

# INF2080

## 2. Regular Expressions and Nonregular languages

Daniel Lupp

Universitetet i Oslo

25th January 2018

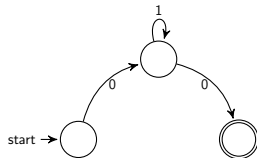


Department of  
Informatics

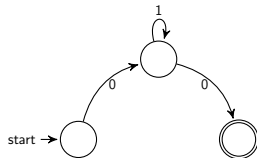


University of  
Oslo

- Deterministic finite automata (DFA)

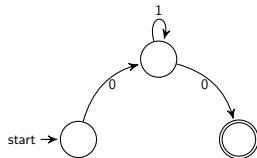


- Deterministic finite automata (DFA)

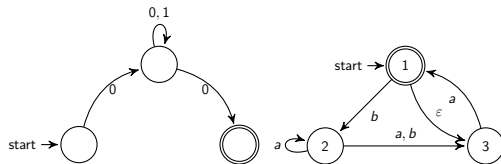


- Regular languages are those languages accepted by DFA's

- Deterministic finite automata (DFA)

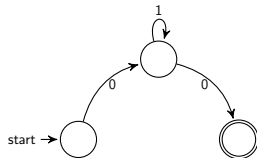


- Regular languages are those languages accepted by DFA's
- Nondeterministic automata (NFA)

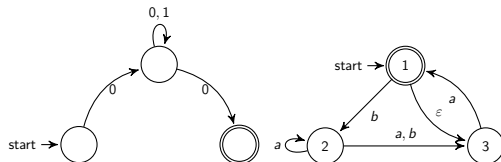


## Last week

- Deterministic finite automata (DFA)



- Regular languages are those languages accepted by DFA's
- Nondeterministic automata (NFA)



- DFA  $\leftrightarrow$  NFA

# Regular Expressions

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,



## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

# Regular Expressions

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

→ Regular expressions represent languages!

## Regular Expressions - Examples

What languages do the following regular expressions (RE) represent?

- $0^*$

## Regular Expressions - Examples

What languages do the following regular expressions (RE) represent?

- $0^*$
- $10^*1$

## Regular Expressions - Examples

What languages do the following regular expressions (RE) represent?

- $0^*$
- $10^*1$
- $(1(0 \cup 1)^*1) \cup (0(0 \cup 1)^*0) \cup 0 \cup 1$

What is the connection between RE and DFA/NFA?

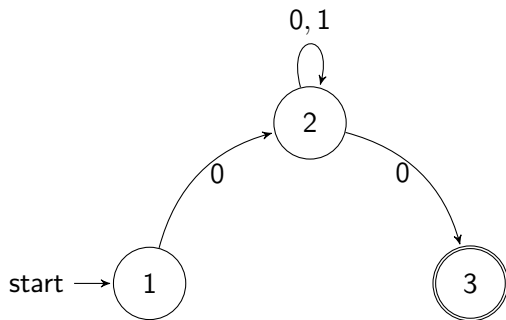


What is the connection between RE and DFA/NFA?

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

What is the connection between RE and DFA/NFA?

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



What is the connection between RE and DFA/NFA?

- Can all RE be represented using DFA/NFA?
- Can all DFA/NFA be described by RE?

# Regular Expressions and Automata

What is the connection between RE and DFA/NFA?

- Can all RE be represented using DFA/NFA?
- Can all DFA/NFA be described by RE?

Yes!

# Regular Expressions and Automata

## Proposition

*Every language described by an RE is regular.*

Proof based on inductive definition of RE!

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

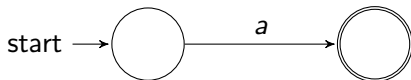
- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

if  $R = a$  for  $a \in \Sigma$ , then  $L(R) = \{a\}$  is accepted by

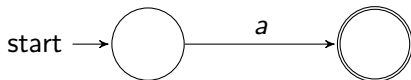


## Definition (Regular Expression)

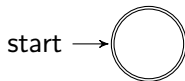
Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

if  $R = a$  for  $a \in \Sigma$ , then  $L(R) = \{a\}$  is accepted by



If  $R = \varepsilon$ , then  $L(R) = \{\varepsilon\}$  is accepted by



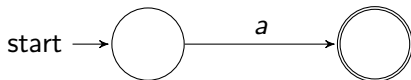


## Definition (Regular Expression)

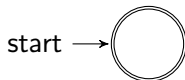
Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

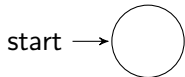
if  $R = a$  for  $a \in \Sigma$ , then  $L(R) = \{a\}$  is accepted by



If  $R = \varepsilon$ , then  $L(R) = \{\varepsilon\}$  is accepted by



If  $R = \emptyset$ , then  $L(R) = \emptyset$  is accepted by



## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

The rest is union, concatenation and Kleene star of regular languages, as discussed last week!

