

FORMAL LANGUAGES, AUTOMATA AND COMPUTATION

REGULAR EXPRESSIONS

Carnegie Mellon University in Qatar

SUMMARY

- Nondeterminism
 - Clone the FA at choice points

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify

SUMMARY

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify
- Nondeterministic FA

SUMMARY

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify
- Nondeterministic FA
 - Multiple transitions from a state with the same input symbol

SUMMARY

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify
- Nondeterministic FA
 - Multiple transitions from a state with the same input symbol
 - ϵ -transitions

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify
- Nondeterministic FA
 - Multiple transitions from a state with the same input symbol
 - ϵ -transitions
- NFAs are equivalent to DFAs

- Nondeterminism
 - Clone the FA at choice points
 - Guess and verify
- Nondeterministic FA
 - Multiple transitions from a state with the same input symbol
 - ϵ -transitions
- NFAs are equivalent to DFAs
 - Determinization procedure builds a DFA with up to 2^k states for an NFA with k states.

CLOSURE THEOREMS

THEOREM

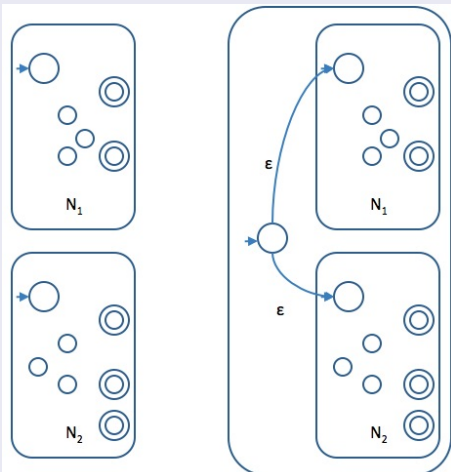
*The class of regular languages is closed under the **union** operation.*

CLOSURE THEOREMS

THEOREM

*The class of regular languages is closed under the **union** operation.*

PROOF IDEA BASED ON NFAS



CLOSURE THEOREMS

THEOREM

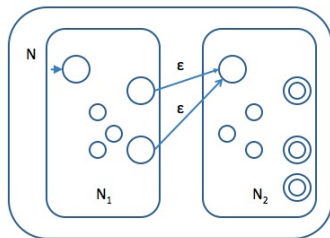
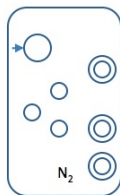
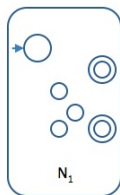
*The class of regular languages is closed under the **concatenation** operation.*

CLOSURE THEOREMS

THEOREM

*The class of regular languages is closed under the **concatenation** operation.*

PROOF IDEA BASED ON NFAS



CLOSURE THEOREMS

THEOREM

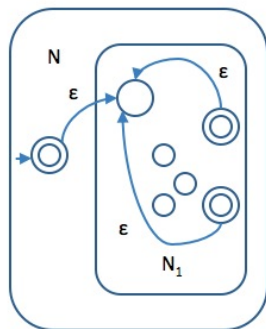
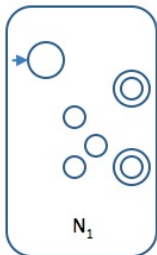
*The class of regular languages is closed under the **star** operation.*

CLOSURE THEOREMS

THEOREM

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

PROOF IDEA BASED ON NFAS



REGULAR EXPRESSIONS

- DFAs are **finite** descriptions of (finite or infinite) sets of strings

REGULAR EXPRESSIONS

- DFAs are **finite** descriptions of (finite or infinite) sets of strings
 - Finite number of symbols, states, transitions

REGULAR EXPRESSIONS

- DFAs are **finite** descriptions of (finite or infinite) sets of strings
 - Finite number of symbols, states, transitions
- **Regular Expressions** provide an algebraic expression framework to describe the same class of strings

REGULAR EXPRESSIONS

- DFAs are **finite** descriptions of (finite or infinite) sets of strings
 - Finite number of symbols, states, transitions
- **Regular Expressions** provide an algebraic expression framework to describe the same class of strings
- Thus, DFAs and Regular Expressions are equivalent.

