

# COMP2022: Formal Languages and Logic

2017, Semester 1, Week 6

Joseph Godbehere

Adapted from slides by A/Prof Kalina Yacef

April 11, 2017



THE UNIVERSITY OF  
**SYDNEY**

# COMMONWEALTH OF AUSTRALIA

## Copyright Regulations 1969

### WARNING

This material has been reproduced and communicated to you by or on behalf of the University of Sydney pursuant to part VB of the Copyright Act 1968 (**the Act**).

The material in this communication may be subject to copyright under the Act. Any further copying or communication of this material by you may be subject of copyright protect under the Act.

**Do not remove this notice.**

# ANNOUNCEMENTS

Friday 14th April is a public holiday

- ▶ No scheduled office hour / advanced seminar

Tuesday 25th April is a public holiday

- ▶ Replacement lecture:
  - ▶ 4pm – 6pm on Wednesday 26th April
  - ▶ PNR LT1 (i.e. the LT *next* to this one)
- ▶ Tuesday tutorials:
  - ▶ Join Wednesday / Thursday tutorials *if there is space*
    - ▶ See Ed for times/locations
  - ▶ Otherwise work through the exercises independently
- ▶ No quiz week 7
- ▶ Quiz week 8 will assess content from weeks 6 and/or 7

# OUTLINE

- ▶ LL(k) Table-Descent Parsing (continued)
- ▶ Grammar transformations
  - ▶ Finding equivalent LL(1) grammars
  - ▶ Chomsky Normal Form
  - ▶ Greibach Normal Form

# FIRST AND FOLLOW SETS

In order to fill in the entries of the table-driven parser, we need to compute some FIRST and FOLLOW sets.

$FIRST(\alpha)$  is the set of all terminals which could start strings derived from  $\alpha$ . We will need to calculate these for every production of  $G$  (i.e. the *right hand side* of each rule).

$FOLLOW(V)$  is the set of all terminals which which could follow the variable  $V$  at any stage of the derivation. Needed whenever  $V$  can derive  $\varepsilon$ .

# WHY DO WE NEED THEM?

Let the current input symbol be  $b$ , and the top of the stack be  $X$ .

There are two ways we might try to derive a string starting with  $b$ :

1. Derive a string starting with  $b$  from the variable  $X$ . Look for a rule  $X \rightarrow \alpha$  where  $b \in FIRST(\alpha)$
2. If  $X$  can be derived to  $\varepsilon$ , then we might be able to derive a string starting with  $b$  by using the symbol(s) *following*  $X$  on the stack.  
i.e. If any of the production rules  $X \rightarrow \alpha$  had  $\varepsilon \in FIRST(\alpha)$ , then we also look at  $FOLLOW(X)$

## ANOTHER WAY TO CALCULATE FIRST SETS

Let  $A \rightarrow X_1 \dots X_n$  be a production of  $A$ , where  $X_i$  could be a terminal or variable.

We recursively compute  $FIRST(X_1 \dots X_n)$  by looking at  $X_1$ :

- ▶ If  $X_1$  is a terminal symbol, then  $FIRST(X_1 \dots X_n) = \{X_1\}$
- ▶ If  $X_1$  is a variable, then  $FIRST(X_1 \dots X_n)$  contains  $FIRST(X_1) \setminus \{\varepsilon\}$
- ▶ If  $X_1$  is a variable  $\varepsilon \in FIRST(X_1)$  then  $FIRST(X_1 \dots X_n)$  also contains  $FIRST(X_2 \dots X_n)$

Don't forget that  $FIRST(\varepsilon) = \{\varepsilon\}$ , so if every  $X_i$  can generate  $\varepsilon$ , then rule 3 will (eventually) give us  $\varepsilon \in FIRST(X_1 \dots X_n)$

# ANOTHER WAY TO CALCULATE FOLLOW SETS

If  $\varepsilon \in FIRST(\alpha)$  for some production rule  $A \rightarrow \alpha$  then we need to compute  $FOLLOW(A)$

Consider each production rule where  $A$  appears in the *right hand side*.

