

Compiler Components & Generators

Lexical Analysis

Guido Wachsmuth

Overview

today's lecture

Overview today's lecture

lexical analysis

Overview today's lecture

lexical analysis

regular languages

- regular grammars
- regular expressions
- finite state automata

Overview today's lecture

lexical analysis

regular languages

- regular grammars
- regular expressions
- finite state automata

equivalence of formalisms

- constructive approach

Overview today's lecture

lexical analysis

regular languages

- regular grammars
- regular expressions
- finite state automata

equivalence of formalisms

- constructive approach

tool generation

I

regular grammars

Recap: A Theory of Language

formal languages



Recap: A Theory of Language

formal languages

vocabulary Σ

finite, nonempty set of elements (words, letters)

alphabet



Recap: A Theory of Language

formal languages

vocabulary Σ

finite, nonempty set of elements (words, letters)

alphabet

string over Σ

finite sequence of elements chosen from Σ

word, sentence, utterance



Recap: A Theory of Language

formal languages

vocabulary Σ

finite, nonempty set of elements (words, letters)

alphabet

string over Σ

finite sequence of elements chosen from Σ

word, sentence, utterance

formal language λ

set of strings over a vocabulary Σ

$\lambda \subseteq \Sigma^*$



Recap: A Theory of Language formal grammars



Recap: A Theory of Language formal grammars

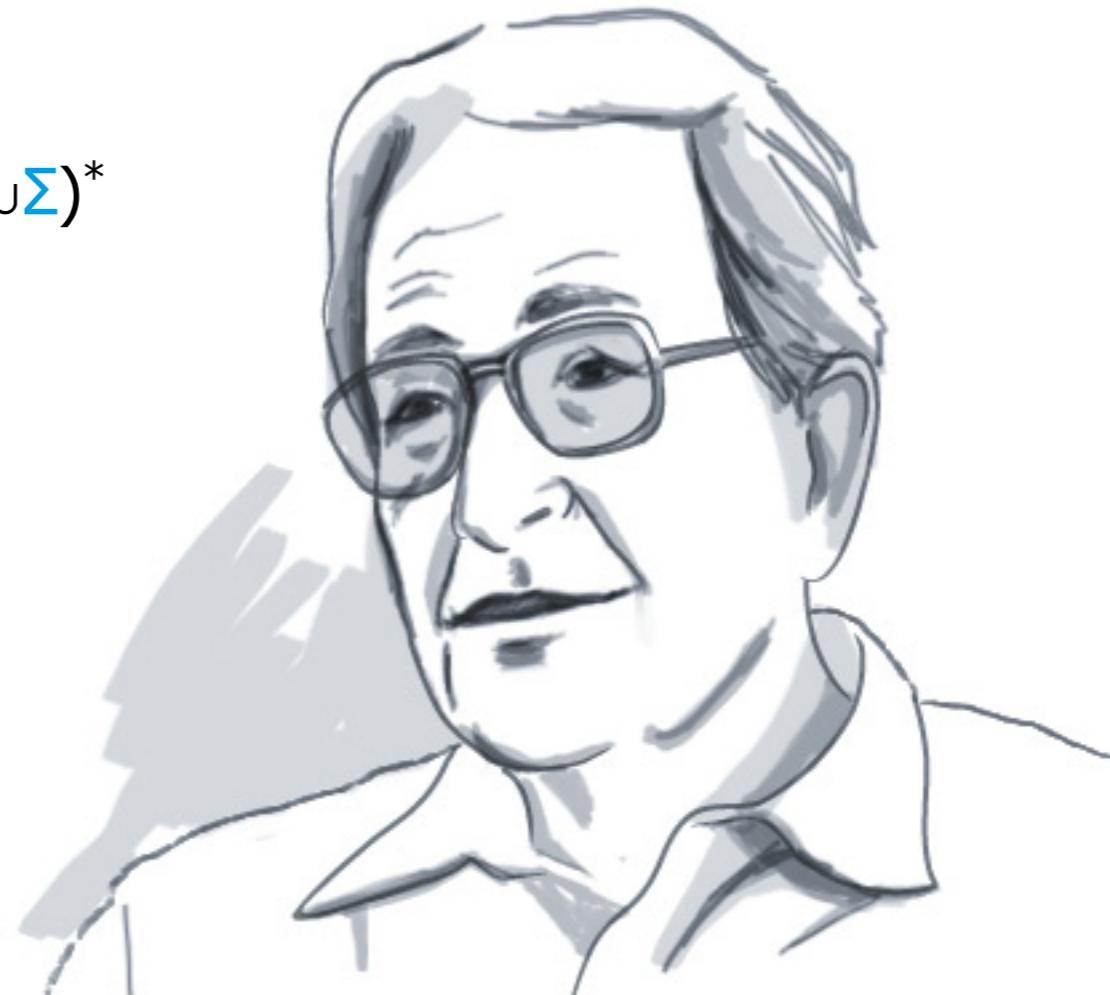
formal grammar $G = (N, \Sigma, P, S)$

nonterminal symbols N

terminal symbols Σ

production rules $P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

start symbol $S \in N$



Recap: A Theory of Language formal grammars

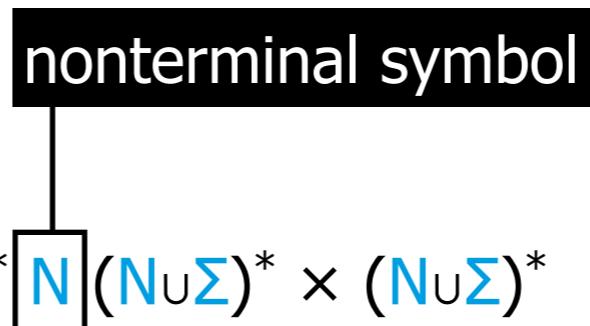
formal grammar $G = (N, \Sigma, P, S)$

nonterminal symbols N

terminal symbols Σ

production rules $P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

start symbol $S \in N$



Recap: A Theory of Language formal grammars

formal grammar $G = (N, \Sigma, P, S)$

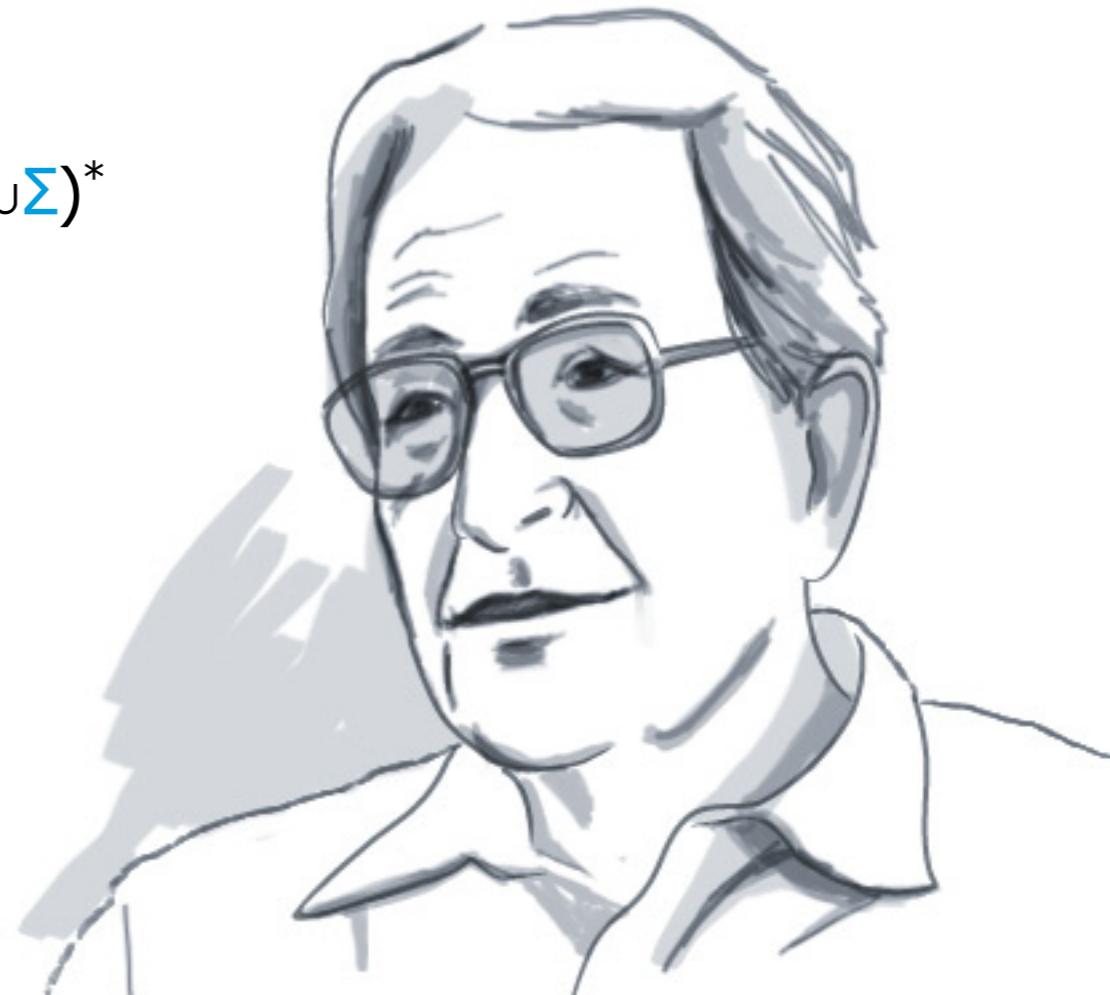
nonterminal symbols N

terminal symbols Σ

production rules $P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

start symbol $S \in N$

context



Recap: A Theory of Language formal grammars

formal grammar $G = (N, \Sigma, P, S)$

nonterminal symbols N

terminal symbols Σ

production rules $P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

start symbol $S \in N$

replacement



Recap: A Theory of Language formal grammars

formal grammar $G = (N, \Sigma, P, S)$

nonterminal symbols N

terminal symbols Σ