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression Regular Set

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
ϕ	$\{\}$

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $\mathbf{a} \in \Sigma$	$\{\mathbf{a}\}$

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
---------------------------	--------------------

ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{\mathbf{a}\}$
ϵ	$\{\epsilon\}$

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 R_2)$	$L(R_1)L(R_2)$

REGULAR EXPRESSIONS

- For every regular expression, there is a corresponding regular set or language
- R, R_1, R_2 are regular expressions; $L(R)$ denotes the corresponding regular set

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 R_2)$	$L(R_1)L(R_2)$
(R^*)	$L(R)^*$

REGULAR EXPRESSIONS— MORE SYNTAX

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
(R^*)	$L(R)^*$

- Also some books use $R_1 + R_2$ to denote union.

REGULAR EXPRESSIONS— MORE SYNTAX

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
(R^*)	$L(R)^*$

- Also some books use $R_1 + R_2$ to denote union.
- In (\dots) , the parenthesis can be deleted

REGULAR EXPRESSIONS— MORE SYNTAX

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
(R^*)	$L(R)^*$

- Also some books use $R_1 + R_2$ to denote union.
- In (\dots) , the parenthesis can be deleted
 - In which case, interpretation is done in the **precedence order**: **star**, **concatenation** and then **union**.

REGULAR EXPRESSIONS— MORE SYNTAX

Regular Expression	Regular Set
ϕ	$\{\}$
\mathbf{a} for $a \in \Sigma$	$\{a\}$
ϵ	$\{\epsilon\}$
$(R_1 \cup R_2)$	$L(R_1) \cup L(R_2)$
$(R_1 \circ R_2)$	$L(R_1) \circ L(R_2)$
(R^*)	$L(R)^*$

- Also some books use $R_1 + R_2$ to denote union.
- In (\dots) , the parenthesis can be deleted
 - In which case, interpretation is done in the **precedence order**: **star**, **concatenation** and then **union**.
- $R^+ = RR^*$ and R^k for k -fold concatenation are useful shorthands.

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

→

Regular Language

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

→

Regular Language

$\{\omega \mid \omega \text{ contains a single } 1\}$

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

Regular Language

$\rightarrow \{\omega \mid \omega \text{ contains a single } 1\}$

\rightarrow

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

Regular Language

→ $\{\omega \mid \omega \text{ contains a single } 1\}$

→ $\{\omega \mid \omega \text{ has at least one } 1\}$

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

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

Regular Language

→ $\{\omega \mid \omega \text{ contains a single } 1\}$

→ $\{\omega \mid \omega \text{ has at least one } 1\}$

→

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

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

Regular Language

→ $\{\omega \mid \omega \text{ contains a single } 1\}$

→ $\{\omega \mid \omega \text{ has at least one } 1\}$

→ $\{\omega \mid \omega \text{ starts and ends with the same symbol}\}$

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

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

$(0^*10^*1)^*0^*$

Regular Language

→ $\{\omega \mid \omega \text{ contains a single } 1\}$

→ $\{\omega \mid \omega \text{ has at least one } 1\}$

→ $\{\omega \mid \omega \text{ starts and ends}$
with the same symbol}

→

REGULAR EXPRESSION EXAMPLES

Regular Expression

0^*10^*

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

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

$(0^*10^*1)^*0^*$

Regular Language

→ $\{\omega \mid \omega \text{ contains a single } 1\}$

→ $\{\omega \mid \omega \text{ has at least one } 1\}$

→ $\{\omega \mid \omega \text{ starts and ends}$
with the same symbol}

→ $\{\omega \mid n_1(\omega) \text{ is even}\}$

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00 (0 \cup 1)^*$

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00(0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00(0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* 1(0 \cup 1)(0 \cup 1)(0 \cup 1)$

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00 (0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* 1 (0 \cup 1) (0 \cup 1) (0 \cup 1)$
- All strings with **no** pair of consecutive 0s

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00 (0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* 1 (0 \cup 1) (0 \cup 1) (0 \cup 1)$
- All strings with **no** pair of consecutive 0s
 - $(1^* 0 1 1^*)^* (0 \cup \epsilon) \cup 1^*$

