

# Lecture 10

## Context Free Grammars

Slides by David Albrecht (2011), modified by Graham Farr (2013).

FIT2014 Theory of Computation

# Overview

- Inductive Definitions
- Context Free Grammars
- Parse Trees
- Derivations

# Arithmetic Expressions

1. All integers are Arithmetic Expressions
2. If A and B are Arithmetic Expressions, so are:
  - (i)  $A + B$
  - (ii)  $A - B$
  - (iii)  $A * B$
  - (iv)  $A / B$
  - (v)  $(A)$

# Production Rules

$AE \rightarrow \text{integer}$

$AE \rightarrow AE + AE$

$AE \rightarrow AE - AE$

$AE \rightarrow AE * AE$

$AE \rightarrow AE / AE$

$AE \rightarrow (AE)$

$S \rightarrow A$

$A \rightarrow \text{integer}$

$A \rightarrow A + A$

$A \rightarrow A - A$

$A \rightarrow A * A$

$A \rightarrow A / A$

$A \rightarrow (A)$

# Backus-Naur Form

(a.k.a. Backus Normal Form)



John Backus (1924-2007)

<http://www-history.mcs.st-and.ac.uk/Biographies/Backus.html>

$$S \rightarrow A$$
$$A \rightarrow \text{integer} \mid A + A \mid A - A \mid A * A \mid A / A \mid (A)$$


Peter Naur (b. 1928)

<http://datamuseum.dk/>

# Historical example: fragment of the BNF of ALGOL 60

## 4.1. COMPOUND STATEMENTS AND BLOCKS

### 4.1.1. Syntax

$\langle \text{unlabelled basic statement} \rangle ::= \langle \text{assignment statement} \rangle |$   
 $\langle \text{go to statement} \rangle | \langle \text{dummy statement} \rangle | \langle \text{procedure statement} \rangle$   
 $\langle \text{basic statement} \rangle ::= \langle \text{unlabelled basic statement} \rangle | \langle \text{label} \rangle :$   
 $\langle \text{basic statement} \rangle$   
 $\langle \text{unconditional statement} \rangle ::= \langle \text{basic statement} \rangle | \langle \text{for statement} \rangle |$   
 $\langle \text{compound statement} \rangle | \langle \text{block} \rangle$   
 $\langle \text{statement} \rangle ::= \langle \text{unconditional statement} \rangle |$   
 $\langle \text{conditional statement} \rangle$   
 $\langle \text{compound tail} \rangle ::= \langle \text{statement} \rangle \text{ end } | \langle \text{statement} \rangle \quad ;$   
 $\langle \text{compound tail} \rangle$   
 $\langle \text{block head} \rangle ::= \text{begin} \langle \text{declaration} \rangle | \langle \text{block head} \rangle \quad ;$   
 $\langle \text{declaration} \rangle$   
 $\langle \text{unlabelled compound} \rangle ::= \text{begin } \langle \text{compound tail} \rangle$   
 $\langle \text{unlabelled block} \rangle ::= \langle \text{block head} \rangle \quad ; \langle \text{compound tail} \rangle$   
 $\langle \text{compound statement} \rangle ::= \langle \text{unlabelled compound} \rangle |$   
 $\langle \text{label} \rangle : \langle \text{compound statement} \rangle$   
 $\langle \text{block} \rangle ::= \langle \text{unlabelled block} \rangle | \langle \text{label} \rangle : \langle \text{block} \rangle$

from: J. W. Backus *et al.*, *Comm. ACM* 3 (5) (May 1960) 299-314.

# Definitions

- A string is in **EQUAL** if it has an equal number of **a**'s and **b**'s.
- E.g.  
 $\{\epsilon, ab, ba, aabb, abab, abba, baba, \dots\}$
- An **a-type** string has one more **a** than **b**.
- A **b-type** string has one more **b** than **a**.

