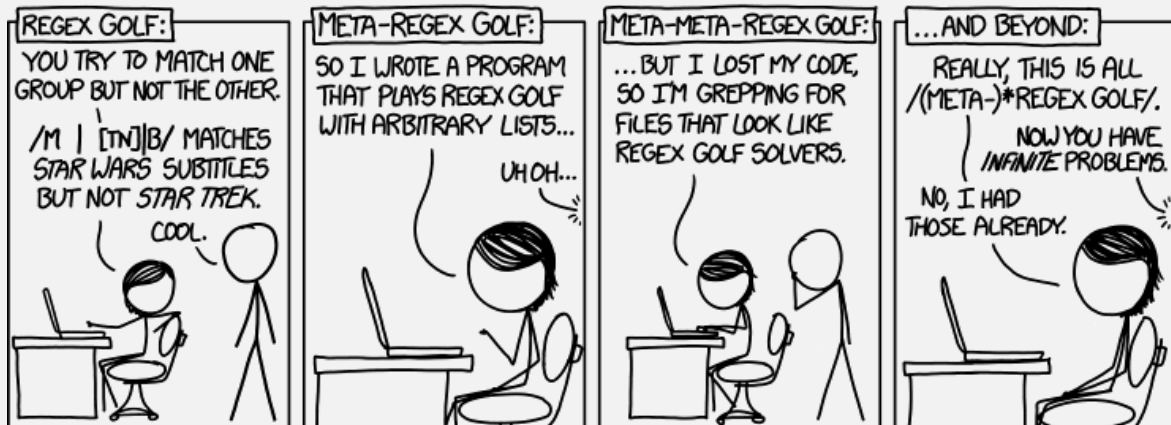


NFA → DFA and Intro to Regular Expressions

Wed, February 10, 2021

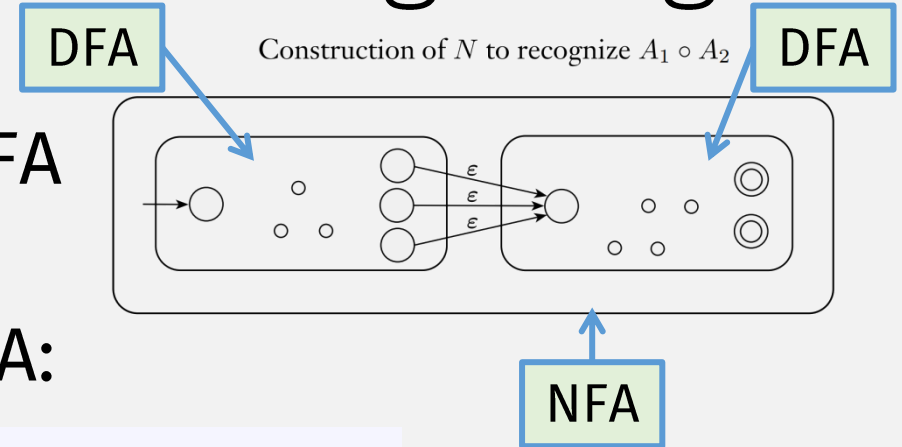


Logistics

- Reminder: no class next Monday 2/15/2021
- Welcome new TA: Benjamin Kwapong
 - See website for additional office hour times
- HW1: solutions posted (to piazza)
 - May use ideas, but not copy (obvi)
- HW2: due Sunday 2/14 11:59pm EST
- HW3: posted, due Sunday 2/21 11:59pm EST
 - Includes a non-coding question
- **Questions?**

Last time: Is Concat Closed for Reg Langs?

- Concatentation of DFAs produces an NFA
- But, regular lang defined using only DFA:



A language is called a *regular language* if some DFA recognizes it.

- To show: *Concatenation is closed for regular languages*, we must prove that NFAs *also* recognize regular languages.
- Specifically:
 - Theorem (1.40): NFAs \Leftrightarrow regular languages

Last time: How to Prove a Theorem: $X \Leftrightarrow Y$

- $X \Leftrightarrow Y$ = “ X if and only if Y ” = X iff Y = $X \Leftrightarrow Y$
- Proof at minimum has 2 parts:
 1. \Rightarrow if X , then Y
 - i.e., assume X , then use it to prove Y
 - “forward” direction
 2. \Leftarrow if Y , then X
 - i.e., assume Y , then use it to prove X
 - “reverse” direction

Proving that NFAs Recognize Reg Langs

- Theorem:
 - A language A is regular **if and only if** some NFA N recognizes it.
- Must prove:
 - \Rightarrow If A is regular, then some NFA N recognizes it
 - Easy
 - We know: if A is regular, then a DFA recognizes it
 - To complete this part of proof: convert DFA to an NFA! (how?)
 - \Leftarrow If an NFA N recognizes A , then A is regular
 - Hard
 - We know: if a DFA recognizes a lang, then it is regular
 - Idea: Convert NFA to DFA

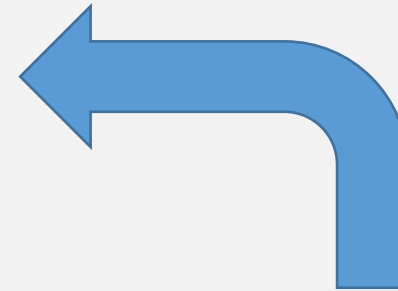
How to convert NFA \rightarrow DFA?

A *finite automaton* is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

1. Q is a finite set called the *states*,
2. Σ is a finite set called the *alphabet*,
3. $\delta: Q \times \Sigma \rightarrow Q$ is the *transition function*,
4. $q_0 \in Q$ is the *start state*, and
5. $F \subseteq Q$ is the *set of accept states*.

