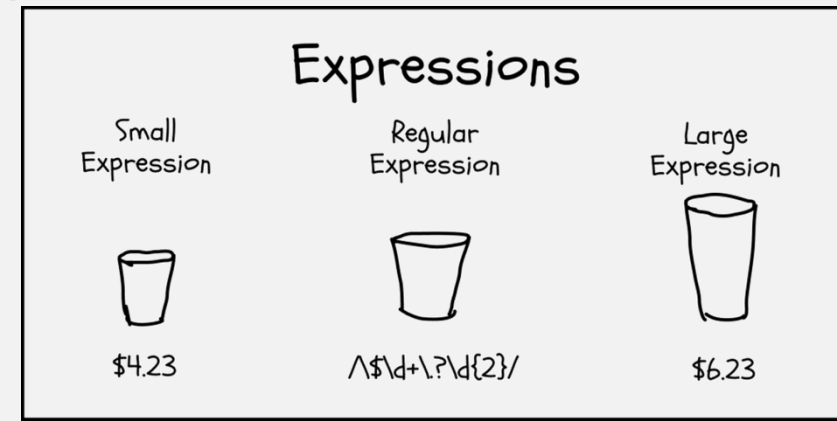


# CS 420 / CS 620

## Regular Expressions

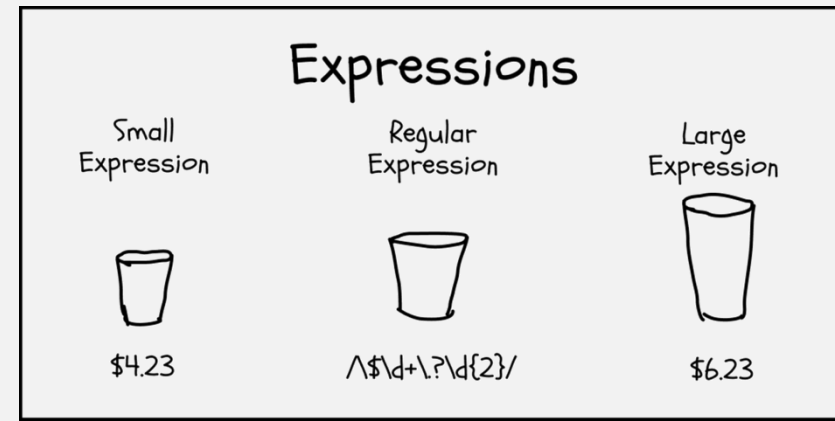
Monday, October 6, 2025

UMass Boston Computer Science



# Announcements

- HW 4
  - ~~Due: Mon 10/6 12pm (noon)~~
- HW 5
  - Out: Mon 10/6 12pm (noon)
  - Due (unofficial): Mon 10/13 12pm (noon) (stay on schedule!)
  - Due (up to): Wed 10/15 12pm (noon)
- HW 6 (most likely)
  - Out: Mon 10/13 12pm (noon)
  - Due: Wed 10/20 12pm (noon)
- **No class:** next Mon 10/13 (Indigenous Peoples)



# In-class question (in gradescope) preview

When used as a string, the epsilon symbol ( $\epsilon$ ) is equivalent to which of the following?

When used as the empty transition, the epsilon symbol ( $\epsilon$ ) is equivalent to which of the following?

When used as a regular expression, the epsilon symbol ( $\epsilon$ ) is equivalent to which of the following?

When used as an input to an NFA's single-step  $\delta$  function, the epsilon symbol ( $\epsilon$ ) is which of the following?

When used as an input to an NFA's multi-step  $\hat{\delta}$  function, the epsilon symbol ( $\epsilon$ ) is which of the following?

When used as a transition label in a GNFA, the epsilon symbol ( $\epsilon$ ) is which of the following?

# List of Closed Ops for Reg Langs (so far)

• Union

• Concatentation

• Kleene Star (repetition) ?

**Star:**  $A^* = \{x_1x_2 \dots x_k \mid k \geq 0 \text{ and each } x_i \in A\}$

# Kleene Star Example

Let the alphabet  $\Sigma$  be the standard 26 letters  $\{a, b, \dots, z\}$ .

If  $A = \{\text{good, bad}\}$

“repeat” zero

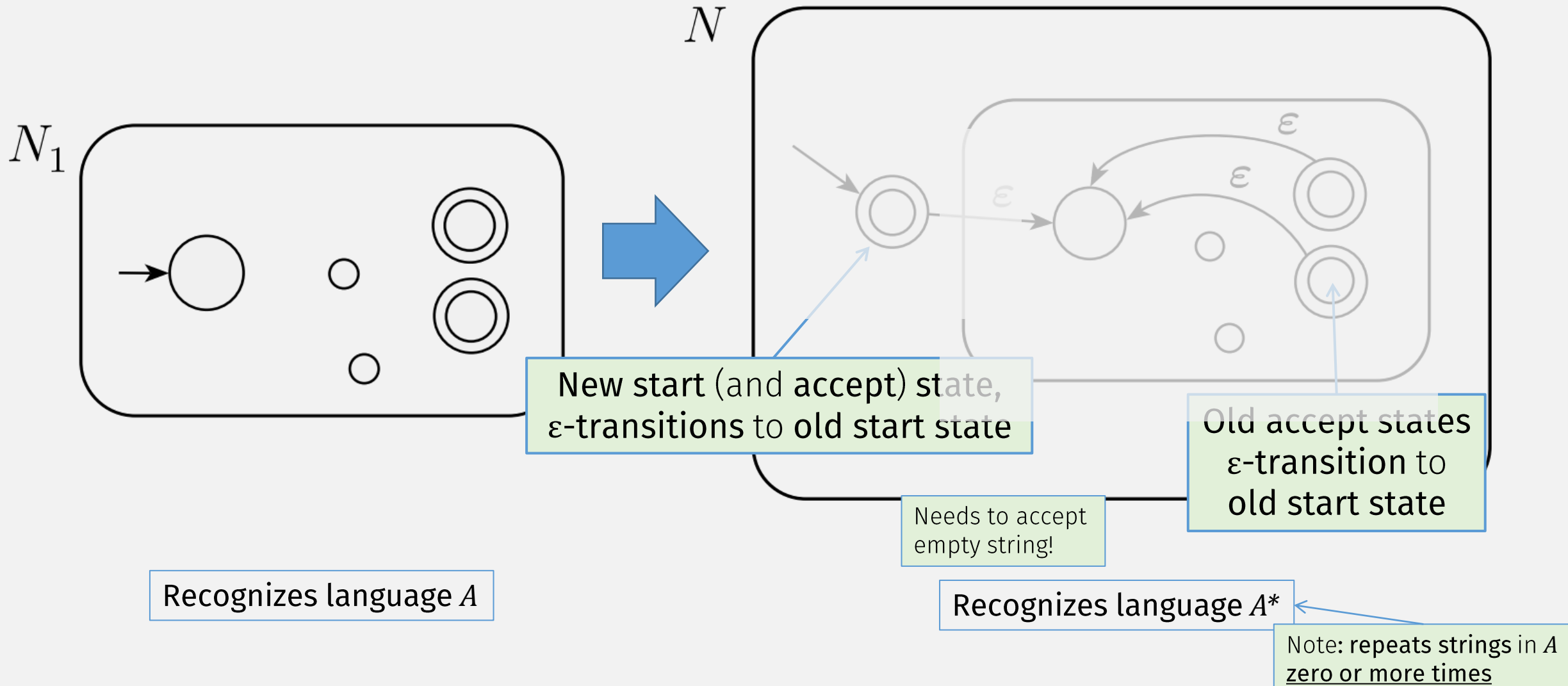
$A^* = \{\epsilon, \text{good, bad, goodgood, goodbad, badgood, badbad, goodgoodgood, goodgoodbad, goodbadgood, goodbadbad, \dots}\}$

Note: repeat strings in  $A$   
zero or more times

(this is an infinite language!)

**Star:**  $A^* = \{x_1x_2 \dots x_k \mid k \geq 0 \text{ and each } x_i \in A\}$

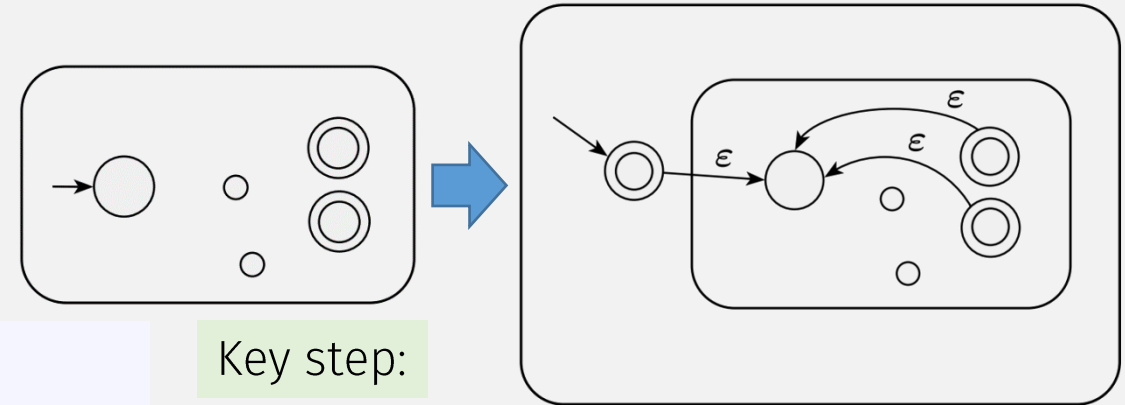
# Kleene Star is Closed for Regular Langs?



# Kleene Star is Closed for Regular Languages

## THEOREM

The class of regular languages is closed under the star operation.



Key step:

$\text{STAR}_{\text{NFA}} : \text{NFA} \rightarrow \text{NFA}$

Where:

$$N = \text{STAR}_{\text{NFA}}(N_1)$$

$$L(N) = L(N_1)^*$$

# Kleene Star is Closed for Regular Languages

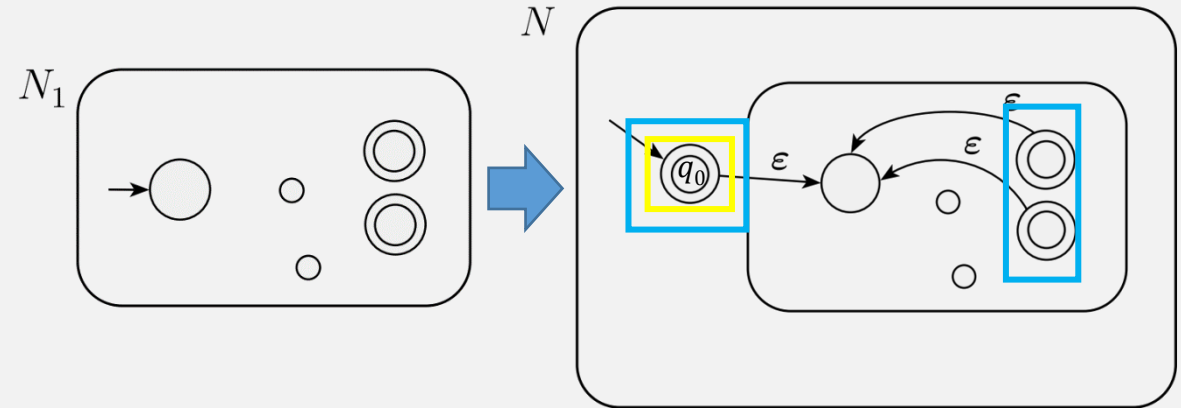
(part of)

**PROOF** Let  $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$  recognize  $A_1$ .

$N = \text{STAR}_{\text{NFA}}(N_1) = (Q, \Sigma, \delta, q_0, F)$  to recognize  $A_1^*$ .

