

Lenguajes y Compiladores

Información: <http://www.cs.famaf.unc.edu.ar/wiki/>

Profesores:

Héctor Gramaglia, Miguel Pagano, Demetrio Vilela

Régimen de regularidad y Promoción

Se tomarán 2 parciales

Promoción: obteniendo al menos 7 en cada uno de los parciales y aprobando el taller.

Regularidad: aprobando 2 parciales (es posible recuperar 1) y aprobando el taller.

Taller: se implementarán intérpretes o compiladores de lenguajes de programación.

Estructura de la materia a grandes rasgos:

Primera Parte: Lenguaje imperativo

Segunda Parte: Lenguaje aplicativo puro, y lenguaje aplicativo con referencias y asignación

Ejes de Contenidos de la primera parte

- 1 Introducción a la sintaxis y la semántica de lenguajes
- 2 El problema de dar significado a la recursión e iteración
- 3 Un Lenguaje Imperativo

¿Cómo se otorga significado a un lenguaje?

El establecer el significado de las frases de un lenguaje de programación es un problema de múltiples aristas en tanto puede tener variados objetivos, que van desde la comprensión humana hasta la necesidad de que una máquina los pueda interpretar o traducir a una secuencia de instrucciones ejecutables.

Del significado trata un manual de usuario, en tanto provee una descripción intuitiva de una acción o una denotación, y también un intérprete, un compilador, o un herramienta teórica destinada a desentrañar principios básicos de diseño.

Primera aproximación

Maneras diferentes de dar significado a un lenguaje:

- *informal, intuitiva:*

explica el funcionamiento de los programas a través de frases comprensibles en el lenguaje natural (ejemplo: manuales, documentación tipo javadoc)

Primera aproximación

Maneras diferentes de dar significado a un lenguaje:

- *informal, intuitiva*
- *axiomática:*

explica el sentido dinámico de manera implícita, estableciendo en el marco de una lógica qué propiedades son asignables a una determinada frase del lenguaje, estableciendo así una manera de razonar sobre programas (ejemplo: $\{P/v \leftarrow e\} v := e \{P\}$)

Primera aproximación

Maneras diferentes de dar significado a un lenguaje:

- *informal, intuitiva*
- *axiomática*
- *operacional:*

explica el sentido dinámico de manera explícita, diciendo de que manera se ejecuta un programa (ejemplo: intérprete)

Primera aproximación

Maneras diferentes de dar significado a un lenguaje:

- *informal, intuitiva*
- *axiomática*
- *operacional*
- *denotacional:*

Asigna a cada programa un significado estático en un universo semántico especialmente definido para representar los fenómenos que el lenguaje describe .

Primera aproximación

Maneras diferentes de dar significado a un lenguaje:

- *informal, intuitiva*
- *axiomática*
- *operacional*
- *denotacional:*

Asigna a cada programa un significado estático en un universo semántico especialmente definido para representar los fenómenos que el lenguaje describe . Este universo podría ser un universo matemático con estructura, o el universo de las frases de otro lenguaje, por ejemplo un compilador.

Semántica Denotacional

Requiere definir:

- Dominio semántico
- Función semántica

L = conjunto de frases del lenguaje

D = dominio semántico formado por objetos abstractos

$[[_]]$ = función semántica

$$L \xrightarrow{[[_]]} D$$

$$e \longrightarrow [[e]]$$

Dominio semántico, función semántica

BIN = secuencias finitas de ceros y unos

$$BIN = \{\alpha_0 \dots \alpha_{n-1} : n \geq 1 \text{ y } \alpha_i \in \{0, 1\}\}$$

\mathbf{D} = los naturales con el cero

$$\llbracket _ \rrbracket : BIN \rightarrow \mathbf{D}$$

$$\llbracket \alpha_0 \dots \alpha_{n-1} \rrbracket = \sum_{i=1}^n \alpha_{i-1} 2^{n-i}$$

Problemas que requieren nuevas herramientas

- El problema de dar la sintaxis del lenguaje, es decir una determinación del conjunto de frases que serán consideradas programas válidos
- El problema de la buena definición del dominio semántico y la función semántica
- El manejo de variables, la ligadura y los problemas de captura
- El problema del significado de la recursión

Gramáticas

Para el estudio de los lenguajes en general, las gramáticas convencionales no son una herramienta adecuada porque sufren del problema de la *ambigüedad*.