Let  $V \rightarrow Y_1 \dots Y_n A Z_1 \dots Z_m$  ( $Y_i, Z_i$  can be terminals or variables)

- ▶ If  $A$  is the start symbol, then  $\$ \in FOLLOW(A)$
- ▶  $FIRST(Z_1 \dots Z_m) \setminus \{\varepsilon\} \subseteq FOLLOW(A)$
- ▶ If  $\varepsilon \in FIRST(Z_1 \dots Z_m)$  then  $FOLLOW(V) \subseteq FOLLOW(A)$



# EXAMPLES

See week 5

# CONSTRUCTING THE PARSE TABLE

Rows: one for each variable of the grammar

Columns: one for each terminal of the grammar, and for the end of string marker \$

Steps to fill the table  $T$ :

1. If there is a rule  $R \rightarrow \alpha$  with  $b \in FIRST(\alpha)$  then put  $\alpha$  in  $T[R, b]$
2. If there is a rule  $R \rightarrow \alpha$  with  $\varepsilon \in FIRST(\alpha)$  and  $b \in FOLLOW(R)$ , then put  $\alpha$  in  $T[R, b]$

# EXAMPLE

See week 5

# IDENTIFY IF A GRAMMAR IS LL(1)

Recall that a grammar is LL(1) if it is sufficient to look at the next symbol to determine which rule to follow next.

i.e. If every cell in the LL(1) parse table contains at most one rule, then the grammar is LL(1)

More formally, a grammar is LL(1) iff for every variable  $A$ :

- ▶ Let  $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$  be the production rules for  $A$
- ▶ Let  $X_i = FIRST(\alpha_i)$  if  $\varepsilon \notin FIRST(\alpha_i)$
- ▶ Let  $X_i = FIRST(\alpha_i) \cup FOLLOW(A)$  otherwise
- ▶ Then  $X_i \cap X_j = \emptyset$  for all  $i \neq j$

# TRANSFORMING NON-LL(1) GRAMMARS

When a grammar is not LL(1) we try to find an equivalent grammar which is, by applying the following techniques:

- ▶ Left factoring
- ▶ Elimination of left recursion

# TRANSFORMING NON-LL(1) GRAMMARS

When a grammar is not LL(1) we try to find an equivalent grammar which is, by applying the following techniques:

- ▶ Left factoring
- ▶ Elimination of left recursion

Recall: grammars are *equivalent* if they generate the same language

# TRANSFORMING NON-LL(1) GRAMMARS

When a grammar is not LL(1) we try to find an equivalent grammar which is, by applying the following techniques:

- ▶ Left factoring
- ▶ Elimination of left recursion

Recall: grammars are *equivalent* if they generate the same language

Such a grammar does not always exist. For example no LL(k) grammar exists for the language

$$\{a^n b^n \mid n \geq 0\} \cup \{a^n b^{2n} \mid n \geq 0\}$$

# LEFT FACTORING: WHY?

Consider the grammar fragment  $S \rightarrow abcC \mid abdD$

- ▶ The two rules both start with the same prefix  $ab$
- ▶ i.e. their FIRST sets both include  $a$
- ▶ The LL(1) parse table will have multiple entries at  $(S, a)$ .



# LEFT FACTORING: WHY?

Consider the grammar fragment  $S \rightarrow abcC \mid abdD$

- ▶ The two rules both start with the same prefix  $ab$
- ▶ i.e. their FIRST sets both include  $a$
- ▶ The LL(1) parse table will have multiple entries at  $(S, a)$ .

We can “factor out” the string  $ab$  to obtain an equivalent grammar, where  $B$  is a new variable:

$$S \rightarrow abB$$

$$B \rightarrow cC \mid dD$$

This grammar fragment is equivalent, but is LL(1)

# LEFT FACTORING: DEFINITION

If a string  $w$  appears on the left of several rules for a variable  $A$ :

$$A \rightarrow wX_1 \mid \dots \mid wX_n$$

Then we can factor out  $w$  and introduce a new variable  $A'$ :

$$\begin{aligned} A &\rightarrow wA' \\ A' &\rightarrow X_1 \mid \dots \mid X_n \end{aligned}$$

Any other rules produced by  $A$  are unaffected.

