

Towards more complex grammar systems

Some basic formal language theory

- Grammars, or: how to specify linguistic knowledge
- Automata, or: how to process with linguistic knowledge
- Levels of complexity in grammars and automata:
The Chomsky hierarchy

Grammars

A grammar is a 4-tuple (N, Σ, S, P) where

- N is a finite set of **non-terminals**
- Σ is a finite set of **terminal symbols**,
with $N \cap \Sigma = \emptyset$
- S is a distinguished **start symbol**, with $S \in N$
- P is a finite set of **rewrite rules** of the form $\alpha \rightarrow \beta$,
with $\alpha, \beta \in (N \cup \Sigma)^*$ and α including at least one
non-terminal symbol.

A simple example

$$N = \{S, NP, VP, V_i, V_t, V_s\}$$

$$\Sigma = \{\text{John, Mary, laughs, loves, thinks}\}$$

$$S = S$$

$$P = \left\{ \begin{array}{l} S \rightarrow NP VP \\ NP \rightarrow \text{John} \\ NP \rightarrow \text{Mary} \\ VP \rightarrow V_i \\ VP \rightarrow V_t NP \\ VP \rightarrow V_s S \\ p \\ V_i \rightarrow \text{laughs} \\ V_t \rightarrow \text{loves} \\ V_s \rightarrow \text{thinks} \end{array} \right\}$$

How does a grammar define a language?

Assume $\alpha, \beta \in (N \cup \Sigma)^*$, with α containing at least one non-terminal.

- A **sentential form** for a grammar G is defined as:
 - The start symbol S of G is a sentential form.
 - If $\alpha\beta\gamma$ is a sentential form and there is a rewrite rule $\beta \rightarrow \delta$ then $\alpha\delta\gamma$ is a sentential form.
- α (directly or immediately) **derives** β if $\alpha \rightarrow \beta \in P$.
One writes:
 - $\alpha \Rightarrow^* \beta$ if β is derived from α in zero or more steps
 - $\alpha \Rightarrow^+ \beta$ if β is derived from α in one or more steps
- A **sentence** is a sentential form consisting only of terminal symbols.
- The **language** $L(G)$ generated by the grammar G is the set of all sentences which can be derived from the start symbol S , i.e., $L(G) = \{\gamma \mid S \Rightarrow^* \gamma\}$

Processing with grammars: automata

An **automaton** in general has three components:

- an **input tape**, divided into squares with a read-write head positioned over one of the squares
- an **auxiliary memory** characterized by two functions
 - fetch: memory configuration \rightarrow symbols
 - store: memory configuration \times symbol \rightarrow memory configuration
- and a **finite-state control** relating the two components.

Different levels of complexity in grammars and automata

Let $A, B \in N$, $x \in \Sigma$, $\alpha, \beta, \gamma \in (\Sigma \cup T)^*$, and $\delta \in (\Sigma \cup T)^+$, then:

Type	Automaton		Grammar	
	Memory	Name	Rule	Name
0	Unbounded	TM	$\alpha \rightarrow \beta$	General rewrite
1	Bounded	LBA	$\beta A \gamma \rightarrow \beta \delta \gamma$	Context-sensitive
2	Stack	PDA	$A \rightarrow \beta$	Context-free
3	None	FSA	$A \rightarrow xB, A \rightarrow x$	Right linear

Abbreviations:

- TM: Turing Machine
- LBA: Linear-Bounded Automaton
- PDA: Push-Down Automaton
- FSA: Finite-State Automaton

Type 3: Right-Linear Grammars and Finite-State Automata

A **right-linear grammar** is a 4-tuple (N, Σ, S, P) with

P a finite set of rewrite rules of the form $\alpha \rightarrow \beta$, with $\alpha \in N$ and $\beta \in \{\gamma\delta \mid \gamma \in \Sigma^*, \delta \in N \cup \{\epsilon\}\}$, i.e.:

- left-hand side of rule: a single non-terminal, and
- right-hand side of rule: a string containing at most one non-terminal, as the rightmost symbol

Right-linear grammars are formally equivalent to left-linear grammars (at most one, leftmost non-terminal).