Proof idea:

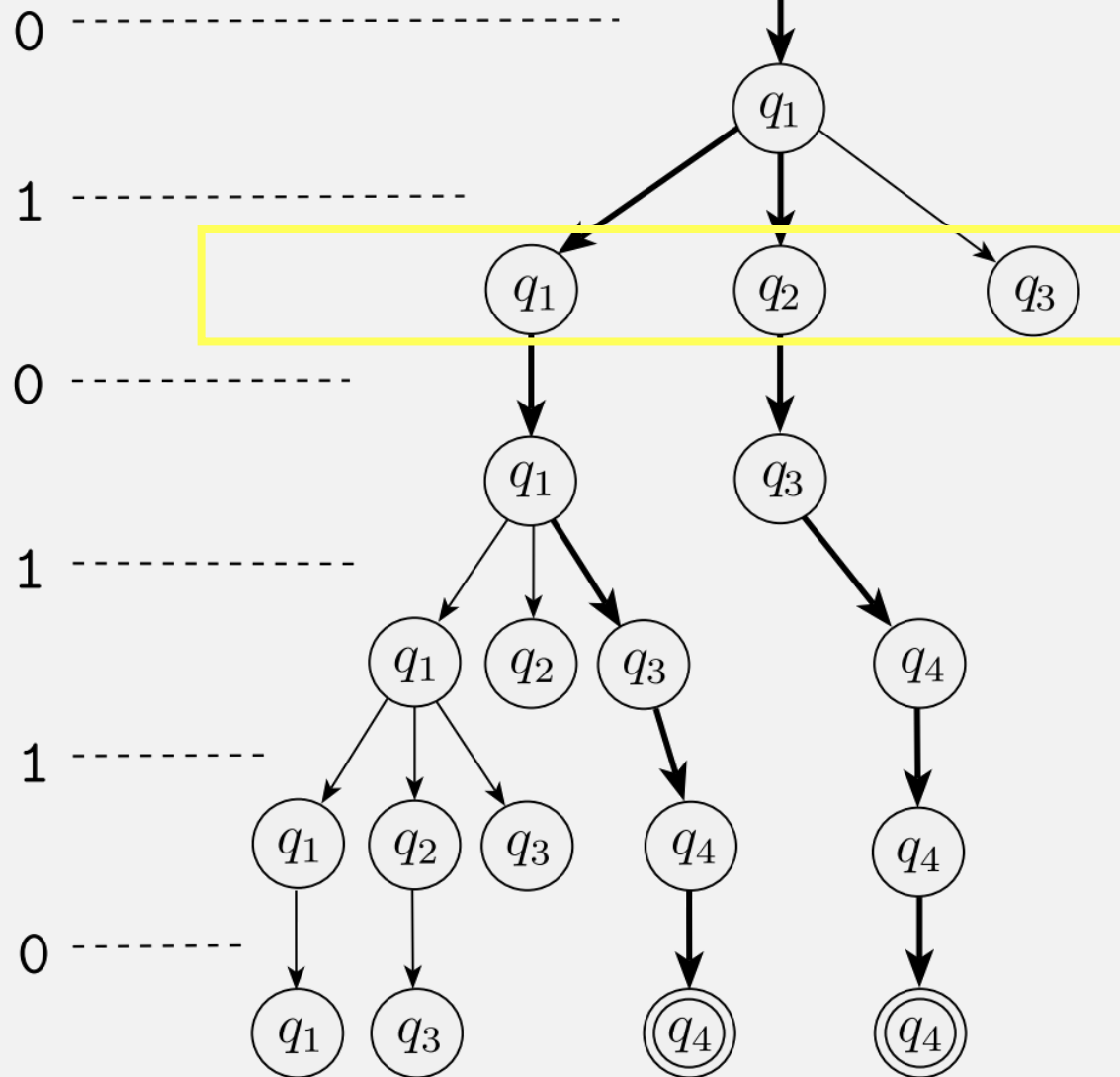
Each “state” of the DFA must be a set of states in the NFA



A *nondeterministic finite automaton* is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, where

1. Q is a finite set of states,
2. Σ is a finite alphabet,
3. $\delta: Q \times \Sigma_\epsilon \rightarrow \mathcal{P}(Q)$ is the transition function,
4. $q_0 \in Q$ is the start state, and
5. $F \subseteq Q$ is the set of accept states.

Symbol read q_1 Start



In a DFA, all these states at each step must be only **one** state

So design a state in the DFA to be a **set of NFA states!**

This is a generalization of the proof strategy from Thm 1.25 (closure of union), where a state = pair of "states"

Convert NFA \rightarrow DFA, Formally

- Let NFA $N = (Q, \Sigma, \delta, q_0, F)$
- An equivalent DFA M has states $Q' = \mathcal{P}(Q)$ (power set of Q)

Example:

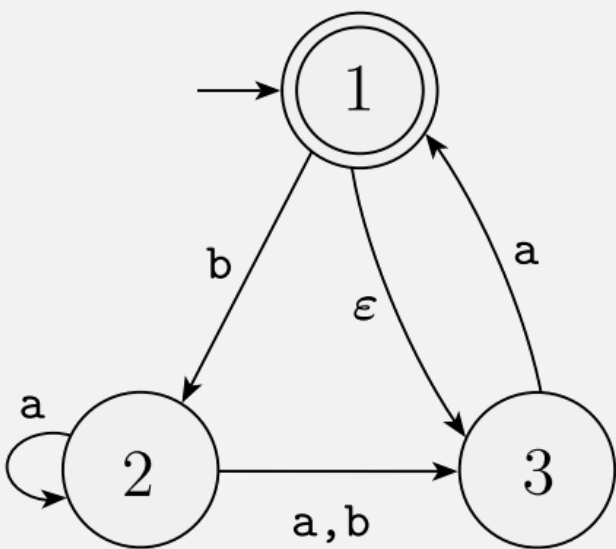


FIGURE 1.42
The NFA N_4

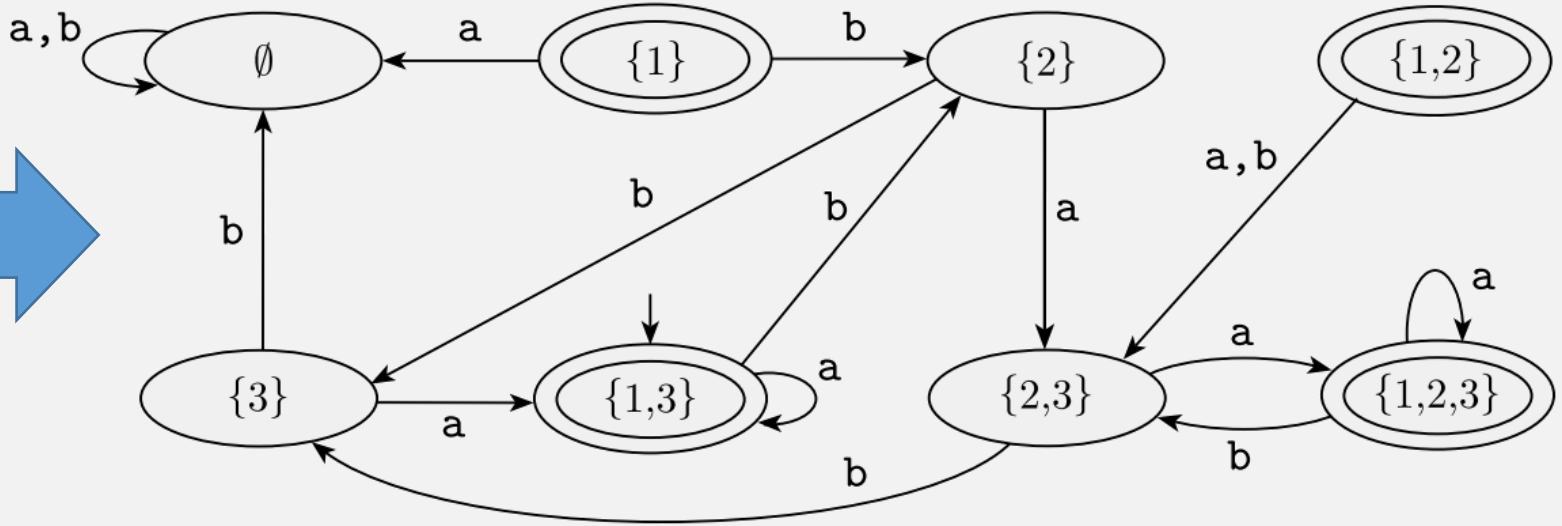
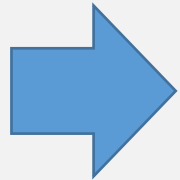


FIGURE 1.43
A DFA D that is equivalent to the NFA N_4

NFA \rightarrow DFA (ignore empty transitions)

- Have: $N = (Q, \Sigma, \delta, q_0, F)$
- Want to: construct a DFA $M = (Q', \Sigma, \delta', q_0', F')$

1. $Q' = \mathcal{P}(Q)$. A state for M is a set of states in N

2. For $R \in Q'$ and $a \in \Sigma$, $R = \text{a state in } M = \text{a set of states in } N$

$$\delta'(R, a) = \bigcup_{r \in R} \delta(r, a)$$

To compute next state for R ,
compute next states of each NFA state r in R ,
then union results into one set

3. $q_0' = \{q_0\}$

4. $F' = \{R \in Q' \mid R \text{ contains an accept state of } N\}$

NFA \rightarrow DFA (with empty transitions)

• Have: $N = (Q, \Sigma, \delta, q_0, F)$

$E(R) = \{q \mid q \text{ can be reached from } R \text{ by traveling along 0 or more } \epsilon \text{ arrows}\}$

• Want to: construct a DFA $M = (Q', \Sigma, \delta', q_0', F')$

1. $Q' = \mathcal{P}(Q)$.

2. For $R \in Q'$ and $a \in \Sigma$,

$$\delta'(R, a) = \bigcup_{r \in R} \cancel{\delta(r, a)} \quad E(\delta(r, a))$$

For each r , do its transition in N , then add states reachable from empty transitions, then union results into one set

3. $q_0' = \cancel{\{q_0\}} \quad E(\{q_0\})$

4. $F' = \{R \in Q' \mid R \text{ contains an accept state of } N\}$

Proving that NFAs Recognize Reg Langs

- Theorem:

- A language A is regular **if and only if** some NFA N recognizes it.

- Must prove:

- \Rightarrow If A is regular, then some NFA N recognizes it

- Easy

- We know: if A is regular, then a DFA recognizes it

- Convert DFA to an NFA

- \Leftarrow If an NFA N recognizes A , then A is regular

- Hard

- We know: if a DFA recognizes a lang, then it is regular

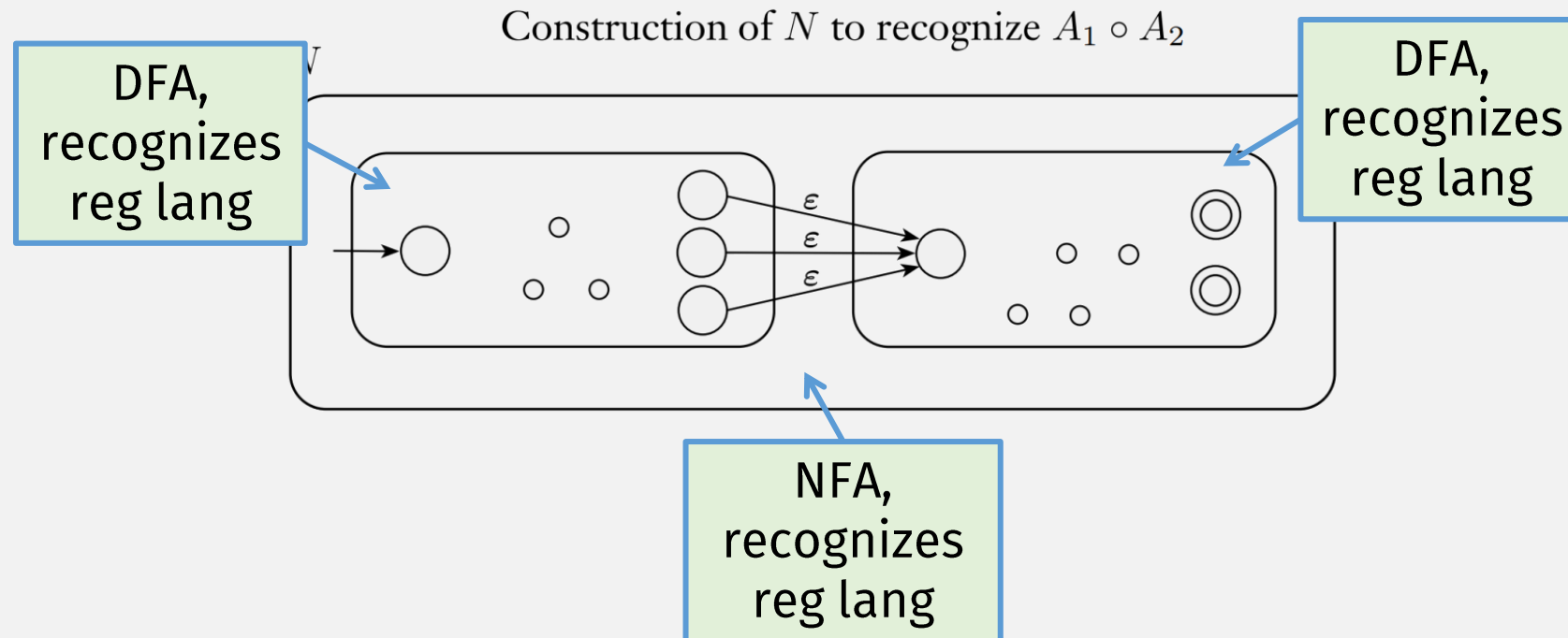
-  Idea: Convert NFA to DFA

- Using NFA \rightarrow DFA algorithm we just created!

■ (Q.E.D.)

So Concatenation is Closed for Reg Langs!

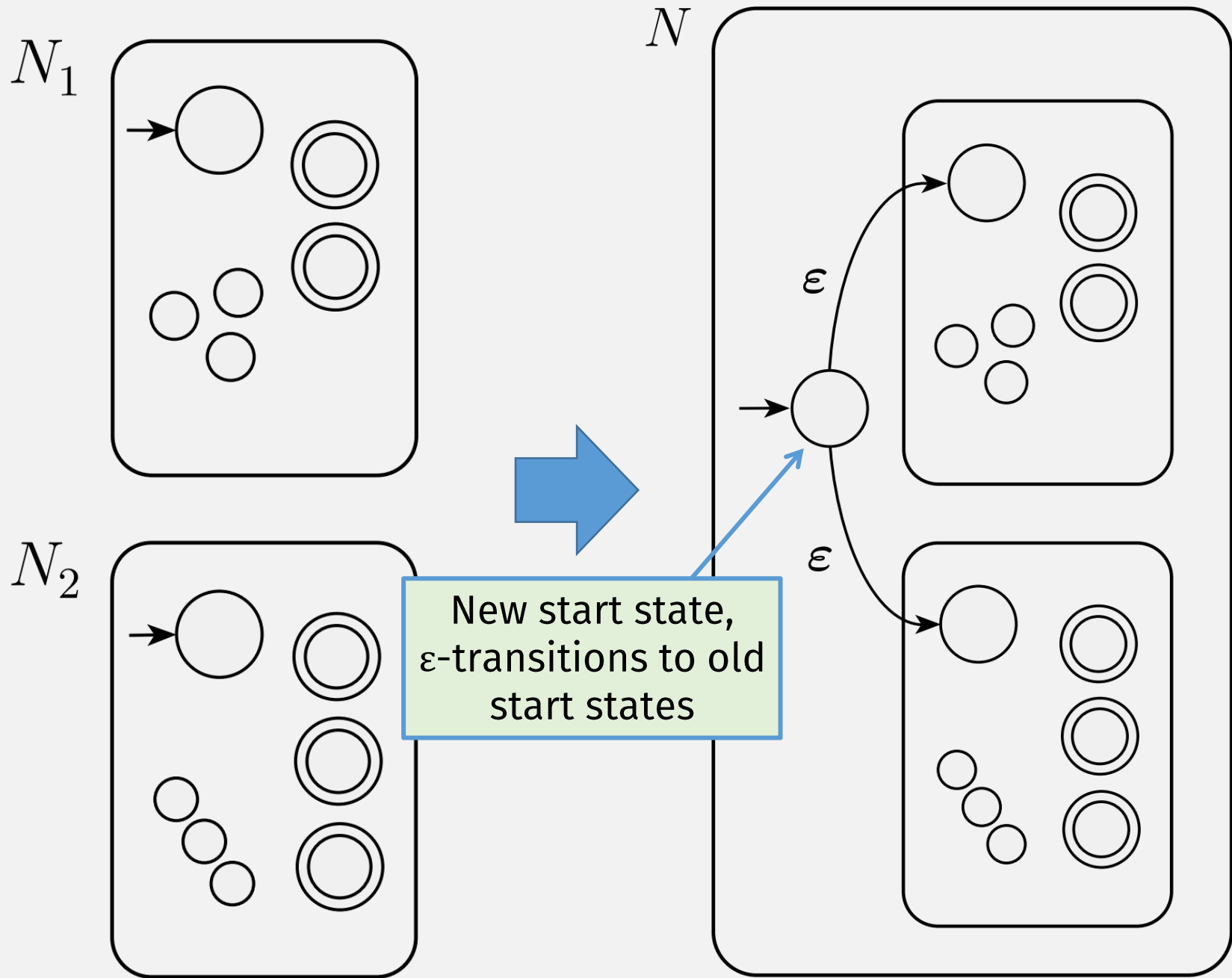
- Concatenation of DFAs produces an NFA



“Regular” Operations

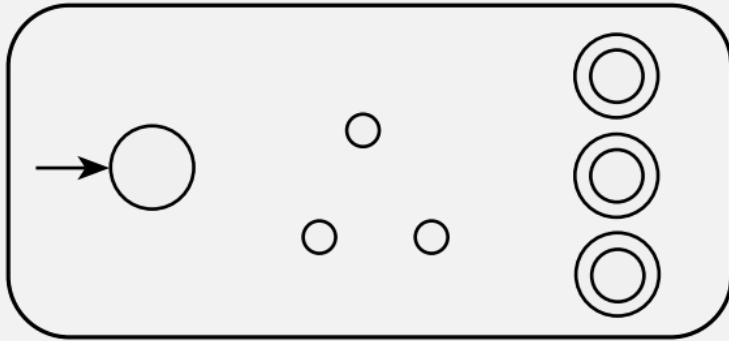
- Regular languages are closed under these operations:
 - Union (already proved with DFAs)
 - Concatenation
 - Kleene Star (repetition, zero or more times)
- Easier to prove closure (by construction) using NFAs

