# ST0244 Lenguajes de Programacion Sintaxis

Andrés Sicard Ramírez

Universidad EAFIT

Semestre 2024-1

#### **Preliminares**

#### Convenciones

- La numeración (capítulos, teoremas, figuras, páginas, etc) en estas diapositivas corresponde a la numeración del texto guía [Lee 2017].
- Los ejemplos que incluyen código fuente están en el repositorio del curso.

Preliminares 2/45

#### Introducción

#### Sintaxis y semántica

- La sintaxis son las reglas para escribir código fuente en un lenguaje de programación (well-formed programs).
- La semántica establece el significado de los programas.

Introducción 3/45

#### Introducción

#### Sintaxis y semántica

- La sintaxis son las reglas para escribir código fuente en un lenguaje de programación (well-formed programs).
- La semántica establece el significado de los programas.

## Pregunta

¿Cuando se está aprendiendo y/o usando un lenguaje de programación son la sintaxis y la semántica igualmente de importantes?

Introducción 4/45

# Problemas sintácticos y semánticos

Tipo	Estático ( <i>compile-time</i> )	Dinámico ( <i>run-time</i> )
Sintaxis	$\checkmark$	
Semántica	$\checkmark$	$\checkmark$

Terminología 5/45

Ejemplo (pág. 32)

¿Es a = b + c; una instrucción correcta en C++?

Terminología 6/45

```
Ejemplo (pág. 32)
```

```
¿Es a = b + c; una instrucción correcta en C++?
```

## Algunas preguntas:

- 1. ¿Tienen b y c valores? (sí es resuelto en *run-time*, es un problema semántico dinámico, pero sí es resuelto en *compile-time*, es un problema semántico estático).
- 2. ¿Han b y c sido declaradas de un tipo que permite la operación +? (resuelto en *compile-time*, problema semántico estático).
- 3. Es la asignación a a compatible con el resultado de la expresión b + c (resuelto en *compile-time*, problema semántico estático).
- 4. ¿Tiene la instrucción de asignación la forma correcta? (resuelto en *compile-time*, problema sintáctico).

Terminología 7/45

#### Definición

Un símbolo terminal (o token) es un símbolo elemental del lenguaje.

# Ejemplo

Palabras reservadas (keywords), tipos simples, operadores, números, identificadores, entre otros.

Terminología 8/45

#### Definición

Un símbolo **no terminal** (o **categoría sintáctica**) representa una secuencia de símbolos terminales.

# Ejemplo

Tipos compuestos, instrucciones, expresiones, sentencias if, aplicación y abstracción de funciones, entre otros.

Terminología 9/45

#### Definición

La **forma de Backus-Naur** (BNF) es un meta-lenguaje formal para especificar la sintaxis de un lenguaje.

Formas de Backus-Naur 10/45

#### Reglas BNF

Una BNF para un lenguaje es un conjunto de reglas de la forma

```
\langle {
m no \ terminal} 
angle \ ::= \ {
m expresión}_1 \ | \ {
m expresión}_2 \ | \ \dots \ | \ {
m expresión}_n
```

# donde

- (i) expresión $_i$  es una cadena de terminales y no terminales,
- (ii) el símbolo ::= significa que el símbolo no terminal a la izquierda «tiene la forma» de la expresión a la derecha,
- (iii) el símbolo | significa «o».

Formas de Backus-Naur 11/45

## Ejemplo

Sea P un conjunto de variables proposicionales (fórmulas atómica) y sea  $p \in P$ . Podemos definir las fórmulas bien formadas (well-formed formulae) de la lógica proposicional por la siguiente BNF:

```
\label{eq:formula} \langle \text{f\'ormula} \rangle \quad ::= \quad \text{p} \\ \quad | \quad \neg \left< \text{f\'ormula} \right> \quad \wedge \left< \text{f\'ormula} \right> ) \\ \quad | \quad \left( \quad \left< \text{f\'ormula} \right> \vee \left< \text{f\'ormula} \right> \right) \\ \quad | \quad \left( \quad \left< \text{f\'ormula} \right> \rightarrow \left< \text{f\'ormula} \right> \right) \\ \quad | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left( \quad \left< \text{f\'ormula} \right> \leftrightarrow \left< \text{f\'ormula} \right> \right) \\ | \quad \left< \text{f\'ormula} \right> \to \left< \text{f\'ormula} \right>
```

#### Observación

Note que la definición de las fórmulas bien formadas es recursiva.