production rules $P \subseteq (N \cup \Sigma)^* N (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

start symbol $S \in N$

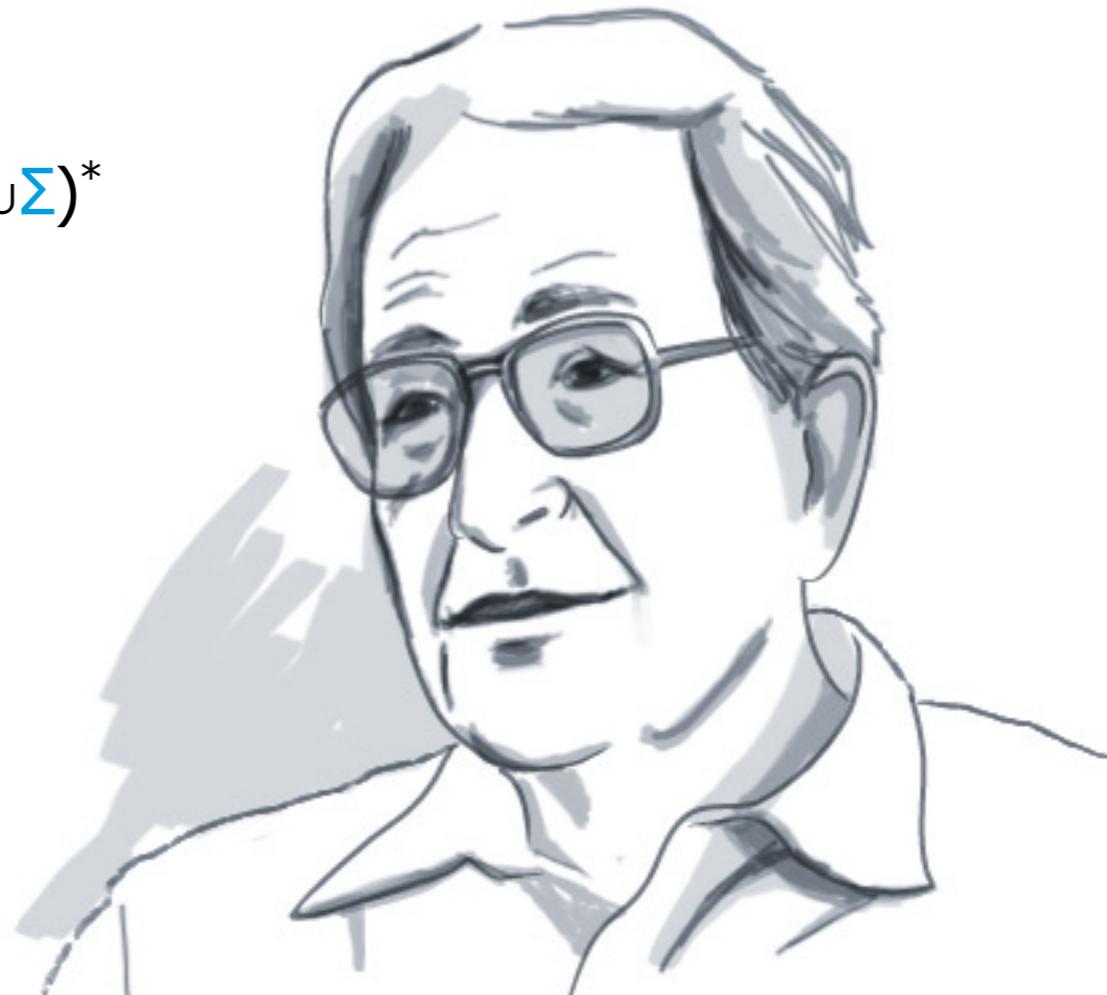
grammar classes

type-0, unrestricted

type-1, context-sensitive: $(a A c, a b c)$

type-2, context-free: $P \subseteq N \times (N \cup \Sigma)^*$

type-3, regular: (A, x) or (A, xB)



Decimal Numbers

right regular grammar

Num → "0" Num
Num → "1" Num
Num → "2" Num
Num → "3" Num
Num → "4" Num
Num → "5" Num
Num → "6" Num
Num → "7" Num
Num → "8" Num
Num → "9" Num

Num → "0"
Num → "1"
Num → "2"
Num → "3"
Num → "4"
Num → "5"
Num → "6"
Num → "7"
Num → "8"
Num → "9"



Identifiers

right regular grammar

$\text{Id} \rightarrow "a" \text{ R}$

...

$\text{Id} \rightarrow "z" \text{ R}$

$\text{R} \rightarrow "a" \text{ R}$

...

$\text{R} \rightarrow "z" \text{ R}$

$\text{R} \rightarrow "0" \text{ R}$

...

$\text{R} \rightarrow "9" \text{ R}$

$\text{Id} \rightarrow "a"$

...

$\text{Id} \rightarrow "z"$

$\text{R} \rightarrow "a"$

...

$\text{R} \rightarrow "z"$

$\text{R} \rightarrow "0"$

...

$\text{R} \rightarrow "9"$



Recap: A Theory of Language

formal languages



Recap: A Theory of Language

formal languages

formal grammar $G = (N, \Sigma, P, S)$



Recap: A Theory of Language

formal languages

formal grammar $G = (N, \Sigma, P, S)$

derivation relation $\Rightarrow_G \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

$w \Rightarrow_G w' \Leftrightarrow$

$\exists (p, q) \in P : \exists u, v \in (N \cup \Sigma)^* :$

$w = u \ p \ v \wedge w' = u \ q \ v$



Recap: A Theory of Language

formal languages

formal grammar $G = (N, \Sigma, P, S)$

derivation relation $\Rightarrow_G \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

$w \Rightarrow_G w' \Leftrightarrow$

$\exists (p, q) \in P: \exists u, v \in (N \cup \Sigma)^*:$

$w = u p v \wedge w' = u q v$

formal language $L(G) \subseteq \Sigma^*$

$L(G) = \{w \in \Sigma^* \mid S \xrightarrow{G}^* w\}$



Recap: A Theory of Language

formal languages

formal grammar $G = (N, \Sigma, P, S)$

derivation relation $\Rightarrow_G \subseteq (N \cup \Sigma)^* \times (N \cup \Sigma)^*$

$w \Rightarrow_G w' \Leftrightarrow$

$\exists (p, q) \in P: \exists u, v \in (N \cup \Sigma)^*:$

$w = u p v \wedge w' = u q v$

formal language $L(G) \subseteq \Sigma^*$

$L(G) = \{w \in \Sigma^* \mid S \xrightarrow{G}^* w\}$

classes of formal languages



II

regular expressions

Recap: Regular Expressions

overview

basics

- symbol from an alphabet
- ϵ

combinators

- alternation: $E_1 \mid E_2$
- concatenation: $E_1 E_2$
- repetition: E^*
- optional: $E? = E \mid \epsilon$
- one or more: $E+ = E E^*$

