

# Regular Languages and Finite State Automata

Data structures and algorithms  
for Computational Linguistics III

Çağrı Çöltekin  
ccoltekin@sfs.uni-tuebingen.de

University of Tübingen  
Seminar für Sprachwissenschaft

Winter Semester 2019–2020

## Why study finite-state automata?

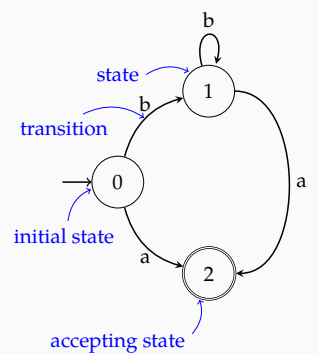
- Unlike some of the abstract machines we discussed, finite-state automata are efficient models of computation
- There are many applications
  - Electronic circuit design
  - Workflow management
  - Games
  - Pattern matching
  - ...
- But more importantly ;-)  
  - Tokenization, stemming
  - Morphological analysis
  - Shallow parsing/chunking
  - ...

## Finite-state automata (FSA)

- A finite-state machine is in one of a finite-number of states in a given time
  - The machine changes its state based on its input
  - Every regular language is generated/recognized by an FSA
  - Every FSA generates/recognizes a regular language
  - Two flavors:
    - Deterministic finite automata (DFA)
    - Non-deterministic finite automata (NFA)
- Note: the NFA is a superset of DFA.

## DFA as a graph

- States are represented as nodes
- Transitions are shown by the edges, labeled with symbols from an alphabet
- One of the states is marked as the *initial state*
- Some states are accepting states



## DFA: formal definition

Formally, a finite state automaton,  $M$ , is a tuple  $(\Sigma, Q, q_0, F, \Delta)$  with

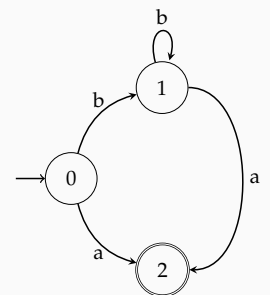
- $\Sigma$  is the alphabet, a finite set of symbols
- $Q$  a finite set of states
- $q_0$  is the start state,  $q_0 \in Q$
- $F$  is the set of final states,  $F \subseteq Q$
- $\Delta$  is a function that takes a state and a symbol in the alphabet, and returns another state ( $\Delta : Q \times \Sigma \rightarrow Q$ )

At any given time, for any input, a DFA has a single well-defined action to take.

## DFA: formal definition

an example

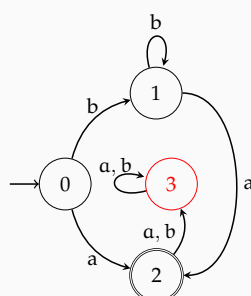
- $\Sigma = \{a, b\}$
- $Q = \{q_0, q_1, q_2\}$
- $q_0 = q_0$
- $F = \{q_2\}$
- $\Delta = \{(q_0, a) \rightarrow q_2, (q_0, b) \rightarrow q_1, (q_1, a) \rightarrow q_2, (q_1, b) \rightarrow q_1\}$



## Another note on DFA

error or sink state

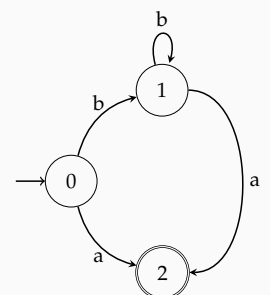
- Is this FSA deterministic?
- To make all transitions well-defined, we can add a sink (or error) state
- For brevity, we skip the explicit error state
  - In that case, when we reach a dead end, recognition fails



## DFA: the transition table

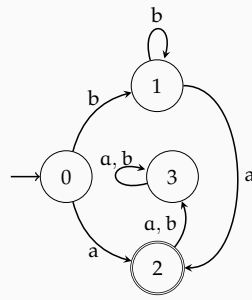
transition table			
		symbol	
		a	b
state	→0	2	1
	1	2	1
	*2	∅	∅

- marks the start state
- \* marks the accepting state(s)



## DFA: the transition table

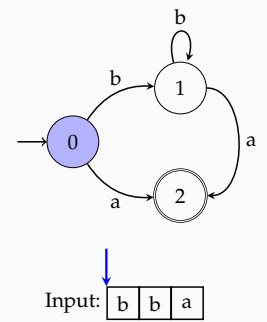
transition table		
state	symbol	
	a	b
→0	2	1
1	2	1
*2	3	3
3	3	3



- marks the start state
- \* marks the accepting state(s)

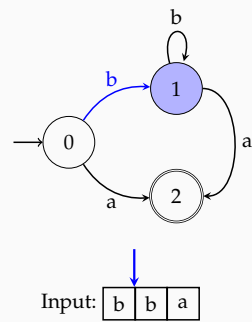
## DFA recognition

1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input



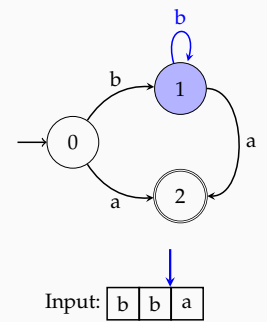
## DFA recognition

1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input



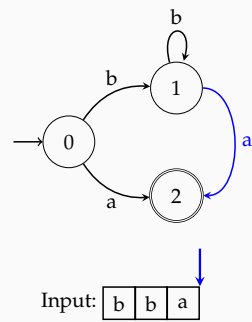
## DFA recognition

1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input



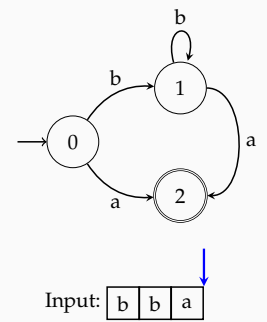
## DFA recognition

1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input



## DFA recognition

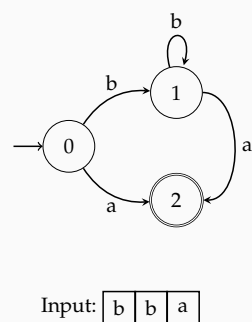
1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input