Formas de Backus-Naur 12/45

# Ejemplo

Una BNF describiendo los números enteros, con o sin signo (p. ej. -344, 56, +9784, 8, 000).

```
 \begin{split} &\langle \mathrm{entero} \rangle &::= \langle \mathrm{signo} \rangle \langle \mathrm{dígitos} \rangle \ | \ \langle \mathrm{dígitos} \rangle \\ &\langle \mathrm{dígitos} \rangle &::= \langle \mathrm{dígito} \rangle \langle \mathrm{dígitos} \rangle \ | \ \langle \mathrm{dígito} \rangle \\ &\langle \mathrm{dígito} \rangle &::= 0 \ | \ 1 \ | \ 2 \ | \ 3 \ | \ 4 \ | \ 5 \ | \ 6 \ | \ 7 \ | \ 8 \ | \ 9 \\ &\langle \mathrm{signo} \rangle &::= + | \ - \end{split}
```

Formas de Backus-Naur

## Ejemplo

La siguiente BNF genera las instrucciones de asignación en un lenguaje de programación.

```
\begin{split} \langle \mathrm{instrucci\'on-asignaci\'on} \rangle &::= \langle \mathrm{var} \rangle := \langle \mathrm{expr-aritm\'etica} \rangle \\ \langle \mathrm{expr-aritm\'etica} \rangle &::= \langle \mathrm{var} \rangle \mid \langle \mathrm{const} \rangle \mid \\ & \qquad \qquad (\langle \mathrm{expr-aritm\'etica} \rangle \langle \mathrm{op-aritm\'etico} \rangle \langle \mathrm{expr-aritm\'etica} \rangle ) \\ \langle \mathrm{op-aritm\'etico} \rangle &::= + \mid - \mid * \mid / \\ \langle \mathrm{const} \rangle &::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\ \langle \mathrm{var} \rangle &::= a \mid b \mid c \mid \dots \mid x \mid y \mid z \end{split}
```

Formas de Backus-Naur 14/45

## BNF extendida (EBNF)

La forma BFN se puede extender con las siguientes definiciones adicionales:

- i) ítem? o [ítem] significa que el ítem es opcional.
- ii) ítem $^*$  o {ítem} significa que cero o más ocurrencias del ítem son permitidas.
- iii) ítem<sup>+</sup> significa que una o más ocurrencias del ítem son permitidas.
- iv) Paréntesis pueden ser usados para agrupación.

Formas de Backus-Naur 15/45

## Ejemplo

Una BNF y una EBNF describiendo los números enteros, con o sin signo.

Formas de Backus-Naur 16/45

## Ejemplo

Una BNF y una EBNF describiendo una lista de identificadores.

```
\langle {\rm lista\text{-}id} \rangle \ ::= \ {\rm id} \ ; \ \langle {\rm lista\text{-}id} \rangle \ | \ {\rm id} \qquad \qquad \langle {\rm lista\text{-}id} \rangle \ ::= \ {\rm id} \ (; \ {\rm id})^* {\rm BNF} \qquad \qquad {\rm EBNF}
```

Formas de Backus-Naur 17/45

# Gramáticas independientes del contexto

### Descripción

Una **gramática independiente del contexto** es un meta-lenguaje formal para especificar la sintaxis de un lenguaje.

# Gramáticas independientes del contexto

#### Definición

Una gramática independiente del contexto (o gramática libre de contexto) (context-free grammar) es una estructura

$$G = (\mathcal{N}, \mathcal{T}, \mathcal{P}, \mathcal{S}),$$

donde

- (i)  $\mathcal N$  es un conjunto finito de símbolos no terminales,
- (ii)  $\mathcal{T}$  es un conjunto finito símbolos terminales,
- (iii)  $\mathcal{P}$  es un conjunto finito de producciones de la forma  $A \to \alpha$  con  $A \in \mathcal{N}$  y  $\alpha \in \{\mathcal{N} \cup \mathcal{T}\}^*$  (conjunto de símbolos terminales y no terminales incluyendo la palabra vacía  $\epsilon$ ),
- (iv)  $S \in \mathcal{N}$  es el símbolo de inicio.