1.  $Q = \{q_0\} \cup Q_1$
2. The state  $q_0$  is the new start state.
3.  $F = \{q_0\} \cup F_1$

Kleene star of a language must accept the empty string!





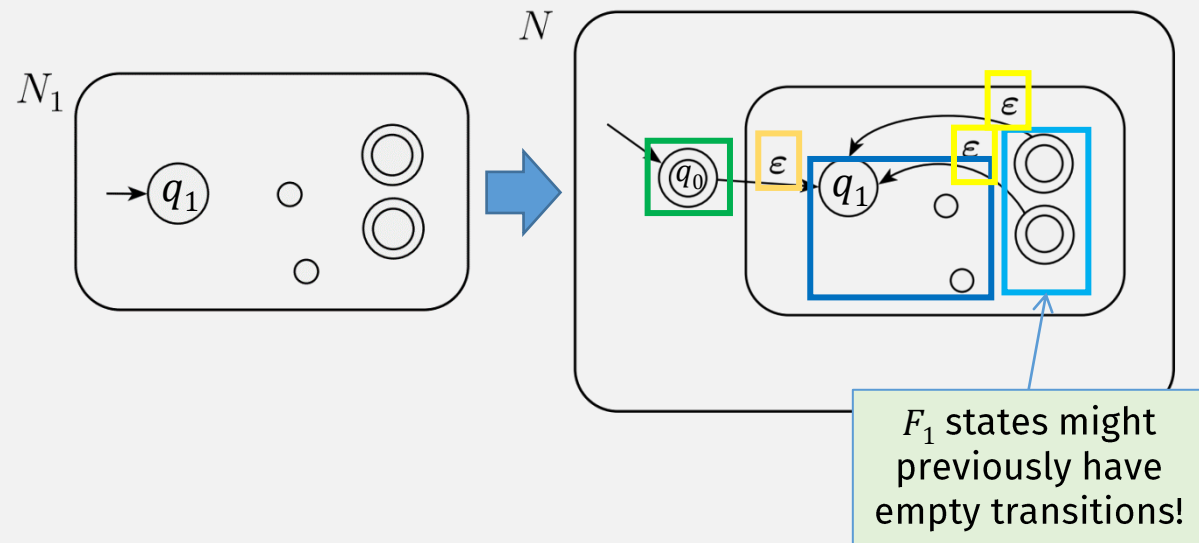
# Kleene Star is Closed for Regular Languages

(part of)

**PROOF** Let  $N_1 = (Q_1, \Sigma, \delta_1, q_1, F_1)$  recognize  $A_1$ .

$N = \text{STAR}_{\text{NFA}}(N_1) = (Q, \Sigma, \delta, q_0, F)$  to recognize  $A_1^*$ .

1.  $Q = \{q_0\} \cup Q_1$
2. The state  $q_0$  is the new start state.
3.  $F = \{q_0\} \cup F_1$
4. Define  $\delta$  so that for any  $q \in Q$  and any  $a \in \Sigma_\epsilon$ ,



$\delta(q, a) =$	}	$\delta_1(q, a?)$	$q \in Q_1$ and $q \notin F_1$	
		$\delta_1(q, a?)$	$q \in F_1$ and $a \neq \epsilon$	
		$\delta_1(q, a?) \cup \{q_1\}$	$q \in F_1$ and $a = \epsilon$	Old accept states $\epsilon$ -transition to old start state $q_1$
		$\{q_1\}$	$q = q_0$ and $a = \epsilon$	New start state $\epsilon$ -transitions to old start state $q_1$
		$\emptyset$	$q = q_0$ and $a \neq \epsilon$ .	New start state has only $\epsilon$ -transitions

Must add transitions for any added states

F<sub>1</sub> states might previously have empty transitions!

# Why These (Closed) Operations?

- Union

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$$

- Concatenation

$$A \circ B = \{xy \mid x \in A \text{ and } y \in B\}$$

- Kleene star (repetition)

$$A^* = \{x_1x_2 \dots x_k \mid k \geq 0 \text{ and each } x_i \in A\}$$

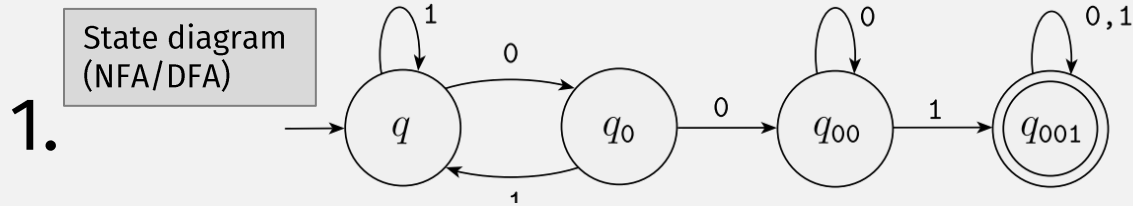
**All regular languages** can be constructed from:

- (language of) **single-char strings** (from some alphabet), ...

- And these **three closed operations!**

e.g., lang {"a"}, lang {"b"}, ...

# So Far: Regular Language Representations



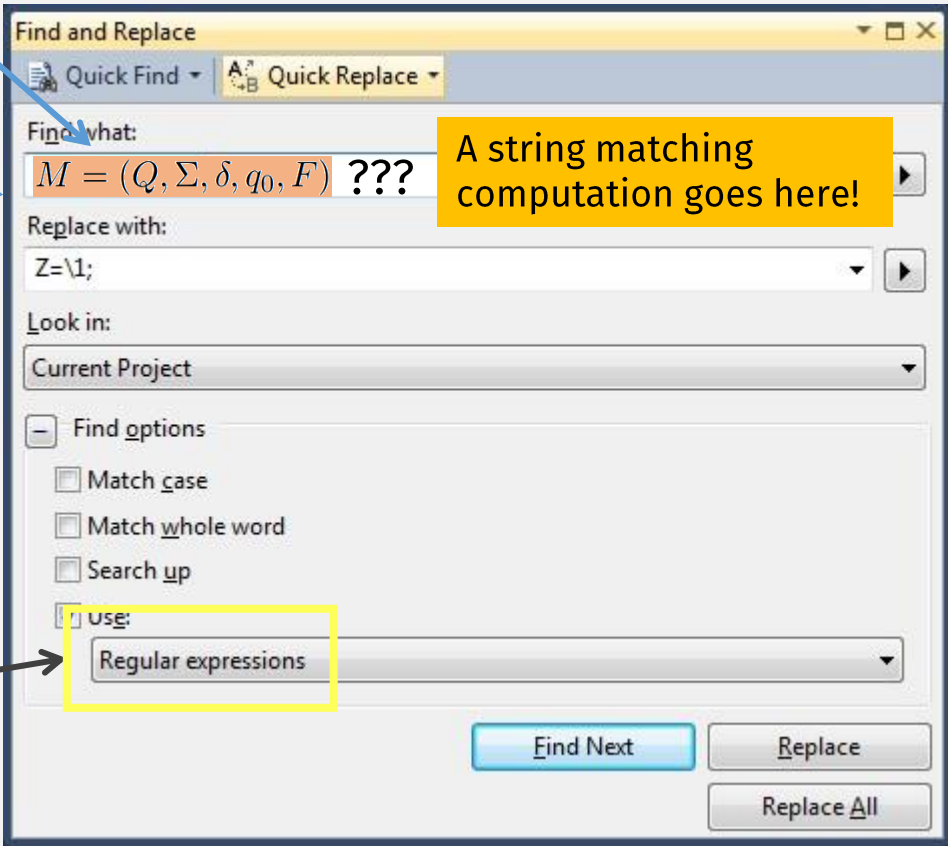
Actually, it's a real programming language, for **text search / string matching** computations

Formal description

1.  $Q = \{q_1, q_2, q_3\}$ ,
2.  $\Sigma = \{0,1\}$ ,
3.  $\delta$  is described as
4.  $q_1$  is the start state
5.  $F = \{q_2\}$

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

(hard to write)



Our Running Analogy:

- Set of all **regular languages** ~ a “programming language”
- One **regular language** (or any equiv representation) ~ a “program”

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

Need a more concise (textual) notation??

# Regular Expressions: A Widely Used Programming Language (usually within tools / languages)

- Unix / Linux
- Java
- Python
- Web APIs

java.util.regex

## Class Pattern

java.lang.Object

java.util.regex.Pattern

```
GREP(1)                                General Commands Manual                                GREP(1)
NAME
    grep, egrep, fgrep, rgrep - print lines matching a pattern
SYNOPSIS
    grep [OPTIONS] PATTERN [FILE...]
    grep [OPTIONS] [-e PATTERN | -f FILE] [FILE...]
DESCRIPTION
    grep searches the named input FILES (or standard input if no files are
    named, or if a single hyphen-minus (-) is given as file name) for lines
    containing a match to the given PATTERN. By default, grep prints the
    matching lines.
```

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

## About regular expressions (regex)

Analytics supports regular expressions so you can create more flexible definitions for things like [view filters](#), [goals](#), [segments](#), [audiences](#), [content groups](#), and [channel groupings](#).

This article covers regular expressions in both Universal Analytics and Google Analytics 4.

In the context of Analytics, regular expressions are specific sequences of characters that broadly or narrowly match patterns in your Analytics data.

For example, if you wanted to create a view filter to exclude site data generated by your own employees, you could use a regular expression to exclude any data from the entire range of IP addresses that serve your employees. Let's say those IP addresses range from 198.51.100.1 - 198.51.100.25. Rather than enter 25 different IP addresses, you could create a regular expression like `198\.\d*\.\d*` that matches the entire range of addresses.

## — Regular expression operations

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

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

# Why These (Closed) Operations?

- Union

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$$

- Concatenation

$$A \circ B = \{xy \mid x \in A \text{ and } y \in B\}$$

- Kleene star (repetition)

$$A^* = \{x_1x_2 \dots x_k \mid k \geq 0 \text{ and each } x_i \in A\}$$

All regular languages can be constructed from:

- (language of) **single-char strings** (from some alphabet), and
- these **three closed operations!**

**They are used to define regular expressions!**

# Regular Expressions: Formal Definition

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

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
2.  $\epsilon$ ,
3.  $\emptyset$ ,
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.