Union

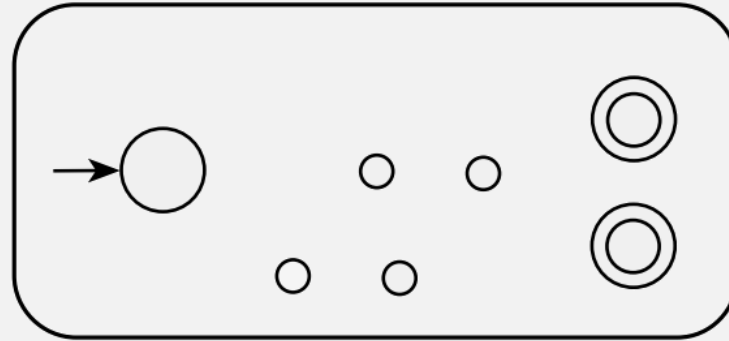


Concat

N_1



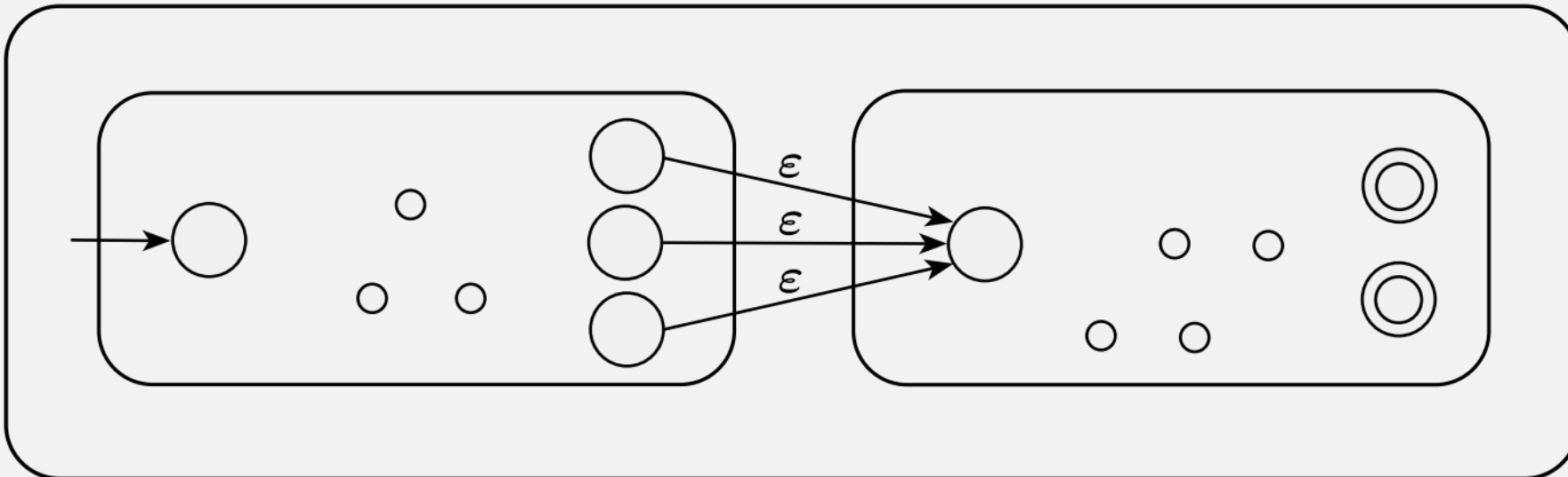
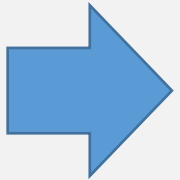
N_2



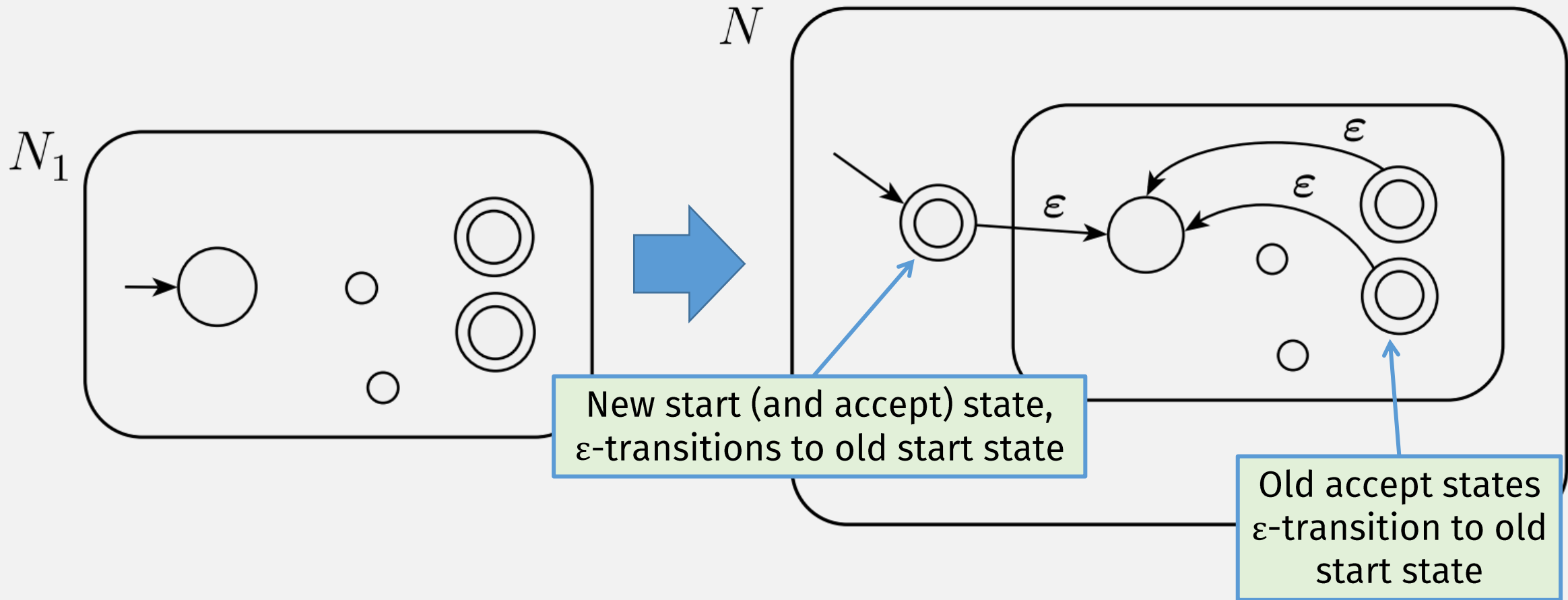
Let N_1 recognize A_1 , and N_2 recognize A_2 .

