

Minicaml, a purely functional, didactical programming language

WORK IN PROGRESS DRAFT

Alessandro Cheli
Course taught by Prof. Gianluigi Ferrari
and Prof. Francesca Levi

January 15, 2020

Abstract

minicaml is a dynamically typed and purely functional interpreted programming language. It is based on the Professor Gianluigi Ferrari and Professor Francesca Levi's minicaml, an evaluation example to show students attending the Programming 2 course at the University of Pisa how interpreters work. It is an interpreted language heavily inspired from the OCaml, Haskell and Scheme languages, with static (lexical scoping), eager and lazy evaluation and a didactical REPL that shows each AST expression and each evaluation step.

1 REPL and command line interface

1.1 Installation

minicaml is available in the opam 2.0 repository. (<https://opam.ocaml.org/>). The easiest way to install minicaml is with the OCaml package manager **opam**. To do so, please check that you have a version of opam $\geq 2.0.0$ and run:

```
opam install minicaml
```

Alternatively, **minicaml** can be installed from source by downloading the source code repository and building it manually. **minicaml** has been tested only on Linux and macOS systems. It has not been tested yet on Windows and BSD derived systems.

```
# download the source code
git clone https://github.com/0x0f0f0f/minicaml
# cd into the source code directory
cd minicaml
# install dependencies
opam install ANSITerminal dune ppx_deriving menhir \
  cmdliner alcotest bisect_ppx ocamline
# compile
make
# test
make test
# execute
make run
# install
make install
```

2 Syntax and Parser

Lexing is achieved with `ocamllex`, the default tool for generating scanners in OCaml. The parser is realized with the **Menhir** parser generator, and is documented using **Obelisk**, which generates a clean text file containing the language grammar, available in Appendix A.

3 Purity Inference

An important feature of the `minicaml` language is the purity inference algorithm, which is performed statically on expressions before evaluation. It is an interpretation of expressions over the domain of purity, meant to prevent side effects by signal an error if they are contained inside the programs written in the language. Expressions are tagged by the algorithm with the `Pure`, `Impure` and `Numerical` labels. An `Impure` expression is an expression that contains calls to primitives that perform I/O operations, mutable variables and/or imperative style assignments. A `Numerical` expression is an expression where only numerical operations are performed. `Pure` expressions are those which do not fall into the previous two categories. To perform impure side effects, the programmer has two constructs available called **purity blocks**. By default, the evaluator is in an `Uncertain` context, which means that it will not allow side effects to be carried on by evaluation, but will allow evaluating purity blocks that change the currently allowed purity context. The `impure` statement takes an expression (the block) and evaluates it in a context where the allowed purity is `Impure`, so that side effects may be performed. The other construct available, the `pure` statement, takes an expression and enforces a `Pure` context, meaning that side effects and nested impure blocks will not be allowed inside of the expression.

4 AST Optimization

After purity inference is performed, and before evaluation, AST expressions are analyzed and optimized by an optimizer function that is recursively called over the tree that is representing the expression. The optimizer simplifies expressions which result is known and therefore does not need to be evaluated. For example, it is known that $5 + 3 \equiv 8$ and $\text{true} \ \&\& \ (\text{true} \ || \ (\text{false} \ \&\& \ \text{false})) \equiv \text{true}$. When a programmer writes a program, she or he may not want to do all the simple calculations before writing the program in which they appear in, we rely on machines to simplify those processes. Reducing constants before evaluation may seem unnecessary when writing a small program, but they do take away computation time, and if they appear inside of loops, it is a wise choice to simplify those constant expressions whose result is already known before it is calculated in all the loop iterations. It is also necessary in optimizing programs before compilation. The optimizer, by now, reduces operations between constants and `if` statements whose guard is always `true` (or `false`). To achieve minimization to an unreducible form, optimizer calls are repeated until it produces an output equal to its input; this way, we get a tree representing an expression that cannot be optimized again. This process is fairly easy:

```
let rec iterate_optimizer e =
```

```
let oe = optimize e in
if oe = e then e (* Bottoms out *)
else iterate_optimizer oe
```

Boolean operations are reduced using laws from the propositional calculus, such as DeMorgan's law, complement, absorption and other trivial ones.

5 Types

6 Evaluation

`minicaml`'s evaluator is heavily inspired by the Metacircular Evaluator defined in the highly acclaimed textbook *Structure and Interpretation of Computer Programs* [1]

6.1 Operational Semantics

Note. The letter e denotes an environment.

The symbol $_$ is used whenever a value exists but is content is irrelevant to the semantical rule, or cannot be determined and therefore is discarded.

Dictionaries

Creation
$\langle e, d \rangle \Rightarrow$ $\{(k, v) \in d \mid \forall i, j \in \mathbb{N} \wedge i, j \in [1, d] \wedge i \neq j \text{ such that } k_i \neq k_j\}$
Insertion
$\langle e, d \rangle \Rightarrow ed$ $\langle e, k \in ed \rangle \Rightarrow \text{false}$ $\langle e, \text{insert } k \ v \ d \rangle \Rightarrow ed \cup (k, v)$
$\langle e, d \rangle \Rightarrow ed$ $\langle e, k \in ed \rangle \Rightarrow \text{true}$ $\langle e, \text{insert } k \ v \ d \rangle \Rightarrow ed \setminus \{(k, _)\} \cup (k, v)$
Deletion
$\langle e, d \rangle \Rightarrow ed$ $\langle e, k \in ed \rangle \Rightarrow \text{true}$ $\langle e, \text{remove } k \ d \rangle \Rightarrow ed \setminus \{(k, v)\}$
$\langle e, d \rangle \Rightarrow ed$ $\langle e, k \in ed \rangle \Rightarrow \text{false}$ $\langle e, \text{remove } k \ d \rangle \Rightarrow \text{error}$
Contains key
$\langle e, d \rangle \Rightarrow ed$ $\langle e, \text{haskey } k \ d \rangle \Rightarrow (k, _) \in ed$