This is a recursive definition

# *Flashback:* Recursive Definitions

**Recursive definitions** are definitions with a self-reference

A valid recursive definition must have:

- **base case** and
- **recursive case** (with a “smaller” self-reference)

# Flashback: Recursive Definitions

```
function factorial( n )  
{  
  if ( n == 0 )  
    return 1;  
  else  
    return n * factorial( n - 1 );  
}
```

Base case

Recursive case

Self-reference

Recursive call with  
"smaller" argument



# *Flashback:* Recursive Definitions

A Natural Number is either:

Self-reference

Base case

• **Zero**, or

Recursive case

• the **Successor** of a **Natural Number**

“smaller” argument

# Flashback: Recursive Definitions

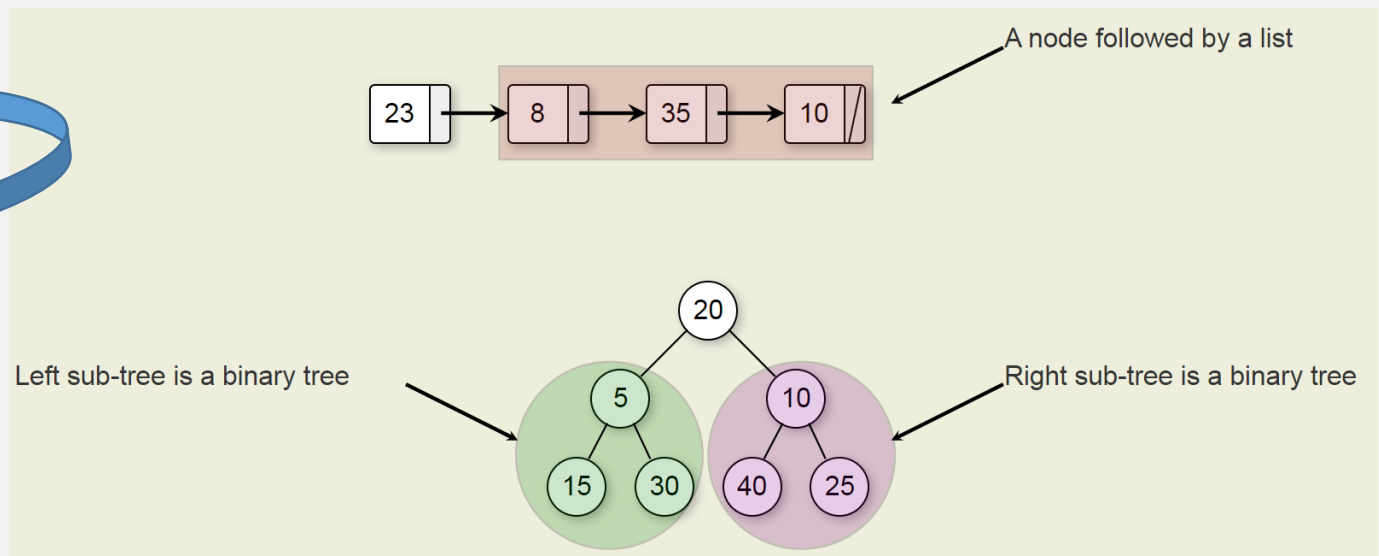
```
/* Linked list Node*/  
class Node {  
    int data;  
    Node next;  
}
```

Smaller self-reference

Q: Where's the base case??

I call it my billion-dollar mistake. It was the invention of the null reference in 1965.

— Tony Hoare —



Data structures are commonly defined recursively

# Regular Expressions: Formal Definition

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

3 Base Cases

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ , (A lang containing  $a$ ) length-1 string

2.  $\epsilon$ , (A lang containing) the empty string (This is the 3<sup>rd</sup> use of the  $\epsilon$  symbol!)

3.  $\emptyset$ , The empty set (i.e., a lang containing no strings)

union

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

concat

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

star

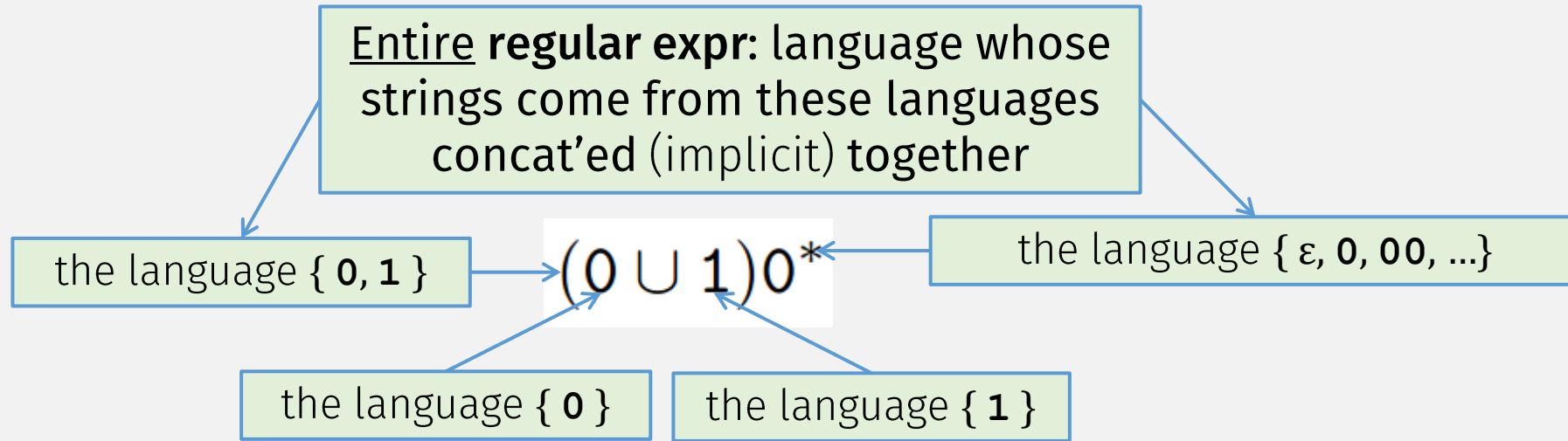
6.  $(R_1^*)$ , where  $R_1$  is a regular expression.

3 Recursive Cases

Note:

- A **regular expression** represents a **language**
- The *set of all regular expressions* represents a *set of languages*

# Regular Expression: Concrete Example



## • Operator Precedence:

- Parentheses
- Kleene Star
- Concat (sometimes use  $\circ$ , sometimes implicit)
- Union

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

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
2.  $\epsilon$ ,
3.  $\emptyset$ ,
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.

# Regular Expression: More Examples

$$0^*10^* = \{w \mid w \text{ contains a single } 1\}$$

$$\Sigma^*1\Sigma^* = \{w \mid w \text{ has at least one } 1\}$$

$\Sigma$  in regular expression = "any char"

$$1^*(01^+)^* = \{w \mid \text{every } 0 \text{ in } w \text{ is followed by at least one } 1\}$$

let  $R^+$  be shorthand for  $RR^*$

$$(0 \cup \epsilon)(1 \cup \epsilon) = \{\epsilon, 0, 1, 01\}$$

$0 \cup \epsilon$  describes the language  $\{0, \epsilon\}$

$$1^*\emptyset = \emptyset$$

$A \circ B = \{xy \mid x \in A \text{ and } y \in B\}$

nothing in  $B =$  nothing in  $A \circ B$

$$\emptyset^* = \{\epsilon\}$$

Star of any lang has  $\epsilon$

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

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
2.  $\epsilon$ ,
3.  $\emptyset$ ,
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.

# Regular Expressions = Regular Langs?

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

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
2.  $\epsilon$ ,
3.  $\emptyset$ ,
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.

We would like to say:

- A **regular expression** represents a **regular language**
- The *set of all regular expressions* represents the *set of all regular languages*

(But we have to prove it)

Thm: A Lang is Regular **iff** Some Reg Expr Describes It

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

⇐ If a language is described by a reg expression, then it's regular  
(Easier)

- Key step: **convert reg expr → equivalent NFA!**
- (Hint: we mostly did this already when discussing closed ops)

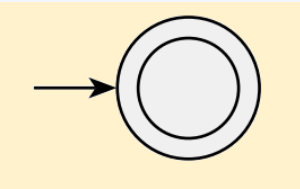
How to show that a language is regular?

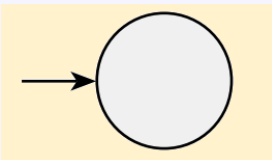
Construct a **DFA or NFA!**

# RegExpr $\rightarrow$ NFA

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

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

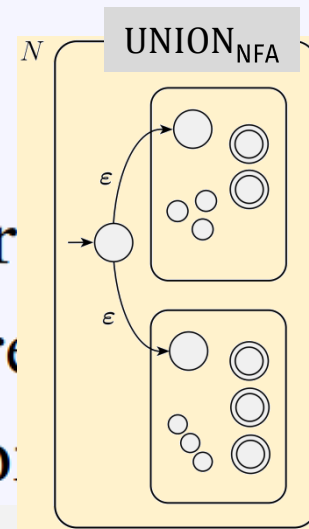
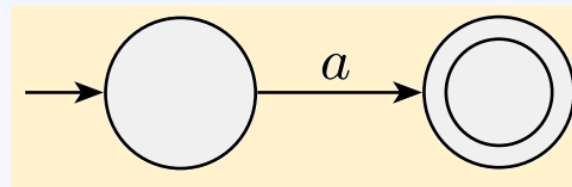
2.  $\epsilon$ , 

3.  $\emptyset$ , 

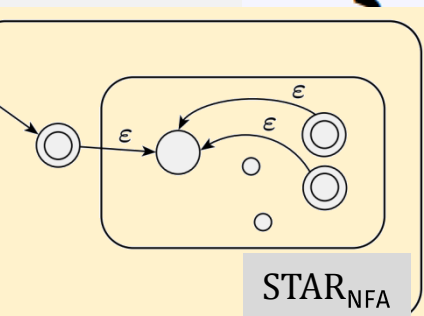
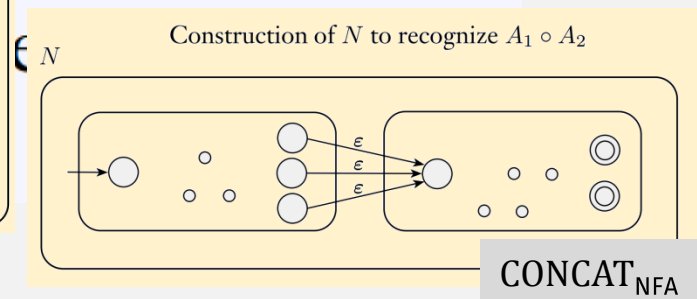
4.  $(R_1 \cup R_2)$ , where  $R_1$  and  $R_2$  are regular expressions,

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.



