



Assignment of master's thesis

Title: SOM Virtual Machine Implementation
Student: Bc. Rudolf Rovňák
Supervisor: Ing. Petr Máj
Study program: Informatics
Branch / specialization: System Programming
Department: Department of Theoretical Computer Science
Validity: until the end of summer semester 2022/2023

Instructions

Familiarize yourself with the details of Simple Object Machine (SOM) and design and implement SOM virtual machine in the C++ language. Prefer clarity of design and code to raw performance and compare your results to the existing SOM++ implementation to determine the overhead.



**FACULTY
OF INFORMATION
TECHNOLOGY
CTU IN PRAGUE**

Master's thesis

SOM Virtual Machine Implementation

Bc. Rudolf Rovňák

Department of Theoretical Computer Science

Supervisor: Ing. Petr Máj

May 6, 2021

Acknowledgements

I would like to thank my supervisor for helping me patiently with this work.

Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No.121/2000 Coll., the Copyright Act, as amended, in particular that the Czech Technical University in Prague has the right to conclude a license agreement on the utilization of this thesis as a school work under the provisions of Article 60 (1) of the Act.

In Prague on May 6, 2021

.....

Czech Technical University in Prague
Faculty of Information Technology
© 2021 Rudolf Rovňák. All rights reserved.

This thesis is school work as defined by Copyright Act of the Czech Republic. It has been submitted at Czech Technical University in Prague, Faculty of Information Technology. The thesis is protected by the Copyright Act and its usage without author's permission is prohibited (with exceptions defined by the Copyright Act).

Citation of this thesis

Rovňák, Rudolf. *SOM Virtual Machine Implementation*. Master's thesis. Czech Technical University in Prague, Faculty of Information Technology, 2021.

Abstrakt

Táto práca popisuje implementáciu virtuálneho stroja pre programovací jazyk Simple Object Machine, založenom na Smalltalku. Takisto obsahuje analýzu existujúcich riešení a analýzu vlastného riešenia. Práca prezentuje syntaktickú analýzu, bajtkód s procesom kompilácie doň a prostredie pre spustenie programov napísaných v programovacom jazyku SOM.

Kľúčová slova behové prostredie, virtuálny stroj, objektovo orientované programovanie

Abstract

This work presents an implementation of a virtual machine for Simple Object Machine programming language, based on Smalltalk. Additionally, it analyzes existing implementations, along with an analysis of the provided solution. Included is a design and implementation of a parser, bytecode with a compiler and runtime environment with garbage collector to allow executing programs written in SOM.

Keywords runtime, virtual machine, object oriented programming

Contents

1 Introduction	1
2 Analysis and design	3
2.1 Existing implementations	3
2.2 Class definition	3
2.2.1 Variables	4
2.2.1.1 Variable name scoping	5
2.2.2 Literals	5
2.3 Primitives	6
2.4 Methods and messages	7
2.5 Blocks	9
2.5.1 Non-local return and block scoping	9
2.6 Expressions	15
2.7 Control structures	16
2.7.1 Conditional branching	16
2.7.2 For loops	17
2.7.3 While loops	17
2.7.4 Class hierarchy	18
2.8 Bytecode	18
2.8.1 Program structure	18
2.8.2 Program entities	19
2.8.3 Instructions	19
2.9 Garbage collection	20
2.9.1 Mark and sweep	21
3 Realisation	23
3.1 Program overview	23
3.1.1 Build	23
3.2 Source code parsing	24

3.3	Abstract Syntax Tree	24
3.3.1	AST Nodes	24
3.3.1.1	Expressions	26
3.3.2	AST construction	28
3.4	Bytecode	30
3.4.1	Values	30
3.4.2	Instructions	30
3.4.3	Compilation	30
3.4.3.1	Method compilation	31
3.5	Interpretation	32
3.5.1	Garbage collection	32
3.5.2	Program counter	32
3.5.3	Execution stack	33
3.5.4	Objects	34
3.5.4.1	Object creation	35
3.5.5	Messages	35
3.6	Core library	36
3.6.1	Primitives	36
3.6.2	Strings	37
3.6.3	Booleans	37
3.6.4	Array	38
3.6.5	Nil	38
3.6.6	Blocks	38
3.7	Block evaluation	39
3.7.1	Argument handling	41
3.7.2	Block restart	41
4	Performance	43
4.1	Fibonacci sequence	43
4.2	While loop	45
4.3	Non local returns	47
Conclusion		49
Future work		49
Bibliography		51
A Acronyms		53
B Contents of enclosed CD		55

List of Figures

2.1	Railroad diagram for <code>classDefinition</code> rule.	4
2.2	Example of a simple class defined in SOM.	4
2.3	Unary message example.	7
2.4	Binary message example.	7
2.5	Array example.	8
2.6	Demonstration of message precedence.	8
2.7	Demonstration of message sends order for one message type.	8
2.8	Demonstration of mathematical operators precedence rules.	9
2.9	Example of blocks usage in SOM.	10
2.10	Example to demonstrate non local return.	11
2.11	Implementation of <code>whileTrue:</code> method in <code>Block</code> class.	11
2.12	Implementation of relevant boolean methods.	12
2.13	Created blocks and their home contexts.	13
2.14	Second example code for non local returns.	13
2.15	Output of the example code from figure 2.14	14
2.16	State of execution before return for example on figure 2.14.	14
2.17	Third example demonstrating non local returns.	15
2.18	Modification of the method from example 2.17.	15
2.19	ANTLR grammar snippet for expressions.	15
2.20	Example of messages functioning as <i>if</i> -control structures.	16
2.21	Example of simple for loops.	17
2.22	Examples of different ways of iterating over an array.	17
2.23	Example of while loops.	18
2.24	Illustration of the object hierarchy before GC algorithm run.	22
2.25	Traversal of the heap objects.	22
2.26	State of the heap after the GC run.	22
3.1	Interface of the program.	23
3.2	Example of some interesting lexer rules from SOM ANTLR grammar.	25
3.3	The syntax of class definition described by ANTR grammar.	25

3.4	The syntax of nested blocks.	26
3.5	Conceptual class diagram for SOM abstract syntax tree.	27
3.6	Example for AST construction.	28
3.7	Expression from example 3.6 as AST.	29
3.8	Conceptual class diagram for compilation environment.	31
3.9	Conceptual class diagram for interpretation environment.	33
3.10	Implementation of Nil class in core library.	38
3.11	Block implementation in SOM core library.	39
3.12	Example to demonstrate block evaluation process.	40
3.13	State of execution environment for execution of example 3.12.	40
4.1	The program for Fibonacci sequence performance test in SOM.	44
4.2	The program for Fibonacci sequence performance test in C++.	44
4.3	The program for Fibonacci sequence performance test in Python.	44
4.4	Execution times of Fibonacci sequence test.	45
4.5	Execution times of Fibonacci sequence test for various SOM im- plementations.	46
4.6	SOM source code for second performance benchmark – while loop evaluation.	46
4.7	Execution times for different SOM implementations for code from 4.6.	47
4.8	Invocation of restart method in Java implementation of SOM (taken from [1]).	47
4.9	SOM source code for third performance benchmark – non local return evaluation.	48
4.10	Execution times for different SOM implementations for code from 4.9.	48

Introduction

In the last decades, a trend of dynamic programming languages ^[1] has been on the rise. As opposed to static programming languages (usually compiled) dynamic ones offer a higher level of abstraction and claim to allow faster and less error-prone development, though validity of such claims is a subject to many discussions. Dynamic languages move a lot of actions traditionally done during compile-time to run-time. This creates the need for another layer, *a runtime environment*.

Dynamic languages are often, but not always, dynamically typed. This is just one aspect – this can be considered just an example of a process done either before program execution (static type checking) or during runtime (dynamic type checking). Other principles that are generally utilised heavily in dynamic languages are late binding, reflection and more.

The goal of this diploma thesis is to implement a process virtual machine for a programming language called SOM, or Simple Object Machine. It is a dynamic, object-oriented programming language based on Smalltalk. It was originally implemented at University of Århus in Denmark to teach object oriented VMs ^[2]. There are several implementations in various programming languages, ranging in speed, optimizations etc.

The main focus of this work will be the clarity of implementation over performance. The solution will provide a basic implementation of a traditional runtime VM that can be built upon in the future. This will include the process of parsing, compiling the bytecode and then providing a runtime environment, along with an implementation for most basic principles (flow control, basic data types, loops).

¹Not to be confused with *dynamically typed programming languages*.

Analysis and design

Simple Object Machine (SOM) is a minimal Smalltalk dialect used primarily for teaching construction of virtual machines. Key characteristics according to official website ([2]) are:

- clarity of implementation over performance,
- common language features such as: objects, classes, closures, non-local returns,
- interpreter optimizations, threading, garbage collectors are different across various implementations.

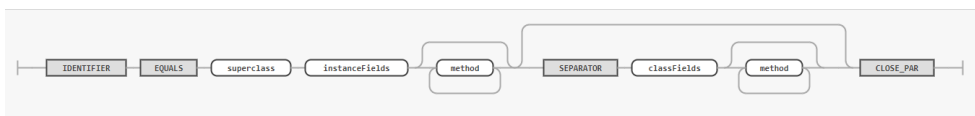
2.1 Existing implementations

There are multiple existing implementations of SOM written in different programming languages, including: Java (SOM), C (CSOM) C++ (SOM++), Python (PySOM) and many others. Some of these are bytecode based, others utilize abstract syntax tree interpretation. They range in implemented optimizations (CSOM offers no optimizations, while TruffleSOM claims to be highly optimised) [2].

Additionally, Read-Eval-Print loop implementation is available on the official site of the project. The REPL accepts only simple expressions and does not support class definitions. Standard library is available. [2]

2.2 Class definition

Syntax for class definition follows the official SOM grammar. The language supports single inheritance as apparent from the use of `subclass` token in the grammar. Not every class has explicitly specified superclass, therefore the actual identifier in the rule is optional.

Figure 2.1: Railroad diagram for `classDefinition` rule.

```
SimpleHello = (
  | name |

  setName: aString (
    name := aString
  )

  printGreeting (
    ('Hello, ', name) print
  )
)
```

Figure 2.2: Example of a simple class defined in SOM.

Declaration of instance side fields follows, denoted by vertical bars. This token itself can be empty. Instance side methods definitions are next. Further details on *methods* and *messages* in SOM are discussed in [2.4](#). Same syntax is used for class side fields and methods separated by a special token.

There are no access modifiers for class variables and methods. By default, all the instance and class variables are accessible only from the given class (or its instances) or its subclass. This corresponds to the behaviour of `protected` access modifier in C++. Every method is invocable by a corresponding message send to an object – this corresponds to `public` access modifier.

2.2.1 Variables

In Smalltalk, a variable is defined as “a *dynamically modifiable association (binding) of either a name or an index to a value. Each distinct variable has exactly one name (or index)*”[\[3\]](#).

A value of a named variable can be any object. Indexed variables are at the core also just an object. They represent an ordered sequence of objects as a single value. Example of those are arrays or Strings. Actual indices are always strictly positive (greater than zero), meaning the first element of an array corresponds to index of value one. This is standard in Smalltalk dialects, although uncommon in C-like languages. Retrieving the values belonging to an index is done via sending a message to the encapsulating object.

When creating a new variable, it is assigned a special value `nil`, meaning the variable is empty. This special value can also be explicitly assigned to a

variable at any point.

SOM is a dynamically typed programming language (as is Smalltalk). As a result, there is no syntax to indicate a data type of a variable. One thing worth pointing out is that in the context of Smalltalk, a *data type* is defined differently than most programming languages. As stated in [3], a class is not a type. A Smalltalk type is defined as “*the power set of messages to which an object can meaningfully respond*” [3]. This is the definition I will be using in the context of SOM. As a result of this, any number of SOM classes can implement one data type.