Retreive a value

$$\begin{array}{c} < e, d > \Rightarrow ed \\ < e, k \in ed > \Rightarrow \text{true} \\ \hline < e, \text{getkey } k \ d > \Rightarrow v \end{array}$$

$$\begin{array}{c} < e, d > \Rightarrow ed \\ < e, k \in ed > \Rightarrow \text{false} \\ \hline < e, \text{getkey } k \ d > \Rightarrow \text{error} \end{array}$$

Filter by keys

$$\begin{array}{c} < e, d > \Rightarrow ed \\ < e, ks > \Rightarrow \{k_1, \dots, k_n\} \\ \hline < e, \text{filterkeys } ks \ d > \Rightarrow \{(k, v) \in ed \mid (k \in ks)\} \end{array}$$

Map

$$\begin{array}{c} < e, d > \Rightarrow ed \\ < e, f > \Rightarrow \lambda(x) \\ \hline < e, \text{map } f \ d > \Rightarrow \{(k, \lambda(v)) \mid (k, v) \in ed\} \end{array}$$

Fold Left

$$\begin{array}{c} < e, f > \Rightarrow \lambda(x, y) \\ < e, d > \Rightarrow \{(k_1, v_1), \dots, (k_n, v_n)\} \\ \hline < e, \text{foldl } f \ a \ d > \Rightarrow \lambda(\dots \lambda(\lambda(a, v_1), v_2), \dots, v_n) \end{array}$$

7 Tests

Unit testing is extensively performed using the alcotest testing framework. Code coverage is provided by the `bisect_ppx` library which yields an HTML page containing the coverage percentage when unit tests are run by the dune build system. After each commit is pushed to the remote version control repository on Github, the package is built and tests are run thanks to the Travis Continuous Integration system.

8 Thanks to

- Prof. Gian-Luigi Ferrari and Francesca Levi for teaching us how to project and develop interpreters in OCaml
- Antonio DeLucreziis for helping me implement lazy evaluation.
- Prof. Alessandro Berarducci for helping me study lambda calculus in depth.
- Giorgio Mossa for helping me polish the lambda-closure mechanism.

A Parsing Grammar

```
<file> ::= EOF
        | <ast_expr> EOF
        | <ast_expr> SEMISEMI <file>
        | <def> SEMISEMI <file>
        | <def> EOF
        | <directive> <file>

<toplevel> ::= <directive>
            | [<ast_expr> (SEMI <ast_expr>)*] [SEMISEMI] EOF
            | <def> [SEMISEMI] EOF
            | <ast_expr> [SEMISEMI] EOF

<assignment> ::= SYMBOL EQUAL <ast_expr>
              | LAZY SYMBOL EQUAL <ast_expr>

<def> ::= LET [<assignment> (AND <assignment>)*]

<directive> ::= DIRECTIVE STRING
             | DIRECTIVE INTEGER
             | DIRECTIVE UNIT

<ast_expr> ::= <ast_app_expr>
            | LPAREN <ast_expr> (SEMI <ast_expr>)* RPAREN
            | <ast_expr> CONS <ast_expr>
            | NOT <ast_expr>
            | <ast_expr> ATSIGN <ast_expr>
            | <ast_expr> CONCAT <ast_expr>
            | <ast_expr> LAND <ast_expr>
            | <ast_expr> OR <ast_expr>
            | <ast_expr> PLUS <ast_expr>
            | <ast_expr> MINUS <ast_expr>
            | <ast_expr> TIMES <ast_expr>
            | <ast_expr> DIV <ast_expr>
            | <ast_expr> EQUAL <ast_expr>
            | <ast_expr> DIFFER <ast_expr>
            | <ast_expr> GREATER <ast_expr>
            | <ast_expr> LESS <ast_expr>
            | <ast_expr> GREATEREQUAL <ast_expr>
            | <ast_expr> LESSEQUAL <ast_expr>
            | IF <ast_expr> THEN <ast_expr> ELSE <ast_expr>
            | <def> IN <ast_expr>
            | LAMBDA SYMBOL+ LARROW <ast_expr>
            | <ast_expr> COMPOSE <ast_expr>
            | <ast_expr> PIPE <ast_expr>

<ast_app_expr> ::= <ast_simple_expr>+

<ast_simple_expr> ::= SYMBOL
                  | UNIT
                  | DOLLAR <ast_expr>
                  | LPAREN <ast_expr> RPAREN
                  | <ast_simple_expr> COLON SYMBOL
                  | PURE <ast_expr>
                  | IMPURE <ast_expr>
                  | LSQUARE [<ast_expr> (SEMI <ast_expr>)*] RSQUARE
                  | LBRACKET [<assignment> (SEMI <assignment>)*] RBRACKET
                  | BOOLEAN
                  | STRING
                  | INTEGER
                  | FLOAT
                  | FLOAT CPLUS FLOAT
                  | FLOAT CMIN FLOAT
```

References

- [1] Harold Abelson, Gerald Jay Sussman, and with Julie Sussman. *Structure and Interpretation of Computer Programs*. MIT Press/McGraw-Hill, Cambridge, 2nd edition, 1996.