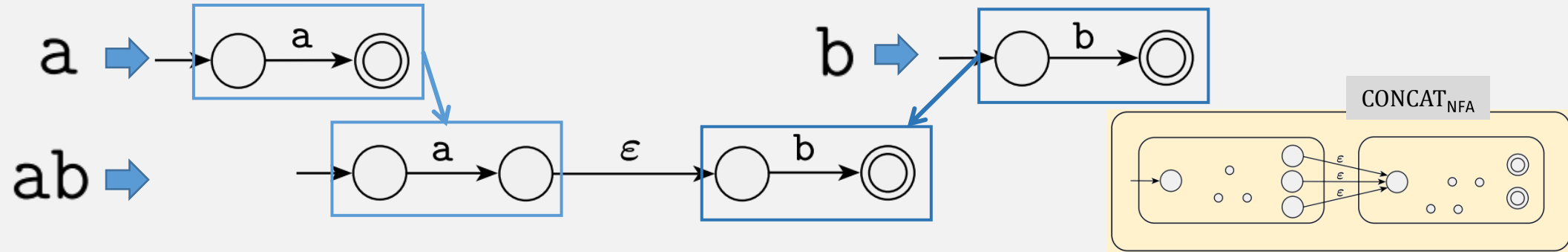
expressions,





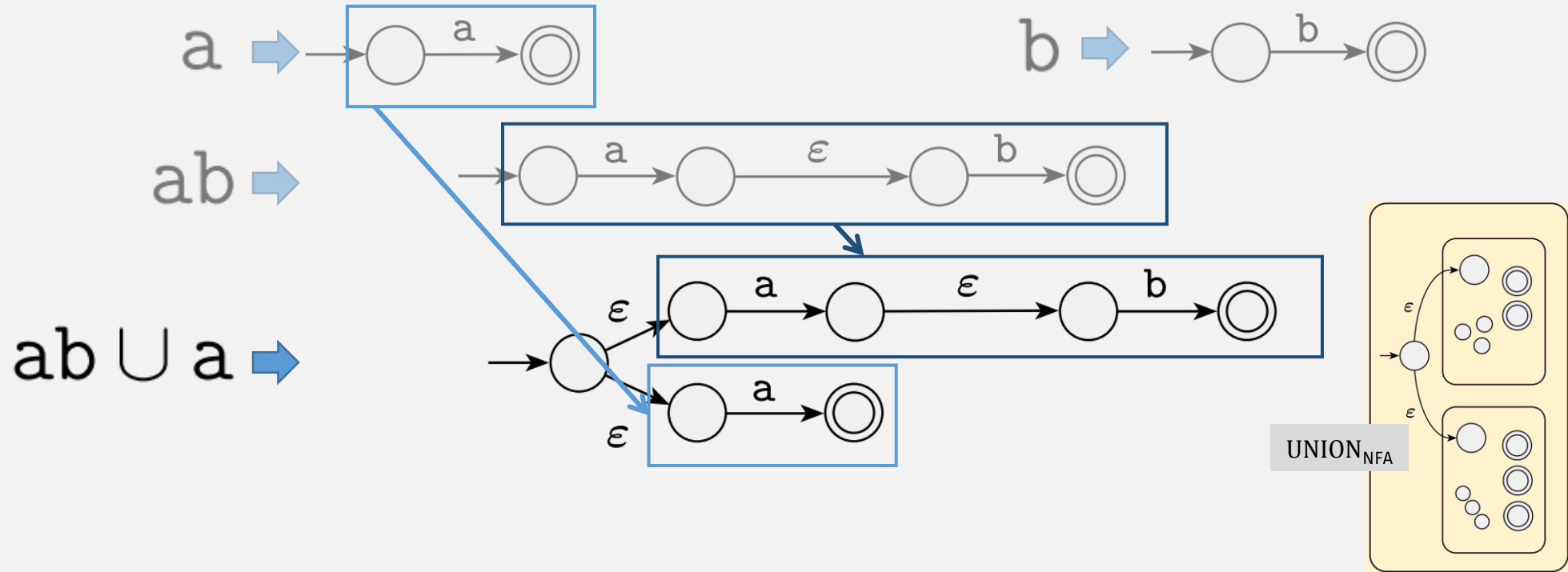
# RegExpr $\rightarrow$ NFA: Example

convert the regular expression  $(ab \cup a)^*$  to an NFA ... step by step



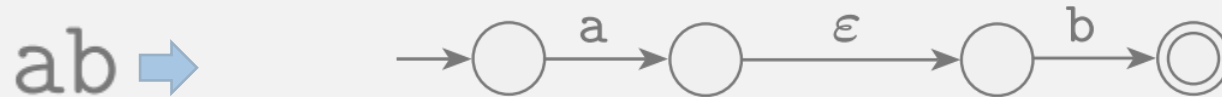
# RegExpr $\rightarrow$ NFA: Example

convert the regular expression  $(ab \cup a)^*$  to an NFA

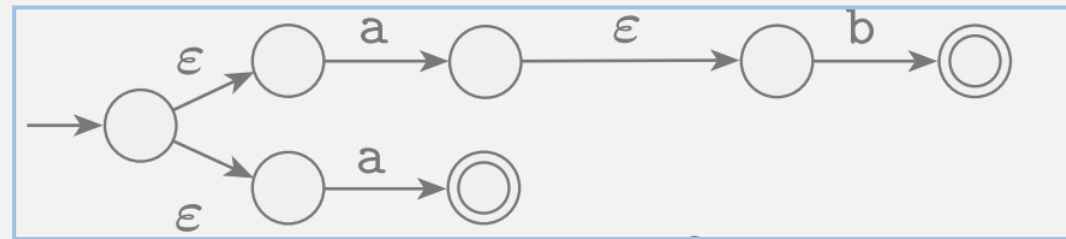


# RegExpr $\rightarrow$ NFA: Example

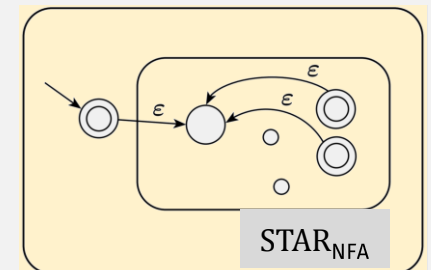
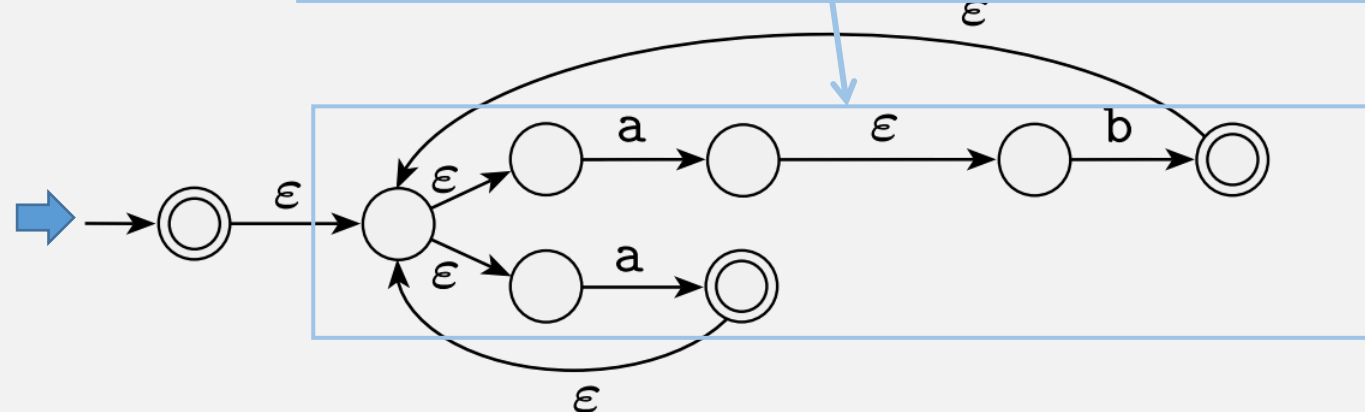
convert the regular expression  $(ab \cup a)^*$  to an NFA



$ab \cup a$   $\rightarrow$



$(ab \cup a)^*$   $\rightarrow$



# Thm: A Lang is Regular iff Some Reg Expr Describes It

⇒ If a language is regular, then it's described by a reg expression  
(Harder)

• Key step: Convert an ~~DFA~~ or **NFA** → equivalent **Regular Expression**

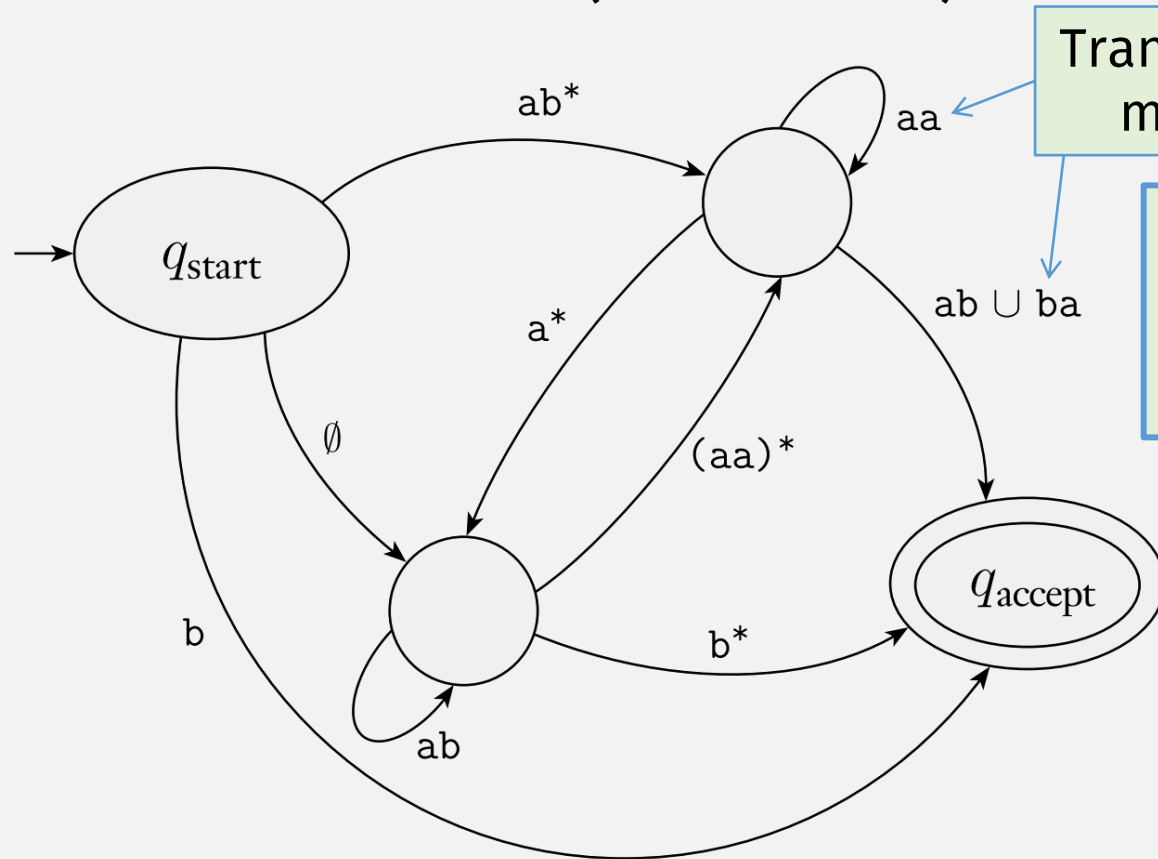
• First, we need another kind of finite automata: a **GNFA**

⇐ If a language is described by a reg expression, then it's regular  
(Easier)

☑ • Key step: Convert the regular expression → an equivalent NFA!

(full proof requires writing Statements and Justifications, and creating an "Equivalence" Table)

# Generalized NFAs (GNFAs)



Transition can read multiple chars

plain NFA = GNFA with single char regular expr transitions

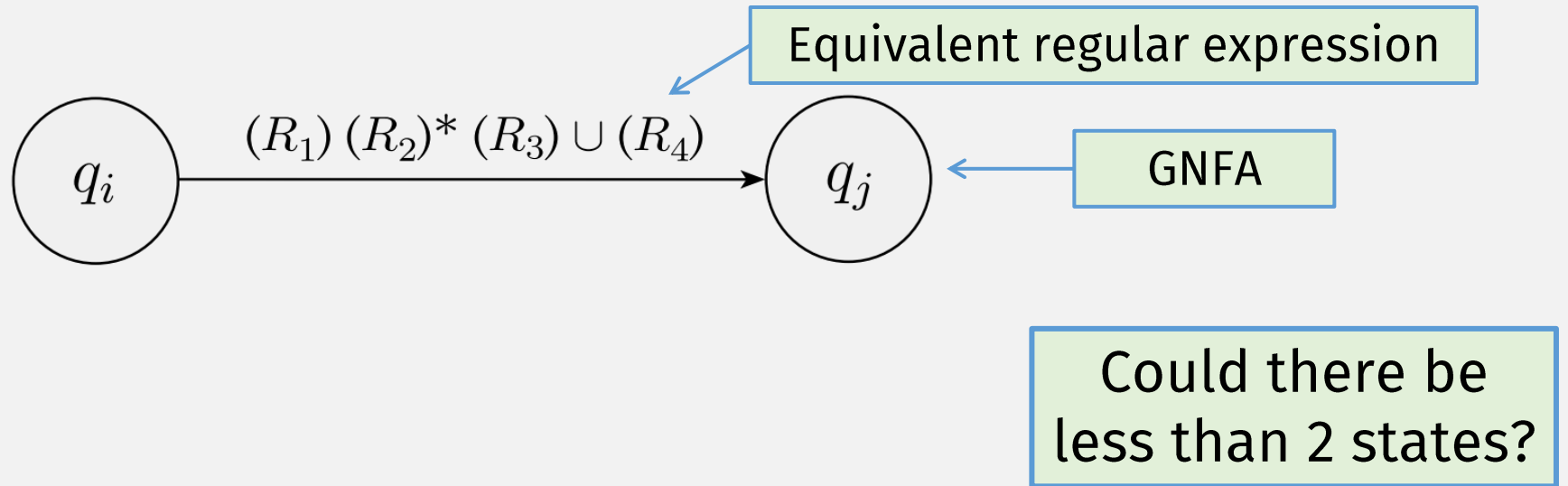
Goal: convert **GNFAs** to equivalent **Regular Exprs**