# Gramáticas independientes del contexto

Ejemplo (gramática para expresiones infijas (§ 2.3.1))

Podemos definir una gramática independiente del contexto  $(\mathcal{N}, \mathcal{T}, \mathcal{P}, E)$  para expresiones infijas por

$$\mathcal{N} = \{E, T, F\},\$$
 $\mathcal{T} = \{\text{id}, \text{num}, +, -, *, /, (,)\},\$ 

y el conjunto  $\mathcal P$  está compuesto por las siguientes producciones:

$$E \rightarrow E + T$$
  $T \rightarrow T * F$   $F \rightarrow (E)$   $E \rightarrow E - T$   $T \rightarrow T / F$   $F \rightarrow \text{id}$   $F \rightarrow \text{num}$ 

#### Definición

Una **sentencia** (sentence) de una gramática G es una cadena de símbolos terminales (tokens) de G.

Derivaciones 21/45

#### Definición

Una **sentencia** (sentence) de una gramática G es una cadena de símbolos terminales (tokens) de G.

# **Ejemplo**

Dos sentencias de la gramática para las expresiones infijas son:

(i) 
$$(5 * x) + y$$
  
(ii)  $)4 + + ($ 

Derivaciones 22/45

#### Definición

Una **forma sentencial** ( $sentential\ form$ ) de una gramática G es una cadena de símbolos terminales y no terminales de G.

Derivaciones 23/45

#### Definición

Una forma sentencial (sentential form) de una gramática G es una cadena de símbolos terminales y no terminales de G.

# **Ejemplo**

Dos formas sentenciales de la gramática para las expresiones infijas son:

(i) 
$$(T * F) + T$$

(ii) (5 \* F F)

Derivaciones 24/45

#### Definición

Una **derivación** de una sentencia S en un gramática G es una secuencia de formas sentenciales de G que comienza en el símbolo inicial de G y termina en S.

Derivaciones 25/45

#### Definición

Una **derivación** de una sentencia S en un gramática G es una secuencia de formas sentenciales de G que comienza en el símbolo inicial de G y termina en S.

#### Observación

Cada forma sentencial en la derivación es obtenida a partir de la anterior reemplazando  $A \in \mathcal{N}$  (símbolo no terminal) por  $\alpha \in \{\mathcal{N} \cup \mathcal{T}\}^*$  (cadena de símbolos terminales y no terminales), si  $A \to \alpha$  es una producción de G.

Derivaciones 26/45

#### Definición

Un sentencia S de una gramática G es **válida** sii existe al menos una derivación para S en G.

Derivaciones 27/45

# Ejemplo

La contancia «/ E \* v ) · · · · de la gramática de las aversaismes inflica es válida margua tiona

La sentencia $((5 * x) + y)$ de la gramation	ca de las expresiones infijas es valida porq	ue tiene
la siguiente derivación:		
$\underline{E} \Rightarrow_1 \underline{E} + T$		
$\Rightarrow_3 \underline{T} + T$	$E  ightarrow E \ + \ T$	(1)
$\Rightarrow_6 \underline{F} + T$	E  ightarrow E~-~T	(2)
$\Rightarrow_7 (\underline{E}) + T$	E  o T	(3)
$\Rightarrow_2 (T) + T$	$T  o T \ * \ F$	(4)

$\Rightarrow_7 (\underline{E}) + T$	E o T	(3)
$\Rightarrow_3 (\underline{T}) + T$	$T  o T \ * \ F$	(4)
$\Rightarrow_4 (\underline{T} * F) + T$	$T  o T \; / \; F$	(5)
$\Rightarrow_6 (\underline{F} * F) + T$	T  o F	(6)
$\Rightarrow_9 (5 * \underline{F}) + T$	$F  o (\ E\ )$	(7)
$\Rightarrow_8 (5 * x) + \underline{T}$	$F o \mathrm{id}$	(8)
$\rightarrow$ (5 $+$ $x$ ) $+$ $F$	F  ightarrow num	(0)

#### Definición

El **lenguaje** de una gramática G, denotado L(G), es el conjunto de sentencias válidas de G.

Derivaciones 29/45

#### Tipos de derivaciones

- (i) **Derivación más a la izquierda** (*left-most derivation*): Siempre reemplaza el símbolo no terminal más a la izquierda).
- (ii) **Derivación más a la derecha** (*right-most derivation*): Siempre reemplaza el símbolo no terminal más a la derecha.