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00 (0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* 1 (0 \cup 1) (0 \cup 1) (0 \cup 1)$
- All strings with **no** pair of consecutive 0s
 - $(1^* 011^*)^* (0 \cup \epsilon) \cup 1^*$
 - Strings consist of repetitions of 1 or 01 or two boundary cases: $(1 \cup 01)^* (0 \cup \epsilon)$

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* 00(0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* 1(0 \cup 1)(0 \cup 1)(0 \cup 1)$
- All strings with **no** pair of consecutive 0s
 - $(1^* 011^*)^*(0 \cup \epsilon) \cup 1^*$
 - Strings consist of repetitions of 1 or 01 or two boundary cases: $(1 \cup 01)^*(0 \cup \epsilon)$
- All strings that **do not end** in 01.

WRITING REGULAR EXPRESSIONS

- All strings with **at least** one pair of consecutive 0s
 - $(0 \cup 1)^* \mathbf{00} (0 \cup 1)^*$
- All strings such that fourth symbol from the end is a 1
 - $(0 \cup 1)^* \mathbf{1} (0 \cup 1) (0 \cup 1) (0 \cup 1)$
- All strings with **no** pair of consecutive 0s
 - $(1^* \mathbf{011}^*)^* (0 \cup \epsilon) \cup 1^*$
 - Strings consist of repetitions of 1 or 01 or two boundary cases: $(1 \cup \mathbf{01})^* (0 \cup \epsilon)$
- All strings that **do not end** in 01.
 - $(0 \cup 1)^* (\mathbf{00} \cup \mathbf{10} \cup \mathbf{11}) \cup \mathbf{0} \cup \mathbf{1} \cup \epsilon$

WRITING REGULAR EXPRESSIONS

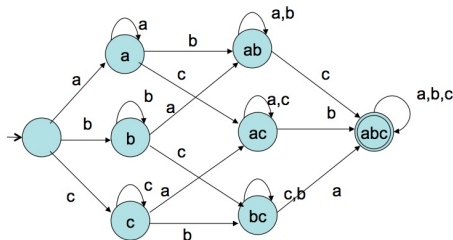
- All strings over $\Sigma = \{a, b, c\}$ that contain every symbol at least once.

WRITING REGULAR EXPRESSIONS

- All strings over $\Sigma = \{a, b, c\}$ that contain every symbol at least once.
- $(a \cup b \cup c)^* a (a \cup b \cup c)^* b (a \cup b \cup c)^* c (a \cup b \cup c)^* \cup$
 $(a \cup b \cup c)^* a (a \cup b \cup c)^* c (a \cup b \cup c)^* b (a \cup b \cup c)^* \cup$
 $(a \cup b \cup c)^* b (a \cup b \cup c)^* a (a \cup b \cup c)^* c (a \cup b \cup c)^* \cup$
 $(a \cup b \cup c)^* b (a \cup b \cup c)^* c (a \cup b \cup c)^* a (a \cup b \cup c)^* \cup$
 $(a \cup b \cup c)^* c (a \cup b \cup c)^* a (a \cup b \cup c)^* b (a \cup b \cup c)^* \cup$
 $(a \cup b \cup c)^* c (a \cup b \cup c)^* b (a \cup b \cup c)^* a (a \cup b \cup c)^*$

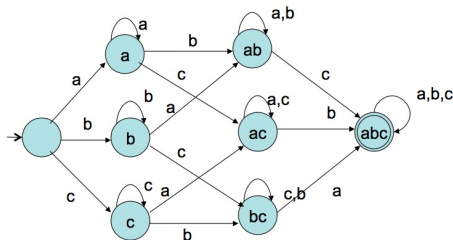
WRITING REGULAR EXPRESSIONS

- All strings over $\Sigma = \{a, b, c\}$ that contain every symbol at least once.



WRITING REGULAR EXPRESSIONS

- All strings over $\Sigma = \{a, b, c\}$ that contain every symbol at least once.