# Definition of EQUAL

- A string is in **EQUAL** if it is
  - $\epsilon$ , or
  - an **a** followed by a string of **b-type**, or
  - it is a **b** followed by a string of **a-type**.
- A string is of **a-type** if it is
  - an **a**, or
  - an **a** followed by a string in **EQUAL**, or
  - a **b** followed by two strings of **a-type**.
- A string is of **b-type** if it is
  - an **b**, or
  - a **b** followed by a string in **EQUAL**, or
  - an **a** followed by two strings of **b-type**.

$$S \rightarrow \epsilon$$

$$S \rightarrow aB$$

$$S \rightarrow bA$$

$$A \rightarrow a$$

$$A \rightarrow aS$$

$$A \rightarrow bAA$$

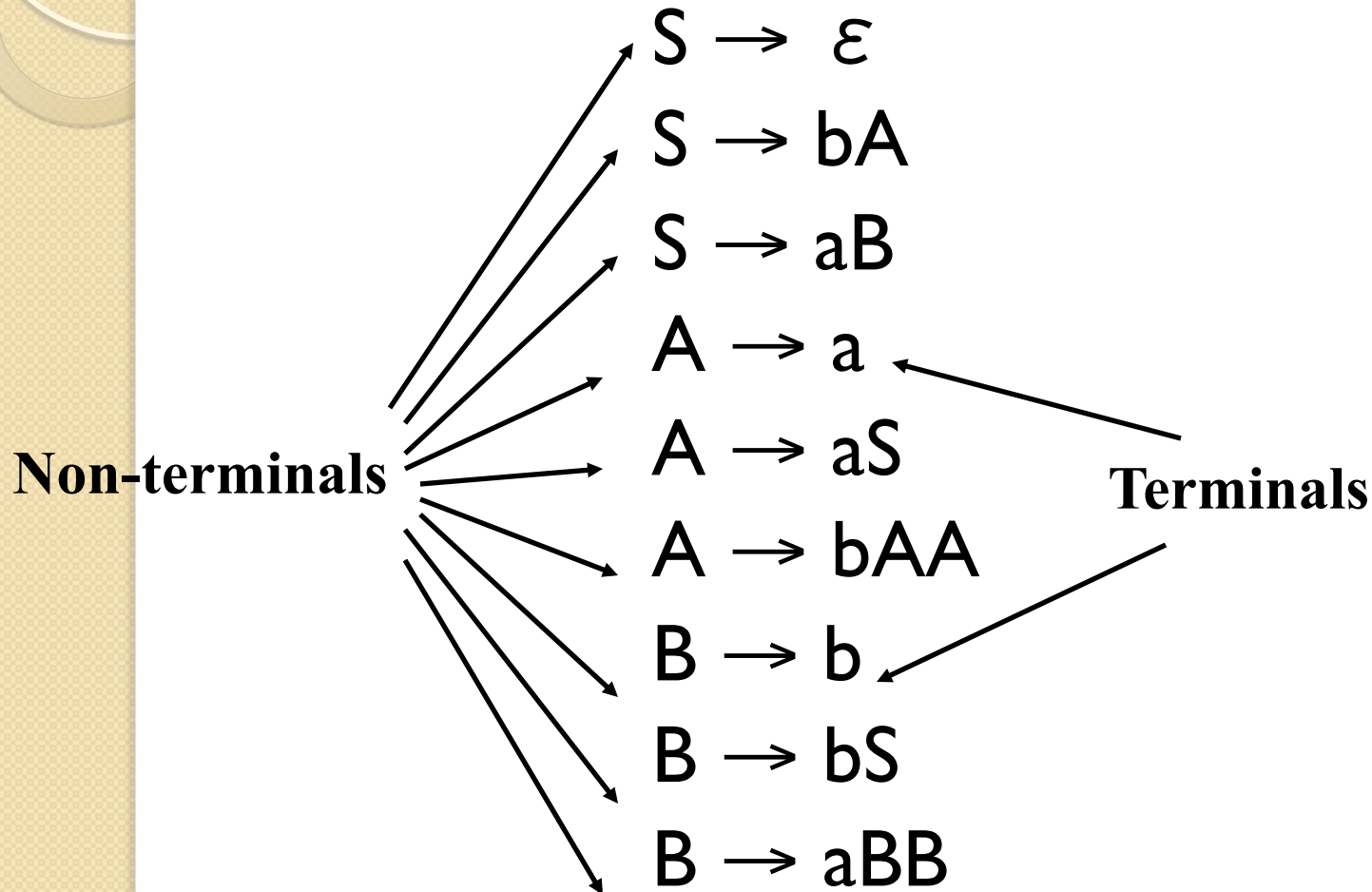
$$B \rightarrow b$$

$$B \rightarrow bS$$

$$B \rightarrow aBB$$



# Production Rules



# A Context Free Grammar consists of:

## 1. An alphabet

- The letters are called **terminals**

## 2. A set of symbols

- We call these symbols **nonterminals**
- One of these symbols is the **Start symbol**
- **S** is often used as the start symbol.

