

Lecture Tue Mar 31  $\text{Diff}(x, i) = \#(' [x_1 \dots x_i] - \#') [x_1 \dots x_i] \quad 0 \leq i \leq n$   
 $n = |x|$

$\text{Bal} = \{x \in \{L, \}^* : \text{Diff}(x, n) = 0 \wedge (\forall i: 0 \leq i < n) \text{Diff}(x, i) \geq 0\}$ .

$G_1: S_1 \rightarrow \epsilon \mid (S_1)S_1$      Prove  $L(G_1) = L(G_2) = \text{Bal}$ .

$G_2: S_2 \rightarrow \epsilon \mid (S_2)S_2$ .     (Both grammars have  $S \rightarrow \epsilon$ )

Note:  $G_1$  simulates  $G_2$  rule-wise via  $S \Rightarrow SS \Rightarrow (S)S$ .

Hence: Soundness of  $G_1 \Rightarrow$  soundness of  $G_2$ . I.e.  $L(G_1) \subseteq \text{Bal} \Rightarrow L(G_2) \subseteq \text{Bal}$ .

But,  $G_2$  does not simulate  $G_1$  rule-wise: in  $G_2$  you cannot get the so-called "sentential form"  $SS$ .  
 any derivable string in  $(V \cup \Sigma)^*$ . (So the "expanded" languages of the two grammars are not the same, but we don't care)

Hence also: If  $G_2$  is comprehensive, then  $G_1$  is comprehensive. I.e.  
 $L(G_2) = \text{Bal} \Rightarrow L(G_1) = \text{Bal}$ . Earlier we proved  $G_1$  is sound, so we only need  $L(G_2) = \text{Bal}$ .

Proof by induction on the lengths  $n$  of strings:  $x \in \text{Bal} \Rightarrow S_2 \Rightarrow^* x$ .  
 Prove  $\forall n P(n)$ , where  $P(n) \equiv$  for each  $x \in \{L, \}^n$ ,  $S_2 \Rightarrow^* x$ .  
 $G: S \rightarrow \epsilon \mid (S)S$      Basis ( $n=0$ )  $\epsilon \in \text{Bal}$ .  $S \Rightarrow \epsilon$  so that checks.

Another Basis ( $n=1$ ): The only strings  $x \in \{L, \}^1$  are  $x = ' and  $x = ''$ .  
 Neither string belongs to  $\text{Bal}$ . So this holds by default - nothing to do.$

Induction ( $n \geq 2$ ): Assume (IH) the statement  $(\forall m < n) P(m)$ .

Given  $n$ , goal: show  $P(n)$ . Let <sup>any</sup>  $x \in \{L, \}^n$  be given.

• If  $x \in \text{Bal}$  (as always happens when  $n$  is odd), there is nothing to do.

• Let  $x \in \text{Bal}$ . Our goal is to show  $S \Rightarrow^* x$  in  $G_2$ . Consider how to parse  $x$ .

We know  $\text{Diff}(x, i) = 0$  when  $i=0$  and when  $i=n=|x|$ , and  $\text{Diff}(x, i) \geq 0$  for all  $i \leq n$ .

We want to isolate a "critical" value of  $i$ .  
 $x = (L) (L) \dots$  Parse as  $S \Rightarrow (S) \epsilon$  ( $i=n \geq 6$ )  
 $x' = (L(L)) (L) \dots$   $i \geq 6$  won't work for  $x'$ : The parse  $S \Rightarrow LS$  leaves  $(L) \epsilon \notin \text{Bal}$ .

Mathematical Khetoni to locate the matching mate of the opening '(': ②

"Take  $i$  to be the least  $i > 0$  such that  $\text{Diff}(X, i) = 0$ ."

By  $\text{Diff}(X, n) = 0$  we know that the loop `while (Diff(X, i) != 0) { i++; }` must exit. [Side remark: when we don't have such a bound, this kind of "Impredicative" math definition is highly controversial.]

Clearer rewrites: • Take  $i$  to be least such that  $\text{Diff}(X, i) = 0$  ( $i > 0$ ).

