Recall

$G_2:$   $S \to \varepsilon \mid AB \mid BA \mid SS$    $L = \{x : \#a(x) = \#b(x)\}$

   $A \to a \mid aS \mid BAA$    $L_A = \{x : \#a(x) = \#b(x) + 1\}$

   $B \to b \mid bS \mid ABB$    $L_B = \{x : \#a(x) = \#b(x) - 1\}$

For $\underline{L(G_2) \subseteq L}$ (soundness) define $P_S, P_A, P_B$, where for instance

$P_A = $ "Every string $x$ that I derive has $\#a(x) = \#b(x) + 1$."

For $\underline{L \subseteq L(G)}$ (comprehensiveness) we do the opposite: $\quad P(n) \equiv P(n)$ $\land Q(n) \land R(n)$

where $Q(n) \equiv$ For each $x \in \Sigma^n$, if $x \in L_A$ then $A \Rightarrow^* x$, which kind of

    says: "I derive every string $x$ that has $\#a(x) = \#b(x) + 1$."

$\underset{\text{Comprehensiveness}}{\text{Proving}} \equiv \underset{\text{parsing}}{\text{Proving an algorithm for}}$   E.g.   Given $x = babaaab$,

                                   how does $x$ get derived?

$S \to \varepsilon \mid AB \mid BA \mid SS$ ← unnecessary   Is $x \in L_S, L_A, L_B$, or none of the above?

$A \to a \mid AS \mid BAA$

$B \to b \mid BS \mid ABB$     $\underline{x \in L_A}$. Proof does subcases for the first letter in $x$.

                   $\underline{x = ay:} \Rightarrow y \in L_S$ so do $A \to \boxed{AS \to aS} \Rightarrow^* ay = x.$

Proof notes that the Diff   $\underline{x = bz:}$    Hence we can parse $x =: buv$

function $\text{diff}(x, i) = \#a(x_1 .. x_i) - \#b(x_1 .. x_i)$   where $|bu| = j$, i.e. $\text{diff}(bu, j) = 0$

• starts at $\boxed{0}$   $\boxed{\text{steps} \pm 1 \text{ at each char}}$

| b | u = a or aba | v = ba²b or aab |
|---|---|---|

• Goes to $\boxed{-1}$ since $x$ begins with $b$

• Ends at $\boxed{+1}$ since $\#a(x) = \#b(x) + 1$   Diff   0   -1            0        +1

∴ By a "Discrete" version of the Intermediate   $x = b \cdot aba \cdot aab$   $A \Rightarrow BAA \Rightarrow bAA$

Value Theorem, $\exists i \; \text{diff}(x_1 .. x_i) = 0.$       $\Rightarrow bAS \Rightarrow ba$

$$A \Rightarrow BAA \Rightarrow \underbrace{bAA}_{aab} \Rightarrow bASA \Rightarrow baSA \Rightarrow buBAA \Rightarrow babAA$$

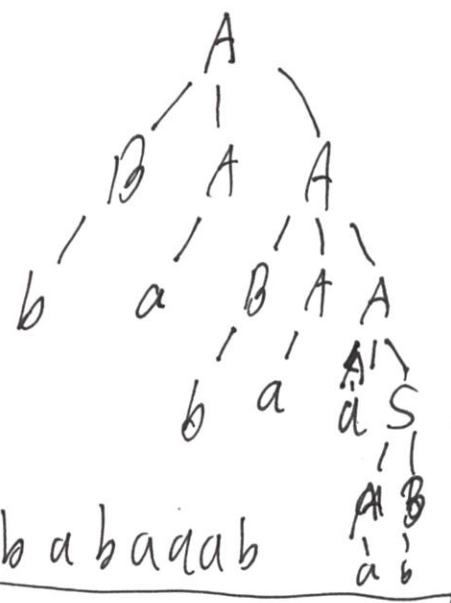$$\overset{aba}{\Rightarrow} \underline{babaA} \Rightarrow babaaS \Rightarrow babaaAB \Rightarrow^2 babaaab$$



