

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

$$A \rightarrow \text{rhs}$$

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

$$A \rightarrow \text{rhs}$$

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

We can also allow rules

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

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

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

Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language $L(G)$ is:

$$\{c_1 \dots c_n \mid \forall i. c_i \in T \wedge S \rightarrow^* c_1 \dots c_n\}$$

Language of a CFG

Let G be a context-free grammar with start symbol S . Then the language $L(G)$ is:

$$\{c_1 \dots c_n \mid \forall i. c_i \in T \wedge S \rightarrow^* c_1 \dots c_n\}$$

- Terminals are so-called because there are no rules for replacing them
- Once generated, terminals are “permanent”
- Terminals ought to be tokens of the language (at least in this course)

Arithmetic Expressions

$E \rightarrow num_token$

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

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

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

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

Arithmetic Expressions

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

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

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

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

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

A CFG is **left-recursive** if it has a nonterminal E such that $E \rightarrow^+ E \cdot \dots$

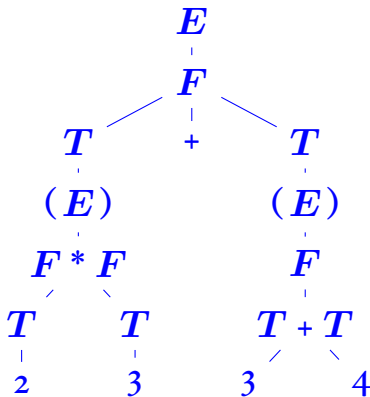
Parse Trees

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

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

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

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



Ambiguous Grammars

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

$$E \rightarrow num_token$$

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

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

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

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

$$1 + 2 * 3 + 4$$

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