- DFA**s and **RE**s may need different ways of looking at the problem.
 - For the DFA, you count symbols
 - For the RE, you enumerate all possible patterns

RE IDENTITIES

- $\mathbf{R} \cup \phi = \mathbf{R}$

RE IDENTITIES

- $\mathbf{R} \cup \phi = \mathbf{R}$
- $\mathbf{R}\epsilon = \epsilon\mathbf{R} = \mathbf{R}$

RE IDENTITIES

- $\mathbf{R} \cup \phi = \mathbf{R}$
- $\mathbf{R}\epsilon = \epsilon\mathbf{R} = \mathbf{R}$
- $\phi^* = \epsilon$

RE IDENTITIES

- $\mathbf{R} \cup \phi = \mathbf{R}$
- $\mathbf{R}\epsilon = \epsilon\mathbf{R} = \mathbf{R}$
- $\phi^* = \epsilon$
- Note that we do not have explicit operators for intersection or complementation!

DIGRESSION: RES IN REAL LIFE

- Linux/Unix Shell, Perl, Awk, Python all have built in regular expression support for pattern matching functionality
- **See** `http://www.wdvl.com/Authoring/Languages/Perl/PerlfortheWeb/perlintro2_table1.html`

DIGRESSION: RES IN REAL LIFE

- Linux/Unix Shell, Perl, Awk, Python all have built in regular expression support for pattern matching functionality
- See http://www.wdvl.com/Authoring/Languages/Perl/PerlfortheWeb/perlintro2_table1.html
- Mostly some syntactic extensions/changes to basic regular expressions with some additional functionality for remembering matches

DIGRESSION: RES IN REAL LIFE

- Linux/Unix Shell, Perl, Awk, Python all have built in regular expression support for pattern matching functionality
- See http://www.wdvl.com/Authoring/Languages/Perl/PerlfortheWeb/perlintro2_table1.html
- Mostly some syntactic extensions/changes to basic regular expressions with some additional functionality for remembering matches
- Substring matches in a string!
- Search for and download *Regex Coach* to learn and experiment with regular expression matching

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA- THE *if* PART

If a language is described by a regular expression, then it is regular

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA- THE *if* PART

If a language is described by a regular expression, then it is regular

PROOF IDEA

Inductively convert a given regular expression to an NFA.

CONVERTING RES TO NFAS: BASIS CASES

Regular Expression Corresponding NFA

ϕ



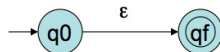
CONVERTING RES TO NFAS: BASIS CASES

Regular Expression Corresponding NFA

ϕ



ϵ



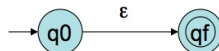
CONVERTING RES TO NFAS: BASIS CASES

Regular Expression Corresponding NFA

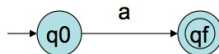
ϕ



ϵ



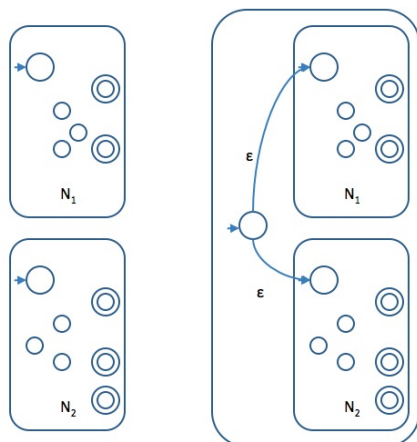
a for $a \in \Sigma$



CONVERTING RES TO NFAS

Union

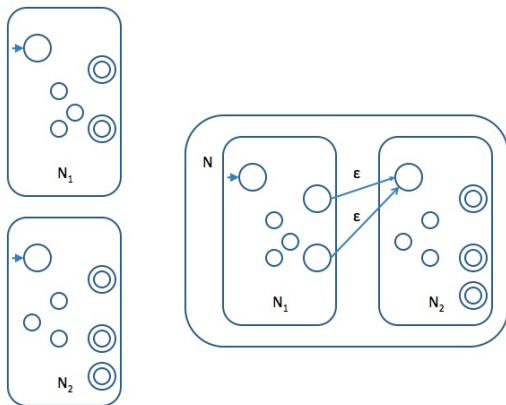
- Let N_1 and N_2 be NFAs for R_1 and R_2 respectively. Then the NFA for $R_1 \cup R_2$ is



