

Compilador para lenguaje COOL

Isabella Sierra

Grupo C412

ISABELLAMARIA.SIERRA@ESTUDIANTES.MATCOM.UH.CU

Adrian Tubal Páez Ruiz

Grupo C412

A.PAEZ@ESTUDIANTES.MATCOM.UH.CU

Eric Martín García

Grupo C411

E.MARTIN@ESTUDIANTES.MATCOM.UH.CU

Tutor(es):

Alejandro Piad, *Universidad de La Habana*

1. Lexer y Parser

Para resolver la sintáctica del compilador se utilizó el módulo de *Python* **PLY** que implementa las populares herramientas de construcción de compiladores *lex* y *yacc*, que utilizamos para generar el primer **AST** del proceso de compilación.

1.1 Tokenización con lex

1. **Creación de los tokens:** Una instancia de objeto **LexToken** (llamemos a este objeto **t**) tiene los siguientes atributos:

- **t.type** que es el tipo de token (por ejemplo: 'INT', 'PLUS', etc.).
- **t.value** que es el lexema
- **t.lineno** que es el número de línea actual

2. **Creación de las reglas:** Las reglas de expresiones regulares pueden definirse como un **string** o como una función. En cualquier caso, el nombre de la variable tiene el prefijo **t_** para denotar que es una regla para hacer coincidir *tokens*.

Para tokens simples, la expresión regular se especificó como **string**: $t_PLUS = r'\+'$

Para tokens mas complejos se definieron funciones:

Algorithm 1 Ejemplo 1

```
def t_STRING(t):
    r"""
    t.lexer.string_start = t.lexer.lexpos
    t.lexer.begin('string')
```

3. **Creación del lexer:** `lexer = lex.lex()`

1.2 Parsing con yacc

Una vez terminado el proceso de tokenización se realiza el proceso de parsing utilizando la herramienta *yacc*. Luego de la creación de una gramática libre del contexto *yacc* generará un parser **LALR**.

La gramática que se utilizó para el proceso de parsing es la siguiente:

```
program -> class_list

class_list -> def_class ; class_list
| def_class ;

def_class -> CLASS TYPE { feature_list }
| CLASS TYPE INHERITS TYPE { feature_list }

feature_list -> def_attr ; feature_list
| def_func ; feature_list
| empty

def_attr -> ID : TYPE <- expr
| ID : TYPE

def_func -> ID ( params ) : TYPE { expr }

params -> param_list

params -> empty

param_list -> param , param_list
| param empty

param -> ID : TYPE

expr -> LET let_attrs IN expr
| CASE expr OF case_list ESAC
| IF expr THEN expr ELSE expr FI
| WHILE expr LOOP expr POOL

expr -> ID <- expr

expr -> expr @ TYPE . ID ( arg_list )
| expr . ID ( arg_list )
| ID ( arg_list )

expr -> expr + expr
| expr - expr
| expr * expr
```

```

| expr \ expr
| expr < expr
| expr <= expr
| expr = expr

expr -> NOT expr
| ISVOID expr
| LNOT expr

expr -> ( expr )

expr -> atom

let_attrs -> def_attr , let_attrs
| def_attr

case_list -> case_elem ; case_list
| case_elem ;

case_elem -> ID : TYPE => expr

arg_list -> arg_list_ne
| empty

arg_list_ne -> expr , arg_list_ne
| expr

atom -> INT
atom -> ID
atom -> NEW TYPE
atom -> block
atom -> BOOL
atom -> STRING

block -> { block_list }

block_list -> expr ; block_list
| expr ;

```

Cada regla de la gramática será definida por una función donde el **string** de documentación de esa función contiene la especificación apropiada. Las declaraciones que conforman el cuerpo de la función implementan las acciones semánticas de la regla junto con las definiciones de los nodos del **AST**.

Algorithm 2 Ejemplo 2

```

def p_program(p):
    """program : class_list"""
    p[0] = ProgramNode(p[1])

```

Al finalizar el proceso de *parsing*, estarán todas las condiciones creadas para iniciar el chequeo semántico.

2. Chequeo semántico

Para realizar el chequeo semántico se utilizó un **visitor** (implementado en `src/semantic`) de una pasada que realiza las siguientes acciones sobre un nodo:

- Aplica el **visitor** correspondiente a las expresiones hijas involucradas y de éste se obtiene el tipo estático inferido de cada una.
- Se aplican las reglas de chequeo de tipos definidas en cada uno de los **visitors** de las expresiones para detectar errores.
- Se infiere el tipo estático de la expresión representada por el nodo y se retorna para ser utilizado por sus padres.

2.1 CoolType

Para auxiliarnos en el chequeo semántico definimos una clase **CoolType** en `src/cool_types` que envuelve la significancia de un tipo en *COOL*.