• Take  $i_0$  to be the least  $i > 0$  such that  $\text{Diff}(X, i) = 0$ .

$$X = \left( \begin{array}{c} \dots \\ \text{Diff}(i) > 0 \quad \text{Diff}(i) = 0 \quad \text{Diff}(i) = 0 \text{ at } i = n \end{array} \right)$$

Then we can parse  $X = (Y)Z$ , where  $|Y| = i - 2$  and:

•  $\text{Diff}(Y, |Y|) = \text{Diff}(X, i-1) - \underline{1} = 1 - 1 = 0$

• For each  $j$ ,  $0 \leq j \leq i-2$ ,  $\text{Diff}(Y, j) = \text{Diff}(X, j+1) - \underline{1}$ .   
taking away the opening '('      ← again, taking away the opening '('

Because  $i$  is least such that  $\text{Diff}(X, i) = 0$  (and  $i > 0$ ), we have that for all  $j$ ,  $j+1$  runs from 1 to  $i-1$  as indices of  $X$ , so that

$$\boxed{\text{Diff}(X, j+1) \geq 1} \quad \therefore \forall j, 0 \leq j \leq |Y| : \text{Diff}(Y, j) \geq 0.$$

$\therefore \text{Diff}(Y, |Y|) = 0$  and  $\forall j, 0 \leq j \leq |Y|, \text{Diff}(Y, j) \geq 0 \therefore \boxed{Y \in \text{Bal}}$

Moreover, for all  $k$ ,  $i \leq k \leq n$   $\text{Diff}(X, k) = \text{Diff}(Z, k-i)$ .

$\therefore \text{Diff}(Z, |Z|) = \text{Diff}(X, n) = 0$ , and for all  $l$ ,  $0 \leq l \leq |Z|, \text{Diff}(Z, l) \geq 0$ .

$\therefore Y \in \text{Bal}$  and  $Z \in \text{Bal}$ . And  $|Y| = i-2 < n$  (even when  $i=n$ ),  $|Z| = n-i < n$

Hence we can apply IH  $P(i-2)$  and  $P(n-i)$  to get

$S \Rightarrow Y$  and  $S \Rightarrow Z$  in  $G_2$ . Thus  $S \Rightarrow (S)S \Rightarrow (Y)S \Rightarrow (Y)Z = X$ .   
since  $i > 0$ !       $\therefore P(n)$       done!

Q: Do we need  $S \rightarrow \epsilon$ ?  $G_2: S \rightarrow \epsilon \mid (S)S$ .  $L(G_2) = \text{Bal}$ . ③

Yes since  $\epsilon \in \text{Bal}$ , but can we get all of  $\text{Bal} - \{\epsilon\}$  without  $\epsilon$ -rules?

How about  $G' = S \rightarrow () \mid (S)S$ .  $G_1$  is sound, but is it comprehensive for  $\text{Bal} - \{\epsilon\}$ ?

Try  $x = (( ))$ .  $S \Rightarrow (S)S \Rightarrow (( ))S$  alas leaves an  $S$  we can no longer get rid of.

Theorem [first part of conversion to Chomsky Normal Form]

For any CFG  $G = (V, \Sigma, R, S)$ , we can build a CFG  $G'$  without  $\epsilon$ -rules such that  $L(G') = L(G) - \{\epsilon\}$ . [optionally, as the text does, if  $\epsilon \in L(G)$  we can finally add a new start symbol  $S'$  with rules  $S' \rightarrow \epsilon \mid S$ ]

Proof: ① Identify the subset  $W \subseteq V$  of nullable variables, meaning  $A \in V$  such that  $A \Rightarrow^* \epsilon$ . In  $G_2: W = V = \{S\}$ .

② For every rule  $B \rightarrow X$  in  $G$ , add to  $R$  all rules obtained by deleting 1 or more nullable variables. [As literally described, this can be an exponential explosion]

$S \rightarrow \epsilon \mid (S)S$ . Add: occurrences of.

$S \rightarrow ( )S$   
 $S \rightarrow (S)$  and  
 $S \rightarrow ( )$  deleting both occurrences.

Now we can do  $S \Rightarrow (S) \Rightarrow ( ( ) )$ .

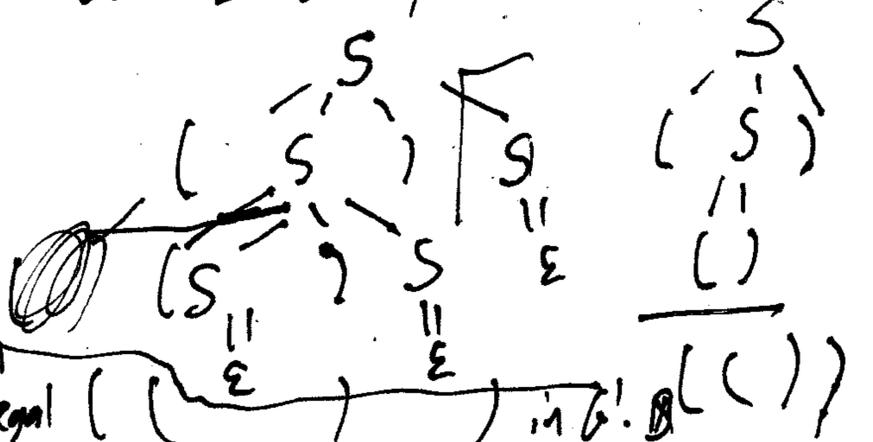
This is sound because all added rules could be simulated by deriving the deleted occurrences to  $\epsilon$ .

