

## Chomsky Grammars

5.1

In general, to describe a language, there are two possible approaches:

1) recognition: describe rules (or a mechanism) to determine whether or not a certain string belongs to a language

e.g. an automaton is such a mechanism

2) generation: define rules to generate all strings of a language

A grammar is a formalism for defining a language in terms of rules that generate all strings of the language.

Since 1920, various formal methods based on the notion of rewriting or derivation have been proposed by Axel Thue, Emil Post, A.A. Markov.

In the mid 1950s the linguist Noam Chomsky introduced the notion of formal grammar with the aim of formalizing natural language. Formal grammars are in fact too simplistic to capture natural language, but they were adopted as the main formal tool to define syntactic properties of artificial languages (e.g. programming languages)

Definition: Given an alphabet  $\Sigma$ , a (formal) grammar  $G$  (5.2)

is a quadruple  $G = (V_N, V_T, P, S)$  where

- $V_T \subseteq \Sigma$  is a finite nonempty set of symbols called terminals
- $V_N$  is a finite nonempty set of symbols s.t.  $V_N \cap \Sigma = \emptyset$ , called variables or nonterminals, or syntactic categories.

Each variable represents a language

-  $S \in V_N$  is called start symbol or axiom, and represents the language being defined by  $G$

-  $P$  is a binary relation over

$$(V_N \cup V_T)^* \cdot V_N \cdot (V_N \cup V_T)^* \times (V_N \cup V_T)^*$$

Each element  $(\alpha, \beta) \in P$  is called a production or rule, and is generally written as  $\alpha \rightarrow \beta$ .

Note:  $\alpha$  ... sequence of terminals and nonterminals with at least one nonterminal

$\beta$  ... sequence of terminals and nonterminals

Definition: The language  $L(G)$  generated by a grammar  $G$

is the set of strings of terminals only that can be generated starting from the axiom by a finite sequence of rule applications

Each application of a rule  $\alpha \rightarrow \beta$  consists in replacing an occurrence of  $\alpha$  with  $\beta$ .

Example: Palindromes:

A palindrome is a word that reads the same both forwards and backwards. (AILATIITALIA, AMOROMA)

$$L_{\text{pal}} = \{w \in \{0,1\}^* \mid w^R = w\}$$

grammar  $G_{\text{pal}} = (V_N, V_T, P, S)$ , where  $P$  consists of

$$\left. \begin{array}{l} 1) S \rightarrow \epsilon \\ 2) S \rightarrow 0 \\ 3) S \rightarrow 1 \end{array} \right\} \text{basis: } \epsilon, 0, 1 \text{ are palindromes}$$

$$\left. \begin{array}{l} 4) S \rightarrow 0S0 \\ 5) S \rightarrow 1S1 \end{array} \right\} \text{induction: if } S \text{ is a palindrome,} \\ \text{so are } 0S0 \text{ and } 1S1$$

Example of derivation:

$$0110 : S \xrightarrow{4} 0S0 \xrightarrow{5} 01S10 \xrightarrow{1} 0110$$

$$11011 : S \xrightarrow{5} 1S1 \xrightarrow{5} 11S11 \xrightarrow{2} 11011$$

**Exercise E5.1** Prove that the above grammar generates all and only palindromes over  $\{0,1\}$ .

Hint: use induction on the length of the derivation

Example: natural language generation

Sentence  $\rightarrow$  NounPhrase VerbPhrase

NounPhrase  $\rightarrow$  Adjective NounPhrase

NounPhrase  $\rightarrow$  Noun

Noun  $\rightarrow$  car

Noun  $\rightarrow$  train

Adjective  $\rightarrow$  big

Adjective  $\rightarrow$  broken

Notation:

1) To denote the set of productions

$$\alpha \rightarrow \beta_1, \alpha \rightarrow \beta_2, \dots, \alpha \rightarrow \beta_m$$

we use

$$\alpha \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

2) We use  $V = V_N \cup V_T$

A production of the form  $\alpha \rightarrow \epsilon$ , with  $\alpha \in V^* \cdot V_N \cdot V^*$  is called  $\epsilon$ -production.

Example:  $L_{eq} = \{w \in \{0, 1\}^* \mid w \text{ has equal number of } 0\text{'s and } 1\text{'s}\}$

We have already seen that this language is not regular.

Idea to define  $G_{eq}$  s.t.  $L(G_{eq}) = L_{eq}$ : use induction

base:  $\epsilon$  is in  $L_{eq}$

induction:  $\therefore 0w_A \in L_{eq}$  if  $w_A$  has one more 1 than 0

$\therefore 1w_B \in L_{eq}$  if  $w_B$  has one more 0 than 1

Characterize also languages for  $w_A$  and  $w_B$  inductively

grammar  $G_{eq} = (\{S, A, B\}, \{0, 1\}, P, S)$  with P

$$S \rightarrow \epsilon \mid 0A \mid 1B$$

$$A \rightarrow 1S \mid 0AA$$

$$B \rightarrow 0S \mid 1BB$$

(A generates strings with one more 1 than 0.)