- GNFA = NFA with regular expression transitions

# GNFA $\rightarrow$ RegExpr function :

On GNFA input  $G$ :

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



# GNFA $\rightarrow$ RegExpr Preprocessing

- Modify input machine to have:

- **New start state:**

- No incoming transitions
- $\epsilon$  transition to old start state

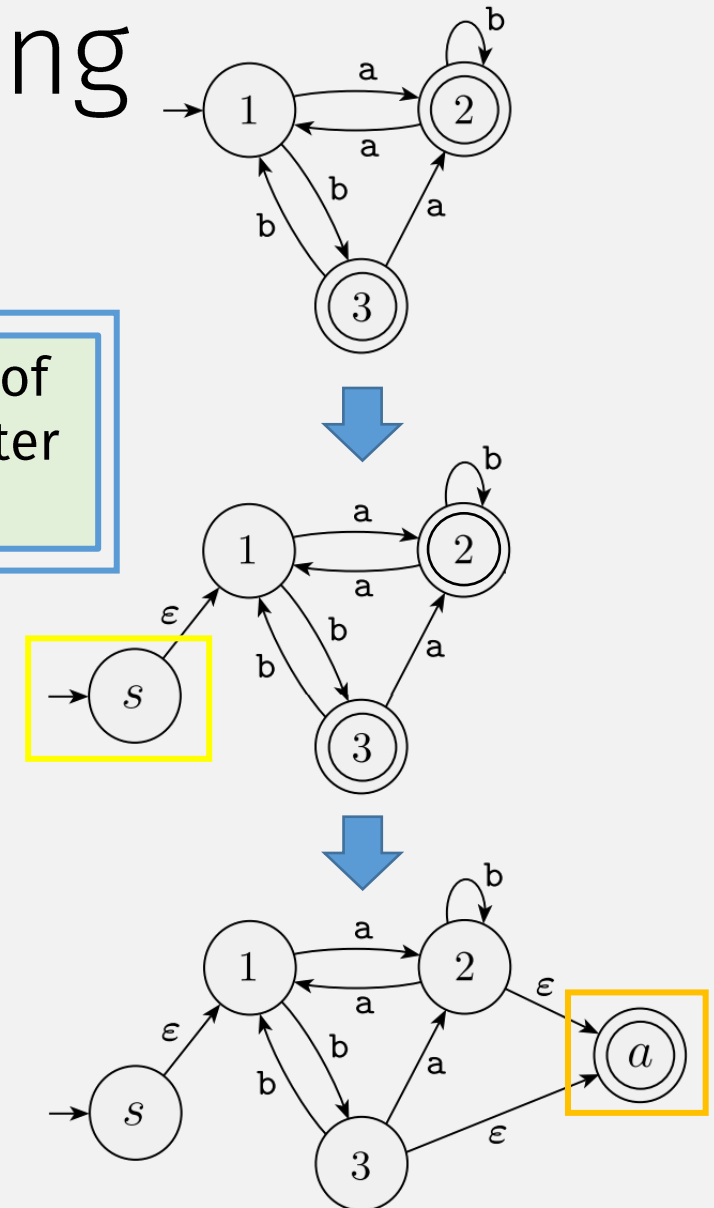
Does this change the language of the machine? i.e., are before/after machines equivalent?

- **New, single accept state:**

- With  $\epsilon$  transitions from old accept states

Modified machine always has 2+ states:

- New start state
- New accept state

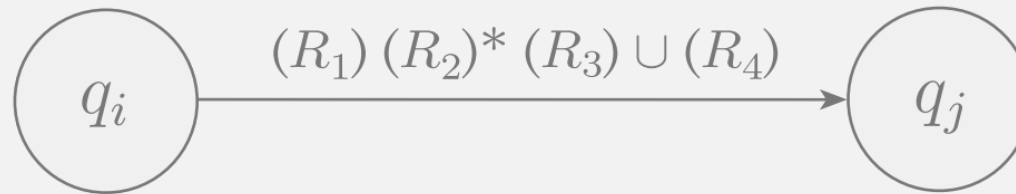


# GNFA $\rightarrow$ RegExpr function (recursive)

On GNFA input  $G$ :

Base  
Case

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



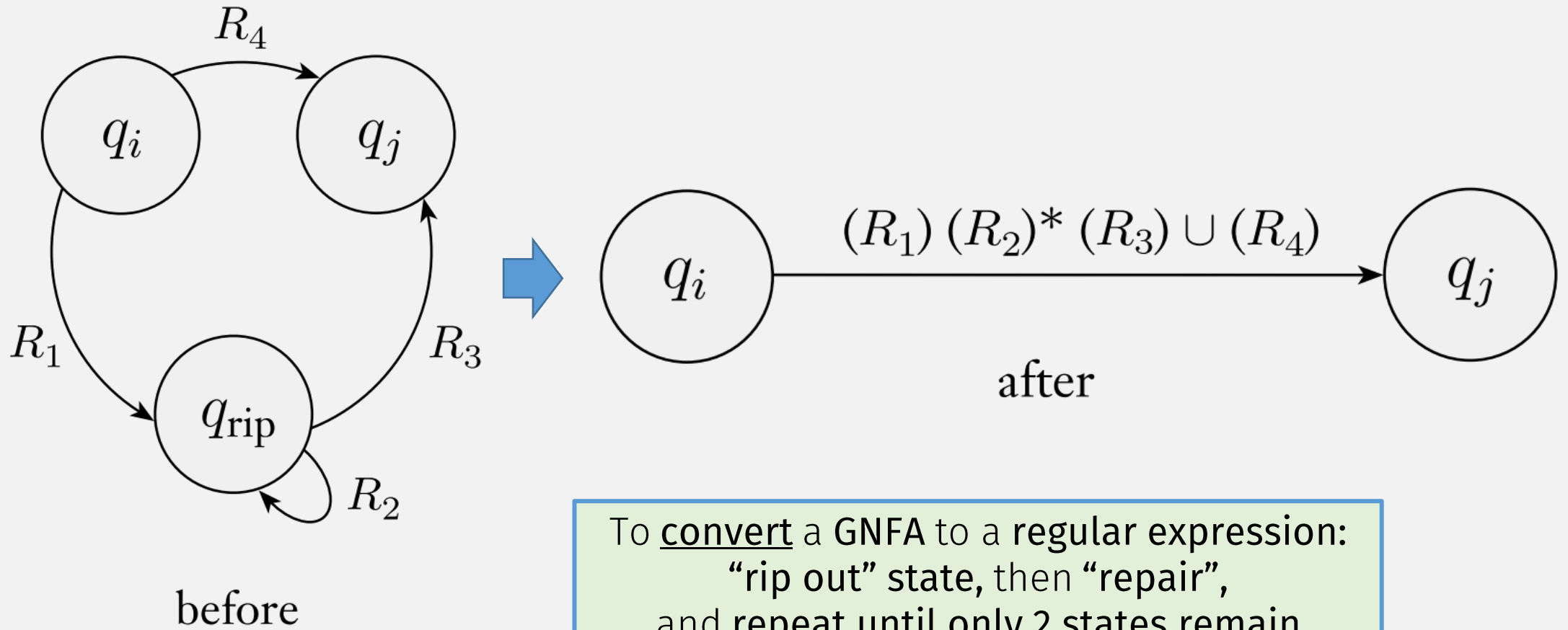
Recursive  
Case

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

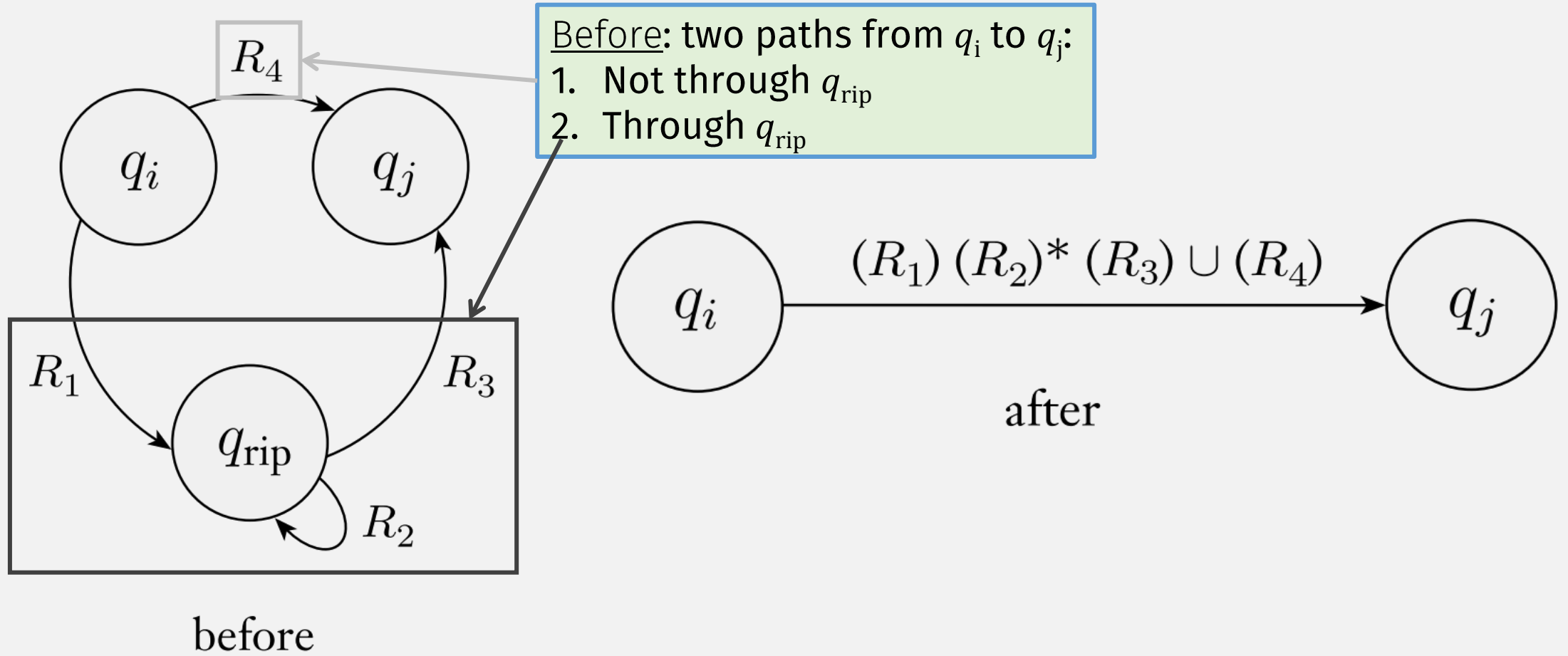
Recursive definitions have:  
- base case and  
- recursive case  
(with “smaller” self-reference)



