

# Tema 4: Gramáticas independientes del contexto

# Bibliografía

- Hopcroft, J. E., Motwani, R., y Ullman, J. D. “Introducción a la Teoría de Autómatas, Lenguajes y Computación”. Addison Wesley. 2002.
  - capítulo 5
- Sudkamp, Thomas A. “Languages and machines : an introduction to the theory of computer science”. Addison-Wesley Publishing Company, 1998.
  - capítulos 3 y 4

# Introducción

- Lenguajes independientes del contexto (LIC)
- Gramáticas independientes del contexto (GIC)
  - implementación de analizadores sintácticos
  - descripción de formatos de documentos: XML
- Árboles sintácticos: representan gráficamente la estructura de la gramática
- Autómatas con pila

# Un ejemplo de GIC

- Palíndromo:  $w = w^R$
- Este lenguaje no es regular. Se comprueba aplicando el lema del bombeo a la cadena  $0^n10^n$
- Para generar el lenguaje de los palíndromos con el alfabeto  $\{0, 1\}$ :
  - Base:  $\varepsilon, 0, 1$  son palíndromos
  - Paso inductivo: si  $w$  es palíndromo, también lo son  $0w0$  y  $1w1$
- Reglas que definen la gramática que genera el lenguaje de los palíndromos con alfabeto  $\{0, 1\}$ :
  - $P \rightarrow \varepsilon$
  - $P \rightarrow 0$
  - $P \rightarrow 1$
  - $P \rightarrow 0P0$
  - $P \rightarrow 1P1$
  - En notación compacta:  $P \rightarrow \varepsilon \mid 0 \mid 1 \mid 0P0 \mid 1P1$

# Definición de gramática

- Las GIC están formadas por cuatro componentes:
  - el alfabeto de símbolos terminales, que son los símbolos finales del lenguaje
  - el conjunto finito de símbolos no terminales o variables, que permiten representar subconjuntos del lenguaje o estados intermedios en la generación de las palabras del lenguaje
  - el símbolo inicial o axioma de la gramática (una de las variables), a partir de cual se obtiene cualquier palabra del lenguaje
  - un conjunto finito de producciones o reglas, que indican las transformaciones posibles desde los símbolos no terminales a las palabras del lenguaje. Las reglas están formadas por:
    - una variable, cabeza de la producción
    - el símbolo de producción  $\rightarrow$
    - una cadena de 0 o más símbolos terminales y variables, que son el cuerpo de la producción
    - las producciones son de la forma  $B \rightarrow x$ , donde  $B \in V$ , y  $x \in (V \cup T)^*$
- $G = (V, T, P, S)$ 
  - ejemplo:  $G_{\text{palíndromo}} = (\{P\}, \{0, 1\}, A, P)$ , donde  $A$  son las producciones o reglas:  $P \rightarrow \varepsilon \mid 0 \mid 1 \mid 0P0 \mid 1P1$

# Definición de gramática (II)

- Ejemplo:  $G = (\{E, I\}, \{+, *, (, ), a, b, 0, 1\}, P, E)$ , donde  $P$  es el conjunto de producciones:

1.  $E \rightarrow I$   
2.  $E \rightarrow E + E$   
3.  $E \rightarrow E * E$   
4.  $E \rightarrow (E)$   
5.  $I \rightarrow a$   
6.  $I \rightarrow b$   
7.  $I \rightarrow Ia$   
8.  $I \rightarrow Ib$   
9.  $I \rightarrow I0$   
10.  $I \rightarrow I1$

# Derivaciones de una gramática

- Aplicando las producciones de una GIC se infieren las cadenas que están en el lenguaje
- Sea  $G = (V, T, P, S)$ , y  $\alpha A \beta$  una cadena de terminales y variables, donde  $A$  está en  $V$ , y  $\alpha$  y  $\beta$  están en  $(V \cup T)^*$ . Sea  $A \rightarrow \gamma$  una producción de  $G$ . Entonces:

$$\alpha A \beta \xRightarrow[G]{} \alpha \gamma \beta$$

- Extensión a varios pasos de derivación

$$\text{– si } \alpha \xRightarrow[G]^* \beta \text{ y } \beta \xRightarrow[G]^* \gamma, \text{ entonces } \alpha \xRightarrow[G]^* \gamma$$