B generates strings with one more 0 than 1)

Exercise E5.2 Prove that  $L(G_{eq}) = L_{eq}$  (by induction)

Definition: Given  $G$ , the direct derivation for  $G$  is the binary relation on  $(V^* \circ V_N \circ V^*) \times V^*$  defined as follows:

$(\varphi, \psi)$  is in the relation if there are  $\alpha, \beta, \gamma, \delta \in V^*$  such that  $\varphi = \gamma\alpha\delta$ ,  $\psi = \gamma\beta\delta$  and  $\alpha \rightarrow \beta \in P$ .

We write  $\varphi \Rightarrow \psi$  and say that  $\psi$  directly derives from  $\varphi$  by  $G$ .

Definition: We call derivation the reflexive, transitive closure of direct derivation. In other words,  $\psi$  derives from  $\varphi$  by  $G$ , written  $\varphi \xRightarrow{*} \psi$  if

- a)  $\varphi = \psi$ , or
- b) there are  $\varphi_1, \dots, \varphi_m \in V^*$  such that  $\varphi_1 = \varphi$ ,  $\varphi_m = \psi$ , and  $\varphi_i \Rightarrow \varphi_{i+1}$ ,  $\forall i, 1 \leq i < m$

Definition: Given a grammar  $G$ , the language generated by  $G$  is

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

Notice: words in  $L(G)$  are constituted by terminals only.

Terminology:

- sentence: any word  $w \in V_T^*$  s.t.  $S \xRightarrow{*} w$ , i.e.  $w \in L(G)$
- sentential form: any  $\alpha \in V^* = (V_T \cup V_N)^*$  s.t.  $S \xRightarrow{*} \alpha$

Notation: terminals:  $a, b, c, \dots$   
 nonterminals:  $A, B, C, \dots$   
 strings of terminals:  $u, v, w, x, y, z, \dots$   
 symbols of  $V = V_T \cup V_N$ :  $X, Y, Z, \dots$   
 sentential forms:  $\alpha, \beta, \gamma, \dots$

Example: Productions for  $G_{eq}$ :

- $S \rightarrow \epsilon \mid 0A \mid 1B$
- $A \rightarrow 1S \mid 0AA$
- $B \rightarrow 0S \mid 1BB$

derivation:

- 1)  $001SA \Rightarrow 001S1S$  (using  $A \rightarrow 1S$ )
- 2)  $001S1S \Rightarrow 0011S$  (using  $S \rightarrow \epsilon$ )
- 3)  $001SA \xrightarrow{*} 0011S$  (using (1) and (2))
- 4)  $S \xrightarrow{*} 001110$

Example: Grammar for  $L_{3n} = \{e^n b^n c^n \mid n \geq 1\}$

18/12/2008

$G_{3n} = (\{A, B, C, S\}, \{e, b, c\}, P, S)$

with P

- 1)  $S \rightarrow eSBC$
  - 2)  $S \rightarrow eBC$
  - 3)  $CB \rightarrow BC$
  - 4)  $eB \rightarrow eb$
  - 5)  $bB \rightarrow bb$
  - 6)  $bC \rightarrow bc$
  - 7)  $cC \rightarrow cc$
- } generate  $ee \dots eBCBC \dots BC$
- } moves the C's to the end
- } generate the terminals from left to right
- Note: we cannot simply have  $B \rightarrow b, C \rightarrow c$  because this would generate many words not in  $L_{eq}$

Example of derivation of  $eaabbbccc$ :

$$\begin{aligned}
 S &\xrightarrow{1} eSBC \xrightarrow{2} eeSBCBC \xrightarrow{2} eeeBCBCBC \\
 &\xrightarrow{3} eeeBCBBCC \xrightarrow{3} eeeBBCBCC \\
 &\xrightarrow{3} eeeBBBCCC \xrightarrow{4} eeebBBCCC \\
 &\xrightarrow{5} eeebbBBCCC \xrightarrow{5} eeebbbCCC \\
 &\xrightarrow{6} eeebbbC \xrightarrow{2} eeebbbCC \\
 &\xrightarrow{7} eeebbbccc
 \end{aligned}$$

Note: not each sequence of direct derivations leads to a sentence in  $L(G_{3m})$

e.g. with the previous grammar we could generate

$$\begin{aligned} S &\Rightarrow \underline{a} SBC \Rightarrow aa \underline{S}BCBC \Rightarrow aaa \underline{B}C \underline{B}C \underline{B}C \\ &\Rightarrow aaa \underline{B}C \underline{B}BCC \Rightarrow aaa \underline{b}C \underline{B}BCC \\ &\Rightarrow aaa bC \underline{B}BCC \end{aligned}$$

we cannot apply any other production

Also, the application of productions could go on forever (e.g. rule 1 in the previous example)

Classification of Chomsky grammars into 4 groups, depending on the form of the productions:

- grammars of type 0 : no limitations
- " " " " 1 : context-sensitive
- " " " " 2 : context-free
- " " " " 3 : regular (or right linear)

Definition: grammar of type 0.

