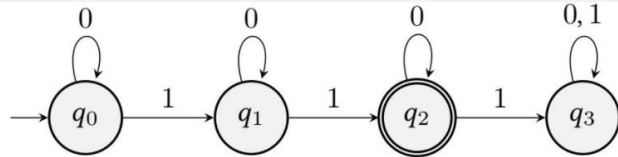
$$\{w \mid w \text{ has exactly two 1s}\}$$

**Q:**

In other words:

1. Design a DFA that recognizes the language; and
2. give an inductive proof that the DFA does indeed recognize the language.

Assume the language contains strings from alphabet  $\Sigma = \{0, 1\}$



**A:**

*Claim.*  $\forall w \in \Sigma^* P(w)$ , where  $P(w) = w \in L(M) \leftrightarrow w \in A$ .

*Proof.* By induction on  $w$ .

*Basis:*  $P(\epsilon)$  holds true as  $\epsilon \notin L(M)$  (the start state  $q_0$  is not accepting) and  $\epsilon \notin A$  ( $\epsilon$  does not have two 1s).

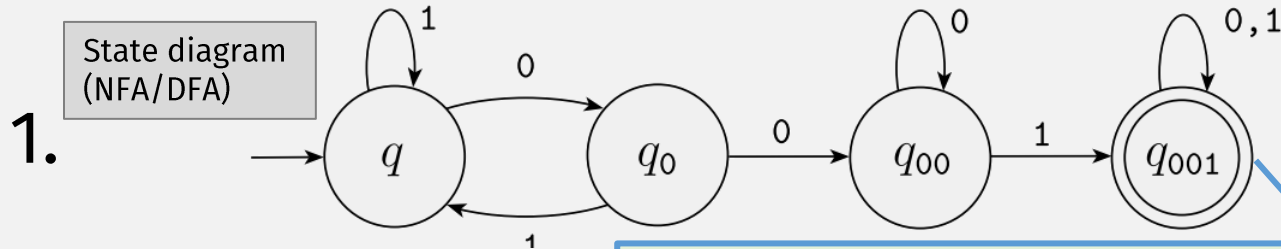
*Induction step:* Let  $w = w'a$  where  $a \in \Sigma$ . Assume  $P(w')$ , and consider  $P(w)$  throughout the following case analysis.

- If  $w'$  has zero 1s, then  $M$  is in state  $q_0$ .
  - Let  $a = 0$ :  $M$  stays in  $q_0$  and rejects with zero 1s.
  - Let  $a = 1$ :  $M$  enters  $q_1$  and rejects with one 1.
- If  $w'$  has one 1, then  $M$  is in state  $q_1$ .
  - Let  $a = 0$ :  $M$  stays in  $q_1$  and rejects with one 1.
  - Let  $a = 1$ :  $M$  enters  $q_2$  and accepts with two 1s.
- If  $w'$  has two 1s, then  $M$  is in state  $q_2$ .
  - Let  $a = 0$ :  $M$  stays in  $q_2$  and accepts with two 1s.
  - Let  $a = 1$ :  $M$  enters  $q_3$  and rejects with three 1s.

Not strong enough!  
(needs to say what each state represents)

These need justification  
(should come from IH)

# So Far: Regular Language Representations



A practical application:  
**text search** ... it doesn't fit!

2. Formal description
- $Q = \{q_1, q_2, q_3\}$ ,
  - $\Sigma = \{0,1\}$ ,
  - $\delta$  is described as

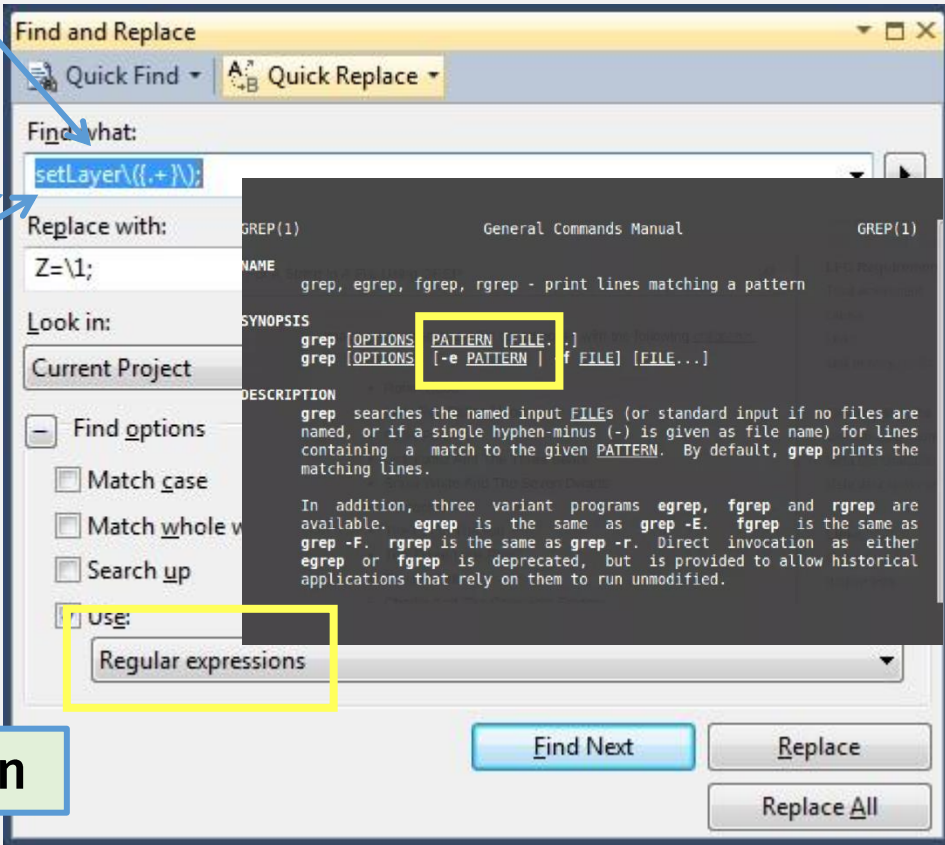
These define a computer (program) that finds strings containing 001

	0	1
$q_1$	$q_1$	$q_2$
$q_2$	$q_3$	$q_2$
$q_3$	$q_2$	$q_2$ ,

- $q_1$  is the start state, and
- $F = \{q_2\}$ .