## Definition (Regular Expression)

Given an alphabet  $\Sigma$ , a *regular expression* is

- $a$  for some  $a \in \Sigma$ ,
- $\varepsilon$ ,
- $\emptyset$ ,
- $(R_1 \cup R_2)$  for regular expressions  $R_1, R_2$ ,
- $(R_1 R_2)$  for regular expressions  $R_1, R_2$ ,
- $R_1^*$  for a regular expression  $R_1$ .

The rest is union, concatenation and Kleene star of regular languages, as discussed last week!

(recall: the union/concatenation/Kleene star of regular languages is itself regular)

# Regular Expressions and Automata

So we've just proven

## Proposition

*Every language described by a RE is regular.*

# Regular Expressions and Automata

So we've just proven

## Proposition

*Every language described by a RE is regular.*

Next:

## Proposition

*Every regular language can be described using a RE.*

# GNFA

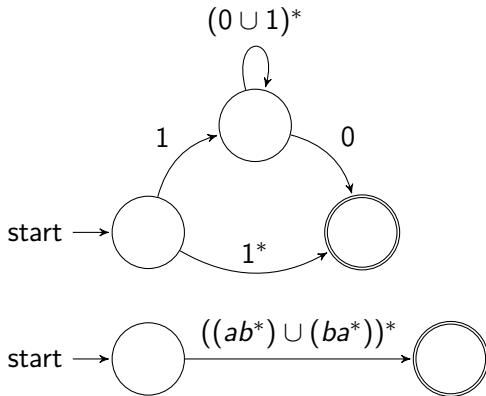
Generalized Nondeterministic Finite Automaton  
(GNFA):

- NFA where the transitions are RE, not only symbols from  $\Sigma$ .

# GNFA

Generalized Nondeterministic Finite Automaton (GNFA):

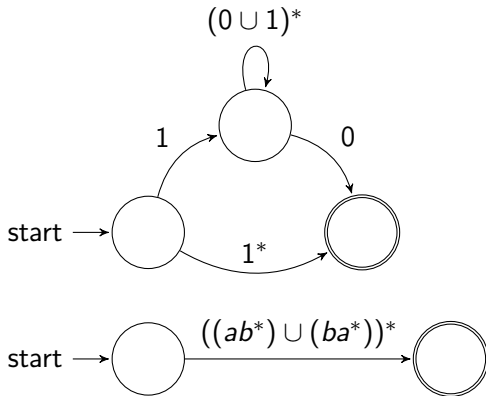
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .



# GNFA

Generalized Nondeterministic Finite Automaton (GNFA):

- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:

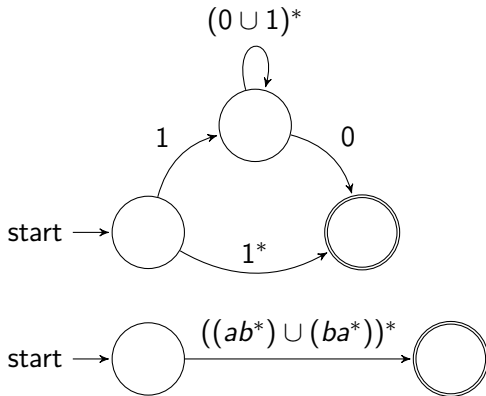




# GNFA

Generalized Nondeterministic Finite Automaton (GNFA):

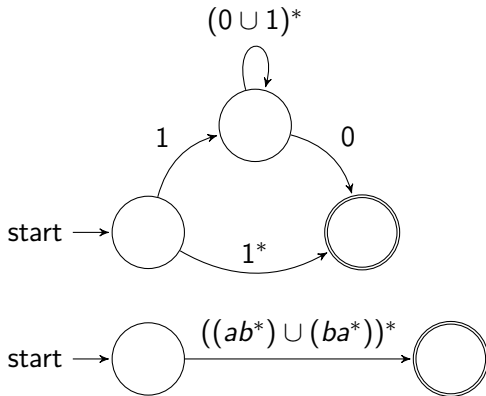
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states



# GNFA

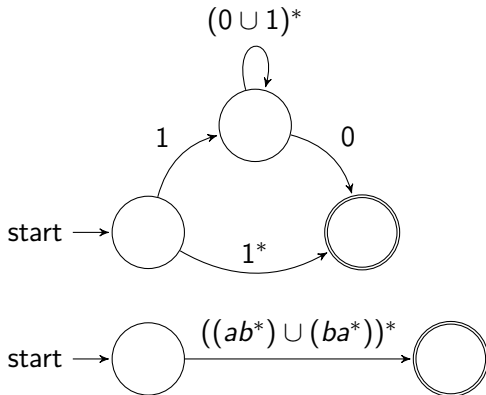
Generalized Nondeterministic Finite Automaton (GNFA):

- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.



## Generalized Nondeterministic Finite Automaton (GNFA):

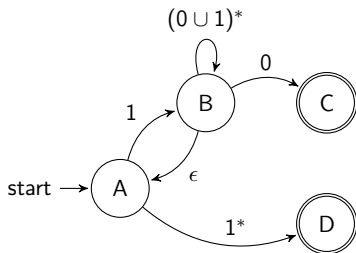
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.