Decimal Numbers & Identifiers

regular expressions

Num: $(0|1|2|3|4|5|6|7|8|9)^+$

Id: $(a|...|z)(a|...|z|0|...|9)^*$



Regular Expressions

formal languages

basics

- $L(a) = \{"a"\}$
- $L(\epsilon) = \{''\}$

combinators

- $L(E_1 \mid E_2) = L(E_1) \cup L(E_2)$
- $L(E_1 E_2) = L(E_1) \cdot L(E_2)$
- $L(E^*) = L(E)^*$

III

finite automata

Finite Automata

formal definition

Finite Automata

formal definition

finite automaton $M = (Q, \Sigma, T, q_0, F)$

states Q

input symbols Σ

transition function T

start state $q_0 \in Q$

final states $F \subseteq Q$

Finite Automata

formal definition

finite automaton $M = (Q, \Sigma, T, q_0, F)$

states Q

input symbols Σ

transition function T

start state $q_0 \in Q$

final states $F \subseteq Q$

transition function

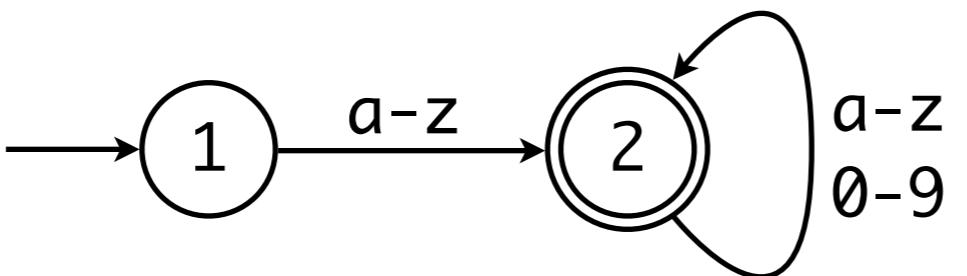
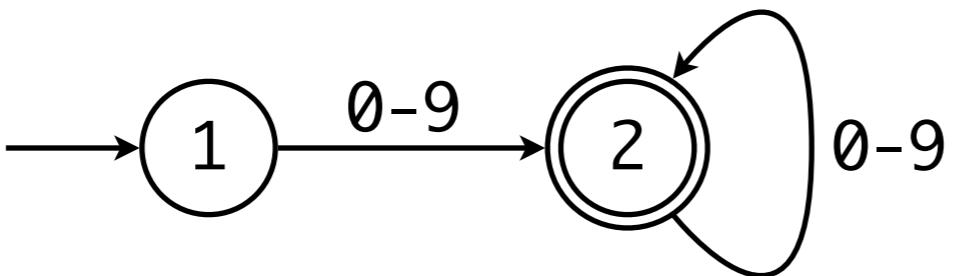
nondeterministic FA $T : Q \times \Sigma \rightarrow P(Q)$

NFA with ϵ -moves $T : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow P(Q)$