- Ejemplo: inferencia de “a \* (a + b00)”

$$\begin{aligned} E &\Rightarrow E * E \Rightarrow I * E \Rightarrow a * E \Rightarrow \\ a * (E) &\Rightarrow a * (E + E) \Rightarrow a * (I + E) \Rightarrow a * (a + E) \Rightarrow \\ a * (a + I) &\Rightarrow a * (a + I0) \Rightarrow a * (a + I00) \Rightarrow a * (a + b00) \end{aligned}$$

# Derivaciones izquierda y derecha

- Derivación más a la izquierda:  $\xRightarrow{mi}^*$
- Derivación más a la derecha:  $\xRightarrow{md}^*$
- Ejemplo:

$$E \xRightarrow{mi}^* a * (a + b00)$$

$$E \xRightarrow{mi} E * E \xRightarrow{mi} I * E \xRightarrow{mi} a * E \xRightarrow{mi} a * (E) \xRightarrow{mi} a * (E + E) \xRightarrow{mi} a * (I + E) \xRightarrow{mi} a * (a + E) \xRightarrow{mi} a * (a + I) \xRightarrow{mi} a * (a + I0) \xRightarrow{mi} a * (a + I00) \xRightarrow{mi} a * (a + b00)$$

$$E \xRightarrow{md}^* a * (a + b00)$$

$$E \xRightarrow{md} E * E \xRightarrow{md} E * (E) \xRightarrow{md} E * (E + E) \xRightarrow{md} E * (E + I) \xRightarrow{md} E * (E + I0) \xRightarrow{md} E * (E + I00) \xRightarrow{md} E * (E + b00) \xRightarrow{md} E * (I + b00) \xRightarrow{md} E * (a + b00) \xRightarrow{md} I * (a + b00) \xRightarrow{md} a * (a + b00)$$



# Lenguaje de una gramática

- Si  $G = (V, T, P, S)$  es una GIC, el lenguaje de  $G$  será:

- $L(G) = \{w \text{ que están en } T^* \mid S \xRightarrow[G]{*} w\}$

- Notación:

- $a, b, \dots, +, (, \dots$ : símbolos terminales
  - $A, B, \dots$ : variables
  - $\dots, w, x, y, z$ : cadenas de terminales
  - $\dots, X, Y, Z$ : terminales o variables
  - $\alpha, \beta, \gamma, \dots$ : cadenas de terminales y variables

# Formas sentenciales

- Si  $G = (V, T, P, S)$  es una GIC, cualquier cadena  $\alpha$  en  $(V \cup T)^*$  tal que  $S \Rightarrow \alpha$  es una forma sentencial
- El lenguaje  $L(G)$  está formado por las formas sentenciales que están en  $T^*$ : se denominan sentencias
- Ejemplo:

$$E \Rightarrow E * E \Rightarrow E * (E) \Rightarrow E * (E + E) \Rightarrow E * (I + E)$$

$$E \underset{mi}{\Rightarrow} E * E \underset{mi}{\Rightarrow} I * E \underset{mi}{\Rightarrow} a * E$$

$$E \underset{md}{\Rightarrow} E * E \underset{md}{\Rightarrow} E * (E) \underset{md}{\Rightarrow} E * (E + E)$$

# Problemas

1. Diseñar las GIC para los siguientes lenguajes
  1.  $\{0^n 1^n \mid n \geq 1\}$
  2.  $\{a^i b^j c^k \mid i \neq j \text{ o } j \neq k\}$
2. La siguiente gramática genera el lenguaje de expresiones regulares  $0^* 1 (0+1)^*$ :  $S \rightarrow A1B$ ,  $A \rightarrow 0A \mid \varepsilon$ ,  $B \rightarrow 0B \mid 1B \mid \varepsilon$ . Obtener las derivaciones más a la izquierda y más a la derecha para la cadena: 00101
3. Sea  $T = \{0, 1, (, ), +, *, \emptyset, e\}$ . Construir una GIC con el conjunto de terminales  $T$  que genere expresiones regulares válidas con el alfabeto  $\{0, 1\}$  (considerar  $e$  como  $\varepsilon$ ).