2.2.1.1 Variable name scoping

Every variable has its scope, which determines the visibility of the variable. SOM follows the rules of Smalltalk when it comes to scoping, as defined in [3]:

- **Local variables** are accessible within the method or code block in which they are defined.
- **Formal method arguments** are accessible by the method wherein they are defined.
- **Formal block arguments** are accessible by the block wherein they are defined.
- **Instance variables** are accessible within all methods of a given object. Each object has its own instances of these variables.
- **Class variables** are accessible by all objects that are instances of the class or its subclasses. All the objects share the same instance of this variable.
- **Global variables** are accessible everywhere. These are primarily special variables (such as `nil`). Users cannot define their own global variables. This functionality can be achieved using class variables.

There is some more details on variable scoping, especially when talking about block closures. Further detail is discussed in [2.5.1](#).

2.2.2 Literals

Outside of variables, there is also a need to represent fixed values in a SOM source code.

Integer literal specifies a value of a decimal whole number, positive or negative. In provided implementation, every integer literal is a representation of an object of class `Integer`. To achieve simple manipulation with integer numbers, all integers are internally represented by a C++ data type `int32_t`

– therefore 32 bit signed integer number. Overflow and underflow is not addressed in the implementation, therefore the behaviour copies that of the underlying C++ type.

Floating point literal approximates a value of a real number. Syntactically, it consists of a decimal, possibly negative, integer literal representing the non-fractional part of the number. It is followed by a decimal point and another decimal (non-negative) integer representing the fractional part of the number. Precision is implicitly given and there is no way to change it. My implementation uses double precision (as defined in C++). Edge cases (such as rounding errors) are not addressed – the behaviour copies that of standard C++ double data type.

String literal represents a sequence of characters. String literals are objects of class `String`. Syntactically, they are delimited by single quotes (`'`). To include a single quote in a string, it needs to be escaped by another single quote.

In the provided implementation of SOM, String objects are not treated as a collection. It is an encapsulated object and every String object is unmutable. Every message, that somehow uses and modifies the value of its receiver creates a new object. This behaviour corresponds with Smalltalk and other SOM implementations.

Array literals specify a sequence of values encapsulated by a single object (that is an instance of class `Array`). Syntactically, the values of an array are surrounded by parentheses and preceded by a hash sign. Note that because of the dynamic typing, elements of an array do not have to be instances of the same class.

Every Array literal results in an Array object instantiation. Arrays are mutable, unlike Strings.

2.3 Primitives

Even though SOM is purely object oriented, in order to get any actual computations done, there is a point where some virtual machine primitives must be invoked. Following things are therefore implemented as primitives (among others):

- memory allocation (`new` message),
- bitwise operations,
- integer arithmetics (`+`, `-`, `=` etc.),
- array accessing (`at:`, `at:put:`)

Primitives are implemented directly in the VM runtime, though not breaking the syntax or core principles of the language. Details on implementation are in section [3.6.1](#).

2.4 Methods and messages

As the SOM language is based on Smalltalk, the concept of messages (and the link to methods) is crucial to understand. *“The only way to invoke a method is to send a message – which necessarily involves dynamic binding (by name) of message to method at runtime (and never at compile time). The internals of an object are not externally accessible, ever – the only way to access or modify an object’s internal state is to send it a message”* [3].

Execution of an invoked method ends with the execution of the last expression in it. Every method implicitly returns `self` (a reference to the object on which the method is invoked). Explicit return of a value is done with a special token `^`. Execution of an expression preceded by this token will exit the method.

The [4] defines a helpful terminology for message passing:

- A message is composed of the message *selector* and the optional message arguments.
- Every message must be sent to its *receiver*.
- Message and its receiver together will be referred to as *message send*.

There are three types of messages (as defined in other Smalltalk dialects, Pharo as an example of one).

Unary messages are sent to an object without any additional information (argument). In the following example, a unary message `size` is sent to a string object.

```
'hello' size "Evaluates to 5"
```

Figure 2.3: Unary message example.

Binary messages are a special type of messages that require exactly one argument. The selector of a binary message can only consist of a sequence of one or more characters from the set: `+`, `-`, `*`, `/`, `&`, `=`, `<`, `>`, `—`, and `@`. A simple example of usage of binary message are arithmetic operations.

```
3 + 4 "Evaluates to 7"
```

Figure 2.4: Binary message example.

Keyword messages require one or more arguments. From the syntactic standpoint, they consist of multiple keywords, each ending in colon (`:`). When sending a message, each keyword is followed by an argument. Note, that a keyword message taking one argument is different to a binary message.

```
| numbers |
numbers := #(1 2 3 4 5). "Simple array"
"Sending a keyword message at:put: to an object of class Array"
numbers at: 1 put: 6 "numbers is now #(6 2 3 4 5)"
```

Figure 2.5: Array example.

When composing messages of various types, there are precedence rules (as defined for Pharo in [\[4\]](#)):

- Unary messages are sent first, followed by binary messages. Keyword messages are sent last.
- Messages in parentheses are sent before other messages.
- Messages of the same kind are evaluated from left to right.

These simple rules permit a very natural way of sending messages, as demonstrated on the next example. First, a simple array is created. Then, a unary message `last` is evaluated, returning the last element of the array. After that, binary message `+` is evaluated (to 2 in this example). Finally, keyword message `at:put:` is sent to an array, putting number 5 on the second position in an array.

```
| numbers |
numbers := #(1 2 3 4 5).
numbers at: 1 + 1 put: numbers last.
"numbers at: (1 + 1) put: (numbers last)"
```

Figure 2.6: Demonstration of message precedence.

Next example demonstrates sending messages from left to right when all of them are of the same type.

```
| numbers |
numbers := #(1 2 3 4 5)
numbers last asString print
"This is equivalent to the following message sends"
((numbers last) asString) print
```

Figure 2.7: Demonstration of message sends order for one message type.

There is a downfall to the simplicity of these rules. Arithmetic operations are all just a simple binary message sends, therefore to ensure proper precedence, it is necessary to use parentheses.

```
"Evaluated as (3 + 2) * 5
3 + 2 * 5
"Parentheses required to achieve mathematical
operators precedence"
3 + (2 * 5)
```

Figure 2.8: Demonstration of mathematical operators precedence rules.

2.5 Blocks

Blocks provide a mechanism to defer the execution of expressions [4]. Blocks can be treated as an object – they can be assigned to variables and passed as arguments.

Blocks can also accept parameters – they are denoted with a leading colon. Parameters are separated from the body of the block by a vertical bar. Local variables can also be declared inside a block. Those are accessible only inside the block and are initialized each time a block is evaluated.

Block is executed by sending it a message `value`. However, this is a unary message and there is no way to pass parameters to a block. To solve this problem, a keyword message `value:` is implemented. So far, this gives a user ability to pass only one parameter to a block. To mitigate this issue, there are two possibilities. The first one is to implement a keyword message for every number of parameters (for example `value:value:`, `value:value:value:`). While this approach is simple, readable and relatively easy to implement for low numbers of parameters, it is impossible for this solution to be exhaustive and the code using very long keyword messages would be bloated.

Another approach would be to implement a keyword message `value:` with an argument of array type. This would permit to use arbitrary number of arguments, though it would require to create arrays of objects before passing them to a block, which could impact readability and clarity of the code. In order to combine pros and cons of these two approaches, the solution follows the implementation in Pharo according to [4, p. 65]. There are keyword methods implemented for up to three parameters (`value:`, `value:value:`). For more than four parameters, a special keyword message `valueWithArguments:` is implemented, where an array of parameters is expected.

2.5.1 Non-local return and block scoping

Block closures are an essential feature to SOM. They allow the implementation of conditionals and loops as messages rather than them being baked in the language syntax. Blocks however bring some dynamic runtime semantics that is not straightforward. When used to the extreme, blocks can introduce some confusion and generally ugly code. However when used correctly, they offer

Figure 2.9: Example of blocks usage in SOM.

```
| b0 b1 b2 b3 |
b0 := [ 1 + 2 ].
b1 := [ :x | x * x ].
b2 := [ :x :y | x * y ].
b3 := [ :x :y :z | x + y + z ].
"Evaluating the blocks"
b0 value. "Returns 3"
b1 value: 3. "Returns 9"
b2 value: 2 value: 8. "Returns 16"
"Message valueWithArguments: can be used with
any number of parameters"
b3 valueWithArguments: #(1 2 3). "Returns 6"
"The next expression is functionally
identical to the previous one"
b3 value: 1 value: 2 value: 3.
```

great way to improve readability and reusability of the code.

Every method has its defined context – a set of variables and objects accessible from the method at given point in an execution. Variables accessible by blocks are bound during runtime, in the context of where the block is *defined*, rather than executed. This is reflected in creation and handling of frames (which can be considered a representation of the context of given method). This also ties in to how returns from blocks function.

The context in which the block is created (and evaluated) is commonly called as *home context* of the block. The block home context is basically a representation of a particular point in the program execution. When a return statement is executed, the execution steps out of the current context and returns to the caller. This can be an implicit return (every method or block implicitly returns the receiver of the message). User can decide to return a different value, denoted by explicit return statement, (^ token).

The behaviour of explicit return in blocks is where the term non-local return comes in. *Non local return returns to the sender of the block home context, i.e., to the method execution point that called the one that created the block* [5]. The important thing is that the home context is tied to the creation of the block, not its evaluation. The pitfall here is there can be a situation where home context of a block being evaluated could end before the block attempts a return to it. This will result in runtime error.

Consider the example on figure [2.10](#)

In the run method, a local variable `x` is created, than assigned an Integer of value 0. Then a loop is executed. To better understand the example, figure [2.11](#) shows the implementation of method `whileTrue:` in the `Block` class.

As apparent from the two examples, there are multiple blocks used and


```

NLReturn = (
  run = (
    | x |
    x := 0.
    [ x < 5 ] whileTrue: [
      x println.
      x := x + 1.
    ]
  )
)

```

Figure 2.10: Example to demonstrate non local return.

```

Block = (
  whileTrue: aBlock = (
    self value ifFalse: [ ↑nil ].
    block value.
    self restart.
  )
)

```

Figure 2.11: Implementation of `whileTrue:` method in `Block` class.

non local returns are crucial to achieve the functionality of the loop (as implemented). The method `restart` is a primitive block method. It could be considered a jump – it jumps to the beginning of the current context. The details of this method are not relevant to non local returns, therefore it will be discussed elsewhere.

Taking a look at the example, the loop consists of the following:

- Local variable `x` is defined and assigned a value 0.
- A block `[x < 5]` is created. The home context of this block is the context of `run` method. The block is therefore able to access the local variable of the method (variable `x` in this example).
- Before the actual message send, the argument is created. This results in another block being created with the same home context. Again, the block has access to the `x` variable each time it is evaluated.
- The message `whileTrue:` is sent to the first block, with the second one as an argument.

For reference, the implementation of relevant boolean methods used in the example are provided on figure [2.12](#)

The execution then continues in the corresponding method:

```
True = (  
  ifFalse: aBlock = ( ↑nil )  
)  
  
False = (  
  ifFalse: aBlock = ( ↑aBlock value )  
)
```

Figure 2.12: Implementation of relevant boolean methods.

- A block [`x < 5`] (`self`) is evaluated, returning the value `true`.
- The instance of `true` is sent a message `ifFalse:`. This returns the `nil` value.
- Execution continues by evaluating the argument block. This prints out the current value of `x` and increments its value by one. Note that even though the block is evaluated from the context of method `whileTrue:`, it is actually evaluated in its home context and still has access to the locals of the home context. The closure captures the variable upon creation, not upon evaluation.
- The `restart` message is then sent to the first block, starting the execution of the method again.

This gets executed until the variable `x` holds a value 5. By then, the `self value` expression returns `false`, therefore `ifFalse:` is sent to a different object. The argument is the block type, its home context is the `whileTrue:` method context.

State of the execution environment upon this `ifFalse:` message send is depicted on figure [2.13](#).

