# Talen en Compilers

## 2009/2010, periode 2

Andres Löh

Department of Information and Computing Sciences
Utrecht University

December 17, 2009

# 10. Regular languages

# This lecture

## Regular languages

Regular languages

Finite state automata

NFAs vs. DFAs

Regular grammars

Regular grammars vs. finite state automata

Universiteit Utrecht

# 10.1  Regular languages

# Context-free languages

The languages we dealt with until now were mostly **context-free** languages:

- ▶ can be described using a context-free grammar,
- ▶ can be parsed relatively easily (for instance, using parser combinators),
- ▶ resulting parsers need polynomial time and space (often not much worse than linear).

Universiteit Utrecht

# Context-free languages

The languages we dealt with until now were mostly **context-free** languages:

- ▶ can be described using a context-free grammar,
- ▶ can be parsed relatively easily (for instance, using parser combinators),
- ▶ resulting parsers need polynomial time and space (often not much worse than linear).

The rest of the course: classes of languages and/or grammars that allow more efficient parsing.

Universiteit Utrecht

# Regular languages

A proper subset of the context-free languages:

- ► can be described using finite state automata,
- ► can be described using regular grammars,
- ► can be described using regular expressions,
- ► can be parsed very easily, in linear time and constant space.

Universiteit Utrecht

# Regular languages

A proper subset of the context-free languages:

- ▶ can be described using finite state automata,
- ▶ can be described using regular grammars,
- ▶ can be described using regular expressions,
- ▶ can be parsed very easily, in linear time and constant space.

We will look at the different formalisms, their respective advantages and disadvantages, and show their equivalence.
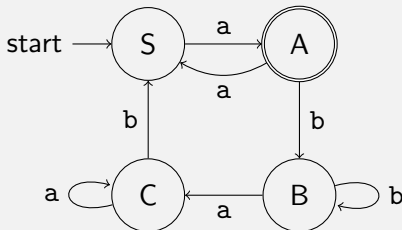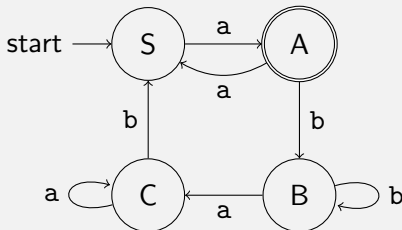
Universiteit Utrecht

# 10.2 Finite state automata
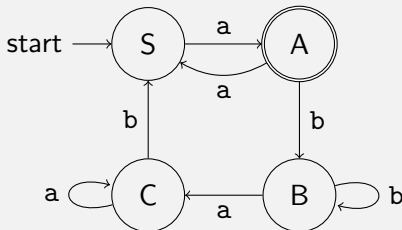
# Deterministic finite state automata (DFA)

Universiteit Utrecht

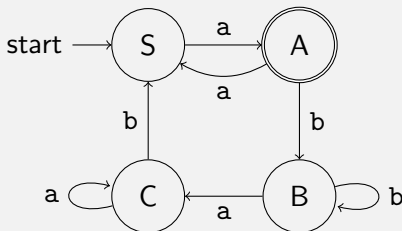# Deterministic finite state automata (DFA)

▶ Input alphabet:
$X = \{a, b\}$

Universiteit Utrecht

# Deterministic finite state automata (DFA)

▶ Input alphabet:
$X = \{a, b\}$

▶ States:
$Q = \{S, A, B, C\}$

**Universiteit Utrecht**

# Deterministic finite state automata (DFA)



- Input alphabet:
  $X = \{\, a, b \,\}$
- States:
  $Q = \{\, S, A, B, C \,\}$
- Transitions:
  $d :: Q \rightarrow X \rightarrow Q$