## GNFA: Convenient assumptions

Generalized Nondeterministic Finite Automaton (GNFA):

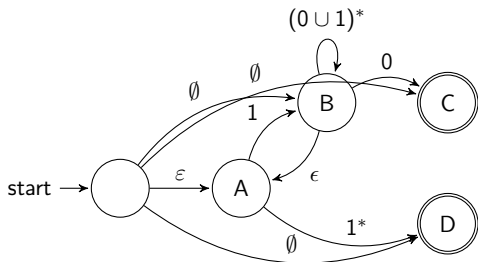
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.



## GNFA: Convenient assumptions

Generalized Nondeterministic Finite Automaton (GNFA):

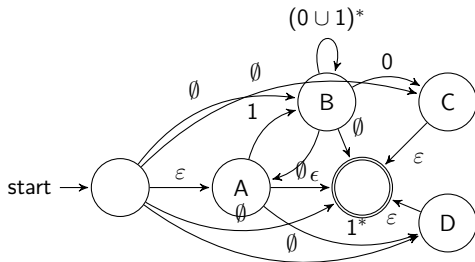
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states ✓
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.



## GNFA: Convenient assumptions

Generalized Nondeterministic Finite Automaton (GNFA):

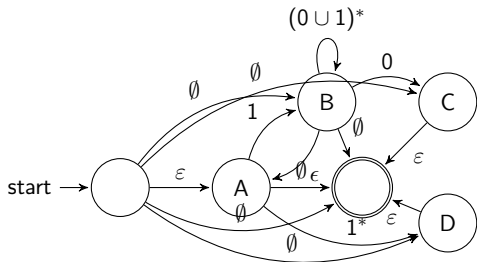
- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states ✓
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows. ✓
- all other states have one transition to all other states, including themselves. (✓)



## GNFA: Convenient assumptions

### Generalized Nondeterministic Finite Automaton (GNFA):

- NFA where the transitions are RE, not only symbols from  $\Sigma$ .
- some other assumptions for convenience:
- start state goes to every other state, but has no incoming states ✓
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows. ✓
- all other states have one transition to all other states, including themselves. (✓)



for the last point, add  $\emptyset$  transitions between any two non-accepting/starting states that were not previously connected (e.g.,  $(B, D)$ )

# Generalized Nondeterministic Finite Automata

## Definition

A generalized nondeterministic finite automaton (GNFA) is a 5-tuple  $(Q, \Sigma, \delta, q_{start}, q_{accept})$  where

- 1  $Q$  is a finite set of states
- 2  $\Sigma$  is the input alphabet
- 3  $\delta : (Q \setminus \{q_{accept}\}) \times (Q \setminus \{q_{start}\}) \rightarrow \mathcal{R}$  is the transition function, where  $\mathcal{R}$  is the set of all RE's over  $\Sigma$ ,
- 4  $q_{start}$  is the start state, and
- 5  $q_{accept}$  is the accept state.



# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Proof idea: take DFA and transform into a GNFA that accepts the same language. Iteratively remove (non-starting and non-accepting) states so that the same language is accepted, until only the starting and accepting state remain. Then the RE along the transition between the two states describes the regular language.

## Proposition

*Every regular language can be described using a RE.*

Proof:

- Given a DFA  $M$ , we construct an equivalent GNFA  $G$  by adding a new start state  $q_{start}$  with an  $\varepsilon$  transition to the old start state  $q_0$ , as well as a new accepting state  $q_{accept}$ , with  $\varepsilon$  transitions from all old accept states.
- add  $\emptyset$  transitions for all state pairs that do not have a transition in  $M$ .

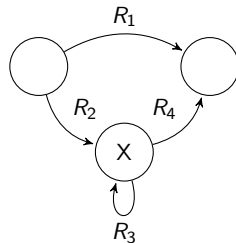
# Regular Expressions and Automata

- Recall the “convenient” properties of GNFA:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.

# Regular Expressions and Automata

- Recall the “convenient” properties of GNFA:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.

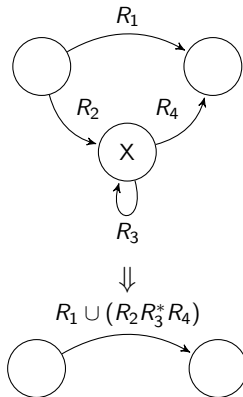
⇒ When removing  $X$ , we only need to consider situations like this:



# Regular Expressions and Automata

- Recall the “convenient” properties of GNFA:
- start state goes to every other state, but has no incoming states
- every state goes to the unique accepting state, which is different from the starting state. The accepting state does not have any outgoing arrows.
- all other states have one transition to all other states, including themselves.

⇒ When removing  $X$ , we only need to consider situations like this:



# Regular Expressions and Automata

Let's formalize this! Let us define a procedure  $\text{CONVERT}(G)$ :

- 1 If  $k = 2$  then  $G$  only has one start and one accept state, so return the regular expression  $R$  of the transition connecting them.

