



Normal Forms for Context-Free Grammars

Robert M. Keller
Harvey Mudd College
28 September 2003



Uses of Normal Forms

- Certain proofs and constructions are simpler if the grammar can be assumed to be in a specific form.
- An example is the proof of the pumping lemma, which used the Chomsky Normal Form.

Chomsky Normal Form

- Each production has one of the following forms:
 - $A \rightarrow BC$, where A , B , and C are auxiliaries (not necessarily distinct)
 - $A \rightarrow \alpha$, where α is a terminal symbol
- Noam Chomsky is a professor
- of linguistics at M.I.T.



Greibach Normal Form

- Each production has the following form:
 - $A \rightarrow \alpha \beta$ where α is a terminal symbol and β is a string of auxiliary symbols (possibly empty).
- Sheila Greibach is a professor of computer science at UCLA.
- [Sheila A. Greibach, A New Normal-Form Theorem for Context-Free Phrase Structure Grammars, Journal of the ACM \(JACM\), v.12 n.1, p.42-52, Jan. 1965](#)



Example: Normal Forms for the balanced parenthesis language

- Chomsky:

$S \rightarrow LT$

$T \rightarrow SR$

$S \rightarrow LR$

$S \rightarrow SS$

$L \rightarrow ($

$R \rightarrow)$

- Greibach:

$S \rightarrow (T$

$S \rightarrow (ST$

$S \rightarrow (TS$

$S \rightarrow (STS$

$T \rightarrow)$



The ϵ Problem

- Neither Chomsky nor Greibach forms can generate ϵ , so we only apply them to languages not containing ϵ .
- This is not a major issue, since if a language has ϵ as an element, we can always generate it by the single production
$$S \rightarrow \epsilon$$
where S is the start symbol and S does not appear on the RHS of any production.
- In the following discussion, we'll assume languages don't contain ϵ unless noted.



Chomsky Normal Form Lemma

- Every context-free language (not containing ϵ) is generated by a grammar in Chomsky Normal Form.



Grammar Cleanup Lemma

- For every context-free grammar (generating a language not containing ϵ), there is an equivalent grammar **not** containing either of the following two types of productions:
 - $A \rightarrow \epsilon$ (called **ϵ productions**)
 - $A \rightarrow B$ (called **unit productions**)




Proof of the Grammar Cleanup Lemma (following Dexter Kozen's book)

- Starting with an arbitrary grammar (without $S \rightarrow \epsilon$), iteratively add additional productions using the following two rules:
 - If $B \rightarrow \epsilon$, then for any production of the form $A \rightarrow \alpha B \beta$, add a production $A \rightarrow \alpha \beta$.
 - If $B \rightarrow C$, then for any production of the form $A \rightarrow \alpha B \beta$, add a production $A \rightarrow \alpha C \beta$.
- The productions added have RHS's no longer than those of the original productions. Hence this process must terminate, because the number of different such strings is finite.
- Adding these productions does not change the language generated.

Continuation of the Proof of the Grammar Cleanup Lemma

- Claim: For any derivation of a terminal string in the modified grammar, there is a derivation not using ϵ - or unit productions.
- Proof: Suppose there is a derivation of a terminal string using a ϵ -production $B \rightarrow \epsilon$:

$$S \xrightarrow{*} \epsilon B \xrightarrow{\epsilon} \epsilon \epsilon \xrightarrow{*} x.$$
 Then B is not S and must have been introduced by a production of the form $A \rightarrow \epsilon B$. But by our rules for augmenting the grammar, we get an equivalent, but shorter, derivation using the rule $A \rightarrow \epsilon \epsilon$ instead.



Continuation (2) of the Proof of the Grammar Cleanup Lemma

- Suppose there is a derivation of a terminal string using a unit production $B \rightarrow C$:

$$S \xrightarrow{*} \alpha B \beta \xrightarrow{\rightarrow C} \alpha C \beta \xrightarrow{*} x.$$

Then B must either be S or have been introduced by a production of the form $A \rightarrow \alpha B \beta$. But by our rules for augmenting the grammar, we get an equivalent derivation of the same length using the rule $A \rightarrow \alpha C \beta$ instead.

- Since neither \rightarrow - or unit productions are needed in derivations in the augmented grammar, all such productions can now be eliminated without changing the language generated.



Example of Grammar Cleanup

- Consider the balanced-parenthesis grammar

$$S \rightarrow (S)$$
$$S \rightarrow SS$$
$$S \rightarrow \epsilon$$

- Augmenting the grammar introduces:

$$S \rightarrow ()$$
$$S \rightarrow S$$

- Removing the ϵ - and unit productions gives:

$$S \rightarrow (S)$$
$$S \rightarrow ()$$
$$S \rightarrow SS$$

This generates the same language without ϵ .



Proof of the Chomsky Normal Form Lemma

- Start with a cleaned-up grammar.
- For every symbol α in the terminal alphabet occurring on some RHS of a production, introduce an auxiliary A_α and a production $A_\alpha \rightarrow \alpha$. Replace α on the RHS with A_α .
- For any production with a RHS longer than 2:
 $A \rightarrow B_1 B_2 B_3 \dots B_n,$
add new auxiliaries C_2, C_3, \dots, C_{n-1} and replace the production with productions
 $A \rightarrow B_1 C_2, C_2 \rightarrow B_2 C_3, \dots, C_{n-1} \rightarrow B_{n-1} B_n$
- The resulting grammar generates the same language and is in Chomsky Normal Form.

Example of the Chomsky Normal Form Lemma

- Start with the cleaned-up grammar:

$$\begin{aligned} S &\rightarrow (S) \\ S &\rightarrow () \\ S &\rightarrow SS \end{aligned}$$

- Add new symbols L and R for (and):

$$\begin{aligned} S &\rightarrow LSR & L &\rightarrow (\\ S &\rightarrow LR & R &\rightarrow) \\ S &\rightarrow SS \end{aligned}$$

- Replace $S \rightarrow LSR$:

$$\begin{aligned} S &\rightarrow LT & L &\rightarrow (\\ T &\rightarrow SR & R &\rightarrow) \\ S &\rightarrow LR \\ S &\rightarrow SS \end{aligned}$$