deterministic FA $T : Q \times \Sigma \rightarrow Q$

Decimal Numbers & Identifiers

finite automata



Nondeterministic Finite Automata

formal languages

Nondeterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

Nondeterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

transition function $T : Q \times \Sigma \rightarrow P(Q)$

$T(\{q_1, \dots, q_n\}, x) := T(q_1, x) \cup \dots \cup T(q_n, x)$

$T^*(\{q_1, \dots, q_n\}, \varepsilon) := \{q_1, \dots, q_n\}$

$T^*(\{q_1, \dots, q_n\}, xw) := T^*(T(\{q_1, \dots, q_n\}, x), w)$

Nondeterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

transition function $T : Q \times \Sigma \rightarrow P(Q)$

$T(\{q_1, \dots, q_n\}, x) := T(q_1, x) \cup \dots \cup T(q_n, x)$

$T^*(\{q_1, \dots, q_n\}, \varepsilon) := \{q_1, \dots, q_n\}$

$T^*(\{q_1, \dots, q_n\}, xw) := T^*(T(\{q_1, \dots, q_n\}, x), w)$

formal language $L(M) \subseteq \Sigma^*$

$L(M) = \{w \in \Sigma^* \mid T^*(\{q_0\}, w) \cap F \neq \emptyset\}$

Deterministic Finite Automata

formal languages

Deterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

Deterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

transition function $T : Q \times \Sigma \rightarrow P(Q)$

$$T^*(q, \varepsilon) := q$$

$$T^*(q, xw) := T^*(T(q, x), w)$$

Deterministic Finite Automata

formal languages

finite automaton $M = (Q, \Sigma, T, q_0, F)$

transition function $T : Q \times \Sigma \rightarrow P(Q)$

$$T^*(q, \varepsilon) := q$$

$$T^*(q, xw) := T^*(T(q, x), w)$$

formal language $L(M) \subseteq \Sigma^*$

$$L(M) = \{w \in \Sigma^* \mid T^*(q_0, w) \in F\}$$