## DFA recognition

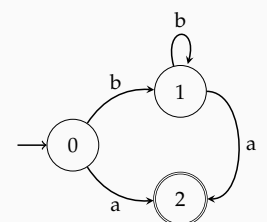
1. Start at  $q_0$
2. Process an input symbol, move accordingly
3. Accept if in a final state at the end of the input

- What is the complexity of the algorithm?
- How about inputs:
  - bbbb
  - aa



## A few questions

- What is the language recognized by this FSA?
- Can you draw a simpler DFA for the same language?
- Draw a DFA recognizing strings with even number of 'a's over  $\Sigma = \{a, b\}$



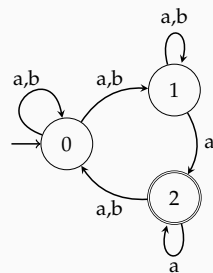
# Non-deterministic finite automata

## Formal definition

A non-deterministic finite state automaton,  $M$ , is a tuple  $(\Sigma, Q, q_0, F, \Delta)$  with

- $\Sigma$  is the alphabet, a finite set of symbols
- $Q$  a finite set of states
- $q_0$  is the start state,  $q_0 \in Q$
- $F$  is the set of final states,  $F \subseteq Q$
- $\Delta$  is a function from  $(Q, \Sigma)$  to  $P(Q)$ , power set of  $Q$  ( $\Delta : Q \times \Sigma \rightarrow P(Q)$ )

# An example NFA



transition table			
		symbol	
		a	b
state	→0	0,1	0,1
	1	1,2	1
	*2	0,2	0

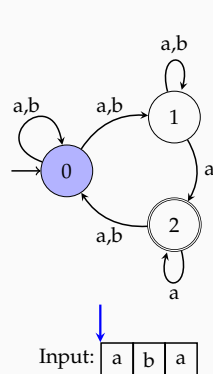
- We have nondeterminism, e.g., if the first input is  $a$ , we need to choose between states 0 or 1
- Transition table cells have *sets* of states

# Dealing with non-determinism

- Follow one of the links, store alternatives, and *backtrack* on failure
- Follow all options in parallel
- Use dynamic programming (e.g., as in chart parsing)

# NFA recognition

as search (with backtracking)

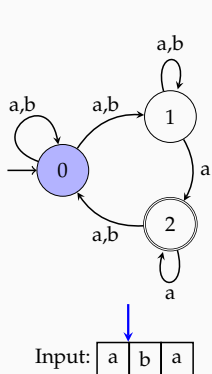


Agenda

- Start at  $q_0$
- Take the next input, place all possible actions to an agenda
- Get the next action from the agenda, act
- At the end of input  
Accept if in an accepting state  
Reject not in accepting state & agenda empty  
Backtrack otherwise

# NFA recognition

as search (with backtracking)

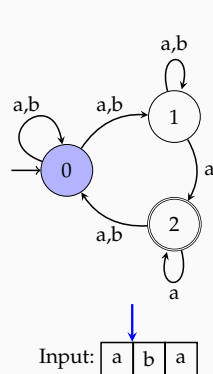


Agenda
(q <sub>0</sub> , 1)
(q <sub>1</sub> , 1)

- Start at  $q_0$
- Take the next input, place all possible actions to an agenda
- Get the next action from the agenda, act
- At the end of input  
Accept if in an accepting state  
Reject not in accepting state & agenda empty  
Backtrack otherwise

# NFA recognition

as search (with backtracking)

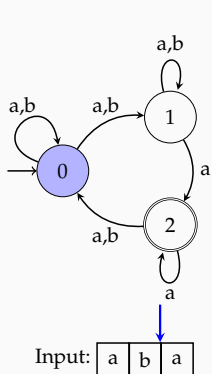


Agenda
(q <sub>0</sub> , 1)
(q <sub>1</sub> , 1)

- Start at  $q_0$
- Take the next input, place all possible actions to an agenda
- Get the next action from the agenda, act
- At the end of input  
Accept if in an accepting state  
Reject not in accepting state & agenda empty  
Backtrack otherwise

# NFA recognition

as search (with backtracking)

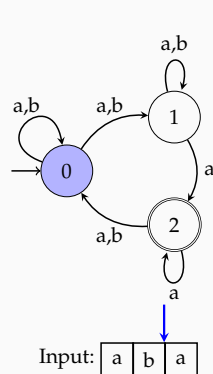


Agenda
(q <sub>0</sub> , 2)
(q <sub>1</sub> , 2)
(q <sub>1</sub> , 1)

- Start at  $q_0$
- Take the next input, place all possible actions to an agenda
- Get the next action from the agenda, act
- At the end of input  
Accept if in an accepting state  
Reject not in accepting state & agenda empty  
Backtrack otherwise

# NFA recognition

as search (with backtracking)

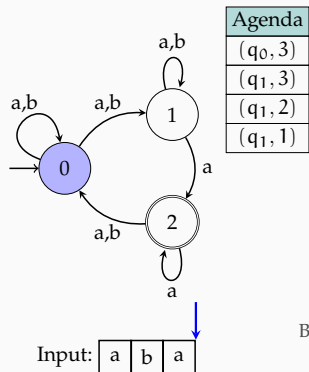


Agenda
(q <sub>0</sub> , 2)
(q <sub>1</sub> , 2)
(q <sub>1</sub> , 1)

- Start at  $q_0$
- Take the next input, place all possible actions to an agenda
- Get the next action from the agenda, act
- At the end of input  
Accept if in an accepting state  
Reject not in accepting state & agenda empty  
Backtrack otherwise

