

giaco.ml

an interpreted EDSL written in OCAML

Mario A. Barbara
Andrea Cancellieri

1 Introduction

giaco.ml is an imperative (but also functional) EDSL interpreter written in OCAML.

The interpreter is capable of handling expressions, commands and declarations all within a program.

A program can be parsed into valid code from a file/string thanks to reflection.

Static taint analysis can be performed, guards have been also added to the `Reflect` command as a form of protection, similar to perl.

The main execution functions are: `interpret` (or `interpret'` which returns also the declarations' output), `reflect` and `taint_analysis` which all operate on a `Prog(ds,cs)` construct.

2 Design Choices

The interpreter works with 3 domains:

generic this domain is shared amongst the syntax (external) and semantic (internal) domains

syntax this domain contains all the domains accessible to the user. it features 3 sub-domains

expressions these are entities which may be mapped directly to a value

commands these are entities which allow modifying the global state of the program, though allocating new memory is not allowed

declarations these are entities which allow allocation of new memory in the global state

semantic this domain is used for representing the internal state of the interpreter.

evaluable values these are values that may be directly evaluated from an expressions. they represent the base internal type for the interpreter. Basic allowed types are Integer, Boolean, Float, String and Subprograms.

environment this is a mapping representing internal non-mutable state. It maps variables (identifiers) to values.

store this is a mapping representing internal mutable state. It allows aliasing of 2 identifiers for one value. It maps addresses (pointers or locations) to values.

denotable values these are values contained in the environment. Basic allowed types are Integer, Boolean, Float, String and Locations.

mutable memory values these are values contained in the store. Basic allowed types are Integer, Boolean, Float, String.

The interpreter itself is implemented by 3 main functions:

eval this function basically maps expressions to evaluable values

cval this function maps commands to a changed internal state (store)

dval this function maps declarations to a new internal state (environment + store)

Some important choices:

- memory is viewed as a mapping. it is therefore impossible to modify in ocaml. changing the type signature of the interpreter would allow this
- because memory is a mapping, new locations are simulated by randomly selecting a new location based on given store size
- regular expressions are used in reflection, so the `Str` module must be available
- fixed points are calculated through ocaml's `function` and `let rec` operations
- to ease debugging, the interpreter is not compiled but interpreted in an ocaml shell with the `#use "file.ml"` construct
- since the program is interpreted, `let and` between multiple files is not available, this has been fixed by using mutable state pointers to future functions. These pointers will be updated when the function is created.
- many utility functions have been created to ease debugging. most notable are `emptyenv emptystore tenv' tstore' tnew'` which simplify dealing with memory
- due to the non-recursive nature of `Procedure` and `Call`, they have been omitted from reflection and taint-analysis. The syntax restriction is due to the difficulty of implementing a parser for the complete OCAML list syntax, and the static analysis of commands that cannot be recursively unpacked (one lies in the *expressions*, the other in the *commands*)

3 Usage

the interpreter can be loaded up with

```
>> #use "giaco.ml";;
```

and subsequent testing may be performed:

```
>> (* have a functional sample: *)
>> Plus(Int(3), Int(4)) in eval e emptyenv emptystore;;
- : eval = EInt 7

>> (* now let's try a complete program *)
>> let e = Lambda("x", Multiply(Val("global"), Var("x"))) in
>> let d = New("global", Str("yes")) in
>> let c = Assign("global", Apply(e, Int(5))) in
>> let program = Prog(d,c) in
>> let result_environment, result_store = interpret' program emptyenv emptystore in
>> eval (Val("global")) result_environment result_store;;
- : eval = EStr "yesyesyesyesyes"
```

The most useful interpreter functions are:

`eval: expr -> env -> store -> eval` evaluates expressions

`cval: com -> env -> store -> store` converts commands into a modified mutable memory

`dval: dec -> env -> store -> env*store` allows extending mutable memory

`interpret: prog -> env -> store -> store` combines all of the above. cannot return expression values like `eval`, though, as we do not have print functionality

interpret': `prog -> env -> store -> env*store` just like **interpret** but return the last evaluated environment as well, which allows for further analysis of the output

emptyenv and **emptystore** already initialized empty environment and store

env' and **store'** allow extending environments and stores outside of syntax. **new'** combines this and allows doing something like the **New(...)** command outside of syntax.

ereffect and **creffect** and **dreflect** and **reflect** reflection of expressions, commands, declarations, full blown programs.

etaint and **ctaint** and **dtaint** and **taint_anlaysia** static taint analysis of expressions, commands, declarations and full blown programs.