3.  $\Sigma^* 001 \Sigma^*$

Need a more concise notation



# Regular Expressions Are Widely Used

- Perl
- Python
- Java
- Every lang!

**NAME**  
perlre - Perl regular expressions

Python » English » 3.8.6rc1 » Documentation » The Python Standard Library » Text Processing Services »

**re** — Regular expression operations

Source code: [Lib/re.py](#)

This module provides regular expression matching operations similar to those found in Perl.

java.util.regex

**Class Pattern**

java.lang.Object  
java.util.regex.Pattern

# Regular Expressions: Formal Definition

$R$  is a *regular expression* if  $R$  is

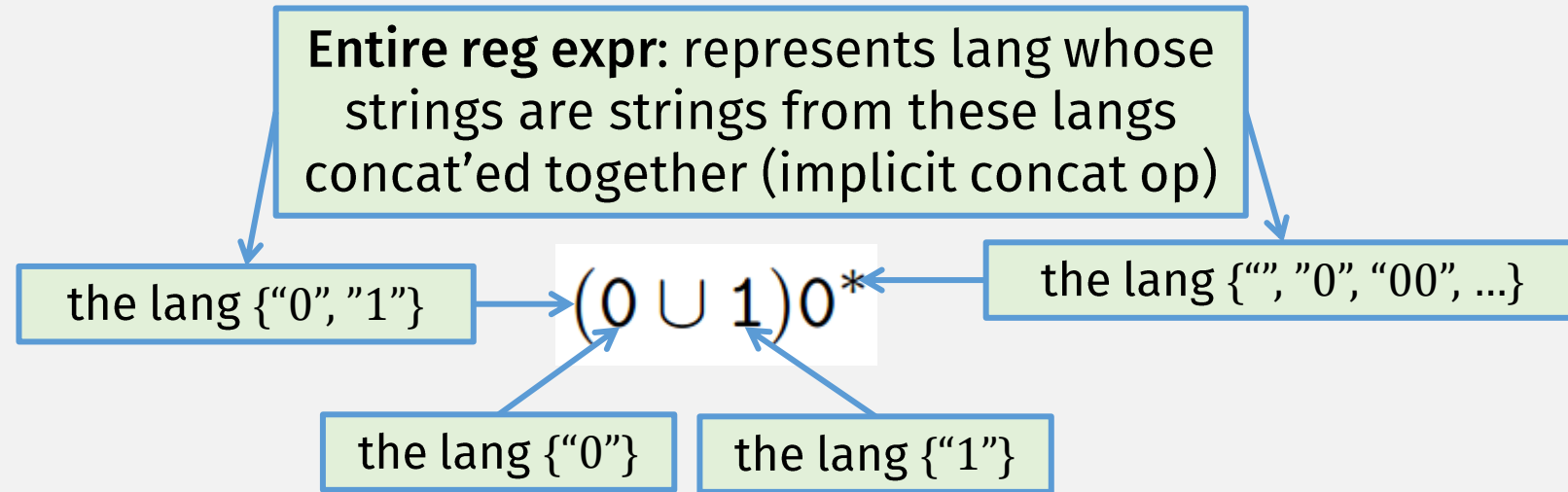
1.  $a$  for some  $a$  in the alphabet  $\Sigma$ , (A lang containing  $a$ ) length-1 string
2.  $\epsilon$ , (A lang containing) the empty string
3.  $\emptyset$ , The empty set (i.e., a lang containing no strings)
4.  $(R_1 \cup R_2)$ , where  $R_1$  and  $R_2$  are regular expressions, union
5.  $(R_1 \circ R_2)$ , where  $R_1$  and  $R_2$  are regular expressions, or concat
6.  $(R_1^*)$ , where  $R_1$  is a regular expression. star

Base cases plus union, concat, and Kleene star can express any regular language!

(But we have to prove it)



# Regular Expression: Concrete Example



- **Operator Precedence:**
  - Parens
  - Star
  - Concat (sometimes implicit)
  - Union

Thm: A lang is regular iff some reg expr describes it

⇒ If a language is regular, it is described by a reg expression

⇐ If a language is described by a reg expression, it is regular

- Easy!
- For a given regexp, construct the equiv NFA!
- (we mostly did it already when discussing closed ops)

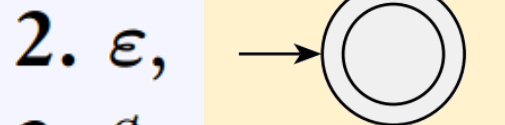
How to show that a lang is regular?

