

# Normal Forms for CFG's

Eliminating Useless Variables

Removing Epsilon

Removing Unit Productions

Chomsky Normal Form

# Variables That Derive Nothing

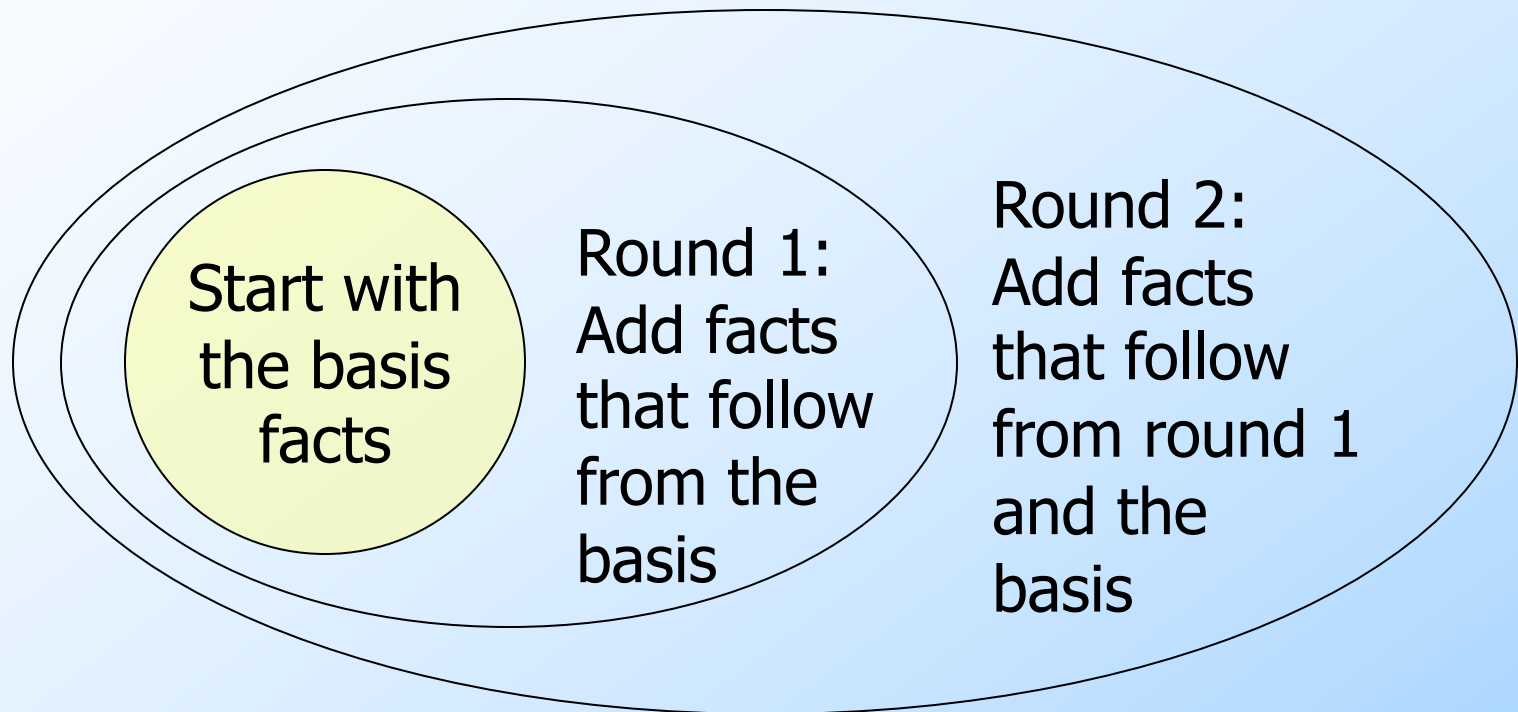
- ◆ Consider:  $S \rightarrow AB$ ,  $A \rightarrow aA \mid a$ ,  $B \rightarrow AB$
- ◆ Although A derives all strings of a's, B derives no terminal strings.
  - ◆ Why? The only production for B leaves a B in the sentential form.
- ◆ Thus, S derives nothing, and the language is empty.

# *Discovery* Algorithms

- ◆ There is a family of algorithms that work inductively.
- ◆ They start discovering some facts that are obvious (the basis).
- ◆ They discover more facts from what they already have discovered (induction).
- ◆ Eventually, nothing more can be discovered, and we are done.