3.1 EXPRESSIONS

EXAMPLE	DESCRIPTION
<code>Int(3)</code>	basic integer
<code>Str("hello world")</code>	basic ASCII string
<code>Bool(true)</code>	basic boolean
<code>Float(4.5)</code>	basic float
<code>Lambda("x", <exp containing x>)</code>	typical function
<code>RecLambda("f", "x", <exp containing f and x>)</code>	typical recursive function
<code>Rec("f", Lambda(...))</code>	just another way to define recursive lambdas
<code>Proc(["x";"y";"z";...], Block(...))</code>	this is a procedure, check the commands section
<code>IfThenElse(Bool(true), .., ..)</code>	control flow element
<code>Var("x")</code>	this is a way to retrieve an immutable variable's content
<code>LetIn("x", e1, e2)</code>	this is a way to nest functional blocks and scopes
<code>Val("x")</code>	this is a way to retrieve a mutable variable's content
<code>Plus(e1, e2)</code>	plus function, applies to: Int, Str, Float
<code>Multiply(e1, e2)</code>	multiply function, appliest to: Int, Str, Float
<code>Apply(e1, e2)</code>	typical function application, e1 is of type: Lambda, RecLambda, Rec
<code>Equals(e1, e2)</code>	like C's ==
<code>Greater(e1, e2)</code>	like C's >
<code>Not(e)</code>	like C's !
<code>Or(e1, e2)</code>	like C's
<code>And(e1, e2)</code>	like C's &&
<code>Len(Str(...))</code>	gets the length of a St
<code>Sub(Str(...), i, j)</code>	gets a substring. i and j of type Int.
<code>Lower(Str(..))</code>	reduces a string to lowercase, like Python's <code>lower()</code>
<code>Upper(Str(...))</code>	reduces a string to uppercase, like Python's <code>upper()</code>
<code>Trim(Str(...))</code>	trims whitespace from a string, like Python's <code>s.trim()</code>
<code>Replace(<string to be replace>,<replacer string>,<string>)</code>	replaces a string with another string in a string, like Python's <code>s.replace()</code>

3.2 COMMANDS

EXAMPLE	DESCRIPTION
Assign("x", e)	this changes the mutable value for the variable "x". e is an expression
Block(d, c)	this is an imperative block with nested scope. d is a declaration, see its section for more detail
Call(p, [e1;e2;e3;..])	this is an application of an imperative procedure. p is of type Proc (check the expressions section)
While(e, c)	like C's <code>while(e){c}</code> , e is an expression and c a command
CIfThen(e, c)	like C's <code>if(e){c}</code>
CIfThenElse(e, c1, c2)	like C's <code>if(e){c1}else{c2}</code>
CSeq(c1, c2)	like C's <code>;</code> it allows concatenation of commands
CSkip	like C's <code>void</code> and Python's <code>pass</code> , it does nothing
Reflect(Str(...))	reflection, see the reflection section

3.3 DECLARATIONS

EXAMPLE	DESCRIPTION
New("x", e)	this allocates a new mutable variable of value e (an expression)
DSeq(d1, d2)	allows concatenation of declarations
DSkip	does nothing

3.4 PROGRAM

NOTE :: to use the functional language only, please use expressions along with `eval`. Programs are only imperative.

EXAMPLE	DESCRIPTION
Prog(<declaration>, <command>)	handy syntax that puts it all together

4 String extension

A few functions have been added to deal with the domain of strings. Functions such as these are taken from the Python language, which has a **very** extensive and popular standard library.

- length comparison (**Greater**)
- concatenation (**Plus** has been extended to allow this)
- substring (**Sub**)
- repetition (**Multiply** has been extended to allow this)
- length (**Len**)
- lowercase (**Lower**)
- uppercase (**Upper**)
- trim (**Trim**), trims all whitespace
- Replace (**Replace**)

Check the examples section for some examples

5 Reflection extension

Reflection consists of allowing any string to be evaluated by the interpreter on the fly. In Python this is akin to the `eval` function. This is also the most essential step to having a good interpreter: the interactive console for Python, one of the most popular interpreted languages, is often called **R.E.P.L.** (Read Eval Print Loop).

The syntax of *giaco.ml* has been extended with the **Reflect** command, which allows on the fly "evaluation" of commands. Unfortunately our language's command syntax is recursive, and furthermore the **ClfThenElse** command uses expressions as boolean conditions, which are also recursive. Therefore, a full blown parser needed to be built to give a string some depth (such as that of an AST).