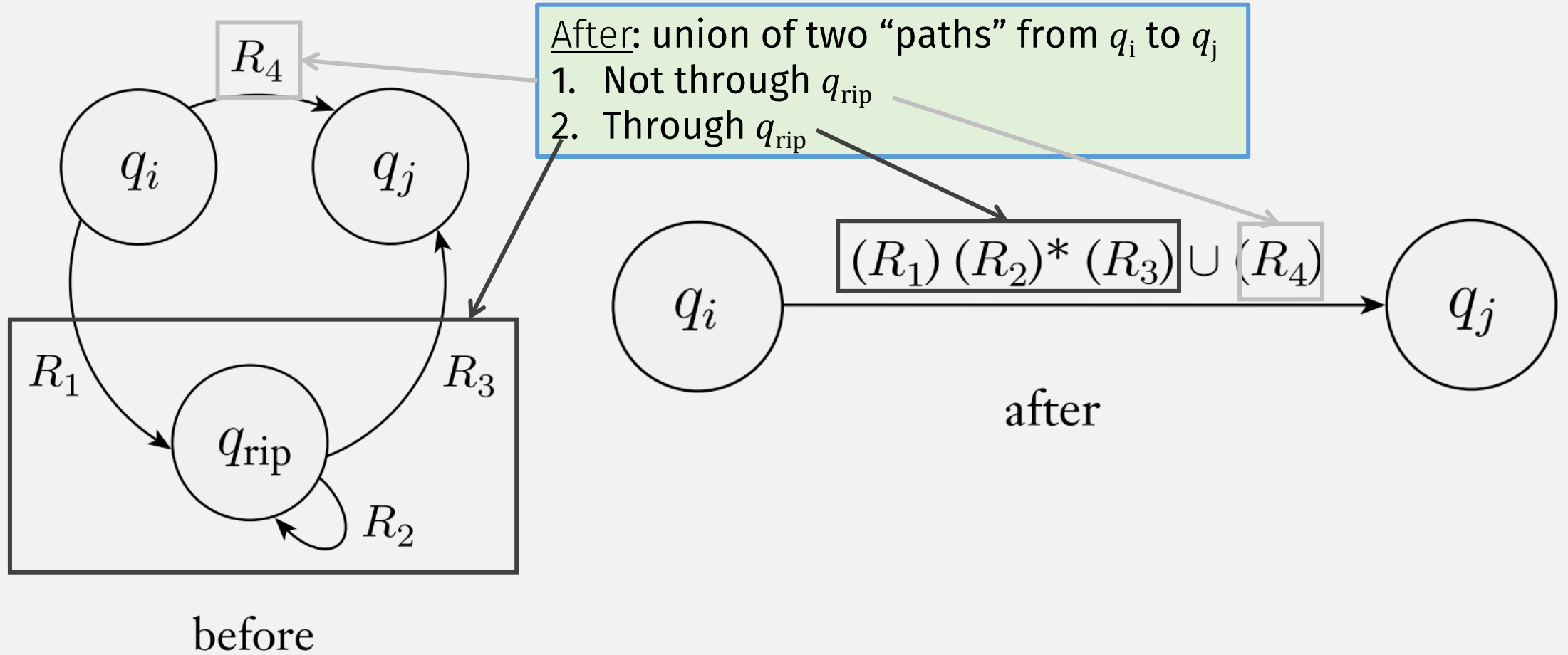
# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step



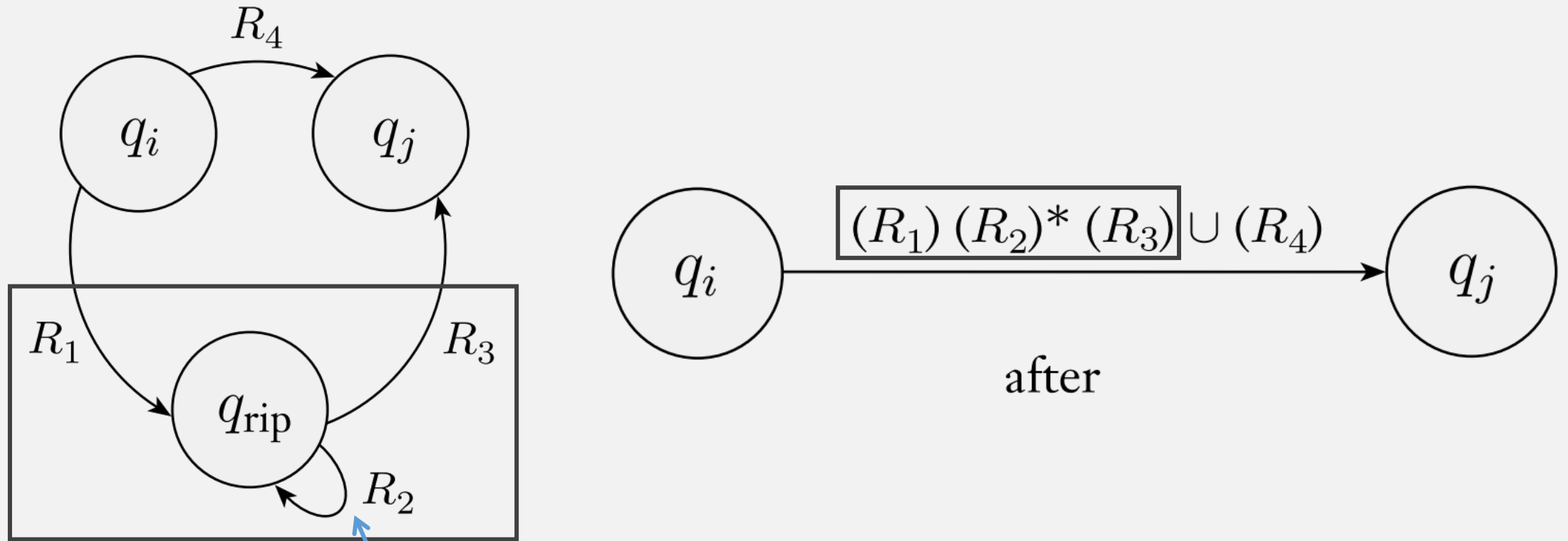
# GNFA $\rightarrow$ RegExpr: “Rip/Repair” step



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



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

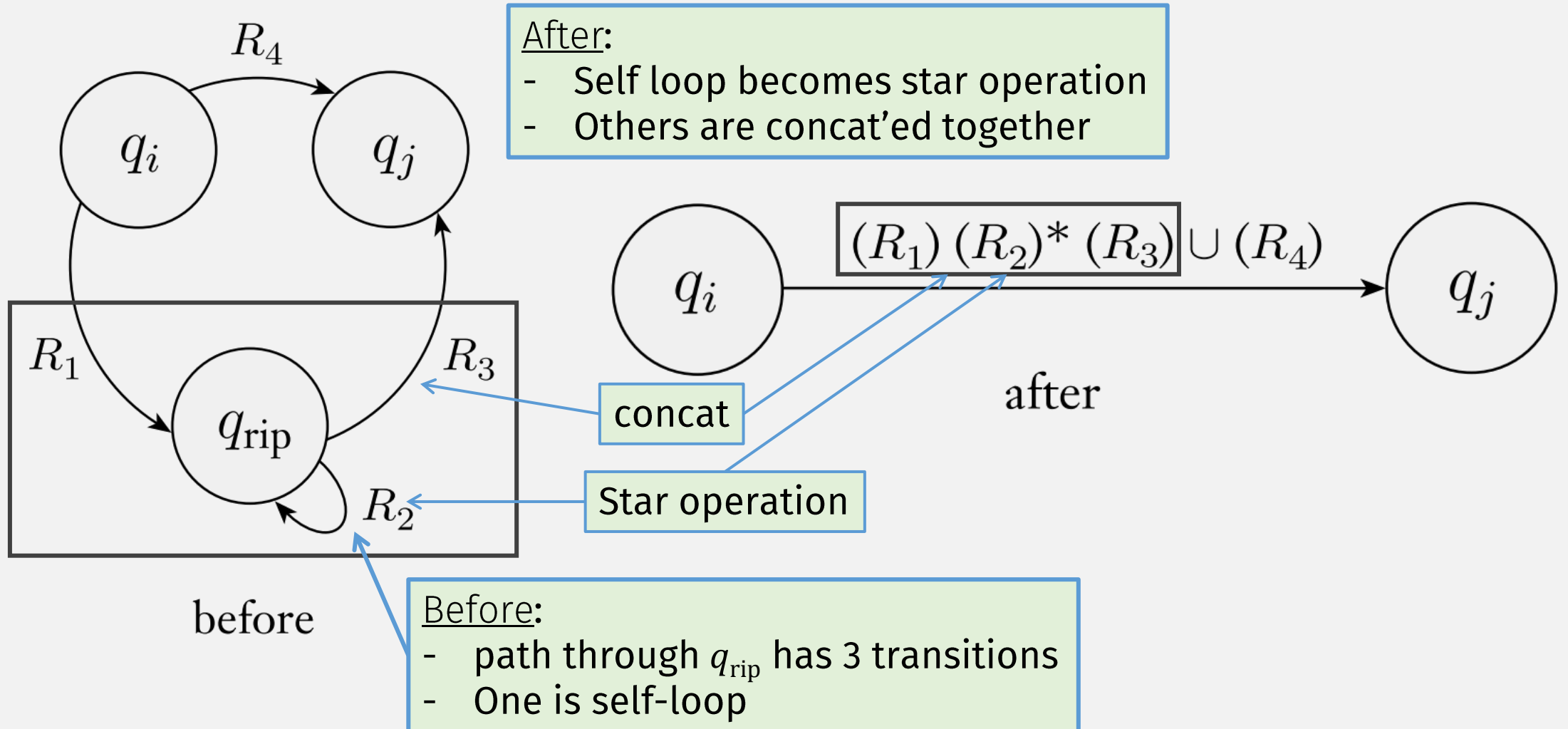


before

Before:

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

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



# Thm: A Lang is Regular iff Some Reg Expr Describes It

⇒ If a language is regular, then it's described by a regular expr

Need to convert DFA or NFA to Regular Expression ...

- Use GNFA → RegExpr to convert GNFA → equiv regular expression!



???

This time, let's really prove  
equivalence!  
(we previously "proved" it  
with an Examples Table)

⇐ If a language is described by a regular expr, then it's regular

- ✓ • Convert regular expression → equiv NFA!

# GNFA $\rightarrow$ RegExpr Correctness

- **Correct** = input and output are **equivalent**
- **Equivalent** = the language does not change (same strings)!