# Regular Expressions and Automata

Let's formalize this! Let us define a procedure CONVERT( $G$ ):

- 1 If  $k = 2$  then  $G$  only has one start and one accept state, so return the regular expression  $R$  of the transition connecting them.
- 2 if  $k > 2$  select a state  $q' \notin \{q_{accept}, q_{start}\}$ . Define  $G' = \{Q', \Sigma, \delta', q_{start}, q_{accept}\}$  with  $Q' = Q \setminus \{q'\}$  and

$$\delta'(q_i, q_j) = R_1 \cup R_2 R_3^* R_4$$

where  $R_1 = \delta(q_i, q_j)$ ,  $R_2 = \delta(q_i, q')$ ,  $R_3 = \delta(q', q')$ ,  $R_4 = \delta(q', q_j)$ , and .



# Regular Expressions and Automata

Let's formalize this! Let us define a procedure CONVERT( $G$ ):

- 1 If  $k = 2$  then  $G$  only has one start and one accept state, so return the regular expression  $R$  of the transition connecting them.
- 2 if  $k > 2$  select a state  $q' \notin \{q_{accept}, q_{start}\}$ . Define  $G' = \{Q', \Sigma, \delta', q_{start}, q_{accept}\}$  with  $Q' = Q \setminus \{q'\}$  and

$$\delta'(q_i, q_j) = R_1 \cup R_2 R_3^* R_4$$

where  $R_1 = \delta(q_i, q_j)$ ,  $R_2 = \delta(q_i, q')$ ,  $R_3 = \delta(q', q')$ ,  $R_4 = \delta(q', q_j)$ , and .

- 3 Return the result of CONVERT( $G'$ ).

# Regular Expressions and Automata

Let's formalize this! Let us define a procedure  $\text{CONVERT}(G)$ :

- 1 If  $k = 2$  then  $G$  only has one start and one accept state, so return the regular expression  $R$  of the transition connecting them.
- 2 if  $k > 2$  select a state  $q' \notin \{q_{\text{accept}}, q_{\text{start}}\}$ . Define  $G' = \{Q', \Sigma, \delta', q_{\text{start}}, q_{\text{accept}}\}$  with  $Q' = Q \setminus \{q'\}$  and

$$\delta'(q_i, q_j) = R_1 \cup R_2 R_3^* R_4$$

where  $R_1 = \delta(q_i, q_j)$ ,  $R_2 = \delta(q_i, q')$ ,  $R_3 = \delta(q', q')$ ,  $R_4 = \delta(q', q_j)$ , and .

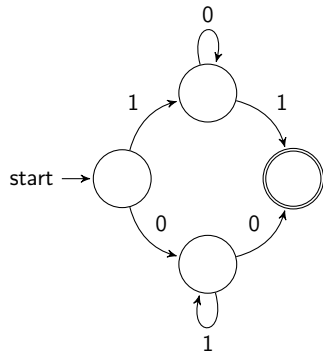
- 3 Return the result of  $\text{CONVERT}(G')$ .
- 4 correctness still remains to be shown! See book for details! (Claim 1.65)

# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Example:  
DFA:

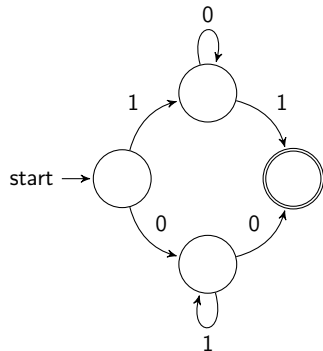


# Regular Expressions and Automata

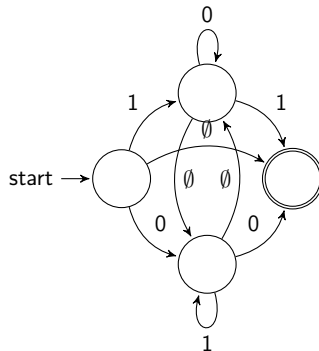
## Proposition

*Every regular language can be described using a RE.*

Example:  
DFA:



GNFA:



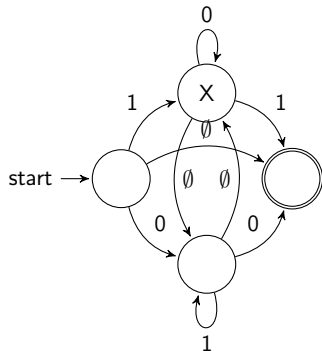
# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Example:

Remove state X:



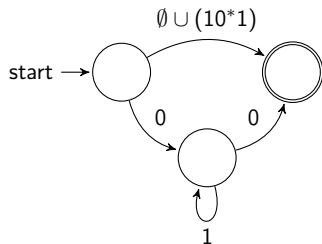
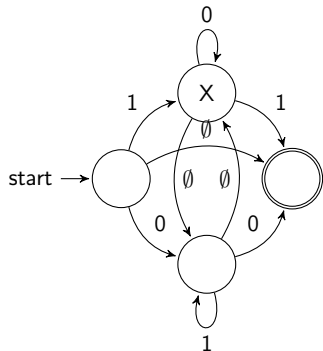
# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Example:

Remove state X:

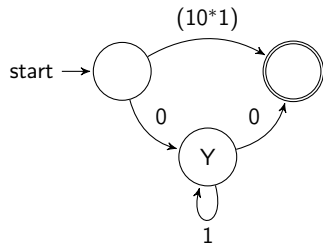


# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Example:  
Remove state Y:

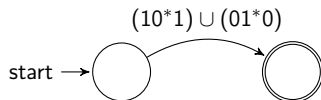
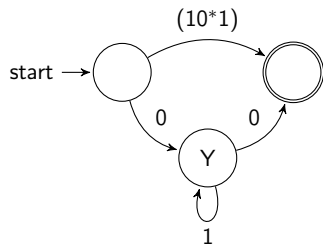


# Regular Expressions and Automata

## Proposition

*Every regular language can be described using a RE.*

Example:  
Remove state Y:





# Summary

So  $RE = GNFA = DFA = NFA = \text{Regular languages...}$

# Summary

So  $RE = GNFA = DFA = NFA = \text{Regular languages...}$   
But when is a language *nonregular*? How can we check?

# Summary

So RE = GNFA = DFA = NFA = Regular languages...

But when is a language *nonregular*? How can we check?

⇒ Pumping Lemma!

# Pumping Lemma

- DFAs only have *finite* memory, aka states.

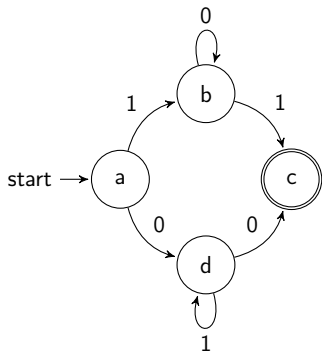
# Pumping Lemma

- DFAs only have *finite* memory, aka states.
- Pumping lemma gives a *pumping length*: if a string is longer than the pumping length, it can be *pumped*, i.e., there is a substring that can be repeated arbitrarily often such that the string remains in the language

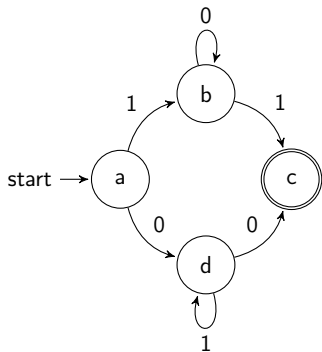
# Pumping Lemma

- DFAs only have *finite* memory, aka states.
- Pumping lemma gives a *pumping length*: if a string is longer than the pumping length, it can be *pumped*, i.e., there is a substring that can be repeated arbitrarily often such that the string remains in the language
- If a DFA has  $p$  states, and a string has length  $\geq p$ , then the accepting path in the DFA must visit at least  $p + 1$  states. In other words, at least one state appears twice.  $\Rightarrow$  loop!
- This loop can be repeated while staying in the language.

# Pumping Lemma - Example



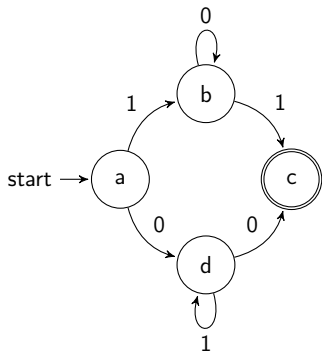
# Pumping Lemma - Example



- Language  $(10^*1) \cup (01^*0)$

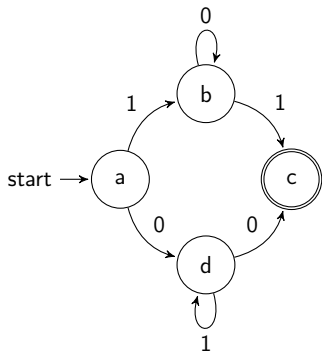


## Pumping Lemma - Example



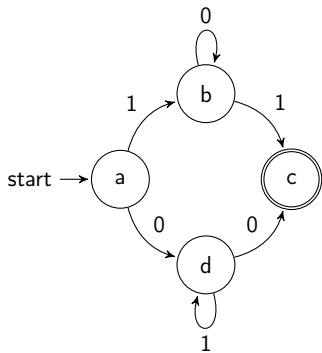
- Language  $(10^*1) \cup (01^*0)$
- DFA has 4 states

## Pumping Lemma - Example



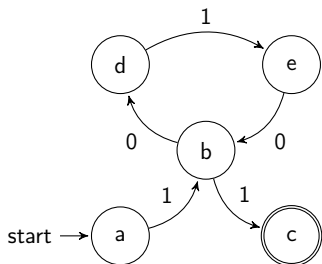
- Language  $(10^*1) \cup (01^*0)$
- DFA has 4 states
- consider string 10001, length 5