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

- ▶ left-recursive: it occurs at the start of the string  $X \Rightarrow^+ X\beta$

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

- ▶ left-recursive: it occurs at the start of the string  $X \Rightarrow^+ X\beta$
- ▶ right-recursive: it occurs at the end of the string  $X \Rightarrow^+ \alpha X$

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

- ▶ left-recursive: it occurs at the start of the string  $X \Rightarrow^+ X\beta$
- ▶ right-recursive: it occurs at the end of the string  $X \Rightarrow^+ \alpha X$
- ▶ self-embedding: it occurs in between:  $X \Rightarrow^+ \alpha X\beta$

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

- ▶ left-recursive: it occurs at the start of the string  $X \Rightarrow^+ X\beta$
- ▶ right-recursive: it occurs at the end of the string  $X \Rightarrow^+ \alpha X$
- ▶ self-embedding: it occurs in between:  $X \Rightarrow^+ \alpha X \beta$

A grammar is recursive if any of its variables is recursive

# RECURSION

If a variable  $X$  can generate a string containing  $X$  itself, then it is recursive

- ▶ left-recursive: it occurs at the start of the string  $X \Rightarrow^+ X\beta$
- ▶ right-recursive: it occurs at the end of the string  $X \Rightarrow^+ \alpha X$
- ▶ self-embedding: it occurs in between:  $X \Rightarrow^+ \alpha X\beta$

A grammar is recursive if any of its variables is recursive

A grammar for an infinite language must contain at least one recursive variable



# ELIMINATE LEFT RECURSION: WHY?

Consider this simple grammar:

$$A \rightarrow c$$

$$A \rightarrow Ab$$

$$FIRST(c) = \{c\}$$

$$FIRST(Ab) = \{c\}$$

# ELIMINATE LEFT RECURSION: WHY?

Consider this simple grammar:

$$A \rightarrow c$$

$$A \rightarrow Ab$$

$$FIRST(c) = \{c\}$$

$$FIRST(Ab) = \{c\}$$

If we try to construct the parse table:

	$b$	$c$	$\$$
$A$		$c$ or $Ab$	

# ELIMINATE LEFT RECURSION: WHY?

Consider this simple grammar:

$$A \rightarrow c$$

$$A \rightarrow Ab$$

$$FIRST(c) = \{c\}$$

$$FIRST(Ab) = \{c\}$$

If we try to construct the parse table:

	$b$	$c$	$\$$
$A$		$c$ or $Ab$	

The base cases for the recursion must have FIRST sets which intersect with the left recursive rule!

# ELIMINATING LEFT RECURSION

Let  $\alpha, \beta$  be arbitrary strings of terminals and/or variables.

Let  $A$  be a variable, and  $R$  a new variable

If  $A$  has left recursive rules:

$$A \rightarrow A\alpha \mid \beta$$

It can be replaced with:

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \varepsilon$$

# ELIMINATING LEFT RECURSION

Let  $\alpha, \beta$  be arbitrary strings of terminals and/or variables.

Let  $A$  be a variable, and  $R$  a new variable

If  $A$  has left recursive rules:

$$A \rightarrow A\alpha \mid \beta$$

It can be replaced with:

$$A \rightarrow \beta R$$

$$R \rightarrow \alpha R \mid \varepsilon$$

What do the parse trees look like for  $\beta\alpha\alpha\alpha$  using the original and transformed grammar?

# SIMPLE EXAMPLE

$$A \rightarrow c$$

$$A \rightarrow Ab$$

Then  $\alpha =$

# SIMPLE EXAMPLE

$$A \rightarrow c$$

$$A \rightarrow Ab$$

Then  $\alpha = b$ ,  $\beta =$

# SIMPLE EXAMPLE

$$A \rightarrow c$$

$$A \rightarrow Ab$$

Then  $\alpha = b$ ,  $\beta = c$ , which gives us:

$$A \rightarrow cR$$

$$R \rightarrow bR \mid \varepsilon$$



# SIMPLE EXAMPLE

$$A \rightarrow c$$

$$A \rightarrow Ab$$

Then  $\alpha = b$ ,  $\beta = c$ , which gives us:

$$A \rightarrow cR$$

$$R \rightarrow bR \mid \varepsilon$$

	$b$	$c$	\$
$A$		$cR$	
$R$	$bR$		$\varepsilon$

# COMPLEX EXAMPLE

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow a \mid b \mid c$$

Then  $\alpha =$

# COMPLEX EXAMPLE

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow a \mid b \mid c$$

Then  $\alpha = +T \mid -T$ ,  $\beta =$

# COMPLEX EXAMPLE

$$E \rightarrow E + T$$

$$E \rightarrow E - T$$

$$E \rightarrow T$$

$$T \rightarrow a \mid b \mid c$$

Then  $\alpha = +T \mid -T$ ,  $\beta = T$ , which gives us:

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) =$$

# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) = \{a, b, c\}$$

$$FIRST(+TR) =$$

# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) = \{a, b, c\}$$

$$FIRST(+TR) = \{+\}$$

$$FIRST(-TR) =$$

# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) = \{a, b, c\}$$

$$FIRST(+TR) = \{+\}$$

$$FIRST(-TR) = \{-\}$$

Because  $\varepsilon \in FIRST(\varepsilon)$ , we calculate  $FOLLOW(R) =$



# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) = \{a, b, c\}$$

$$FIRST(+TR) = \{+\}$$

$$FIRST(-TR) = \{-\}$$

Because  $\varepsilon \in FIRST(\varepsilon)$ , we calculate  $FOLLOW(R) = \{\$, \}$

# COMPLEX EXAMPLE

$$E \rightarrow TR$$

$$R \rightarrow +TR \mid -TR \mid \varepsilon$$

$$T \rightarrow a \mid b \mid c$$

$$FIRST(TR) = \{a, b, c\}$$

$$FIRST(+TR) = \{+\}$$

$$FIRST(-TR) = \{-\}$$

Because  $\varepsilon \in FIRST(\varepsilon)$ , we calculate  $FOLLOW(R) = \{\$ \}$

	$a$	$b$	$c$	$+$	$-$	$\$$
$E$	$TR$	$TR$	$TR$			
$R$				$+TR$	$-TR$	$\varepsilon$
$T$	$a$	$b$	$c$			

# PROVING THAT A GRAMMAR IS *not* LL(1)

It is sufficient to show any one of the following:

- ▶ The grammar is left recursive
- ▶ The grammar needs left factoring
- ▶ The first sets of the production rules for a variable are not disjoint

# TYPICAL EXAM QUESTION

Consider the grammar  $G$ :

$$S \rightarrow ST \mid ab$$

$$T \rightarrow aTbb \mid ab$$

Show that the grammar  $G$  is not LL(1)

Transform  $G$  to obtain a grammar  $G'$  which is LL(1)

Give the LL(1) parse table for  $G'$

# CHOMSKY NORMAL FORM

Noam Chomsky, 1959

A grammar  $G$  is in Chomsky Normal Form iff every rule is of one of these forms:

- ▶  $A \rightarrow BC$  (neither  $B$  nor  $C$  are the start symbol)
- ▶  $A \rightarrow a$
- ▶  $S \rightarrow \varepsilon$  (iff  $S$  is the start symbol and  $\varepsilon \in L(G)$ )

Any CFG can be transformed into an equivalent Chomsky Normal Form.

We will use an algorithm which increases the size of the grammar to at most  $|G|^2$

# CHOMSKY NORMAL FORM: BUT WHY??

Useful properties of CNF:

- ▶ No rules (except  $S$ ) can produce  $\varepsilon$ .
- ▶ The derivations are always *binary trees*

The *CYK bottom up parsing algorithm* takes advantage of these

- ▶ It has a running time of  $O(n^3 \cdot |G|)$ , if  $G$  is in CNF
- ▶ i.e. effectively  $O(n^3 \cdot |G|^2)$  if  $G$  was not in CNF (ignoring the one-off cost of performing the grammar transformations.)
- ▶ One of the most efficient parsers in terms of *worst case* asymptotic complexity.

# CHOMSKY NORMAL FORM: ALGORITHM

1. **Eliminate the start symbol from all production rules**
2. Eliminate rules with non-solitary terminals
3. Eliminate rules with more than two variables
4. Eliminate epsilon productions
5. Eliminate unit rules

Add a new start symbol and rule  $S_0 \rightarrow S$

# CHOMSKY NORMAL FORM: ALGORITHM

1. Eliminate the start symbol from all production rules
2. **Eliminate rules with non-solitary terminals**
3. Eliminate rules with more than two variables
4. Eliminate epsilon productions
5. Eliminate unit rules

Create a new rule  $N_a \rightarrow a$  for every terminal which appears in one or more rules which are not of the form  $A \rightarrow a$ , then replace  $a$  with  $N_a$  in every such production rule.

