

HULK — Havana University Language Kompiler

Diego A. Martínez Jiménez

Introducción

HULK es un lenguaje de programación imperativo, funcional, estática y fuertemente tipado. Casi todas las instrucciones en HULK son expresiones. En particular, el subconjunto de HULK implementado se compone solamente de expresiones que pueden escribirse en una línea.

Un Intérprete

A la hora de interpretar un lenguaje el código fuente pasa por varias fases. A continuación una breve descripción de cada una.

+ Análisis léxico

El análisis léxico constituye la primera fase, aquí se lee el programa fuente de izquierda a derecha y se agrupa en componentes léxicos (tokens), que son secuencias de caracteres que tienen un significado. Además, todos los espacios en blanco, líneas en blanco, y demás información innecesaria se elimina del programa fuente. También se comprueba que los símbolos del lenguaje (palabras clave, operadores, etc.) se han escrito correctamente.

+ Análisis sintáctico/Parser

En esta fase los caracteres o componentes léxicos se agrupan jerárquicamente en frases gramaticales que el compilador utiliza para sintetizar la salida. Se comprueba si lo obtenido de la fase anterior es sintácticamente correcto (obedece a la gramática del lenguaje). Por lo general, las frases gramaticales del programa fuente se representan mediante un árbol de análisis sintáctico.

+ Análisis semántico

La fase de análisis semántico revisa el programa fuente para tratar de encontrar errores semánticos y reúne la información sobre los tipos para la fase posterior de generación de código. En ella se utiliza la estructura jerárquica determinada por la fase de análisis sintáctico para identificar los operadores y operandos de expresiones y proposiciones.

Un componente importante del análisis semántico es la verificación de tipos. Aquí, el compilador verifica si cada operador tiene operandos permitidos por la especificación del lenguaje fuente, entre otros.

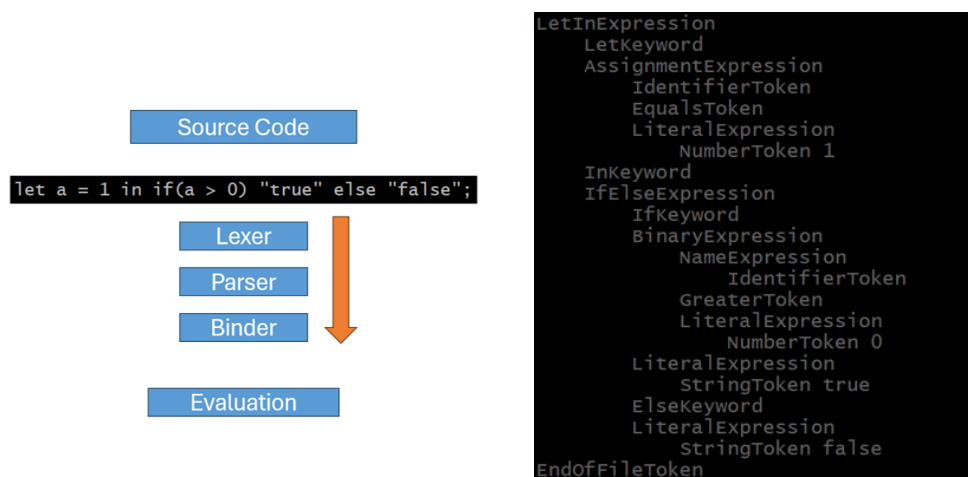
+ Evaluación

En la fase de evaluación, el intérprete ya tiene toda la información necesaria para saber de qué forma evaluar cada expresión. Por lo que en esta fase ya ejecutará el código fuente de la forma esperada o, en caso de haber algún error en el código fuente, los reportará.

Intérprete de HULK

+ Análisis del código

El código fuente se analiza como se explicó anteriormente. Para esto existen tres clases principales **Lexer**, **Parser** y **Binder**. Estas clases se encargan del análisis léxico, sintáctico y semántico respectivamente. A continuación una breve descripción de su funcionamiento.



Estructura del intérprete y AST

+ Clase Lexer

Esta clase se encarga de tokenizar el código. Primero almacena el código en la variable estática luego con la función `Lex()` retorna el próximo `SyntaxToken`, siendo `SyntaxToken` la clase padre que representa todos los tokens. En esta clase se usan algunos métodos auxiliares de la clase `SyntaxFacts`; por ejemplo `SyntaxFacts.GetText(SyntaxKind kind)` y `SyntaxFacts.GetKeywordKind(string text)`, mas adelante se explicará esta clase. Haciendo uso de la clase `DiagnosticsBag` se reportan los errores léxicos que se encuentren a la hora de analizar la cadena.

+ Clase Parser

La clase Parser se encarga de crear el AST. El constructor de esta clase recibe el código, luego hace uso de la clase Lexer para almacenar todos los tokens ignorando los espacios en blanco y los tokens desconocidos; en el constructor también son almacenados los errores encontrados en el análisis léxico para luego unirlos a los que errores que se encuentren en esta parte del análisis. Una vez construida una instancia se hace uso del método público `ParseCompilationUnit()` que se encarga de iniciar el proceso de análisis sintáctico. En este proceso primero se analiza si se la expresión es la declaración de una función, en caso de que no se analiza si es una expresión binaria; es decir que sea una expresión tipo *left + right* por ejemplo, para lo cual se analiza la parte *left* la cual será una expresión primaria (expresión primaria se refiere a las expresiones que deben parsearse antes de una expresión binaria o una expresión unaria) o una expresión unaria, en el caso de una expresión unaria se analiza la precedencia del operador de la expresión unaria y a partir de la precedencia se parseará antes o después de la expresión binaria (en principio todos los operadores unarios implementados en el proyecto tienen la misma precedencia y también precedencia mayor que todos los operadores binarios). Luego de analizar *left*, se analiza el operador binario en cuanto a su precedencia, luego analiza *right* de similar forma a *left* y así se obtiene la expresión binaria y en caso de no serlo *left* será el tipo de expresión que luego habrá que se tendrá que luego analizar semánticamente.

+ Clase Binder

Esta clase se encarga del análisis semántico del AST. El método `BindExpression()` de esta clase se encarga de crear una `BoundExpression` la cual será la expresión con toda la información necesaria para evaluar el código correctamente y haciendo la menor cantidad de verificaciones. Este método recibirá el root del AST y se irá analizando recursivamente. Las clases `BoundUnaryOperator` y `BoundBinaryOperator` son esenciales para el asegurar que los operadores de una expresión unaria o binaria estén siendo correctamente usados; con el método `Bind()` cada una de las clases asegura la existencia de el operador con los tipos especificados. De forma similar se verifican las funciones predefinidas mediante la clase `BoundPredefinedFunction`. Las expresiones de asignación infieren el tipo a partir de lo que se encuentra a la derecha del `=`.

+ Clase Evaluator

Luego del análisis semántico se evalúa el código mediante la clase `Evaluator`. Similar a la clase `Binder` esta clase analiza recursivamente el AST pero ahora en vez de analizar la semántica evalúa las expresiones, hasta tener un literal que será el resultado de la expresión evaluada.

*

Para este proyecto se uso como guía una **playlist de YouTube** acerca de como crear un compilador.