## NFA recognition

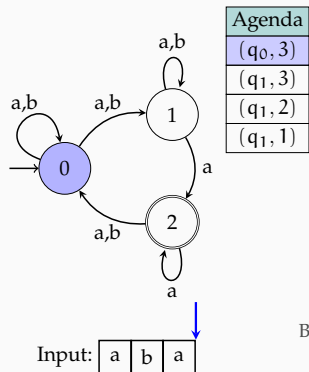
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

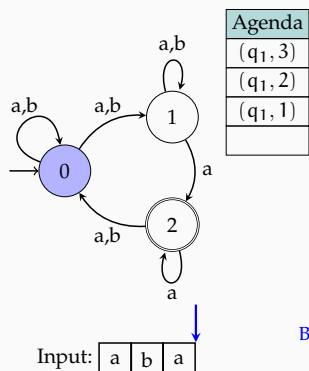
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

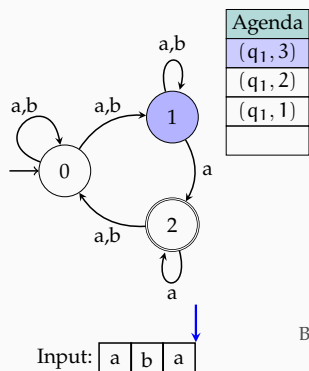
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

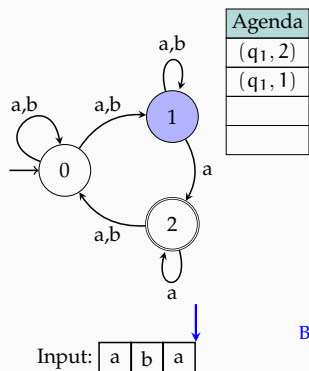
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

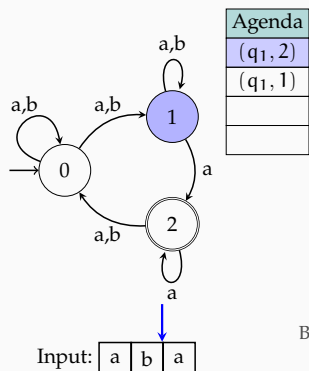
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

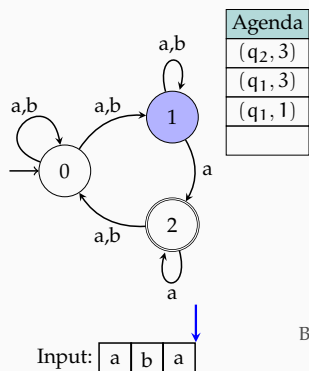
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

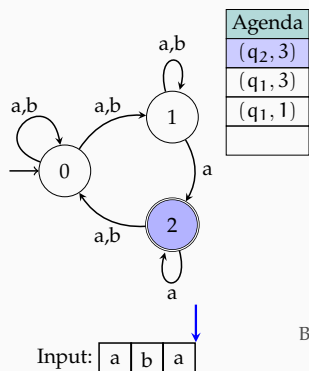
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

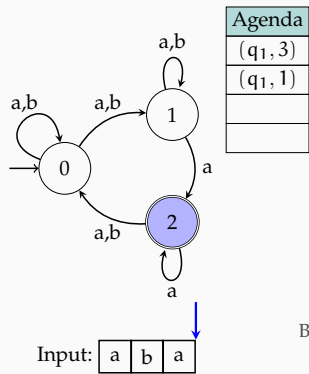
as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input  
 Accept if in an accepting state  
 Reject not in accepting state & agenda empty  
 Backtrack otherwise

## NFA recognition

as search (with backtracking)



1. Start at  $q_0$
2. Take the next input, place all possible actions to an agenda
3. Get the next action from the agenda, act
4. At the end of input
  - Accept if in an accepting state
  - Reject not in accepting state & agenda empty
  - Backtrack otherwise

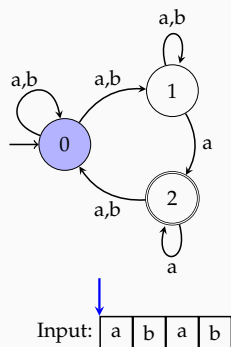
## NFA recognition as search

summary

- Worst time complexity is exponential
  - Complexity is worse if we want to enumerate all derivations
- We used a stack as *agenda*, performing a depth-first search
- A queue would result in breadth-first search
- If we have a reasonable heuristic A\* search may be an option
- Machine learning methods may also guide finding a fast or the best solution

## NFA recognition

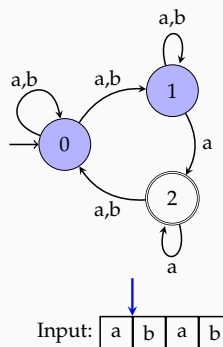
parallel version



1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

## NFA recognition

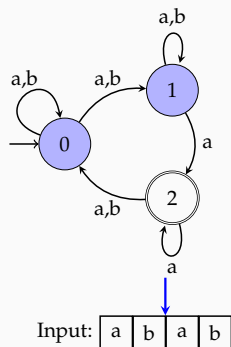
parallel version



1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

## NFA recognition

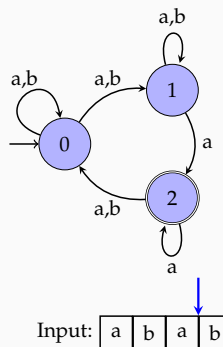
parallel version



1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

## NFA recognition

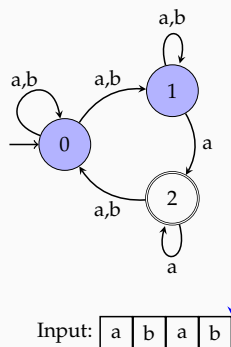
parallel version



1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

## NFA recognition

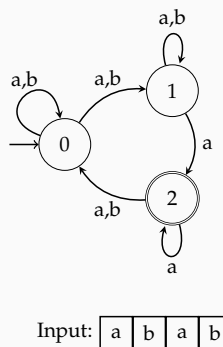
parallel version



1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

## NFA recognition

parallel version

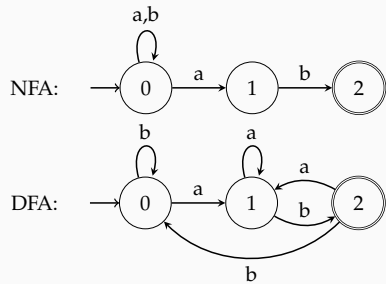


1. Start at  $q_0$
2. Take the next input, mark all possible next states
3. If an accepting state is marked at the end of the input, accept

Note: the process is *deterministic*, and *finite-state*.

