

Compiladores

Facultad de Ciencias UNAM

Representaciones Intermedias

Lourdes Del Carmen González Huesca *

22 de abril de 2024

Resumen

En este punto del proceso de compilación se ha realizado el análisis léxico, que devolvió un *stream* de *tokens*, un análisis sintáctico que definió si el código era correcto de acuerdo a la gramática y un análisis semántico que determinó si el código era semánticamente válido bajo ciertas características. Por tanto se han cubierto todas las fases del *front-end* de un compilador. Para pasar a las fases siguientes es necesario transformar el código fuente en uno que facilite su manejo. A lo largo de esta nota se presentan algunas representaciones intermedias que puede tener el código y los tipos de éstas para cubrir diferentes objetivos.

Términos clave:

Es: Generación de código intermedio, Código de tres direcciones, Gráficas de dependencias de atributos, Gráficas de flujo.

En: Intermediate-code generation, Three-address code, Directed Acyclic Graphs for Expressions, Flow graphs.

Después de las fases de análisis en el front-end, el compilador ha obtenido suficiente información del programa fuente para realizar una traducción que facilite el manejo del programa en las fases de back-end. Esta información será utilizada para tener una **representación intermedia** (IR *Intermediate Representation*) que sea fiel al código fuente y se traduzca a código de bajo nivel.

Una representación intermedia es una estructura de datos que sólo existe en tiempo de compilación. Dependiendo del diseño e implementación del compilador es posible que se usen diferentes representaciones intermedias para mejorar las traducciones entre fases así como utilizar optimizaciones a la representación obtenida por el front-end para que el programa traducido y equivalente al fuente se ejecute de forma eficiente.

Las representaciones intermedias tienen como objetivos principales:

- simplificar el tratamiento de código y separar las fases de front-end y back-end;
- mantener el compilador modularizado para mejor implementación y manutención;
- facilitar optimizaciones independientes de la máquina.

Y tienen las siguientes propiedades:

- hacer más fácil la generación y manipulación de código intermedio;
- proveer una mejor abstracción del programa fuente.

*Material revisado por el servicio social de Apoyo a la Docencia y Asesoría Académica en la Facultad de Ciencias de la UNAM con clave 2023-12/12-292, de Diana Laura Nicolás Pavia.

Existen tres categorías de representaciones intermedias:

Estructurales o gráficas

son representaciones orientadas a gráficas utilizadas en su mayoría por traductores entre lenguajes fuente.

Lineales

representaciones que están dirigidas a máquinas abstractas y que son simples y compactas. El código compilado es representado por una secuencia ordenada de operaciones.

Híbridas

estas son combinación de las anteriores para obtener los beneficios de ambas.

Código de tres direcciones

Esta representación intermedia o lenguaje intermedio divide expresiones en instrucciones más sencillas y cuya sintaxis es más parecida a las instrucciones de bajo nivel. Cada instrucción incluye asignaciones y operaciones básicas, contiene a lo más 4 items: tres operandos y un operador. Las asignaciones se realizan en direcciones relativas que serán más fáciles de traducir a direcciones físicas en una fase del back-end llamada asignación de registros. Las instrucciones siguen un orden lineal (estructurado) pero puede haber saltos (**goto**), bifurcaciones, condicionales, llamadas a funciones o subrutinas, accesos a memoria y asignaciones, etc.

Por ejemplo, consideremos un condicional sencillo y su traducción a código de tres direcciones:

```

if (x + y*z > x*y +z)      t1 = y*z
    a = 0;                   t2 = x+t1
                              t3 = x*y
                              t4 = t3+z
                              if (t2 <= t4) goto L
                              a = 0
                              L:

```

Este código secciona las operaciones del condicional en instrucciones sencillas usando direcciones temporales (**t1**) para almacenar los resultados intermedios. La etiqueta **L** indica la siguiente instrucción al término del condicional y permite continuar con la ejecución del resto del programa. Esta representación puede ser obtenida a partir de la gramática del lenguaje, sólo es cuestión de extender las reglas semánticas mediante acciones de traducción.

