

CS:4330 Theory of Computation
Spring 2018

Context-Free Languages
Context-Free Grammars

Haniel Barbosa



Readings for this lecture

Chapter 2 of [Sipser 1996], 3rd edition. Section 2.1.

Context-Free Grammars (CFG)

- ▷ There are languages, such as $\{0^n 1^n \mid n \geq 0\}$ that cannot be described by finite automata (or regexps)

- ▷ Context-free grammars provide a more powerful mechanism for language specification.

- ▷ Context-free grammars can describe features that have a recursive structure, making them useful beyond finite automata.

Formal definition of a CFG

A context-free grammar is a 4-tuple (V, Σ, R, S) in which:

- ▷ V is a finite set of symbols called the *variables* or *nonterminals*
- ▷ Σ is a finite set of symbols, disjoint from V , called *terminals*
- ▷ R is a finite set of *rules* of the form $lhs \rightarrow rhs$, in which $lhs \in V$ and $rhs \in (V \cup \Sigma)^*$
- ▷ $S \in V$ is the start nonterminal

Example

- ▷ CFG G_1 has the following rules:

$$A \rightarrow 0A1$$

$$A \rightarrow B$$

$$B \rightarrow \#$$

- ▷ Nonterminals of G_1 are $\{A, B\}$ and A is the start symbol

- ▷ Terminals of G_1 are $\{0, 1, \#\}$

Language specification

A grammar is used for a language specification by generating each string of the language in the following manner:

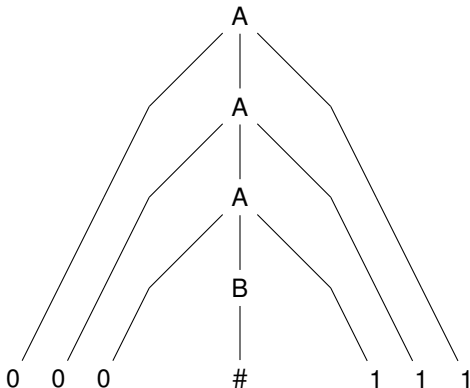
1. Write down the start variable; it is the *lhs* of the first rule, unless specified otherwise
2. Find a variable that is written down and a rule whose *lhs* is that variable. Replace the written down variable with the *rhs* of that rule.
3. Repeat step 2 until no variables remain in the string thus generated.

Note

The sequence of substitutions used to obtain a string using a CFG is called a *derivation* and may be represented by a tree called *derivation tree* or a *parse tree*

Example derivation tree

The derivation tree of the string 000#111 using CFG G_1 is:



Note

- ▷ All strings of terminals generated in this way constitute the language specified by the grammar

- ▷ We write $\mathcal{L}(G)$ for the language generated by the grammar G . Thus,
$$\mathcal{L}(G_1) = \{0^n\#1^n \mid n \geq 0\}$$

- ▷ A language generated by a context-free grammar (CFG) is called a *Context-Free Language* (CFL).

CFG G_2

The CFG G_2 specifies a fragment of English:

$\langle \text{SENTENCE} \rangle$	\rightarrow	$\langle \text{NOUN-PHRASE} \rangle \langle \text{VERB-PHRASE} \rangle$
$\langle \text{NOUN-PHRASE} \rangle$	\rightarrow	$\langle \text{CP-NOUN} \rangle \mid \langle \text{CP-NOUN} \rangle \langle \text{PREP-PHRASE} \rangle$
$\langle \text{VERB-PHRASE} \rangle$	\rightarrow	$\langle \text{CP-VERB} \rangle \mid \langle \text{CP-VERB} \rangle \langle \text{PREP-PHRASE} \rangle$
$\langle \text{PREP-PHRASE} \rangle$	\rightarrow	$\langle \text{PREP} \rangle \langle \text{CP-NOUN} \rangle$
$\langle \text{CP-NOUN} \rangle$	\rightarrow	$\langle \text{ARTICLE} \rangle \langle \text{NOUN} \rangle$
$\langle \text{CP-VERB} \rangle$	\rightarrow	$\langle \text{VERB} \rangle \mid \langle \text{VERB} \rangle \langle \text{NOUN-PHRASE} \rangle$
$\langle \text{ARTICLE} \rangle$	\rightarrow	a the
$\langle \text{NOUN} \rangle$	\rightarrow	boy girl flower
$\langle \text{VERB} \rangle$	\rightarrow	touches likes sees
$\langle \text{PREP} \rangle$	\rightarrow	with