Statement to Prove:

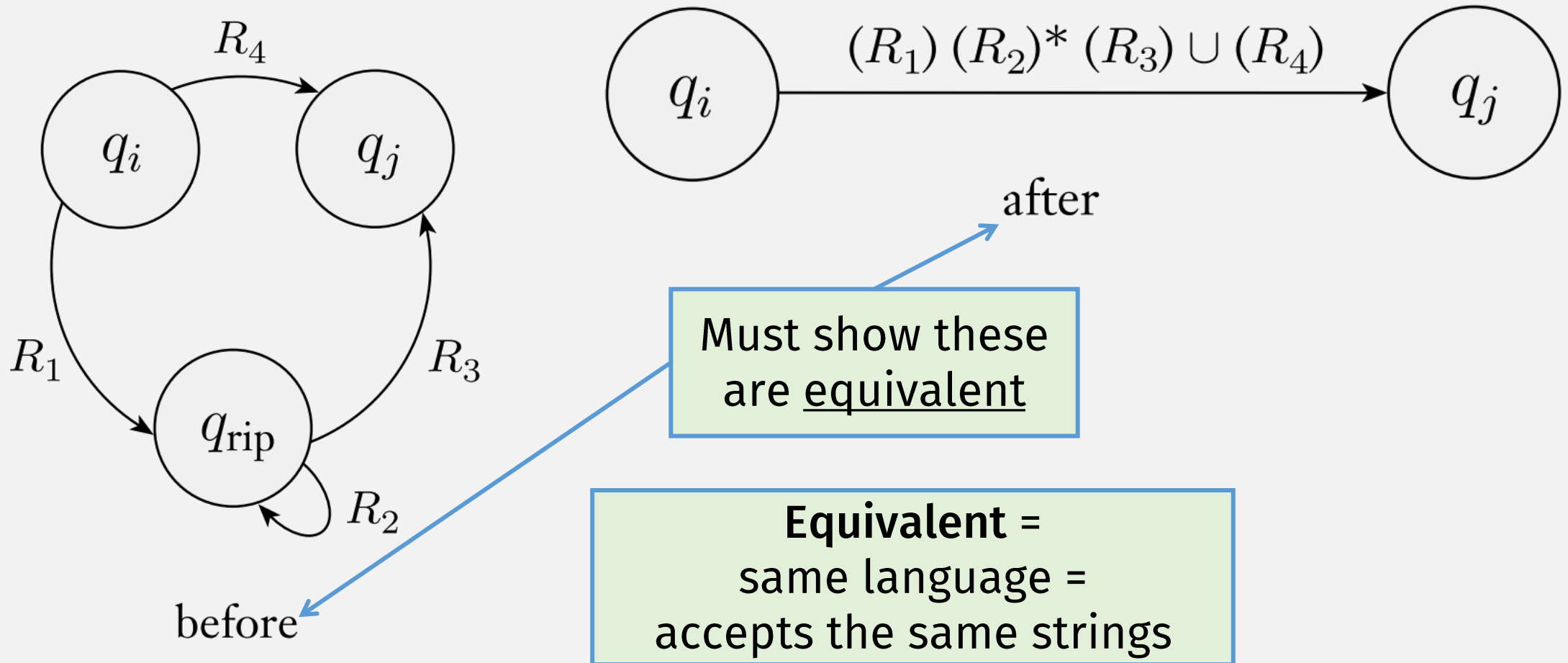
???

$$\text{LANGOF} ( G ) = \text{LANGOF} ( R )$$

We are ready to really  
prove equivalence!  
(we previously “proved” it  
with some examples)

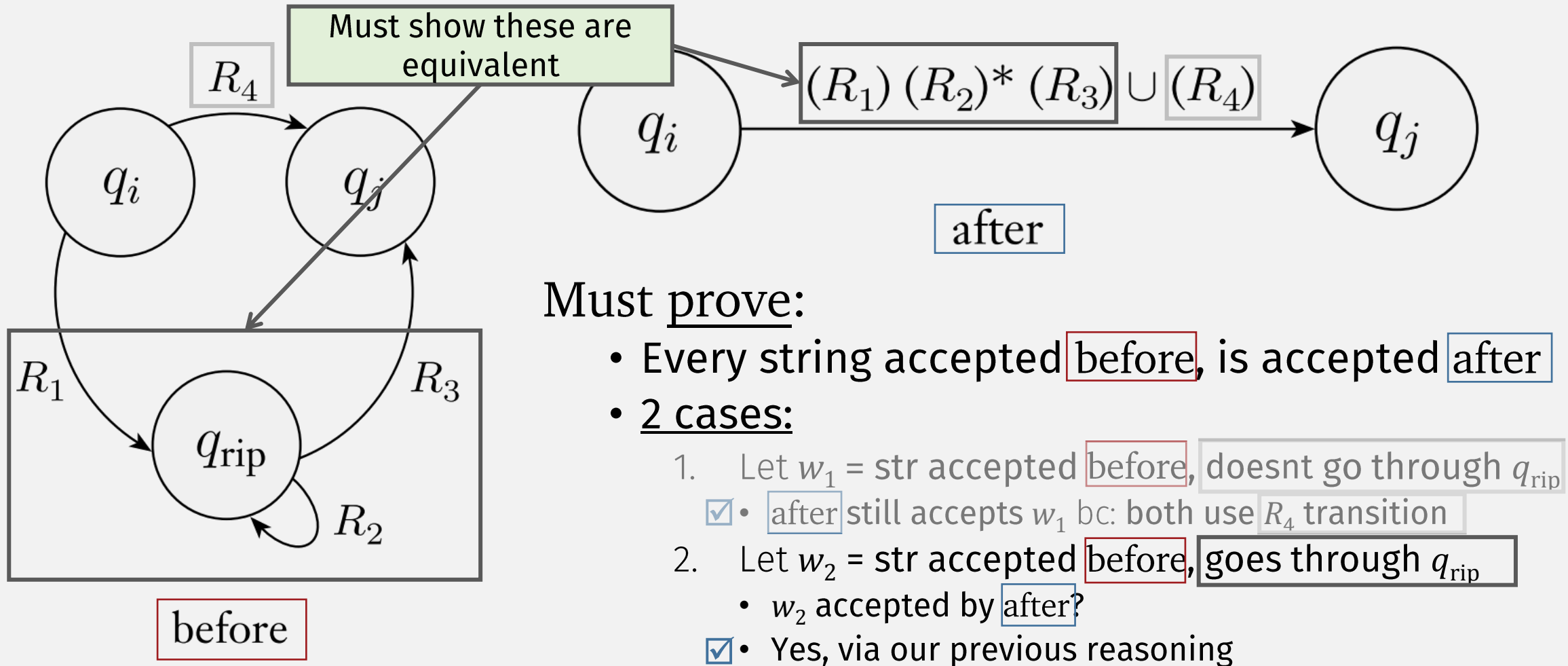
- where:
  - $G$  = a GNFA
  - $R$  = a Regular Expression = GNFA $\rightarrow$ RegExpr( $G$ )
- Key step: the rip/repair step

# GNFA $\rightarrow$ RegExpr: Rip/Repair Correctness





# GNFA $\rightarrow$ RegExpr: Rip/Repair Correctness



# GNFA $\rightarrow$ RegExpr Equivalence

- **Equivalent** = the language does not change (i.e., same set of strings)!

Statement to Prove: input output ???

$$\text{LANGOF} ( G ) = \text{LANGOF} ( R )$$

This time, let's really prove equivalence!  
(we previously "proved" it with some examples)

- where:
  - $G$  = a GNFA
  - $R$  = a Regular Expression = GNFA $\rightarrow$ RegExpr( $G$ )

Language could be infinite set of strings!

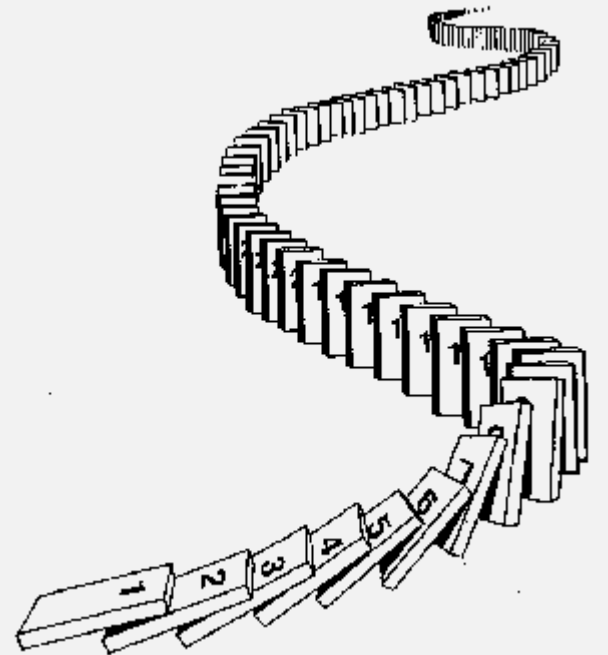
(how can we show equivalence for a possibly infinite set of strings?)

Recursion!

*Next Time*

# Inductive Proofs

(recursive)



# Proof by Induction

*Previously:*

## Recursive function

- Use it when: writing a function involving a recursive definition

*Now:*

**Proof by induction** (recursion) = “a recursive proof”

- Use it when: proving something involving a recursive definition

The recursive definition is (always) the key!

A valid recursive definition has:

- **base case(s)** and
- **recursive case(s)** (with “smaller” self-reference)

# Proof by Induction

(A proof for **each case** of some recursive definition)

To Prove: *Statement* for recursively defined “thing”  $x$ :

1. Prove: *Statement* for **base case** of  $x$

2. Prove: *Statement* for **recursive case** of  $x$ :

- Assume: **induction hypothesis (IH)**

i.e., *Statement* is true for **some  $x_{\text{smaller}}$**  (This is just the recursive part from the recursive definition!)

- E.g., if  $x$  is number, then “smaller” = lesser number

- Prove: *Statement* for  $x$ , using IH (and known definitions, theorems ...)

- Typically: show that going from  $x_{\text{smaller}}$  to larger  $x$  is true! i.e., a normal proof

A valid recursive definition has:

- **base case(s)** and

- **recursive case(s)** (with “smaller” self-reference)

# Natural Numbers Are Recursively Defined

A Natural Number is:

Base Case

• **0**

Self-reference

Recursive Case

• Or  **$k + 1$** , where  $k$  is a Natural Number

Recursive definition is valid because self-reference is “smaller”

So, proving things about:  
recursive Natural Numbers requires  
recursive proof,  
i.e., **proof by induction!**

A valid recursive definition has:

- **base case(s)** and
- **recursive case(s)** (with “smaller” self-reference)

# Proof By Induction Example (Sipser Ch 0)

Prove true:  $P_t = PM^t - Y \left( \frac{M^t - 1}{M - 1} \right)$

- $P_t$  = loan balance after  $t$  months
- $t$  = # months
- $P$  = principal = original amount of loan
- $M$  = interest (multiplier)
- $Y$  = monthly payment

(Details of these variables not too important here)

# Proof By Induction Example (Sipser Ch 0)

Prove true:  $P_t = PM^t - Y \left( \frac{M^t - 1}{M - 1} \right)$

A proof by induction follows the cases of the recursive definition (here, natural numbers) that the induction is “on”

Proof: by **induction** on natural number  $t$

**Base Case,  $t = 0$ :**

A Natural Number is:

- 0
- Or  $k + 1$ , where  $k$  is a natural number

$$P_0 = PM^0 - Y \left( \frac{M^0 - 1}{M - 1} \right) = P$$

Plug in  $t = 0$

Simplify

$P_0 = P$  is a true statement!  
(amount owed at start = loan amount)



# Proof By Induction Example (Sipser Ch 0)

Prove true:  $P_t = PM^t - Y \left( \frac{M^t - 1}{M - 1} \right)$

A **proof by induction** follows cases of recursive definition (here, natural numbers) that the induction is "on"

**Inductive Case:**  $t = k + 1$ , for some natural num  $k$

A Natural Number is:  
 • 0   
 •  $k + 1$ , for some nat num  $k$

- **Inductive Hypothesis (IH)**, assume statement is true for some  $t =$  (smaller)  $k$

IH plugs in "smaller"  $k$

$$P_k = PM^k - Y \left( \frac{M^k - 1}{M - 1} \right)$$

Goal statement to prove, for  $t = k + 1$ :

$$P_{k+1} = PM^{k+1} - Y \left( \frac{M^{k+1} - 1}{M - 1} \right)$$

Simplify, to get to goal statement

Write  $t = k + 1$  case in terms of "smaller"  $k$

Plug in IH for  $P_k$

- Proof of Goal:

$$P_{k+1} = P_k M - Y$$

Definition of Loan:  
 amt owed in month  $k + 1 =$   
 amt owed in month  $k * interest  $M -$  amt paid  $Y$$

# In-class Exercise: Proof By Induction

A **proof by induction** follows cases of recursive definition (here, natural numbers) that the induction is “on”

Prove: ( $z \neq 1$ )

$$\sum_{i=0}^m z^i = \frac{1 - z^{m+1}}{1 - z}$$

A Natural Number is:

- 0
- $k + 1$ , for some nat num  $k$

**Use Proof by Induction.**

Make sure to: clearly **state what the induction is “on”**

i.e., which recursively defined value (and its **type**) will the proof focus on

# Proof by Induction: CS 420 Example

*Statement to prove:*

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

- Where:
  - $G$  = a GNFA
  - $R$  = a Regular Expression  $\text{GNFA} \rightarrow \text{RegExpr}( G )$
- i.e.,  $\text{GNFA} \rightarrow \text{RegExpr}$  must not change the language!