Productions have the most general form  $\alpha \rightarrow \beta$ ,

with  $\alpha \in V^+ = V_N^+ = V^*$   $\beta \in V^*$

Grammars of type 0 allow for derivations that shorten the sentential form.

A language generated by a grammar of type 0 is called of type 0.

Definition: grammar of type 1, or context sensitive

Productions have the form  $\alpha \rightarrow \beta$ , with

$$\alpha \in V^* \cdot V_N \cdot V^*, \quad \beta \in V^+, \quad |\alpha| \leq |\beta|$$

These productions cannot shorten the length of the sentential form to which they are applied.

A language generated by a grammar of type 1 is called of type 1, or context sensitive.

Example:  $G_{3n}$  is context sensitive. Obviously, it is also of type 0.

Definition: grammar of type 2, or context-free

Productions have the form  $A \rightarrow \beta$ , with  $A \in V_N, \beta \in V^+$ .

These productions are productions of type 1, with the additional requirement that on the left there is a single nonterminal.

A language generated by a grammar of type 2 is called of type 2, or context free.

Example  $L_{2n} = \{a^n b^n \mid n \geq 1\}$  is of type 1, since the

following grammar  $G'_{2n}$  generates  $L_{2n}$

$$S \rightarrow aB \mid SAB$$

$$BA \rightarrow AB$$

$$aA \rightarrow aa$$

$$aB \rightarrow ab$$

$$bB \rightarrow bb$$

$L_{2n}$  is also of type 2, since it is generated by

$$S \rightarrow aSb \mid ab$$



↓ We said that grammars of type 1 are also called context-sensitive (in contrast to context-free grammar). This is justified by the original definition by Chomsky for context-sensitive grammars:

Definition: Chomsky CS-grammar

Productions have the form  $\varphi_1 A \varphi_2 \rightarrow \varphi_1 \beta \varphi_2$   
with  $\varphi_1, \varphi_2 \in V^*$ ,  $A \in V_N$ ,  $\beta \in V^+$

Intuitively,  $A$  is replaced by  $\beta$  only if it appears "in the context" of  $\varphi_1$  and  $\varphi_2$

Theorem: Grammars of type 1 and Chomsky CS grammars generate the same class of languages

Proof: We show that, for every language  $L$ :

There is a type-1 grammar  $G_1$  s.t.  $L = \mathcal{L}(G_1)$  iff

there is a Chomsky CS grammar  $G_c$  s.t.  $L = \mathcal{L}(G_c)$

"if" immediate, since each Chomsky CS grammar is of type 1

(in  $\varphi_1 A \varphi_2 \rightarrow \varphi_1 \beta \varphi_2$  we have  $\beta \in V^+$  and hence

$$|\varphi_1 A \varphi_2| \leq |\varphi_1 \beta \varphi_2|)$$

"only-if": let  $G_1$  be a type-1 grammar for  $L$ .

We construct from  $G_1$  a Chomsky CS grammar  $G_c$  as follows:

- 1) for each  $a \in V_T$ , add a new nonterminal  $N_a$ ,
- 2) replace in each production of  $G_1$ , each  $a \in V_T$  by  $N_a$

Now all productions have the form

$$A_1 A_2 \dots A_m \rightarrow B_1 B_2 \dots B_m \text{ with } m \leq n$$

and all  $A_i, B_j \in V_N$

3) For each such production  $A_1 \dots A_m \rightarrow B_1 \dots B_n$ , introduce a new nonterminal  $N$ , and replace the production by the following ones:

$$\begin{aligned}
A_1 A_2 \dots A_m &\rightarrow N A_2 \dots A_m \\
N A_2 \dots A_m &\rightarrow N B_2 A_3 \dots A_m \\
N B_2 A_3 \dots A_m &\rightarrow N B_2 B_3 A_4 \dots A_m \\
&\vdots \\
N B_2 \dots B_{m-1} A_m &\rightarrow N B_2 \dots B_{m-1} B_m \dots B_n \\
N B_2 \dots B_m &\rightarrow B_1 B_2 \dots B_n
\end{aligned}$$

(note that, due to the presence of  $N$ , these productions will not "interfere" with other ones)

Observe that all such productions are of the form

$$\gamma_1 A \gamma_2 \rightarrow \gamma_1 \beta \gamma_2 \quad \text{with } \gamma_1, \gamma_2 \in V^*, A \in V_N, \beta \in V^+$$

4) For each  $a \in V_T$ , add the production

$$N_a \rightarrow a \quad (\text{where } N_a \text{ is the new non-terminal associated to } a)$$