## Pumping Lemma - Example

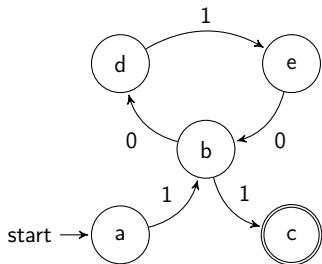


- Language  $(10^*1) \cup (01^*0)$
- DFA has 4 states
- consider string 10001, length 5
- $\Rightarrow$  path must contain a loop (in this case, at node b)

# Pumping Lemma - Example

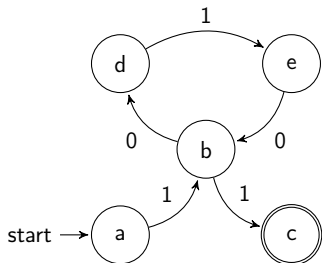


# Pumping Lemma - Example



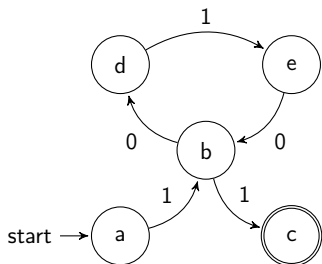
- Language  $1(010)^*1$

# Pumping Lemma - Example



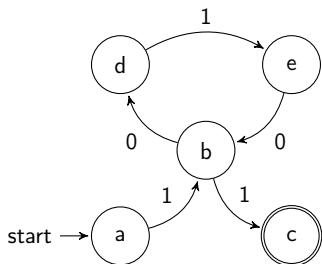
- Language  $1(010)^*1$
- DFA has 5 states

## Pumping Lemma - Example



- Language  $1(010)^*1$
- DFA has 5 states
- consider string 10101, length 5

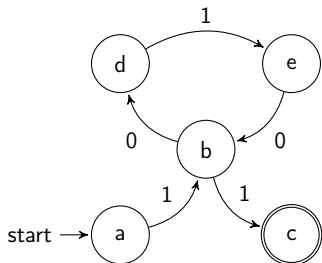
## Pumping Lemma - Example



- Language  $1(010)^*1$
- DFA has 5 states
- consider string 10101, length 5
- $\Rightarrow$  path must contain a loop (in this case, at nodes b,d,e)



## Pumping Lemma - Example



- Language  $1(010)^*1$
- DFA has 5 states
- consider string 10101, length 5
- $\Rightarrow$  path must contain a loop (in this case, at nodes b,d,e)
- $\Rightarrow$  10100101 is also a word!

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

Proof: We formalize our intuition.

- Let  $M = (Q, \Sigma, \delta, q_0, F)$  be a DFA,  $A = L(M)$  the language accepted by  $M$ , and  $p$  be the number of states in  $M$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

Proof: We formalize our intuition.

- Let  $M = (Q, \Sigma, \delta, q_0, F)$  be a DFA,  $A = L(M)$  the language accepted by  $M$ , and  $p$  be the number of states in  $M$ .
- Let  $w = w_1 \cdots w_n$  be a word in  $A$  of length  $n \geq p$ . Since  $w \in A$ , it is accepted by  $M$ , i.e., there exists a sequence of states  $s_1, s_2, \dots, s_{n+1}$  of length  $n + 1$ , where  $s_{n+1}$  is an accept state and  $\delta(s_i, w_{i+1}) = s_{i+1}$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $w = w_1 \cdots w_n$  be a word in  $A$  of length  $n \geq p$ . Since  $w \in A$ , it is accepted by  $M$ , i.e., there exists a sequence of states  $s_1, s_2, \dots, s_{n+1}$  of length  $n + 1$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $w = w_1 \cdots w_n$  be a word in  $A$  of length  $n \geq p$ . Since  $w \in A$ , it is accepted by  $M$ , i.e., there exists a sequence of states  $s_1, s_2, \dots, s_{n+1}$  of length  $n + 1$ .
- Since  $n + 1 \geq p + 1$ , one state must occur twice within the first  $p + 1$  elements of the sequence (pigeonhole principle).

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $w = w_1 \cdots w_n$  be a word in  $A$  of length  $n \geq p$ . Since  $w \in A$ , it is accepted by  $M$ , i.e., there exists a sequence of states  $s_1, s_2, \dots, s_{n+1}$  of length  $n + 1$ .
- Since  $n + 1 \geq p + 1$ , one state must occur twice within the first  $p + 1$  elements of the sequence (pigeonhole principle).
- Let these occurrences be  $s_j$  and  $s_l$ . Since these occur in the first  $p + 1$  elements of the sequence, we have  $l \leq p + 1$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let these occurrences be  $s_j$  and  $s_l$ . Since these occur in the first  $p + 1$  elements of the sequence, we have  $l \leq p + 1$ .



# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let these occurrences be  $s_j$  and  $s_l$ . Since these occur in the first  $p + 1$  elements of the sequence, we have  $l \leq p + 1$ .
- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let these occurrences be  $s_j$  and  $s_l$ . Since these occur in the first  $p + 1$  elements of the sequence, we have  $l \leq p + 1$ .
- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .
- Then  $|y| > 0$  and  $|xy| \leq p$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^i z \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .
- $x$  takes  $M$  from  $s_1$  to  $s_j$ ,  $y$  takes  $M$  from  $s_j$  to  $s_l$ ,  $z$  takes  $M$  from  $s_l$  to  $s_{n+1}$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that

- 1  $xy^i z \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .
- $x$  takes  $M$  from  $s_1$  to  $s_j$ ,  $y$  takes  $M$  from  $s_j$  to  $s_l$ ,  $z$  takes  $M$  from  $s_l$  to  $s_{n+1}$ .
- Thus the word  $xy^i z$  takes  $M$  from the start state  $s_1$  to  $s_j$ , follows the path from  $s_j$  to  $s_l$   $i$  times (recall that  $s_j = s_l$ ), then takes  $M$  from  $s_l$  to  $s_{n+1}$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that

- 1  $xy^i z \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .
- $x$  takes  $M$  from  $s_1$  to  $s_j$ ,  $y$  takes  $M$  from  $s_j$  to  $s_l$ ,  $z$  takes  $M$  from  $s_l$  to  $s_{n+1}$ .
- Thus the word  $xy^i z$  takes  $M$  from the start state  $s_1$  to  $s_j$ , follows the path from  $s_j$  to  $s_l$   $i$  times (recall that  $s_j = s_l$ ), then takes  $M$  from  $s_l$  to  $s_{n+1}$ .
- Thus  $M$  accepts any word  $xy^i z$  for  $i \geq 0$ .

# Pumping Lemma

## Lemma (Pumping Lemma)

If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that

- 1  $xy^i z \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- define  $x = w_1 \cdots w_{j-1}$ ,  $y = w_j \cdots w_{l-1}$ ,  $z = w_l \cdots w_n$ .
- $x$  takes  $M$  from  $s_1$  to  $s_j$ ,  $y$  takes  $M$  from  $s_j$  to  $s_l$ ,  $z$  takes  $M$  from  $s_l$  to  $s_{n+1}$ .
- Thus the word  $xy^i z$  takes  $M$  from the start state  $s_1$  to  $s_j$ , follows the path from  $s_j$  to  $s_l$   $i$  times (recall that  $s_j = s_l$ ), then takes  $M$  from  $s_l$  to  $s_{n+1}$ .
- Thus  $M$  accepts any word  $xy^i z$  for  $i \geq 0$ . □

# Pumping Lemma

- very useful for determining if a language is nonregular

# Pumping Lemma

- very useful for determining if a language is nonregular
- $\rightarrow$  find a string with length  $\geq p$  such that the pumping lemma does not hold



# Pumping Lemma

- very useful for determining if a language is nonregular
- $\rightarrow$  find a string with length  $\geq p$  such that the pumping lemma does not hold
- *not* very useful for proving a language is regular

# Pumping Lemma

- very useful for determining if a language is nonregular
- $\rightarrow$  find a string with length  $\geq p$  such that the pumping lemma does not hold
- *not* very useful for proving a language is regular
- $\rightarrow$  not an if and only if statement!

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- ❶  $xy^iz \in A$  for every  $i \geq 0$ ,
  - ❷  $|y| > 0$ ,
  - ❸  $|xy| \leq p$ .
- Let  $A = \{0^n1^n \mid n \geq 0\}$ .

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $A = \{0^n1^n \mid n \geq 0\}$ .
- Is  $A$  regular?

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $A = \{0^n1^n \mid n \geq 0\}$ .
- Is  $A$  regular?
- If it is, then the pumping lemma gives us a pumping length  $p$ .
- Let  $s = 0^p1^p$ .

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $A = \{0^n1^n \mid n \geq 0\}$ .
- Let  $s = 0^p1^p$ .

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- ❶  $xy^iz \in A$  for every  $i \geq 0$ ,
- ❷  $|y| > 0$ ,
- ❸  $|xy| \leq p$ .

- Let  $A = \{0^n1^n \mid n \geq 0\}$ .
- Let  $s = 0^p1^p$ .
- Condition 3 tells us that  $y$  consists of only 0s.

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

*If  $A$  is a regular language, then there is a number  $p$ , called the pumping length, where if  $w$  is a word in  $A$  of length  $\geq p$  then  $w$  can be divided into three parts,  $w = xyz$ , such that*

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $A = \{0^n1^n \mid n \geq 0\}$ .
- Let  $s = 0^p1^p$ .
- Condition 3 tells us that  $y$  consists of only 0s.
- $\Rightarrow$  then  $xy^iz$  for  $i \geq 2$  has more 0s than 1s. Contradiction!  $\Rightarrow A$  is nonregular.



