

Experience Report: OCaml for an industrial-strength static analysis framework

Pascal Cuoq* Julien Signoles
CEA LIST, Software Reliability Labs,
91191 Gif-sur-Yvette Cedex, France
First.Last@cea.fr

with Patrick Baudin, Richard Bonichon,
Géraud Canet, Loïc Correnson, Benjamin
Monate, Virgile Prevosto, Armand Puccetti

Abstract

This experience report describes the choice of OCaml as the implementation language for Frama-C, a framework for the static analysis of C programs. OCaml became the implementation language for Frama-C because it is expressive. Most of the reasons listed in the remaining of this article are secondary reasons, features which are not specific to OCaml (modularity, availability of a C parser, control over the use of resources...) but could have prevented the use of OCaml for this project if they had been missing.

Categories and Subject Descriptors D1.1 [*Programming techniques*]: Applicative (Functional) Programming

General Terms Design, Languages

Keywords OCaml, software architecture, plug-ins, static analysis

1. Introduction

Frama-C is a framework that allows static analyzers, implemented as plug-ins, to collaborate towards the study of a C program. Although it is distributed as Open Source, Frama-C is very much an industrial project, both in the size it has already reached and in its intended use for the certification, quality assurance, and reverse-engineering of industrial code. Frama-C is written in OCaml, and this article reports on the most noticeable consequences of this choice on the human (section 2) and technical (section 3) levels, as well as providing an overview of the implementation of Frama-C (appendix A).

2. Human Context

Frama-C is developed collaboratively between the ProVal team (a joint laboratory of INRIA Saclay Île-de-France and LRI) and CEA LIST. This article describes our (CEA LIST) own analysis of this collaborative development. The Open Source nature of the software and other partnerships involving CEA LIST mean that Frama-C is in fact developed at three different sites, by around ten full-time programmers, with infrequent inter-site eye-to-eye meetings.

*This work has been supported by the french RNTL project CAT ANR05RNTL00301

2.1 Recruiting OCaml programmers

It is typical for an article of this nature to include a few words to the effect that it is harder to find people who can program in functional language Y (Minsky 2007; Nanavati 2008) than in C++, sometimes nuanced by more words pointing out that this is balanced by the higher quality of Y candidates. The first proposition did not apply for us in the case of Frama-C and OCaml. CEA LIST is an applied research laboratory that recruits almost exclusively PhDs. When the choice is restricted to candidates with a PhD in the field of formal methods, it is not harder to find a candidate with the motivation to program in OCaml than in C++.

2.2 Objectives of the Frama-C project

Although it is developed by research institutes, Frama-C tries to fulfill focused, specific needs expressed by industrial partners. It aims past the R&D departments and into the hands of the engineers who develop embedded code in any industry with criticality issues¹. Frama-C is structured as a kernel to which different analysis plug-ins are connected. It is composed as a whole of 100 to 200 kloc. All this OCaml code provides a wide range of functionalities, but the plug-ins do not just sit side by side. It is more accurate to think of them as built on top of each other. To give an example, a value analysis plug-in computes supersets of possible values for the variables of the program, indexed by statement of the original program. Unions, structs, arrays and casts thereof are handled with the precision necessary for embedded code (Cuoq 2008). These values, especially the values of pointers and expressions used as indices in arrays, are used by another plug-in to compute synthetic functional dependencies between the inputs and the outputs of each analyzed function. These synthetic dependencies are in turn used by a slicer plug-in which produces simplified, compilable C programs that are guaranteed to be equivalent to the original for the slicing criterion. And the building blocks of this slicer are used by one of Frama-C's most sophisticated plug-ins to date, a security-aware slicer that preserves the confidentiality of information: the (functional) confidentiality is guaranteed to be exactly the same in the sliced program as in the original program (Monate and Signoles 2008). This means that it is safe to study functional confidentiality on the (smaller) sliced program, for instance for a security audit of the source code. This would not be possible with a traditional slicer, because a traditional slicer might remove information leaks — in particular, a malicious programmer could insert information leaks that he knows the traditional slicer used for the audit will remove.

