



Apunte de cátedra: Optimizando la Máquina Virtual

En el apunte anterior definimos la BVM, una máquina virtual de pila, con las siguientes reglas de ejecución:

$$\begin{array}{l}
 \langle \text{CONST } n; c \mid e \mid s \rangle \longrightarrow \langle c \mid e \mid n : s \rangle \\
 \langle \text{ADD}; c \mid e \mid n : m : s \rangle \longrightarrow \langle c \mid e \mid m + n : s \rangle \\
 \langle \text{SUB}; c \mid e \mid n : m : s \rangle \longrightarrow \langle c \mid e \mid m - n : s \rangle \\
 \langle \text{ACCESS } i; c \mid e \mid s \rangle \longrightarrow \langle c \mid e \mid e!i : s \rangle \\
 \langle \text{CALL}; c \mid e \mid v : (e_f, c_f) : s \rangle \longrightarrow \langle c_f \mid v : e_f \mid (e, c)_{RA} : s \rangle \\
 \langle \text{FUNCTION}(c_f); c \mid e \mid s \rangle \longrightarrow \langle c \mid e \mid (e, c_f) : s \rangle \\
 \langle \text{RETURN}; - \mid - \mid v : (e, c)_{RA} : s \rangle \longrightarrow \langle c \mid e \mid v : s \rangle \\
 \langle \text{PRINTN}; c \mid e \mid n : s \rangle \xrightarrow{n} \langle c \mid e \mid n : s \rangle \\
 \langle \text{PRINT}; x_0; \dots x_n; \text{NULL}; c \mid e \mid s \rangle \xrightarrow{x_0 \dots x_n} \langle c \mid e \mid s \rangle \\
 \langle \text{FIXPOINT}(c_f); c \mid e \mid s \rangle \longrightarrow \langle c \mid e \mid (e_{\text{fix}}, c_f) : s \rangle \\
 \langle \text{SHIFT}; c \mid e \mid v : s \rangle \longrightarrow \langle c \mid v : e \mid s \rangle \\
 \langle \text{DROP}; c \mid v : e \mid s \rangle \longrightarrow \langle c \mid e \mid s \rangle
 \end{array}$$

Estas reglas deben ser completadas con la implementación de IFZ pedida en el apunte anterior.

1. La BVM Implementada en C

Como prometimos en el apunte anterior, la máquina virtual puede implementarse en un lenguaje de más bajo nivel donde cada instrucción tiene un costo muy bajo. El archivo `bvm.c` en el repo contiene la implementación, y usa el mismo formato y opcodes que la versión en Haskell, por lo cual deberían ser equivalentes.

El código está documentado como para ser leído, pero exponemos las ideas principales de la máquina aquí. En la función `run` hay 3 variables: `c`, `e` y `s`, punteros a código, entorno y pila respectivamente. El entorno `e` es una lista enlazada, inicialmente vacía. La pila está representada por un arreglo `stack`, y un puntero `s` que apunta a la primera posición libre (es decir, una posición adelante de la cima). Por ello, el pushear un valor `v` a la pila se hace simplemente con `*s++ = v`, y quitar un valor se hace con `v = *--s`. La máquina interpreta la instrucción apuntada por el puntero de código, es decir el valor `*c`, incrementa el puntero `c` y luego salta al código que implementa esa instrucción (via el `switch`).

La implementación de la instrucción `ADD`, por ejemplo, es entonces esencialmente:

```
case ADD: {
    uint32_t y = (*--s).i;
    uint32_t x = (*--s).i;
    (*s++).i = x+y;
    break;
}
```

En la implementación en `bvm.c`, la pila no es de enteros, sino de `union value` una unión (no etiquetada) de los distintos tipos de valores que soporta la máquina.

```
union value {
    int i;
    struct clo clo;
};
typedef union value value;
```

Es decir, también incluye clausuras y direcciones de retorno que definimos con la estructura `clo`.¹

```
struct clo {
    env clo_env;
    code clo_body;
};
```

Pregunta. ¿Por qué no usamos un arreglo para los entornos?