## Pumping Lemma - Applied

- Even if a language is nonregular, it might contain strings for which the pumping lemma is true!

## Pumping Lemma - Applied

- Even if a language is nonregular, it might contain strings for which the pumping lemma is true!
- We have to be careful!

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

$|w| \geq p, w = xyz, s.t.$

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{w \mid w \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $w = (01)^p$ .

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

$|w| \geq p, w = xyz, s.t.$

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $w = (01)^p$ .
- $x = \varepsilon, y = 01, z = (01)^{p-1}$

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

$|w| \geq p, w = xyz, s.t.$

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $w = (01)^p$ .
- $x = \varepsilon, y = 01, z = (01)^{p-1}$
- all conditions are met!

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^i z \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .
- $x = \varepsilon, y = 0^p 1^p, z = \varepsilon$

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .
- $x = \varepsilon, y = 0^p 1^p, z = \varepsilon$
- looks like it can be pumped



# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .
- $x = \varepsilon, y = 0^p 1^p, z = \varepsilon$
- looks like it can be pumped, but are all conditions met?

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .
- $x = \varepsilon, y = 0^p 1^p, z = \varepsilon$
- looks like it can be pumped, but are all conditions met?
- condition 3  $\Rightarrow y$  must contain only 0s, so it cannot be pumped

# Pumping Lemma - Applied

## Lemma (Pumping Lemma)

- 1  $xy^iz \in A$  for every  $i \geq 0$ ,
- 2  $|y| > 0$ ,
- 3  $|xy| \leq p$ .

- Let  $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$ .
- Let  $s = 0^p 1^p$ .
- $x = \varepsilon, y = 0^p 1^p, z = \varepsilon$
- looks like it can be pumped, but are all conditions met?
- condition 3  $\Rightarrow y$  must contain only 0s, so it cannot be pumped  $\Rightarrow B$  nonregular!

## Pumping Lemma - Applied

- $A = \{0^n 1^n \mid n \geq 0\}$ .
- $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$
- Another way of showing  $B$  is nonregular is to reduce it to the nonregularity of  $A$ :

## Pumping Lemma - Applied

- $A = \{0^n 1^n \mid n \geq 0\}$ .
- $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$
- Another way of showing  $B$  is nonregular is to reduce it to the nonregularity of  $A$ :
- regular languages are closed under intersection

## Pumping Lemma - Applied

- $A = \{0^n 1^n \mid n \geq 0\}$ .
- $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$
- Another way of showing  $B$  is nonregular is to reduce it to the nonregularity of  $A$ :
- regular languages are closed under intersection
- and  $A = B \cap 0^*1^*$

## Pumping Lemma - Applied

- $A = \{0^n 1^n \mid n \geq 0\}$ .
- $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$
- Another way of showing  $B$  is nonregular is to reduce it to the nonregularity of  $A$ :
- regular languages are closed under intersection
- and  $A = B \cap 0^*1^*$
- if  $B$  is regular and since  $0^*1^*$  is regular, then  $A$  must be as well

## Pumping Lemma - Applied

- $A = \{0^n 1^n \mid n \geq 0\}$ .
- $B = \{\omega \mid \omega \text{ contains an equal number of 0s and 1s}\}$
- Another way of showing  $B$  is nonregular is to reduce it to the nonregularity of  $A$ :
- regular languages are closed under intersection
- and  $A = B \cap 0^*1^*$
- if  $B$  is regular and since  $0^*1^*$  is regular, then  $A$  must be as well, contradiction!



# Summary

- regular expressions are shorthand notations for languages

# Summary

- regular expressions are shorthand notations for languages
- $RE = GNFA = DFA = NFA$ , i.e., regular expressions are shorthand for *regular* languages

# Summary

- regular expressions are shorthand notations for languages
- $RE = GNFA = DFA = NFA$ , i.e., regular expressions are shorthand for *regular* languages
- proof involved transforming a DFA to a GNFA then reducing the number of states to 2 while accepting the same language

# Summary

- regular expressions are shorthand notations for languages
- $RE = GNFA = DFA = NFA$ , i.e., regular expressions are shorthand for *regular* languages
- proof involved transforming a DFA to a GNFA then reducing the number of states to 2 while accepting the same language
- $\rightarrow$  the regular expressions describe the paths in the DFA

# Summary

- regular expressions are shorthand notations for languages
- $RE = GNFA = DFA = NFA$ , i.e., regular expressions are shorthand for *regular* languages
- proof involved transforming a DFA to a GNFA then reducing the number of states to 2 while accepting the same language
- $\rightarrow$  the regular expressions describe the paths in the DFA
- every regular language has a pumping length

# Summary

- regular expressions are shorthand notations for languages
- $RE = GNFA = DFA = NFA$ , i.e., regular expressions are shorthand for *regular* languages
- proof involved transforming a DFA to a GNFA then reducing the number of states to 2 while accepting the same language
- $\rightarrow$  the regular expressions describe the paths in the DFA
- every regular language has a pumping length
- useful for determining if a language is *nonregular*