coffee break



IV

equivalence

Regular Languages

formalisms

left regular
grammars

right regular
grammars

regular
expressions

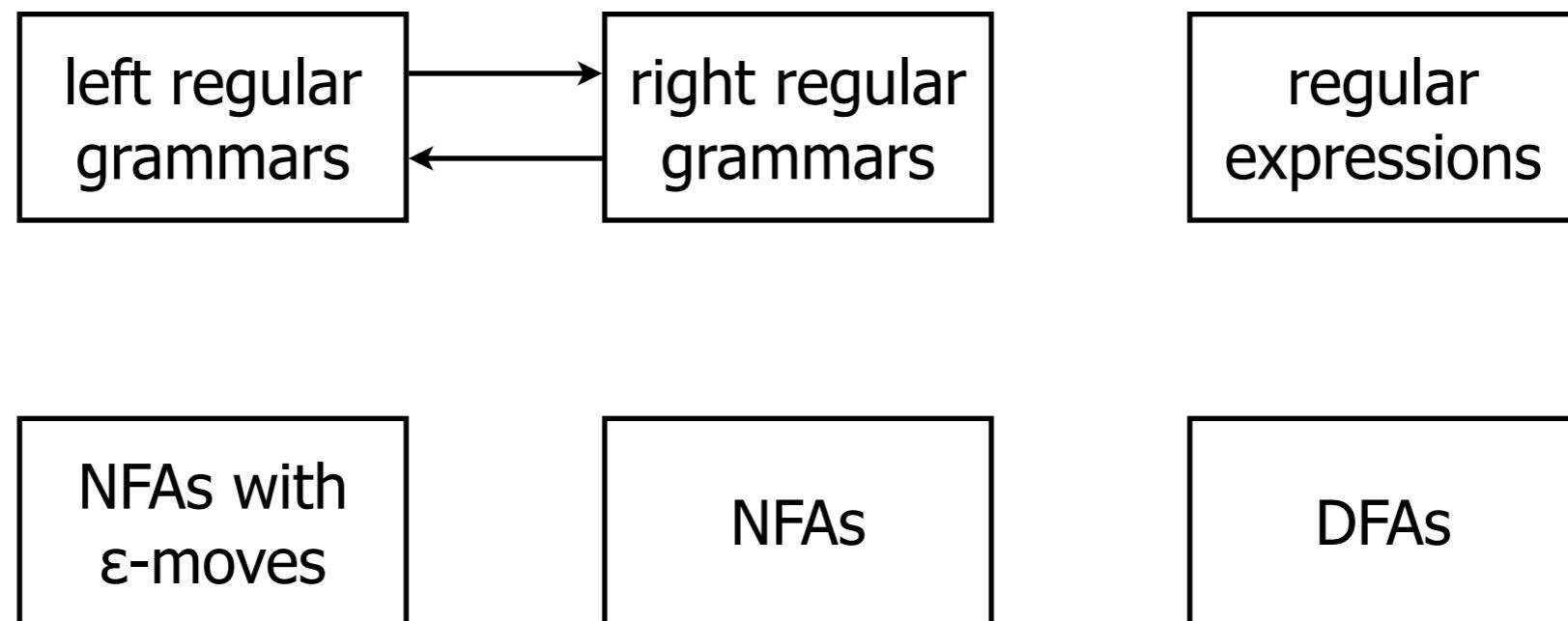
NFAs with
 ϵ -moves

NFAs

DFAs

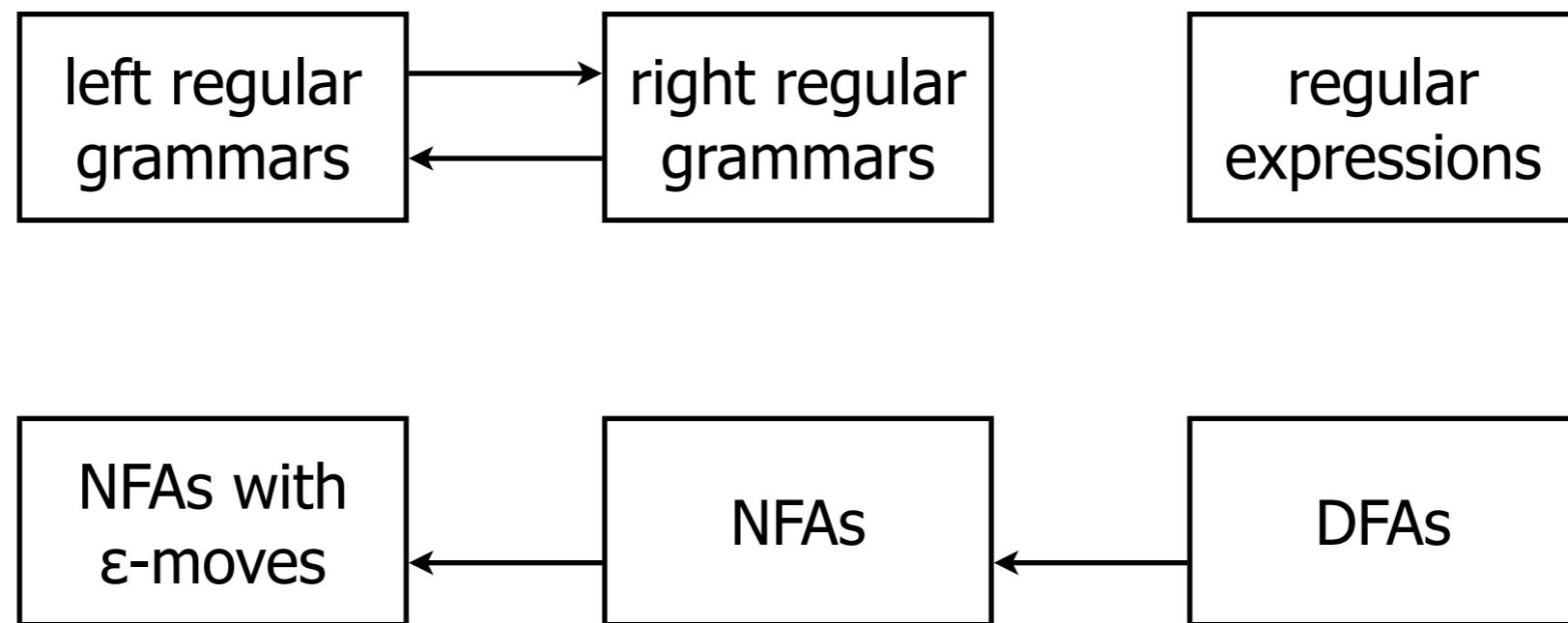
Regular Languages

formalisms



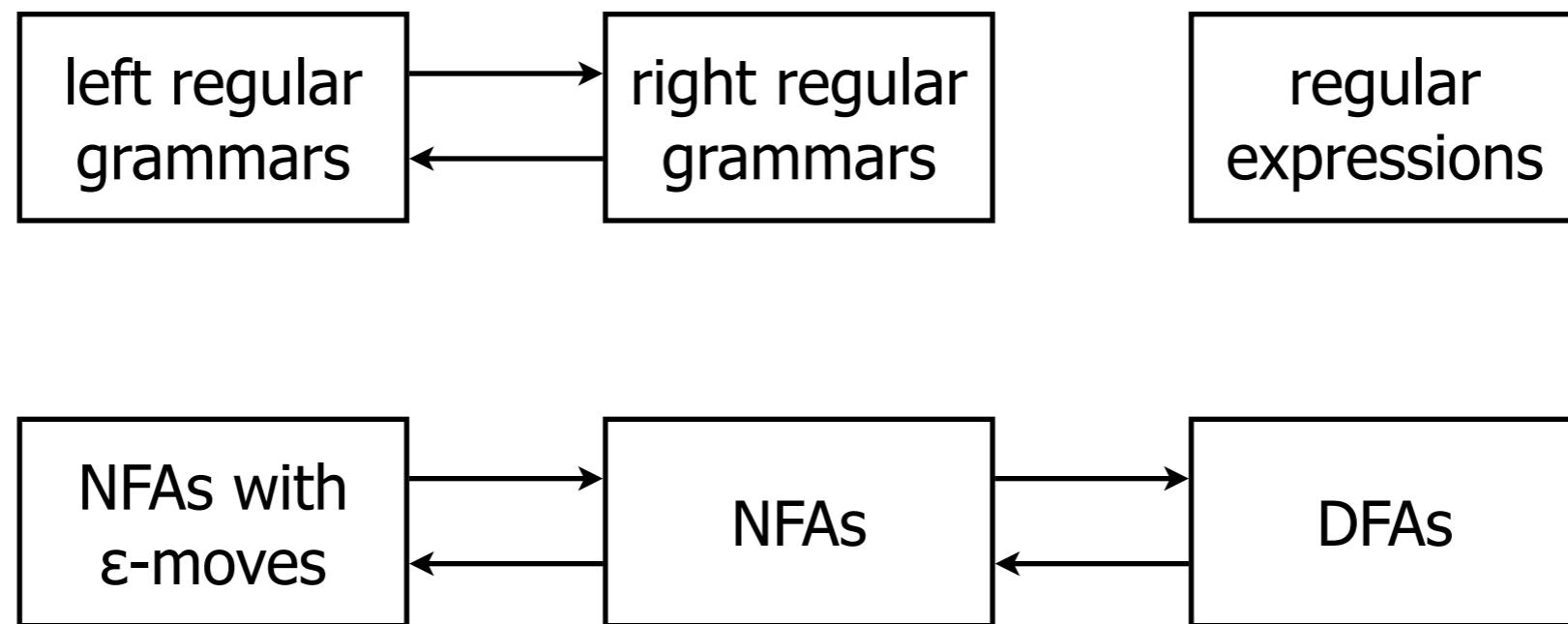
Regular Languages

formalisms



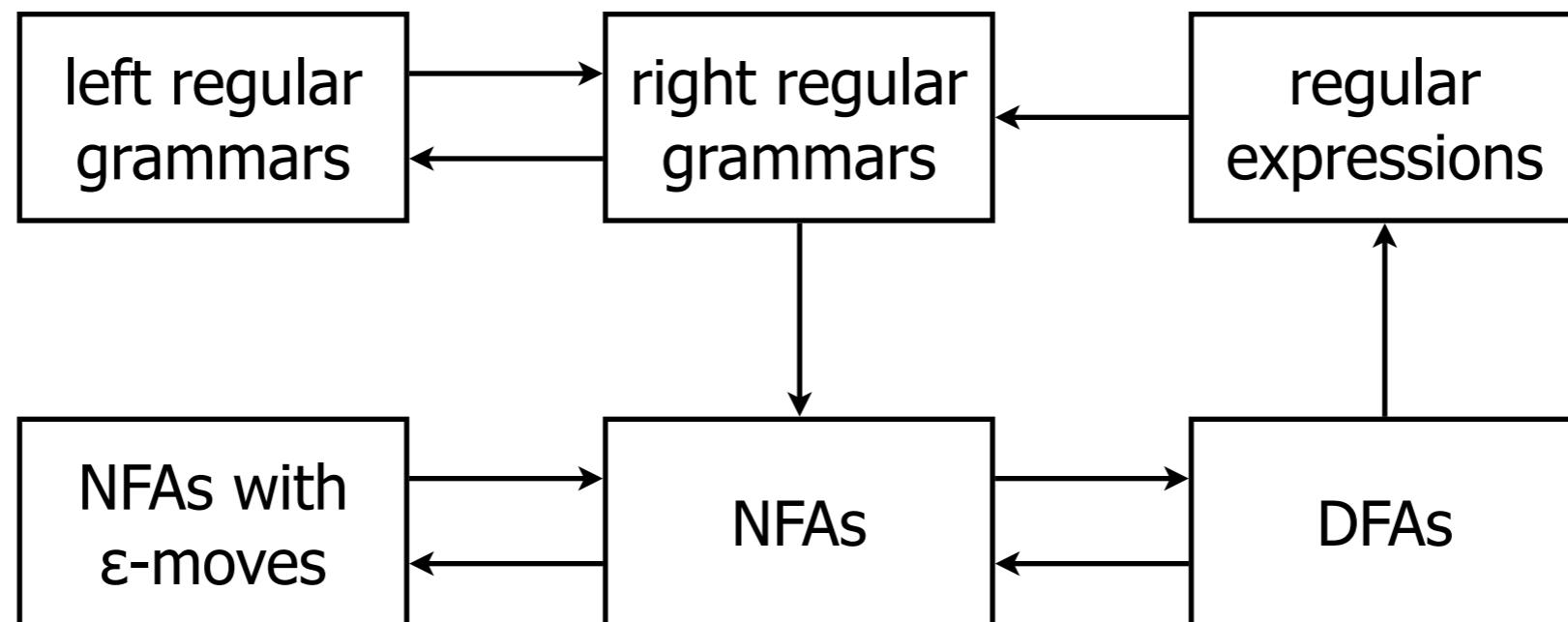
Regular Languages

formalisms



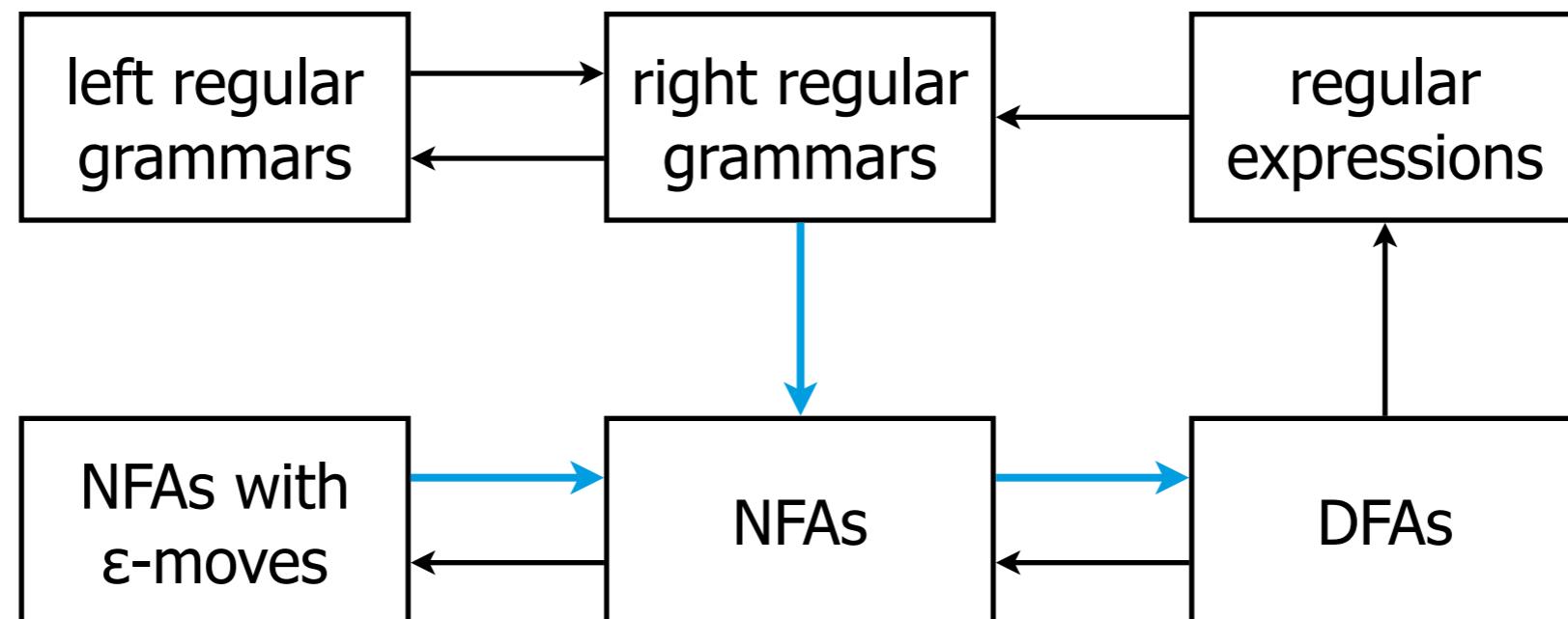
Regular Languages

formalisms



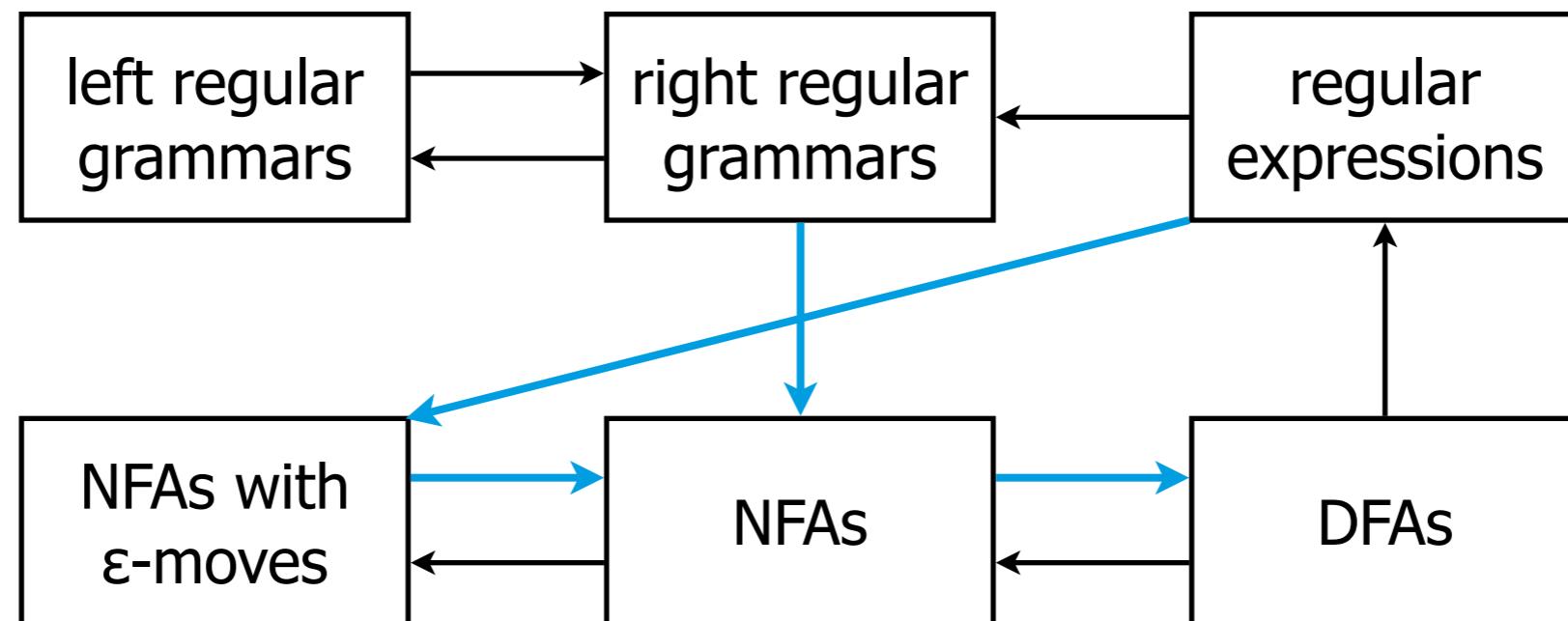
Regular Languages

formalisms



Regular Languages

formalisms



NFA construction

right regular grammar

NFA construction right regular grammar

formal grammar $G = (N, \Sigma, P, S)$

NFA construction right regular grammar

formal grammar $G = (N, \Sigma, P, S)$

finite automaton $M = (N \cup \{f\}, \Sigma, T, S, F)$

NFA construction right regular grammar

formal grammar $G = (N, \Sigma, P, S)$

finite automaton $M = (N \cup \{f\}, \Sigma, T, S, F)$

transition function T

$(X, aY) \in P : (X, a, Y) \in T$

$(X, a) \in P : (X, a, f) \in T$

NFA construction right regular grammar

formal grammar $G = (N, \Sigma, P, S)$

finite automaton $M = (N \cup \{f\}, \Sigma, T, S, F)$

transition function T

$(X, aY) \in P : (X, a, Y) \in T$

$(X, a) \in P : (X, a, f) \in T$

final states F

$(S, \epsilon) \in P : F = \{S, f\}$

else: $F = \{f\}$

NFA construction example

Num → "0" Num

Num → "1" Num

Num → "2" Num

Num → "3" Num

Num → "4" Num

Num → "5" Num

Num → "6" Num

Num → "7" Num

Num → "8" Num

Num → "9" Num

Num → "0"

Num → "1"

Num → "2"

