

# FSA and regular languages

Data Structures and Algorithms for Computational Linguistics III  
(ISCL-BA-07)

Çağrı Çöltekin  
ccoltekin@ifa.uni-tuebingen.de

University of Tübingen  
Seminar für Sprachwissenschaft

Winter Semester 2025/26

version: 00230a-00230a-00-00

## Three ways to define a regular language

- A language is regular if there is regular grammar that generates/recognizes it
- A language is regular if there is an FSA that generates/recognizes it
- A language is regular if we can define a regular expressions for the language

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 1 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Regular languages: some properties/operations

- $\mathcal{L}_1 \mathcal{L}_2$  Concatenation of two languages  $\mathcal{L}_1$  and  $\mathcal{L}_2$ : any sentence of  $\mathcal{L}_1$  followed by any sentence of  $\mathcal{L}_2$
- $\mathcal{L}^*$  Kleene star of  $\mathcal{L}$ :  $\mathcal{L}$  concatenated with itself 0 or more times
- $\mathcal{L}^R$  Reverse of  $\mathcal{L}$ : reverse of any string in  $\mathcal{L}$
- $\bar{\mathcal{L}}$  Complement of  $\mathcal{L}$ : all strings in  $\Sigma^*$  except the ones in  $\mathcal{L}$  ( $\Sigma^* - \mathcal{L}$ )
- $\mathcal{L}_1 \cup \mathcal{L}_2$  Union of languages  $\mathcal{L}_1$  and  $\mathcal{L}_2$ : strings that are in any of the languages
- $\mathcal{L}_1 \cap \mathcal{L}_2$  Intersection of languages  $\mathcal{L}_1$  and  $\mathcal{L}_2$ : strings that are in both languages

Regular languages are closed under all of these operations.

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 2 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Regular expressions and some extensions

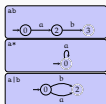
- Kleene star ( $*$ ), concatenation ( $ab$ ) and union ( $a|b$ ) are the basic operations
- Parentheses can be used to group the sub-expressions. Otherwise, the priority of the operators are as listed above:  $a|bc* = a|(b(c*))$
- In practice some short-hand notations are common
  - $\cdot = (a_1) \dots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

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 3 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Converting regular expressions to FSA



- 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

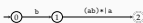
Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 4 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Exercise

convert  $b|(ab)^*|a$  to an NFA



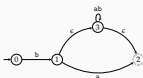
Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 5 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Exercise

convert  $b|(ab)^*|a$  to an NFA



Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 6 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Regular expressions

- Every regular language can be expressed by a regular expression, and every regular expressions defines a regular language
- A regular expression  $e$  defines a regular language  $\mathcal{L}(e)$ 
  - Relations between regular expressions and regular languages
    - $\mathcal{L}(\emptyset) = \emptyset$
    - $\mathcal{L}(a) = \{a\}$
    - $\mathcal{L}(ab) = \mathcal{L}(a)\mathcal{L}(b)$
    - $\mathcal{L}(a^*) = \mathcal{L}(a)^*$
  - where,  $\epsilon$  is the empty string,  $\emptyset$  is the language that accepts nothing (e.g.,  $\Sigma^* - \Sigma^*$ )
- Note: no complement and intersection operators in common regex libraries

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 7 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Some properties of regular expressions

Useful identities for simplifying regular expressions

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

An exercise	
Simplify $a ab^*$	
$a ab^* = ac ab^*$	
$= a(c b^*)$	
$= ab^*$	

Note: some of these are direct statements of Kleene algebra, others can be derived from them.

Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 8 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Exercise

convert  $b|(ab)^*|a$  to an NFA



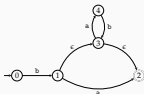
Ç. Çöltekin, INF | University of Tübingen

Winter Semester 2025/26 9 / 27

Recap: regular languages and automata Regular expressions Operations on FSA

## Exercise

convert  $b|(ab)^*|a$  to an NFA



Ç. Çöltekin, INF | University of Tübingen

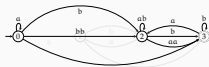
Winter Semester 2025/26 10 / 27

# Converting FSA to regular expressions



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

# Converting FSA to regular expressions



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

# Converting FSA to regular expressions



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

# Converting FSA to regular expressions



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

# Converting FSA to regular expressions



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

# Converting FSA to regular expressions



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

An exercise: simplify the resulting regular expressions

# Two example FSA

what languages do they accept?

$L_1 = \mathcal{L}(M_1)$



Odd number of a's over {a, b}.

$L_2 = \mathcal{L}(M_2)$



Odd number of b's over {a, b}.

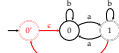
We will use these languages and automata for demonstration.

# Kleene star

$L_1$



$L_1^*$



- What if there were more than one accepting states?

# Reversal

$L_1$



$L_1^R$



# Complement

$L_1$

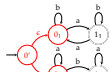


$L_1^c$

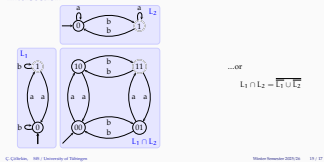


# Union

$L_1 \cup L_2$



## Intersection



## Closure properties of regular languages

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

## Wrapping up

- PSA and regular expressions express regular languages
- Regular languages and PSA are closed under
  - Concatenation
  - Kleene star
  - Complement
  - Reversal
  - Union
  - Intersection
- To prove a language is regular, it is sufficient to find a regular expression or PSA for it
- To prove a language is not regular, we can use pumping lemma (see Appendix)

Next:

- PSTs

## Acknowledgments, credits, references

- The classic reference for PSA, regular languages and regular grammars is Hopcroft and Ullman (1979) (there are recent editions).

- Hopcroft, John E., Rajeev Motwani, and Jeffrey D. Ullman (2007). *Introduction to Automata Theory, Languages, and Computation*. 3rd. Pearson/Addison Wesley. isbn: 9780521462251.
- 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.

## Another exercise on intersection

Construct the intersection of the automata below (adapted from Hopcroft, Motwani, and Ullman (2007), Fig. 4.4)



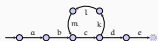
## Is a language regular?

— or not

- To show that a language is regular, it is sufficient to find a PSA 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



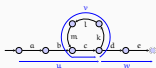
- What is the length of longest string generated by this PSA?
- Any PSA generating an infinite language has to have a loop (application of recursive rule(s) in the grammar)
- Part of every string longer than some number will include repetition of the same substring ('cklm' above)

## 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
  - $uv^i w \in L, (\forall i \geq 0)$
  - $v \neq \epsilon$
  - $|uv| \leq p$
- Pick the string  $a^p b^p$
- For the sake of example, assume  $p = 5, x = aaaaaabbbb$
- Three different ways to split



Pumping Lemma

C. Collaris

SR | University of Tübingen

Week

Winter Semester 2020/21

A.10

Pumping Lemma

C. Collaris

SR | University of Tübingen

Week

Winter Semester 2020/21

A.11