Derivaciones 30/45

# Ejemplo

Derivaciones

Derivaciones más a la izquierda y más a la derecha para la sentencia «( 5 \* x ) + y ».

$$\underline{E} \Rightarrow_{1} \underline{E} + T 
\Rightarrow_{3} \underline{T} + T 
\Rightarrow_{6} \underline{E} + T 
\Rightarrow_{7} (\underline{E}) + T$$

$$\underline{E} \Rightarrow_{1} E + \underline{T} 
\Rightarrow_{6} E + \underline{F} 
\Rightarrow_{8} \underline{E} + \underline{y} 
\Rightarrow_{8} \underline{E} + \underline{y} 
\Rightarrow_{3} T + \underline{y}$$

$$E \Rightarrow_{E} T$$

 $\Rightarrow_6 F + y$ 

 $\Rightarrow_7 (E) + y$ 

 $\Rightarrow_3 (T) + y$ 

 $\Rightarrow_A (T * F) + y$ 

 $\Rightarrow_8 (T * x) + y$ 

 $\Rightarrow_6 (F * x) + y$ 

 $\Rightarrow_0 (5 * x) + y$ 

Más a la derecha

(1)

(2)

(3)

(4)

(5)

(6)

(7)

(8)

(9)

31/45

 $T \to T * F$ 

 $T \to T / F$ 

 $F \rightarrow (E)$ 

 $T \to F$ 

 $F \to \mathrm{id}$ 

 $F \to \text{num}$ 

**Producciones** 

$$\Rightarrow_{6} \underline{F} + T \qquad \Rightarrow_{8} \underline{E} + y$$
$$\Rightarrow_{7} (\underline{E}) + T \qquad \Rightarrow_{3} T + y$$

$$\Rightarrow_8 \underline{E} + y$$
$$\Rightarrow_3 \underline{T} + y$$

$$\geq_1 \underline{E} + T$$

 $\Rightarrow_3 (T) + T$ 

 $\Rightarrow_4 (T * F) + T$ 

 $\Rightarrow_6 (F * F) + T$ 

 $\Rightarrow_9 (5 * F) + T$ 

 $\Rightarrow_8 (5 * x) + T$ 

 $\Rightarrow_6 (5 * x) + F$ 

 $\Rightarrow_8 (5 * x) + y$ 

Más a la izquierda

## Expresiones prefijas

En las expresiones prefijas el operador aparece antes que los operandos.

# Ejemplo

$$4 + (a - b) * x$$
 (expresión infija)  
  $+ 4 * - a b x$  (expresión prefija)

Derivaciones 32/45

Ejemplo (prefix expresiones gramática (§ 2.4.3))

Podemos definir una gramática independiente del contexto  $(\mathcal{N}, \mathcal{T}, \mathcal{P}, E)$  para expresiones infijas por

$$\mathcal{N} = \{E\},$$

$$\mathcal{T} = \{\text{id}, \text{num}, +, -, *, /\},$$

y el conjunto  ${\mathcal P}$  está compuesto por las siguientes producciones:

$$E \rightarrow + E E \mid -E E \mid *E E \mid /E E \mid id \mid num$$

Derivaciones 33/45

# Árboles de análisis sintáctico

#### Definición

Sea G una gramática. Un **árbol de análisis sintáctico** (parser tree) es una representación gráfica de la derivación de una sentencia de L(G).

Árboles de análisis sintáctico 34/45

# Árboles de análisis sintáctico

#### Definición

Sea G una gramática. Un **árbol de análisis sintáctico** (parser tree) es una representación gráfica de la derivación de una sentencia de L(G).

# **Propiedades**

Un árbol de análisis sintáctico de una sentencia de L(G) tiene las siguientes propiedades [Aho, Lam, Sethi y Ullman 2006, pág. 45]:

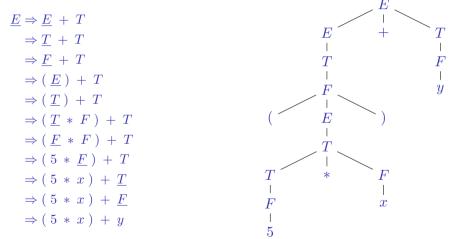
- (i) «La raíz se etiqueta con el símbolo inicial.»
- (ii) «Cada hoja se etiqueta con un terminal, o con  $\epsilon$ .»
- (iii) «Cada nodo interior se etiqueta con un no terminal.»
- (iv) «Si A es el no terminal que etiqueta a cierto nodo interior, y  $X_1, X_2, \ldots, X_n$  son las etiquetas de los hijos de ese nodo de izquierda a derecha, entonces debe haber una producción  $A \to X_1, X_2, \ldots, X_n$ .»