It is not difficult to see that  $\mathcal{L}(G_1) = \mathcal{L}(G_c)$

(the proof is by induction on the length of the derivation of a string  $w \in \mathcal{L}(G_1)$  - (resp.,  $\mathcal{L}(G_c)$ )

↑ END OPTIONAL

(5.11)

Definition: grammar of type 3, or regular, or right linear

Productions have the form  $A \rightarrow \delta$  with  $A \in V_N$   
 $\delta \in V_T \cup (V_T \circ V_N)$

(i.e.,  $A \rightarrow eB$  or  $A \rightarrow e$ , with  $A, B \in V_N, e \in V_T$ )

A language generated by a grammar of type 3 is called of type 3 or regular

Example:  $\{e^n b \mid n \geq 0\}$  is of type 3, since it is generated by the grammar

$$\begin{aligned} S &\rightarrow eS \\ S &\rightarrow b \end{aligned}$$

Note: a grammar of type 3 is called linear, because on the right hand side of a production there is at most one non-terminal. It is called right-linear because the non-terminal is on the right of the terminal

↓ OPTIONAL

Exercise: E5.3: Show that grammars of type 3 generate the class of regular languages that do not contain  $\epsilon$ .

Hint: given  $G = (V_N, V_T, P, S)$ , construct an NFA

$A_G = (V_N \cup \{F\}, V_T, \delta, S, \{F\})$  with

$B \in \delta(A, e)$  iff  $A \rightarrow eB$  and

$F \in \delta(A, e)$  iff  $A \rightarrow e$

Show by induction on  $|w|$  that  $w \in \mathcal{L}(A_G)$  iff  $w \in \mathcal{L}(G)$ .

Conversely, given an NFA  $A$ , construct a grammar  $G_A$  by having again nonterminals correspond to states of  $A$ .

↑ END OPTIONAL

Note on  $\epsilon$ -productions (for grammars of type 1, 2, 3)

As we have defined them, grammars of type 1 (resp. 2, 3) cannot generate the empty string  $\epsilon$ .

We could extend the definition by allowing also the generation of  $\epsilon$ :

- if the start symbol  $S$  does not appear on the right-hand side of productions, we allow also for a production  $S \rightarrow \epsilon$  ( $\epsilon$ -production)
- if the start symbol  $S$  appears on the right-hand side of productions, we introduce a new non-terminal  $S_{new}$ , make it the new start symbol, add a production  $S_{new} \rightarrow S$ , and allow for  $S_{new} \rightarrow \epsilon$ .

Hence, an  $\epsilon$ -production used just to generate  $\epsilon$  is harmless.

Note that, allowing for  $\epsilon$ -productions for every non-terminal is not that harmless.

OPTIONAL  
↓

Exercise: E5.4: Show that, for every language  $L$  of type 0

there is a grammar of type 1 extended with  $\epsilon$ -productions on arbitrary non-terminals that generates  $L$ .

Hint: introduce a new nonterminal  $N_\epsilon$  that is eliminated through an  $\epsilon$ -production  $N_\epsilon \rightarrow \epsilon$ , and use  $N_\epsilon$  to make the right-hand side of productions as long as the left-hand side.

↑ END OPTIONAL

# Context-free grammars (CFGs)

5.13

In a CFG, the productions have the form  $A \rightarrow B$  with  $A \in V_N$ ,  $B \in V^*$  (note: we allow for  $\epsilon$ -productions)

Example: CFG for arithmetic expressions over variable  $i$

$G = (\{E, T, F\}, \{i, +, *, (, )\}, P, E)$ , where  $P$  is

$E \rightarrow T \mid E + T$        $E$  ... expression

$T \rightarrow F \mid T * F$        $T$  ... term

$F \rightarrow i \mid (E)$        $F$  ... factor

This grammar generates, e.g.,  $i + i * i$

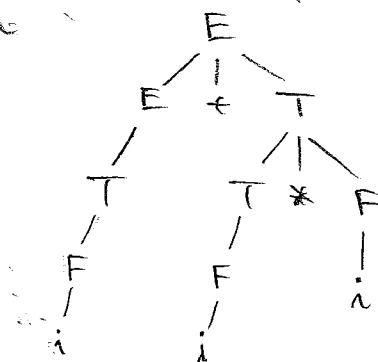
$\underline{E} \rightarrow \underline{E} + T \rightarrow \underline{T} + T \rightarrow \underline{F} + T \rightarrow i + \underline{T} \rightarrow$   
 $\rightarrow i + \underline{T} * F \rightarrow i + i * \underline{E} \rightarrow i + i * i$

We can also represent a derivation of a string by a CFG by means of a tree, called parse-tree:

Is a tree whose nodes are labeled by elements of  $V \cup \{\epsilon\}$  satisfying:

- 1) each interior node is labeled by a non-terminal
- 2) each leaf is labeled by a non-terminal, a terminal, or  $\epsilon$ . If it is labeled by  $\epsilon$ , then it is the only child of its parent
- 3) If an interior node is labeled  $A$ , and its children from left to right are labeled  $X_1, X_2, \dots, X_k$ , then there is a production  $A \rightarrow X_1 X_2 \dots X_k$  in  $P$ .