**Construct DFA or NFA!**

# RegExpr $\rightarrow$ NFA

$R$  is a *regular expression* if  $R$  is

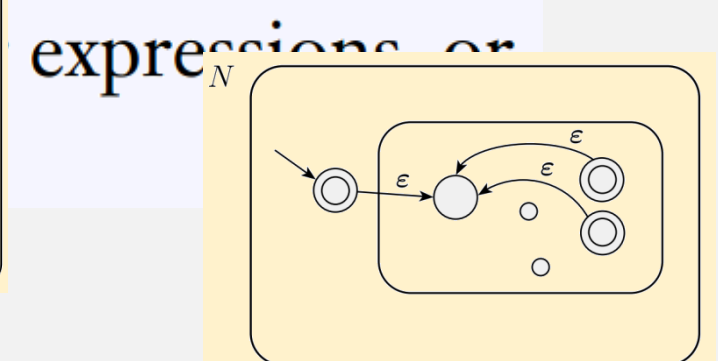
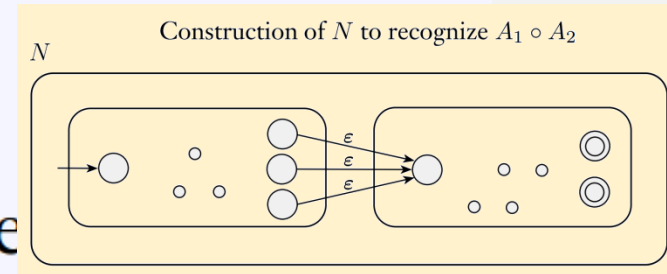
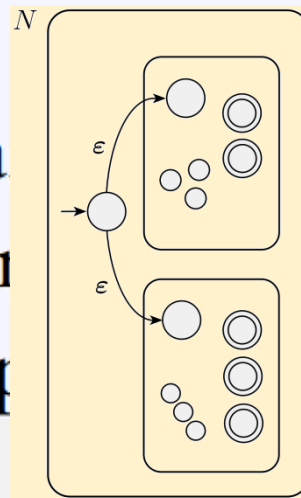
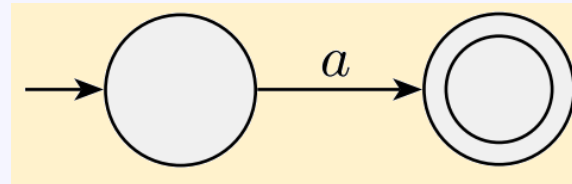
1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,



4.  $(R_1 \cup R_2)$ , where  $R_1$  and  $R_2$  are regular expressions,

5.  $(R_1 \circ R_2)$ , where  $R_1$  and  $R_2$  are regular expressions,

6.  $(R_1^*)$ , where  $R_1$  is a regular expression.



Thm: A lang is regular iff some reg expr describes it

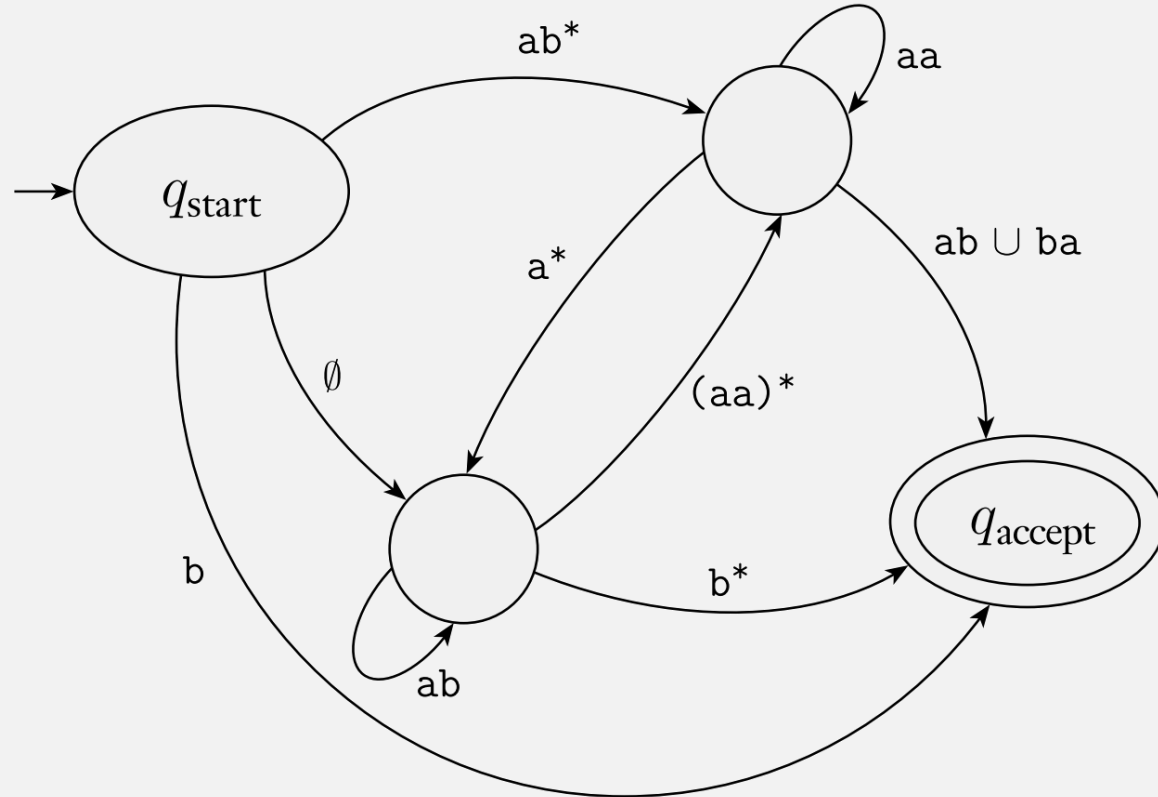
⇒ If a language is regular, it is described by a reg expression

- Harder!
- Need to convert DFA or NFA to Regular Expression
- To do so, need new kind of machine: a GNFA

⇐ If a language is described by a reg expression, it is regular

- Easy!
- Construct the NFA! (**Done**)