Contiene los atributos:

- **name**: Representa el nombre del tipo.
- **inherits**: Determina si la clase es heredera de alguna clase definida en el programa.
- **parent**: Es una instancia **CoolType** que representa el padre de la clase definida.
- **methods**: Contiene todos los métodos definidos por la clase (De tipo **CoolTypeMethod**).
- **attributes**: Contiene los atributos definidos por la clase (De tipo **CoolTypeAttribute**).
- **childs**: Contiene todas los **CoolType** que heredan directamente de la clase.
- **order_interval**: Atributo que define un intervalo (x, y) tal que si $A \leq B$, entonces (x_A, y_A) estará contenido en (x_B, y_B) . Más adelante explicaremos para qué se utiliza este atributo.

y los métodos:

- **get_method**: Obtiene el método determinado por el argumento **id**, teniendo en cuenta si pudo haber sido definido por la clase guardada en **parent**.
- **get_all_methods**: Obtiene todos los métodos definidos por la clase, así como los definidos por alguna clase superior en su línea de herencia.
- **add_method**: Agrega el método con nombre **id** en el conjunto de métodos de la clase. Esta función tiene implícitas las reglas de definición de métodos.
- **add_attr**: Agrega el atributo con nombre **id** en el conjunto de métodos de la clase. Esta función tiene implícitas las reglas de definición de atributos.
- **get_all_attributes**: Obtiene todos los atributos, heredados y propios.
- **get_self_attributes**: Obtiene todos los atributos propios.

Tipos *built-in*: Los tipos *built-in* `ObjectType`, `IObjectType`, `StringType`, `IntType`, `BoolType` son declarados en el archivo `src/cool_types/types.py`, así como sus métodos de instancia.

3. Generación de código

3.1 Código Intermedio: CIL

Para facilitar el paso desde la sintaxis de *COOL* hacia MIPS, nos servimos de CIL como lenguaje intermedio.

Para generar un AST de CIL utilizamos el patrón `visitor` sobre los nodos del AST de *COOL*. La selección del visitor adecuado se realiza a partir de un diccionario de la forma `{type:visitor}`.

3.1.1 AST

Un programa CIL se representa con un `ProgramNode` y consta de cuatro partes:

1. Una sección `.data`, representada por una lista de `DataNode`, que a su vez consiste en la declaración de una constante involucrada en el programa.
2. Una sección `.types`, representada por una lista de `TypeNode`, que a su vez consiste en la declaración de un tipo con los identificadores de sus métodos y sus atributos.
3. Una sección `.text`, representada por una lista de `DefFuncNode`, que a su vez consiste en la declaración de una función.

DataNode: Un nodo de tipo `.data` consiste en una constante de tipo *string* y una etiqueta que lo mapea:

```
data_1: "Hello world\n"
```

TypeNode: Un nodo de tipo `.type` consiste en una etiqueta que representa el nombre del tipo y un conjunto de *features* que pueden ser atributos con su tipo o métodos. Un método será mapeado al método real a ejecutar de ser llamado éste. Se representa de la forma:

```
Type:
  attr1:   Type_1
  method_1: actual_method_1
  method_2: actual_method_2
```

DefFuncNode: Un nodo de tipo `deffunc` constituye la declaración de un cuerpo de función, dentro de ella se ejecutará un bloque de código CIL, siguiendo las convenciones usuales.

3.1.2 CILBLOCK

Para representar un bloque abstracto de código CIL utilizamos una clase llamada `CILBlock` que contiene un `body`, un `locals` y un `value`.

Para convertir un AST de *COOL* a un AST de CIL visitaremos cada uno de los nodos del primer AST y generaremos un `CILBlock` correspondiente, donde cada uno de los elementos del `body` consiste en un nodo del AST de CIL.

3.1.3 EJEMPLO: ASIGNACIÓN

Sea el código en *COOL* `c<-a+b`, parseado a un árbol de la forma `AssignNode(VarNode(a),...expression...)` y con chequeo semántico perfecto.

El nodo más externo creará un `CILBlock` de la forma:

```
body = concat(expression.body, [
    CILAssignNode(a, expression.value)
])
value = expression.value

locals = concat(expression.locals, [
    a
])
```

3.2 Código MIPS

Para generar código MIPS utilizamos el mismo patrón `visitor` sobre los nodos del AST de CIL generado, para entonces generar un AST de código MIPS de forma similar a la generación de código intermedio.

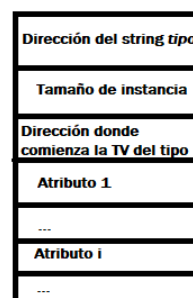
3.2.1 UTILIZACIÓN DE REGISTROS

Los registros se utilizaron de la siguiente forma:

- `$fp` para referenciar el inicio de la memoria correspondiente a la función actual.
- `$sp` para referenciar la posición donde se encuentra el tope de la pila.
- `$t0-t9` para realizar operaciones intermedias.
- `$ra` para guardar la dirección de retorno después de la ejecución de una función.
- `$v0` para guardar el valor de retorno de una función.
- `$a0-$a3` para guardar los argumentos necesarios en la ejecución de los `syscall`

3.2.2 UTILIZACIÓN DE MEMORIA

Representación de instancia: Una instancia se representa en memoria de la siguiente forma:



Paso de argumentos: En el paso de argumentos utilizamos la pila, liberando la memoria utilizada al terminar la ejecución de la función.