Example: parse tree for  $i + i * i$

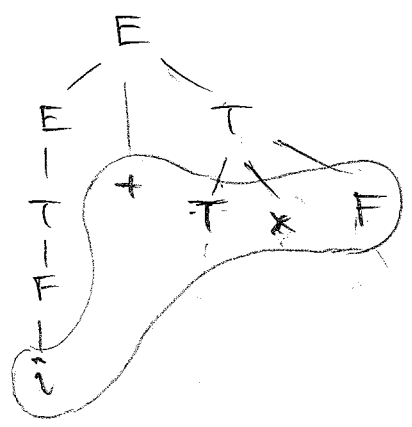


We call A-tree a subtree of the parse-tree rooted at non-terminal A.

Yield (or frontier) of a tree:

is the sequence of labels of the leaves from left-to-right.

Example:



Theorem:  $\alpha \in V^+$  is the yield of an A-tree  $\iff A \Rightarrow^* \alpha$

Proof: by induction on the height of the tree  
(see textbook)

8/1/2003

Note: a parse tree does not specify a unique way to derive  $\alpha$  from A. (the order in which non-terminals are expanded is not specified).

The parse tree specifies, however, which rule is applied for each non-terminal.

Specific derivation orders:

- leftmost derivation: obtained by traversing the tree depth-first, by first going to the left subtree, and then to the right one.

e.g.  $E \xrightarrow{lm} E + T \xrightarrow{lm} i + T \xrightarrow{lm} i + T * F \xrightarrow{lm} i + T * F \dots$

- rightmost derivation: defined similarly:  $E \xrightarrow{rm} E + T \xrightarrow{rm} E + T * F \dots$

Theorem: the following are all equivalent statements for (5.15)

a CFG  $G = (V, T, P, S)$  and a string  $w \in T^*$

1)  $w \in \mathcal{L}(G)$  (or  $S \xRightarrow{*} w$ )

2)  $S \xrightarrow{lm}^* w$

3)  $S \xrightarrow{rm}^* w$

4) There exists an  $S$ -tree with yield  $w$ .

Proof: the equivalence of (1) and (4) follows from the previous theorem. The other equivalences are obvious.

Thus, we could always use  $lm$ -derivation as a canonical way to derive any  $w \in \mathcal{L}(G)$ ; i.e. as a canonical way to interpret a parse tree for  $w$ .

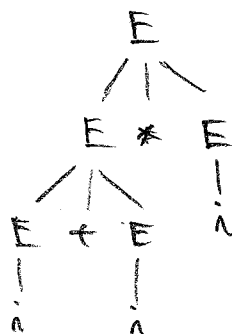
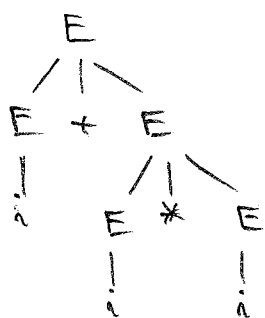
Ambiguous grammars:

$\exists w \in \mathcal{L}(G)$  could have two distinct parse trees, and hence two distinct  $lm$ -derivations

Example: another grammar for arithmetic expressions

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

$w = i + i * i$



These parse trees correspond to two different  $lm$ -derivations, and also to two ways of interpreting  $w$ .

Definition: A CFG  $G$  is ambiguous if for some  $w \in L(G)$  there exist two distinct parse trees.

Ambiguity has to be avoided in compilers, since it corresponds to different ways of interpreting a string.

Sometimes grammars can be redesigned to remove ambiguity. (e.g., for arithmetic expressions)

This is not always possible:

Definition: A CF language is (inherently) ambiguous if all its grammars are ambiguous

Example:  $L = \{a^n b^m c^m d^m \mid n, m \geq 1\} \cup \{a^m b^n c^n d^n \mid n, m \geq 1\}$

$L$  is CF (show for exercise)

Consider strings of the form  $a^k b^k c^k d^k$ .

We cannot tell whether they come from first or second types of strings in  $L$ , and any CFG must allow for both possibilities.



# Properties of context-free languages

5.17

We will study

- 1) Normal forms for CFGs (useful for proving properties of CFLs)
- 2) Expressive power  $\Rightarrow$  pumping lemma for CFLs
- 3) Closure and decision properties

## Normal forms for CFGs

We look at how to simplify CFGs, while preserving the generated language.