2. Llamadas (y Recursión) de Cola

Consideremos el siguiente fragmento de código FD4.

```
let g (x:Nat) : Nat = ...
let f (x:Nat) : Nat = g (1 + x)
  in f 0
```

Analizamos qué ocurre al evaluar la aplicación `f 0`. Esta aplicación nos llevará a ejecutar una aplicación `g 1`, donde tendremos en la pila el valor del argumento (1), una clausura para `g`, y la dirección de retorno de `f`. Luego de llamada, `f` retorna el resultado con **RETURN**. La traza es algo como:

$$\begin{array}{llll}
& \langle \text{CALL; RETURN} \mid e \mid 1 : (e_g, c_g) : (e', c)_{RA} : s \rangle & (e_0) \\
\longrightarrow & \langle c_g \mid 1 : e_g \mid (e, c_{RET})_{RA} : (e', c)_{RA} : s \rangle & (e_1) \\
\longrightarrow^* & \langle \text{RETURN} \mid 1 : e_g \mid v : (e, c_{RET})_{RA} : (e', c)_{RA} : s \rangle & (e_2) \\
\longrightarrow & \langle \text{RETURN} \mid e \mid v : (e', c)_{RA} : s \rangle & (e_3) \\
\longrightarrow & \langle c \mid e' \mid v : s \rangle & (e_4)
\end{array}$$

donde c_{RET} apunta al **RETURN** y v es el resultado de evaluar g 1.

En el estado (e_0) , cuando `f` está por llamar a `g`, tiene su propia dirección de retorno guardada en la pila, después de los dos elementos necesarios para efectuar el **CALL**. Al hacer la llamada, en (e_1) , agrega una nueva dirección de retorno para el **RETURN** de `f`, y salta a `g`. La función `g` hace alguna cantidad de pasos, y llega a un **RETURN** (estado (e_2)) habiendo agregado exactamente un elemento a la pila. Finalmente, en (e_3) , `g` retorna a `f` y en (e_4) `f` retorna a su llamante. Resumiendo, `f` queda como “intermediario” de la llamada a `g`, guardando una dirección de retorno a su **RETURN**.

En este caso, cuando una función directamente retorna el resultado de una llamada *sin operar sobre él*, decimos que es una llamada *de cola*. En otras palabras, la llamada a `g` desde `f` es de cola porque es lo *último* que hace `f`. En estos casos, podemos optimizar la ejecución mediante la *omisión* de la dirección de retorno intermedia, y pasando la propia de la función llamante en su lugar. Para ello, introducimos una nueva instrucción **TAILCALL** a tal efecto. Si `f` fuera compilada con una **TAILCALL** a `g`, la traza sería algo como:

$$\begin{array}{llll}
& \langle \text{TAILCALL} \mid e \mid 1 : (e_g, c_g) : (e', c)_{RA} : s \rangle \\
\longrightarrow & \langle c_g \mid 1 : e_g \mid (e', c)_{RA} : s \rangle \\
\longrightarrow^* & \langle \text{RETURN} \mid 1 : e_g \mid v : (e', c)_{RA} : s \rangle \\
\longrightarrow & \langle c \mid e' \mid v : s \rangle
\end{array}$$

Con lo que ahorramos un elemento en la pila y un paso. Aunque a simple vista esto no parece ser una gran optimización, la verdadera ganancia se manifiesta en las funciones recursivas de cola.

Las funciones *recursivas de cola* son aquellas funciones en las cuales sus llamadas recursivas ocurren sólo en posición de cola. Por ejemplo:

```
let rec suma (m n : Nat) : Nat =
  ifz n then m
  else suma (1 + m) (n - 1)
```

Aquí, la optimización de llamadas de cola reduce *asintóticamente* la memoria consumida por esta función, de $O(n)$ a $O(1)$. Notar la diferencia con la versión de `suma` vista anteriormente, que no es recursiva de cola. Ser consciente de la recursión de cola y poder escribir código en esa forma es una habilidad esencial para los programadores funcionales, aunque también es importante en otros lenguajes².

¹Para las direcciones de retorno, usamos el caso `clo`. Si bien existe una diferencia conceptual entre clausuras y direcciones de retorno, es totalmente inútil crear un nuevo caso idéntico. El caso `clo` bien podría estar nombrado `clo_or_ra`.