Universiteit Utrecht

# Deterministic finite state automata (DFA)



- ▶ Input alphabet:
  $X = \{\, a, b \,\}$
- ▶ States:
  $Q = \{\, S, A, B, C \,\}$
- ▶ Transitions:
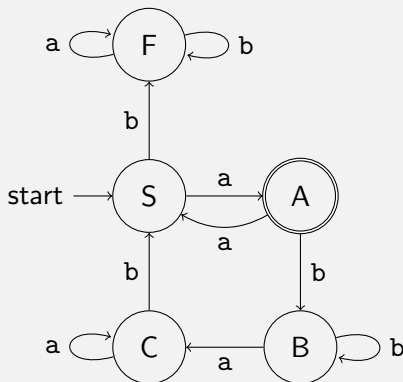  $d :: Q \rightarrow X \rightarrow Q$
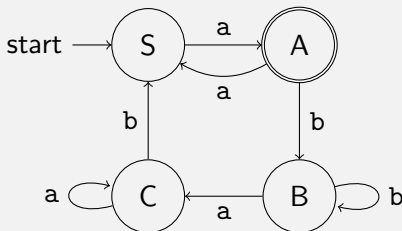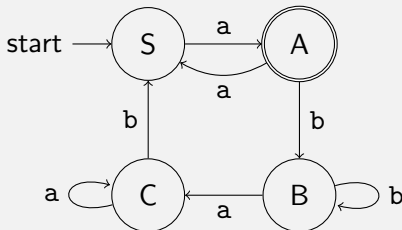
# Deterministic finite state automata (DFA)



- ▶ Input alphabet:
  $X = \{a, b\}$
- ▶ States:
  $Q = \{S, A, B, C\}$
- ▶ Transitions:
  $d :: Q \rightarrow X \rightarrow Q$

Universiteit Utrecht

# Deterministic finite state automata (DFA)



- ▶ Input alphabet:
  $X = \{a, b\}$
- ▶ States:
  $Q = \{S, A, B, C\}$
- ▶ Transitions:
  $d :: Q \to X \to Q$
- ▶ Start state:
  $S$ (where $S \in Q$)
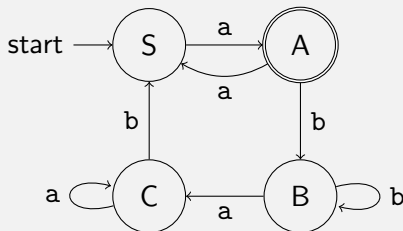
Universiteit Utrecht

# Deterministic finite state automata (DFA)



- ▶ Input alphabet:
  $X = \{a, b\}$
- ▶ States:
  $Q = \{S, A, B, C\}$
- ▶ Transitions:
  $d :: Q \to X \to Q$
- ▶ Start state:
  S (where $S \in Q$)
- ▶ Accepting states:
  $F = \{A\}$
  (where $F \subseteq Q$)

Universiteit Utrecht

# Definition of a DFA

A DFA is given by

- an input alphabet $X$,
- a set of states $Q$,
- a transition function d of type $Q \rightarrow X \rightarrow Q$,
- a start state $S \in Q$,
- a set of accepting states $F \subseteq Q$.

Universiteit Utrecht

# Definition of a DFA

A DFA is given by

- an input alphabet $X$,
- a set of states $Q$,
- a transition function $d$ of type $Q \to X \to Q$,
- a start state $S \in Q$,
- a set of accepting states $F \subseteq Q$.

Often, a DFA is simply written as a tuple $(X, Q, d, S, F)$.

Universiteit Utrecht

# Definition of a DFA

A DFA is given by

- an input alphabet X,
- a set of states Q,
- a transition function d of type $Q \rightarrow X \rightarrow Q$,
- a start state $S \in Q$,
- a set of accepting states $F \subseteq Q$.

Often, a DFA is simply written as a tuple $(X, Q, d, S, F)$.

Sometimes, when X and Q are clear from the context, $(d, S, F)$ is sufficient to specify a DFA.

# Running a DFA

dfa :: (Q → X → Q) → Q → [X] → Q
dfa d q []       = q
dfa d q (x : xs) = dfa d (d q x) xs

# Running a DFA

```
dfa :: (Q → X → Q) → Q → [X] → Q
dfa d q []       = q
dfa d q (x : xs) = dfa d (d q x) xs
```

## Question

Does this function look familiar?

