

Is English a finite state language?

## Formal languages

In **formal language theory**, a language is a set of strings. A string is just a sequence of symbols chosen from an agreed-upon set of symbols, called the vocabulary or **lexicon**.

formal language	(non-)examples
any sequence of English words from the OED	$\in$ “at by green maybe”, “vanity, all vanity” $\notin$ “I think que nous allons tromper the lesson”
strings over $\{a, b, c\}$ starting with $a$	$\in$ $abbb$ , $a$ , $a$ followed by a million $b$ 's $\notin$ $bcacc$ , the zero-length string $\epsilon$
strings over $\{a, b, c, d\}$ in alphabetical order	$\in$ $abd$ , $ad$ , $bcd$ , $b$ , $abcd$ $\notin$ $dbcd$ , $ba$

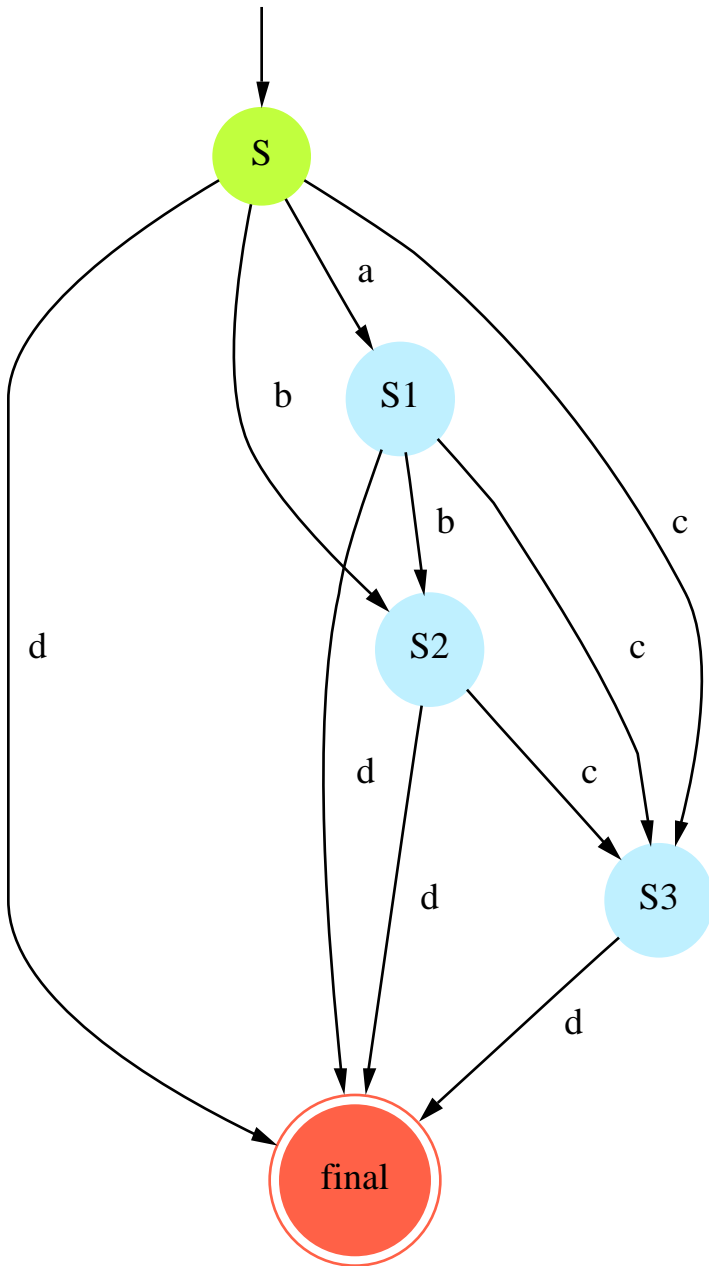
The idea of **generative grammar** is to use grammars to define a set that closely resembles a natural language – for instance, all and only the acceptable English sentences. However, not all sets are definable by all types of grammars. What kind of rules do we require to accurately describe natural languages?