e.g.  $A \rightarrow \alpha a \beta$  would become  $A \rightarrow \alpha N_a \beta$



# CHOMSKY NORMAL FORM: ALGORITHM

1. Eliminate the start symbol from all production rules
2. Eliminate rules with non-solitary terminals
3. **Eliminate rules with more than two variables**
4. Eliminate epsilon productions
5. Eliminate unit rules

For every rule of the form  $A \rightarrow X_1 \dots X_n$  (where  $n > 2$ ), delete it and create new variables  $A_1 \dots A_{n-2}$  and rules:

$$\begin{aligned}
 A &\rightarrow X_1 A_1 \\
 A_1 &\rightarrow X_2 A_2 \\
 &\dots \\
 A_{n-3} &\rightarrow X_{n-2} A_{n-2} \\
 A_{n-2} &\rightarrow X_{n-1} X_n
 \end{aligned}$$

# CHOMSKY NORMAL FORM: ALGORITHM

1. Eliminate the start symbol from all production rules
2. Eliminate rules with non-solitary terminals
3. Eliminate rules with more than two variables
4. **Eliminate epsilon productions**
5. Eliminate unit rules

For every rule of the form  $A \rightarrow \varepsilon$  (except  $S_0 \rightarrow \varepsilon$ )

- Delete the rule
- For each rule  $R \rightarrow \alpha$  containing  $A$ , create every possible rule  $R \rightarrow \alpha'$  where  $\alpha'$  is  $\alpha$  with one or more  $A$ 's removed.

Example: removing  $A \rightarrow \varepsilon$

If there is a rule  $R \rightarrow \alpha A \beta A \gamma$  then we add 3 new rules

$R \rightarrow \alpha \beta \gamma, R \rightarrow \alpha A \beta \gamma, R \rightarrow \alpha \beta A \gamma$  (and keep the original)

# CHOMSKY NORMAL FORM: ALGORITHM

1. Eliminate the start symbol from all production rules
2. Eliminate rules with non-solitary terminals
3. Eliminate rules with more than two variables
4. Eliminate epsilon productions
5. **Eliminate unit rules**

For each rule of the form  $A \rightarrow B$

- ▶ For each rule of the form  $B \rightarrow \alpha$  create a new rule  $A \rightarrow \alpha$
- ▶ Remove the rule  $A \rightarrow B$

# CHOMSKY NORMAL FORM: ALGORITHM

1. Eliminate the start symbol from all production rules
2. Eliminate rules with non-solitary terminals
3. Eliminate rules with more than two variables
4. Eliminate epsilon productions
5. Eliminate unit rules



All done! The grammar should be equivalent to the original one, but in Chomsky Normal Form.

# CHOMSKY NORMAL FORM: EXAMPLE

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

Eliminate start symbol from all production rules:

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid aB$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

Eliminate rules with non-solitary terminals:

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid \mathbf{N_a}B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$\mathbf{N_a} \rightarrow \mathbf{a}$$

# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow ASA \mid N_a B$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$N_a \rightarrow a$$

Eliminate rules with more than two variables:

$$S_0 \rightarrow S$$

$$S \rightarrow \mathbf{AS_1} \mid N_a B$$

$$\mathbf{S_1 \rightarrow SA}$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$N_a \rightarrow a$$

# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B$$

$$S_1 \rightarrow SA$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b \mid \varepsilon$$

$$N_a \rightarrow a$$

Eliminate epsilon production  $B \rightarrow \varepsilon$

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid \mathbf{N_a}$$

$$S_1 \rightarrow SA$$

$$A \rightarrow B \mid S \mid \varepsilon$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$



# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid N_a$$

$$S_1 \rightarrow SA$$

$$A \rightarrow B \mid S \mid \varepsilon$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$

Eliminate epsilon production  $A \rightarrow \varepsilon$

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid N_a \mid \mathbf{S_1}$$

$$S_1 \rightarrow SA \mid \mathbf{S}$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

$$N \rightarrow a$$

# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid N_a \mid S_1$$

$$S_1 \rightarrow SA \mid S$$

$$A \rightarrow B \mid S$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$

Eliminate unit rules  $S \rightarrow N_a$ ,  $A \rightarrow B$ , then  $S \rightarrow S_1$