- gain efficiency in parsing
- simplify proving properties

### 1) Eliminate useless symbols:

We say that  $X \in V$  is useful if

$$S \Rightarrow^* \alpha X \beta \Rightarrow^* w \quad \text{with } w \in V_T^* \\ \alpha, \beta \in V^*$$

Thus, a symbol is useless (not useful) if it does not participate in any derivation of strings of the language.

$\Rightarrow$  can be eliminated

Definition:  $X \in V$  is generating if  $X \Rightarrow^* w$ , for  $w \in V_T^*$

$X \in V$  is reachable if  $S \Rightarrow^* \alpha X \beta$ , for  $\alpha, \beta \in V^*$

Hence,  $X$  is useful, if it is both generating and reachable.

We identify useless symbols by

1a) eliminating non-generating symbols and all their productions

1b) --- unreachable ---

Note: it is important to do these two steps in the above order

Example:  $\begin{cases} S \rightarrow AB \mid b \\ A \rightarrow \epsilon \end{cases}$  Let us consider what happens if we do first (1b) and then (1a)

• we eliminate unreachable symbols: all are reachable  
• --- non-generating ---

we eliminate B and  $S \rightarrow AB$

$\Rightarrow$  we obtain:  $S \rightarrow b$   
 $A \rightarrow \epsilon$

But if we do it in the right order:

1a) eliminate non-generating symbols: B and  $S \rightarrow AB$

1b) --- unreachable --- : A and  $A \rightarrow \epsilon$

$\Rightarrow$  We obtain:  $S \rightarrow b$

9a) Eliminating non-generating symbols:

We construct the set H of generating symbols, and then eliminate the symbols not in H and the productions containing them.

Algorithm to construct the set H of generating symbols

Input: grammar  $G = (V_N, V_T, P, S)$

Output: set H of generating symbols

$H \leftarrow V_T$

while there is a change in H do

for each production  $A \rightarrow X_1 \dots X_k$  in P do

if  $\{X_1, \dots, X_k\} \subseteq H$  (i.e., all of  $X_1, \dots, X_k$  are generating)

then  $H \leftarrow H \cup \{A\}$

return H

Example:  $G_1 = (V_N, V_T, P, S)$

with  $V_N = \{S, A, B, C, D\}$ ,  $V_T = \{a, b, c, d\}$ ,

$P$ :

$$\begin{cases} S \rightarrow AB \mid AC \mid CD \\ A \rightarrow BB \\ B \rightarrow AC \mid ab \\ C \rightarrow Ce \mid CC \\ D \rightarrow Bc \mid b \mid d \end{cases}$$

initialization:  $H = \{e, b, c, d\}$

iteration 1:  $H \leftarrow H \cup \{B, D\}$

- " - 2:  $H \leftarrow H \cup \{A\}$

- " - 3:  $H \leftarrow H \cup \{S\}$

- " - 4:  $H$  does not change

$C$  is not generating  $\Rightarrow$  Remove  $C$  and all productions containing  $C$

### 1b) Eliminating unreachable symbols

We construct the set  $R$  of reachable symbols, and then eliminate the symbols not in  $R$  and the productions containing them.

Algorithm to construct the set  $R$  of reachable symbols

Input: grammar  $G = (V_N, V_T, P, S)$

Output: set  $R$  of reachable symbols

$R \leftarrow \{S\}$

while there is a change in  $R$  do

  for each production  $A \rightarrow X_1 \dots X_k$  in  $P$  do

    if  $A \in R$  (i.e.,  $A$  is reachable)

      then  $H \leftarrow H \cup \{X_1, \dots, X_k\}$

return  $R$

Example:  $G_2 = (V_N, V_T, P, S)$  with  $V_N = \{S, A, B, D\}$ ,  $V_T = \{a, b, c, d\}$

$P$ :

$$\begin{cases} S \rightarrow AB \\ A \rightarrow BB \\ B \rightarrow ab \\ D \rightarrow b \mid d \end{cases}$$

initialization:  $R = \{S\}$

iteration 1:  $R \leftarrow R \cup \{A, B\}$

- " - 2:  $R \leftarrow R \cup \{e, b\}$

- " - 3:  $R$  does not change

$D, c, d$  are unreachable

$\Rightarrow$  Remove  $D, c, d$  and all productions containing them.

## 2) Eliminate $\epsilon$ -productions

$\epsilon$ -production:  $A \rightarrow \epsilon$  slows down parsing

Definition:  $A \in V_N$  is nullable if  $A \Rightarrow^* \epsilon$

We first need to find all nullable symbols

Algorithm to construct the set  $N$  of nullable symbols

Input: grammar  $G = (V_N, V_T, P, S)$

Output: set  $N$  of nullable symbols

$N = \emptyset$

for each production  $A \rightarrow \epsilon$  in  $P$  do  $N \leftarrow N \cup \{A\}$

while there is a change in  $N$  do

for each production  $A \rightarrow X_1 \dots X_k$  in  $P$  do

if  $\{X_1, \dots, X_k\} \subseteq N$  (i.e., all of  $X_1, \dots, X_k$  are nullable)  
then  $N \leftarrow N \cup \{A\}$

return  $N$

Example:  $G_3 = (V_N, V_T, P, S)$  with  $V_N = \{S, A, B, C\}$ ,  $V_T = \{a, b\}$