Universiteit Utrecht

# Running a DFA

$$\text{dfa} :: (Q \to X \to Q) \to Q \to [X] \to Q$$
$$\text{dfa d q } [] \quad\quad = q$$
$$\text{dfa d q } (x : xs) = \text{dfa d } (d\ q\ x)\ xs$$

## Question

Does this function look familiar?

$$\text{dfa} = \text{foldl}$$

# Acceptance by a DFA

A word xs is **accepted** by a DFA if running the DFA on the word, starting in the start state S, yields an accepting state.

Universiteit Utrecht

# Acceptance by a DFA

A word xs is **accepted** by a DFA if running the DFA on the word, starting in the start state S, yields an accepting state.

```
dfaAccept :: [X] → (Q → X → Q, Q, Set Q) → Bool
dfaAccept xs (d, s, fs) = dfa d s xs 'member' fs
```

Universiteit Utrecht

# Language of a DFA

All words that are accepted by the DFA $(d, S, F)$.

$\{\, w \in [X] \mid \text{dfaAccept } w\ (d, S, F)\,\}$

Universiteit Utrecht

# Language of a DFA

All words that are accepted by the DFA $(d, S, F)$.

$\{\, w \in [X] \mid \mathsf{dfaAccept}\ w\ (d, S, F)\,\}$

One language can in general be described by multiple automata.
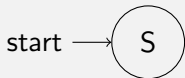
Universiteit Utrecht

### Question

Can the empty language be described by a DFA?

# Exercise

### Question

Can the empty language be described by a DFA?

start $\longrightarrow$ ( S )

Any automaton without accepting states is possible.

Universiteit Utrecht

# Exercise

## Question

Can the language $\{\varepsilon\}$ be described by a DFA?

Universiteit Utrecht

# Exercise

## Question

Can the language $\{\varepsilon\}$ be described by a DFA?

start $\longrightarrow$ (( S ))

Universiteit Utrecht

# Exercise

### Question

Can the language $\{\varepsilon\}$ be described by a DFA?

start $\longrightarrow$ ( ( S ) )

In general, any automaton where the starting state is accepting will accept the empty word (and possibly other words).

Universiteit Utrecht

# Observation

Running a DFA is clearly possible in linear time and constant space.

Universiteit Utrecht

# Nondeterministic finite state automata (NFA)

Similar to DFA, but:

▶ Potentially multiple start states.
▶ Potentially multiple transitions for the same terminal from a given state.

# Nondeterministic finite state automata (NFA)

Similar to DFA, but:

- ▶ Potentially multiple start states.
- ▶ Potentially multiple transitions for the same terminal from a given state.

Formally:

- ▶ an input alphabet X,
- ▶ a set of states Q,
- ▶ a transition function d of type $Q \rightarrow X \rightarrow \text{Set } Q$,
- ▶ a **set of** start states $S \subseteq Q$,
- ▶ a set of accepting states $F \subseteq Q$.

# Interpretation using choices

- ▶ We can choose a start state.
- ▶ When processing a terminal, we can choose one of the possible transitions for that terminal at that state and thereby end up with a new state.
- ▶ A word is accepted if there is a sequence of choices that gets us to an accepting state.

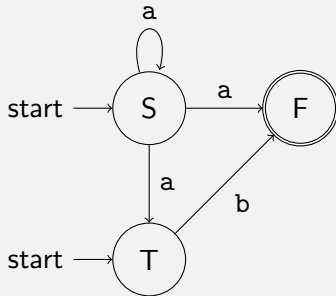# Interpretation using a set of all possible choices

- ▶ At any time, a set of states in the NFA is active. We start with the set of start states.
- ▶ When we process a terminal, we take all possible actions from all current states and thereby end up with a new set of states.
- ▶ A word is accepted if the set of states that is active after processing the word contains at least one accepting state.