Num → "3"

Num → "4"

Num → "5"

Num → "6"

Num → "7"

Num → "8"

Num → "9"

NFA construction example

Num → "0" Num
Num → "1" Num
Num → "2" Num
Num → "3" Num
Num → "4" Num
Num → "5" Num
Num → "6" Num
Num → "7" Num
Num → "8" Num
Num → "9" Num

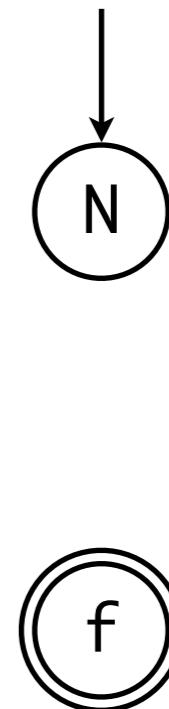
Num → "0"
Num → "1"
Num → "2"
Num → "3"
Num → "4"
Num → "5"
Num → "6"
Num → "7"
Num → "8"
Num → "9"



NFA construction example

Num → "0" Num
Num → "1" Num
Num → "2" Num
Num → "3" Num
Num → "4" Num
Num → "5" Num
Num → "6" Num
Num → "7" Num
Num → "8" Num
Num → "9" Num

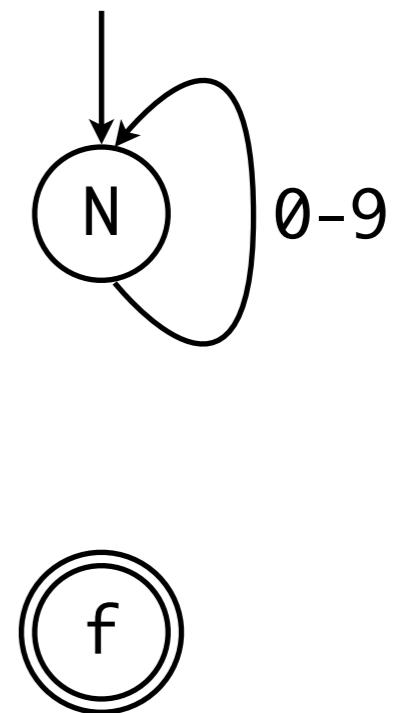
Num → "0"
Num → "1"
Num → "2"
Num → "3"
Num → "4"
Num → "5"
Num → "6"
Num → "7"
Num → "8"
Num → "9"



NFA construction example

Num → "0" Num
Num → "1" Num
Num → "2" Num
Num → "3" Num
Num → "4" Num
Num → "5" Num
Num → "6" Num
Num → "7" Num
Num → "8" Num
Num → "9" Num

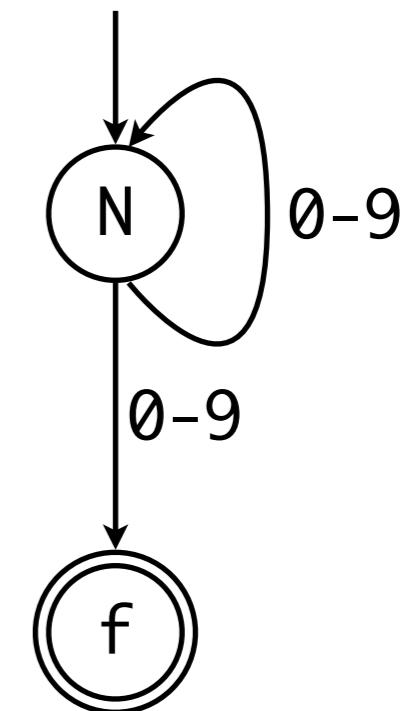
Num → "0"
Num → "1"
Num → "2"
Num → "3"
Num → "4"
Num → "5"
Num → "6"
Num → "7"
Num → "8"
Num → "9"



NFA construction example

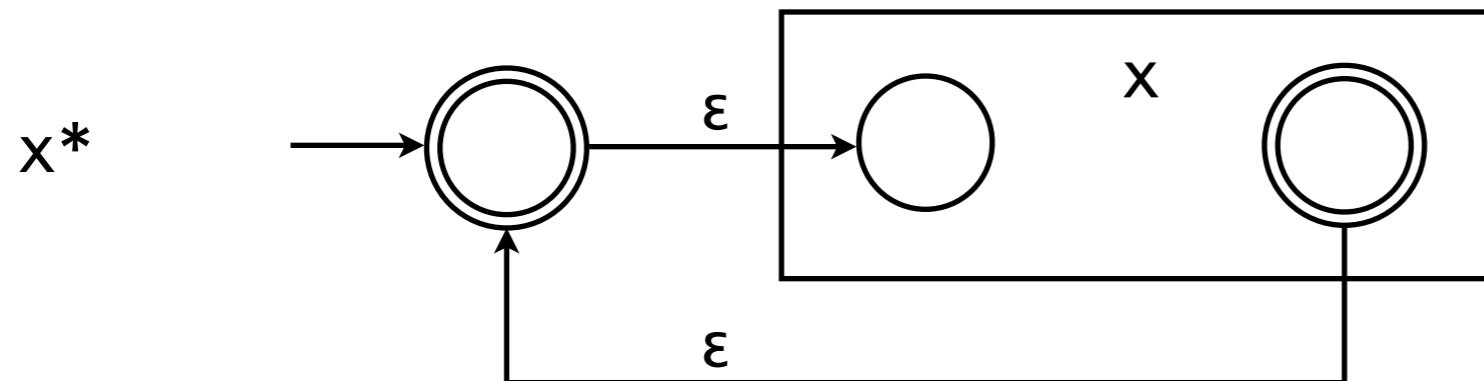
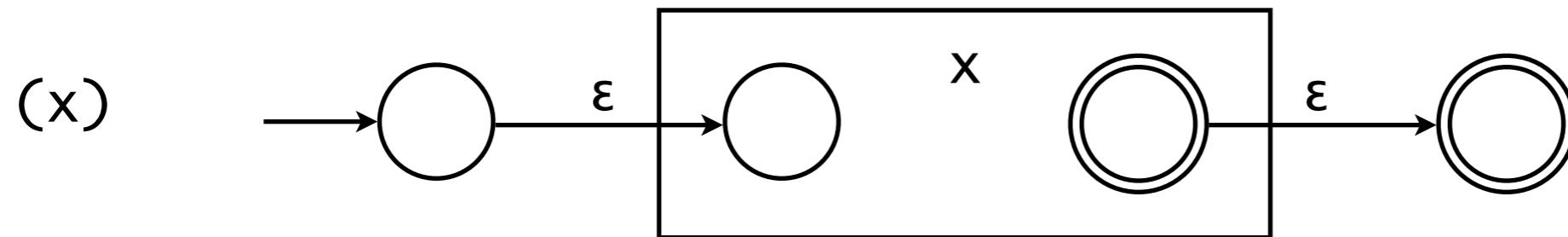
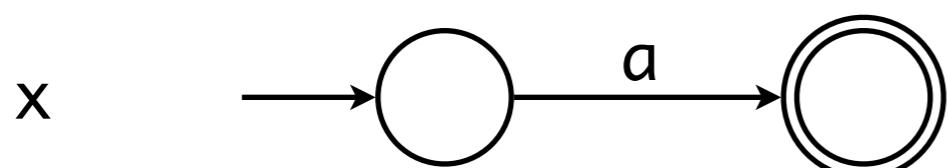
Num → "0" Num
Num → "1" Num
Num → "2" Num
Num → "3" Num
Num → "4" Num
Num → "5" Num
Num → "6" Num
Num → "7" Num
Num → "8" Num
Num → "9" Num

Num → "0"
Num → "1"
Num → "2"
Num → "3"
Num → "4"
Num → "5"
Num → "6"
Num → "7"
Num → "8"
Num → "9"