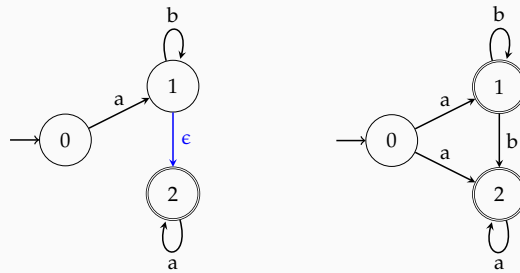
### An exercise

Construct an NFA and a DFA for the language over  $\Sigma = \{a, b\}$  where all sentences end with  $ab$ .

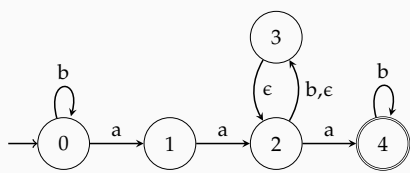


### One more complication: $\epsilon$ transitions

- An extension of NFA,  $\epsilon$ -NFA, allows moving without consuming an input symbol, indicated by an  $\epsilon$ -transition (sometimes called a  $\lambda$ -transition)
- Any  $\epsilon$ -NFA can be converted to an NFA



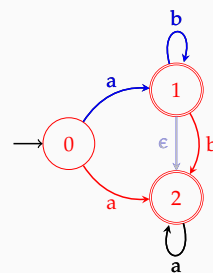
### $\epsilon$ -transitions need attention



- How does the (depth-first) NFA recognition algorithm we described earlier work on this automaton?
- Can we do without  $\epsilon$  transitions?

### $\epsilon$ removal

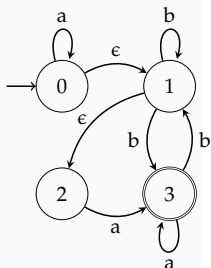
- We start with finding the  $\epsilon$ -closure of all states
  - $\epsilon$ -closure( $q_0$ ) = { $q_0$ }
  - $\epsilon$ -closure( $q_1$ ) = { $q_1, q_2$ }
  - $\epsilon$ -closure( $q_2$ ) = { $q_2$ }
- Replace each arc to each state with arc(s) to all states in the  $\epsilon$ -closure of the state



### $\epsilon$ removal

a(nother) solution with the transition table

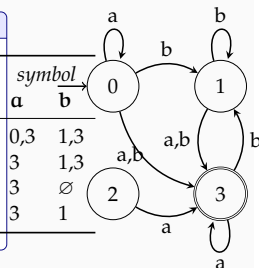
		symbol			
		a	b	$\epsilon$	$\epsilon^*$
state	$\rightarrow 0$	0	$\emptyset$	1	0,1,2
	1	$\emptyset$	1,3	2	1,2
	2	3	$\emptyset$	$\emptyset$	2
	*3	3	1	$\emptyset$	3



### $\epsilon$ removal

a(nother) solution with the transition table

		symbol			
		a	b	$\epsilon$	$\epsilon^*$
state	$\rightarrow 0$	0	$\emptyset$	1	0,1,2
	1	$\emptyset$	1,3	2	1,2
	2	3	$\emptyset$	$\emptyset$	2
	*3	3	1	$\emptyset$	3



### NFA-DFA equivalence

- The language recognized by every NFA is recognized by some DFA
- The set of DFA is a subset of the set of NFA (a DFA is also an NFA)
- The same is true for  $\epsilon$ -NFA
- All recognize/generate regular languages
- NFA can automatically be converted to the equivalent DFA

### Why do we use an NFA then?

- NFA (or  $\epsilon$ -NFA) are often easier to construct
  - Intuitive for humans (cf. earlier exercise)
  - Some representations are easy to convert to NFA rather than DFA, e.g., regular expressions
- NFA may require less memory (fewer states)

A quick exercise – and a not-so-quick one

1. Construct (draw) an NFA for the language over  $\Sigma = \{a, b\}$ , such that 4th symbol from the end is an  $a$
2. Construct a DFA for the same language

## Determinization

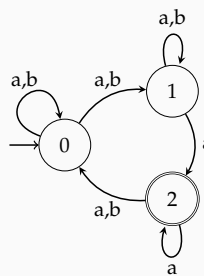
the subset construction

Intuition: remember **the parallel NFA recognition**. We can consider an NFA being a deterministic machine which is at a **set of states** at any given time.

- **Subset construction** (sometimes called power set construction) uses this intuition to convert an NFA to a DFA
- The algorithm can be modified to handle  $\epsilon$ -transitions (or we can eliminate  $\epsilon$ 's as a preprocessing step)

## The subset construction

by example



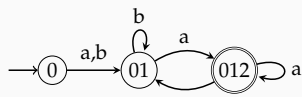
transition table with subsets

	symbol	
	a	b
$\emptyset$	$\emptyset$	$\emptyset$
$\rightarrow \{0\}$	$\{0, 1\}$	$\{0, 1\}$
$\{1\}$	$\{1, 2\}$	$\{1\}$
$* \{2\}$	$\{0, 2\}$	$\{0\}$
$\{0, 1\}$	$\{0, 1, 2\}$	$\{0, 1\}$
$* \{0, 2\}$	$\{0, 1, 2\}$	$\{0, 1\}$
$* \{1, 2\}$	$\{0, 1, 2\}$	$\{0, 1\}$
$* \{0, 1, 2\}$	$\{0, 1, 2\}$	$\{0, 1\}$