Using the "least $j$" rule — Parses breaks babaaab as

$$b \cdot a \cdot babaab$$
$$0 \quad -1 \quad 0 \qquad +1$$

<u>Refined grammar</u>
$$S \to \varepsilon \mid aB \mid bA$$
$$A \to aS \mid \underline{bAA}$$

By symmetry, can use $B \to bS \mid \underline{aBB}$
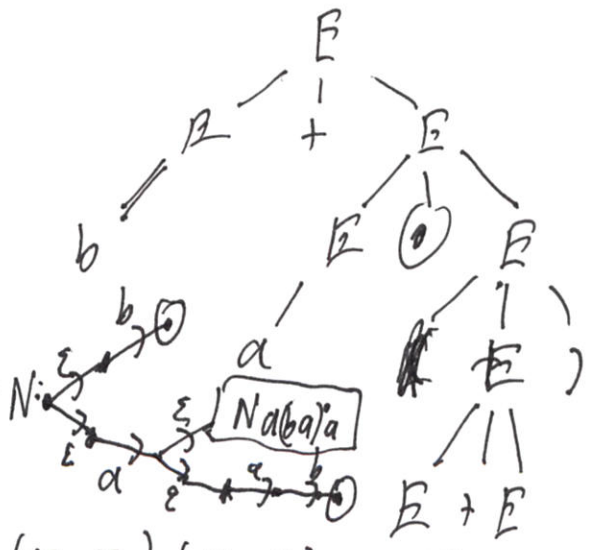
This grammar $G_2'$ remains comprehensive:
$L \subseteq L(G_2')$ by what the proof <u>actually used</u>.
Unfortunately $G_2'$ is still ambiguous, as the case $X = babaaab$ showed.



babaaab

Regular Expression Parsing:

$$b \cup a \big( a(ba)^* a \cup ab \big)$$
$$\underline{b + a \cdot (a \cdot (b \cdot a)^* \cdot a + a \cdot b)}$$



Fully parenthesized: $E \to \emptyset \mid (\varepsilon \mid a \mid b) \mid (E + E) \mid (E \cdot E) \mid (E^*)$

Thus the Regexp-to-NFA proof became one by Structural Induction.
was

## Chomsky Normal Form → Makes the proof of the CFL Pumping Lemma in §2.3 nicer.

__Def^n__: A CFG $G = (V, \Sigma, R, S)$ is in Chomsky Normal

Form (ChNF) if every rule $A \to X$ has $\underline{X \in \Sigma}$ or

Allowed: $A \to \overset{c \in \Sigma}{c}$, $A \to AA$, $A \to BC$          $|X| = 2$

Not Allowed: $\underset{\varepsilon\text{-rule}}{\underline{A \to \varepsilon}}$  $\underset{\text{mixed RHS}}{\underline{A \to aS}} | \underset{\text{unit rule}}{\underline{A \to B}}$  $\underset{|RHS| \geq 3}{\underline{A \to BCD}}$   with $\underline{X \in VV}$.

(Text allows $\underline{S \to \varepsilon}$ provided $S$ is not on the RHS of any rule.)
Can always put $\varepsilon$ in last by a new start var. $S_0$ and rules $S_0 \to \varepsilon | \binom{RHS}{\text{of } S}$

__Theorem__: Given any CFG $G$, we can build $G'$ in ChNF st. $L(G') = L(G)$.

__Proof. By__
__Algorithm__  ① Identify the set NULLABLE of variables $A$ st. $A \overset{*}{\Rightarrow} \varepsilon$.

(not the fastest possible)  ② For every rule $A \to X$ include every rule $A \to X'$

or blends ① & ②   obtained by erasing 1 or more __occurrences__ of vars in NULLABLE

keep
$\Rightarrow BCDB$)   Idea: Suppose $A \to X$ is $A \to B\underline{CDB}$ and $B, C$ are nullable.
$\Rightarrow BD$,   Derivations can do $B \overset{*}{\Rightarrow} \varepsilon$ and $C \overset{*}{\Rightarrow} \varepsilon$, hence they can do all of:
$\quad \underline{A \overset{*}{\Rightarrow} D}$, $A \overset{*}{\Rightarrow} DB$, $A \overset{*}{\Rightarrow} BCD$, $A \overset{*}{\Rightarrow} CD$, $A \overset{*}{\Rightarrow} CDB$, $A \overset{*}{\Rightarrow} BDB$