# Picture of Discovery

And so on ...



# Testing Whether a Variable Derives Some Terminal String

- ◆ **Basis:** If there is a production  $A \rightarrow w$ , where  $w$  has no variables, then  $A$  derives a terminal string.
- ◆ **Induction:** If there is a production  $A \rightarrow \alpha$ , where  $\alpha$  consists only of terminals and variables known to derive a terminal string, then  $A$  derives a terminal string.

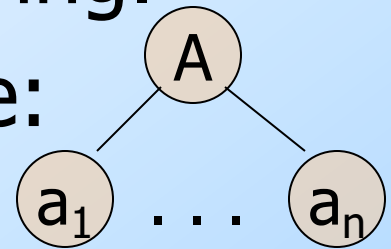
# Testing – (2)

- ◆ Eventually, we can find no more variables.
- ◆ An easy induction on the order in which variables are discovered shows that each one truly derives a terminal string.
- ◆ Conversely, any variable that derives a terminal string will be discovered by this algorithm.

# Proof of Converse

- ◆ The proof is an induction on the height of the least-height parse tree by which a variable  $A$  derives a terminal string.

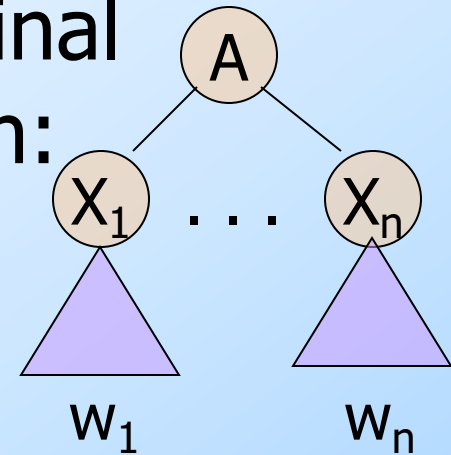
- ◆ **Basis:** Height = 1. Tree looks like:



- ◆ Then the basis of the algorithm tells us that  $A$  will be discovered.

# Induction for Converse

- ◆ Assume IH for parse trees of height  $< h$ , and suppose  $A$  derives a terminal string via a parse tree of height  $h$ :
- ◆ By IH, those  $X_i$ 's that are variables are discovered.
- ◆ Thus,  $A$  will also be discovered, because it has a right side of terminals and/or discovered variables.





# Algorithm to Eliminate Variables That Derive Nothing

1. Discover all variables that derive terminal strings.
2. For all other variables, remove all productions in which they appear in either the head or body.

# Example: Eliminate Variables

$S \rightarrow AB \mid C, A \rightarrow aA \mid a, B \rightarrow bB, C \rightarrow c$

- ◆ **Basis:** A and C are discovered because of  $A \rightarrow a$  and  $C \rightarrow c$ .
- ◆ **Induction:** S is discovered because of  $S \rightarrow C$ .
- ◆ Nothing else can be discovered.
- ◆ **Result:**  $S \rightarrow C, A \rightarrow aA \mid a, C \rightarrow c$

# Unreachable Symbols

- ◆ Another way a terminal or variable deserves to be eliminated is if it cannot appear in any derivation from the start symbol.
- ◆ **Basis**: We can reach  $S$  (the start symbol).
- ◆ **Induction**: if we can reach  $A$ , and there is a production  $A \rightarrow \alpha$ , then we can reach all symbols of  $\alpha$ .

# Unreachable Symbols – (2)

- ◆ Easy inductions in both directions show that when we can discover no more symbols, then we have all and only the symbols that appear in derivations from  $S$ .
- ◆ **Algorithm:** Remove from the grammar all symbols not discovered reachable from  $S$  and all productions that involve these symbols.

# Eliminating Useless Symbols

- ◆ A symbol is *useful* if it appears in some derivation of some terminal string from the start symbol.
- ◆ Otherwise, it is *useless*.  
Eliminate all useless symbols by:
  1. Eliminate symbols that derive no terminal string.
  2. Eliminate unreachable symbols.

## Example: Useless Symbols – (2)

$S \rightarrow AB, A \rightarrow C, C \rightarrow c, B \rightarrow bB$

- ◆ If we eliminated unreachable symbols first, we would find everything is reachable.
- ◆ A, C, and c would never get eliminated.