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid \mathbf{a} \mid \mathbf{SA} \quad (\text{and useless } S \rightarrow S)$$

$$S_1 \rightarrow SA \mid S$$

$$A \rightarrow \mathbf{b} \mid S$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$

# CHOMSKY NORMAL FORM: EXAMPLE

$$S_0 \rightarrow S$$

$$S \rightarrow AS_1 \mid N_a B \mid a \mid SA$$

$$S_1 \rightarrow SA \mid S$$

$$A \rightarrow b \mid S$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$

Eliminate unit rules  $S_0 \rightarrow S$ ,  $S_1 \rightarrow S$ ,  $A \rightarrow S$

$$S_0 \rightarrow \mathbf{AS_1} \mid \mathbf{N_a B} \mid \mathbf{a} \mid \mathbf{SA}$$

$$S \rightarrow AS_1 \mid N_a B \mid a \mid SA$$

$$S_1 \rightarrow \mathbf{AS_1} \mid \mathbf{N_a B} \mid \mathbf{a} \mid SA$$

$$A \rightarrow b \mid \mathbf{AS_1} \mid \mathbf{N_a B} \mid \mathbf{a} \mid \mathbf{SA}$$

$$B \rightarrow b$$

$$N \rightarrow a$$

# CHOMSKY NORMAL FORM: EXAMPLE

All done!

$$S_0 \rightarrow AS_1 \mid N_a B \mid a \mid SA$$

$$S \rightarrow AS_1 \mid N_a B \mid a \mid SA$$

$$S_1 \rightarrow AS_1 \mid N_a B \mid a \mid SA$$

$$A \rightarrow b \mid AS_1 \mid N_a B \mid a \mid SA$$

$$B \rightarrow b$$

$$N_a \rightarrow a$$

# GREIBACH NORMAL FORM

Sheila Greibach, 1965

A grammar is in Greibach Normal Form iff every rule is in the form

- ▶  $A \rightarrow a$
- ▶  $A \rightarrow aB_1 \dots B_n$  (where all  $B_i$  are variables)
- ▶  $S \rightarrow \varepsilon$  (iff  $\varepsilon \in L(G)$ )

Any CFG can be transformed into an equivalent Greibach Normal Form.

# GREIBACH NORMAL FORM: BUT WHY??

Useful properties of GNF:

- ▶ The grammar does not contain left recursion
- ▶ (but there can be many rules starting with the same terminal)
- ▶ *Every derivation step consumes one input symbol*

# GREIBACH NORMAL FORM: BUT WHY??

Useful properties of GNF:

- ▶ The grammar does not contain left recursion
- ▶ (but there can be many rules starting with the same terminal)
- ▶ *Every derivation step consumes one input symbol*

Too good to be true?

# GREIBACH NORMAL FORM: BUT WHY??

Useful properties of GNF:

- ▶ The grammar does not contain left recursion
- ▶ (but there can be many rules starting with the same terminal)
- ▶ *Every derivation step consumes one input symbol*

Too good to be true?

GNF grammars can be very large, i.e.  $O(|G|^4)$  in general, or  $O(|G|^3)$  for languages without  $\epsilon$ .



# GREIBACH NORMAL FORM: BUT WHY??

Useful properties of GNF:

- ▶ The grammar does not contain left recursion
- ▶ (but there can be many rules starting with the same terminal)
- ▶ *Every derivation step consumes one input symbol*

Too good to be true?

GNF grammars can be very large, i.e.  $O(|G|^4)$  in general, or  $O(|G|^3)$  for languages without  $\epsilon$ .

However, for some languages they can result in extremely efficient parsers, because top-down parsers will only need to perform  $n$  derivation steps.