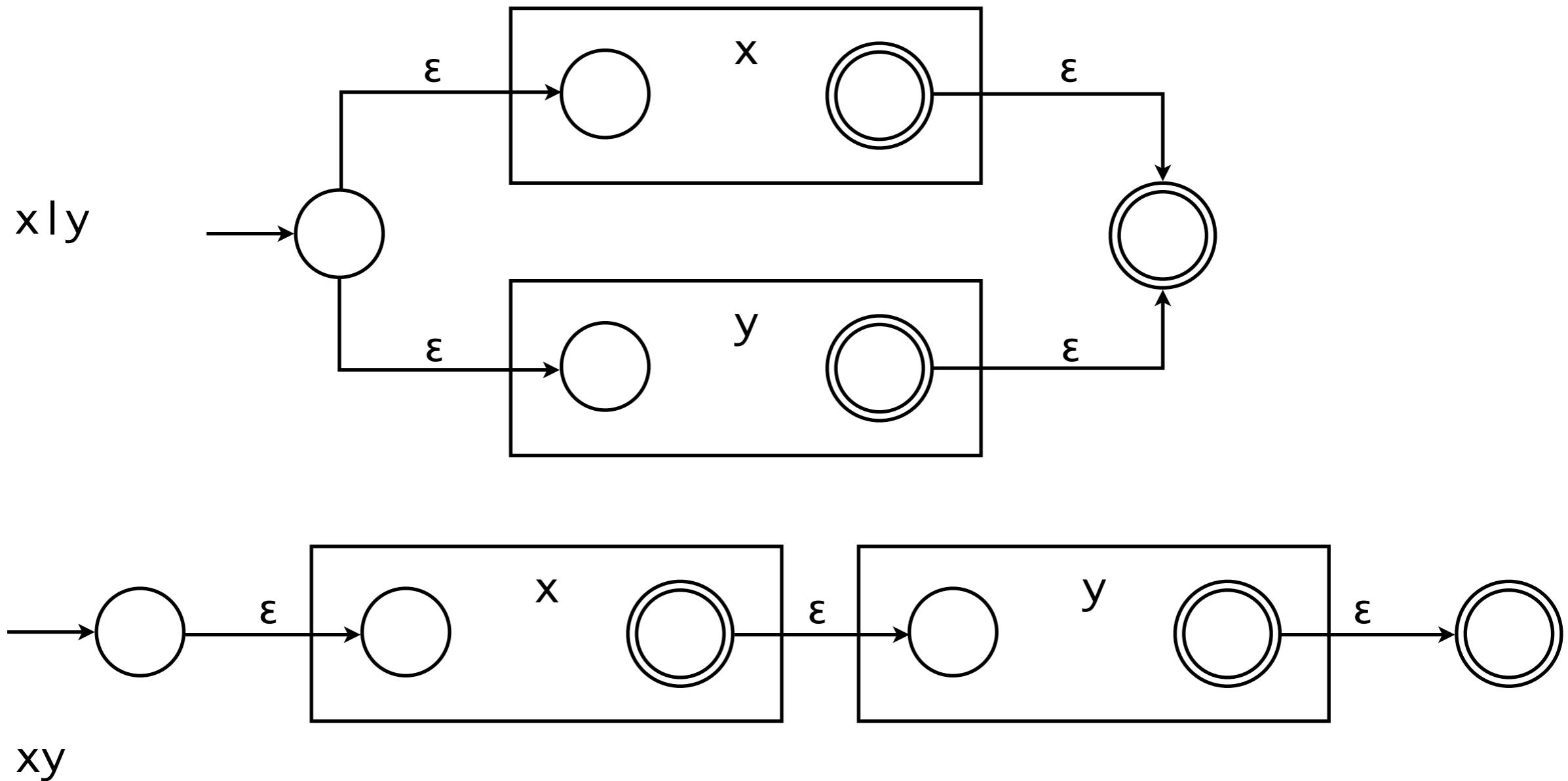
NFA construction

regular expressions



NFA construction

regular expressions



NFA construction

ϵ elimination

additional final states

- states with ϵ -moves into final states
- become final states themselves

additional transitions

- ϵ -move from source to target state
- transitions from target state
- add these transitions to the source state

Powerset construction eliminating nondeterminism

nondeterministic finite automaton $M = (Q, \Sigma, T, q_0, F)$

deterministic finite automaton $M' = (P(Q), \Sigma, T', \{q_0\}, F')$

transition function T'

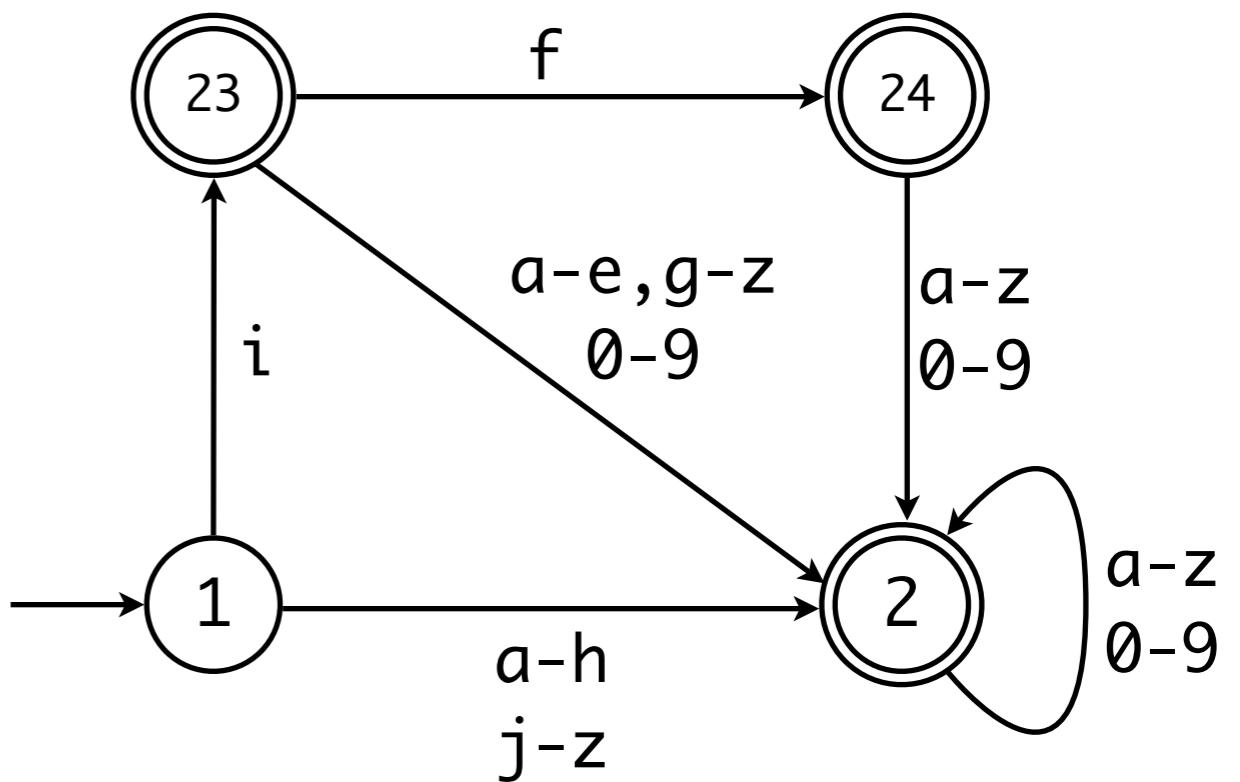
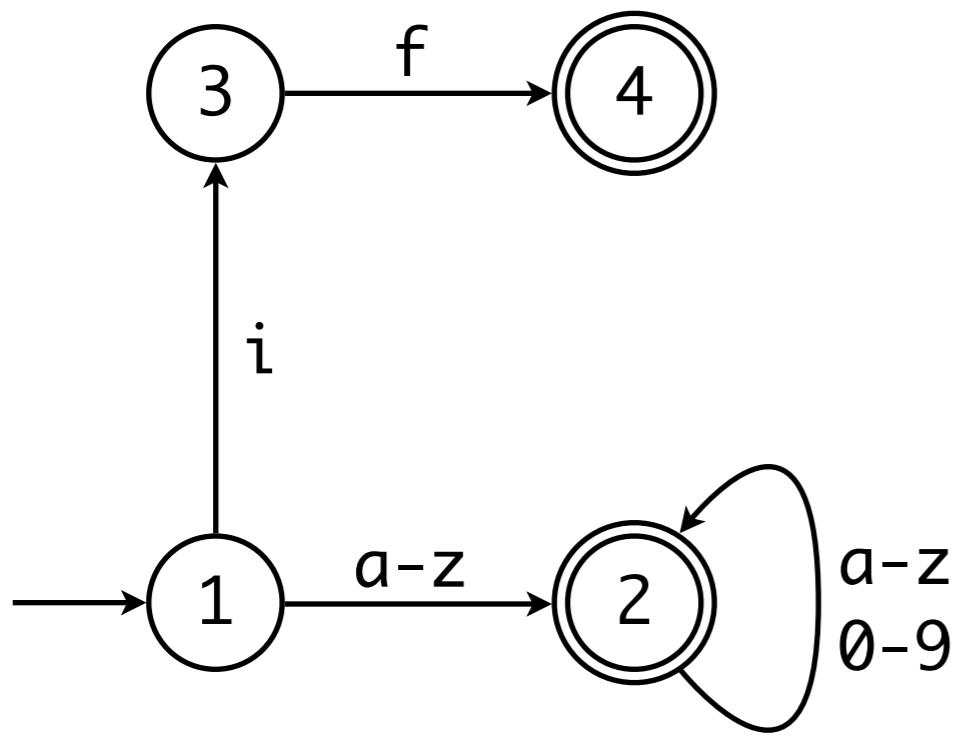
- $T'(\{q_1, \dots, q_n\}, x) = T(\{q_1, \dots, q_n\}, x) = T(q_1, x) \cup \dots \cup T(q_n, x)$

final states $F' = \{S \mid S \subseteq Q, S \cap F \neq \emptyset\}$

- all states that include a final state of the original NFA

Powerset construction

example



V

summary

Summary

lessons learned

Summary

lessons learned

What are the formalisms to describe regular languages?