The argument of the method `ifFalse:` contains a return statement. This is the non local return. When the block is evaluated, the execution *jumps out of the block's home context*. This means the execution ends up in the `run` method, returning the value `nil`.

The contexts can be represented by stack frames. When performing local return from a current context, the execution returns to the point that created it – that means one frame is removed. In non local return, the execution can jump any numbers of contexts – the home context does not even have to exist anymore.

The concept of non local returns can get confusing when dealing with blocks assigned to variable, mainly instance variables. The context in which the block is created and evaluated can be far apart. Consider the example at figure [2.14](#), where two blocks are assigned as instance variables – one with a non local return, one without it.

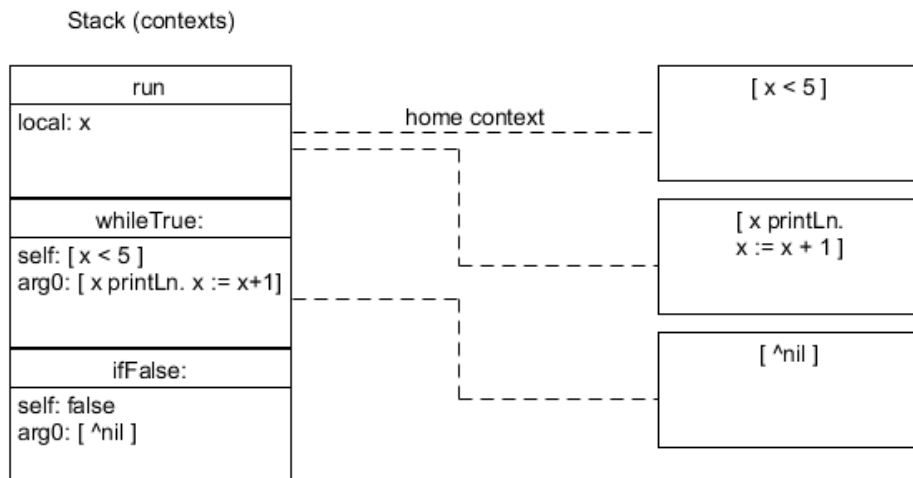


Figure 2.13: Created blocks and their home contexts.

```

NLReturn2 = (
  | lBlock nlBlock |

  run = (
    lBlock := [ 'Local return block' println ].
    nlBlock := [ ↑'Non local return block' println ].
    1 to: 10 do: [ :index |
      index = 3
        ifTrue: [ nlBlock value ]
        ifFalse: [ lBlock value ]
    ].
    'Run method is exiting' println
  )
)

```

Figure 2.14: Second example code for non local returns.

The code assigns the instance variables block objects. Then blocks are evaluated in the loop, depending on the value of the index of the iteration. If the index is equal to three, block with non local return is evaluated. Then the program prints out exiting message. The output of the program is on figure [2.15](#).

As is evident from the program output, the execution stops after the third step of the iteration. The expression following the iteration block is also not executed. When taking a look at the block contexts, it is clear why it is so.

There are multiple methods being called before the block is executed. With

Local return block
 Local return block
 Non local return block

Figure 2.15: Output of the example code from figure 2.14

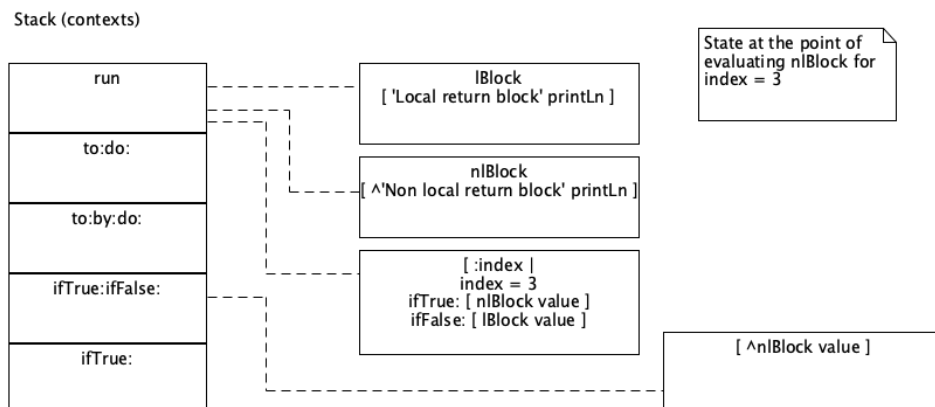


Figure 2.16: State of execution before return for example on figure 2.14.

block in variable `lBlock`, the block’s home context plays little role – after the block is executed, simple return is performed (from the method `value`). With `nBlock`, non local return is performed after evaluation – meaning the execution jumps out of the home context of the block – which means the end of execution of the program.

As stated earlier, the non local returns offer a great way to write more compact and readable code. There are some pitfalls to be aware of. The nature of the non local returns means that it is not always obvious where the execution ends up.

The last example serves to demonstrate how a created block can attempt to exit a context that no longer exists.

Execution of this program will result in a runtime error. Both of the blocks are created in the context belonging to the method `createBlocks`. While it may seem, that at least `block1` could be executed as the evaluation of that block returns to the main method. However, none of the blocks can be evaluated in this case and the reason becomes more clear with a slight modification to the `createBlocks` method, as seen on figure 2.18.

This modification to the method is completely valid – both block are created within the method’s context and thus can access its local variables. However, after the return from this method, its context is discarded. By the time blocks are evaluated, variable `mssg` does not exist.

```

NLReturn3 = (
  | block1 block2 |
  run = (
    self createBlocks.
    block1 value.
    block2 value
  )

  createBlocks = (
    block1 := [ 'Local return' println ].
    block2 := [ ↑'Non local return' println ]
  )
)

```

Figure 2.17: Third example demonstrating non local returns.

```

createBlocks = (
  | mssg |
  mssg := 'Block evaluation'.
  block1 := [ mssg println ].
  block2 := [ ↑mssg println ]
)

```

Figure 2.18: Modification of the method from example [2.17](#).

2.6 Expressions

According to [\[3\]](#), *an expression is a segment of code in a body of executable code that can be evaluated to yield a value as a result of its execution.* Expressions can contain another expressions.

```

expression: assignment | evaluation;
assignment: assignments evaluation;
assignments: assignment+;
assignment: variable ASSIGN;
evaluation: primary messages?;
primary: variable | nestedTerm | nestedBlock | literal;
messages:
  unaryMessage+ binaryMessage* keywordMessage?
| binaryMessage+ keywordMessage?
| keywordMessage;

```

Figure 2.19: ANTLR grammar snippet for expressions.

Syntactically, an expression can consist of [3](#):

- literal,
- variable/constant reference,
- message send,
- nested expression.

2.7 Control structures

In Smalltalk, there are no built-in control structures, unlike for example C++ or Java. SOM follows this principle from Smalltalk, therefore there are no grammatical rules for branching or loops.

The way controlling the flow of program works in SOM is, again, by sending messages. One big advantage of this approach is that the programmer can define their own control structures, simply by implementing classes and methods as needed.

To make working with SOM easier and faster, my implementation provides multiple message implementations, corresponding to the most used control structures in other programming languages. Syntax of these messages corresponds to other Smalltalk dialects.

2.7.1 Conditional branching

There are 3 messages that function as an if control structure. Selectors for these messages are `ifTrue:`, `ifFalse:`, `ifTrue:ifFalse:`. As apparent, they are keyword messages, the receiver is an instance of a Boolean class. All of these messages take blocks as arguments, then evaluating or not evaluating them based on the Boolean value. [Figure 2.20](#) shows a simple example of usage.

```
"Subtracts b from a only if a is greater than b"
a > b ifTrue: [ a - b ].
a <= b ifFalse: [ a - b ].
"Subtracts the smaller number from the bigger one"
a < b
  ifTrue: [ b - a ]
  ifFalse: [ a - b ]
```

Figure 2.20: Example of messages functioning as *if*-control structures.

2.7.2 For loops

The simplest example of a for loop is iterating over a range of integers. There are 2 messages, `to:do:` and `to:by:do:.` The receiver of the message is an integer. The receiver of the message is the lower bound of the iteration, the argument for `to:` keyword is the upper bound, `by:` specifies a step of iteration, `do:` takes a block that is evaluated (note that the block has to have exactly one parameter, so it is possible to capture the value of index in every step).

```
"Prints all numbers from 1 to 10"
1 to: 10 do: [ :index | index asString println ].
"Prints all the even numbers between 1 and 100"
0 to: 100 by: 2 do: [ :index | index asString println ]
```

Figure 2.21: Example of simple for loops.

This way of looping is also usable when iterating over arrays (or any indexable collection). As seen on figure [2.22](#), this method is not very concise, therefore a message `do:` is implemented. Array class implements a method corresponding to this message, iterating over every element of the array. It takes a block as an argument. The block has to have one parameter – that is the element of the array of the given step of the iteration.

```
| array |
array := #(1 2 3).
"Printing the elements by iterating over index"
1 to: array size do: [ :index |
  (array at: index) asString println ].
"Printing the elements by iterating over array"
array do: [ :element | element asString println ]
```

Figure 2.22: Examples of different ways of iterating over an array.

2.7.3 While loops

While loops are implemented as a keyword message sent to a block that returns a boolean value. There are actually two messages, `whileTrue:` and `whileFalse:.` The first one repeats the evaluation of an argument (a block) as long as the receiver returns `true`. The second one, as the name suggests, does the same thing if the receiver returns `false` value. Example in figure [2.23](#) shows printing numbers from 0 to 10 using `whileTrue:` message.

```
| index |  
index := 0.  
[ index < 10 ] whileTrue: [  
    index println.  
    index := index + 1  
]
```

Figure 2.23: Example of while loops.

2.7.4 Class hierarchy

SOM supports single inheritance and it is a vital aspect of the language. As a purely object-oriented programming language, it takes full advantage of polymorphism.

Single inheritance permits only one superclass per class. The superclass can be explicitly defined. If it is not defined, every class is a subclass of special Object class. This therefore means that every object is an instance of the Object class or its subclass.

Tied with inheritance is a principle of late binding. The method to invoke is decided during runtime by name (the message selector). The method lookup is simple – it follows the chain of inheritance, from the subclass to superclass (always ending at Object class).

2.8 Bytecode

The solution uses a bytecode as an intermediary form to represent a program. The bytecode is saved in a binary file that can be interpreted. The structure of bytecode files and semantics and syntax of operation codes is described in the following sections.

2.8.1 Program structure

SOM program has a very simple structure consisting of:

1. **The constant pool:** This is a list of all the entities of the program. The choice of the word *entity* over *object* is intentional to avoid confusion with what objects are in OOP languages. Each entity can be accessed by its index.
2. **Entry point:** An index to a Method that is executed on program start. There can only be one entry point to a program. It is a unary method with selector `run`. It can be a member of any class of the program.

2.8.2 Program entities

All entities in the constant pool are one of these types:

1. **Nil entity** represents an undefined value.
2. **Int entity** represents a 32 bit signed integer number. It is used for LIT instructions.
3. **Double entity** represents a double-precision floating point number.
4. **String entity** represents a value of string of characters of arbitrary length. It is used for constants in the program as well as to store all the identifiers to classes, method selectors and variables.
5. **Field entity** represents a variable in an object. It consists of one index to a string value that represents the name of the slot.
6. **Method entity** represents a method of an object. It holds and index to a string representing the selector, number of arguments (arity of the corresponding message), number of local variables and an array of instructions.
7. **Primitive entity** is a method that needs an implementation in the VM. These are used to handle constructs that cannot be expressed in the SOM language.
8. **Block entity** is a block of code. It holds the number of arguments and an array of instructions (similar to a method).
9. **Class entity** represents the structure of objects. It consists of an array of indices to all the fields of the object. Each one of these fields point either to a Field entity or a Method entity.

