TopHat 3557    $G: S \to SS \mid bA \mid aB \mid \varepsilon$    Target $E =$
$\qquad\qquad\quad A \to aS \mid bAA$    $\{ x \in \{a, b\}^* :$
$\qquad\qquad\quad B \to bS \mid aBB.$    $\#a(x) = \#b(x) \}.$

Note $S \to \varepsilon$ is the basis for all three variables, but we will not need any special treatment.
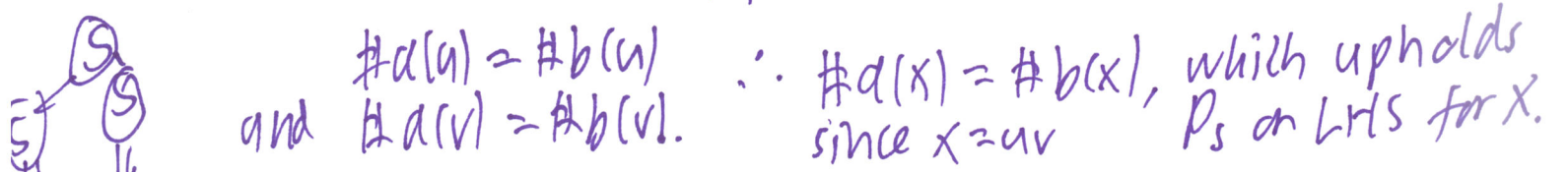
Proof of $L(G) \subseteq E$ by "Structural Induction."
Define the properties

$P_S \equiv$ "Every $x$ {that I derive} such that $S \Rightarrow^* x$ has $\#a(x) = \#b(x)$, belongs to $E$."

$P_A \equiv$ "Every $y$ s.t. $A \Rightarrow^* y$ has $\#a(y) = \#b(y) + 1$.
one more $a$ than $b$.

$P_B \equiv$ "Every $z$ s.t. $B \Rightarrow^* z$ has $\#b(z) = \#a(z) + 1$.
one more $b$ than $a$.

Go through all rules for each variable and show that if the variables on the RHS derive strings that obey their properties, then the resulting string obeys ("upholds") the property of the variable on the LHS.

$\underline{S \to SS}$:  Suppose $S \Rightarrow^* x$ using this rule first. Then $x ::= uv$ where $S \Rightarrow^* u$ and $S \Rightarrow^* v$. By IH $P_S$ on RHS (twice), we have



$\qquad\qquad \#a(u) = \#b(u)$    $\therefore \#a(x) = \#b(x)$, which upholds
$\qquad\qquad$ and $\#a(v) = \#b(v)$.    since $x = uv$    $P_S$ on LHS for $x$.

$\underline{u \quad v = x}$   $\underline{S \to bA}$: Suppose $S \Rightarrow^* utrf$. Then $x ::= by$ where $A \Rightarrow^* y$

By IH $P_A$ on RHS, $y$ has 1 more $a$ than $b$. The leading 'b' in $x$ thus equalizes the count: $\#a(x) = \#b(x)$, and this upholds $P_S$ on LHS.

$\qquad \underline{S \to aB}$: "Similar to $S \to bA$ case." This finishes $S$. Are we done?

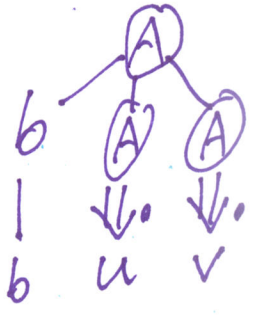$\qquad\underline{S \to \varepsilon}$: Suppose $S \Rightarrow^* x$ utrf. Then $x = \varepsilon$. And $\#a(\varepsilon) = 0 = \#b(\varepsilon)$ and this upholds $P_S$ on LHS.   No . . .

$S \to SS \mid aB \mid bA \mid \varepsilon$
$A \to aS \mid bAA$
$B \to bS \mid aBB$

$\underline{A \to aS}$: Suppose $A \Rightarrow^* y$ utrf. Then
$y =: \underline{a}x$ where $S \Rightarrow^* x$. By IH $P_S$ on RHS
$\#a(x) = \#b(x)$. $\therefore y$ has one more a than b
$\therefore P_A$ on LHS is upheld.

$\underline{A \to bAA}$: Suppose $A \Rightarrow^* x$ utrf. Then $x =: buv$
where $u$ and $v$ each have one more 'a' than 'b'.
The leading b thus brings the count in $x$ down from two
more a's in the uv part to one more 'a' overall. $\therefore P_A$ on LHS.