$$\begin{aligned}\langle \text{intexp} \rangle & ::= 0 \mid 1 \mid 2 \mid \dots \\ & \quad - \langle \text{intexp} \rangle \mid \\ & \quad \langle \text{intexp} \rangle + \langle \text{intexp} \rangle \mid \\ & \quad \langle \text{intexp} \rangle * \langle \text{intexp} \rangle\end{aligned}$$

Observemos ejemplos de frases generadas por esta gramática:

-15

-15+3

2*-4

Ambigüedad

La gramática es ambigua en el siguiente sentido: algunas frases admiten diferentes maneras de generarse. Tal es el caso de $2+3*4$:

(1)

$\langle intexp \rangle$
→ $\langle intexp \rangle * \langle intexp \rangle$
→ $\langle intexp \rangle * 4$
→ $\langle intexp \rangle + \langle intexp \rangle * 4$
→ $\langle intexp \rangle + 3 * 4$
→ $2 + 3 * 4$

(2)

$\langle intexp \rangle$
→ $\langle intexp \rangle + \langle intexp \rangle$
→ $2 + \langle intexp \rangle$
→ $2 + \langle intexp \rangle * \langle intexp \rangle$
→ $2 + 3 * \langle intexp \rangle$
→ $2 + 3 * 4$

Solución que proveen los lenguajes usuales

Paréntesis y convenciones de precedencia: la gramática dice cómo se escriben concretamente las frases del lenguaje.

$$\begin{aligned}\langle \textit{intexp} \rangle & ::= \langle \textit{intexp} \rangle + \langle \textit{termexp} \rangle | \langle \textit{termexp} \rangle \\ \langle \textit{termexp} \rangle & ::= \langle \textit{groundexp} \rangle | - \langle \textit{termexp} \rangle \\ \langle \textit{groundexp} \rangle & ::= 0 | 1 | 2 | \dots | (\langle \textit{intexp} \rangle)\end{aligned}$$

Aquí queda claro que el + asocia a izquierda y que el menos tiene mayor precedencia. Podríamos llamarla **gramática concreta**, y a las frases que genera, **frases concretas**. Podríamos decir que define la **sintaxis concreta** del lenguaje.

La gramática concreta no es apropiada

- La gramática concreta resulta más complicada que la que dimos anteriormente.
- Nos obliga a fijar detalles de la sintaxis del lenguaje que son irrelevantes para nuestros propósitos, y que de hecho cada lenguaje de programación lo resuelve de una manera distinta

Nos interesa describir sólo la estructura de las frases, decir cuáles son las construcciones que hay en el lenguaje.

No interesa cómo se escriben las frases sino qué frases hay.

Gramáticas abstractas

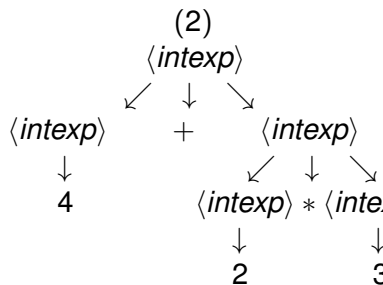
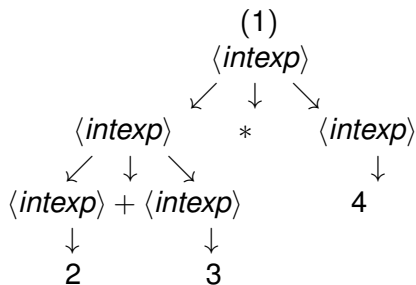
Utilizamos la forma de definición de las gramáticas convencionales:

$$\begin{aligned} \langle \text{intexp} \rangle ::= & 0 \mid 1 \mid 2 \mid \dots \\ & - \langle \text{intexp} \rangle \mid \\ & \langle \text{intexp} \rangle + \langle \text{intexp} \rangle \mid \\ & \langle \text{intexp} \rangle * \langle \text{intexp} \rangle \end{aligned}$$

Para nosotros será una gramática abstracta, que expresa qué construcciones tiene el lenguaje, cuál es la estructura de las frases que hay, y cuáles son las subfrases.

Estructura de las frases abstractas

Estructuras abstractas diferentes se ponen de manifiesto al representar cada derivación a través de un *árbol sintáctico*:



Trabajaremos con gramáticas abstractas

Problema: ¿qué notación que utilizaremos para referirnos a una producción particular de tal gramática?

En las ocasiones en que nos resulte necesario precisar notación específica resolveremos los problemas de ambigüedad mencionados estableciendo precedencias, "stopping symbols" y utilizando los paréntesis que sean necesarios.

Tales símbolos **no serán parte de la gramática** en cuestión, sino sólo convenciones de notación que nos permitirán entender de qué producción de la gramática abstracta estamos hablando.

