

Automata and Formal Languages (7)

Email: christian.urban at kcl.ac.uk
Office: SI.27 (1st floor Strand Building)
Slides: KEATS (also home work is there)

CFGs

A **context-free** grammar (CFG) G consists of:

- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}_1 | \text{rhs}_2 | \dots$$

where **rhs** are sequences involving terminals and nonterminals (can also be empty).

CFGs

A **context-free** grammar (CFG) G consists of:

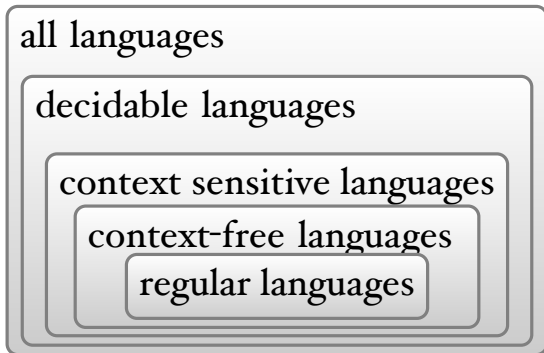
- a finite set of nonterminal symbols (upper case)
- a finite terminal symbols or tokens (lower case)
- a start symbol (which must be a nonterminal)
- a set of rules

$$A \rightarrow \text{rhs}_1 | \text{rhs}_2 | \dots$$

where **rhs** are sequences involving terminals and nonterminals (can also be empty).

Hierarchie of Languages

Recall that languages are sets of strings.



Numbers

A grammar for numbers:

$$N \rightarrow N \cdot N \mid 0 \mid 1 \mid \dots \mid 9$$

Unfortunately left-recursive (and ambiguous).

A problem for **recursive descent parsers**
(e.g. parser combinators).

Numbers

A grammar for numbers:

$$N \rightarrow N \cdot N \mid 0 \mid 1 \mid \dots \mid 9$$

Unfortunately left-recursive (and ambiguous).

A problem for **recursive descent parsers**
(e.g. parser combinators).

A non-left-recursive grammar for numbers

$$N \rightarrow 0 \cdot N \mid 1 \cdot N \mid \dots \mid 0 \mid 1 \mid \dots \mid 9$$

Chomsky Normal Form

All rules must be of the form

$$A \rightarrow a$$

or

$$A \rightarrow B \cdot C$$

CYK Algorithm

$S \rightarrow N \cdot P$

$P \rightarrow V \cdot N$

$N \rightarrow N \cdot N$

$N \rightarrow$ students | Jeff | geometry | trains

$V \rightarrow$ trains

Jeff trains geometry students

CYK Algorithm

- runtime is $O(n^3)$
- grammars need to be transferred into CNF

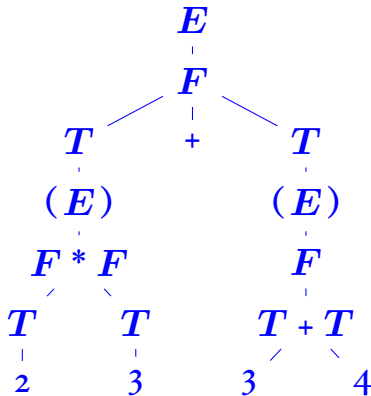
Parse Trees

$E \rightarrow F \mid F \cdot * \cdot F$

$F \rightarrow T \mid T \cdot + \cdot T \mid T \cdot - \cdot T$

$T \rightarrow \text{num_token} \mid (.E.)$

$(2*3)+(3+4)$



Ambiguous Grammars

A CFG is **ambiguous** if there is a string that has at least two parse trees.

$$E \rightarrow \textit{num_token}$$

$$E \rightarrow E \cdot + \cdot E$$

$$E \rightarrow E \cdot - \cdot E$$

$$E \rightarrow E \cdot * \cdot E$$

$$E \rightarrow (\cdot E \cdot)$$

$$1 + 2 * 3 + 4$$

Dangling Else

Another ambiguous grammar:

$$\begin{array}{l} E \rightarrow \text{if } E \text{ then } E \\ \quad | \text{if } E \text{ then } E \text{ else } E \\ \quad | \text{id} \end{array}$$

if a then if x then y else c

A CFG Derivation

- 1 Begin with a string with only the start symbol S
- 2 Replace any non-terminal X in the string by the right-hand side of some production $X \rightarrow \text{rhs}$
- 3 Repeat 2 until there are no non-terminals

$S \rightarrow \dots \rightarrow \dots \rightarrow \dots \rightarrow \dots$