CONVERTING RES TO NFAS

Concatenation

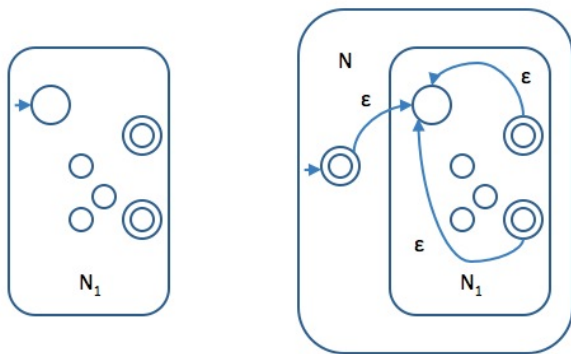
- Let N_1 and N_2 be NFAs for R_1 and R_2 respectively. Then the NFA for R_1R_2 is



CONVERTING RES TO NFAS: STAR

Star

- Let N be NFAs for R . Then the NFA for R^* is



RE TO NFA CONVERSION EXAMPLE

- Let's convert $(\mathbf{a} \cup \mathbf{b})^* \mathbf{aba}$ to an NFA.

RE TO NFA TO DFA

- Regular Expression \rightarrow NFA (possibly with ϵ -transitions)

RE TO NFA TO DFA

- Regular Expression \rightarrow NFA (possibly with ϵ -transitions)
- NFA \rightarrow DFA via determinization

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA – THE *only if* PART

If a language is regular then it is described by a regular expression

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA – THE *only if* PART

If a language is regular then it is described by a regular expression

PROOF IDEA

- **Generalized transitions:** label transitions with regular expressions

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA – THE *only if* PART

If a language is regular then it is described by a regular expression

PROOF IDEA

- **Generalized transitions:** label transitions with regular expressions
- **Generalized NFAs** (GNFA)

EQUIVALENCE WITH FINITE AUTOMATA

THEOREM

A language is regular if and only if some regular expression describes it.

LEMMA – THE *only if* PART

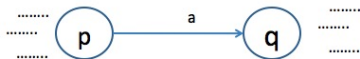
If a language is regular then it is described by a regular expression

PROOF IDEA

- **Generalized transitions:** label transitions with regular expressions
- **Generalized NFAs** (GNFA)
- Iteratively eliminate states of the GNFA one by one, until only two states and a single generalized transition is left.

GENERALIZED TRANSITIONS

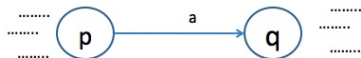
- DFAs have single symbols as transition labels



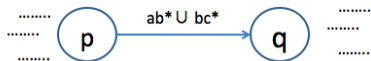
- If you are in state p and the next input symbol matches a , go to state q

GENERALIZED TRANSITIONS

- DFAs have single symbols as transition labels

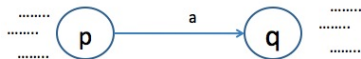


- If you are in state p and the next input symbol matches a , go to state q
- Now consider

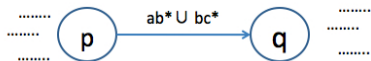


GENERALIZED TRANSITIONS

- DFAs have single symbols as transition labels



- If you are in state p and the next input symbol matches a , go to state q
- Now consider



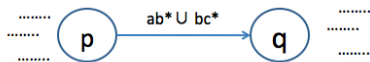
- If you are in state p and **a prefix of the remaining input** matches the regular expression $ab^* \cup bc^*$ then go to state q

GENERALIZED TRANSITIONS AND NFA

- A generalized transition is a transition whose label is a regular expression

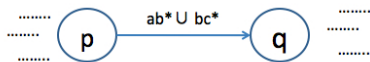
GENERALIZED TRANSITIONS AND NFA

- A generalized transition is a transition whose label is a regular expression



GENERALIZED TRANSITIONS AND NFA

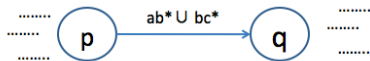
- A generalized transition is a transition whose label is a regular expression



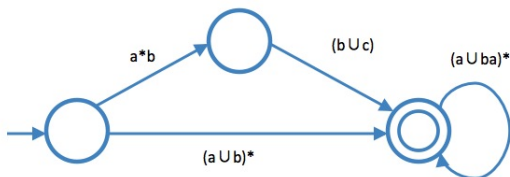
- A Generalized NFA is an NFA with generalized transitions.