Finite-state transition network: states + arcs

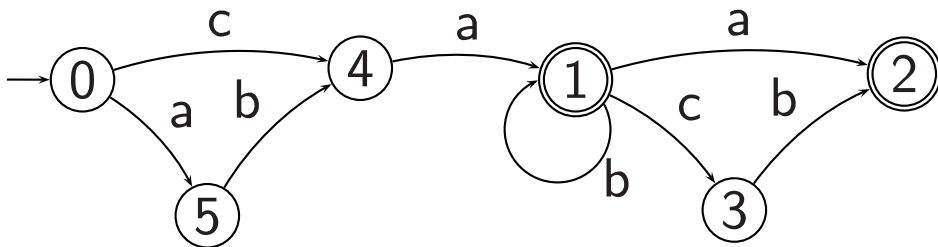
A **finite-state machine** consists of

- a tape
- a finite-state control
- no auxiliary memory

A regular language example:

$$(ab|c)ab^* (a|cb)?$$

Finite-state transition network:



Right-linear grammar:

$$N = \{\text{Expr}, X, Y, Z\}$$

$$\Sigma = \{a, b, c\}$$

$$S = \text{Expr}$$

$$P = \left\{ \begin{array}{l} \text{Expr} \rightarrow ab X \\ \text{Expr} \rightarrow c X \\ X \rightarrow a Y \\ Y \rightarrow b Y \\ Y \rightarrow Z \\ Z \rightarrow a \\ Z \rightarrow cb \\ Z \rightarrow \epsilon \end{array} \right.$$

Thinking about regular languages

- A language is regular iff one can define a FSM (or regular expression) for it.
- An FSM only has a fixed amount of memory, namely the number of states.
- Strings longer than the number of states, in particular also any infinite ones, must result from a loop in the FSM.
- Pumping Lemma: if for an infinite string there is no such loop, the string cannot be part of a regular language.

Type 2: Context-Free Grammars and Push-Down Automata

A **context-free grammar** is a 4-tuple (N, Σ, S, P) with

P a finite set of rewrite rules of the form $\alpha \rightarrow \beta$, with $\alpha \in N$ and $\beta \in (\Sigma \cup N)^*$, i.e.:

- left-hand side of rule: a single non-terminal, and
- right-hand side of rule: a string of terminals and/or non-terminals

A **push-down automaton** is a

- finite state automaton, with a
- stack as auxiliary memory

A context-free language example: $a^n b^n$

Context-free grammar:

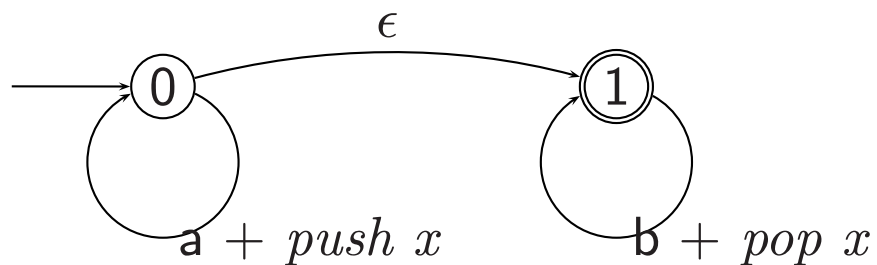
$$N = \{S\}$$

$$\Sigma = \{a, b\}$$

$$S = S$$

$$P = \left\{ \begin{array}{l} S \rightarrow a S b \\ S \rightarrow \epsilon \end{array} \right\}$$

Push-down automaton:



Type 1: Context-Sensitive Grammars and Linear-Bounded Automata

A rule of a **context-sensitive grammar**

- rewrites at most one non-terminal from the left-hand side.
- Contextual restrictions on the occurrence of this non-terminal may be imposed.
- The non-terminal must not rewrite as the empty string ϵ .

A **linear-bounded automaton** is a

- finite state automaton, with an
- auxiliary memory which cannot exceed the length of the input string.

A context-sensitive language example:

$$a^n b^n c^n$$

Context-sensitive grammar:

$$N = \{S, B, C\}$$

$$\Sigma = \{a, b\}$$

$$S = S$$

$$P = \left\{ \begin{array}{l} S \rightarrow a S B C, \\ S \rightarrow a b C, \\ b B \rightarrow b b, \\ b C \rightarrow b c, \\ c C \rightarrow c c, \\ C B \rightarrow B C \end{array} \right\}$$

Type 0: General Rewrite Grammar and Turing Machines

- In a **general rewrite grammar** there are no restrictions on the form of a rewrite rule.
- A **turing machine** has an unbounded auxiliary memory.
- Any language for which there is a recognition procedure can be defined, but recognition problem is not decidable.

Properties of different language classes

Reasoning:

- Languages are defined to be sets of strings.
- One can therefore apply set operations to languages and investigate results for particular language classes.