## The subset construction

by example: the resulting DFA

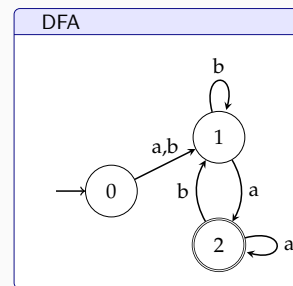
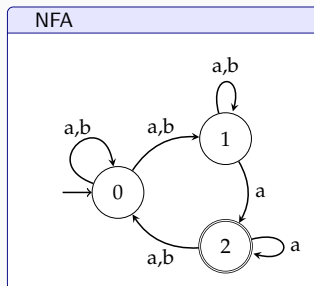
	symbol	
	a	b
$\rightarrow \{0\}$	$\{0, 1\}$	$\{0, 1\}$
$\{0, 1\}$	$\{0, 1, 2\}$	$\{0, 1\}$
$* \{0, 1, 2\}$	$\{0, 1, 2\}$	$\{0, 1\}$



Do you remember the set of states marked during **parallel NFA recognition**?

## The subset construction

by example: side by side

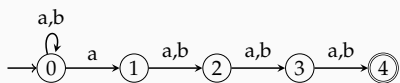


- What language do they recognize?

## The subset construction

wrapping up

- In worst case, resulting DFA has  $2^n$  nodes
- Worst case is rather rare, number of nodes in an NFA and the converted DFA are often similar
- In practice, we do not need to enumerate all  $2^n$  subsets
- We've already seen a typical problematic case:



- We can also skip the unreachable states during subset construction

## Yet another exercise

Determinize the following automaton

## Regular languages: definition

A regular grammar is a tuple  $G = (\Sigma, N, S, R)$  where

$\Sigma$  is an alphabet of terminal symbols

$N$  are a set of non-terminal symbols

$S$  is a special 'start' symbol  $\in N$

$R$  is a set of rewrite rules following one of the following patterns ( $A, B \in N, a \in \Sigma, \epsilon$  is the empty string)

Left regular
1. $A \rightarrow a$
2. $A \rightarrow Ba$
3. $A \rightarrow \epsilon$

Right regular
1. $A \rightarrow a$
2. $A \rightarrow aB$
3. $A \rightarrow \epsilon$

## Regular languages: another definition

A language is regular if there is an FSA that recognizes it

- We denote the language recognized by a finite state automaton  $M$ , as  $\mathcal{L}(M)$
- The above definition reformulated: if a language  $L$  is regular, there is a DFA  $M$ , such that  $\mathcal{L}(M) = L$
- Remember: any NFA (with or without  $\epsilon$  transitions) can be converted to a DFA

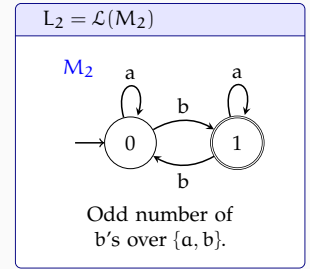
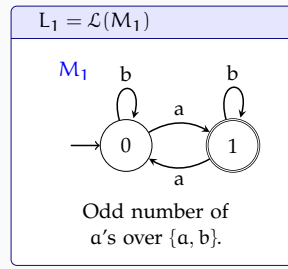
## Some operations on regular languages (and FSA)

- $L_1 L_2$  Concatenation of two languages  $L_1$  and  $L_2$ : any sentence of  $L_1$  followed by any sentence of  $L_2$
- $L^*$  Kleene star of  $L$ :  $L$  concatenated by itself 0 or more times
- $L^R$  Reverse of  $L$ : reverse of any string in  $L$
- $\overline{L}$  Complement of  $L$ : all strings in  $\Sigma_1^*$  except the ones in  $L$  ( $\Sigma_1^* - L$ )
- $L_1 \cup L_2$  Union of languages  $L_1$  and  $L_2$ : strings that are in any of the languages
- $L_1 \cap L_2$  Intersection of languages  $L_1$  and  $L_2$ : strings that are in both languages

Regular languages are closed under all of these operations.

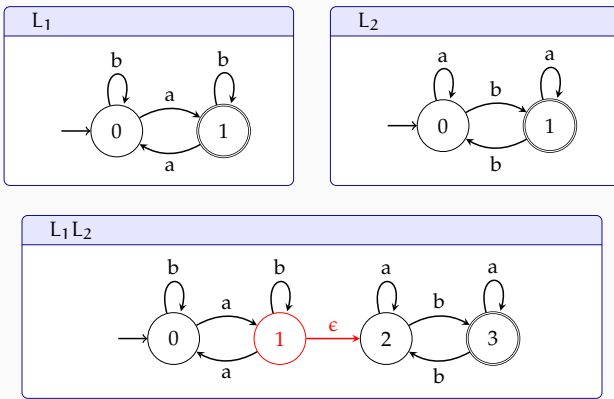
## Two example FSA

what languages do they accept?

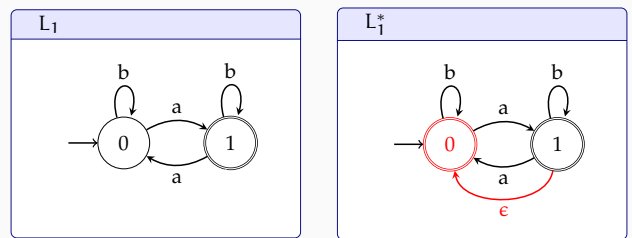


We will use these languages and automata for demonstration.