Árboles de análisis sintáctico 35/45

# Árboles de análisis sintáctico

# Ejemplo

Árbol de análisis sintáctico para la sentencia (5 \* x) + y.



Árboles de análisis sintáctico 36/45

# Árboles sintáctico abstractos

#### Definición

Un **árbol sintáctico abstracto** (abstract syntax tree) es un árbol de análisis sintáctico sin información esencial requerida para evaluar (generar código en la compilación o ejecutar en la interpretación) la sentencia (pág. 38):

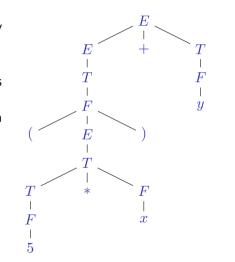
- i) «Non-terminal nodes in the tree are replaced by nodes that reflect the part of the sentencia they represent.»
- ii) «Unit productions in the tree are collapsed.»

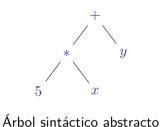
Árboles sintáctico abstractos 37/45

# Árboles sintáctico abstractos

# Ejemplo

Árbol de análisis sintáctico y árbol sintáctico abstracto (los nodos interiores representan operadores y las hojas representan operandos) para la sentencia (5 \* x) + y.





Árbol de análisis sintáctico

# Ambigüedad en las gramáticas

#### Definición

Un gramática G es **ambigua** sii existe una sentencia en L(G) que tiene más de un árbol de análisis sintáctico.

# Ambigüedad en las gramáticas

# Ejemplo

Dada la gramática  $(\mathcal{N}, \mathcal{T}, \mathcal{P}, E)$  donde

$$\mathcal{N} = \{E\},\$$

$$\mathcal{T} = \{*, +, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9\},\$$

$$E \to E \ * E \mid E + E$$

$$E \to 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

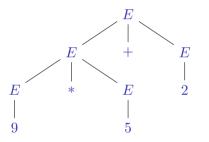
(continua en la próxima diapositiva)

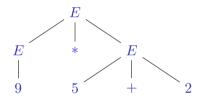
Ambigüedad en las gramáticas 40/45

# Ambigüedad en las gramáticas

Ejemplo (continuación)

La sentencia 9 \* 5 + 2 tiene dos árboles de análisis sintáctico.





# Limitaciones de las definiciones sintácticas

#### Algunas limitaciones

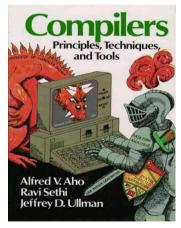
- La sintaxis de un lenguaje de programación es una descripción incompleta de éste (p. ej. 5 + 4/0).
- «The set of programs in any interesting lenguaje is not context-free.» (pág. 50) (p. ej. a + b)
- Un gramática independiente del contexto no especifica la semántica de un lenguaje de programación.

# Limitaciones de las definiciones sintácticas

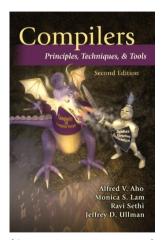
# Ejemplo (context-sensitive issues (págs. 50–1))

- $\bullet$  «In an array declaration in C++, the array size must be a non-negative value.»
- «Operands for the && operation must be boolean in Java.»
- «In a method definition, the return value must be compatible with the return type in the method declaration.»
- «When a method is called, the actual parameters must match the formal parameter types.»

# El libro del drágon



(Primera edición, 1986)



(Segunda edición, 2006)

## Referencias



Aho, Alfred V., Lam, Monica S., Sethi, Ravi y Ullman, Jeffrey D. [1986] (2006). Compilers: Principles, Techniques, & Tools. 2.ª ed. Addison-Wesley (vid. págs. 34, 35).



Lee, Kent D. [2014] (2017). Foundations of Programming Languages. 2.a ed. Undergraduate Topics in Computer Science. Springer (vid. pág. 2).

Referencias 45/45