$P: \begin{cases} S \rightarrow ABC \mid BCB \\ A \rightarrow aB \mid a \\ B \rightarrow CC \mid b \\ C \rightarrow S \mid \epsilon \end{cases}$

initialization:  $N = \{C\}$   
iteration 1:  $N \leftarrow N \cup \{B\}$   
" 2:  $N \leftarrow N \cup \{S\}$   
" 3:  $N$  does not change

Knowing the nullable symbols, allows us to compensate for the elimination of  $\epsilon$ -productions.

Example: in  $G_3$ , since  $B$  and  $C$  are nullable, we can derive:

$S \Rightarrow^* BCB, S \Rightarrow^* CB, S \Rightarrow^* BC, S \Rightarrow^* BB$   
 $S \Rightarrow^* B, S \Rightarrow^* C, S \Rightarrow^* \epsilon$

Hence, if we eliminate  $C \rightarrow \epsilon$  (and  $B, C$  are not nullable anymore), we have to add direct productions for the above derivations.

Algorithm to eliminate  $\epsilon$ -productions

- 1) Identify all nullable symbols
- 2) Replace each production  $A \rightarrow X_1 \dots X_n$  by the set of all productions of the form  $A \rightarrow \alpha_1 \dots \alpha_n$  where  $\alpha_i = X_i$ , if  $X_i$  is not nullable  
 $\alpha_i = X_i$  or  $\epsilon$ , if  $X_i$  is nullable
- 3) If the resulting grammar contains  $S \rightarrow \epsilon$ , introduce a new start symbol  $S'$  and add the productions  $S' \rightarrow S | \epsilon$
- 4) Remove all  $\epsilon$ -productions, except possibly the one for  $S'$ .

Example: by applying steps (1) and (2) to  $G_3$ , we get

$$\left\{ \begin{array}{l} S \rightarrow ABC \mid AB \mid AC \mid A \mid \\ \quad BCB \mid BC \mid BB \mid CB \mid B \mid C \mid \epsilon \\ A \rightarrow aB \mid \epsilon \\ B \rightarrow CC \mid C \mid \epsilon \mid b \\ C \rightarrow S \mid \epsilon \end{array} \right.$$

Since we have  $S \rightarrow \epsilon$ , we add  $S' \rightarrow S | \epsilon$  and remove the remaining  $\epsilon$ -productions.

Eliminate unit productions

Unit-production:  $A \rightarrow B$  slows down parsing

Algorithm to eliminate unit-productions

- 1) Remove  $\epsilon$ -productions
- 2) For all  $A, B \in V_N$   
 if  $A \Rightarrow^* B$  and  $B \rightarrow \alpha$  is not unit  
 then add  $A \rightarrow \alpha$
- 3) Eliminate all unit-productions

How do we determine whether  $A \Rightarrow^* B$  holds?

5.22

Since we have no  $\epsilon$ -productions, we have that

$A \Rightarrow^* B$  if and only if

$$A \Rightarrow B_1 \Rightarrow B_2 \Rightarrow \dots \Rightarrow B_{k-1} \Rightarrow B_k \Rightarrow B$$

where all  $B_i$ 's are pairwise distinct. Hence  $k \leq |V_M|$ .

(if we had a sequence where two  $B_i$ 's are the same, we could eliminate all steps inbetween, and get a new sequence where all  $B_i$ 's are pairwise distinct)

Each single derivation step  $B_i \Rightarrow B_{i+1}$  must correspond to a unit production  $B_i \rightarrow B_{i+1}$  of  $G$ .

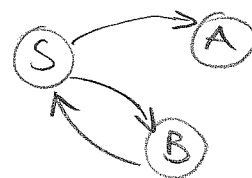
Hence, we can detect whether  $A \Rightarrow^* B$  by checking whether  $B$  is reachable from  $A$  in the graph of the unit productions:

- nodes: non-terminals

- edges: one edge  $(A) \rightarrow (B)$  for each unit prod.  $A \rightarrow B$

Example:  $G_1 = \begin{cases} S \rightarrow A \mid B \\ A \rightarrow S a \mid \epsilon \\ B \rightarrow S \mid b \end{cases}$

Graph of unit productions



reachability:  $S \Rightarrow^* A$        $S \Rightarrow^* B$   
 $B \Rightarrow^* S$        $B \Rightarrow^* A$

we get:  $\begin{cases} S \rightarrow A \mid B \mid S a \mid \epsilon \mid S \mid b \\ A \rightarrow S a \mid a \\ B \rightarrow S \mid b \mid A \mid B \mid S a \mid \epsilon \end{cases}$

Removing unit productions, we get  $\begin{cases} S \rightarrow S a \mid \epsilon \mid b \\ A \rightarrow S a \mid \epsilon \\ B \rightarrow S a \mid \epsilon \mid b \end{cases}$

Note: A and B have become unreachable

We have seen: removal of: useless symbols,  $\epsilon$ -prod, unit-prod. (5.23)

Does the order of the steps matter?