Construction of N to recognize $A_1 \circ A_2$

N



Kleene Star



Why do we care?

- Union, concat, and kleene star represent all regular languages
- I.e., they define **regular expressions**

DEFINITION 1.52

Say that R is a *regular expression* if R is

1. a for some a in the alphabet Σ ,
2. ϵ ,
3. \emptyset ,

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

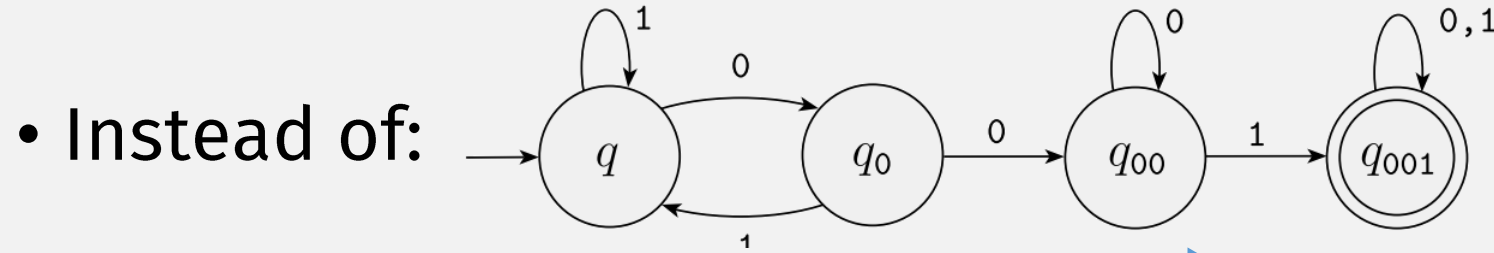
star



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

Poll: Regexes

Ways to Recognize a Regular Language



1. $Q = \{q_1, q_2, q_3\}$,
2. $\Sigma = \{0,1\}$,
3. δ is described as

• Or:

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

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

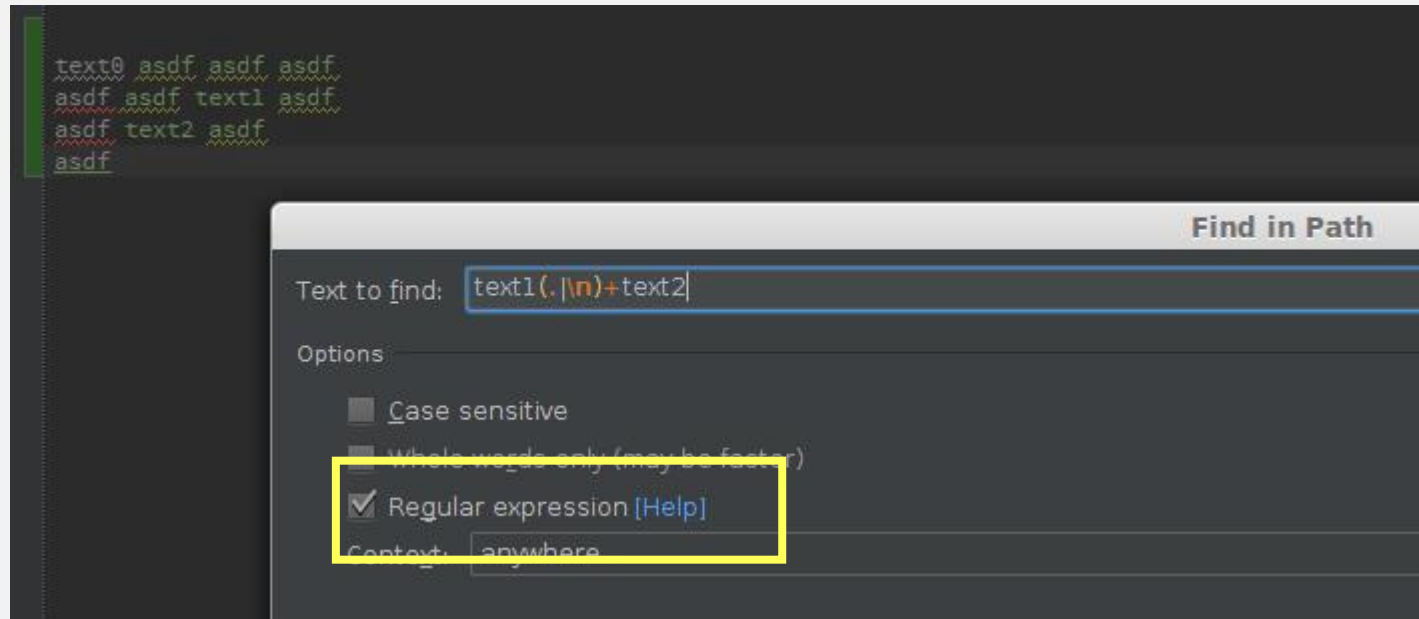
• We can write a regexp: $\Sigma^* 001 \Sigma^*$

These all define a computer (program) that accepts all strings containing 001

Which would you rather write?