GENERALIZED TRANSITIONS AND NFA

- A generalized transition is a transition whose label is a regular expression



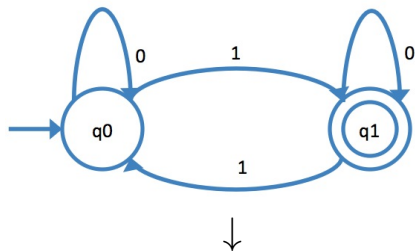
- A Generalized NFA is an NFA with generalized transitions.



- In fact, all standard DFA transitions are generalized transitions with regular expressions of a single symbol!

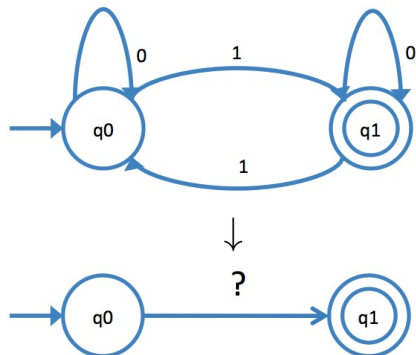
GENERALIZED TRANSITIONS

- Consider the 2-state DFA



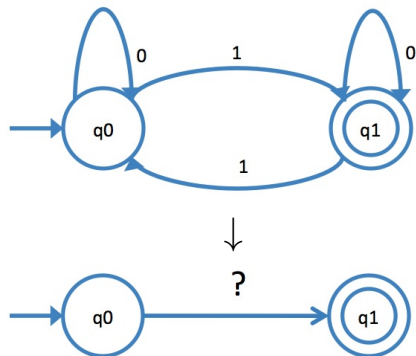
GENERALIZED TRANSITIONS

- Consider the 2-state DFA



GENERALIZED TRANSITIONS

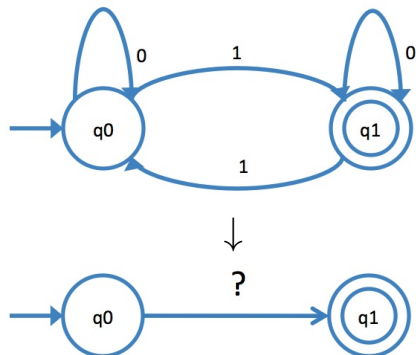
- Consider the 2-state DFA



- 0^*1 takes the DFA from state q_0 to q_1

GENERALIZED TRANSITIONS

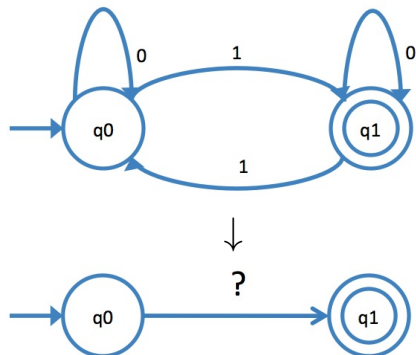
- Consider the 2-state DFA



- 0^*1 takes the DFA from state q_0 to q_1
- $(0 \cup 10^*1)^*$ takes the machine from q_1 back to q_1

GENERALIZED TRANSITIONS

- Consider the 2-state DFA



- 0^*1 takes the DFA from state q_0 to q_1
- $(0 \cup 10^*1)^*$ takes the machine from q_1 back to q_1
- So $? = 0^*1(0 \cup 10^*1)^*$ represents all strings that take the DFA from state q_0 to q_1

GENERALIZED NFAs

- Take any NFA and transform it into a GNFA
 - with only two states: one start and one accept

GENERALIZED NFAs

- Take any NFA and transform it into a GNFA
 - with only two states: one start and one accept
 - with one generalized transition

GENERALIZED NFAs

- Take any NFA and transform it into a GNFA
 - with only two states: one start and one accept
 - with one generalized transition
- then we can “read” the regular expression from the label of the generalized transition (as in the example above)