Some closure properties:

- All language classes are closed under **union with themselves**.
- All language classes are closed under **intersection with regular languages**.
- The class of **context-free languages is not closed under intersection with itself**. Proof:

Assume the two context-free languages L_1 and L_2 :

- $L_1 = \{a^n b^n c^i \mid n \geq 1 \text{ and } i \geq 0\}$
- $L_2 = \{a^j b^n c^n \mid n \geq 1 \text{ and } j \geq 0\}$

Their intersection is not context-free:

- $L_1 \cap L_2 = \{a^n b^n c^n \mid n \geq 1\}$

Criteria under which to evaluate grammar formalisms

There are three kinds of criteria:

- linguistic naturalness
- mathematical power
- computational effectiveness and efficiency

The weaker the type of grammar:

- the stronger the claim made about possible languages
- the greater the efficiency of the parsing procedure

Reasons for choosing a stronger grammar class:

- to capture the bare facts about actual languages
- to provide for elegant analyses capturing more generalizations (→ more “compact” grammars)

Language classes and natural languages

Natural languages are not regular

Example grammar:

$$\left\{ \begin{array}{l} S \rightarrow \text{If } S \text{ then } S \\ S \rightarrow \text{Either } S \text{ or } S \\ S \rightarrow \text{I laugh} | \text{I have to sneeze} | \text{Tim screams} \end{array} \right\}$$

Example analyses:

- [If [I laugh] then [I have to sneeze]]
- [If [either [I laugh] or [I have to sneeze]] then [Tim screams]]
- . . .

A more abstract version of the grammar rules:

$$\left\{ \begin{array}{l} S \rightarrow aSaS \\ S \rightarrow bSbS \\ S \rightarrow \epsilon \end{array} \right\}$$

which accepts $a^n b^m b^m a^n$ which is not a regular language.

Accounting for bare facts vs. linguistically sensible analyses

Looking at grammars from a linguistic perspective, one can distinguish their

- **weak generative capacity**, considering only the set of strings generated by a grammar
- **strong generative capacity**, considering the set of strings and their syntactic analyses generated by a grammar

Two grammars can be strongly or weakly equivalent.

Example for weakly equivalent grammars

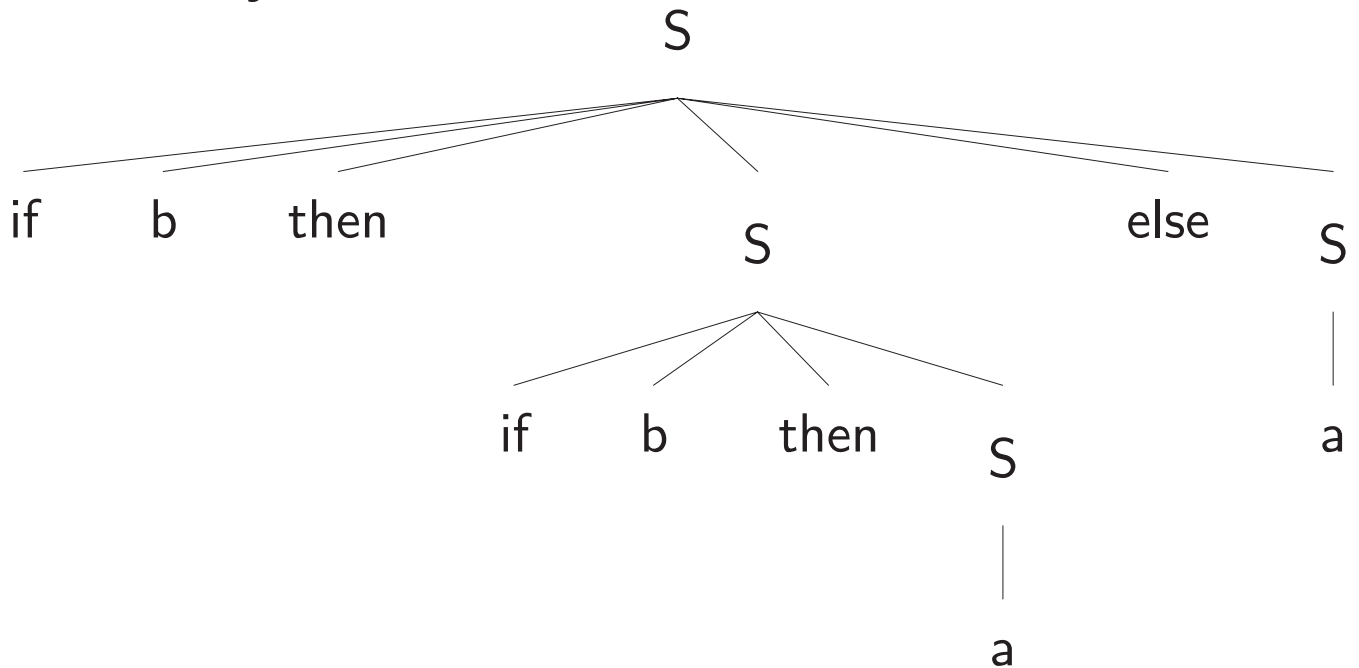
Example string:

if b then if b then a else a p

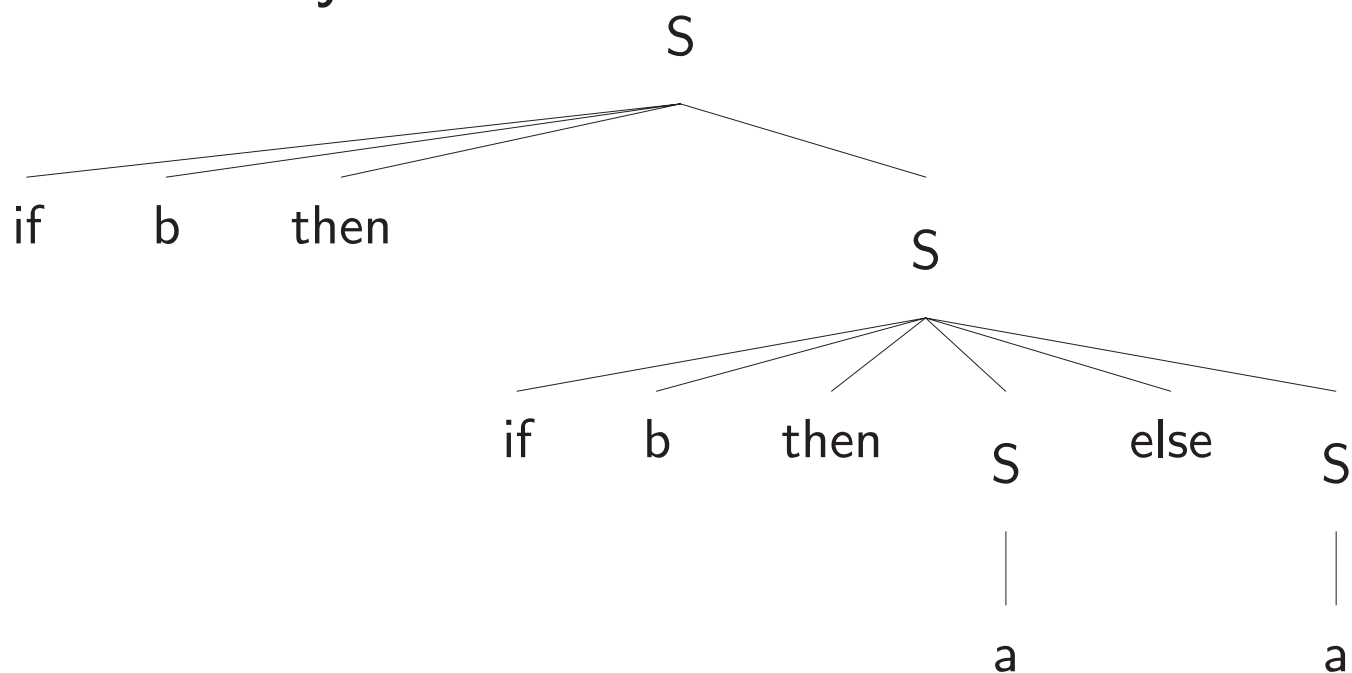
Grammar 1 rules:

$$\left\{ \begin{array}{l} S \rightarrow \text{if b then } S \text{ else } S, \\ S \rightarrow \text{if b then } S, \\ S \rightarrow a \end{array} \right\}$$

First analysis:



Second analysis:



Grammar 2 rules: A weakly equivalent grammar eliminating the ambiguity (only licenses second analysis).

$$\left\{ \begin{array}{l} S1 \rightarrow \text{if } b \text{ then } S1, \\ S1 \rightarrow \text{if } b \text{ then } S2 \text{ else } S1, \\ S1 \rightarrow a, \\ S2 \rightarrow \text{if } b \text{ then } S2 \text{ else } S2, \\ S2 \rightarrow a \end{array} \right\}$$