# Generalized NFAs (GNFAs)



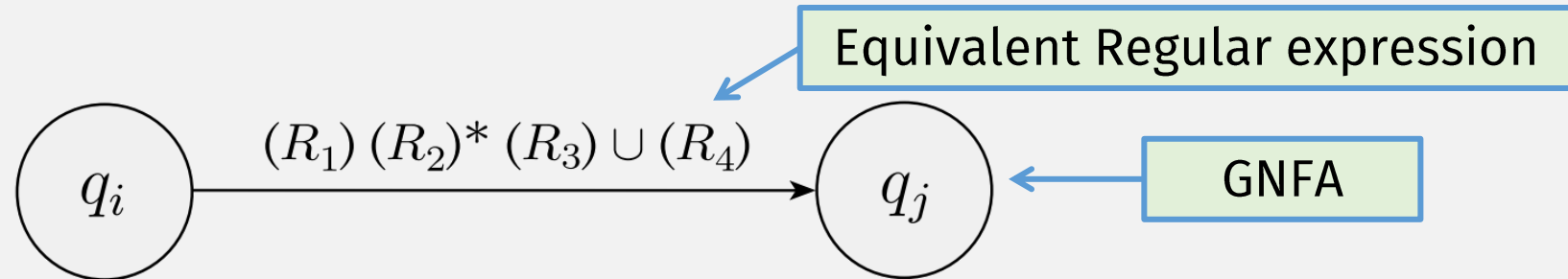
- GNFA = NFA with regular expression transitions

**Want to convert  
GNFAs to Reg Exprs**

# GNFA $\rightarrow$ RegExpr function

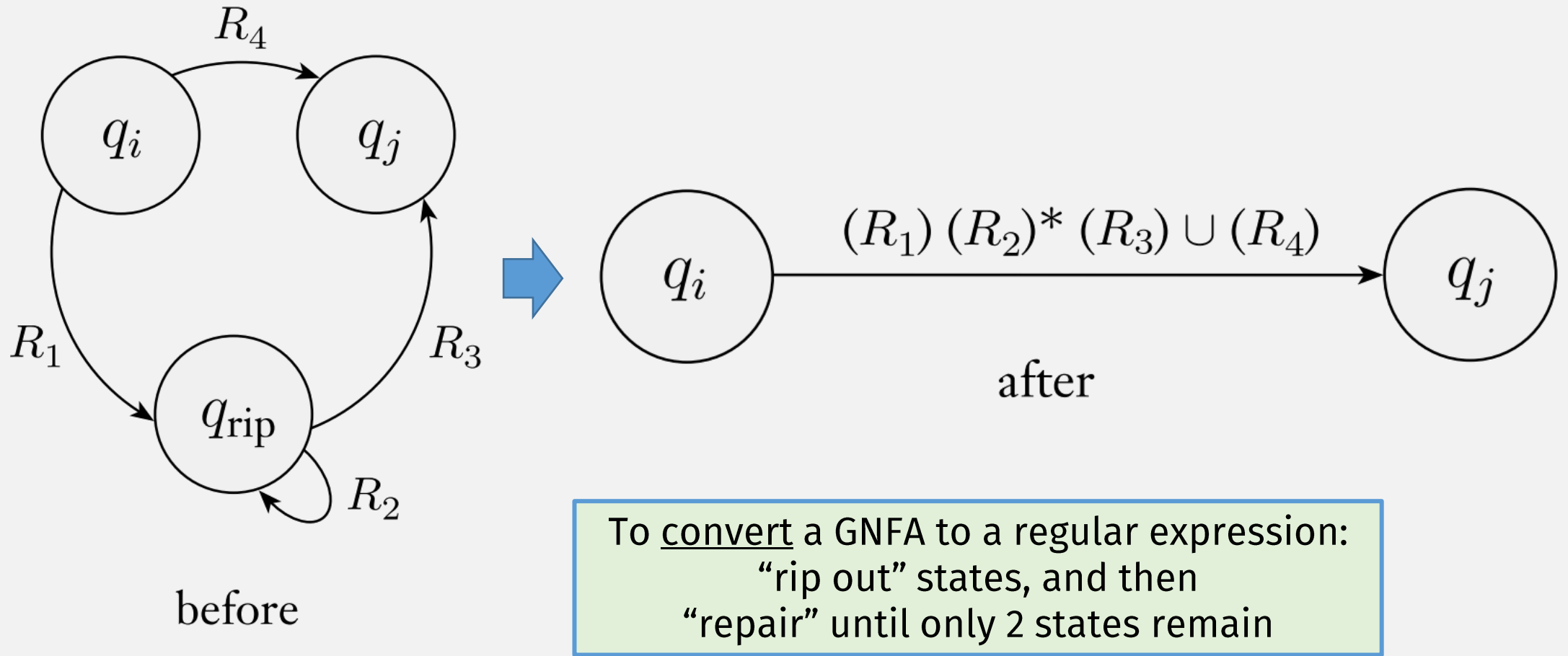
On GNFA input  $G$ :

- If  $G$  has 2 states, return the regular expression transition, e.g.:

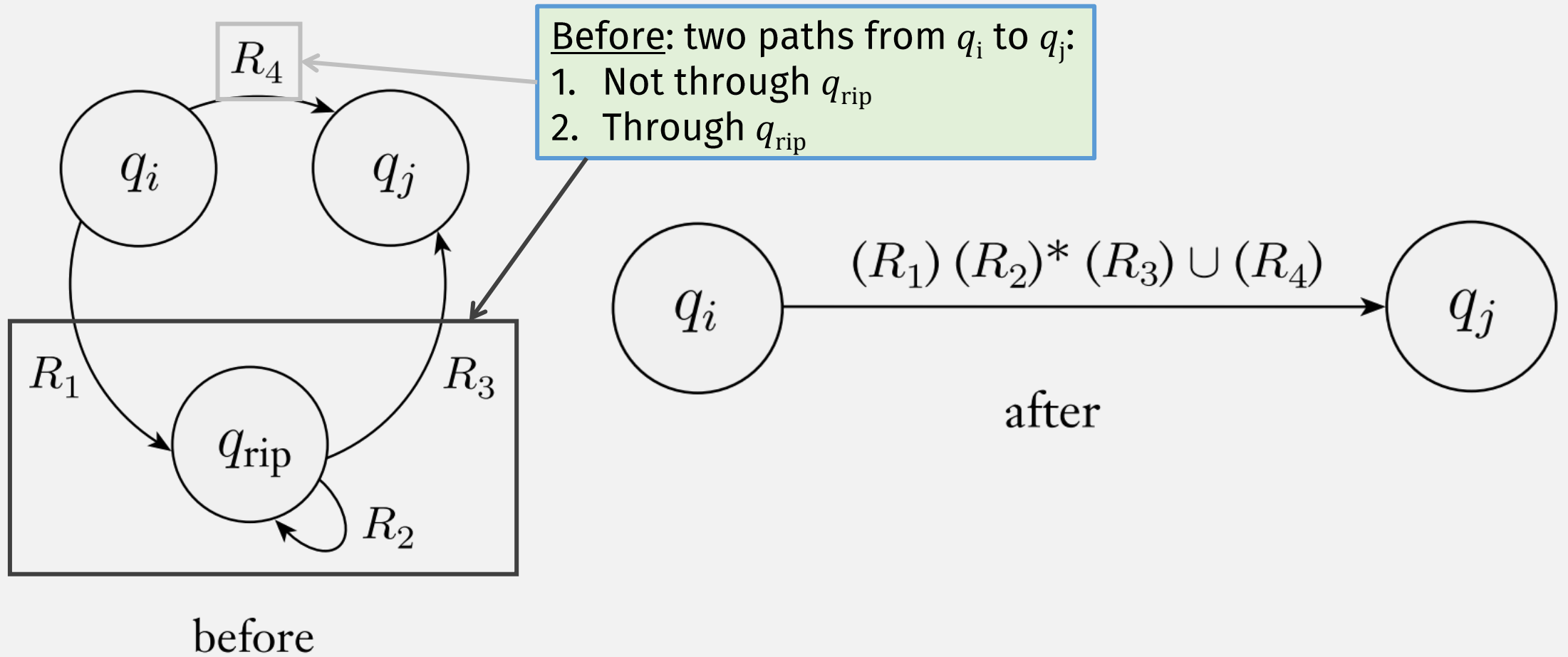


- Else:
  - “Rip out” one state
  - “Repair” the machine to get an equivalent GNFA  $G'$
  - Recursively call  $\text{GNFA} \rightarrow \text{RegExpr}(G')$

# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step

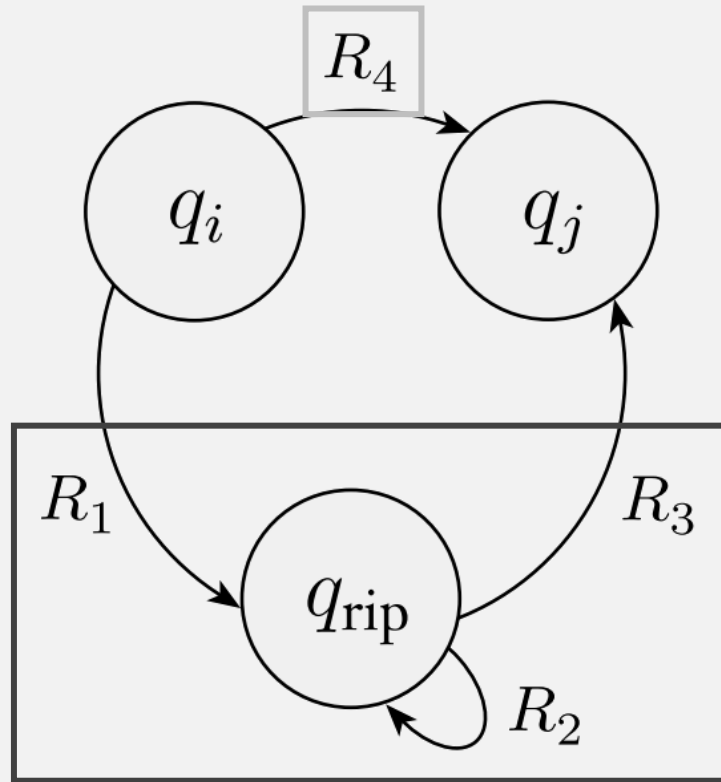


# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step





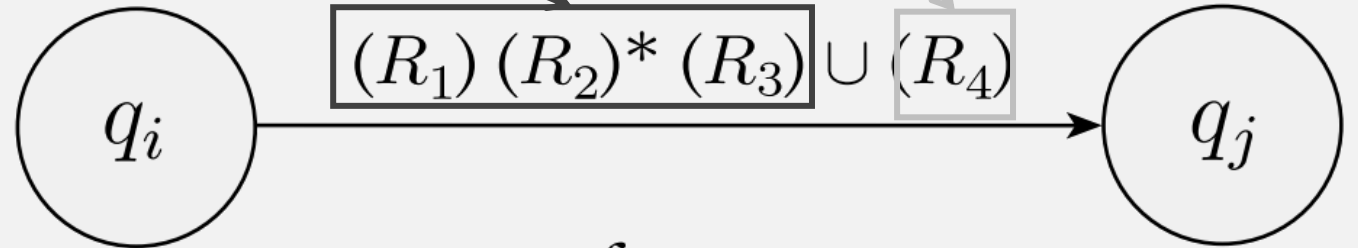
# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step