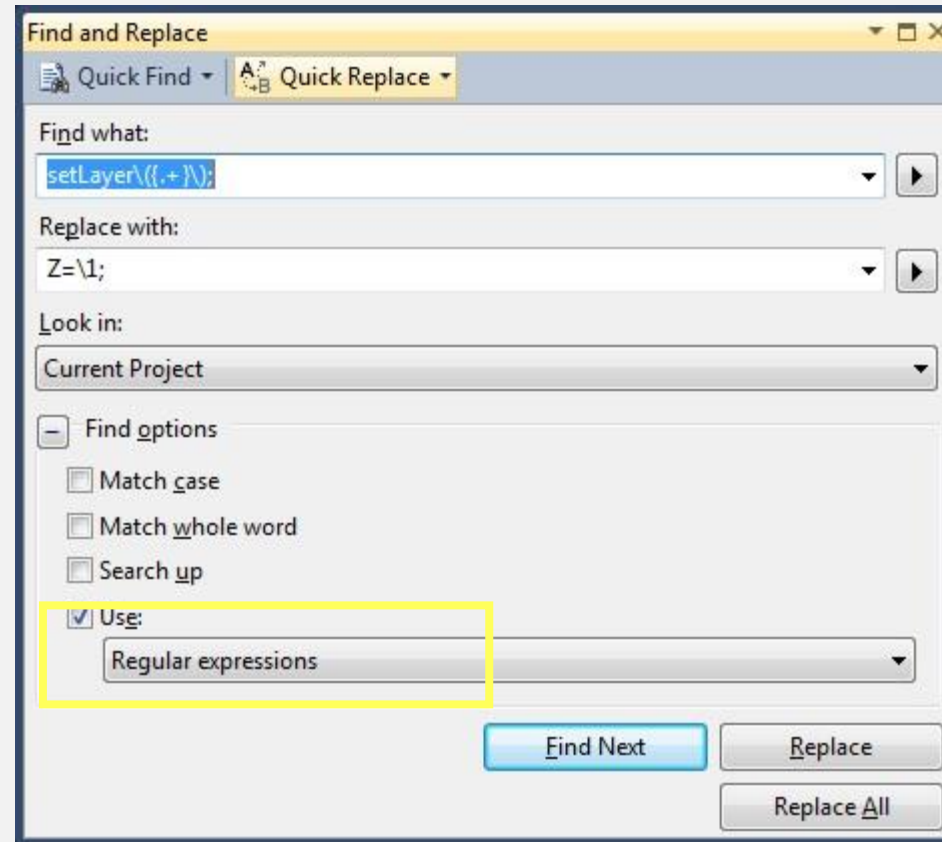
Regular Expressions are Super Useful

- IntelliJ



Regular Expressions are Super Useful

- Visual Studio



Regular Expressions are Super Useful

- Grep (Linux)

```
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.

    In addition, three variant programs egrep, fgrep and rgrep are
    available. egrep is the same as grep -E. fgrep is the same as
    grep -F. rgrep is the same as grep -r. Direct invocation as either
    egrep or fgrep is deprecated, but is provided to allow historical
    applications that rely on them to run unmodified.
```

Regexps supported in every language

- Perl
- Python
- Java
- Every lang!

NAME

perlre - Perl regular expressions

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

Table of Contents

- re — Regular expression operations
 - Regular Expression Syntax
 - Module Contents
 - Regular Expression

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

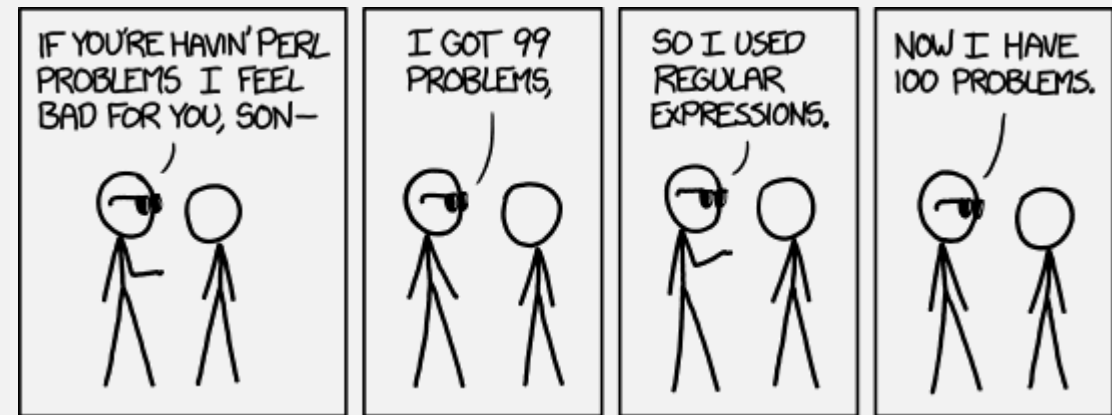
java.lang.Object

java.util.regex.Pattern

Caveat: Regexp are useful, if used correctly



Regexps: potentially **useful** ...



... only if used correctly

Where "used correctly" =
only use it to recognize
regular languages

(To do this, you must know what is,
and is not, a regular language!)

Someone Who Did Not T

RegEx match open tags except XHTML self-con

Asked 10 years, 10 months ago Active 1 month ago Viewed 2.9m times

I need to match all of these opening tags:

1553

```
<p>
<a href="foo">
```

But not these:

6572

You can't parse [X]HTML with regex. Because HTML can't be parse
Regex is not a tool that can be used to correctly parse HTML. As I h
HTML-and-regex questions here so many times before, the use of re
allow you to consume HTML. Regular expressions are a tool that is
sophisticated to understand the constructs employed by HTML. HTM
regular language and hence cannot be parsed by regular expressior
queries are not equipped to break down HTML into its meaningful p
times but it is not getting to me. Even enhanced irregular regular exp
used by Perl are not up to the task of parsing HTML. You will never i