## 3. A finite set of production rules of the form:

One nonterminal  $\rightarrow$  finite string of terminals  
and/or nonterminals

# Definition

- The ***language generated*** by a Context Free Grammar (CFG)
  - consists of those strings which can be produced from the start symbol using the production rules.
- A language generated by a CFG is called a **Context Free Language (CFL)**.

# History

- Pāṇini (c520BC-c460BC)
  - studied *Sanskrit*
- Noam Chomsky (b. 1928)
  - studied *natural languages*
- John Backus
  - studied *programming languages*



<http://www-groups.dcs.st-and.ac.uk/history/Biographies/Panini.html>



Noam Chomsky,  
during visit to Australia in 2011 to  
accept Sydney Peace Prize.

<http://www.abc.net.au/news/2011-06-02/noam-chomsky/2741826>

$$S \rightarrow aS \mid Sa \mid \varepsilon$$

## Derivation of aaaa

1.  $S \rightarrow Sa$

2.  $S \rightarrow aS$

3.  $S \rightarrow \varepsilon$

$$S \Rightarrow Sa \quad (\text{Rule 1})$$

$$\Rightarrow aSa \quad (\text{Rule 2})$$

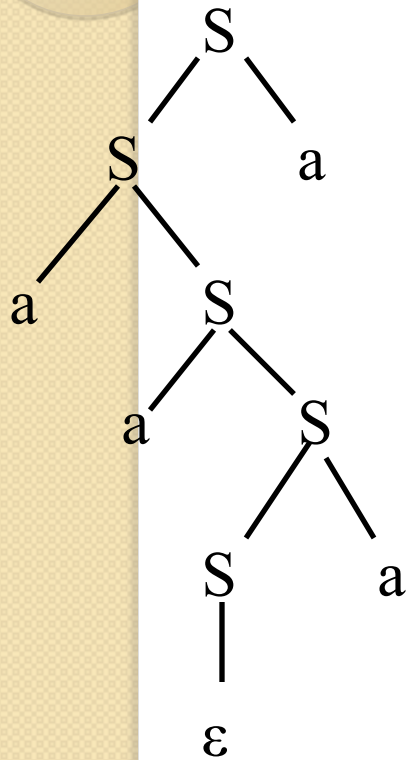
$$\Rightarrow aaSa \quad (\text{Rule 2})$$

$$\Rightarrow aaSaa \quad (\text{Rule 1})$$

$$\Rightarrow aa\varepsilon aa \quad (\text{Rule 3})$$

$$= aaaa$$

# Parse Tree



## Derivation of aaaa

$S \Rightarrow Sa$  (Rule 1)

$\Rightarrow aSa$  (Rule 2)

$\Rightarrow aaSa$  (Rule 2)

$\Rightarrow aaSaa$  (Rule 1)

$\Rightarrow aa\epsilon aa$  (Rule 3)