## Finite-state grammars

To define (“generate”) strings over  $\{a, b, c, d\}$  in alphabetical order, let  $S$  be the start symbol of the grammar in (1).

$$\begin{array}{l}
 S \rightarrow a S1 \\
 S \rightarrow b S2 \\
 S \rightarrow c S3 \\
 S \rightarrow d \\
 (1) \quad S1 \rightarrow b S2 \\
 \quad S1 \rightarrow c S3 \\
 \quad S1 \rightarrow d \\
 \quad S2 \rightarrow c S3 \\
 \quad S2 \rightarrow d \\
 \quad S3 \rightarrow d
 \end{array}$$

Grammar (1) has a corresponding **finite state machine** that recognizes all and only the sentences it generates.



## The general form of finite state grammars

Any grammar having only rules of the form  $A \rightarrow bC$  where  $A, B$  are nonterminals and  $b$  is a terminal has a corresponding finite state machine. Given a string, if a path can be found through the machine, the string is generated by the grammar and vice versa.

There are some languages that cannot be recognized by finite state machines.

formal language	(non-)examples
sequence of $a$ 's followed by an equal number of $b$ 's	$\in$ $ab, aabb, aaabbb, \dots$ $\notin$ $aabbb, aaaaaaaaaaabb$

Call the number of  $a$ 's in the sentence being analyzed  $n$ . A finite state machine would need to 'remember' this number  $n$  while waiting for the end of the  $b$ 's. But, by definition, a finite state machine will only have enough states to remember some fixed number of  $a$ 's. Hence there exists neither a finite-state grammar nor a finite state machine for the language  $a^n b^n$ .

## English is just like $a^n b^n$

Consider

- (2) a. The cat died.  
 b. The cat the dog chased died.  
 c. The cat the dog the rat bit chased died.  
 d. The cat the dog the rat the elephant admired bit chased died  
 ⋮

If  $A = \{\text{the cat, the dog, the rat, the elephant, the kangaroo, ...}\}$  and  $B = \{\text{chased, bit, admired, ate, befriended, ...}\}$  it's clear that (2) have the structure

$$(3) \quad a^n b^{n-1} \text{died}$$

Chomsky and Miller (1963) argued that the obligatory paired dependencies presented by *either...or*, *if...then* or the agreement between verbs and subjects can nest inside one another to an arbitrary depth.

- (4) Anyone<sub>1</sub> who feels that if<sub>2</sub> so-many<sub>3</sub> more<sub>4</sub> students<sub>5</sub> whom we<sub>6</sub> haven't<sub>6</sub> actually admitted are<sub>5</sub> sitting in on the course than<sub>4</sub> ones we have that<sub>3</sub> the room had to be changed, then<sub>2</sub> probably auditors will have to be excluded, is<sub>1</sub> likely to agree that the curriculum needs revision.

This is analogous to the reversal language  $\{xx^R \mid x \in \{a, b\}^*\}$ , which also cannot be generated by a finite automaton.

Grammar (5) does generate the language  $a^n b^n$ . It is a **context-free** grammar, and as such is capable of deriving any number of center-embeddings.

$$(5) \quad \begin{array}{l} S \rightarrow a b \\ S \rightarrow a S b \end{array}$$

There is a corresponding machine that accepts the context-free languages: the **push down automaton**.

## The Chomsky Hierarchy

Let  $A, B$  be nonterminals,  $b$  a terminal,  $\alpha$  a nonempty sequence of either kind of symbol, and  $\gamma, \delta$  possibly empty sequences of either kind of symbol. The **Chomsky Hierarchy** is a classification of languages in a subset relationship. Each language has corresponding class of machine that recognizes it.

language	rule type	machine
finite-state	$A \rightarrow bC$	finite state machine
context-free	$A \rightarrow \alpha$	push down automaton
context-sensitive	$\gamma A \delta \rightarrow \gamma \alpha \delta$	linear bounded automaton
unrestricted	no restriction	Turing machine

The cardinality of the set of languages definable by a grammar formalism is called its **generative capacity**. Natural languages are believed to reside somewhere between context-free and context-sensitive.