2.8.3 Instructions

- **LIT *i*** retrieves a constant value from the constant pool at the index *i* and pushes it on the stack. The item can be either integer, double or string value.
- **GET SLOT *i*** pops a value from the operand stack, assuming it is an object. Then it retrieves a value with index *i* from the constants pool, assuming it is a string. It then retrieves the value stored in the slot with the name specified by the string and pushes it onto the stack.
- **SET SLOT *i*** pops a value from the stack. This value is then assigned to an instance variable with identifier at index *i* in the constants pool.

- `SEND i n` sends a message to an object, which in most cases results in calling a method. A new frame is created on the execution stack, arguments are pushed and the execution jumps to the first instruction of the method.
- `GET LOCAL i` retrieves a local variable with an index `i` and pushes it to the top of the stack.
- `SET LOCAL i` pops a value `x` from the top of the stack and then assigns the `x` into a local variable with the index `i`.
- `GET SELF` retrieves the callee of executed method. The object is pushed to the top of the stack.
- `GET ARG i` retrieves the `i`-th argument of the current message from the stack and pushes it on top.
- `BLOCK i` creates a code block object. The argument `i` points to a block value in the constant pool. The block object is instantiated on the heap and pushed to the top of the stack.
- `RET` is used to return from a method call. The value from the top of the stack is returned. The address to return to is retrieved from the current frame, then the frame is popped and execution jumps to an instruction after the `CALL` that invoked the method.
- `RETNL i` - non local return. The value at the top of the current frame is used as the return value. Argument `i` specifies the type of non local return.

2.9 Garbage collection

The process of *garbage collection* performed by *garbage collector (GC)* is the process of allocating and freeing memory during application runtime. The main advantage of this mechanics is to prevent *memory leaks* – parts of a program that allocate memory without freeing it when it is not needed [6]. Most modern high-level programming languages implement some form of garbage collection.

There are multiple possible algorithms that solve this problem. I have decided to implement a mark and sweep algorithm, due to its simplicity. The expectation is that this algorithm will not be particularly fast, though it leaves a lot of room for possible improvements and it can be used as a demonstration of performance effects of a garbage collector.

2.9.1 Mark and sweep

According to [7], the algorithm consists of two main phases:

- Mark phase – discovery of all the reachable objects.
- Sweep phase – clearing the heap of all unreachable objects.

Every object allocated on the heap holds a mark bit. This represents the reachability of the object – 0 (false) for unreachable, 1 (true) for reachable. Upon object creation, the mark bit is set to zero. In the sweep phase a simple graph traversal algorithm can be used (such as DFS – Depth First Search) to mark all the reachable objects. Every object can be considered a node and variables serve as neighbour lists. Root nodes would be local variables and fields that are directly accessible.

In the sweep phase, the whole heap can be traversed linearly and objects with mark bit set to zero are deleted. As the [7] states, main advantages of this algorithm are:

- algorithm handles cyclic references and therefore cannot end up in an infinite loop,
- no additional overhead during the execution of the algorithm (such as extra data structures etc.).

On the other hand, the simplicity of the algorithm means there are some disadvantages:

- normal program execution is suspended during the garbage collection process,
- the algorithm does not address memory fragmentation.

The problem of memory fragmentation means that after several runs of the algorithm, the reachable objects will be separated by chunks of free memory. This can be solved by shuffling the objects, though at a cost of further performance degradation.

To demonstrate the execution of the algorithm, consider the following object hierarchy allocated during a program execution. Each node represents an object allocated on the heap, with marked bit set to false. Arrows represent variable references from one object to another. In reality, there will be multiple root objects – therefore the algorithm can run for every root object.

The heap is then traversed from the root object (figure 2.25). Each colour represents one step of the traversal. This example represents breadth first search.

After the run of the mark phase, all the objects reachable from the root objects (or multiple root objects after a run for each one) are marked. One linear traversal over objects can then remove all of the unmarked ones (figure 2.26).

2. ANALYSIS AND DESIGN

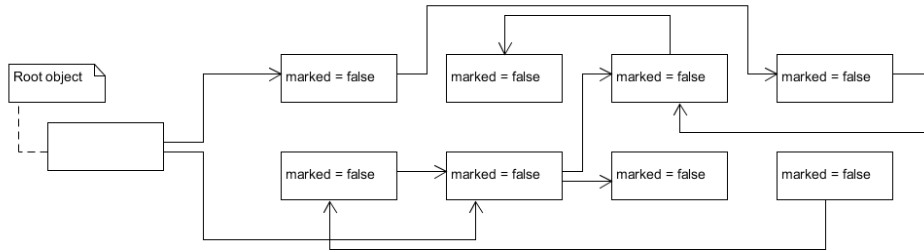


Figure 2.24: Illustration of the object hierarchy before GC algorithm run.

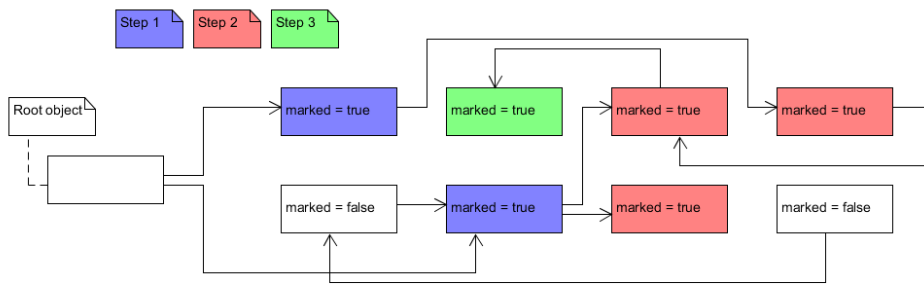


Figure 2.25: Traversal of the heap objects.

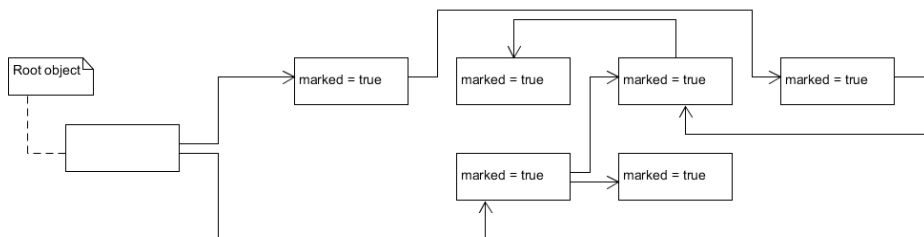


Figure 2.26: State of the heap after the GC run.

Realisation

3.1 Program overview

The implemented program provides a way to compile SOM source code and execute it.

```
<som_executable> [OPTION] [SOURCE]
```

Figure 3.1: Interface of the program.

The interface of the application is simple and consists of two user provided arguments.

The argument `OPTION` can have two values and alters the mode the app will function in:

- `-c` is compile mode. The argument `SOURCE` is a folder containing the source code of the program. This folder is searched for SOM source files – those are recognized by their file extension, which should be `.som`. The folder is searched non recursively. Every SOM source file is compiled and one binary file containing bytecode is created. The name of the file is the same as the provided folder name.
- `-r` loads and runs a compiled bytecode. The `SOURCE` argument is a file name of the compiled bytecode.

The SOM source files lookup uses standard C++ `std::filesystem` implementation. This is the reason the compilation requires standard C++17 or newer.

3.1.1 Build

The project uses CMake to automate the building process. ANTLR executable is needed, either available at the system's `PATH` variable or with its path spec-

ified in the root `CMakeLists.txt` file. The build process then handles the following:

- downloads the ANTLR runtime from GitHub,
- uses the ANTLR executable to generate the parser and visitors from grammar files,
- compiles and links the project into an executable.

3.2 Source code parsing

The first step of executing the source code is lexical and syntactical analysis. There are multiple solutions to streamline the implementation of parsers. For this project, ANTLR is used (ANother Tool for Language Recognition). It is a tool that is able to generate a parser from a formal grammar defined for a language.

Grammar for SOM is defined in two files – `SOMLexer.g4` and `SOMParser.g4`. These contain rules definitions for lexer and parser respectively. Every rule in the grammar has its name and definition.

Lexer rule names always start with an uppercase letter. These are the rules that form a foundation for the subsequent parser rules. Every lexer rule represents a token of a language, such as keywords, special characters etc.

Figure 3.2 shows definitions of some of those lexer rules. The simplest rules define one or multiple character tokens (rules like `NewTerm` or `Primitive`). Comments are also defined here – the rule tells the lexer not to tokenize any characters between quotation marks. Format of identifiers is also defined here (for classes, variables, messages etc.). String and numerical values are also defined here.

Parser rules are where the structure of the source code comes in. Parser rules are defined in file `SOMParser.g4`. They help build the abstract syntax tree (or at least reflect the structure of the AST to some degree) [8]. For reference, excerpts of parser rules are demonstrated on figures 3.3 and 3.4.

3.3 Abstract Syntax Tree

After the parsing is complete, Abstract Syntax Tree (AST) is constructed. AST is, by definition, stripped of many syntactic detail. It mainly represents the structural and content-related aspects of the code.

The conceptual design of the AST is depicted on figure 3.5.

3.3.1 AST Nodes

Class node represents a class in the program, while the program itself is basically an array of different classes. A class holds its name, its member

```

NewTerm: '(';
EndTerm: ')';
Primitive: 'primitive';

Comment:  ''' ~["]* ''' -> skip;
Identifier: [\p{Alpha}] [\p{Alpha}0-9_]*;
STString:
    '\''
    (
        '\\t'
        | '\\b'
        | '\\n'
        | '\\r'
        | '\\f'
        | '\\0'
        | '\\\\'
        | '\\\\\\'
        | ~('\''| '\\')
    )*
    '\'';
Integer: [0-9]+;
Double: [0-9]+ '.' [0-9]+;

```

Figure 3.2: Example of some interesting lexer rules from SOM ANTLR grammar.

```

classdef:
    Identifier Equal superclass
    instanceFields method*
    ( Separator classFields classMethod* )?
    EndTerm;
superclass: Identifier? NewTerm;
instanceFields: ( Or variable* Or )?;
classFields: ( Or variable* Or )?;
classMethod: method;

```

Figure 3.3: The syntax of class definition described by ANTR grammar.

fields (member variables in C++ terminology), instance-side methods (member functions), class-side fields (static member variables) and class-side methods (static member functions).

Method node represents a method – instance or class side. It consists of a pattern, local variable definitions and a block to be executed.

```
nestedBlock: NewBlock blockPattern? blockContents? EndBlock;  
blockPattern: blockArguments Or;  
blockArguments: ( Colon argument )+;
```

Figure 3.4: The syntax of nested blocks.

Pattern represents a message corresponding to the method. There are 3 types of messages in SOM, therefore there are 3 distinct types of patterns:

- The simplest one is **unary pattern** – consisting of only one identifier as there are no arguments.
- **Binary pattern** is treated as a separate pattern. It consists of identifier and exactly one argument. There are special requirements for binary pattern identifier – there is a special set of characters permitted that can form a binary pattern.
- **Keyword pattern** then consists of one or more keywords and same number of arguments, each corresponding to one keyword. Concatenation of keywords form a selector of the method. **Keywords** holds the string value of the keyword, always ending in colon (:).

Variable node represents instance/class side variables, arguments to messages or blocks. It holds the identifier of the variable as a string value.

Block represents a block of executable code with its own scope. The simplest block consists of local variable definitions and an array of expressions to be evaluated. While this is enough to represent a method block, other uses may require more information, therefore there is another similar node discussed later.

3.3.1.1 Expressions

Expression is an abstract term in the context of the AST - there are two types. The common thing is they can be evaluated – therefore forming the actual executable code of the program.

Evaluation is the first form of expression - it represents a message sends to an object, thus returning a single value when evaluated. This node consists of messages (optional) and a *primary*.

Primary is another abstract concept. In its core, a primary represents an object, though there are multiple ways to reference an object in SOM. There are four AST nodes that can be classified as a primary:

- **Literal** – a constant basic value (of integer, floating point, string or array type). Each of these have their dedicated literal node holding the value as seen on figure [3.5](#).