$= aaaa$

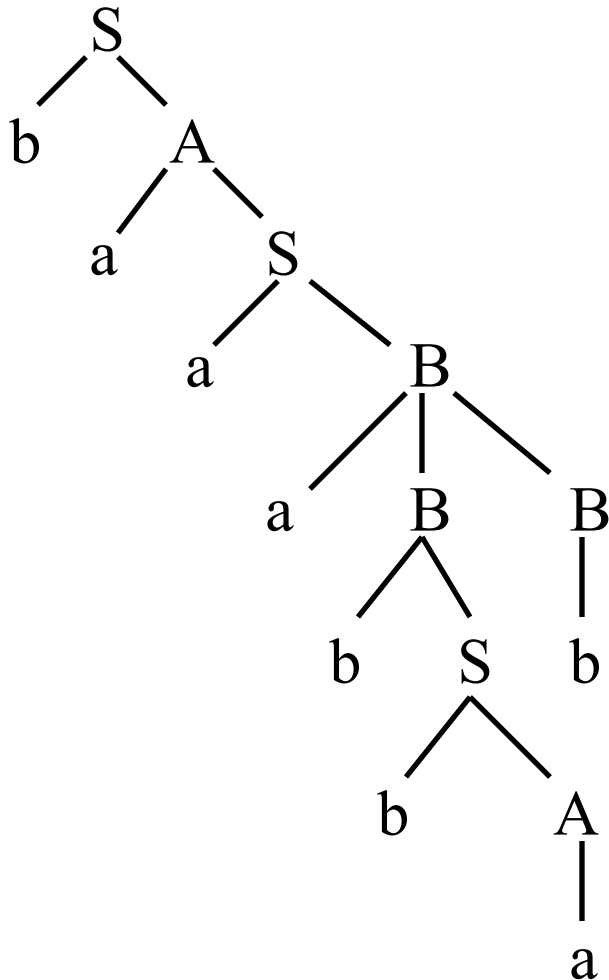
# EQUAL

## Derivation of baaabbab

1.  $S \rightarrow \varepsilon$
2.  $S \rightarrow bA$
3.  $S \rightarrow aB$
4.  $A \rightarrow a$
5.  $A \rightarrow aS$
6.  $A \rightarrow bAA$
7.  $B \rightarrow b$
8.  $B \rightarrow bS$
9.  $B \rightarrow aBB$

$S \Rightarrow bA$  (Rule 2)  
 $\Rightarrow baS$  (Rule 5)  
 $\Rightarrow baaB$  (Rule 3)  
 $\Rightarrow baaaBB$  (Rule 9)  
 $\Rightarrow baaaBb$  (Rule 7)  
 $\Rightarrow baaabSb$  (Rule 8)  
 $\Rightarrow baaabbAb$  (Rule 2)  
 $\Rightarrow baaabbab$  (Rule 4)

# Parse Tree



$S \Rightarrow bA$

$\Rightarrow baS$

$\Rightarrow baaB$

$\Rightarrow baaaBB$

$\Rightarrow baaaBb$

$\Rightarrow baaabSb$

$\Rightarrow baaabbAb$

$\Rightarrow baaabbbab$



# PARENTHESES (a.k.a. the Dyck Language)

1.  $S \rightarrow \varepsilon$
2.  $S \rightarrow (S)$
3.  $S \rightarrow SS$

Derivation of  $()()$

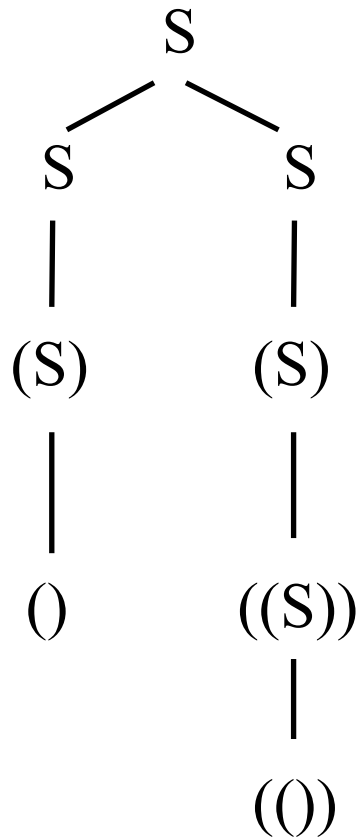
$S \Rightarrow SS$	(Rule 3)
$\Rightarrow (S)S$	(Rule 2)
$\Rightarrow (S)(S)$	(Rule 2)
$\Rightarrow ()(S)$	(Rule 1)
$\Rightarrow ()((S))$	(Rule 2)
$\Rightarrow ()()$	(Rule 1)



Walther von Dyck (1856-1934)

[http://www-history.mcs.st-and.ac.uk/Biographies/Von\\_Dyck.html](http://www-history.mcs.st-and.ac.uk/Biographies/Von_Dyck.html)

# Parse Tree



$S \Rightarrow SS$

$\Rightarrow (S)S$

$\Rightarrow (S)(S)$

$\Rightarrow ()(S)$

$\Rightarrow ()((S))$

$\Rightarrow ()()$

# Exercises

- Suppose we have two types of brackets, such as round and square:  $()$  and  $[]$ . Find a context-free language for the set of all valid strings of such brackets.
- Find a context-free grammar for **PALINDROMES**
- For other languages we have met:
  - find context-free grammars for them, OR
  - if you think they don't have one, think about *why*.

# A simple property of derivations

At any stage, the string to the left of the first nonterminal must be a prefix of the final (derived) string.

.....

$x_1 \dots x_k \mathbf{A} \mathbf{B} \dots$

$\Rightarrow x_1 \dots x_k \mathbf{aXYB} \dots$  (using  $\mathbf{A} \rightarrow \mathbf{aXY}$ )

.....