Semántica denotacional

Es necesario recurrir a herramientas matemáticas para la construcción de los dominios semánticos y para la definición misma de la función.

- La semántica denotacional debe ser una herramienta que aporte claridad conceptual, y no que sea un instrumento de traducir algo poco comprensible en otra cosa incomprensible,
- la semántica denotacional debe ser útil para comparar otras posibles semánticas,
- se deben poder utilizar resultados típicos de la teoría de lenguajes para estudiar las características del lenguaje en cuestión.

Función semántica y Ecuaciones semánticas

Función semántica:

$$[[_]] : \langle \text{intexp} \rangle \rightarrow \mathbf{Z}$$

Ecuaciones semánticas:

$$[[0]] = 0$$

$$[[1]] = 1$$

$$\vdots$$

$$[[- e]] = - [[e]]$$

$$[[e + e']] = [[e]] + [[e']]$$

$$[[e * e']] = [[e]] * [[e']]$$

Aspectos destacables de la definición

- Los mismos símbolos aparecen con dos sentidos distintos. El 0 de la derecha es una frase del lenguaje de las expresiones enteras, y el de la izquierda es el objeto cero perteneciente a los números enteros. Esta aparente circularidad tiene un nombre: **metacircularidad**.
- Cada una de las tres últimas ecuaciones representa, en realidad, infinitas ecuaciones. Por ejemplo, $\llbracket -e \rrbracket = -\llbracket e \rrbracket$ establece una propiedad que vale cualquiera sea la expresión e . Aquí e NO es una expresión, sólo es un objeto que representa cualquier expresión del lenguaje. Es una variable del **metalenguaje**, o sea una **metavariable**.

Dirección por sintaxis

Las ecuaciones semánticas, ¿definen una función?

Un conjunto de ecuaciones es *dirigido por sintaxis* cuando se satisfacen las siguientes condiciones:

- hay 1 ecuación por cada producción de la gramática abstracta
- cada ecuación que expresa el significado de una frase compuesta, lo hace puramente en función de los significados de sus subfrases inmediatas

Composicionalidad

La dirección por sintaxis garantiza existencia y unicidad del significado, pero además garantiza la propiedad de **composicionalidad de una función**.

Se dice que una semántica es **composicional**, cuando el significado de una frase no depende de ninguna propiedad de sus subfrases, salvo de sus significados.

Efectos de la composicionalidad: podemos reemplazar una subfrase e_0 de e por otra de igual significado que e_0 sin alterar el significado de la frase e .

Los conceptos de variable y ligadura

Gramática abstracta para el lenguaje de los predicados

$\langle \text{intexp} \rangle ::= 0 \mid 1 \mid 2 \mid \dots$	$\langle \text{assert} \rangle ::= \mathbf{true} \mid \mathbf{false}$
$\langle \text{var} \rangle$	$\langle \text{intexp} \rangle = \langle \text{intexp} \rangle$
$-\langle \text{intexp} \rangle$	$\langle \text{intexp} \rangle < \langle \text{intexp} \rangle$
$\langle \text{intexp} \rangle + \langle \text{intexp} \rangle$	$\langle \text{intexp} \rangle \leq \langle \text{intexp} \rangle$
$\langle \text{intexp} \rangle * \langle \text{intexp} \rangle$	$\langle \text{intexp} \rangle > \langle \text{intexp} \rangle$
$\langle \text{intexp} \rangle - \langle \text{intexp} \rangle$	$\langle \text{intexp} \rangle \geq \langle \text{intexp} \rangle$
$\langle \text{intexp} \rangle / \langle \text{intexp} \rangle$	$\neg \langle \text{assert} \rangle$
$\langle \text{intexp} \rangle \% \langle \text{intexp} \rangle$	$\langle \text{assert} \rangle \vee \langle \text{assert} \rangle$
	$\langle \text{assert} \rangle \wedge \langle \text{assert} \rangle$
	$\exists \langle \text{var} \rangle . \langle \text{assert} \rangle$
	$\forall \langle \text{var} \rangle . \langle \text{assert} \rangle$

Particularidades de la sintaxis

- Disponemos ahora de 3 tipos de frases, las cuales por supuesto tendrán 3 dominios distintos de significado.
- $\langle var \rangle$ carece de producciones. No especificaremos qué es el conjunto de frases $\langle var \rangle$, pero asumimos que es un conjunto infinito de símbolos. Usualmente utilizaremos las metavARIABLES u, v, w, z para referirnos a elementos de $\langle var \rangle$. (Este es otro aspecto que contribuye a la condición de "abstracta" de la gramática.)
- Aparecen construcciones (para todo, existe) que poseen subfrases de distinto tipo.