Ejemplo 1. Definición dirigida por la sintaxis para construir código de tres direcciones para expresiones que representan asignaciones, los atributos a usar son **code** para el código y **addr** para guardar la dirección que almacenará el valor de una expresión:

PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(id.lexeme) \neq E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr \neq E_1.addr \neq E_2.addr)$
$ - E_1$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel$ $gen(E.addr \neq \text{'minus'} E_1.addr)$
$ (E_1)$	$E.addr = E_1.addr$ $E.code = E_1.code$
$ \text{id}$	$E.addr = top.get(id.lexeme)$ $E.code = ''$

Ejemplo 2. La siguiente traducción obtiene una representación intermedia de tres direcciones en un arreglo de referencias. El arreglo se define mediante la gramática

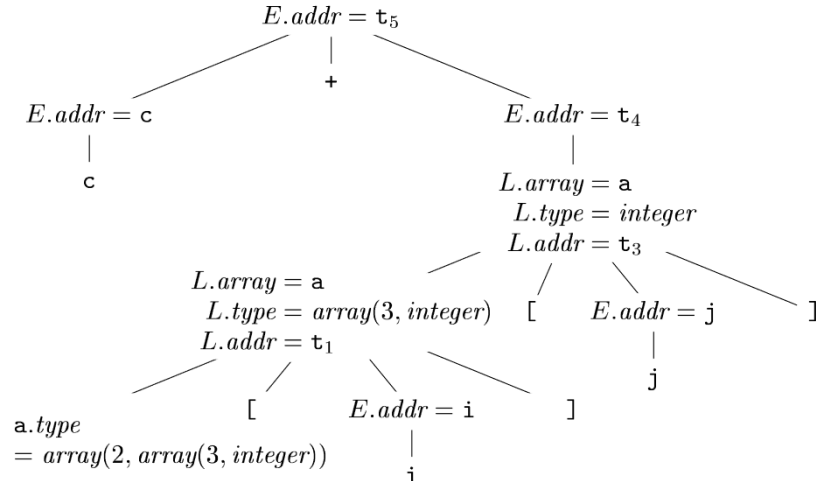
$$L \rightarrow L[E] \mid \mathbf{id}[E]$$

que extiende a $S \rightarrow \mathbf{id} = E \mid L = E \quad E \rightarrow E + E \mid \mathbf{id} \mid L$.

Los atributos usados son: **addr** para almacenar temporalmente el desplazo en el arreglo, **array** es el apuntador a la tabla de símbolos donde se almacena el arreglo y **type** es el tipo del subarreglo correspondiente:

$$\begin{aligned}
S &\rightarrow \mathbf{id} = E ; \quad \{ \text{gen}(\text{top.get}(\mathbf{id.lexeme}) \text{'=' } E.\text{addr}); \} \\
&\mid L = E ; \quad \{ \text{gen}(L.\text{array.base} \text{'[' } L.\text{addr} \text{'}]' \text{'=' } E.\text{addr}); \} \\
E &\rightarrow E_1 + E_2 \quad \{ E.\text{addr} = \mathbf{new Temp}(); \\
&\quad \text{gen}(E.\text{addr} \text{'=' } E_1.\text{addr} \text{'+' } E_2.\text{addr}); \} \\
&\mid \mathbf{id} \quad \{ E.\text{addr} = \text{top.get}(\mathbf{id.lexeme}); \} \\
&\mid L \quad \{ E.\text{addr} = \mathbf{new Temp}(); \\
&\quad \text{gen}(E.\text{addr} \text{'=' } L.\text{array.base} \text{'[' } L.\text{addr} \text{'}]'); \} \\
L &\rightarrow \mathbf{id} [E] \quad \{ L.\text{array} = \text{top.get}(\mathbf{id.lexeme}); \\
&\quad L.\text{type} = L.\text{array.type.elem}; \\
&\quad L.\text{addr} = \mathbf{new Temp}(); \\
&\quad \text{gen}(L.\text{addr} \text{'=' } E.\text{addr} \text{'*' } L.\text{type.width}); \} \\
&\mid L_1 [E] \quad \{ L.\text{array} = L_1.\text{array}; \\
&\quad L.\text{type} = L_1.\text{type.elem}; \\
&\quad t = \mathbf{new Temp}(); \\
&\quad L.\text{addr} = \mathbf{new Temp}(); \\
&\quad \text{gen}(t \text{'=' } E.\text{addr} \text{'*' } L.\text{type.width}); \\
&\quad \text{gen}(L.\text{addr} \text{'=' } L_1.\text{addr} \text{'+' } t); \}
\end{aligned}$$

Y el siguiente árbol es un parse tree decorado para la expresión $\mathbf{c} + \mathbf{a}[\mathbf{i}][\mathbf{j}]$ donde **a** es un arreglo de enteros de dimensión 3×2 de tamaño 24 considerando que un entero es de tamaño 4:



El código de tres direcciones que representa esa misma expresión es:

$$\begin{aligned} t_1 &= i * 12 \\ t_2 &= j * 4 \\ t_3 &= t_1 + t_2 \\ t_4 &= a[t_3] \\ t_5 &= c + t_4 \end{aligned}$$

Ejemplo 3. Three-Address Code ¹

PRODUCTION	SEMANTIC RULES
$P \rightarrow S$	$S.next = newlabel()$ $P.code = S.code \parallel label(S.next)$
$S \rightarrow \text{assign}$	$S.code = \text{assign}.code$
$S \rightarrow \text{if} (B) S_1$	$B.true = newlabel()$ $B.false = S_1.next = S.next$ $S.code = B.code \parallel label(B.true) \parallel S_1.code$
$S \rightarrow \text{if} (B) S_1 \text{ else } S_2$	$B.true = newlabel()$ $B.false = newlabel()$ $S_1.next = S_2.next = S.next$ $S.code = B.code$ $\parallel label(B.true) \parallel S_1.code$ $\parallel gen('goto' S.next)$ $\parallel label(B.false) \parallel S_2.code$
$S \rightarrow \text{while} (B) S_1$	$begin = newlabel()$ $B.true = newlabel()$ $B.false = S.next$ $S_1.next = begin$ $S.code = label(begin) \parallel B.code$ $\parallel label(B.true) \parallel S_1.code$ $\parallel gen('goto' begin)$
$S \rightarrow S_1 S_2$	$S_1.next = newlabel()$ $S_2.next = S.next$ $S.code = S_1.code \parallel label(S_1.next) \parallel S_2.code$

PRODUCTION	SEMANTIC RULES
$B \rightarrow B_1 \parallel B_2$	$B_1.true = B.true$ $B_1.false = newlabel()$ $B_2.true = B.true$ $B_2.false = B.false$ $B.code = B_1.code \parallel label(B_1.false) \parallel B_2.code$
$B \rightarrow B_1 \&\& B_2$	$B_1.true = newlabel()$ $B_1.false = B.false$ $B_2.true = B.true$ $B_2.false = B.false$ $B.code = B_1.code \parallel label(B_1.true) \parallel B_2.code$
$B \rightarrow ! B_1$	$B_1.true = B.false$ $B_1.false = B.true$ $B.code = B_1.code$
$B \rightarrow E_1 \text{ rel } E_2$	$B.code = E_1.code \parallel E_2.code$ $\parallel gen('if' E_1.addr \text{ rel.op } E_2.addr 'goto' B.true)$ $\parallel gen('goto' B.false)$
$B \rightarrow \text{true}$	$B.code = gen('goto' B.true)$
$B \rightarrow \text{false}$	$B.code = gen('goto' B.false)$