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a⟨NOUN⟩⟨VERB-PHRASE⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a boy⟨VERB-PHRASE⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a boy⟨VERB-PHRASE⟩
⇒ a boy⟨CP-VERB⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a boy⟨VERB-PHRASE⟩
⇒ a boy⟨CP-VERB⟩
⇒ a boy⟨VERB⟩

Example derivation with G_2

⟨SENTENCE⟩	→	⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⟨NOUN-PHRASE⟩	→	⟨CP-NOUN⟩ ⟨CP-NOUN⟩⟨PREP-PHRASE⟩
⟨VERB-PHRASE⟩	→	⟨CP-VERB⟩ ⟨CP-VERB⟩⟨PREP-PHRASE⟩
⟨PREP-PHRASE⟩	→	⟨PREP⟩⟨CP-NOUN⟩
⟨CP-NOUN⟩	→	⟨ARTICLE⟩⟨NOUN⟩
⟨CP-VERB⟩	→	⟨VERB⟩ ⟨VERB⟩⟨NOUN-PHRASE⟩
⟨ARTICLE⟩	→	a the
⟨NOUN⟩	→	boy girl flower
⟨VERB⟩	→	touches likes sees
⟨PREP⟩	→	with

⟨SENTENCE⟩ ⇒ ⟨NOUN-PHRASE⟩⟨VERB-PHRASE⟩
⇒ ⟨CP-NOUN⟩⟨VERB-PHRASE⟩
⇒ ⟨ARTICLE⟩⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a⟨NOUN⟩⟨VERB-PHRASE⟩
⇒ a boy⟨VERB-PHRASE⟩
⇒ a boy⟨CP-VERB⟩
⇒ a boy⟨VERB⟩
⇒ a boy sees

Direct derivation

- ▷ If $u, v, w \in (V \cup \Sigma)^*$ (i.e. are strings of variables and terminals) and $A \rightarrow w \in R$ (i.e. is a rule of the grammar), then we say that uAv yields uwv , written

$$uAv \Rightarrow uwv$$

- ▷ We may also say that uwv is directly derived from uAv using the rule $A \rightarrow w$

Derivation

▷ We write $u \xRightarrow{*} v$ if $u = v$ or if a sequence $u_1, \dots, u_k \in (V \cup \Sigma)^*$ exists, for $k \geq 0$, and $u \Rightarrow u_1 \Rightarrow \dots \Rightarrow u_k \Rightarrow v$

▷ We may also say that u, u_1, \dots, u_k, v is a derivation of v from u

Language specified by G

If $G = (V, \Sigma, R, S)$ is a CFG then the language specified by G (or the language of G) is a CFL

$$\mathcal{L}(G) = \{w \in \Sigma^* \mid S \xRightarrow{*} w\}$$

More examples of CFGs

- ▷ Consider the grammar

$$G_3 = (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid SS \mid \epsilon\}, S)$$

More examples of CFGs

- ▷ Consider the grammar

$$G_3 = (\{S\}, \{a, b\}, \{S \rightarrow aSb \mid SS \mid \epsilon\}, S)$$

- ▷ $\mathcal{L}(G_3)$ contains strings such as

abab, aaabbb, aababb

Note

If one thinks of a and b as the symbols '(' and ')' then we can see that $\mathcal{L}(G_3)$ is the language of all strings of properly nested parenthesis