$\Rightarrow x_1 \dots x_k \mathbf{a} \dots \dots$  (derived string)

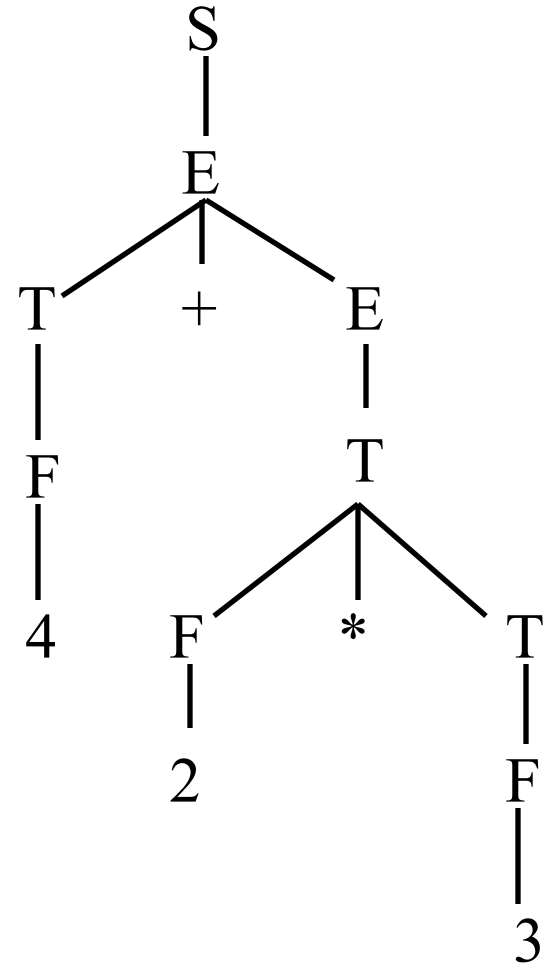
$4 + 2 * 3$

$S \rightarrow E$

$E \rightarrow T + E \mid T - E \mid T$

$T \rightarrow F * T \mid F / T \mid F$

$F \rightarrow \text{integer} \mid (E)$



$$4 + 2 * 3$$

$$S \rightarrow E$$

$$E \rightarrow T + E \mid T - E \mid T$$

$$T \rightarrow F * T \mid F / T \mid F$$

$$F \rightarrow \text{integer} \mid (E)$$

**Leftmost Derivation**

$$S \Rightarrow E$$

$$\Rightarrow T + E$$

$$\Rightarrow F + E$$

$$\Rightarrow 4 + E$$

$$\Rightarrow 4 + T$$

$$\Rightarrow 4 + F * T$$

$$\Rightarrow 4 + 2 * T$$

$$\Rightarrow 4 + 2 * F$$

$$\Rightarrow 4 + 2 * 3$$

**4 + 2\*3**

**$S \rightarrow E$**

**$E \rightarrow T + E \mid T - E \mid T$**

**$T \rightarrow F * T \mid F / T \mid F$**

**$F \rightarrow \text{integer} \mid (E)$**

**Rightmost Derivation**

**$S \Rightarrow E$**

**$\Rightarrow T + E$**

**$\Rightarrow T + T$**

**$\Rightarrow T + F * T$**

**$\Rightarrow T + F * F$**

**$\Rightarrow T + F * 3$**

**$\Rightarrow T + 2 * 3$**

**$\Rightarrow F + 2 * 3$**

**$\Rightarrow 4 + 2 * 3$**

# Leftmost and rightmost derivations

- In a **Leftmost derivation**, the leftmost non-terminal is always replaced first.
- In a **Rightmost derivation**, the rightmost non-terminal is always replaced first.

## Theorem

Whenever a string has a derivation, it also has a leftmost derivation of the same length.

**Proof:** see Tute 3.

Does the same hold for rightmost derivations?



# A simple property of leftmost derivations

Whenever a production

**X**  $\rightarrow$  **terminals** **Non-terminal** **theRest**

is applied, the **terminal letters on the left** are appended to the current prefix to give a larger prefix of the derived string

.....

$x_1 \dots x_k \mathbf{A} \mathbf{B} \dots$

$\Rightarrow x_1 \dots x_k \mathbf{a} \mathbf{X} \mathbf{Y} \mathbf{B} \dots$  (using  $\mathbf{A} \rightarrow \mathbf{a} \mathbf{X} \mathbf{Y}$ )

.....

$\Rightarrow x_1 \dots x_k \mathbf{a} \dots \dots$  (derived string)

# Revision

- Context Free Grammars
  - Definition. How to use them.
- Parse Trees
  - Definition. How to make them.
- Be able to construct leftmost and rightmost derivations.
- Read
  - M. Sipser, , “***Introduction to the Theory of Computation***”, Chapter 2.