Convenciones para la representación de la estructura abstracta de las frases:

Utilizaremos las convenciones usuales de precedencia de operadores aritméticos, sumadas a la convención de que el alcance de un cuantificador \forall o \exists se extiende hasta el final de la frase, o hasta la aparición de un paréntesis que cierra cuyo alcance contiene al cuantificador.

Estados y Función semántica

Σ será el conjunto de todos los estados posibles. Definimos el *conjunto de estados posibles* como la familia de todas las funciones totales de $\langle var \rangle$ en \mathbf{Z} :

$$\Sigma = \langle var \rangle \rightarrow \mathbf{Z}$$

Función Semántica

$$[[_]]^{intexp} : \langle intexp \rangle \rightarrow \Sigma \rightarrow \mathbf{Z}$$

$$[[e]]^{intexp} : \Sigma \rightarrow \mathbf{Z} \quad (\text{para cualquier expresión } e)$$

$$[[e]_{\sigma}]^{intexp} \in \mathbf{Z} \quad (\text{para cualquier expresión } e)$$

Operaciones auxiliares sobre los estados

Sea $\sigma \in \Sigma$ un estado, v una variable y n un entero.

$[\sigma|v : n]$ será el estado que coincide con σ en todas las variables salvo posiblemente en v , donde este nuevo estado tiene asignado n .

$$[\sigma|v : n]w = \begin{cases} n & \text{si } (w = v) \\ \sigma w & \text{sino} \end{cases}$$

Función semántica de las exp. booleanas

$$\llbracket _ \rrbracket^{assert} : \langle assert \rangle \rightarrow \Sigma \rightarrow \{V, F\}$$

$$\llbracket p \rrbracket^{assert} : \Sigma \rightarrow \{V, F\} \quad (\text{para cualquier predicado } p)$$

Notar que son dos funciones, una para expresiones enteras $\llbracket _ \rrbracket^{intexp}$ y otra para predicados $\llbracket _ \rrbracket^{assert}$, pero para no recargar la notación usaremos indistintamente la notación $\llbracket _ \rrbracket$ para ambas.

Ecuaciones semánticas de las exp. booleanas

$$\llbracket \text{true} \rrbracket_{\sigma} = V$$

$$\llbracket \text{false} \rrbracket_{\sigma} = F$$

$$\llbracket e = e' \rrbracket_{\sigma} = (\llbracket e \rrbracket_{\sigma} = \llbracket e' \rrbracket_{\sigma})$$

$$\vdots$$

$$\llbracket p \wedge q \rrbracket_{\sigma} = (\llbracket p \rrbracket_{\sigma} \wedge \llbracket q \rrbracket_{\sigma})$$

$$\vdots$$

$$\llbracket \forall v. p \rrbracket_{\sigma} = (\forall n \in \mathbf{Z}. \llbracket p \rrbracket_{\sigma|v:n})$$

$$\llbracket \exists v. p \rrbracket_{\sigma} = (\exists n \in \mathbf{Z}. \llbracket p \rrbracket_{\sigma|v:n})$$

¿Cuál es el significado de $1/0$?

Acá debemos cuidarnos de no trasladar al lenguaje los vicios del metalenguaje, por ejemplo los problemas relacionados con la división por 0.

En esta primera aproximación a la semántica denotacional vamos a asumir que todas las funciones son totales. En particular, la división está siempre definida (como función del dominio semántico), incluso si el divisor es 0.

Más adelante abordaremos la manera en que los lenguajes gestionan situaciones de error o no definición.

Variables y metavariables

Noción de metavariable: es una variable del metalenguaje.

Ejemplos e y e' (corren sobre expresiones enteras), p y q (corren sobre los predicados), n (corre sobre enteros) y σ (corre sobre estados), v, u, w (corre sobre las variables del lenguaje)

Usaremos u, v, w para metavariables y x, y, z para variables del lenguaje

Ligadura

- **Ocurrencia ligadora:** una ocurrencia ligadora de una variable es la que se encuentra inmediatamente después de un cuantificador (\forall o \exists).
- **Alcance de una ocurrencia ligadora:** En $\forall v. p$ o $\exists v. p$, el predicado p es el alcance de la ocurrencia ligadora de v .
- **Ocurrencia ligada:** cualquier ocurrencia de v en el alcance de una ocurrencia ligadora de v es una ocurrencia ligada de v (se toma la de menor alcance).
- **Ocurrencia libre:** una ocurrencia de una variable que no es ligadora ni ligada es una ocurrencia libre.
- **Variable libre:** una variable que tiene ocurrencias libres es una variable libre.
- **Expresión cerrada:** no tiene variables libres.