²Por ejemplo, en C, la recursión de cola puede implementarse como un salto al inicio de la función.

2.1. Compilando Llamadas de Cola

Por suerte, tenemos una definición muy concreta de lo que es una llamada de cola: un `CALL` seguido de un `RETURN`. La idea para mejorar la compilación es dividirla en dos funciones $\mathcal{C}(-)$ y $\mathcal{T}(-)$, mutuamente recursivas. La primera es la compilación de un término general a bytecode, mientras que la segunda es usada para la compilación de un término *en posición de cola*. La compilación en posición de cola $\mathcal{T}(-)$ tiene una regla especial para una aplicación, generando una llamada de cola. En todo otro caso se interpreta solamente como $\mathcal{C}(-)$, pero agregando un `RETURN` al final.

$$\begin{aligned}\mathcal{T}(ab) &= \mathcal{C}(a); \mathcal{C}(b); \text{TAILCALL} \\ \mathcal{T}(t) &= \mathcal{C}(t); \text{RETURN} && (\text{en otro caso}) \\ \mathcal{C}(\lambda t) &= \text{FUNCTION}(\mathcal{T}(t)) \\ \mathcal{C}(ab) &= \mathcal{C}(a); \mathcal{C}(b); \text{CALL} \\ &\dots \text{otras reglas de } \mathcal{C}(-) \dots\end{aligned}$$

Con estas reglas, vemos que la compilación (con $\mathcal{C}(-)$) de un término $\lambda.fe$ resulta:

$$\text{FUNCTION}(\mathcal{C}(f); \mathcal{C}(e); \text{TAILCALL})$$

como es de esperar.

Pregunta. Una idea para implementar esta optimización es tomar el bytecode compilado de la forma original y reemplazar las ocurrencias de `[CALL; RETURN]` por `[TAILCALL]`. ¿Qué problemas puede haber con este enfoque?

Hay que notar que la mayoría de las llamadas recursivas ocurren bajo un `ifz`, por lo que esta optimización como está escrita no surte mucho efecto. Pero podemos notar que si un `ifz` está en posición de cola, entonces sus ramas también lo están. En otras palabras, si el `ifz` es lo último a evaluar antes de retornar, entonces sus ramas también tienen esa propiedad (cuando son ejecutadas). Entonces, podemos agregar una regla como la siguiente a $\mathcal{T}(-)$:

$$\mathcal{T}(\text{ifz } c \text{ then } t \text{ else } e) = \mathcal{C}(c); \text{IFZ } n; \mathcal{T}(t); \mathcal{T}(e)$$

donde la instrucción `IFZ` chequea si el valor de la pila es 0 o no. En el caso de que sea 0, no hace nada. Esto causará que se ejecute el código $\mathcal{T}(t)$ y se retorne antes de llegar al código de $\mathcal{T}(e)$. En el caso de que el valor de la pila no sea 0, la instrucción `IFZ` salta n posiciones. Tomando $n = |\mathcal{T}(t)|$ conseguimos que salte a la rama del `else`.

Por otro lado, hay otra instrucción que podemos obviar antes de un `RETURN`: la instrucción `DROP`. Esta instrucción solo afecta el entorno actual de la máquina, que se ve descartado al ejecutar `RETURN`. Como `DROP` solo se genera al final de un `letbinding`, podemos optimizar los `lets` en posición de cola agregando otro caso a la función $\mathcal{T}(-)$.

$$\mathcal{T}(\text{let } x = m \text{ in } n) = \mathcal{C}(m); \text{SHIFT}; \mathcal{T}(n)$$

De esta forma, evitamos el `DROP`.

$$\mathcal{C}(\lambda. \text{let } f = g \text{ in } fx) = \text{FUNCTION}(\mathcal{C}(g); \text{SHIFT}; \text{ACCESS } 0; \mathcal{C}(x); \text{TAILCALL})$$

2.2. Tareas

- Completar la ejecución de `IFZ` en la máquina.
- Implementar `TAILCALL` en la BVM, tanto en la implementación en Haskell y como en la implementación en C.
- Cambiar la compilación a bytecode para que use llamadas de cola en lo posible, siguiendo §2. Intente verificar la mejora en uso de memoria para alguna función apropiada.