# Why It Works

- ◆ After step (1), every symbol remaining derives some terminal string.
- ◆ After step (2) the only symbols remaining are all derivable from  $S$ .
- ◆ In addition, they still derive a terminal string, because such a derivation can only involve symbols reachable from  $S$ .

# Epsilon Productions

- ◆ We can almost avoid using productions of the form  $A \rightarrow \epsilon$  (called  *$\epsilon$ -productions* ).
  - ◆ The problem is that  $\epsilon$  cannot be in the language of any grammar that has no  $\epsilon$ -productions.
- ◆ **Theorem:** If  $L$  is a CFL, then  $L - \{\epsilon\}$  has a CFG with no  $\epsilon$ -productions.



# Nullable Symbols

- ◆ To eliminate  $\epsilon$ -productions, we first need to discover the *nullable symbols* = variables  $A$  such that  $A \Rightarrow^* \epsilon$ .
- ◆ **Basis**: If there is a production  $A \rightarrow \epsilon$ , then  $A$  is nullable.
- ◆ **Induction**: If there is a production  $A \rightarrow \alpha$ , and all symbols of  $\alpha$  are nullable, then  $A$  is nullable.

# Example: Nullable Symbols

$S \rightarrow AB, A \rightarrow aA \mid \epsilon, B \rightarrow bB \mid A$

- ◆ **Basis:** A is nullable because of  $A \rightarrow \epsilon$ .
- ◆ **Induction:** B is nullable because of  $B \rightarrow A$ .
- ◆ Then, S is nullable because of  $S \rightarrow AB$ .

# Eliminating $\epsilon$ -Productions

- ◆ **Key idea:** turn each production  $A \rightarrow X_1 \dots X_n$  into a family of productions.
- ◆ For each subset of nullable  $X$ 's, there is one production with those eliminated from the right side "in advance."
  - ◆ Except, if all  $X$ 's are nullable (or the body was empty to begin with), do not make a production with  $\epsilon$  as the right side.

# Example: Eliminating $\epsilon$ -Productions

$S \rightarrow ABC, A \rightarrow aA \mid \epsilon, B \rightarrow bB \mid \epsilon, C \rightarrow \epsilon$

◆  $A, B, C,$  and  $S$  are all nullable.

◆ New grammar:

$S \rightarrow \cancel{ABC} \mid AB \mid \cancel{AC} \mid \cancel{BC} \mid A \mid B \mid \cancel{C}$

$A \rightarrow aA \mid a$

$B \rightarrow bB \mid b$

Note:  $C$  is now useless.  
Eliminate its productions.

# Why it Works

- ◆ **Prove** that for all variables  $A$ :
  1. If  $w \neq \epsilon$  and  $A \Rightarrow_{\text{old}}^* w$ , then  $A \Rightarrow_{\text{new}}^* w$ .
  2. If  $A \Rightarrow_{\text{new}}^* w$  then  $w \neq \epsilon$  and  $A \Rightarrow_{\text{old}}^* w$ .
- ◆ Then, letting  $A$  be the start symbol proves that  $L(\text{new}) = L(\text{old}) - \{\epsilon\}$ .
- ◆ (1) is an induction on the number of steps by which  $A$  derives  $w$  in the old grammar.

# Proof of 1 – Basis

- ◆ If the old derivation is one step, then  $A \rightarrow w$  must be a production.
- ◆ Since  $w \neq \epsilon$ , this production also appears in the new grammar.
- ◆ Thus,  $A \Rightarrow_{\text{new}} W$ .

# Proof of 1 – Induction

- ◆ Let  $A \Rightarrow_{\text{old}}^* w$  be a  $k$ -step derivation, and assume the IH for derivations of fewer than  $k$  steps.
- ◆ Let the first step be  $A \Rightarrow_{\text{old}} X_1 \dots X_n$ .
- ◆ Then  $w$  can be broken into  $w = w_1 \dots w_n$ , where  $X_i \Rightarrow_{\text{old}}^* w_i$ , for all  $i$ , in fewer than  $k$  steps.

# Induction – Continued

- ◆ By the IH, if  $w_i \neq \epsilon$ , then  $X_i \Rightarrow_{\text{new}}^* w_i$ .
- ◆ Also, the new grammar has a production with  $A$  on the left, and just those  $X_i$ 's on the right such that  $w_i \neq \epsilon$ .
  - **Note:** they all can't be  $\epsilon$ , because  $w \neq \epsilon$ .
- ◆ Follow a use of this production by the derivations  $X_i \Rightarrow_{\text{new}}^* w_i$  to show that  $A$  derives  $w$  in the new grammar.



