

Extended Free Grammars

Translated and adapted by L. Breveglieri

FREE GRAMMARS EXTENDED BY MEANS OF REGULAR EXPRESSIONS

REGULAR EXPRESSIONS ARE MORE READABLE AND PERSPICUOUS than free grammars, thanks to the iteration operators (star and cross) and to the choice operator (union)

To make free grammars more readable as well, it is often permitted to use regexps also in the right sides of the rules of the grammar. One thus obtains the so-called EXTENDED FREE GRAMMARS (Extended BNF or EBNF). The grammars of the technical languages (programming languages, description languages, etc) are usually presented in an extended form

The EBNF rules can be represented as syntax diagrams. A syntax diagram can be conceived as the flow graph of a syntax analysis algorithm

The family LIB of free languages is closed with respect to regular operators (union, star, concatenation); therefore the extended free grammars have the same generating power as the non-extended ones

EXAMPLE: list of declarations of variable names

character testo1, testo2; *real* temp, result;

integer alfa, beta2, gamma;

$\Sigma = \{c, i, r, v, ', ' , ';'\}$ where v is a variable name

$\left((c | i | r) v (, v)^* ; \right)^+$

$S \rightarrow SE | E \quad E \rightarrow AF; \quad A \rightarrow c | i | r \quad F \rightarrow v, F | v$

The grammar is longer than the regexp, and makes less evident the existence of a hierarchy of two lists. Moreover, the names A, E and F are arbitrary

A FREE EXTENDED GRAMMAR (EBNF) $G = \{ V, \Sigma, P, S \}$ contains exactly $|V|$ rules, each of the form $A \rightarrow \eta$, where η is a regexp over the alphabet $V \cup \Sigma$. For better clarity it is allowed to use derived regular operators, like for instance cross, optionality, repetition, etc

EXAMPLE: a simplified Algol-like language

A more compact but less readable version. The non-terminal B cannot be eliminated, as it is indispensable to generate nested structures

$$B \rightarrow b[D]Ie$$

$$D \rightarrow ((c | i | r)v(,v)^*;)^+$$

$$I \rightarrow F(;F)^*$$

$$F \rightarrow a | B$$

$$B \rightarrow b \left[((c | i | r)v(,v)^*;)^+ \right] F(;F)^* e$$

$$B \rightarrow b((c | i | r)v(,v)^*;)^* F(;F)^* e$$

$$B \rightarrow b((c | i | r)v(,v)^*;)^* (a | B)(; (a | B))^* e$$

DERIVATION AND SYNTAX TREE IN THE EXTENDED GRAMMARS

The right side of an extended rule is a regexp that defines an infinite set of strings

Such strings could be imagined as the right sides of a grammar G' , equivalent to G , except that it contains infinitely many rules

$$A \rightarrow (aB)^+$$

$$A \rightarrow aB \mid aBaB \mid \dots$$

DERIVATION RELATION IN AN EXTENDED GRAMMAR

Similarly, one can define multiple step derivations and the generated language

given the strings $\eta_1 \quad \eta_2 \in (\Sigma \cup V)^*$

say that η_2 *derives* (immediately) from η_1

$\eta_1 \Rightarrow \eta_2$ if the two strings are factored as follows

$\eta_1 = \alpha A \gamma, \quad \eta_2 = \alpha \mathcal{G} \gamma$ and there is a rule

$$A \rightarrow e : \quad e \stackrel{*}{\Rightarrow} \mathcal{G}$$

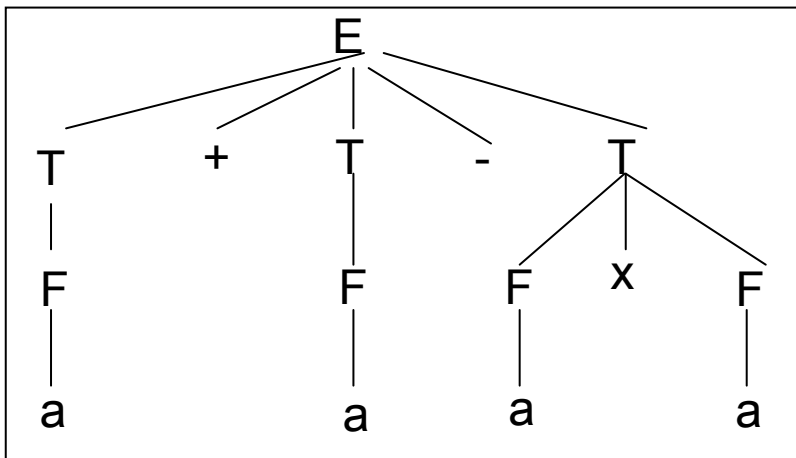
note that $\eta_1 \quad \eta_2 (\quad \alpha A \gamma \quad \alpha \mathcal{G} \gamma)$ do not contain either regular operators or parentheses, and e is a regexp