This time, let's really prove equivalence!  
(we previously “proved” it with some examples)

# Proof by Induction: CS 420 Example

*Statement to prove:*

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

Recursively defined “thing”

Proof: by Induction on # of states in  $G$

1. Prove Statement is true for base case

$G$  has 2 states

Why is this an ok  
base case  
(instead of zero)?

(Modified) Recursive definition:

- A “NatNumber > 1” is:
- 2
  - Or  $k + 1$ , where  $k$  is a “NatNumber > 1”

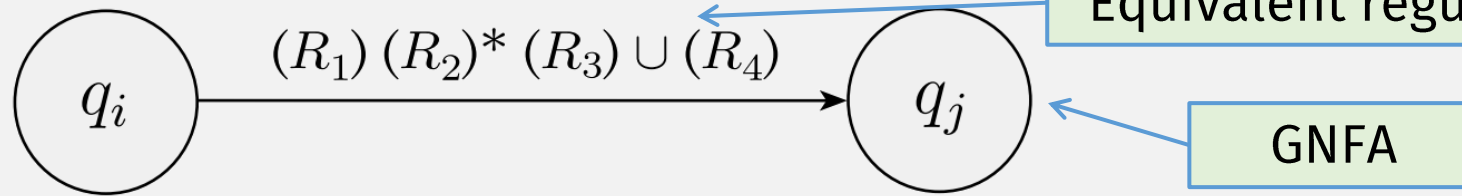
Last Time

# GNFA $\rightarrow$ RegExpr (recursive) function

On GNFA input  $G$ :

Base Case

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



# Proof by Induction: CS 420 Example

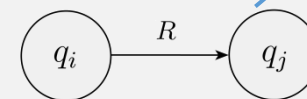
Statement to prove:

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

Proof: by Induction on # of states in  $G$

✓ 1. Prove Statement is true for base case

$G$  has 2 states



Plug in

## Statements

1.  $\text{LANGOF} ( \text{GNFA} ( q_i \xrightarrow{R} q_j ) ) = \text{LANGOF} ( R )$
  2.  $\text{GNFA} \rightarrow \text{RegExpr} ( \text{GNFA} ( q_i \xrightarrow{R} q_j ) ) = R$
- $$\text{LANGOF} ( \text{GNFA} ( q_i \xrightarrow{R} q_j ) ) = \text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( \text{GNFA} ( q_i \xrightarrow{R} q_j ) ) )$$

Plug in  $R$

## Justifications

1. Definition of GNFA
2. Definition of  $\text{GNFA} \rightarrow \text{RegExpr}$  (base case)
3. From (1) and (2)

Goal

Don't forget the Statements / Justifications !

# Proof by Induction: CS 420 Example

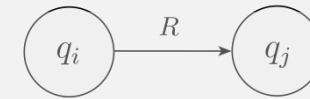
Statement to prove:

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

Proof: by Induction on # of states in  $G$

✓ 1. Prove Statement is true for base case

$G$  has 2 states



2. Prove Statement is true for recursive case:  $G$  has  $> 2$  states

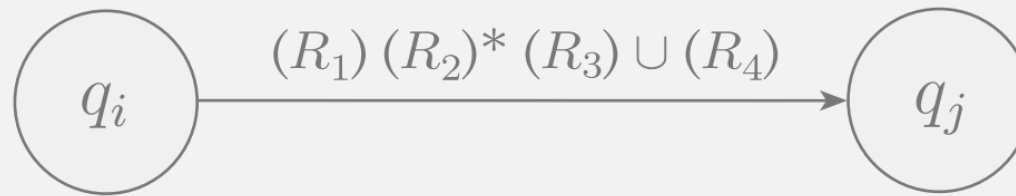
Last Time

# GNFA $\rightarrow$ RegExpr (recursive) function

On GNFA input  $G$ :

Base  
Case

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



- Else:

Recursive  
Case

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

Recursive call  
(with a “smaller”  $G'$ )



# Proof by Induction: CS 420 Example

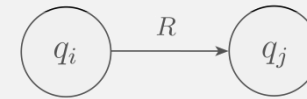
Statement to prove:

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

Proof: by Induction on # of states in  $G$

✓ 1. Prove Statement is true for base case

$G$  has 2 states



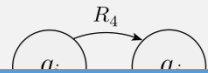
2. Prove Statement is true for recursive case:

$G$  has > 2 states

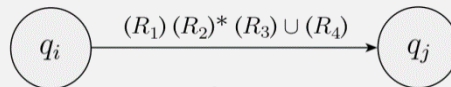
IH Assumption

- Assume the induction hypothesis (IH):
  - *Statement* is true for smaller  $G'$
- Use it to prove *Statement* is true for  $G > 2$  states
  - Show that going from  $G$  to smaller  $G'$  is true!

$$\begin{aligned} &\text{LANGOF} ( G' ) \\ &= \\ &\text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( G' ) ) \\ &\text{(Where } G' \text{ has less states than } G \text{)} \end{aligned}$$



Don't forget the Statements / Justifications !



smaller  $G'$

Show that "rip/repair" step ✓ converts  $G$  to smaller, equivalent  $G'$

before

$G$

# Proof by Induction: CS 420 Example

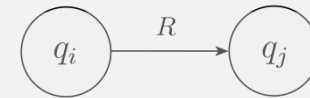
Statement to prove:

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

Proof: by Induction on # of states in  $G$

✓ 1. Prove Statement is true for base case

$G$  has 2 states



✓ 2. Prove Statement is true for recursive case:

$G$  has > 2 states

- Assume the inductive hypothesis:
  - Statement is true for  $n$  states
- Use it to prove the statement for  $n+1$  states
- Show that going from  $n$  to  $n+1$  is true!

Known "facts" available to use:

- ✓ IH
- ✓ Equiv of Rip/Repair step
- ✓ Def of GNFA  $\rightarrow$  RegExpr

$$\begin{aligned} & \text{LANGOF} ( G' ) \\ & = \\ & \text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( G' ) ) \\ & \text{(Where } G' \text{ has less states than } G \text{)} \end{aligned}$$

## Statements

1.  $\text{LANGOF} ( G' ) = \text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( G' ) )$
2.  $\text{LANGOF} ( G ) = \text{LANGOF} ( G' )$
3.  $\text{GNFA} \rightarrow \text{RegExpr} ( G ) = \text{GNFA} \rightarrow \text{RegExpr} ( G' )$  Plug in
4.  $\text{LANGOF} ( G ) = \text{LANGOF} ( \text{GNFA} \rightarrow \text{RegExpr} ( G ) )$

## Justifications

1. IH
2. Equivalence of Rip/Repair step (prev)
3. Def of  $\text{GNFA} \rightarrow \text{RegExpr}$  (recursive call)
4. From (1), (2), and (3)

Goal

Thm: A Lang is Regular iff Some Reg Expr Describes It

⇒ If a language is regular, then it's described by a regular expr

☑ • Use GNFA → RegExpr to convert GNFA → equiv regular expression!

⇐ If a language is described by a regular expr, then it's regular

☑ • Convert regular expression → equiv NFA! ■

Now: we can use regular expressions to represent regular langs!

So we also have another way to prove things about regular languages!

So a regular language has these equivalent representations:

- DFA
- NFA
- Regular Expression

# *So Far:* How to Prove A Language Is Regular?

Key step, either:

- Construct DFA
- Construct NFA
- Create Regular Expression

Slightly different because of recursive definition

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

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
2.  $\epsilon$ ,
3.  $\emptyset$ ,
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.

# Proof by Induction

To Prove: a ***Statement*** about a recursively defined “thing”  $x$ :

1. Prove: *Statement* for base case of  $x$
2. Prove: *Statement* for recursive case of  $x$ :
  - Assume: **induction hypothesis (IH)**
    - i.e., *Statement* is true for some  $x_{\text{smaller}}$ 
      - E.g., if  $x$  is number, then “smaller” = lesser number
      - ➔• E.g., if  $x$  is regular expression, then “smaller” = ...
  - Prove: *Statement* for  $x$ , using IH (and known definitions, theorems ...)
    - Usually, must show that going from  $x_{\text{smaller}}$  to larger  $x$  is true!

1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
  2.  $\epsilon$ ,
  3.  $\emptyset$ ,
  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.
- 
- Whole reg expr
- "smaller"

# Thm: Reverse is Closed for Regular Langs

Example string:  $\mathbf{abc}^{\mathcal{R}} = \mathbf{cba}$

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  $A^{\mathcal{R}} = \{w^{\mathcal{R}} \mid w \in A\}$

Example language:

$\{\mathbf{a, ab, abc}\}^{\mathcal{R}} = \{\mathbf{a, ba, cba}\}$

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

Proof: by induction on the regular expression of  $A$

# Thm: Reverse is Closed for Regular Langs

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

Proof: by Induction on regular expression of  $A$ : (6 cases)

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.

Need to Prove: if  $A$  is a regular language, described by reg expr  $R_1 \cup R_2$ , then  $A^{\mathcal{R}}$  is regular

IH1: if  $A_1$  is a regular language, described by reg expr  $R_1$ , then  $A_1^{\mathcal{R}}$  is regular

IH1: if  $A_2$  is a regular language, described by reg expr  $R_2$ , then  $A_2^{\mathcal{R}}$  is regular

"smaller"



# Thm: Reverse is Closed for Regular Langs

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

Proof: by Induction on regular expression of  $A$ : (Case # 4)

## Statements

1. Language  $A$  is regular, with reg expr  $R_1 \cup R_2$
2.  $R_1$  and  $R_2$  are regular expressions
3.  $R_1$  and  $R_2$  describe regular langs  $A_1$  and  $A_2$
4. If  $A_1$  is a regular language, then  $A_1^{\mathcal{R}}$  is regular
5. If  $A_2$  is a regular language, then  $A_2^{\mathcal{R}}$  is regular
6.  $A_1^{\mathcal{R}}$  and  $A_2^{\mathcal{R}}$  are regular
7.  $A_1^{\mathcal{R}} \cup A_2^{\mathcal{R}}$  is regular
8.  $A_1^{\mathcal{R}} \cup A_2^{\mathcal{R}} = (A_1 \cup A_2)^{\mathcal{R}}$
9.  $A = A_1 \cup A_2$
10.  $A^{\mathcal{R}}$  is regular

Goal

## Justifications

1. Assumption of IF in IF-THEN
2. Def of Regular Expression
3. Reg Expr  $\Leftrightarrow$  Reg Lang (Prev Thm)
4. IH
5. IH
6. By (3), (4), and (5)
7. Union Closed for Reg Langs
8. Reverse and Union Ops Commute
9. By (1), (2), and (3)
10. By (7), (8), (9)

# Thm: Reverse is Closed for Regular Langs

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

Proof: by Induction on regular expression of  $A$ : (6 cases)

Base cases

- ✓ 1.  $a$  for some  $a$  in the alphabet  $\Sigma$ ,
- ✓ 2.  $\epsilon$ ,
- ✓ 3.  $\emptyset$ ,

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.

Remaining cases will use similar reasoning

# Non-Regular Languages?

- Are there languages that are not regular languages?
- How can we **prove** that a language is not a regular language?