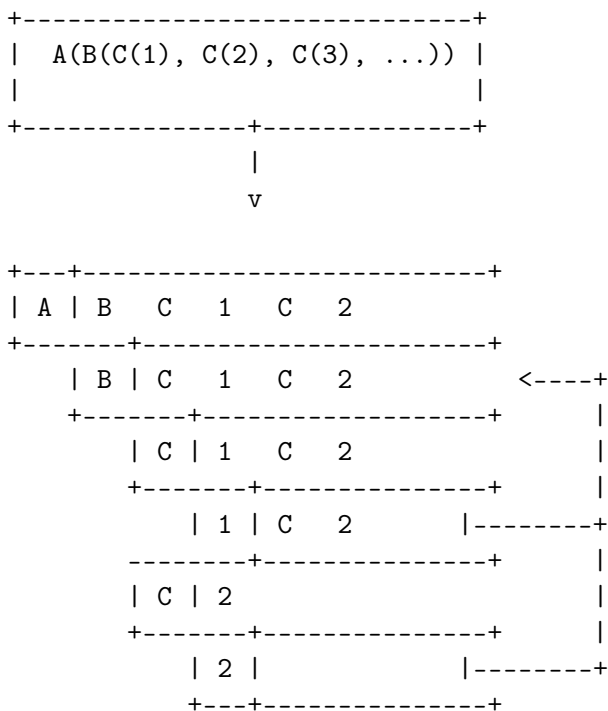
The reasoning is as follows:

1. a function called `next_unit` is charged with grabbing the first word up until a (or) or , or multiple consecutive repetitions.
2. to get the command to match against, `next_unit` is called on the string and the result is matched against some constants, taken from the language's syntax
3. to get a command's arguments (which may be recursive and contain any amount of () , , caution must be taken to correctly identify the argument boundaries, which are all separated by a , comma. 2 options are given:

iterative by counting the number of open parentheses matched thus far, and decreasing each time a closed parentheses is found, it is possible to correctly identify the recursive structure of the syntax.

recursive (but faster) since we know the amount of parameter each command needs, it is simply required to recursively reflect upon the arguments' string as many times as needed. Of course, each time a command is consumed, it shall return the arguments' string, so as to allow its father to continue looking for arguments.

Our interpreter implements the recursive and faster technique. Here is a simple diagram drawing to illustrate the flow of this technique:



6 Taint-Analysis extension

Static taint analysis consists of understanding how much damage some unsafe elements (of undefined value but defined nature) will yield. A classic example is an unsanitized input on a HTML form, which may result in an SQL

Injection attack and damage your company's most valuable assets.

In our simple language, we have no operations that deal with the outside world (yet). We are therefore forced to ask the user to label some variables in the environment and store as **Clean** or **Dirty**. Afterwards, we will analyze a program and check the **Taint** for every possible variable assignment. The semantic domains have been revisited, allowing memory (environment and store) to only contain tainted values (or store locations, in the environment's case).

Our analysis is based on 2 simple concepts:

pure evaluation a **tor** function will take 2 taints and return **Dirty** if one of them is as well, otherwise **Clean**. This process can be lazy.

- All constants are **Clean**
- If a function is involved (such as a **Lambda**) then the formal parameters are identified as **Clean** (as they cannot be expressions) and then the body is analyzed. If the body is clean, the function is clean
- A function application requires a **tor** amongst the analysis of the function itself and the passed parameter.
- If a condition is involved, then 2 outputs are possible. If the condition is **Dirty**, that means the attacker may choose either output and (regardless of the output's default taint) will result in a **Dirty** value. If the condition is **Clean**, then either output may occur during execution, so they must be passed to **tor**.

imperative state change all possible assignments in a command are gathered. Only the latest possible assignments matter (if i set **x** to **Dirty** and then **Clean** it is **Clean**).

- Afterwards, we check whether 2 branches are possible: if they are, a **tor** function must be applied to all assignments of same key, merging the 2 branches.
- If the branches are subject to a condition (such as in a **CIfThenElse**) then a **Dirty** condition will mean an attacker may choose amongst any of the 2 branches, therefore dirtying all assignments of shared key (amongst the 2 branches). If the condition is **Clean**, then the normal merge has already evaluated taint with **tor**.

7 Examples

check test.ml for some code examples.