# Unit Productions

- ◆ A *unit production* is one whose body consists of exactly one variable.
- ◆ These productions can be eliminated.
- ◆ **Key idea:** If  $A \Rightarrow^* B$  by a series of unit productions, and  $B \rightarrow \alpha$  is a non-unit-production, then add production  $A \rightarrow \alpha$ .
- ◆ Then, drop all unit productions.

# Unit Productions – (2)

- ◆ Find all pairs  $(A, B)$  such that  $A \Rightarrow^* B$  by a sequence of unit productions only.
- ◆ **Basis**: Surely  $(A, A)$ .
- ◆ **Induction**: If we have found  $(A, B)$ , and  $B \rightarrow C$  is a unit production, then add  $(A, C)$ .

# Proof That We Find Exactly the Right Pairs

- ◆ By induction on the order in which pairs  $(A, B)$  are found, we can show  $A \Rightarrow^* B$  by unit productions.
- ◆ Conversely, by induction on the number of steps in the derivation by unit productions of  $A \Rightarrow^* B$ , we can show that the pair  $(A, B)$  is discovered.

# Proof The the Unit-Production-Elimination Algorithm Works

- ◆ **Basic idea:** there is a leftmost derivation  $A \Rightarrow_{lm}^* w$  in the new grammar if and only if there is such a derivation in the old.
- ◆ A sequence of unit productions and a non-unit production is collapsed into a single production of the new grammar.

# Cleaning Up a Grammar

- ◆ **Theorem:** if  $L$  is a CFL, then there is a CFG for  $L - \{\epsilon\}$  that has:
  1. No useless symbols.
  2. No  $\epsilon$ -productions.
  3. No unit productions.
- ◆ I.e., every body is either a single terminal or has length  $\geq 2$ .

# Cleaning Up – (2)

- ◆ **Proof:** Start with a CFG for L.
- ◆ Perform the following steps in order:

1. Eliminate  $\epsilon$ -productions.
2. Eliminate unit productions.
3. Eliminate variables that derive no terminal string.
4. Eliminate variables not reached from the start symbol.

Must be first. Can create unit productions or useless variables.

# Chomsky Normal Form

- ◆ A CFG is said to be in *Chomsky Normal Form* if every production is of one of these two forms:
  1.  $A \rightarrow BC$  (body is two variables).
  2.  $A \rightarrow a$  (body is a single terminal).
- ◆ **Theorem:** If  $L$  is a CFL, then  $L - \{\epsilon\}$  has a CFG in CNF.

# Proof of CNF Theorem

- ◆ **Step 1:** “Clean” the grammar, so every body is either a single terminal or of length at least 2.
- ◆ **Step 2:** For each body  $\neq$  a single terminal, make the right side all variables.
  - ◆ For each terminal  $a$  create new variable  $A_a$  and production  $A_a \rightarrow a$ .
  - ◆ Replace  $a$  by  $A_a$  in bodies of length  $\geq 2$ .



## Example: Step 2

- ◆ Consider production  $A \rightarrow BcDe$ .
- ◆ We need variables  $A_c$  and  $A_e$ . with productions  $A_c \rightarrow c$  and  $A_e \rightarrow e$ .
  - **Note:** you create at most one variable for each terminal, and use it everywhere it is needed.
- ◆ Replace  $A \rightarrow BcDe$  by  $A \rightarrow BA_cDA_e$ .

# CNF Proof – Continued

- ◆ **Step 3:** Break right sides longer than 2 into a chain of productions with right sides of two variables.
- ◆ **Example:**  $A \rightarrow BCDE$  is replaced by  $A \rightarrow BF$ ,  $F \rightarrow CG$ , and  $G \rightarrow DE$ .
  - ◆ F and G must be used nowhere else.

## Example of Step 3 – Continued

- ◆ Recall  $A \rightarrow BCDE$  is replaced by  $A \rightarrow BF$ ,  $F \rightarrow CG$ , and  $G \rightarrow DE$ .
- ◆ In the new grammar,  $A \Rightarrow BF \Rightarrow BCG \Rightarrow BCDE$ .
- ◆ **More importantly**: Once we choose to replace  $A$  by  $BF$ , we must continue to  $BCG$  and  $BCDE$ .
  - ◆ Because  $F$  and  $G$  have only one production.