Important application

Context-free grammars are used as basis for compiler design and implementation

- ▷ Context-free grammars are used as specification mechanisms for programming languages

- ▷ Designers of compilers use such grammars to implement compiler's components, such as scanners, parsers, code generators, *code synthesizers*

- ▷ The implementation of almost any programming languages is preceded by a context-free grammar that specifies it

Example

▷ Consider the grammar $G_4 = (\{E, T, F\}, \{a, +, *, (,)\}, R, E)$ in which R is:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid a$$

Example

- ▷ Consider the grammar $G_4 = (\{E, T, F\}, \{a, +, *, (,)\}, R, E)$ in which R is:

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid a$$

- ▷ $\mathcal{L}(G_4)$ is the language of arithmetic expressions

Application: program synthesis

Consider the grammar $G_S = (\{E, B\}, \{0, 1, x, y, +, \text{ite}, \leq, \wedge, \neg, , , (,)\}, R, E)$ in which R is:

$$E \rightarrow 0 \mid 1 \mid x \mid y \mid (E + E) \mid \text{ite}(B, E, E)$$

$$B \rightarrow (\neg B) \mid (B \wedge B) \mid (E \leq E)$$

What is a program generated with this grammar that solves the following problem:

$$\text{prog}(x, y) \geq x \wedge \text{prog}(x, y) \geq y$$

Application: program synthesis

Consider the grammar $G_S = (\{E, B\}, \{0, 1, x, y, +, \text{ite}, \leq, \wedge, \neg, , , (,)\}, R, E)$ in which R is:

$$E \rightarrow 0 \mid 1 \mid x \mid y \mid (E + E) \mid \text{ite}(B, E, E)$$

$$B \rightarrow (\neg B) \mid (B \wedge B) \mid (E \leq E)$$

What is a program generated with this grammar that solves the following problem:

$$\text{prog}(x, y) \geq x \wedge \text{prog}(x, y) \geq y$$

A solution is $\text{prog}(x, y) = \text{ite}(x \leq y, y, x)$, i.e. the max function.

Designing CFGs

- ▷ As with the design of automata, the design of CFGs requires creativity

- ▷ CFGs are even trickier to construct than finite automata since “we are more accustomed to programming a machine than to specify programming languages.”

Design Techniques

- ▷ Many CFGs are unions of simpler CFGs. Hence the suggestion is to construct smaller, simpler grammars first and then to join them into larger grammar
- ▷ The mechanism of grammar combination consists of putting all their rules together and adding the new rules

$$S \rightarrow S_1 \mid \cdots \mid S_k$$

where the nonterminals S_i , for $1 \leq i \leq k$, are the start variables of the individual grammars and S is the new variable

Example Grammar Design

Design a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$

Example Grammar Design

Design a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$

1. Construct the grammar $S_1 \rightarrow 0S_11 \mid \epsilon$ for the language

$$\{0^n 1^n \mid n \geq 0\}$$

Example Grammar Design

Design a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$

1. Construct the grammar $S_1 \rightarrow 0S_11 \mid \epsilon$ for the language

$$\{0^n 1^n \mid n \geq 0\}$$

Construct the grammar $S_2 \rightarrow 1S_20 \mid \epsilon$ for the language

$$\{1^n 0^n \mid n \geq 0\}$$

Example Grammar Design

Design a grammar for the language $\{0^n 1^n \mid n \geq 0\} \cup \{1^n 0^n \mid n \geq 0\}$

1. Construct the grammar $S_1 \rightarrow 0S_1 1 \mid \epsilon$ for the language

$$\{0^n 1^n \mid n \geq 0\}$$

Construct the grammar $S_2 \rightarrow 1S_2 0 \mid \epsilon$ for the language

$$\{1^n 0^n \mid n \geq 0\}$$

3. Put them together adding the rule $S \rightarrow S_1 \mid S_2$, obtaining

$$\begin{aligned} S &\rightarrow S_1 \mid S_2 \\ S_1 &\rightarrow 0S_1 1 \mid \epsilon \\ S_2 &\rightarrow 1S_2 0 \mid \epsilon \end{aligned}$$