- regular grammars
- regular expressions
- finite state automata

Summary lessons learned

What are the formalisms to describe regular languages?

- regular grammars
- regular expressions
- finite state automata

Why are these formalisms equivalent?

- constructive proofs

Summary lessons learned

What are the formalisms to describe regular languages?

- regular grammars
- regular expressions
- finite state automata

Why are these formalisms equivalent?

- constructive proofs

How can we generate compiler tools from that?

- implement DFAs
- generate transition tables

Literature

learn more

Literature

[learn more](#)

formal languages

Noam Chomsky: Three models for the description of language. 1956

J. E. Hopcroft, R. Motwani, J. D. Ullman: Introduction to Automata Theory, Languages, and Computation. 2006

Literature

[learn more](#)

formal languages

Noam Chomsky: Three models for the description of language. 1956

J. E. Hopcroft, R. Motwani, J. D. Ullman: Introduction to Automata Theory, Languages, and Computation. 2006

lexical analysis

Andrew W. Appel, Jens Palsberg: Modern Compiler Implementation in Java, 2nd edition. 2002

Outlook

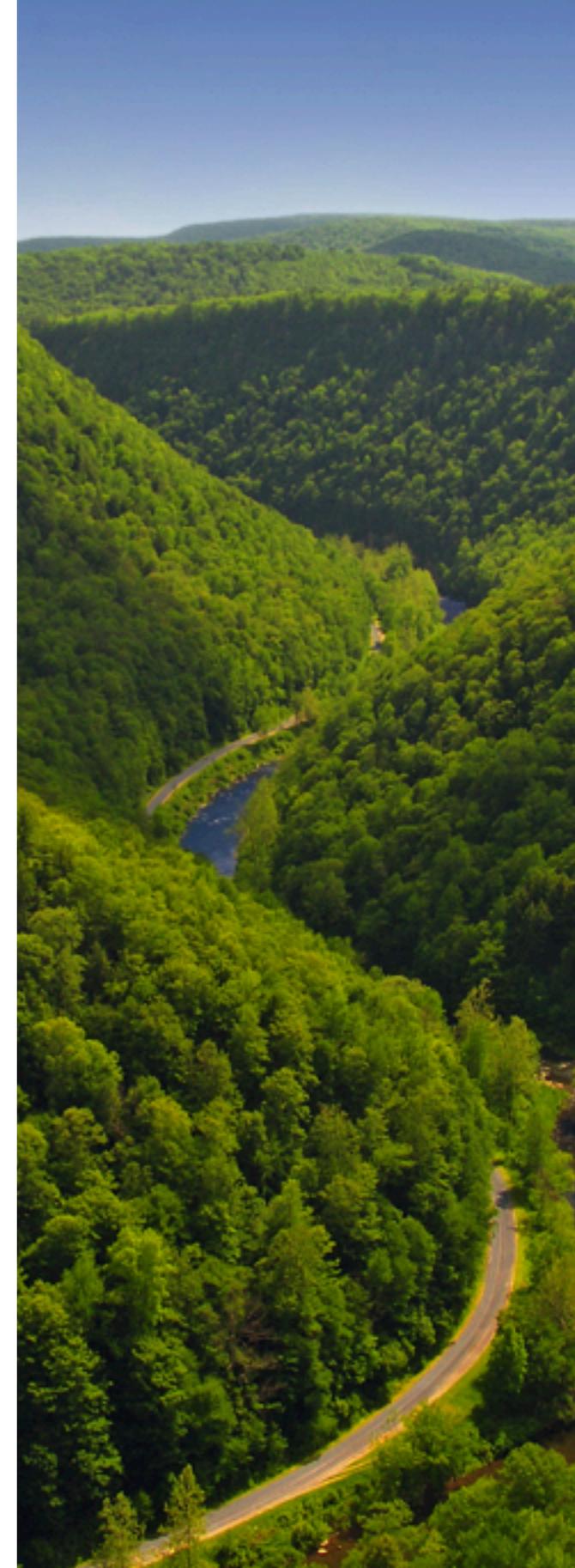
coming next

compiler components and their generators

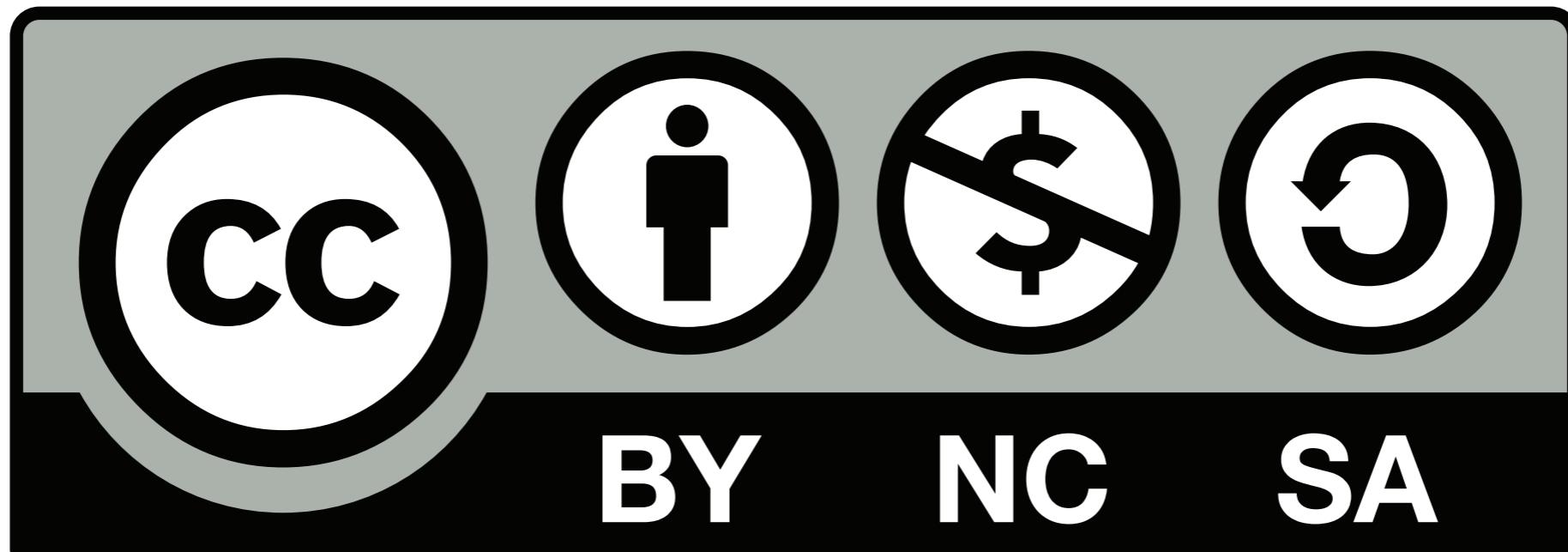
- Lecture 12: LL parsing
- Lecture 13: LR parsing
- guest lecture Eelco Visser
- exam preparation lecture

Lab Nov 30

- Java & Java Bytecode
- Jasmin editor
- initial code generator



copyrights



Pictures copyrights

Slide 1:

Book Scanner by Ben Woosley, some rights reserved

Slides 4, 5, 8:

Noam Chomsky by Fellowsisters, some rights reserved

Slide 6, 7, 11, 15:

Tiger by Bernard Landgraf, some rights reserved

Slide 18:

Coffee Primo by Dominica Williamson, some rights reserved

Slide 32:

Pine Creek by Nicholas, some rights reserved