before

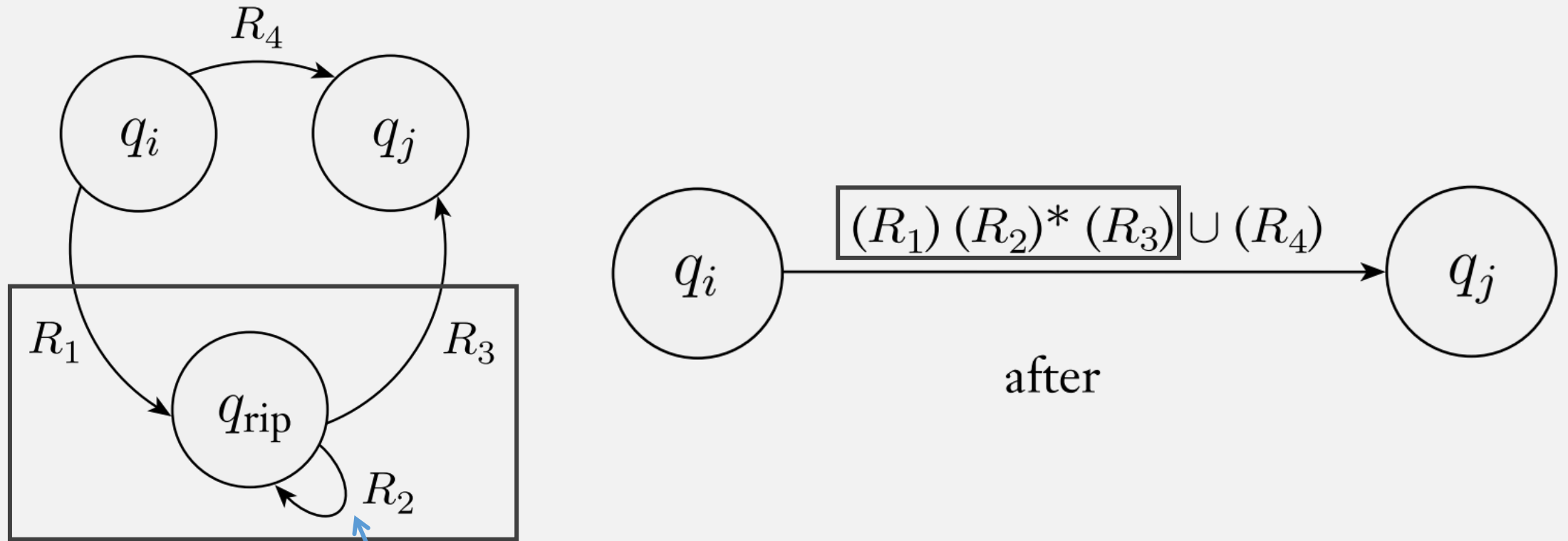
After: still two “paths” from  $q_i$  to  $q_j$

1. Not through  $q_{rip}$
2. Through  $q_{rip}$



after

# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step

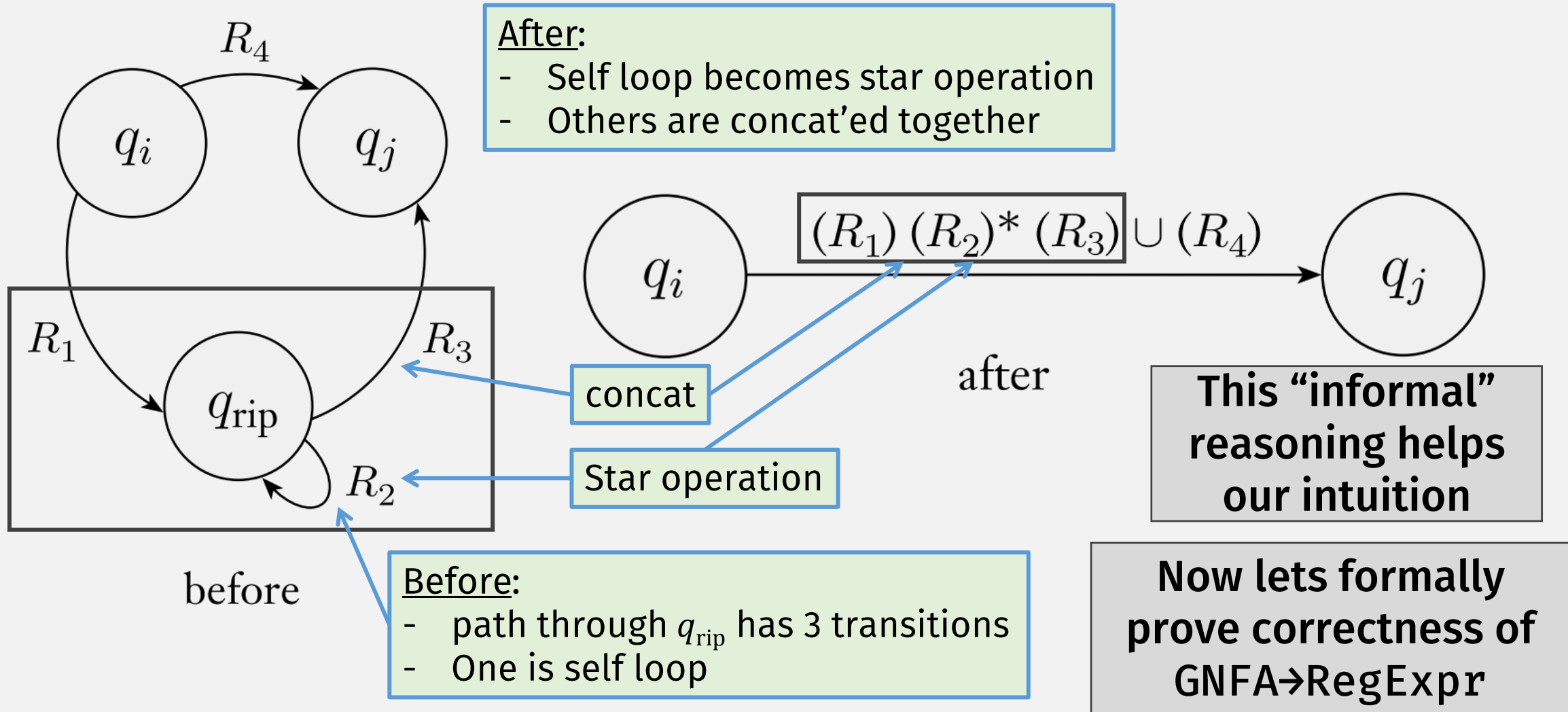


before

Before:

- path through  $q_{rip}$  has 3 transitions
- One is self loop

# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step



# GNFA $\rightarrow$ RegExpr “Correctness”

- Where “Correct” means:

Use Proof by induction ... on size of G

$$\text{LANGOF} ( G ) = \text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( G ) )$$

This is the property we want to prove

- i.e., GNFA $\rightarrow$ RegExpr must not change the language!