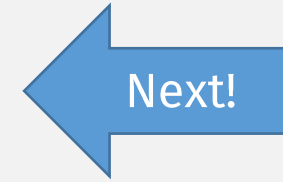
4414

HTML is a language of sufficient complexity that it cannot be parsed by regular expressions. Even Jon Skeet cannot parse HTML using regular expressions. Every time you attempt to parse HTML with regular expressions, the unholy child weeps the blood of virgins, and Russian hackers pwn your webapp. Parsing HTML with regex summons tainted souls into the realm of the living. HTML and regex go together like love, marriage, and ritual infanticide. The <center> cannot hold it is too late. The force of regex and HTML together in the same conceptual space will destroy your mind like so much watery putty. If you parse HTML with regex you are giving in to Them and their blasphemous ways which doom us all to inhuman toil for the One whose Name cannot be expressed in the Basic Multilingual Plane, he comes. HTML-plus-regex will liquify the nerves of the sentient whilst you observe, your psyche withering in the onslaught of horror. Regēx-based HTML parsers are the cancer that is killing StackOverflow *it is too late it is too late we cannot be saved* the transgression of a child ensures regex will consume all living tissue (except for HTML which it cannot, as previously prophesied) *dear lord help us how can anyone survive this scourge* using regex to parse HTML has doomed humanity to an eternity of dread torture and security holes *using regex* as a tool to process HTML establishes a breach *between this world* and the dread realm of corrupt entities (like SGML entities, but *more corrupt*) a mere glimpse of the world of regex parsers for HTML will instantly transport a programmer's consciousness into a world of ceaseless screaming, he comes, the pestilent slithy regex-infection will devour your HTML parser, application and existence for all time like Visual Basic only worse he comes he comes do not fight he comes, his unholy radiancé destroying all enlightenment, HTML tags *leaking from your eyes like liquid pain*, the song of regular expression parsing will extinguish the voices of mortal man from the sphere I can see it can you see if it is beautiful the final snuffing of the lies of Man ALL IS LOST ALL IS LOST the pony he comes he comes he comes the anchor permeates all MY FACE MY FACE oh god no NO! NO! NO! NO! stop the angles are not real ZALGO IS TONY THE PONY, HE COMES

Have you tried using an XML parser instead?

Big Picture Road Map

- We ultimately want to prove:
 - Regular Languages \Leftrightarrow Regular Expressions
- First, we need to show these operations are closed for reglang:
 - Union (**done!**)
 - Concatentation (**done!**)
 - Kleene star (**done!**)



Thm: A lang is regular iff some regexp describes it

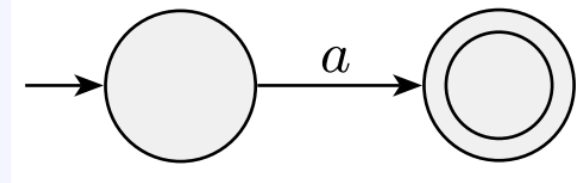
- \Rightarrow If a language is regular, it is described by a regexp
- \Leftarrow If a language is described by a regexp, it is regular
 - Easy!
 - Construct the NFA!
 - See Lemma 1.55

Regexp \rightarrow NFA

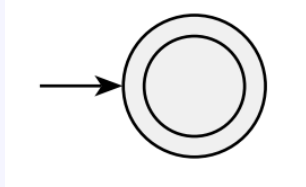
DEFINITION 1.52

Say that R is a *regular expression* if R is

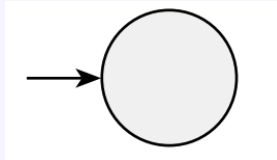
1. a for some a in the alphabet Σ ,



2. ϵ ,



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,

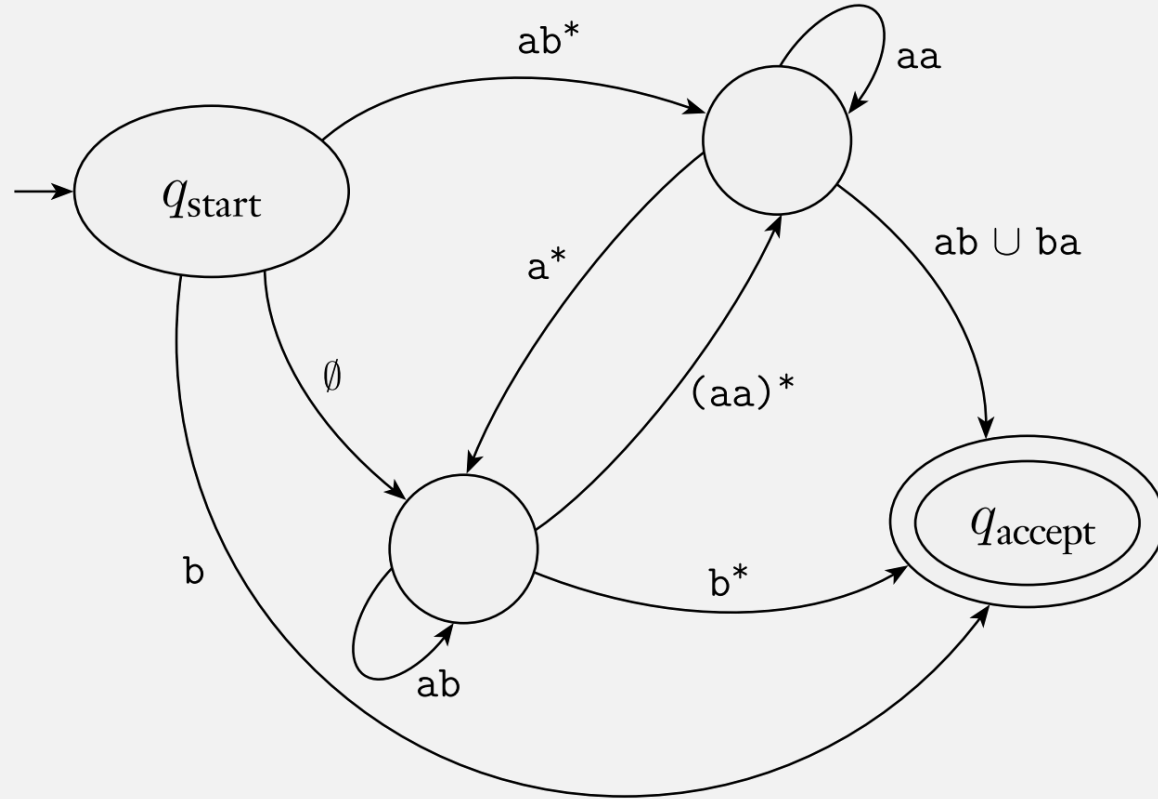
6. (R_1^*) , where R_1 is a regular expression.

Constructions from before!

Thm: A lang is regular iff some regexp describes it

- \Rightarrow If a language is regular, it is described by a regexp
 - Hard!
 - Need something new: a GNFA
- \Leftarrow If a language is described by a regexp, it is regular
 - Easy!
 - Construct the NFA! (**Done**)

GNFA = NFA with regexp transitions



- To convert to regexp, keep “ripping out” states until only 2 are left

Check-in Quiz 2/10

On gradescope