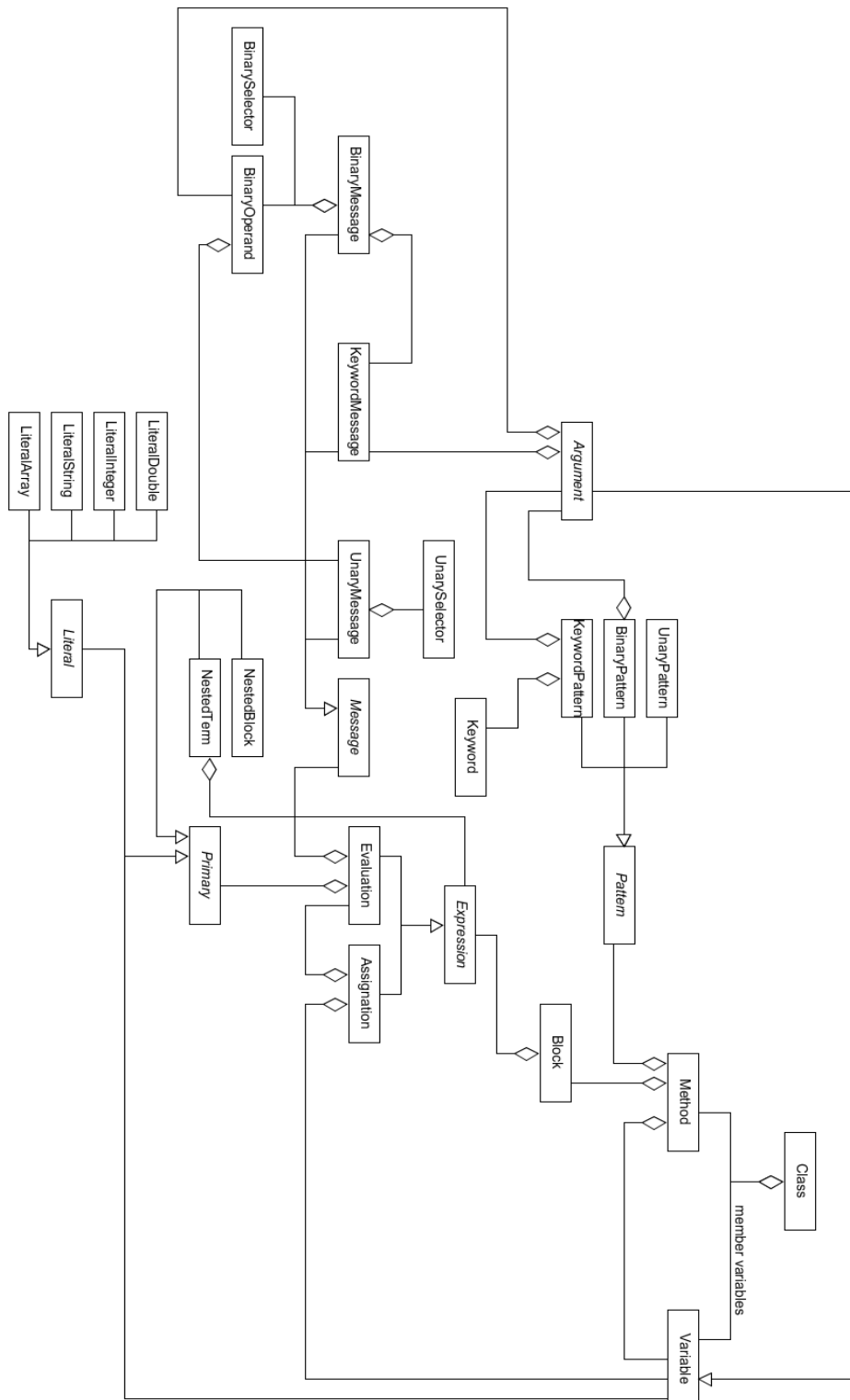


Figure 3.5: Conceptual class diagram for SOM abstract syntax tree.

- **Variable** is self explanatory – a reference to an object accessed via the identifier.
- **Nested term** is an expression that needs to be evaluated to retrieve the reference to an object. Syntactically, the nested terms are enclosed in parentheses.
- **Nested block** is a block of expressions returning the reference to an object. It is enclosed in square brackets in the syntax. Nested blocks consist of the same elements as block discussed with methods with addition of a block pattern – nested block can have their arguments.

The second part of the evaluation node is the message sends to the primary. There are three types corresponding to three types of messages in SOM. **UnaryMessage** node is self explanatory – there are no arguments, only the message selector. **BinaryMessage** holds its selector too with addition of the argument. The argument of the binary message send can be a primary, along with unary message sends (because unary message sends take precedence). **KeywordMessage** is made up of the keywords (forming the selector) and something called *formulas*. Formulas are binary (and also unary) message sends, that take precedence over keyword messages.

Assignment is the second form of an expression. The name suggests this node represents assigning a value into a variable. Therefore the node consists of the **Variable** node to assign to and an **Evaluation** node returning the value to assign.

Expression with an explicit return statement is encapsulated in an **Result** node. The information about explicit return being used is needed to handle non local returns.

Consider the following expressions:

```
"Expression 1"  
receiver1 foo: x asString bar: (x + y toInteger)
```

Figure 3.6: Example for AST construction.

Expression 1 is a keyword message send to a variable `receiver1`, with arguments as sub-expressions. The AST produced for this expression is on figure [3.7](#).

3.3.2 AST construction

The AST is constructed by visiting over the ANTLR-generated parse tree. The visitor is implemented in class `CParseTreeConverter`. This is a subclass of `SOMParserBaseVisitor`, which is a base visitor implementation provided by ANTLR that performs depth-first traversal over the parse tree. Some

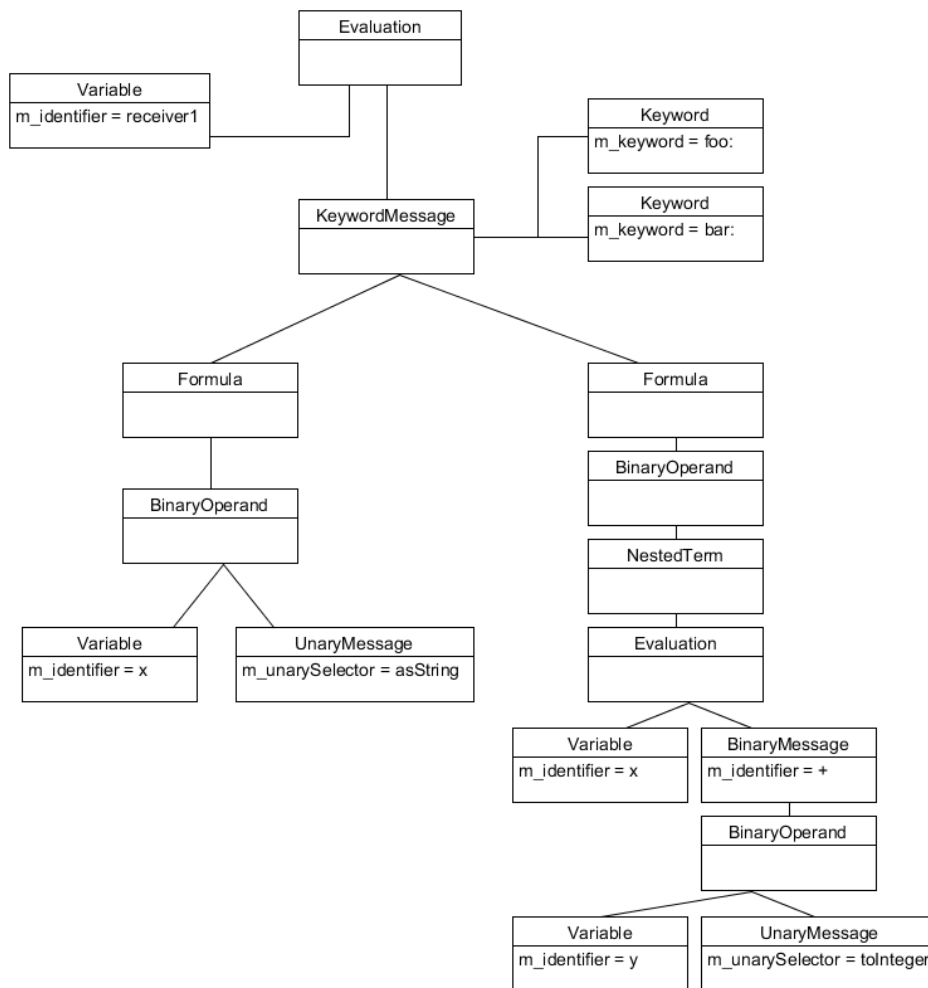


Figure 3.7: Expression from example 3.6 as AST.

member functions in `CParseTreeConverter` are not overridden and make use of this default behaviour (which is just iterating over child nodes and visiting them).

ANTLR also provides an abstract class that defines the interface of a visitor over a parse tree. Every member function of the visitor returns a special value – `antlr::Any`. It is a special class defined to hold any data type. The implementation therefore has to convert these values (to pointers to AST nodes in this case).

3.4 Bytecode

After constructing the AST it is compiled into bytecode. The bytecode definitions are located in source files `Bytecode.h/Bytecode.cpp`. The actual compilation is implemented as a depth-first traversal of the AST, therefore a visitor pattern is used. There is an abstract class `ASTVisitor` that defines the interface of any AST visitor. This can be used to further implement visualisations of the AST or to add support for different bytecode instructions sets, for example Java bytecode to add support for running inside a Java runtime.

3.4.1 Values

Every constant value in the code is saved as a `Value` struct. Each value contains the one byte tag and the actual value to hold. There are method implementations to print every value into human readable format for better visualisation of the compiled bytecode. Additionally, every instruction struct is able to serialize itself – write the data needed in binary format to a file.

3.4.2 Instructions

Similar to values, every instruction is represented by a struct holding the relevant information (such as operation codes and arguments). Every instruction is capable of printing itself in human readable format and serialize itself to binary format, the same as all the values.

3.4.3 Compilation

The process of compilation consists of a walk over the AST and constructing the program represented by the `Program` class. This is the class that holds the constants pool and all of the instructions. It also holds the information on the entry point of the program (address of the first instruction to execute).

The AST visiting and compilation is handled by the class `CBytecodeCompiler`, a subclass of abstract AST visitor. The process traverses every class definition and visits all its children nodes. The compiler keeps track of all the scopes to help resolve every encountered identifier.

The high level look at a compilation of a class is:

- The class identifier is registered as a string constant in the constants pool.
- Instance side and class side variable identifiers are registered as string constant in the constants pool.
- Every method defined in the class is visited (and compiled). For every method, method value is registered in the constant pool.

- Class value is then added into the program. All of the previously compiled variables and methods are registered as slots.

3.4.3.1 Method compilation

When a method node is encountered, a new scope is created on the scope stack. The scope stack contains `CMethodCompilationCtx` objects. Every such object holds:

- a vector of local variable identifiers,
- a vector of formal argument identifiers.