# Árboles de derivación

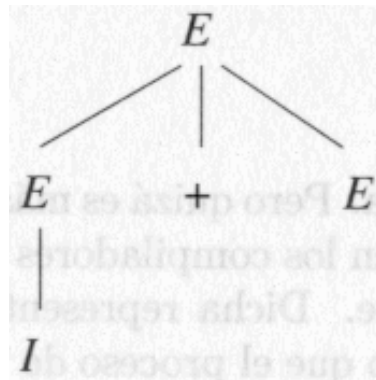
- Representación de las derivaciones en forma de árbol
  - muestra la agrupación de los símbolos de una cadena terminal en subcadenas
  - utilización en los compiladores para representar la estructura sintáctica del programa fuente
- Terminología
  - los árboles son conjuntos de nodos unidos entre sí por la relación *padre-hijo*. Un nodo tiene como máximo un padre y cero o más hijos
  - nodo raíz: no tiene padre
  - hojas: nodos sin hijos
  - nodos interiores: nodos que no son hojas
  - el hijo de un hijo de ...: descendiente
  - el padre de un padre de ...: antepasado
  - los hijos de un nodo se ordenan de izquierda a derecha

# Construcción de árboles de derivación

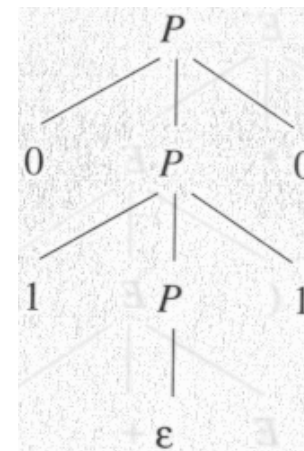
- Sea  $G = (V, T, P, S)$ . El árbol de derivación para  $G$  tendrá las siguientes características:
  - cada nodo interior está etiquetado con una variable
  - cada hoja está etiquetada con una variable, un terminal o  $\varepsilon$ . Si es  $\varepsilon$ , tiene que ser el único hijo de su padre
  - si un nodo interior está etiquetado con  $A$ , y sus hijos están etiquetados con  $X_1, X_2, \dots, X_k$  (de izquierda a derecha), entonces  $A \rightarrow X_1X_2\dots X_k$  es una producción de  $P$

- Ejemplos:

$$E \xRightarrow{*} I + E$$



$$P \xRightarrow{*} 0110$$



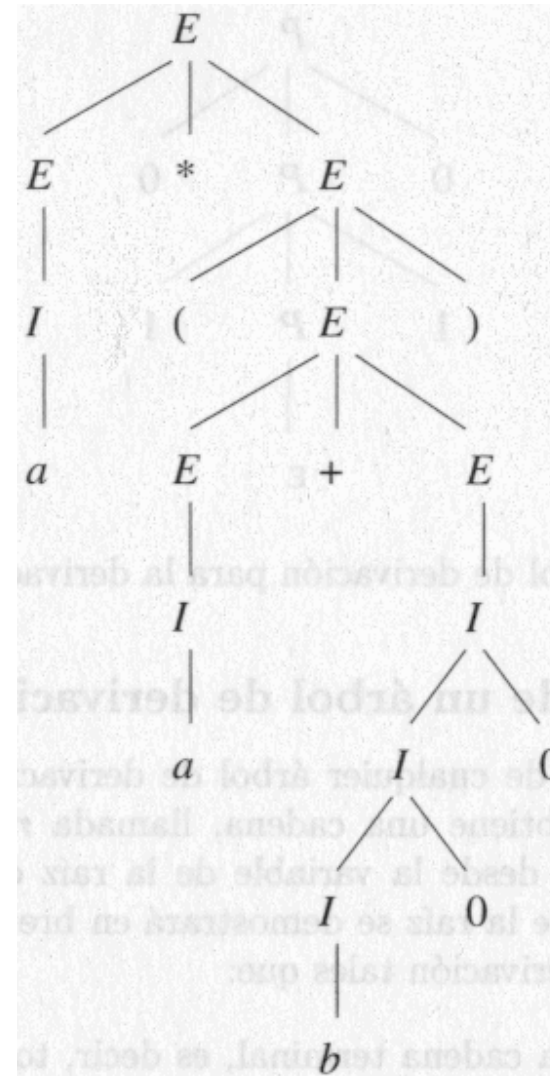
# Resultado de un árbol de derivación

- Concatenando la hojas de un árbol desde la izquierda se obtiene la cadena resultado del árbol, que se deriva desde la raíz
- Son especialmente importantes los árboles de derivación tales que:
  - su resultado es una cadena terminal
  - la raíz está etiquetada con el símbolo inicial
- Estos árboles tienen como resultado las cadenas del lenguaje generado por esa gramática

# Un ejemplo

$$E \Rightarrow^* a^*(a + b00)$$

1.  $E \rightarrow I$
2.  $E \rightarrow E + E$
3.  $E \rightarrow E * E$
4.  $E \rightarrow (E)$
5.  $I \rightarrow a$
6.  $I \rightarrow b$
7.  $I \rightarrow Ia$
8.  $I \rightarrow Ib$
9.  $I \rightarrow I0$
10.  $I \rightarrow I1$



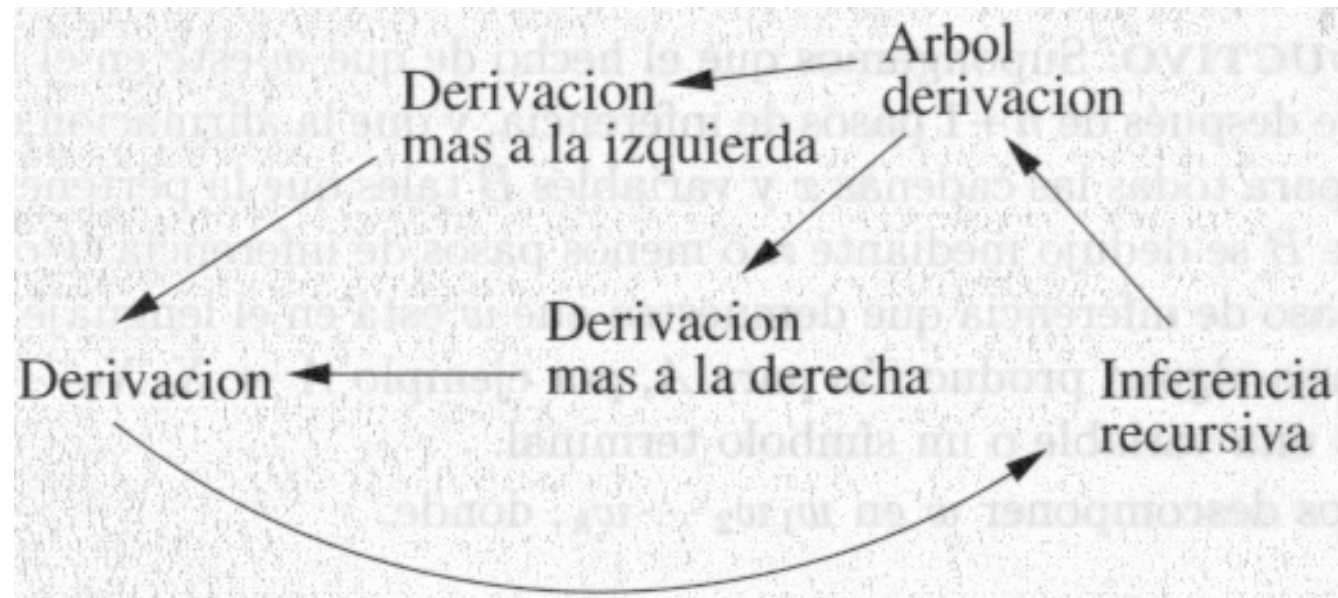
# Inferencia, derivaciones y árboles

- Dada  $G = (V, T, P, S)$ , las siguientes afirmaciones son equivalentes:
  - el procedimiento de inferencia recursiva determina que la cadena terminal  $w$  pertenece al lenguaje de la variable  $A$