Creación de una instancia: Para la creación de una instancia reservamos memoria utilizando `sbrk`. Cada una de las instancias tendrá como atributos especiales su tipo y su tamaño, en la primera y última posición respectivamente; por tanto, al reservar memoria tenemos que hacerlo con un tamaño igual a la cantidad de atributos más dos.

Localización de variables: Para la localización de variables utilizamos la pila, y el acceso a cada una se controla a partir de un *offset* que las representa dentro de `$fp` para cada función. Este control se realiza a través de un diccionario en tiempo de compilación.

3.2.3 LLAMADO A FUNCIÓN

Tabla Virtual: Al inicio de la ejecución del programa se construye una Tabla Virtual de la siguiente forma:

1. Se reserva un espacio para cada tipo en la sección *data*.
2. Se guarda, para la función *X*, la dirección de la etiqueta de la función *Y* adecuada, definida desde la generación de código intermedio. El *offset* respecto al inicio de la tabla virtual se prefija en compilación de manera tal que si *A* implementa *X* y *B* lo redefine, entonces *X* se encuentra en el mismo *offset* en ambas *Tablas Virtuales*.

Decisión de la función a llamar: Para decidir qué función se debe llamar, basta localizar el *offset* predefinido en compilación de la función, a partir de la dirección de *Tabla Virtual* del tipo de la instancia, guardada como atributo en ésta.

Pasos del llamado a función:

1. Se salvan los registros utilizados: `$fp`, `$ra`. (son los únicos que trascienden un llamado de función).
2. Se pasan los argumentos a la pila.
3. Se realiza la instrucción `Jal function_name` en el caso de saber la función precisa a llamar; de lo contrario, se utiliza la *Tabla Virtual* y el salto se realiza utilizando la instrucción `Jalr $ti`, donde `$ti` guarda la dirección de la función correcta. En cualquier caso que guarda la dirección de retorno en el registro `$ra`.
4. Se sacan los argumentos de la pila.
5. Se restauran los registros.
6. Se toma el valor de retorno de `$v0`.

La función llamada debe:

1. Mover `$fp` a donde se encuentra `$sp`.
2. Tomar sus argumentos de la pila.
3. Reservar memoria para todas sus variables locales.
4. Realizar las instrucciones que le corresponden.

5. Liberar el espacio de la pila correspondiente a todas las variables locales.
6. Realizar un salto al contenido del registro `$ra`.

3.2.4 ORDEN DE HERENCIA EN MIPS PARA INSTRUCCIÓN CASE

La relación de herencia se determina a través de unos atributos llamados `order` y `min_order`, creado desde la generación de código intermedio como atributo especial. Si *A* hereda de *B*, entonces seguramente $\text{min_order}(B) < \text{min_order}(A) < \text{order}(A) < \text{min_order}(B)$.

De esta forma garantizamos que en la ejecución de una instrucción *case*, podamos decidirnos por la primera rama cuyo intervalo definido por $(\text{min_order}(x), \text{order}(x))$ esté contenido en el mismo intervalo para el tipo de la expresión. Para garantizar que la primera sea la correcta, las ramas del *case* se ordenan en tiempo de compilación por niveles de herencia.

3.2.5 CONSTRUCTOR DE INSTANCIA EN MIPS

Para cada uno de los tipos, definimos una función especial llamada `__init__`, que se crea en la fase de generación de código intermedio y se encarga de manejar la inicialización de los atributos heredados y los propios.

3.2.6 MANEJO DE strings EN MEMORIA

- Un *string* constante se reserva en la sección *.data* en el inicio del programa, la estructura que contiene dichos *strings* está construida desde la fase de generación de código intermedio.
- El espacio para guardar un *string* que introduce el usuario se reserva utilizando el `syscall` definido por `sbrk` con un tamaño de *buffer* de 1024 bytes.
- El espacio para un *string* resultante de una operación de cadena se reserva con el `syscall` definido por `sbrk` con el tamaño exacto necesario, que se conoce en ejecución a través de los argumentos del llamado.

3.2.7 DETERMINACIÓN DEL TIPO EN EJECUCIÓN

Dentro de cada instancia, como representamos anteriormente, se guarda un atributo especial que consiste en la dirección en memoria de la tabla virtual del tipo de ésta. Para acceder a esta dirección construimos una instrucción apropiada desde la fase de código intermedio, que se traduce en MIPS al acceso directo a la dirección en memoria predefinida donde sabemos que se encuentra este atributo.