Second Design Technique

- ▷ Constructing a CFG for a regular language is easy if one can first construct a DFA for the language
- ▷ Conversion procedure:
 1. Make a variable R_i for each state q_i of the DFA
 2. Add rules $R_i \rightarrow aR_j$ to the CFG if $\delta(q_i, a) = q_j$ is a transition in the DFA
 3. Add the rule $R_i \rightarrow \epsilon$ if q_i is an accept state of the DGA
 4. If q_0 is the start state of the DGA make R_0 the start variable of the CFG.

Theorem

Every regular language is context-free.

Third Design Technique

- ▷ Certain CFLs contain strings with two related substrings such as 0^n and 1^n in $\{0^n 1^n \mid n \geq 0\}$
- ▷ Example of relationship: to recognize such a language a machine would need to remember an unbounded amount of information about one of the substrings
- ▷ A CFG that handles this situation uses a rule of the form $R \rightarrow uRv$ which generates strings wherein the portion containing u 's corresponds to the portion containing v 's.

Example Application

Consider the CFG $G = (\{S, B\}, \{a, b\}, \{S \rightarrow aSb \mid B \mid \epsilon, B \rightarrow bB \mid b\}, S)$

- ▷ The following are derivations with G :

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSbBB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSBbB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaBB$$

which shows that derivations in this grammar can be quite complex

- ▷ When rewriting the strings $aaSBB$ we can consider further derivations of each of its symbols in isolation

Example Application

Consider the CFG $G = (\{S, B\}, \{a, b\}, \{S \rightarrow aSb \mid B \mid \epsilon, B \rightarrow bB \mid b\}, S)$

- ▷ The following are derivations with G :

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSbBB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSBbB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaBB$$

which shows that derivations in this grammar can be quite complex

- ▷ When rewriting the strings $aaSBB$ we can consider further derivations of each of its symbols in isolation
- ▷ Derivations from B are $B \Rightarrow bB \Rightarrow bbB \xRightarrow{*} b^{k-1}B \Rightarrow b^k, k \geq 1$

Example Application

Consider the CFG $G = (\{S, B\}, \{a, b\}, \{S \rightarrow aSb \mid B \mid \epsilon, B \rightarrow bB \mid b\}, S)$

- ▷ The following are derivations with G :

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSbBB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSBbB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaBB$$

which shows that derivations in this grammar can be quite complex

- ▷ When rewriting the strings $aaSBB$ we can consider further derivations of each of its symbols in isolation
- ▷ Derivations from B are $B \Rightarrow bB \Rightarrow bbB \xRightarrow{*} b^{k-1}B \Rightarrow b^k, k \geq 1$
- ▷ Therefore $S \Rightarrow aSB \xRightarrow{*} aSb^k, k \geq 1$

Example Application

Consider the CFG $G = (\{S, B\}, \{a, b\}, \{S \rightarrow aSb \mid B \mid \epsilon, B \rightarrow bB \mid b\}, S)$

▷ The following are derivations with G :

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSbBB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSBbB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaSB,$$

$$S \Rightarrow aSb \Rightarrow aaSBB \Rightarrow aaBB$$

which shows that derivations in this grammar can be quite complex

- ▷ When rewriting the strings $aaSBB$ we can consider further derivations of each of its symbols in isolation
- ▷ Derivations from B are $B \Rightarrow bB \Rightarrow bbB \xRightarrow{*} b^{k-1}B \Rightarrow b^k, k \geq 1$
- ▷ Therefore $S \Rightarrow aSB \xRightarrow{*} aSb^k, k \geq 1$
- ▷ Hence, $\mathcal{L}(G) = \{a^n b^m \mid n \leq m\}$

Ambiguity

- ▷ If a CFG generates the same string in several different ways, we say that the string is derived *ambiguously* in that grammar
- ▷ If a CFG generates some string we say that the grammar is *ambiguous*