## Concatenation

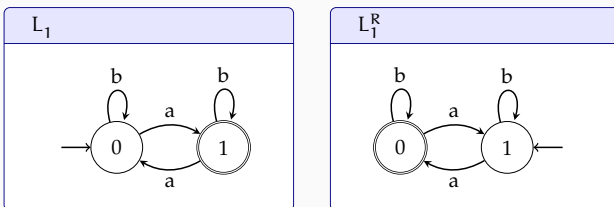


## Kleene star

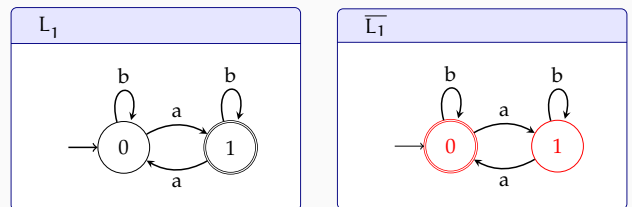


- What if there were more than one accepting states?

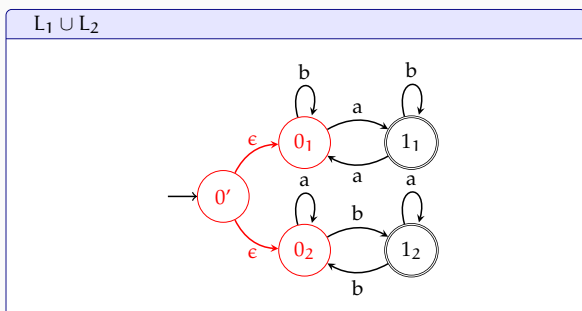
## Reversal



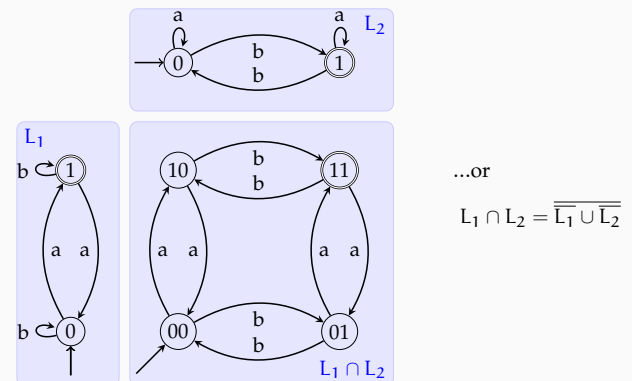
## Complement



## Union



## Intersection





## Closure properties of regular languages

- Since results of all the operations we studied are FSA: Regular languages are closed under
  - Concatenation
  - Kleene star
  - Reversal
  - Complement
  - Union
  - Intersection

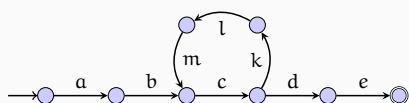
## Is a language regular?

— or not

- To show that a language is regular, it is sufficient to find an FSA that recognizes it.
- Showing that a language is *not* regular is more involved
- We will study a method based on *pumping lemma*

## Pumping lemma

intuition



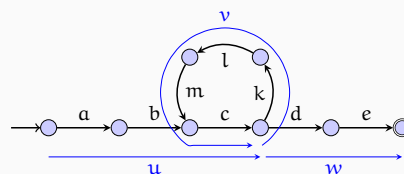
- What is the length of longest string generated by this FSA?
- Any FSA generating an infinite language has to have a loop (application of recursive rule(s) in the grammar)
- Part of every string longer than some number will include repetition of the same substring ('cklm' above)

## Pumping lemma

definition

For every regular language  $L$ , there exist an integer  $p$  such that a string  $x \in L$  can be factored as  $x = uvw$ ,

- $uv^i w \in L, \forall i \geq 0$
- $v \neq \epsilon$
- $|uv| \leq p$



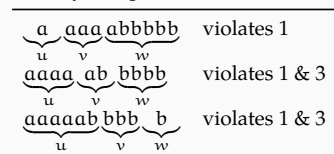
## How to use pumping lemma

- We use pumping lemma to prove that a language is not regular
- Proof is by contradiction:
  - Assume the language is regular
  - Find a string  $x$  in the language, for all splits of  $x = uvw$ , at least one of the pumping lemma conditions does not hold
    - $uv^i w \in L (\forall i \geq 0)$
    - $v \neq \epsilon$
    - $|uv| \leq p$

## Pumping lemma example

prove  $L = a^n b^n$  is not regular

- Assume  $L$  is regular: there must be a  $p$  such that, if  $uvw$  is in the language
  1.  $uv^i w \in L (\forall i \geq 0)$
  2.  $v \neq \epsilon$
  3.  $|uv| \leq p$
- Pick the string  $a^p b^p$
- For the sake of example, assume  $p = 5, x = aaaaaabbbb$
- Three different ways to split

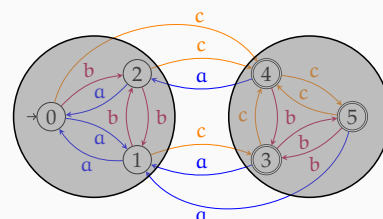


## DFA minimization

- For any regular language, there is a unique *minimal* DFA
- By finding the minimal DFA, we can also prove equivalence (or not) of different FSA
- In general the idea is:
  - Throw away unreachable states (easy)
  - Merge equivalent states
- There are two well-known algorithms for minimization:
  - Hopcroft's algorithm: find and eliminate equivalent states by partitioning the set of states
  - Brzowski's algorithm: 'double reversal'

## Finding equivalent states

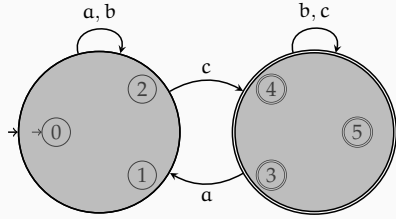
Intuition



The edges leaving the group of nodes are identical. Their *right* languages are the same.

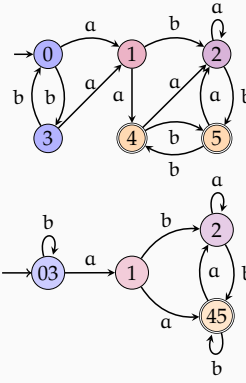
## Finding equivalent states

Intuition



The edges leaving the group of nodes are identical.  
Their right languages are the same.