involves a   ③ Then delete all $\varepsilon$-rules. Got $G_1$ st. $L(G_1) = L(G) \setminus \{\varepsilon\}$.
transitive
closure of   ④ Determine which $A, B \in V$ allow $A \underset{G_1}{\overset{*}{\Rightarrow}} B$  (note: we may get more unit rules from step 2)
$\to B$ relation. ⑤ For each such $A, B$ $(B \neq A)$, make every RHS of $B$ a RHS of $A$.
⑥ Then __delete__ all unit rules to get $G_2$ st. $L(G_2) = L(G) \setminus \{\varepsilon\}$.
Alias every terminal $a$ to a variable $X_a$, and ⑧ Use dummy variables to reduce RHS to 2.

Example: $A \rightarrow BCdB$, then add variables $X_d, Y_1, Y_2$ and do:

$A \rightarrow BY_1, \quad Y_1 \rightarrow C.Y_2, \quad Y_2 \rightarrow X_dB, \quad X_d \rightarrow d.$

## Added — for Thursday — Full Examples:

$$S \rightarrow \varepsilon \mid (S)S \qquad L = \{x \in \{(,)\}^* : X \text{ is } \underline{balanced}\}$$

1. NULLABLE $= \{S\}$ since $S \Rightarrow^* \varepsilon$

2. Add rules: $S \rightarrow () \mid (S) \mid ()S$

3. Delete $S \rightarrow \varepsilon$: $G_1 = S \rightarrow (S)S \mid () \mid (S) \mid ()S$ $\qquad L(G_1) = L \setminus \{\varepsilon\}$

4. No unit rules were introduced or were there previously. 5.6. can skip.

7. Alias $L \rightarrow ($ , $R \rightarrow )$. $\quad G_2 = S \rightarrow LSRS \mid LR \mid LSR \mid LRS \mid \begin{array}{l} L \rightarrow ( \\ R \rightarrow ) \end{array}$

8. Add variables $Y_1, Y_2, Y_3, Y_4$.

$$\boxed{\begin{array}{l} G_3 = S \rightarrow LY_1 \mid LR \mid LY_3 \mid LY_4, \quad L \rightarrow (, R \rightarrow ) \\ \quad Y_1 \rightarrow SY_2, \quad Y_2 \rightarrow RS, \quad Y_3 \rightarrow SR, \quad Y_4 \rightarrow RS \end{array}}$$

Yes we could economize with $Y_4 \equiv Y_2$, but who cares? $G_3$ is UGLY!!!

9. Now $L(G_3) = L \setminus \{\varepsilon\}$. If we want to put $\varepsilon$ back in the language, do:

$$G' = S_0 \rightarrow \varepsilon \mid LY_1 \mid LR \mid LY_3 \mid LY_4 \qquad Y_1 \rightarrow SY_2, \quad Y_2 \rightarrow RS, \quad Y_3 \rightarrow SR, \quad Y_4 \rightarrow RS$$
$$S \rightarrow LY_1 \mid LR \mid LY_3 \mid LY_4 \qquad L \rightarrow '(', \quad R \rightarrow ')! \quad \boxed{\text{OK in text ChNF}}$$

We couldn't add $S_0 \rightarrow S$ since that would be a unit rule. Then $\underline{L(G') = L(G) = L}$

---

The algorithm to define NULLABLE is an important kind of iterative loop akin to Breadth-First Search. Note that if we initialize NULLABLE $= \emptyset$ instead, then since $\emptyset^* = \varepsilon$, we shall get the $\varepsilon$-rule variables on the first iteration.

```
NULLABLE := {A ∈ V: A → ε is a rule};
bool foundNew = true;
while (foundNew) {
    foundNew = false;
    foreach (rule A → X in R) {
        if (X ∈ NULLABLE*) {
            NULLABLE = NULLABLE ∪ {A};
            foundNew = true;
        } }
```