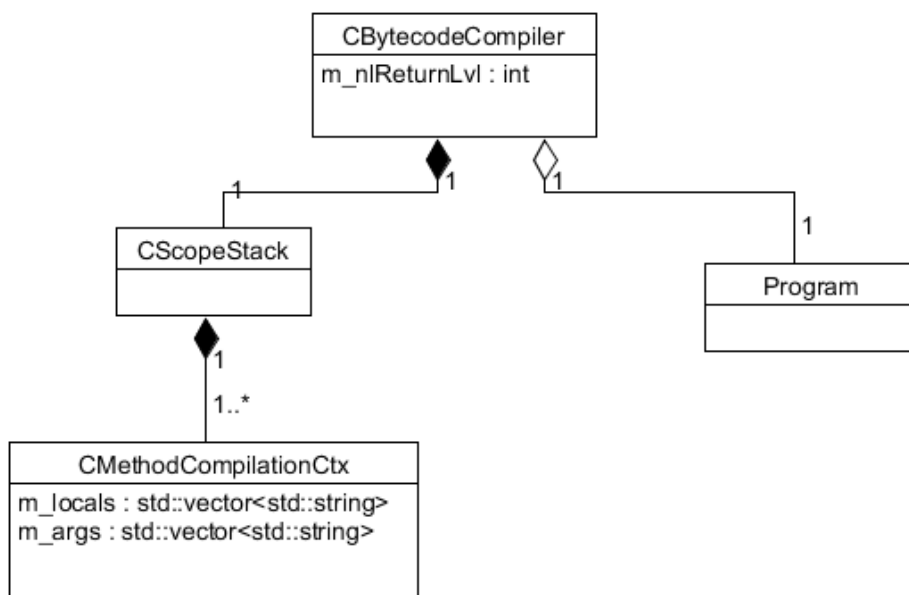


Figure 3.8: Conceptual class diagram for compilation environment.

The position of the identifier in the vector then determines the index that is assigned to said argument / local variable. That index is then used for `GET ARG` / `GET LOCAL` instructions. When the method compilation finishes, its context is popped from the stack.

After the arguments and locals are registered in the method scope, the method selector is then added to the constants pool as string value. For unary and binary methods, the selector is saved as defined. For keyword messages, the selector is defined as a concatenation of all of the keywords (including the colons). This means that for a method defined as `foo: fooArg` `bar: barArg`, the selector in the constants pool ends up being `foo:bar:.`

After that, the method block is visited. An empty vector of instructions is created, then every expression is compiled and added to the vector. Whenever and identifier is encountered during this, it is resolved as:

- if the identifier is of value `self`, the result is `GET SELF` instruction,
- if the identifier is registered as the current method's local variable with index `i`, the result is `GET LOCAL i` instruction,
- if the identifier is registered as the current method's argument with index `i`, the result is `GET ARG i` instruction,
- otherwise, the identifier's value is either looked up in the constants pool (if the value exists) or registered (if it does not exist); the index `i` of the value is then used in the resulting `GET SLOT i` instruction.

Handling of literal values is straightforward – the value is registered in the constants pool and `LIT` instruction is returned. String delimiters are stripped from the quotations in this step.

3.5 Interpretation

The process of interpreting the loaded bytecode is handled by class `CInterpret`. The class diagram of the interpretation environment is in the figure [3.9](#). The diagram serves as a reference point, all the parts are described in further detail in the following sections.

3.5.1 Garbage collection

The process of garbage collection is a part of the interpretation loop. The algorithm used is mark and sweep (described in subsection [2.9.1](#)). There is no mechanism implemented to explicitly trigger garbage collection. As running the process after each executed instruction would have a potential of slowing the execution down, the garbage collection is triggered after every context switch.

3.5.2 Program counter

The class `CProgramCounter` is an implementation for a program counter for the VM. During initialization of this object, the entry point of the program is loaded. The program counter also saves the "end" address. The addresses are implemented as iterators to a vector of instructions – that is possible simply because all of the executable code will be contained within a method, and all the method instructions are loaded into a vector of instructions. The entry point of the program is therefore the first instruction of the `run` method, while

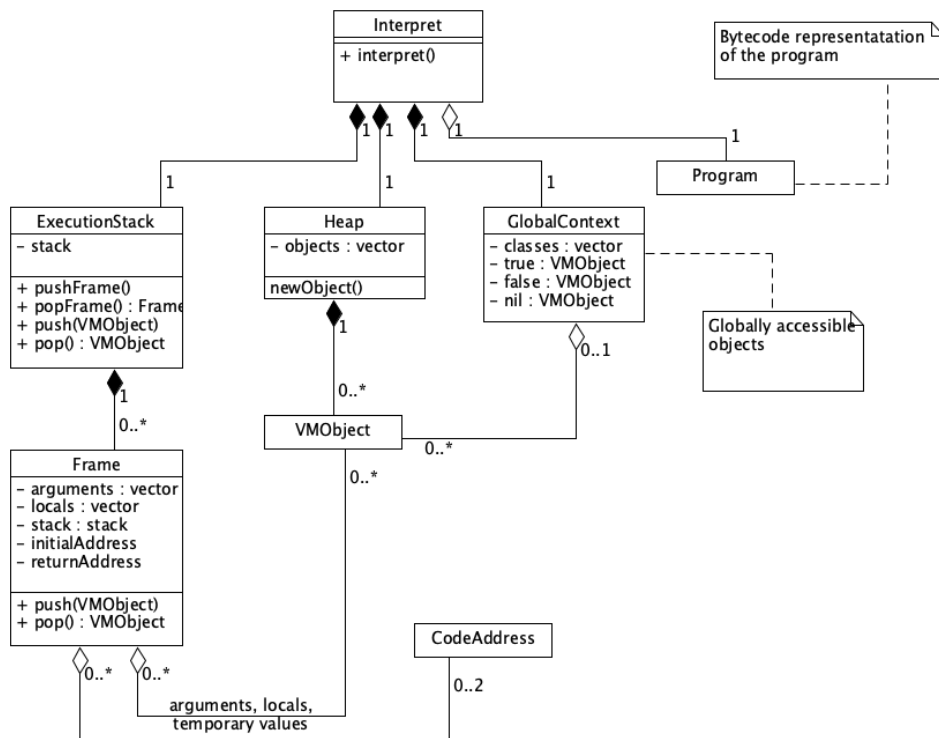


Figure 3.9: Conceptual class diagram for interpretation environment.

”end” address is the end iterator of the same vector. This means that the end address is not an executable instruction (nor it is a dereferenceable iterator).

3.5.3 Execution stack

The execution stack is responsible for management of method calls (and nested block evaluations). It is implemented as a LIFO (last in, first out) structure. For every method call or block evaluation, a new frame is created and pushed on the stack. A frame contains:

- Arguments accessible by their index, assigned at compilation. For the sake of stack, the callee is also considered an argument and is always the last element of the arguments array.
- Array of local variables, accessible by their index assigned during compilation.
- Return address. The program counter is set to this address when a RET instruction is executed.

3. REALISATION

- The address of the beginning of the current context. This is used only to achieve iteration without the need of recursion.
- Local data for the method/block execution. This part is used for all the temporary object creation (e.g. literal values) and argument passing.

Sending a message or evaluating a block therefore can be split into multiple steps:

- The receiver of the message is pushed to the stack.
- The arguments are then pushed to the stack, in the order from left to right.
- When the **SEND** instruction is executed, new stack frame is created. The return address is initialized to the address of the instruction following the **SEND** instruction. Array of arguments is initialized with the values from the top of the stack, number of arguments is provided as an argument of the instruction.
- The new frame is pushed to the stack. The receiver is then accessed, added to the end of the argument array and given the responsibility to invoke the method corresponding to the sent message.

When returning, the top of the stack contains the return value of the method or block. The frame is popped from the stack, program counter is set to the return address and the top value is pushed to the new top of the stack.

3.5.4 Objects

Every object created during runtime (implicitly or explicitly) is represented by the **VMObject** class. Every value is represented by an instance of this class or its subclass. Every object holds a pointer to its class and a map of the instance values – the identifier as a key, object pointers as a value. Upon creation, instance field values are initialize to the **nil** value.

Every class is represented by a **VMClass** instance, which is a subclass of **VMObject**. This allows every class to be manipulated as an object. There is only one class instance per class definition during runtime and instantiating new class object is not possible. Before the execution of the code, the byte-code is traversed and for every **CLASS** value, a singleton object is created and initialized. These objects are accessible globally by the identifier – the class name.

Every class holds all the information needed to create and manipulate its instance. This includes the method dispatch.

3.5.4.1 Object creation

One of the main responsibilities of the class object is to handle dynamic creation of its instances. This is done via a `new` message send. Upon the `new` message send:

- `new VMObject` is created on the heap and its class is set,
- instance fields of the new object are initialized,
- pointer to the new object is returned.

3.5.5 Messages

In SOM, almost everything is handled via message sends. Message sends invoke a method, implemented either directly in SOM, or as a primitive. The process of invoking a method on an object depends on what kind of message is sent. Additionally, a small number of special cases needs to be handled for code blocks, to allow for iteration.

When a `SEND` instruction is encountered, the receiver of the message is retrieved from the stack. As the receiver is pushed to the stack first, before the argument, it needs to be retrieved using the message arity. To achieve the late binding, the selector (a string value) is then retrieved from the constants pool.

Reference to the class of the object than handles the method dispatch.

Every method invocation results in a new frame being created on the execution stack. Because the execution stack is implemented as a stack of separate objects representing the frames (and not a continuous array), every new frame needs to be initialized. This handles the argument passing, as well as scoping of local variables when sending messages to code blocks.

When the new frame is created, multiple actions take place:

- Return address on the frame is set to the instruction following the `SEND` instruction. As every block of code must end with return (local or non local), the existence of the address is guaranteed.
- The array of arguments is initialized. Arguments had been pushed to the stack in the order of corresponding keywords (therefore left to right). As a result, the arguments are reversed in the actual array – member functions access them from the back. The last argument is set to the receiver – that is the object accessed by `self` keyword (comparable to `this` pointer in C++).
- Array of local variable is initialized. This step is skipped for every object outside of code blocks – instances of `Block` class. This is due to the scoping rules – code blocks can access the local variables of the method, in which they are evaluated.

After the new stack frame is created, the execution of the method can take place. For primitive values, the corresponding member function is called directly in the VM. For native methods, program counter is set to the address of the method and interpreting loop continues.

Method execution should always end with the return value at the top of the stack – for primitives and native methods. Every method should also end with the `RET` instruction. When this instruction is executed, the top value of the current stack is taken. The top frame of the stack is removed, program counter is set to its return address and the top value is pushed to the top of the underlying frame. Interpreting loop can then continue.

3.6 Core library

SOM programming language is very light on features. In order for it to be usable, there needs to be an implementation of some fundamental principles provided. For example, the language itself does not provide a way to manipulate numbers, character strings, standard input and output or control structures. All of these can be implemented in the VM itself while preserving the consistency of rules of the language.

In order to achieve that, a keyword `primitive` is defined. This keyword is used for methods that require an implementation in the VM runtime. From the outside, calling a primitive method is no different than calling a method implemented directly in SOM.

While the user is able to implement their own primitives, a lot of them are already provided in the core library. This library is a set of classes loaded every time a SOM code is interpreted. This library provides implementations for:

- strings,
- numbers – integers and doubles,
- boolean values with messages that provided control flow features,
- code blocks,
- arrays.

3.6.1 Primitives

Every class that marks a method primitive has to provide its implementation in the VM. This is done by creating a new subclass of `VMClass`. This subclass then has to implement the primitive methods (in the form of `void` member functions that take `CInterpret*` as an argument). Resolving of method selector (which is a simple string) and the function to invoke is then done via member function `dispatchPrimitive`.

To demonstrate the functioning of the primitive method, consider the following implementation of the `concatenate:` method in `String` class.

```
void VMString::concatenate(CInterpret* interpret) {
    auto argument =
        interpret->executionStack().getArgument(0)->getValue().asString();
    auto receiver =
        interpret->executionStack().getSelf()->getValue().asString();
    interpret->executionStack().push(
        std::make_shared<VMObject>(
            interpret->globalContext().getClass("String"),
            VMValue(receiver.append(argument)))
        );
}
```

The method takes one string argument and returns the concatenation of `self` and argument. The member function itself does not return anything, it only manipulates the interpreter environment. The objects are retrieved from the current frame of the execution stack and the result is returned as the top value on the current frame.

3.6.2 Strings

Every literal string value is represented by an instance of `String` class during runtime. The corresponding class implementing the primitives is `VMClass`. In SOM, every string object is immutable, therefore every string manipulation results in creating a new string object.