Example

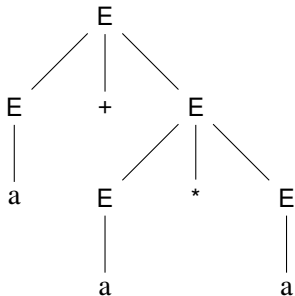
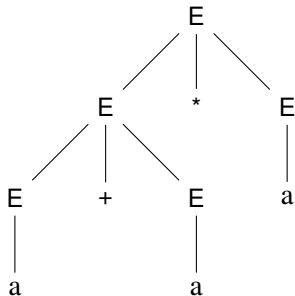
The grammar G_5 , whose rules are

$$E \rightarrow E + E \mid E * E \mid (E) \mid a$$

generates ambiguously some arithmetic expressions

Ambiguous expressions

Two different derivation trees for $a + a * a$



Note

- ▷ The grammar G_5 does not capture the casual precedence relations and so groups the $+$ before $*$ and vice versa
- ▷ In contrast, the grammar G_4 generates the same language, but every generated string has a unique derivation tree
 $G_4 = (\{E, T, F\}, \{a, +, *, (,)\}, R, E)$ in which R is:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid a \end{aligned}$$

- ▷ Hence, G_5 is ambiguous and G_4 is not, i.e. G_4 is *unambiguous*

Note

- ▷ When a grammar generates a string ambiguously it means that the string has two different derivation trees

- ▷ However, two different derivations may produce the same derivation tree because they may differ in the order in which they replace nonterminals, not in the rules they use

- ▷ To concentrate on the structure of the derivations we need to fix the order of rule application

Fixing rule application order

Definition (Leftmost derivation)

A derivation of a string w in a grammar G is a *leftmost derivation* if at every step the leftmost nonterminal is replaced.

Definition (Rightmost derivation)

A derivation of a string w in a grammar G is a *rightmost derivation* if at every step the rightmost nonterminal is replaced.

Inherent Ambiguity

- ▷ Some CFLs can have both ambiguous and unambiguous grammars.
- ▷ Some CFLs, however, can be generated only by an ambiguous grammar.
- ▷ A CFL that can be generated only by ambiguous grammars is called *inherently ambiguous*.

Example of inherently ambiguous language

$$\{0^i 1^j 2^k \mid i = j \vee j = k\}$$

Chomsky Normal Form

- ▷ It is often convenient to simplify CFGs so we can reason about them

- ▷ One of the simplest and most useful simplified forms of CFGs is called the Chomsky Normal Form

- ▷ Another normal form usually used in algebraic specifications is Greibach normal form

Definition

A context-free grammar G is in Chomsky normal form if every rule is of the form

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a \end{aligned}$$

where a is a terminal, A, B, C are nonterminals, and B, C may not be the start variable

Note

The rule $S \rightarrow \epsilon$, where S is the start variable, is not excluded

Chomsky Normal Form characterizes CFLs

Theorem

Any context-free language is generated by a context-free grammar in Chomsky normal form

Proof ideas

- ▷ Show that any grammar G can be converted into Chomsky normal form
- ▷ Conversion procedure has several states where the rules that violate Chomsky normal form conditions are replaced with equivalent ones that satisfy these conditions
- ▷ Order of transformations:
 1. add a new start variable
 2. eliminate all ϵ -rules
 3. eliminate unit rules
 4. convert rules
- ▷ Check that the obtained grammar defines the same language as the initial one.