¹ or even without criticality issues, but so far the industrial interest has come from critical embedded systems

2.3 History of Frama-C

At the Software Reliability Labs, Frama-C was initially seen as a major evolution to Caveat (Baudin et al. 2002; Randimbivololona et al. 1999; Delmas et al. 2008), a software verification tool for C programs based on Hoare logic. Caveat is used by Airbus for part of the verification of part of the software embedded in the A380. As the DO-178B standard mandates, Caveat has been qualified by Airbus as a verification tool to be used for the certification of this particular software. Caveat is supported by the Software Reliability Labs but new ideas are now developed within Frama-C.

OCaml was pushed as the implementation language to choose for Frama-C by the new hires, but the actual reason it was accepted is that OCaml was not completely unheard of to the senior researchers there. Indeed, OCaml was already used in the (predominantly C++) Caveat project as the scripting language that allows an interactive validation process to be re-played in batch mode.

3. Technical context

3.1 Expressivity

OCaml's expressivity was crucial in the adoption phase of the language. An initial one-person internal prototype was able to produce results that convinced management to extend the experiment to two programmers. Eventually, this prototype was able to persuade industrial users to get involved in a national project.

To be precise, industrial partners agreed to be part of the project because of their previous experiences with the providers of static analyzers in the project. Some of the tools that they were familiar with were written in OCaml (Caduceus by the ProVal team), and some of them weren't (Caveat). The time it takes to build these relationships should not be underestimated, and we are not saying that the choice of any programming language can shorten it.

However, the progress made by Frama-C after the project had started, which can at least partly be attributed to OCaml, convinced the industrial participants to become involved beyond expectations. At each phase of the bootstrap process, OCaml's expressivity was important in quickly writing the proof-of-concept that took the project to the next stage.

3.2 Control over the use of resources

One of Frama-C's first plug-ins was a value analysis based on abstract interpretation. This plug-in computes supersets of possible values for expressions of the program. Among other things, these over-approximated sets are useful to exclude the possibility of a run-time error. In contrast with the heuristic techniques used in other static analysis tools, which may be very efficient but solve a different problem, shortcuts do not work when the question is to correctly — that is, without false negatives — find all possible run-time errors in large embedded programs. The analysis has to start from the entry point of the analyzed program and unroll function calls (and often loops). In addition, the modular nature of Frama-C and the interface that the value analysis aimed at providing to other plug-ins meant that abstract states had to be memorized at each statement, which dictated the choice of persistent data structures, with sharing between identical sub-parts of states (Cuoq and Doligez 2008). This meant at the very least that a garbage-collected language had to be used. While there are popular imperative languages that are garbage-collected nowadays, and some of these languages have huge libraries of data structures ready to use, persistent data structures are often under-represented in these libraries, and are slightly annoying to write in these languages. For writing a static analyzer, one is not worse off with OCaml and a few select libraries (Conchon et al. 2008; Filiâtre and Conchon 2006) than with any of Python, the .NET framework or the Java platform, al-

though it is by no means impossible to write static analyzers in any of these.

3.3 Existence of CIL

CIL (Necula et al. 2002) is an OCaml library that provides a parser and Abstract Syntax Tree (AST)-level linker for C code. CIL is well documented and provides ready-made generic analyses that are very useful to get a prototypal analyzer started. Unlike the earlier mentioned data structures that are nice to find in OCaml but could be written quickly if they were missing, having to write a C parser when one's goal is to write a static analyzer would be a major time-sink. It would have been a significant counter-argument to the choice of OCaml if such a library had not existed. We were probably a little bit lucky, and OCaml is certainly at a disadvantage with respect to other languages from this point of view. We weren't perhaps lucky to find a C parser so much as to be working in a field that is also of interest to academia. So many researchers in software engineering use OCaml for their experiments that despite the amount of work involved, such a parser could be expected to get written someday. On the other hand, finding a library, or even bindings to an existing library, for an OCaml project in a field that does not interest students and researchers could be a problem.