# Previously: Recursive (Inductive) Definitions

- Have (at least) two parts:
  - Base case
  - Inductive case
    - Self-reference must be “smaller”

- Example:

**Def:  $\text{GNFA} \rightarrow \text{RegExp}$ :** input  $G$  is a GNFA with  $n$  states:

Base case

If  $n = 2$ : return the regular expression on the transition

Inductive case

Else ( $G$  has  $n > 2$  states):

- “Rip” out one state and “Repair” to get  $G'$
- Recursively Call  $\text{GNFA} \rightarrow \text{RegExp}(G')$

“smaller” self-reference

This is exactly the structure of an inductive proof!

# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$

Def: GNFA $\rightarrow$ RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call **GNFA $\rightarrow$ RegExpr**( $G'$ )

➤ Proof (by induction on size of  $G$ ):

# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA} \rightarrow \text{RegExpr}(G))$

Def: GNFA  $\rightarrow$  RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

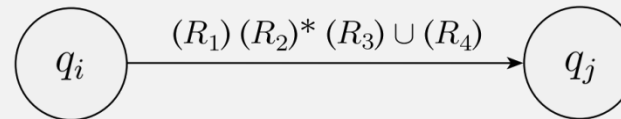
Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call  $\text{GNFA} \rightarrow \text{RegExpr}(G')$

• Proof (by induction on size of  $G$ ):

➤ Base case:  $G$  has 2 states



•  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA} \rightarrow \text{RegExpr}(G))$  is true, by def of GNFA!

# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$

Def: GNFA $\rightarrow$ RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call **GNFA $\rightarrow$ RegExpr**( $G'$ )

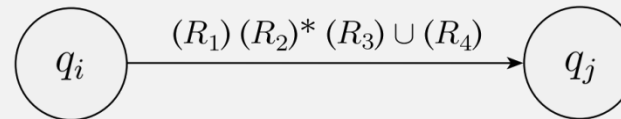
• Proof (by induction on size of  $G$ ):

• Base case:  $G$  has 2 states

•  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$  is true!

➤ IH: Assume  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$

• For some  $G'$  with  $n-1$  states





# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$

Def: GNFA $\rightarrow$ RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call GNFA $\rightarrow$ RegExpr( $G'$ )

• Proof (by induction on size of  $G$ ):

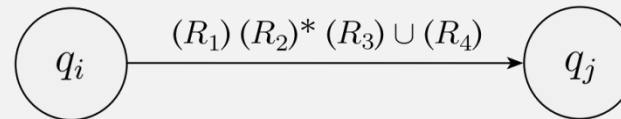
• Base case:  $G$  has 2 states

•  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$  is true!

• IH: Assume  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$

• For some  $G'$  with  $n-1$  states

➤ Induction Step: Prove it's true for  $G$  with  $n$  states



# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$

Def: GNFA $\rightarrow$ RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call GNFA $\rightarrow$ RegExpr( $G'$ )

• Proof (by induction on size of  $G$ ):

• Base case:  $G$  has 2 states

•  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$  is true!

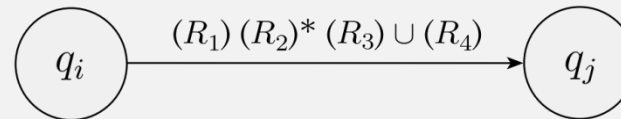
• IH: Assume  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$

• For some  $G'$  with  $n-1$  states

➤ Induction Step: Prove it's true for  $G$  with  $n$  states

• After “rip/repair” step, we have exactly a GNFA  $G'$  with  $n-1$  states

• And we know  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$  from the IH!



# GNFA $\rightarrow$ RegExpr is correct

Want to prove:  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$

Def: GNFA $\rightarrow$ RegExpr: input  $G$  is a GNFA with  $n$  states:

If  $n = 2$ : return the regular expression on the transition

Else ( $G$  has  $n > 2$  states):

“Rip” out one state and “Repair” to get  $G'$

Recursively Call GNFA $\rightarrow$ RegExpr( $G'$ )

• Proof (by induction on size of  $G$ ):

• Base case:  $G$  has 2 states

•  $\text{LANGOF}(G) = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G))$  is true!

