

CSCE 355
3/3/2025

Context-free languages

Used to describe programming language syntax.

Def: A context-free grammar (CFG)

is a tuple $\langle V, T, S, P \rangle$ where

V - a finite set. Elements are called variables, nonterminals, syntactic categories

(also Σ) T - an ~~finite~~ alphabet (elements are called terminals or tokens)

~~$V \cap T = \emptyset$~~ $V \cap T = \emptyset$

$S \in V$ (the start symbol)

P is a finite set of productions (or rewrite rules), each one has the form

$$\underbrace{A}_{\text{head}} \rightarrow \underbrace{\alpha}_{\text{body}}$$

where $A \in V$ and α is any string in $(V \cup T)^*$.

Idea: $A \rightarrow \alpha$ means "an A could be an α "
or "A can be replaced by α "
 S - the broadest syntactic category,

productions that describe syntax of expressions^②
over the arithmetic operators
 $+, -, *, /$ & constants & parens

Ex:
 $E \rightarrow E + E$
 $E \rightarrow E - E$
 $E \rightarrow E * E$
 $E \rightarrow E / E$
 $E \rightarrow (E)$
 $E \rightarrow c$

Think
 E — "expression"
("constant")

string

$(c + c) * c$

$E \Rightarrow \underline{E} * E \Rightarrow (E) * E \Rightarrow (E + E) * E$

$\Rightarrow (c + E) * E \Rightarrow (c + c) * E \Rightarrow \boxed{(c + c) * c}$

So $(c + c) * c$ is a syntactically correct arithmetic expression.

Formally, the grammar is

$G = \langle \{E\}, \{+, -, *, /, '(', ')', c\}, E, \{E \rightarrow E + E, \dots, E \rightarrow c\} \rangle$

Usually just give the list of productions:

- nonterminals are the left-hand sides (heads) of the productions
- terminals are any symbols appearing on the right-hand sides (bodies) ~~of~~ that are not nonterminals.
- start symbol is head of the first production in the list

Nonterminals — upper case roman letters

lang. of statements (S) in a prog. lang.

Ex: $S \rightarrow \{ L \}$ ($L = "$ list of statements

$S \rightarrow \underline{\text{while}} E \underline{\text{do}} S$

$S \rightarrow \underline{\text{if}} E \underline{\text{then}} S$

$S \rightarrow \underline{\text{if}} E \underline{\text{then}} S \underline{\text{else}} S$

$S \rightarrow \underline{\text{assignment}} ;$

$E \rightarrow$

.

;

as before

$E \rightarrow$

$L \rightarrow S L$

$L \rightarrow \epsilon$

(ϵ is metasymbol denoting the empty string)

$L \Rightarrow S L \Rightarrow S S L \Rightarrow S S S L \Rightarrow S S S$

Def: Let $G = \langle V, \Sigma, S, P \rangle$ be a CFG.

A derivation in G of a string $w \in \Sigma^*$ is a sequence of strings

$A_0 \Rightarrow A_1 \Rightarrow A_2 \Rightarrow \dots \Rightarrow A_n$

where all $A_i \in (V \cup \Sigma)^*$

- $A_0 = S$

- $A_n = w$

- For all $i, 0 \leq i < n, A_{i+1}$ is obtained from A_i by replacing an occurrence

of some $B \in V$ with the body of a production in P , whose head is B . (4)

Ex: $abBcd \Rightarrow abxyzcd$

where $B \rightarrow xyz$ is a production in P .

A_0, A_1, \dots, A_n — sentential forms

$A_0 \Rightarrow \dots \Rightarrow w$ is derivation of w

say that w is derivable from G

Def: Let G be a CFG with terminal alphabet Σ
 $L(G) \subseteq \Sigma^*$:

$$L(G) = \{ w \in \Sigma^* : w \text{ is derivable from } G \}$$

Ambiguity Def: A derivation is leftmost (resp. rightmost) if in each step, the leftmost (respectively rightmost) nonterminal of the sentential form is replaced.

Fact: If w is derivable, then it has a leftmost and a rightmost derivation.

Ambiguity: A CFG is ambiguous if there is some string in $L(G)$ that has more than one leftmost derivation.

Recall $E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow (E)$
 $E \rightarrow c$ } ambiguous

Let $w = c + c * c$

† leftmost deriv: $E \rightarrow E + E \Rightarrow c + E$
 $\Rightarrow c + E * E$
 $\Rightarrow c + c * E$
 $\Rightarrow c + c * c$

another

$E \Rightarrow E * E \Rightarrow E + E * E$
 $\Rightarrow c + E * E \Rightarrow c + c * E \Rightarrow c + c * c$

Equivalent
 Unambiguous grammars are preferable (if they exist)

An equiv. unambiguous grammar for arith exprs:

$E \rightarrow E + T$	$E = \text{"expression"}$
$E \rightarrow T$	$T = \text{"term"}$
$T \rightarrow T * F$	$F = \text{"factor"}$
$T \rightarrow F$	
$F \rightarrow c$	
$F \rightarrow (E)$	

$E \Rightarrow E + T \Rightarrow E + \overset{T}{\cancel{E}} + T \Rightarrow E + T + T + T \Rightarrow T + T + T + T$
 $T \Rightarrow T * F \Rightarrow T * F * F \Rightarrow F * F * F$

$c + c * c$ — has only 1 leftmost derivation;

⑥

$$E \Rightarrow E + T \Rightarrow T + T \Rightarrow F + T \Rightarrow C + T$$

$$\Rightarrow C + T * F \Rightarrow C + F * F$$

$$\Rightarrow C + C * F \Rightarrow C + C * C$$

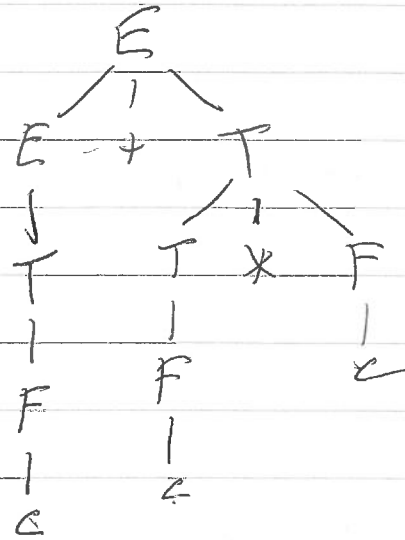
Parse Trees (alternative to derivation)

Ex: for $C + C * C$:

Def: Given a CFG G ,

a parse tree of G

is a rooted ordered tree where



- internal nodes have nonterminal labels
- leaves have terminal labels (or ϵ)
- root labeled by the start symbol
- the child of each internal node (label A , say) read left to right from the body of some production with head A

The ~~yield~~ yield of a parse tree T is the ~~the~~ in-order concatenation of its leaves.