This approach can however lead to inefficiencies when constructing strings. To combat this, a class handling dynamic string creation could be implemented, similar to streams.

When looking at `String` implementation, there is a big difference to how they are implemented in SOM and Smalltalk. When a string object is created in Smalltalk, it is handled as an object with indexed fields, each containing a character. Provided implementation of SOM does not conform to this, as strings are simply handled in the VM runtime. Possible extension of the language could be developed, allowing for creation of objects with indexed fields.

This could be then applied to arrays, thus creating a unified ways to handle collections.

3.6.3 Booleans

The core library also provides an implementation of boolean values. The two values are implemented as separate classes (`True` and `False`) with a com-

```
Nil = (  
  asString = ( ↑'nil' )  
  isNil = ( ↑true )  
  notNil = ( ↑false )  
  
  ifNil: aBlock = ( ↑aBlock value )  
  ifNotNil: aBlock = ( ↑self )  
  ifNil: goBlock ifNotNil: noGoBlock = (  
    ↑goBlock value )  
)
```

Figure 3.10: Implementation of Nil class in core library.

mon superclass (`Boolean`). The runtime then provides a globally accessible instances of these two classes under the identifiers `true` and `false`.

3.6.4 Array

Array object implementation is also provided in the core library. Array class implements messages for element access and iteration over elements. Internally, every array object holds a `std::vector` containing pointer to the element.

3.6.5 Nil

This class implements a special type of object, signifying an empty value. There is only one instance of this class during runtime, accessible globally under `nil` identifier. This object is used to initialize variables (locals, instance fields etc.) that are not otherwise initialized.

3.6.6 Blocks

As blocks are treated as any other objects in SOM, there needs to be a class implemented for them. Every block in the program is an instance of the `Block` class. Given the special nature of blocks (as discussed in section [2.5.1](#)), they require more implementation details compared to other classes.

The actual implementation for blocks is pretty simple – methods for evaluating and methods for iterating (simulating while loops known from other programming languages).

Every block created is an object instance – therefore it is implemented as an instance of `VMObject` in the VM. Every block object is initialized with a special primitive value (struct `BlockContext`), holding information needed to evaluate the block. That is:

- address of the first instruction of the block,

```
Block = (  
  value = primitive  
  value: arg = primitive  
  "More evaluation methods omitted"  
  
  whileTrue: block = (  
    self value ifFalse: [ ↑nil ].  
    block value.  
    self restart  
  )  
  
  whileFalse: block = (  
    [ self value not ] whileTrue: block  
  )  
)
```

Figure 3.11: Block implementation in SOM core library.

- a pointer to the home context of the block (pointer to the stack frame corresponding with the method that created the block).

3.7 Block evaluation

What happens upon the evaluation of the block is as follows:

- The message `value` (or any variant with arguments) is sent to a block object. This is dispatched using standard method dispatch. As a part of standard method dispatch, new stack frame is created and the execution of the primitive method starts.
- Upon the start of the execution of the method, the new stack frame is discarded and replaced by the block's home context. In standard situations, this will result in multiple pointers pointing to a single frame (as execution stack stores pointers to frames). The return address from the discarded frame is saved, as that is the return address in the case of local return.
- Program counter is set to the first instruction of the block.
- The local return address is set in a separate part of the context. It cannot overwrite the return address already in the context, as that is needed in case execution moves to the context of the executed block again.

3. REALISATION

- The interpret loop is called (creating a nested call inside of the main interpret loop). On one of the return instructions, interpret jumps out of the nested interpret call.

Every block ends with `RETNL i` instruction. The integer from an argument decides the behaviour of the return.

- If $i = 1$, the return behaves as a local return. Interpretation jumps out of the nested interpret loop, then removes the top stack frame (as is standard for every return from a method). As a return address, the special value from the initially discarded stack frame is used.
- If $i > 1$, the return is a non local return. Stack frames are popped until the block's home context is at the top. The standard return from that context ensures the execution returns from the block's home context. The execution of the program ends if the home context is the bottom-most frame on the stack.

Consider the following example of a simple block evaluation on figure 3.12:

```
BlockEval = (
  run = ( [ 'Hello' println ] value )
)
```

Figure 3.12: Example to demonstrate block evaluation process.

The state of the interpret environment after the `value` message send (right before the first block instruction is evaluated) is depicted on figure 3.13.

There are few things to note. First, the execution stack contains one stack frame two times. This is to ensure that:

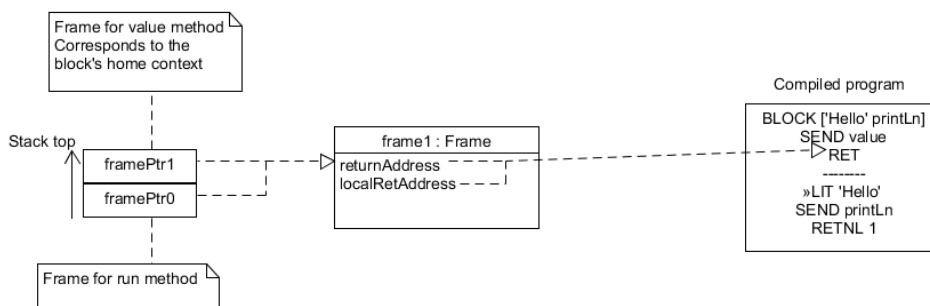


Figure 3.13: State of execution environment for execution of example 3.12.

- block can be evaluated in its home context without the need to copy values on initialization and again on return,
- standard mechanism for returning from methods can be used when performing local return from the block.

In this case, local return address and standard return address are the same, because the block is evaluated in its home context. This does not always have to be the case. In the example, the local return address is used as the block does not perform non local return.

3.7.1 Argument handling

Every block can require arbitrary number of arguments in order to be evaluated. Given that blocks are executed inside already existing contexts (with possible existing arguments), this needs to be handled on two levels:

- During compilation. The arguments are not resolved via their identifier, but they are assigned an index in their scope. The indices of the block cannot clash with already existing ones (arguments of the method creating the block).
- During interpretation. At the beginning of block evaluation, the values of arguments have to be initialized inside the home context. They then have to be removed as to not be accessible from the original method.

3.7.2 Block restart

In order to achieve a possibility of iteration without the need of recursive calls, `Block` class implements a special `restart` method. The method is implemented as primitive and is used in core library. When writing programs in SOM, the programmer should use existing messages aimed at providing iteration functionality, thus the use of this message outside of core library is not recommended.

The method simulates a recursive execution of a method, without the need of allocating and initializing new stack frames on each message send. Implementation itself is pretty simple:

- New stack frame is created as a part of standard message send mechanism. This stack frame is discarded.
- The program counter is set to the address of the initial instruction of current context (stack frame).
- Interpret loop is then called to avoid context exiting mechanisms.

Performance

To see, how a basic SOM runtime implementation compares to other SOM implementations and some fully featured programming languages, a small set of benchmarks was performed. With each one, an implementation code snippet is provided.

4.1 Fibonacci sequence

The first test is calculating a specific position in the Fibonacci sequence. Even though linear complexity is easy to achieve, for the purposes of the benchmark, recursive approach is used. This will result in lot of redundant computations, yielding $\mathcal{O}(2^n)$ time complexity and resulting in easily measurable execution times.

The conditions of the time measurements are as follows:

- Times stated are an average of 100 runs.
- The execution time is measured using `Measure-Command` on Windows Powershell.

The execution time is measured as a time it takes for a command that executes the program to finish. This includes all the potential runtime initializations, procedures after the execution is finished etc.

The tested program in SOM is on figure [4.1](#). It calculates the thirteenth number in the Fibonacci sequence using recursive calls and assigns it to a variable.

Following are programs in C++ (figure [4.2](#)) and Python(figure [4.3](#)), doing the exact same thing.

The results are displayed at figure . As expected, my interpreter of SOM is slower than the traditional programming languages by an order of magnitude.

4. PERFORMANCE

```
Fibonacci = (  
    run = (  
        | result |  
        result := self fibonacci: 13  
    )  
  
    fibonacci: n = (  
        ↑(n <= 1)  
        ifTrue: 1  
        ifFalse: [  
            (self fibonacci: n - 1) + (self fibonacci: n - 2)  
        ]  
    )  
)
```

Figure 4.1: The program for Fibonacci sequence performance test in SOM.

```
int fibonacci(int n) {  
    if (n <= 1)  
        return 1;  
    return fibonacci(n - 1) + fibonacci(n - 2);  
}  
  
int main() {  
    int res = fibonacci(13);  
}
```

Figure 4.2: The program for Fibonacci sequence performance test in C++.

```
def fibonacci(n):  
    if n <= 1:  
        return 1  
    else:  
        return (fibonacci(n-1) + fibonacci(n-2))  
  
res = fibonacci(13)
```

Figure 4.3: The program for Fibonacci sequence performance test in Python.

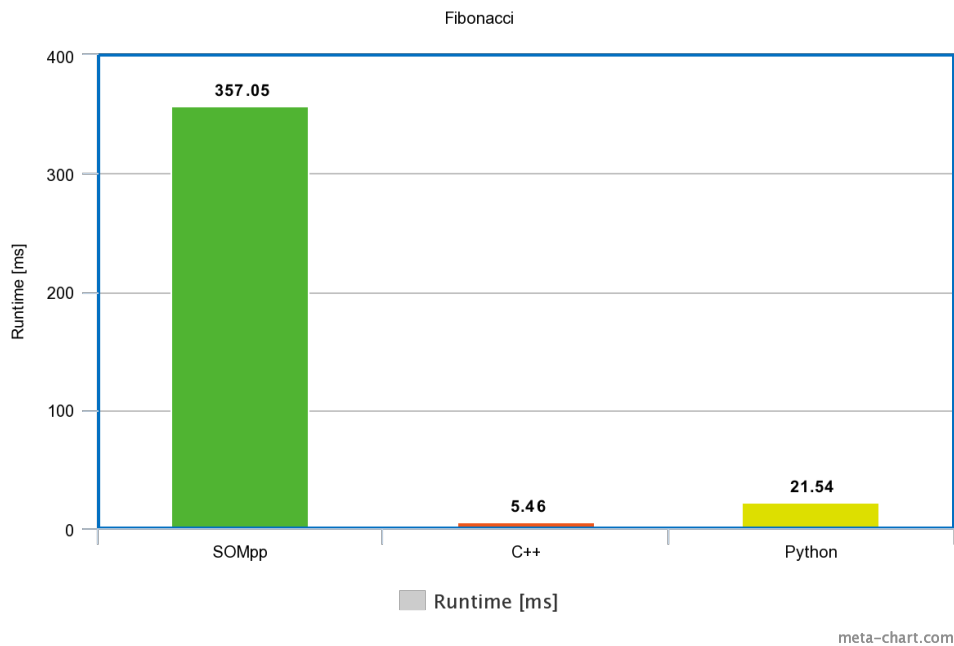


Figure 4.4: Execution times of Fibonacci sequence test.

When compared to other SOM implementations, the differences are smaller, though CppSOM implementation is still the slowest. The assumption is that this is due to the fact that no optimizations (during compilation or during runtime) are implemented.

Figure 4.5 shows comparison with the first Java implementation of SOM, along with the Python one (PySOM). The latter has an option of two types of interpreters – bytecode interpreter and AST interpreter. Both options have been tested and the results are shown. I have also included a bytecode interpreter written in Rust – SOM-RS along with SOM++ (reference C++ implementation).

I have decided not to include CSOM (C implementation) in the comparison, as it failed to parse SOM source files other implementations did.

4.2 While loop

It is pretty apparent that further comparisons of SOM performance to traditional, fully featured programming languages have no purpose. In the next test, I will compare how a very basic language feature performs across the tested SOM implementations.

The program consists of a simple while loop, iterating and incrementing a local variable two million times. The SOM source code (same for all implementations) is on figure 4.6.