EXAMPLE: extended derivation for arithmetic expressions

The extended grammar G shown below generates the arithmetic expressions with the four standard operators (in infix notation), round brackets and the symbol a to represent generic variables or constants

$$E \Rightarrow [+|-]T((+|-)T)^* \quad T \rightarrow F((\times|/)F)^* \quad F \rightarrow (a|'('E')')$$

leftmost derivation

$$\begin{aligned} E &\Rightarrow T + T - T \Rightarrow F + T - T \Rightarrow a + T - T \Rightarrow a + F - T \Rightarrow \\ &\Rightarrow a + a - T \Rightarrow a + a - F \times F \Rightarrow a + a - a \times F \Rightarrow a + a - a \times a \end{aligned}$$


In EBNF, the arity of a tree node may be unlimited, in general

However, as the tree gets larger the depth gets lower

AMBIGUITY IN EXTENDED GRAMMARS

A non-extended but ambiguous grammar is ambiguous also when conceived as extended (as BNF is a special case of EBNF)

However, the presence of regexps in the rules may give rise to specific ambiguity forms

$a^*b \mid ab^*$ numbered as $a_1^*b_2 \mid a_3b_4^*$
is ambiguous because ab is
derivable as a_1b_2 or as a_3b_4
 $S \rightarrow a^*b \mid ab^*$ is ambiguous as well

MORE GENERAL GRAMMARS – CHOMSKY CLASSIFICATION

<i>Grammar</i>	<i>Rule Type</i>	<i>Lang. Family</i>	<i>Recognizer Device</i>
Type 0 recursively enumerable	$\beta \rightarrow \alpha \quad \text{where}$ $\alpha, \beta \in (\Sigma \cup V)^+ ;$ $\beta \neq \varepsilon$	R ecursively enumerable languages	T uring machine
Type 1 context- dependent	$\beta \rightarrow \alpha \quad \text{where}$ $\alpha, \beta \in (\Sigma \cup V)^+ ;$ $ \beta \leq \alpha $	C ontext- dependent languages	T uring machine with a memory tape of length equal to that of the string to be recognized

CHOMSKY CLASSIFICATION (continued)

<i>Grammar</i>	<i>Rule Type</i>	<i>Lang. Family</i>	<i>Recognizer Device</i>
Type 2 context- free or BNF form	<div style="border: 1px solid black; padding: 10px;"> $A \rightarrow \alpha \text{ where}$ $A \text{ is nonterm.}$ $\alpha \in (\Sigma \cup V)^*$ </div>	Context-free languages or BNF languages or algebraic languages	Pushdown automaton (non-deterministic)
Type 3 right or left uni-linear	<div style="border: 1px solid black; padding: 10px;"> R-lin: $: A \rightarrow uB$ L-lin: $: A \rightarrow Bu$ where A is nonterm. $u \in \Sigma^* \quad B \in (V \cup \varepsilon)$ </div>	Regular or rational languages	Finite state automaton or sequential machine

The language families mentioned before are in a hierarchic relation of strict containment

Differences: form of the rules and properties of the recognizer devices

type 0, 1 and 2 automaton with unbounded memory

type 3 automaton with finite memory

CLOSURE WITH RESPECT TO:

UNION, CONCAT., STAR, MIRROR, INTERS. WITH REG. LANG.: 0, 1, 2, 3

CLOSURE WITH RESPECT TO COMPLEMENT: 1, 3

TYPE 0: there does not exist any general algorithm to decide whether a string belongs to a given language (the membership problem is semi-decidable)

TYPE 1, 2: there exists a general algorithm to decide whether a string belongs to a given language (the membership problem is semi-decidable)

TYPE 3: it is decidable whether two grammars are weakly equivalent

TYPE 0 e 1: the derivation cannot be represented as a syntax tree; more complex graph structures are necessary

EXAMPLE OF A TYPE 1 GRAMMAR (just for completeness)

$$L = \{a^n b^n c^n \mid n \geq 1\}$$

$$G(\text{ type 1 }): \begin{array}{lll} 1. S \rightarrow aSBC & 3. CB \rightarrow BC & 5. bC \rightarrow bc \\ 2. S \rightarrow abC & 4. bB \rightarrow bb & 6. cC \rightarrow cc \end{array}$$

Differently to what happens in the free grammars, the language generated by means of leftmost derivations is different from that generated by means of rightmost derivations

REPLICAS OR LIST
OF MATCHINGS

$$L_{\text{replic}} = \{uu \mid u \in \Sigma^+\}$$
$$\Sigma = \{a, b\} \quad x = abbbabbb$$

Type 1 grammars are seldom used to design artificial languages for compiler construction (one exception: Algol-68 was defined by means of two-level grammars)

Structures like replicas and matching lists are defined by resorting to semantic methods (attribute grammars or syntax-driven translation)