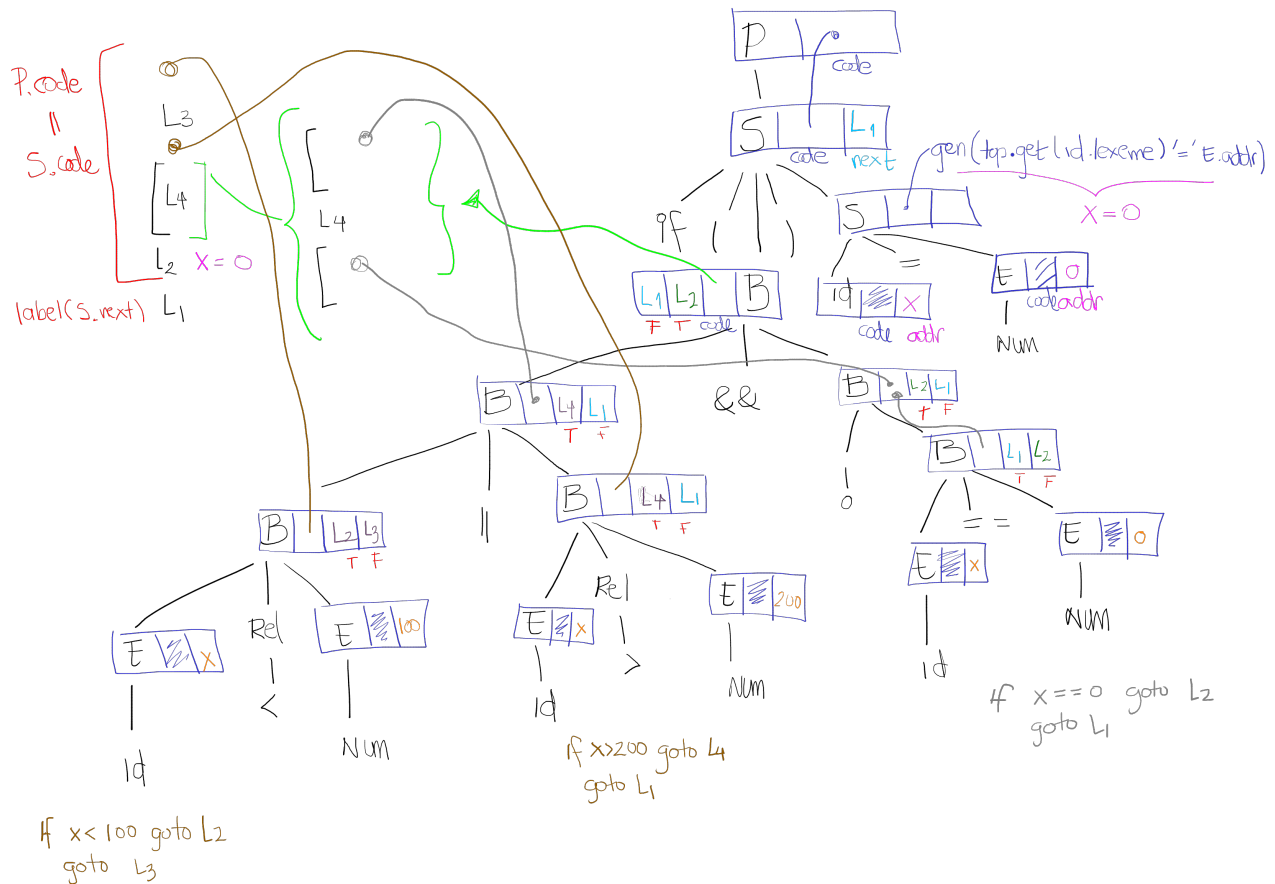
PRODUCTION	SEMANTIC RULES
$S \rightarrow \text{id} = E ;$	$S.code = E.code \parallel$ $gen(top.get(\text{id.lexeme}) ' = ' E.addr)$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $gen(E.addr ' = ' E_1.addr ' + ' E_2.addr)$
$ - E_1$	$E.addr = \text{new Temp}()$ $E.code = E_1.code \parallel$ $gen(E.addr ' = ' \text{minus} ' E_1.addr)$
$ (E_1)$	$E.addr = E_1.addr$ $E.code = E_1.code$
$ \text{id}$	$E.addr = top.get(\text{id.lexeme})$ $E.code = ''$

¹Figuras y ejemplo tomado del libro “Compilers, Principles, Techniques and Tools”, Aho A. et al., Capítulo 6.

Código fuente y código de tres direcciones usando las funciones semánticas:

```
if ( x<100 || x>200 && x!=0 ) x=0;
```

```
if x < 100 goto L2
goto L3
L3: if x > 200 goto L4
goto L1
L4: if x != y goto L2
goto L1
L2: x = 0
L1:
```



Gráficas de dependencias de atributos

Representación intermedia que modela el flujo de la información entre las instancias de los atributos en un parse-tree. Es una gráfica (alternativa) donde los nodos son los diferentes atributos asociados a cada símbolo y las aristas dirigidas representan las restricciones de las reglas semánticas en el cálculo de los atributos.

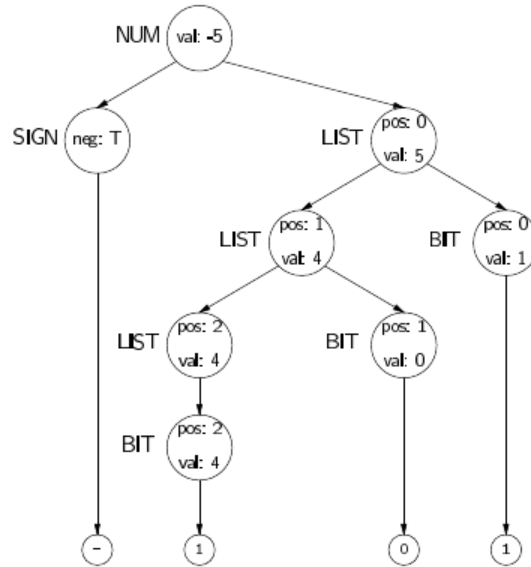
Ejemplo 4. Considera la siguiente gramática para representar números binarios con signo

$$\text{num} \rightarrow \text{sign list} \quad \text{sign} \rightarrow + \mid - \quad \text{list} \rightarrow \text{bit} \mid \text{list bit} \quad \text{bit} \rightarrow 0 \mid 1$$