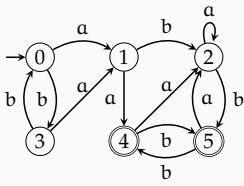
## Minimization by partitioning



- Accepting & non-accepting states form a partition  
 $Q_2 = \{0, 1, 2, 3\}$ ,  $Q_2 = \{4, 5\}$
- If any two nodes go to different sets for any of the symbols split
- $Q_1 = \{0, 3\}$ ,  $Q_3 = \{1\}$ ,  $Q_4 = \{2\}$ ,  $Q_2 = \{4, 5\}$
- Stop when we cannot split any of the sets, merge the indistinguishable states

## Minimization by partitioning

tabular version

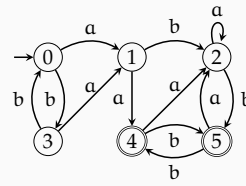


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

1					
2					
3					
4					
5					
	0	1	2	3	4

## Minimization by partitioning

tabular version

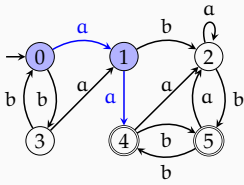


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

1					
2					
3					
4	•	•	•	•	
5	•	•	•	•	
	0	1	2	3	4

## Minimization by partitioning

tabular version

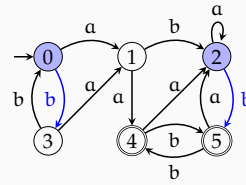


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

1					
2					
3					
4	•	•	•	•	
5	•	•	•	•	
	0	1	2	3	4

## Minimization by partitioning

tabular version

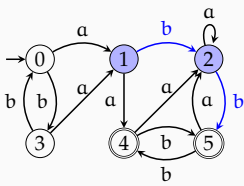


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

1	•				
2					
3					
4	•	•	•	•	
5	•	•	•	•	
	0	1	2	3	4

## Minimization by partitioning

tabular version

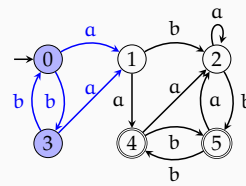


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

1	•				
2	•				
3					
4	•	•	•	•	
5	•	•	•	•	
	0	1	2	3	4

## Minimization by partitioning

tabular version

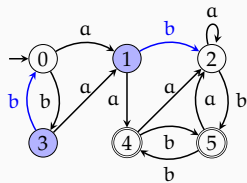


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

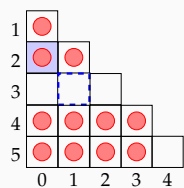
1	•	•			
2	•	•			
3					
4	•	•	•	•	
5	•	•	•	•	
	0	1	2	3	4

## Minimization by partitioning

tabular version

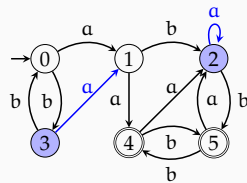


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

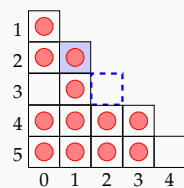


## Minimization by partitioning

tabular version

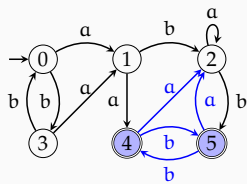


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

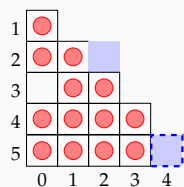


## Minimization by partitioning

tabular version

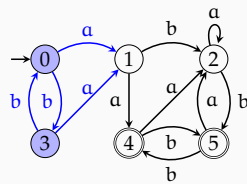


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

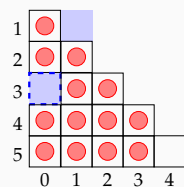


## Minimization by partitioning

tabular version

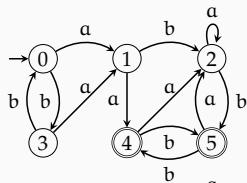


- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$

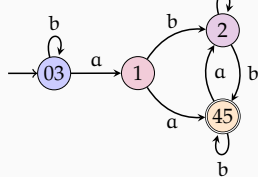
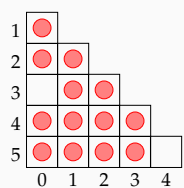


## Minimization by partitioning

tabular version



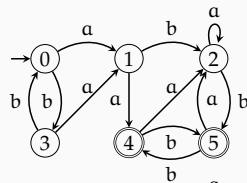
- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$



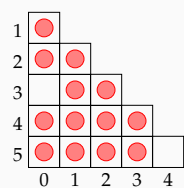
- Merge indistinguishable states

## Minimization by partitioning

tabular version



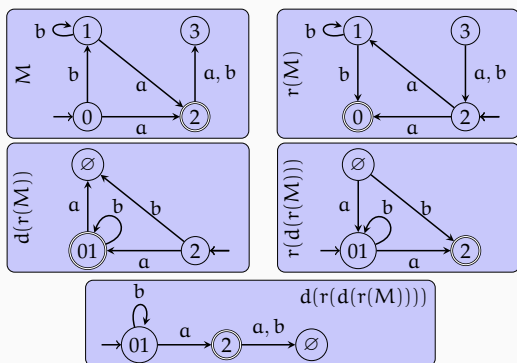
- Create a state-by-state table, mark *distinguishable* pairs:  $(q_1, q_2)$  such that  $(\Delta(q_1, x), \Delta(q_2, x))$  is a distinguishable pair for any  $x \in \Sigma$



- Merge indistinguishable states
- The algorithm can be improved by choosing which cell to visit carefully

## Brzowski's algorithm

double reverse (r), determinize (d)