$\underline{B \to bS}$, $\underline{B \to aBB}$: Similar to the last two cases." $\therefore L(G) \subseteq E$ by SR.



How about $E \subseteq L(G)$? Since $G$ is sound, if $G$ is comprehensive
then we get $L(G) = E$.

$\underline{Sketch}$: Idea is induction on the length of strings, also using the
languages $E_1 = \{x : \#a(x) = \#b(x) + 1\}$ and $E_2 = \{x : \#a(x) = \#b(x) - 1\}$

Think of this as a task of parsing when you're in "S mode", A mode, or B mode.
These "modes" stand for routines that build parse trees from the given variable.

Example:: $x = abbbabaa$

$S \to aB \mid bA \mid \varepsilon$   $B \to bS$
$A \to aS \mid bAA$   $B \to aBB$

$S \Rightarrow aB$.   $x = ay$ with $y = bbbabaa$.


$S \Rightarrow^* x$.

Key point: $n' = |y|$ is $< n = |x|$. So we can use an IH
for $E_2$ saying that for all $y$ of length $\leq n$, $\underline{if}$
$\#b(y) = \#a(y) + 1$   $\underline{then}$ $B \Rightarrow^* y$. "B is comprehensive for $E_2$"
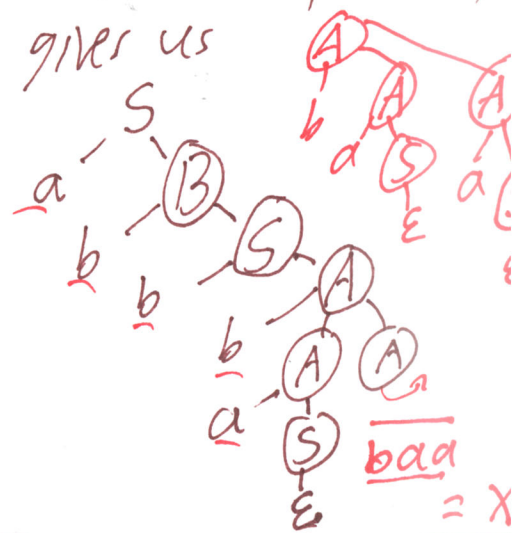
We also need to verify othe variables recursively. So with $y = \underline{bbbabaa}$,
we have $y \in E_2$ so in 'B mode' we need to parse it from B.
We can use $B \to bS$ and go back to S mode on $x' = bbabaa$.
Next we do $S \Rightarrow bA$ and we're in A-mode on $y' = babaa$. What now
This is OK because we can parse $y'$ as $b \cdot \underline{a} \cdot \underline{baa}$. The parts u=a,
v=baa are both in E

Hence by induction hypothesis on comprehensiveness of A for $F_1$, [3]
we can derive both of them from A. This gives us whole parse so far is:

$$A \Rightarrow bAA$$
$$\Downarrow_\bullet \quad \Downarrow_\bullet$$
$$\overline{\phantom{xx}u \quad v \phantom{xx}}$$
$$b \, u \, v = \gamma'$$



The full proof that the other rules $B \Rightarrow a BB$ and $A \Rightarrow aS$ and $S \Rightarrow aB$ make all three variables comprehensive is similar and never needs the rule $S \rightarrow Ss$. So we can delete it from G.

$$G' = \begin{array}{l} S \rightarrow aB \mid bA \mid \varepsilon \\ A \rightarrow aS \mid bAA \\ B \rightarrow bS \mid aBB \end{array}$$

Deleting the rule left G' still sound and we showed G' is still comprehensive.

Because $\gamma'$ could have been broken as $b \cdot aba \cdot a$ G and G' are ambiguous!

## Chomsky normal form

<u>Def$^n$</u>: A grammar G is in ChNF if all of its rules have the form $A \rightarrow c \ (c \in \Sigma)$ or $A \rightarrow BC$ with $B, C \in V$, (either can be A itself or C = B allowed)

The original def$^n$ (as in the text) forbids B or C = S.
It also disallows $S \rightarrow \varepsilon$. But both can be tolerated by re-defining a new start symbol $S_0$ with rules $S_0 \rightarrow \varepsilon \mid$ any r.h.s. of S.
Hence in lecture, ChNF will allow S on right-hand sides.

[ADDED]: I have decided this time not to cover the proof that every CFG G can be converted to G' in ChNF (strict ChNF if $\varepsilon \notin L(G)$), except that the step of bypassing $\varepsilon$-rules will be covered where relevant in Ch. 4. The <u>fact</u> of the conversion will still be used to make the CFL Pumping Lemma (§2.3) easier to visualize on T...