• IH: Assume  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$

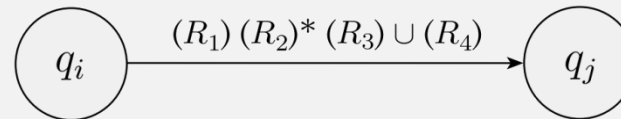
• For some  $G'$  with  $n-1$  states

▪ Induction Step: Prove it's true for  $G$  with  $n$  states

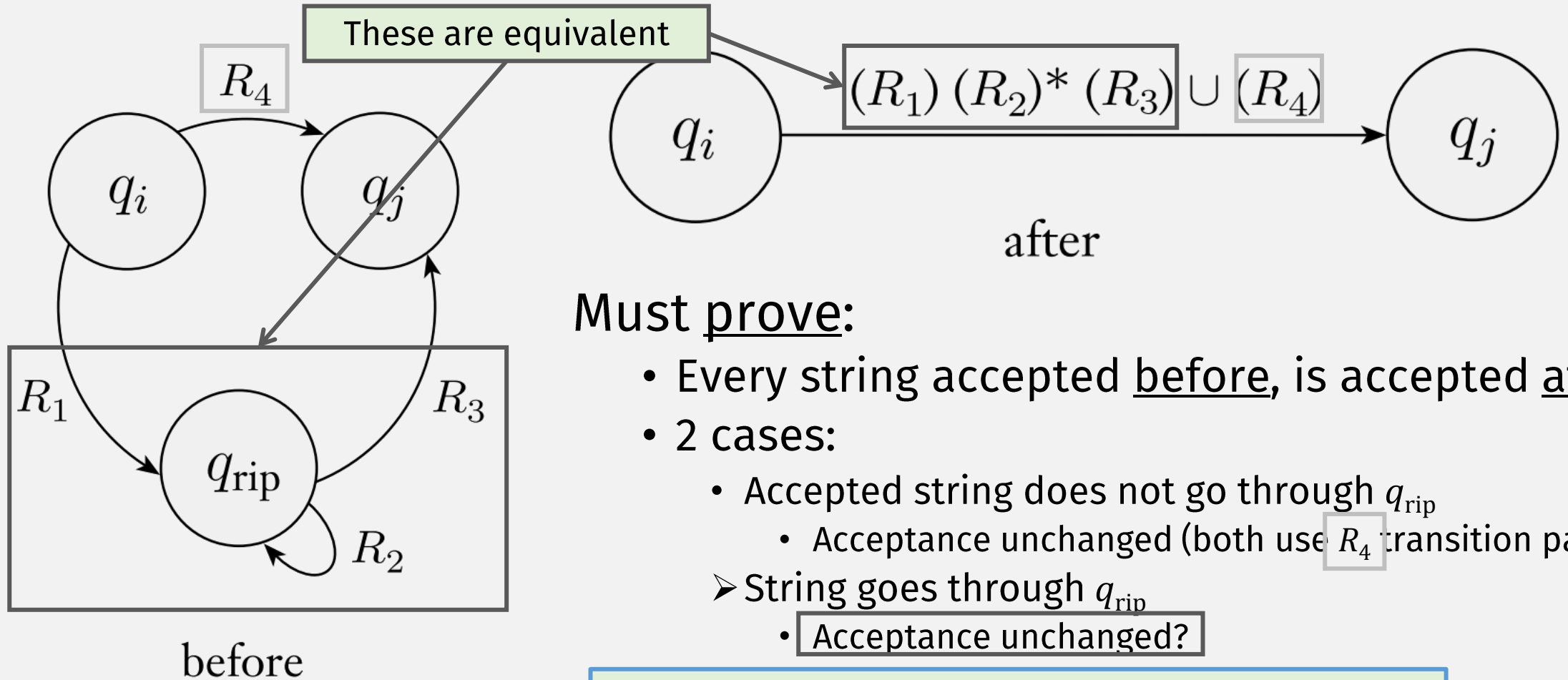
• After “rip/repair” step, we have exactly a GNFA  $G'$  with  $n-1$  states

• And we know  $\text{LANGOF}(G') = \text{LANGOF}(\text{GNFA}\rightarrow\text{RegExpr}(G'))$  from the IH!

➤ To go from  $G$  to  $G'$ : just need to prove correctness of “rip/repair” step



# GNFA $\rightarrow$ RegExpr: “rip/repair” correctness



**Mostly done this already!**  
**Just need to state more formally**

Thm: A lang is regular iff some reg expr describes it

⇒ If a language is regular, it is described by a reg expr

- Harder!
- Need to convert DFA or NFA to Regular Expression
- Use GNFA→RegExpr to convert GNFA to regular expression! (Done!)

⇐ If a language is described by a reg expr, it is regular

- Construct the NFA! (**Done**)

Now we may use regular expressions to represent regular langs.

I.e., we have another way to prove things about reg langs!

So a regular language has these equivalent representations:

- DFA
- NFA
- Regular Expression

# Thm: Reverse is Closed for Regular Langs

- For any string  $w = w_1w_2 \cdots w_n$ , the *reverse* of  $w$ , written  $w^{\mathcal{R}}$ , is the string  $w$  in reverse order,  $w_n \cdots w_2w_1$ .

For any language  $A$ , let  $A^{\mathcal{R}} = \{w^{\mathcal{R}} \mid w \in A\}$

- Theorem: if  $A$  is regular, so is  $A^{\mathcal{R}}$
- Proof (by induction on regular expressions):

Remember: A language is regular iff it has a regular expression representation

# Thm: Reverse is Closed for Regular Langs

if  $A$  is regular, so is  $A^{\mathcal{R}}$

Case Analysis, assume some regular language  $A$  is represented with the regular expression ...

Base cases

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ , same reg. expr. represents  $A^{\mathcal{R}}$  so it is regular

2.  $\epsilon$ , same reg. expr. represents  $A^{\mathcal{R}}$  so it is regular

3.  $\emptyset$ , same reg. expr. represents  $A^{\mathcal{R}}$  so it is regular

Inductive cases

4.  $(R_1 \cup R_2)$ , where  $R_1$  and  $R_2$  are regular expressions, ←

5.  $(R_1 \circ R_2)$ , where  $R_1$  and  $R_2$  are regular expressions, or

6.  $(R_1^*)$ , where  $R_1$  is a regular expression.

Other cases will use similar reasoning

Need to show: if  $A_1 \cup A_2$  is a regular language, then  $(A_1 \cup A_2)^{\mathcal{R}}$  is regular

IH: if  $A_1$  and  $A_2$  are the regular languages represented by  $R_1$  and  $R_2$ , then  $A_1^{\mathcal{R}}$  and  $A_2^{\mathcal{R}}$  are regular too

Proof:  $(A_1 \cup A_2)^{\mathcal{R}} = A_1^{\mathcal{R}} \cup A_2^{\mathcal{R}}$ , because reversal and union don't affect each other and are interchangeable

$A_1^{\mathcal{R}}$  and  $A_2^{\mathcal{R}}$  are regular (from IH) and union is closed for regular langs (class thm), so  $A_1^{\mathcal{R}} \cup A_2^{\mathcal{R}}$  is regular

# **In-Class quiz 9/22**

See gradescope