7.1 Numbers

INPUT	OUTPUT
Int(5)	EInt 5
Float(133.7)	EFloat 133.7
Plus(Int(1), Int(2))	EInt 3
Multiply(Float(2.5),Float(10.0))	EFloat 25
Greater(Int(3),Int(5))	EBool false

7.2 Booleans

INPUT	OUTPUT
Bool(true)	EBool true
Not(Bool(true))	EBool false
And(Equals(Float(4.5),Float(4.6)),Equals(Float(0.1),Float(0.1)))	EBool false
Or(Equals(Float(4.5),Float(4.6)),Equals(Float(0.1),Float(0.1)))	EBool true

7.3 Strings

INPUT	OUTPUT
Str("hello world")	EStr "hello world"
Plus(Str("hello "),Str("world!"))	EStr "hello world!"
Multiply(Str("abc"),Int(10))	EStr "abcabcabcabcabcabcabcabcabc"
Len(Multiply(Str("abc"),Int(10)))	EInt 30
Greater(Str("two"),Str("three"))	EBool false
Sub(Str("threeeeeeee"),Int(2),Int(10))	EStr "reeeeeeeee"
Upper(Str("im so lonely"))	EStr "IM SO LONELY"
Lower(Upper(Str("im so lonely")))	EStr "im so lonely"
Trim(Str(" italia "))	EStr "italia"
Replace(Str("hello"),Str("goodbye"),Str("hello world!"))	EStr "goodbye world!"

7.4 Functional Control Flow

INPUT	OUTPUT
IfThenElse(Bool(true), Int(1337), Str("i am"))	EInt 1337
IfThenElse(Not(Greater(Str("bob"),Str("mouse"))),Str("ciao mondo"),Int(5))	EStr "ciao mondo"

7.5 Functional Blocks

INPUT	OUTPUT
Var("x")	EInt 20
xxx = LetIn("a",Int(3),Multiply(Var("a"),Var("a")))	EInt 9
LetIn("a",Int(5),(LetIn("b",xxx,LetIn("c",Int(6),Plus(Var("a"),Plus(Var("b"),Var("c"))))))))	EInt 20

7.6 Functional Subprograms

INPUT	OUTPUT
Apply(Lambda("x", Plus(Var("x"), Int(1))), Int(99))	EInt 100
Apply(RecLambda("fact", "x", IfThenElse(Equals(Var("x"), Int(0)), Int(1), Multiply(Var("x"), Apply(Var("fact"), Plus(Var("x"), Int(-1)))))), Int(10))	EInt 3628800

7.7 Imperative State Change

INPUT	VARIABLE OUTPUT
Val("y")	EInt 10
Assign("y", Plus(Val("y"), Val("y")))	EInt 20

7.8 Imperative Control Flow

INPUT	VARIABLE OUTPUT
Val("y"), Val("z")	EInt 10, EInt 0
CIfThenElse(Not(Equals(Val("y"),Int(11))), Assign("y", Int(50)))	EInt 50, EInt 0
While(Not(Equals(Val("y"), Int(100))), CSeq(Assign("y", Plus(Val("y"), Int(1))), Assign("z", Plus(Val("z"), Int(1)))))	EInt 50, EInt 50

7.9 Imperative Blocks

INPUT	OUTPUT
Val("y"), Val("z")	EInt 10, EInt 0
Block(New("z", Int(1000)), Assign("y", Plus(Val "y", Val "z")))	EInt 1010, EInt 0

7.10 Imperative Subprograms

INPUT	OUTPUT
Val("y"), Val("z")	EInt 10, EInt 0
f = Proc(["z"], Block(DSkip, Assign("y", Val("z")))) in Call (Val "f", [Val "z"])	EInt 0, EInt 0

7.11 Declarations

INPUT	OUTPUT
Val("y"), Val("z")	Failure 'y' not in environment, Failure 'z' not in environment
DSeq(New("y", Int(10)), New("z", Int(0)))	EInt 10, EInt 0

7.12 Reflection

INPUT	OUTPUT
Val("y")	EInt 10
ereffect ("Plus(Plus(Int(1), Int(2)), Plus(Int(3), Int(4)))")	EInt 10
Reflect(Str("Assign(\"y\", Int(5))"))	EInt 5

7.13 Taint Analysis

"dirty" is Dirty, "clean" is Clean

INPUT	VALUE
"dirty"	Dirty
"clean"	Clean
e	Equals(Plus(Val("x"), Val("y")), Int(6))
assign1	CSeq(Assign("x", Val("dirty")), Assign("y", Val("clean")))
assign2	CSeq(Assign("x", Val("clean")), Assign("y", Val("dirty")))
d	DSeq(New("x", Val("dirty")), New("y", Val("clean")))
c	CIfThenElse(e, assign1, assign2)

INPUT	OUTPUT
taint_analysis Prog(d,c)	[("clean", TLoc 15n); ("dirty", TLoc 27n); ("x", TLoc 76n); ("y", TLoc 41n)] [(15n, Clean); (27n, Dirty); (41n, Dirty); (76n, Dirty)]

8 Noted Bugs:

- reflection on the **Reflect** command is supposedly allowed (since it is implemented recursively) but OCAML's escape of strings within strings is weird and needs some fixing...