Variables libres

$$FV [n] = \emptyset$$

$$FV v = \{v\}$$

$$FV(-e) = FV e$$

$$FV(e + e') = FV e \cup FV e'$$

⋮

$$FV \text{ true} = \emptyset$$

$$FV \text{ false} = \emptyset$$

$$FV(e = e') = FV e \cup FV e'$$

⋮

$$FV(\neg p) = FV p$$

$$FV(p \wedge q) = FV p \cup FV q$$

⋮

$$FV(\forall v. p) = (FV p) - \{v\}$$

$$FV(\exists v. p) = (FV p) - \{v\}$$

Operador Sustitución

Las expresiones con variables libres pueden instanciarse sustituyendo sus variables libres por términos.

Conjunto de sustituciones:

$$\Delta = \langle var \rangle \rightarrow \langle intexp \rangle$$

El operador sustitución "opera" sobre expresiones enteras (términos) y expresiones booleanas (predicados).

$$_/_ \in \langle intexp \rangle \times \Delta \rightarrow \langle intexp \rangle$$

$$_/_ \in \langle assert \rangle \times \Delta \rightarrow \langle assert \rangle$$

Operador Sustitución

$$0/\delta = 0\dots$$

$$v/\delta = \delta v$$

$$(-e)/\delta = -(e/\delta)$$

$$(e + f)/\delta = (e/\delta) + (f/\delta)\dots$$

$$(e = f)/\delta = (e/\delta) = (f/\delta)$$

$$(\neg p)/\delta = \neg(p/\delta)\dots$$

$$(p \vee q)/\delta = (p/\delta) \vee (q/\delta)\dots$$

El problema de la captura

Aplicación "ingenua" de la sustitución:

si $e = \delta y$ entonces

$$(\exists x. x > y) / \delta = \exists x. x > e$$

Pero si e es x , quedaría $\exists x. x > x$ que es falsa, a diferencia de $\exists x. x > y$ que es verdadera.

El problema de la captura

Visto de otra manera:

$$\exists x. x > y$$

$$\exists z. z > y$$

sólo se diferencian en el nombre de la variable ligada, de manera que al sustituir en uno u otro me debería dar resultados equivalentes. Pero para $\delta y = x + 1$ obtenemos predicados no equivalentes.

Problema: se está "capturando" x , que era libre, y ahora pasa a ser ligada.

Solución al problema de la captura

Renombramos la variable ligadora y ligada antes de efectuar una susutitución que producirá una captura.

Renombramos eligiendo una variable "nueva", que llamaremos v_{new} . Esto es, una variable no "capturable". Las variables "capturables" son las que pueden aparecer al aplicar la sustitución.

$$\begin{aligned}(\forall v . b)/\delta &= \forall v_{new} . (b/[\delta|v : v_{new}]) \\(\exists v . b)/\delta &= \exists v_{new} . (b/[\delta|v : v_{new}])\end{aligned}$$

donde

$$v_{new} \notin \bigcup_{w \in FV(b) - \{v\}} FV(\delta w)$$

Propiedades de la semántica

Hay dos propiedades de la semántica que resultan relevantes para los lenguajes de programación.

Teorema de Coincidencia (TC)

Si dos estados σ y σ' coinciden en las variables libres de p , entonces da lo mismo evaluar p en σ o σ' . En símbolos:

$$(\forall w \in FV(p) . \sigma w = \sigma' w) \implies \llbracket p \rrbracket \sigma = \llbracket p \rrbracket \sigma'.$$

Teorema de Renombre (TR)

Los nombres de las variables ligadas no tienen importancia. En símbolos,

$$u \notin FV(q) - \{v\} \implies \llbracket \forall u . q/v \rightarrow u \rrbracket = \llbracket \forall v . q \rrbracket$$

Teorema de Sustitución

Si aplico la sustitución δ a p y luego evalúo en el estado σ , puedo obtener el mismo resultado a partir de p sin sustituir, si evalúo en un estado que hace el trabajo de δ y de σ (en las variables libres de p).

En símbolos:

$$(\forall w \in FV(p) . \llbracket \delta w \rrbracket \sigma = \sigma' w) \implies \llbracket p/\delta \rrbracket \sigma = \llbracket p \rrbracket \sigma'.$$