–  $A \xRightarrow{*} w$

–  $A \xRightarrow{mi} w$

–  $A \xRightarrow{md} w$



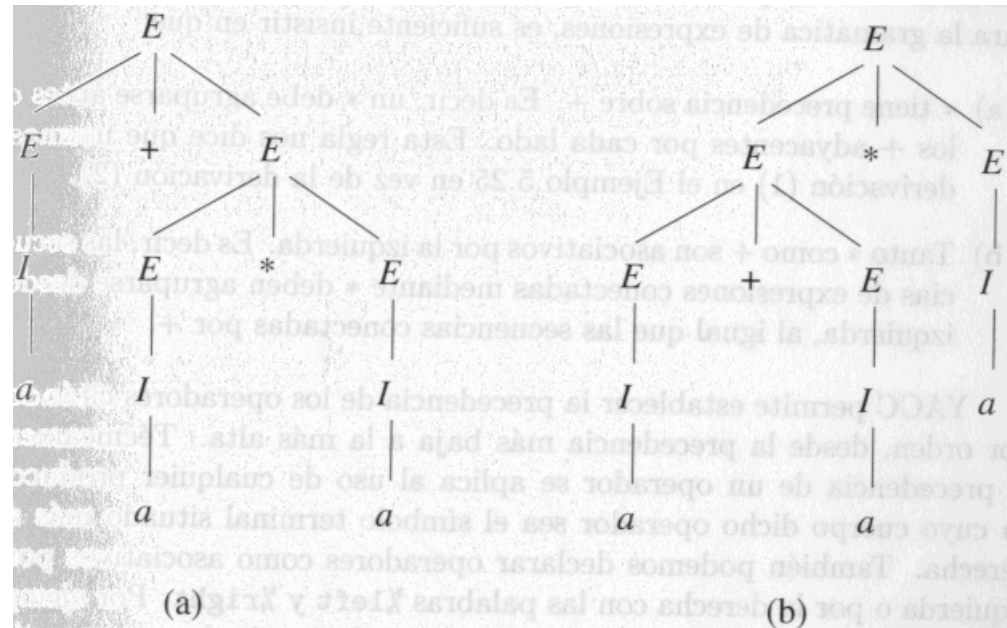
- Existe un árbol de derivación con la raíz etiquetada con  $A$  y resultado  $w$



# Ambigüedad

- Una GIC  $G = (V, T, P, S)$  es ambigua si existe al menos una cadena  $w$  en  $T^*$  para la que podemos encontrar dos árboles de derivación distintos con la raíz etiquetada con  $S$  y cuyo resultado es  $w$

1.	$E \rightarrow I$
2.	$E \rightarrow E + E$
3.	$E \rightarrow E * E$
4.	$E \rightarrow (E)$
5.	$I \rightarrow a$
6.	$I \rightarrow b$
7.	$I \rightarrow Ia$
8.	$I \rightarrow Ib$
9.	$I \rightarrow IO$
10.	$I \rightarrow II$



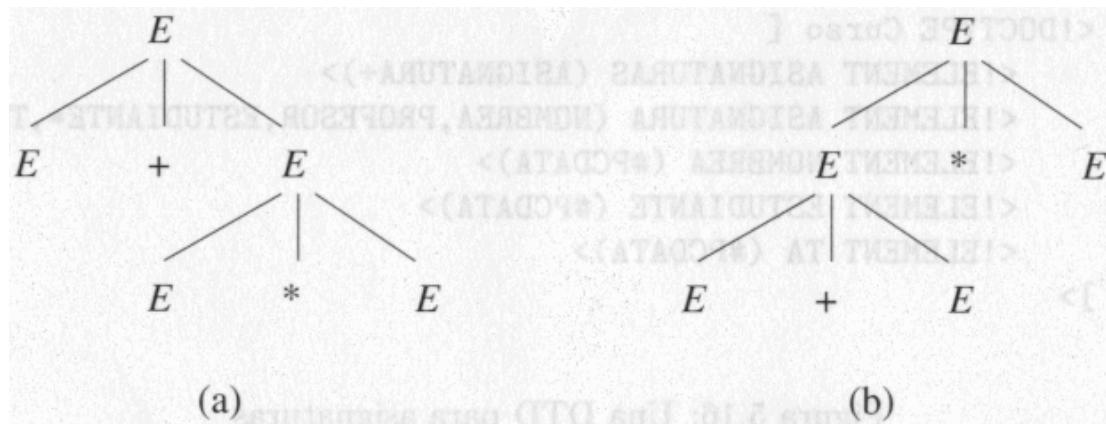
- La existencia de derivaciones diferentes para una cadena no supone un defecto en la gramática, *la existencia de árboles de derivación diferentes sí supone un problema*