Universiteit Utrecht

# Running an NFA

nfa :: $(Q \rightarrow X \rightarrow \text{Set } Q) \rightarrow \text{Set } Q \rightarrow [X] \rightarrow \text{Set } Q$
nfa d qs [] $\quad$ = qs
nfa d qs $(x : xs)$ = nfa d (join (map (flip d x) qs)) xs

where join is the concat for sets and computes the union of a set of sets:

join :: $\text{Set } (\text{Set } Q) \rightarrow \text{Set } Q$

Universiteit Utrecht

# Acceptance by an NFA

$$\text{nfaAccept} :: [X] \rightarrow (Q \rightarrow X \rightarrow \text{Set } Q, \text{Set } Q, \text{Set } Q) \rightarrow \text{Bool}$$
$$\text{nfaAccept xs } (d, ss, fs) = \text{not } (\text{null } (\text{nfa } d \text{ ss xs `intersect` fs}))$$

Universiteit Utrecht

# 10.3  NFAs vs. DFAs

# From DFA to NFA

Every DFA $(d, S, F)$ is trivially an NFA.

Universiteit Utrecht

# From DFA to NFA

Every DFA $(d, S, F)$ is trivially an NFA.

The start state $S$ becomes the one-element set of start states $\{S\}$.

Universiteit Utrecht

# From DFA to NFA

Every DFA $(d, S, F)$ is trivially an NFA.

The start state $S$ becomes the one-element set of start states $\{S\}$.

The transition function is changed such that it returns singleton sets:

$d' :: Q \rightarrow X \rightarrow Set\ Q$
$d'\ q\ x = singleton\ (d\ x)$

# From DFA to NFA

Every DFA $(d, S, F)$ is trivially an NFA.

The start state $S$ becomes the one-element set of start states $\{S\}$.

The transition function is changed such that it returns singleton sets:

```
d' :: Q → X → Set Q
d' q x = singleton (d x)
```

It is quite easy to show that the resulting NFA accepts the same language.

Universiteit Utrecht

# From NFA to DFA

We can also make a DFA from an NFA.

Question

How?

Universiteit Utrecht

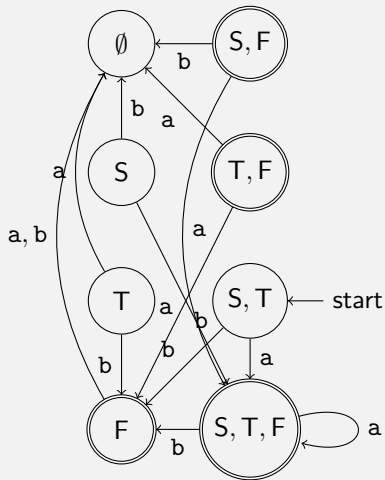# From NFA to DFA

We can also make a DFA from an NFA.

Question

How?

The construction is called the **powerset construction**:

- For each **set of states** in the NFA, we get **one state** in the DFA.
- The set of start states in the NFA thus corresponds to a single state in the DFA.
- Since the transition function for the NFA takes sets of states to sets of states, we can then reuse it for the DFA.
- All states that contain an accepting state of the NFA become accepting states in the DFA.

# From NFA to DFA – example

**Universiteit Utrecht**

# From NFA to DFA – example

**Universiteit Utrecht**

# 10.4 Regular grammars

# Regular grammar

A context-free grammar G is called **regular** if all productions are of one of the following two forms:

$A \rightarrow xB$
$A \rightarrow x$

where $x$ is a (possibly empty) sequence of terminals, and A and B are nonterminals.

**Universiteit Utrecht**

# Regular grammar

A context-free grammar G is called **regular** if all productions are of one of the following two forms:

$A \rightarrow xB$
$A \rightarrow x$

where $x$ is a (possibly empty) sequence of terminals, and A and B are nonterminals.

In other words: Every right hand side has at most one nonterminal that must occur in the end.

Universiteit Utrecht

# Regular language

A language is called **regular** if it can be described by a regular grammar.

Universiteit Utrecht

# Regular vs. context-free

From the definition, it is immediately clear that every regular language is context-free.

Universiteit Utrecht

# Regular vs. context-free

From the definition, it is immediately clear that every regular language is context-free.

## Question

Is every context-free language regular?

Universiteit Utrecht

# Regular vs. context-free

From the definition, it is immediately clear that every regular language is context-free.

## Question

Is every context-free language regular?

No. The standard example is the language $\{a^n b^n \mid n \in \mathbb{N}\}$.

Universiteit Utrecht

# Regular vs. context-free

From the definition, it is immediately clear that every regular language is context-free.

## Question

Is every context-free language regular?