Conversion: 1 - introduce new start state

Add a new start symbol S_0 and the rule $S_0 \rightarrow S$ where S was the original start symbol

Note

This change guarantees that the start symbol does not occur on the *rhs* of any rule

Conversion: 2 - eliminate ϵ -rules

Repeat

1. Eliminate the ϵ rule $A \rightarrow \epsilon$ where A is not the start symbol
2. For each occurrence of A on the *rhs* of a rule, add a new rule with that occurrence of A deleted
Example: To delete $A \rightarrow \epsilon$, replace $B \rightarrow uAv$ by $B \rightarrow uAv \mid uv$; replace $R \rightarrow uAvAw$ by $B \rightarrow uAvAw \mid uvAw \mid uAvw \mid uwv$
3. Replace the rule $B \rightarrow A$, (if it is present) by $B \rightarrow A \mid \epsilon$ unless the rule $B \rightarrow \epsilon$ has not been previously eliminated

until all ϵ rules are eliminated.

Conversion: 3 - remove unit rules

Repeat

1. Remove a unit rule $A \rightarrow B$
2. For each rule $B \rightarrow u$ that appears, add the rule $A \rightarrow u$, unless it was a unit rule previously removed

until all unit rules are eliminated.

Note

u is a string of variables and terminals

Conversion: 4 - Convert all remaining rules

Repeat

1. Replace a rule $A \rightarrow u_1, \dots, u_k$, $k \geq 3$, where each u_i , $1 \leq i \leq k$, is a variable or terminal, by:

$$A \rightarrow u_1 A_1, A_1 \rightarrow u_2 A_2, \dots, A_{k-2} \rightarrow u_{k-1} u_k$$

where A_1, \dots, A_{k-2} are new variables

2. If $k \geq 2$ replace any terminal u_i with a new variable U_i and add the rule $U_i \rightarrow u_i$

until no rules of the form $A \rightarrow u_1, \dots, u_k$ with $k \geq 3$ remain.

Example CFG conversion

Consider the grammar

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \epsilon$$

After the first step of transformation we get

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \epsilon$$

Removing ϵ rules

Removing $B \rightarrow \epsilon$:

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASA \mid aB \mid a \\ A &\rightarrow B \mid S \mid \epsilon \\ B &\rightarrow b \end{aligned}$$

Removing ϵ rules

Removing $B \rightarrow \epsilon$:

$$\begin{aligned}S_0 &\rightarrow S \\S &\rightarrow ASA \mid aB \mid a \\A &\rightarrow B \mid S \mid \epsilon \\B &\rightarrow b\end{aligned}$$

Removing $A \rightarrow \epsilon$

$$\begin{aligned}S_0 &\rightarrow S \\S &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S \\A &\rightarrow B \mid S \\B &\rightarrow b\end{aligned}$$

Removing unit rule

Removing $S \rightarrow S$

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

Removing unit rule

Removing $S \rightarrow S$

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \\ A &\rightarrow B \mid S \\ B &\rightarrow b \end{aligned}$$

Removing $S_0 \rightarrow S$

$$\begin{aligned} S_0 &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \\ S &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \\ A &\rightarrow B \mid S \\ B &\rightarrow b \end{aligned}$$

More unit rules

Removing $A \rightarrow B$

$S_0 \rightarrow S$

$S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$

$A \rightarrow S \mid b$

$B \rightarrow b$

More unit rules

Removing $A \rightarrow B$

$$\begin{aligned} S_0 &\rightarrow S \\ S &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \\ A &\rightarrow S \mid b \\ B &\rightarrow b \end{aligned}$$

Removing $A \rightarrow S$

$$\begin{aligned} S_0 &\rightarrow ASA \mid aB \mid SA \mid AS \\ S &\rightarrow ASA \mid aB \mid a \mid SA \mid AS \\ A &\rightarrow b \mid ASA \mid aB \mid a \mid SA \mid SA \\ B &\rightarrow b \end{aligned}$$

Converting the remaining rules

$S_0 \rightarrow AA_1 \mid UB \mid a \mid SA \mid AS$

$S \rightarrow AA_1 \mid UB \mid a \mid SA \mid AS$

$A \rightarrow b \mid AA_1 \mid UB \mid a \mid SA \mid AS$

$A_1 \rightarrow SA$

$U \rightarrow a$

$B \rightarrow b$