- $E \rightarrow E + E \rightarrow I + E \rightarrow a + E \rightarrow a + I \rightarrow a + b$
- $E \rightarrow E + E \rightarrow E + I \rightarrow E + b \rightarrow I + b \rightarrow a + b$

# Eliminación de la ambigüedad

- No existe un algoritmo que nos indique si una GIC es ambigua
- Existen LIC que sólo tienen GIC ambiguas: inherentemente ambiguos
- Para las construcciones de los lenguajes de programación comunes existen técnicas para la eliminación de la ambigüedad
- Ejemplo: causas de ambigüedad en la siguiente gramática
  - no se respeta la precedencia de operadores
  - una secuencia de operadores idénticos puede agruparse desde la izquierda y desde la derecha. Lo convencional es agrupar desde la izquierda

```
1.  E → I
2.  E → E + E
3.  E → E * E
4.  E → (E)
5.  I → a
6.  I → b
7.  I → Ia
8.  I → Ib
9.  I → I0
10. I → I1
```

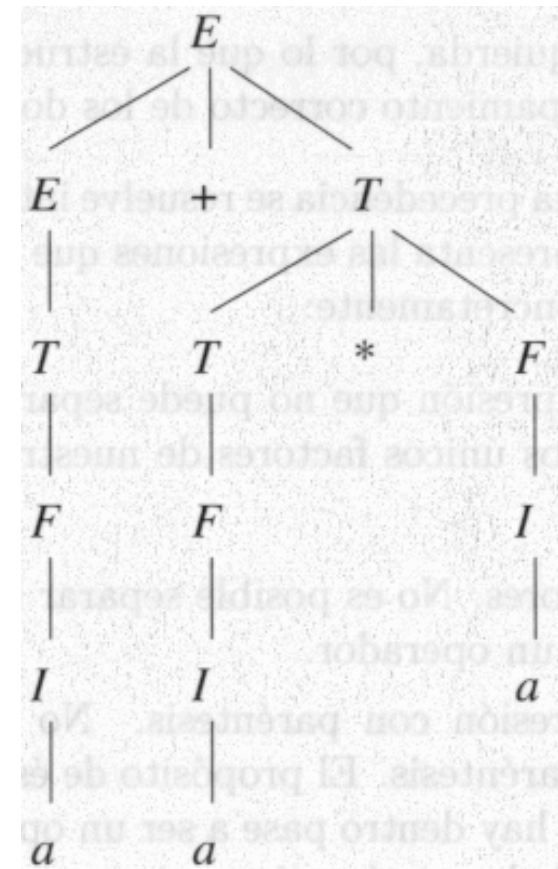


# Eliminación de la ambigüedad (II)

- Ejemplo: modificamos la gramática para forzar la precedencia

1.  $E \rightarrow I$   
 2.  $E \rightarrow E + E$   
 3.  $E \rightarrow E * E$   
 4.  $E \rightarrow (E)$   
 5.  $I \rightarrow a$   
 6.  $I \rightarrow b$   
 7.  $I \rightarrow Ia$   
 8.  $I \rightarrow Ib$   
 9.  $I \rightarrow I0$   
 10.  $I \rightarrow I1$

$I \rightarrow a \mid b \mid Ia \mid Ib \mid I0 \mid I1$   
 $F \rightarrow I \mid (E)$   
 $T \rightarrow F \mid T * F$   
 $E \rightarrow T \mid E + T$