No. The standard example is the language $\{a^n b^n \mid n \in \mathbb{N}\}$.

We will investigate later how such a negative statement (not belonging to the class of regular languages) can be proved.

# Closure properties

As context-free languages, regular languages are closed under

- union (corresponds to the $<|>$ combinator),
- sequencing (corresponds to the $<\!\!*\!\!>$ combinator),
- the star operator (corresponds to the many combinator).

# Closure properties

As context-free languages, regular languages are closed under

- union (corresponds to the $<|>$ combinator),
- sequencing (corresponds to the $<\!\!*\!\!>$ combinator),
- the star operator (corresponds to the many combinator).

While context-free languages are not closed under intersection, regular languages are (werkcollege).

Universiteit Utrecht

# Simplifying regular grammars

For every regular language, there is a regular grammar that has no productions of the form

$A \rightarrow B$
$C \rightarrow \varepsilon$

where A, B, and C are nonterminals, except that for the start symbol there may be a production

$S \rightarrow \varepsilon$

Universiteit Utrecht

# Simplifying regular grammars – contd.

The grammar transformation works in two phases:

- ▶ first all productions of the form A → B are removed;
- ▶ then all epsilon-productions are removed.

Universiteit Utrecht

# Simplifying regular grammars – contd.

Consider all pairs of nonterminals Y and Z.

If $Y \Rightarrow^* Z$:

- for every production $Z \rightarrow z$ (with z a sequence of symbols, but not a single nonterminal), add a production $Y \rightarrow z$.

If $Y \rightarrow Z$ is in the grammar, remove it.

# Simplifying regular grammars – contd.

Consider all pairs of nonterminals Y and Z.

If $Y \Rightarrow^* Z$:

- ▶ for every production $Z \rightarrow z$ (with z a sequence of symbols, but not a single nonterminal), add a production $Y \rightarrow z$.

If $Y \rightarrow Z$ is in the grammar, remove it.

The only problematic productions left are now epsilon-productions.

Universiteit Utrecht

# Simplifying regular grammars – contd.

For each production $Y \rightarrow \varepsilon$, consider all productions $Z \rightarrow xY$ (where $x$ now can consist only of terminals) and add a production $Z \rightarrow x$.

Then remove all epsilon-productions but $S \rightarrow \varepsilon$ if it exists.

# Simplifying regular grammars – contd.

We can simplify a regular grammar even further and require that all productions are of one of the following two forms

$$Y \rightarrow xZ$$
$$Y \rightarrow x$$

where $x$ is thus a single terminal, except for the start symbol S, for which a production of $S \rightarrow \varepsilon$ is allowed.

Universiteit Utrecht

# Simplifying regular grammars – contd.

We can simplify a regular grammar even further and require that all productions are of one of the following two forms

$$Y \rightarrow xZ$$
$$Y \rightarrow x$$

where $x$ is thus a single terminal, except for the start symbol S, for which a production of $S \rightarrow \varepsilon$ is allowed.

The transformation works by introducing new nonterminals. For example

$$A \rightarrow xyC$$

is transformed into

$$A \rightarrow xB$$
$$B \rightarrow yC$$

Universiteit Utrecht

[Faculty of **Science**
Information and Computing Sciences]

# 10.5 Regular grammars vs. finite state automata
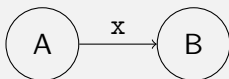
# From NFA to regular grammars

For each NFA, there exists a regular grammar that describes the same language.

Universiteit Utrecht

# From NFA to regular grammars

For each NFA, there exists a regular grammar that describes the same language.

- ▶ the states become nonterminals,
- ▶ the start state becomes the start symbol,
- ▶ for each transition

$$A \xrightarrow{\quad x \quad} B$$

we introduce a production

$$A \rightarrow xB$$

- ▶ for each accepting state F we introduce a production

$$A \rightarrow \varepsilon$$

Universiteit Utrecht

We can also produce an automaton for every regular grammar.

Universiteit Utrecht

# From regular grammars to an NFA

We can also produce an automaton for every regular grammar.

We first simplify the grammar. Then, all the hard work is done.

**Universiteit Utrecht**