# GREIBACH NORMAL FORM: IDEA

The full algorithm is a little tedious, but the general idea is:

1. Eliminate all left recursions
2. Eliminate null productions (except on  $S'$ )
3. Make substitutions to transform the grammar into the proper form:
  - ▶ Expand the first variable in each rule until we get a terminal on the left
  - ▶ Cut cycles which cannot reach a terminal

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow AB \mid B$$

$$A \rightarrow Aa \mid b$$

$$B \rightarrow Ab \mid c$$

Eliminate left recursion

$$S \rightarrow AB \mid B$$

$$\mathbf{A} \rightarrow \mathbf{bR}$$

$$\mathbf{R} \rightarrow \mathbf{aR} \mid \varepsilon$$

$$B \rightarrow Ab \mid c$$

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow AB \mid B$$

$$A \rightarrow bR$$

$$R \rightarrow aR \mid \varepsilon$$

$$B \rightarrow Ab \mid c$$

Eliminate null productions

$$S \rightarrow AB \mid B$$

$$A \rightarrow bR \mid \mathbf{b}$$

$$R \rightarrow aR \mid \mathbf{a}$$

$$B \rightarrow Ab \mid c$$

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow AB \mid B$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

Substitute variables on the left side of the productions of  $S$

$$S \rightarrow \mathbf{bRB} \mid \mathbf{bB} \mid \mathbf{Ab} \mid \mathbf{c}$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow bRB \mid bB \mid Ab \mid c$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

Substitute variables on the left side of the productions of  $S$ , again

$$S \rightarrow bRB \mid bB \mid \mathbf{bRb} \mid \mathbf{bb} \mid c$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow bRB \mid bB \mid bRb \mid bb \mid c$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

Substitute variables on the left side of the productions of  $B$

$$S \rightarrow bRB \mid bB \mid bRb \mid bb \mid c$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow \mathbf{bRb} \mid \mathbf{bb} \mid c$$

# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow bRB \mid bB \mid bRb \mid bb \mid c$$

$$A \rightarrow bR \mid b$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow Ab \mid c$$

$A$  is unreachable now, delete it

$$S \rightarrow bRB \mid bB \mid bRb \mid bb \mid c$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow bRb \mid bb \mid c$$



# GREIBACH NORMAL FORM: EXAMPLE

$$S \rightarrow bRB \mid bB \mid bRb \mid bb \mid c$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow bRb \mid bb \mid c$$

Make a new variable  $N_b$ , to substitute in for the  $b$ 's that are not at the start of a production

$$S \rightarrow bRB \mid bB \mid bRN_b \mid bN_b \mid c$$

$$R \rightarrow aR \mid a$$

$$B \rightarrow bRN_b \mid bN_b \mid c$$

$$N_b \rightarrow b$$

All done!

# REVIEW

- ▶ Identifying when a grammar is LL(1)
  - ▶ If left factoring is needed, then it is *not* LL(1)
  - ▶ If left recursion exists, then it is *not* LL(1)
  - ▶ It is LL(1) iff all positions in the table contain at most one production rule

# REVIEW

- ▶ Identifying when a grammar is LL(1)
  - ▶ If left factoring is needed, then it is *not* LL(1)
  - ▶ If left recursion exists, then it is *not* LL(1)
  - ▶ It is LL(1) iff all positions in the table contain at most one production rule
- ▶ Transforming grammars to LL(1)
  - ▶ Not always possible
  - ▶ Try left factoring
  - ▶ Try removing left recursion

# REVIEW

- ▶ Identifying when a grammar is LL(1)
  - ▶ If left factoring is needed, then it is *not* LL(1)
  - ▶ If left recursion exists, then it is *not* LL(1)
  - ▶ It is LL(1) iff all positions in the table contain at most one production rule
- ▶ Transforming grammars to LL(1)
  - ▶ Not always possible
  - ▶ Try left factoring
  - ▶ Try removing left recursion
- ▶ Chomsky Normal Form
  - ▶ Efficient parsing
  - ▶ Easy to build
  - ▶ All CFL have a CNF

# REVIEW

- ▶ Identifying when a grammar is LL(1)
  - ▶ If left factoring is needed, then it is *not* LL(1)
  - ▶ If left recursion exists, then it is *not* LL(1)
  - ▶ It is LL(1) iff all positions in the table contain at most one production rule
- ▶ Transforming grammars to LL(1)
  - ▶ Not always possible
  - ▶ Try left factoring
  - ▶ Try removing left recursion
- ▶ Chomsky Normal Form
  - ▶ Efficient parsing
  - ▶ Easy to build
  - ▶ All CFL have a CNF
- ▶ Greibach Normal Form
  - ▶ Derivation depth is the length of input
  - ▶ Can be excellent for parsing
  - ▶ Sometimes prohibitively large
  - ▶ All CFL have a GNF