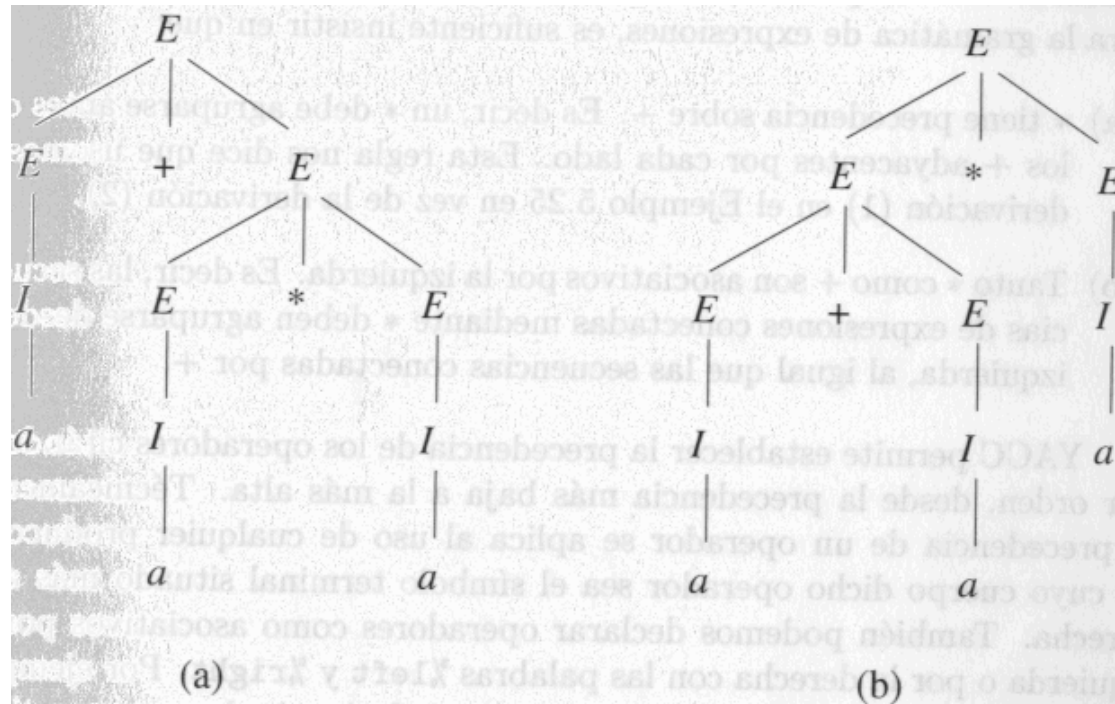


# Derivaciones como forma de expresar ambigüedad

- En una gramática no ambigua las derivaciones más a la derecha y más a la izquierda serán únicas
- Teorema: para toda gramática  $G = (V, T, P, S)$  y toda cadena  $w$  en  $T^*$ ,  $w$  tiene dos árboles de derivación distintos si y sólo si  $w$  tiene dos derivaciones más a la izquierda distintas desde  $S$ 
  - lo mismo para las derivaciones más a la derecha

# Derivaciones como forma de expresar ambigüedad (II)

- Ejemplo:



a)  $E \Rightarrow E + E \Rightarrow I + E \Rightarrow a + E \Rightarrow a + E * E \Rightarrow a + I * E \Rightarrow a + a * E \Rightarrow a + a * I \Rightarrow a + a * a$

b)  $E \Rightarrow E * E \Rightarrow E + E * E \Rightarrow I + E * E \Rightarrow a + E * E \Rightarrow a + I * E \Rightarrow a + a * E \Rightarrow a + a * I \Rightarrow a + a * a$

# Problemas finales

1. Un palíndromo es una cadena que se lee igual de izquierda a derecha que de derecha a izquierda. Por ejemplo, las palabras “radar”, “oso” y “abba” son palíndromos. Dado el alfabeto  $\Sigma = \{a, b\}$ , determinar una gramática que describa palíndromos. La gramática debería generar palabras como “abba”, “aba”, “bb”, “babab”, “a”, “b”, ..., y  $\varepsilon$ .
2. Determinar una gramática que genere el lenguaje  $L = \{a^i b^j c^k \mid i=j \text{ o } j=k, i, j, k > 0\}$ .
3. Determinar una gramática que genere el lenguaje  $L = \{a^i b^j a^j b^i \mid i, j > 0\}$ ,