Observation:

- removing useless: does not add productions at all symbols (and therefore not  $\epsilon$ -prod. or unit-prod.)
- removing  $\epsilon$ -prod: could add unit-prod
- removing unit-prod: - needs removing  $\epsilon$ -prod first  
- could create useless symbols  
- cannot create  $\epsilon$ -prod.

$\Rightarrow$  The right order for removal is

- 1)  $\epsilon$ -productions
- 2) unit-productions
- 3) useless symbols: first non-generating then unreachable

### Chomsky Normal Form

Definition: A CFG  $G$  is in Chomsky Normal (CNF) if all its productions are of the form

$$\begin{array}{l} A \rightarrow a \\ A \rightarrow BC \end{array} \quad \text{with} \quad \begin{array}{l} a \in V_T \\ A, B, C \in V_N \end{array}$$

Given a CFG  $G$ , we can always construct a CFG  $G_C$  that is in CNF and such that  $L(G_C) = L(G) \setminus \{\epsilon\}$ .

Note: since a CFG in CNF cannot generate  $\epsilon$ , if  $G$  generates  $\epsilon$ , then we cannot have that  $L(G_C) = L(G)$ . However, apart from  $\epsilon$ , the two languages are equal.

Starting from  $G_1$ , we construct  $G_c$  in several steps:

- 1) Eliminate  $\epsilon$ -productions (without introducing the new start symbol  $S'$  with  $S' \rightarrow S | \epsilon$ )
- 2) Eliminate unit-productions

$\Rightarrow$  All productions are of the form

$$A \rightarrow a$$

$$A \rightarrow X_1 \dots X_k \quad \text{with } k \geq 2$$

with  $A \in V_N$ ,  $a \in V_T$ ,  $X_1, \dots, X_k \in V$

- 3) Remove non-generating symbols
- 4) Remove unreachable symbols

- 5) Remove "injected bodies"

for each  $a \in V_T$ , add a new nonterminal  $N_a$  and production  $N_a \rightarrow a$

in each production  $A \rightarrow X_1 \dots X_k$ , replace  $a$  with  $N_a$

$\Rightarrow$  All productions are of the form

$$A \rightarrow a$$

$$A \rightarrow A_1 \dots A_k \quad (k \geq 2)$$

with  $a \in V_T$ ,  $A, A_1, \dots, A_k \in V_N$

- 6) "Factor" long productions

for each  $A \rightarrow A_1 \dots A_k$  with  $k \geq 3$

- add new nonterminals  $B_1, \dots, B_{k-2}$

- replace  $A \rightarrow A_1 \dots A_k$

with  $A \rightarrow A_1 B_1$

$$B_1 \rightarrow A_2 B_2$$

$\vdots$

$$B_{k-2} \rightarrow A_{k-1} A_k$$

The grammar we get is in CNF by construction.

It is easy to show that the language is preserved, apart possibly for the empty string  $\epsilon$ , which cannot be generated by a grammar in CNF.



Example:  $G_1$   $\left\{ \begin{array}{l} S \rightarrow ABB \mid \epsilon b \\ A \rightarrow Be \mid be \\ B \rightarrow \epsilon A b B \mid b \end{array} \right.$

Steps 1-4: nothing to do

Step 5:  $\left\{ \begin{array}{l} N_e \rightarrow \epsilon \\ N_b \rightarrow b \\ S \rightarrow ABB \mid N_e N_b \\ A \rightarrow B N_e \mid N_b N_e \\ B \rightarrow N_e A N_b B \mid b \end{array} \right.$

Step 6:  $\left\{ \begin{array}{l} N_e \rightarrow \epsilon \\ N_b \rightarrow b \\ S \rightarrow A B_1 \mid N_e N_b \\ B_1 \rightarrow BB \\ A \rightarrow B N_e \mid N_b N_e \\ B \rightarrow N_e C_1 \mid b \\ C_1 \rightarrow A C_2 \\ C_2 \rightarrow N_b B \end{array} \right.$

Note: If the original grammar generated  $\epsilon$ , and we want that the grammar in CNF also generates  $\epsilon$ , we can execute step 1 by introducing the new start symbol  $S'$  and the productions:  $S' \rightarrow S \mid \epsilon$

Step 2 will then replace the unit production  $S' \rightarrow S$  by other productions, but none of the transformations 2-5 will introduce a production where  $S'$  is on the right side.

Therefore, in the end, we will have a grammar in CNF, apart for the production  $S' \rightarrow \epsilon$ .