## Minimization algorithms

final remarks

- There are many versions of the 'partitioning' algorithm. General idea is to form equivalence classes based on *right-language* of each state.
- Partitioning algorithm has  $O(n \log n)$  complexity
- 'Double reversal' algorithm has exponential worst-time complexity
- Double reversal algorithm can also be used with NFAs (resulting in the minimal equivalent DFA – NFA minimization is intractable)
- In practice, there is no clear winner, different algorithms run faster on different input

## Regular expressions

- Another way to specify a regular language (RL) is use of *regular expressions* (RE)
- Every RL can be expressed by a RE, and every RE defines a RL
- A RE  $x$  defines a RL  $\mathcal{L}(x)$
- Relations between RE and RL

$$\begin{aligned} & - \mathcal{L}(\emptyset) = \emptyset, & - \mathcal{L}(a|b) = \mathcal{L}(a) \cup \mathcal{L}(b) \\ & - \mathcal{L}(\epsilon) = \epsilon, & & \text{(some author use the} \\ & - \mathcal{L}(a) = a, & & \text{notation } a+b, \text{ we will use} \\ & - \mathcal{L}(ab) = \mathcal{L}(a)\mathcal{L}(b) & & \text{notation } a|b \text{ as in many practical} \\ & - \mathcal{L}(a^*) = \mathcal{L}(a)^* & & \text{implementations)} \end{aligned}$$

where,  $a, b \in \Sigma$ ,  $\epsilon$  is empty string,  $\emptyset$  is the language that accepts nothing (e.g.,  $\Sigma^* - \Sigma^*$ )

- Note: no standard complement operation in RE

## Regular

some extensions

- Kleene star ( $a^*$ ), Concatenation ( $ab$ ) and union ( $a|b$ ) are the common operations
- Parentheses can be used to group the sub-expressions. Otherwise, the priority of the operators as specified above  $a|bc^* = a|(b(c^*))$
- In practice some short-hand notations are common
  - $\cdot = (a_1 | \dots | a_n)$ , for  $\Sigma = \{a_1, \dots, a_n\}$
  - $a^+ = aa^*$
  - $[a-c] = (a|b|c)$
  - $[\hat{a-c}] = \cdot - (a|b|c)$
  - $\backslash d = (0|1|\dots|8|9)$
  - ...
- And some non-regular extensions, like  $(a^*)b \setminus 1$  (sometimes the term *regex* is used for expressions with non-regular extensions)

## Some properties of regular expressions

Kleene algebra

These identities are often used to simplify regular expressions.

- $\epsilon u = u$
- $\emptyset u = \emptyset$
- $u(vw) = (uv)w$
- $\emptyset^* = \epsilon$
- $\epsilon^* = \epsilon$
- $(u^*)^* = u^*$
- $u|v = v|u$
- $u|u = u$
- $u|\emptyset = u$
- $u|\epsilon = u$
- $u|(v|w) = (u|v)|w$
- $u|(v|w) = uv|uw$
- $(u|v)^* = (u^*|v^*)^*$

An exercise

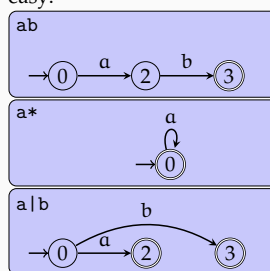
Simplify  $a|ab^*$

$$\begin{aligned} a|ab^* &= a\epsilon|ab^* \\ &= a(\epsilon|b^*) \\ &= ab^* \end{aligned}$$

Note: most of these follow from set theory, and some can be derived from others.

## Converting regular expressions to FSA

Converting to NFA is easy:

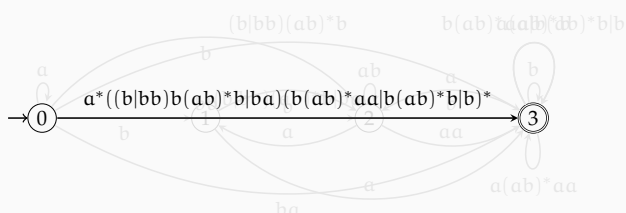


Note the similarity with operations on regular languages discussed earlier.

- For more complex expressions, one can replace the paths for individual symbols with corresponding automata
- Using  $\epsilon$  transitions may ease the task
- The reverse conversion (from automata to regular expressions) is also easy:
  - identify the patterns on the left, collapse paths to single transitions with regular expressions

## Converting FSA to regular expressions

Converting an FSA to a regular expression is also easy:



- The general idea: remove (intermediate) states, replacing edge labels with regular expressions

An exercise: simplify the resulting regular expressions

## Wrapping up

- FSA and regular expressions express regular languages
- FSA have two flavors: DFA, NFA (or maybe three:  $\epsilon$ -NFA)
- DFA recognition is linear
- Any NFA can be converted to a DFA (in worst case number of nodes increase exponentially)
- Regular languages and FSA are closed under
  - Concatenation
  - Kleene star
  - Complement
  - Reversal
  - Union
  - Intersection
- Every FSA has a unique minimal DFA

Next:

- Finite state transducers (FSTs)
- Applications of FSA and FSTs

## References / additional reading material

- Hopcroft and Ullman (1979, Ch. 2&3) (and its successive editions) covers (almost) all topics discussed here
- Jurafsky and Martin (2009, Ch. 2)
- Other textbook references include:
  - Sipser (2006)
  - Kozen (2013)

## References / additional reading material (cont.)

Hopcroft, John E. and Jeffrey D. Ullman (1979). *Introduction to Automata Theory, Languages, and Computation*. Addison-Wesley Series in Computer Science and Information Processing. Addison-Wesley. ISBN: 9780201029888.

Jurafsky, Daniel and James H. Martin (2009). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. second. Pearson Prentice Hall. ISBN: 978-0-13-504196-3.

Kozen, Dexter C. (2013). *Automata and Computability*. Undergraduate Texts in Computer Science. Berlin Heidelberg: Springer.

Sipser, Michael (2006). *Introduction to the Theory of Computation*. second. Thomson Course Technology. ISBN: 0-534-95097-3.