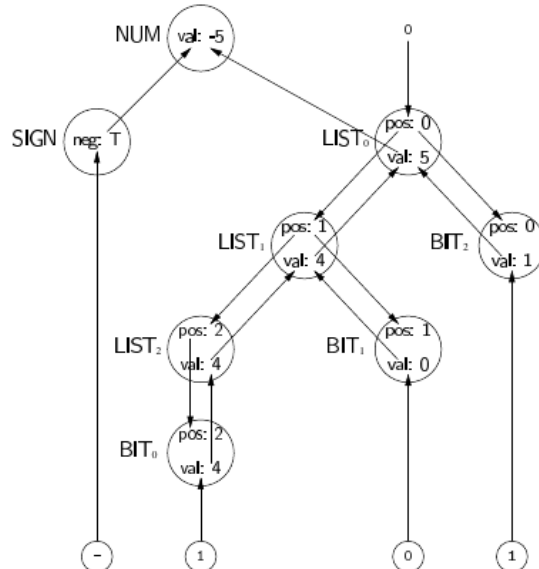
La gramática extendida con reglas semánticas que obtienen el valor decimal que representa el número binario es la siguiente:

PRODUCTION	SEMANTIC RULES
NUM \rightarrow SIGN LIST	LIST.pos := 0 if SIGN.neg NUM.val := -LIST.val else NUM.val := LIST.val
SIGN \rightarrow +	SIGN.neg := false
SIGN \rightarrow -	SIGN.neg := true
LIST \rightarrow BIT	BIT.pos := LIST.pos LIST.val := BIT.val
LIST \rightarrow LIST ₁ BIT	LIST ₁ .pos := LIST.pos + 1 BIT.pos := LIST.pos LIST.val := LIST ₁ .val + BIT.val
BIT \rightarrow 0	BIT.val := 0
BIT \rightarrow 1	BIT.val := 2 ^{BIT.pos}

El siguiente árbol se obtiene al procesar la cadena -101 usando la gramática anterior:



Después de obtener el parse-tree, se pueden (re)visitar los nodos para calcular los atributos:



El primer atributo en calcularse es el signo seguido del atributo `pos` en los nodos `LIST0`, `LIST1` y `LIST2` en ese orden. Con estos atributos se puede calcular el atributo `pos` para los nodos `BIT0`, `BIT1` y `BIT2`. Después se calcula el valor de los bits y de las listas en el siguiente orden: `BIT0`, `LIST2`, `BIT1`, `LIST1`, `BIT2`, `LIST0`. Este orden está definido por la regla semántica `LIST.val:=LIST2.val+BIT.val`. Finalmente se calcula el valor decimal en el nodo raíz.

Gráficas de control de flujo

Representación que modela la transferencia de control en el programa. Es una gráfica donde los nodos son bloques básicos y las aristas dirigidas representan el flujo del control de las estructuras de control (ciclos, casos, saltos, etc.).

Esta es una representación que respeta el significado original del código fuente, pero que será fácil de traducir al lenguaje máquina.

Los nodos de la gráfica de flujo son nodos básicos. Hay un arista del bloque B al bloque C si y sólo si es posible que la primera instrucción del bloque C se siga a partir de la última instrucción del bloque B. Lo que significa que existen dos formas en que un arista puede ser justificado:

1. Existe un salto condicional o incondicional desde la última instrucción de B hacia la primera de C.
2. La primera instrucción de C se sigue inmediatamente a la última de B en el orden original del código de tres direcciones y B no termina en un salto incondicional.

Referencias

- [1] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers, Principles, Techniques and Tools*. Pearson Education Inc., Second edition, 2007.
- [2] Torben Ægidius Mogensen. *Basics of Compiler Design*. Lulu Press, 2010.
- [3] Hanne Riis Nielson and Flemming Nielson. *Semantics with Applications: An Appetizer (Undergraduate Topics in Computer Science)*. Springer-Verlag, Berlin, Heidelberg, 2007.
- [4] Frank Pfenning. Notas del curso (15-411) Compiler Design. <https://www.cs.cmu.edu/~fp/courses/15411-f14/>, 2014.
- [5] Michael Lee Scott. *Programming Language Pragmatics*. Morgan-Kaufman Publishers, Third edition, 2009.
- [6] Yunlin Su and Song Y. Yan. *Principles of Compilers, A New Approach to Compilers Including the Algebraic Method*. Springer-Verlag, Berlin Heidelberg, 2011.
- [7] Linda Torczon and Keith Cooper. *Engineering A Compiler*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2nd edition, 2011.