4. PERFORMANCE

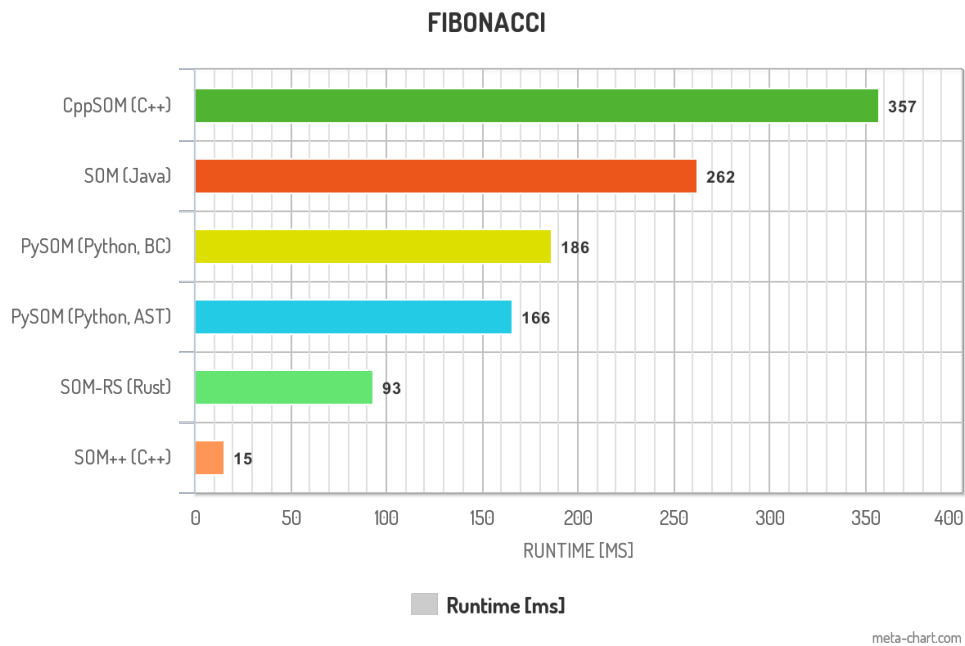


Figure 4.5: Execution times of Fibonacci sequence test for various SOM implementations.

```
WhileLoopBench = (  
  run = (  
    | x |  
    x := 0.  
    (x < 2000000) whileTrue: [  
      x := x + 1.  
    ]  
  )  
)
```

Figure 4.6: SOM source code for second performance benchmark – while loop evaluation.

The results of the test (figure 4.7) are unexpected. The provided C++ implementation of SOM is very close to (highly optimized) PySOM bytecode interpreter implementation and consistently faster than Java one. This could be caused by multiple things:

- The optimizations in PySOM and SOM (Java) did not cover the specific code in the benchmark. This could mean that relative speedup of my implementation is simply due to underlying programming languages used (compiled C++ interpreter against interpreted Java/Python ones).

- The implementation of `restart` method (used in `whileTrue:`) in `Block` class in different virtual machines.

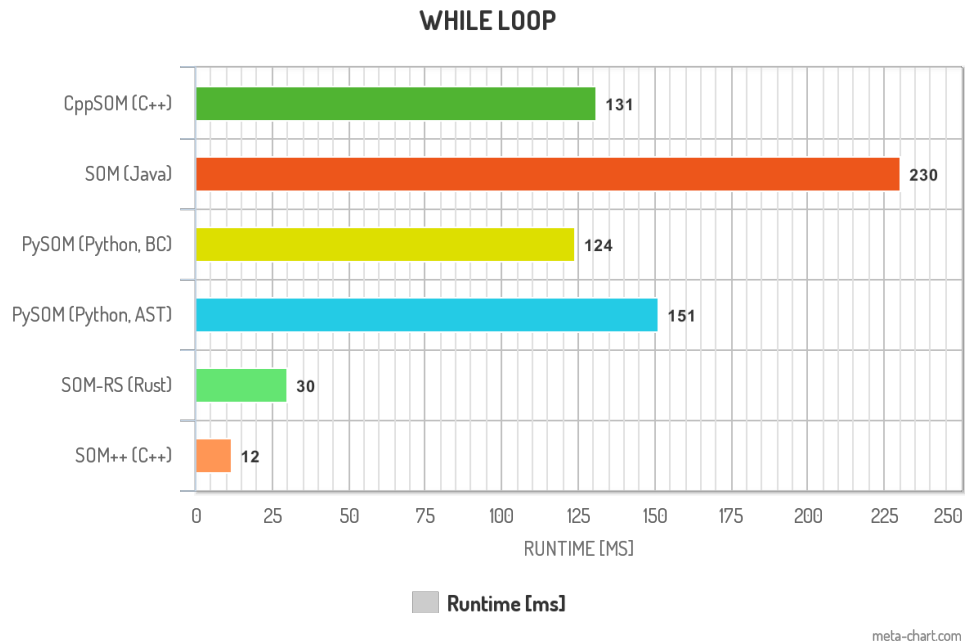


Figure 4.7: Execution times for different SOM implementations for code from [4.6](#).

When examining the Java implementation, I found that the `restart` primitive method for block is implemented very similarly to mine. From the code snippet () it is safe to assume that no extra message sends or any different actions are performed.

```
public void invoke(Frame frame, Interpreter interpreter) {
    frame.setBytecodeIndex(0);
    frame.resetStackPointer();
}
```

Figure 4.8: Invocation of `restart` method in Java implementation of SOM (taken from [\[1\]](#)).

4.3 Non local returns

In the third test, I take a look at performance of non local returns. The code that is run in all implementation is on figure [4.9](#).

4. PERFORMANCE

```
NLReturn = (  
  
    first: a = ( self second: a )  
    second: a = (self third: a )  
    third: a = ( a value )  
  
    nlr = (  
        self first: [ ↑1 ]  
    )  
  
    run = (  
        | x |  
        x := 0.  
        [ x < 200 ] whileTrue: [  
            x := x + self nlr  
        ].  
        x println  
    )  
)
```

Figure 4.9: SOM source code for third performance benchmark – non local return evaluation.

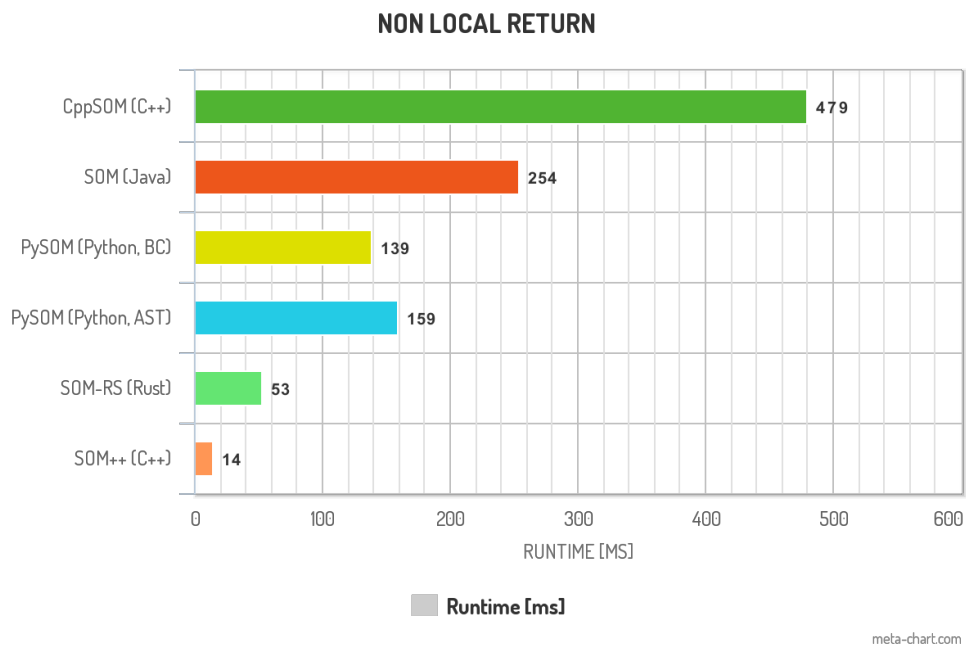


Figure 4.10: Execution times for different SOM implementations for code from [4.9](#).

Conclusion

The thesis provided an overview of the syntax and core principles of the Simple Object Machine programming language. It put the language in the context of Smalltalk, on which SOM is based.

It provides a basic implementation of a virtual machine to execute SOM programs. The solution parses the source files, compiles them into bytecode and is then able to load the compiled bytecode and execute it. The thesis provides the description of abstract syntax tree design, as well as bytecode design and its implementation.

The project provides basic functionality, compared to existing SOM implementations. Especially from the performance point of view, there is a lot of potential for further work and improvement. Optimizations over AST and bytecode or more advanced garbage collection algorithms are the first things that come to mind.

On the other hand, this relatively small implementation offers an advantage of providing an easy to get into playground for people that want to experiment with runtime systems. The clear separation of every step of compilation/execution allows faster understanding of the functioning, as well as ability to easily modify or completely change the parser, compiler or interpreter.

Future work

The provided solution offers basic functionality to execute SOM programs. It can be further built upon and improved in several ways. As demonstrated in the benchmarking section, the performance is lackluster compared to existing solutions. I will propose a few solutions, that could be implemented in the future to remedy this.

Provided solution utilises a custom design bytecode, based on SOM/Smalltalk bytecode. An alternate compiler could be built, utilising an existing bytecode with existing VM implementation.

CONCLUSION

Implemented garbage collection algorithm is a fairly basic one. An interesting addition would be the implementation of more advanced GC algorithms (such as reference counting or generational collection). Other possible improvements to GC could be the implementation of concurrent mark and sweep to shorten the execution pauses.

Bibliography

- [1] Marr, S. SOM. [cit. 2021-5-5]. Available from: <https://github.com/SOM-st/som-java>
- [2] SOM. SOM: A minimal Smalltalk for teaching and research on Virtual Machines. 2020, [cit. 2020-11-17]. Available from: <https://som-st.github.io/>
- [3] Lovejoy, A. L. Smalltalk: Getting The Message. 2007, [cit. 2021-1-19]. Available from: <http://devrel.zoomquiet.top/data/20080627141054/index.html>
- [4] Ducasse, S.; Chloupis, D.; et al. Pharo By Example 5. 2017.
- [5] Bera, C. Blocks: A Detailed Analysis. [cit. 2021-4-30]. Available from: <http://pharobooks.gforge.inria.fr/PharoByExampleTwo-Eng/latest/Block.pdf>
- [6] Boersma, E. Memory leak detection - How to find, eliminate, and avoid. January 2020, [cit. 2020-11-5]. Available from: <https://raygun.com/blog/memory-leak-detection/>
- [7] Agarwal, C. Mark-and-Sweep: Garbage Collection Algorithm. [cit. 2021-5-1]. Available from: <https://www.geeksforgeeks.org/mark-and-sweep-garbage-collection-algorithm/>
- [8] Tahir, N. An introduction to language lexing and parsing with ANTLR. 2017, cit. [2021-5-6]. Available from: <https://willowtreeapps.com/ideas/an-introduction-to-language-lexing-and-parsing-with-antlr>

Acronyms

- AST** Abstract syntax tree
- GC** Garbage collector
- OOP** Object Oriented Programming
- SOM** Simple Object Machine
- VM** Virtual machine

Contents of enclosed CD

readme.txt	the file with CD contents description
src	the directory of source codes
├─ antlr_grammar	ANTLR grammar definitions
├─ cmake	CMake files for ANTLR building
├─ core_lib	SOM core library implementation
├─ examples	examples of SOM source code
├─ src	implementation source files
thesis	the directory of L ^A T _E X source codes of the thesis
├─ media	the thesis figures directory
├─ *.tex	the L ^A T _E X source code files of the thesis
text	the thesis text directory
├─ thesis.pdf	the Diploma thesis in PDF format