Finally, define  $G'$  by deleting all  $\epsilon$ -rules.

$L(G') \subseteq L(G) - \{\epsilon\}$ . Need  $L(G') \supseteq L(G) - \{\epsilon\}$

Here we show that  $L(G') = L(G) - \{\epsilon\}$ , we have to show that all nonempty strings  $x \in L(G)$  can be derived in  $G'$ .

Consider any parse tree for  $x$  in the original grammar  $G$ . Just delete all subtrees that yielded  $\epsilon$ .  $G'$  covers all such deletions with its rules, so the new parse is legal.



Defn: A CFG  $G$  is in Chomsky Normal Form (ChNF) <sup>(4)</sup> if all rules in  $R$  have the form  $A \rightarrow c$  or  $A \rightarrow BC$  with  $A, B, C \in V$ ,  $c \in \Sigma$ .  $\underbrace{c}_{\in \Sigma}$ . possibly  $B, C = A$

Theorem: For every CFG  $G$  we can build a CFG  $G''$  in ChNF such that  $L(G'') = L(G) \setminus \{\epsilon\}$ , [optionally add  $S'' \rightarrow \epsilon$  | right-hand sides of  $S$  to allow  $L(G'') = L(G)$ .]

Note: ChNF involves the step of aliasing every terminal  $c \in \Sigma$  to a variable  $X_c$  and adding the rules  $X_c \rightarrow c$  as the only productions with terminals on the RHS.  $S \rightarrow X_L S X_R S \mid X_L X_R S \mid X_L S X_R \mid X_L X_R$   
 $X_R \rightarrow ' ) ' \quad X_L \rightarrow '('$

Extra The rest of the algorithm:

- ① Eliminate  $\epsilon$ -rules as above (don't do the optional  $S' \rightarrow S \mid \epsilon$  yet: save it for the end).
- ② Now find all pairs of variables  $(A, B)$  such that  $A \Rightarrow^* B$ . This might happen by transitivity:  $A \rightarrow C, C \rightarrow D, D \rightarrow B$ . Add all right-hand sides of  $B$  as rule options for  $A$ . Finally delete all "unit productions"  $A \rightarrow C$  etc. Call the result  $G'$ .
- ③ Finally, the "silly steps": Alias every  $c \in \Sigma$  to  $X_c$  as above, except in rules  $A \rightarrow b$ .
- ④ For every right-hand side of length  $\geq 3$ , make variables  $Y_1, \dots, Y_{l-1}$ . If the rule is  $A \rightarrow V_1 V_2 \dots V_l$ , the new rules are  $A \rightarrow V_1 Y_1, Y_1 \rightarrow V_2 Y_2, Y_2 \rightarrow V_3 Y_3, \dots$ , up to  $Y_{l-1} \rightarrow V_{l-1} V_l$ . Use different "Y"-variables for each rule. Call that  $G''$  too. Optionally, add  $S'' \rightarrow \epsilon$  | right-hand sides of  $S$  in  $G''$ . Then  $L(G'') = L(G)$  in ChNF.

Example of step ④ for our parentheses grammar:

$S \rightarrow X_L Y_{11} \mid X_L Y_{21} \mid X_L Y_{31} \mid X_L X_R, X_R \rightarrow ' ) ', X_L \rightarrow '('$   
 $Y_{11} \rightarrow S Y_{12}, Y_{12} \rightarrow X_R S, Y_{21} \rightarrow X_R S, Y_{31} \rightarrow S X_R$ . Option:  $S'' \rightarrow \epsilon \mid X_L Y_{11} \mid X_L Y_{21} \mid X_L Y_{31} \mid X_L X_R$

Then  $G''$  is in Chomsky NF.

(Can combine  $Y_{12}$  and  $Y_{21}$ , but this is hideously unreadable anyway.)