This drawback (less developers implies less available libraries) is mitigated by higher code re-usability when you do find code to re-use, the existence of the Caml Hump² and the fact that the grapevine works well between OCaml developers. OCaml does not have anything that begins to compare for instance with CPAN, the Comprehensive Perl Archive Network, but it does have a healthy community.

3.4 Portability

It would have been an annoyance if the development platform for Frama-C had allowed it to be used only on Unix variants, because many potential users only have access to Windows stations. Unix being the system used by the majority of the researchers at the Software Reliability Labs, the switch to a Windows-only platform was not considered. Motivated users — and at this time actually deploying formal methods requires motivation anyway — have found their way past this limitation for previous projects developed at the Labs.

From this point of view, the choice of OCaml (and later of the GTK+ toolkit for the graphical interface, through the `lablgtk` bindings) was an excellent compromise, with Unix clearly being the primary platform of the compiler, and Win32 robustly supported through either Cygwin or Visual C++. Compiling a large OCaml project on Windows+Cygwin is slow. This is probably caused one way or the other by the use of Unix compilation idioms (configuration script, makefile) on an OS where things are usually done differently, and is not a limitation in this context.

It should be noted that many OCaml developments are referenced in the source distribution `GODI`, with dependency lists and automated compilation scripts. All of those Frama-C dependencies that are written in OCaml are referenced in `GODI`, and Frama-C itself is. Some Frama-C users who have no interest in OCaml outside of Frama-C have found that this was the most convenient installation path for them.

Frama-C has been tested under Windows, Mac OS X, Solaris (32-bit), OpenBSD and Linux. Binaries are also distributed for some of these platforms.

64-bit readiness A 64-bit address space is available to OCaml programs for many 64-bit platforms. Frama-C compiles in both 32- and 64-bit mode, and the resulting versions are functionally identical. This was not a big effort, as OCaml encourages to write

²The Caml Hump is an informal central repository for OCaml libraries

high-level code, for instance by providing bignums. For some 64-bit-aware platforms (Mac OS X), it is a simple `configure` option to choose between 32-bit or 64-bit pointers at the time of compiling OCaml. For others, it is troublesome to go against the default size (Linux). With Linux, getting an OCaml compiler with a word size different from the distribution default is akin to building a cross-compiler. However, efforts are under way in the OCaml community to improve support for cross-compilation, including from Linux to Win32. We are looking forward to the maturation of such initiatives.

Availability of a graphical toolkit Frama-C uses the GTK+ toolkit for its graphical user interface. This section does not discuss the merits or demerits of this toolkit with respect to others. The question it tries to answer is “If I choose OCaml for a software project, will I find one satisfactory toolkit to design the user interface with?”. The choice of the GTK+ toolkit for Frama-C was somewhat arbitrary, but it allows to give a positive answer to the question above, without prejudice to other available toolkits.

Our experience is that for some Unix variants (Solaris, Mac OS X, very old Linux distributions), it is necessary to obtain and compile the missing GTK+ libraries manually, or semi-manually with the help of source distribution systems such as GARNOME or MacPorts. In this case, retrieving and installing the dependencies of the `gtksourceview1` library (a GTK+ widget for displaying source code) is a pain. This is not directly an OCaml problem, but another development platform could have made it possible to use the modern `gtksourceview2` (which solves the dependencies problem) or provided more toolkits to choose from initially (to the best of our knowledge, OCaml only offers Tk and GTK+ at this time). Now that `gtksourceview2` has become stable, there is talk on the `lablgtk` development list about including it in `lablgtk`. This is anyway a very minor quibble. It should be kept in mind for comparison that Java or .NET/Mono do not come pre-installed on every platform either. Again, we have no reason to regret the choices of OCaml and GTK+ from the standpoint of portability.

OCaml as a scripting language It should be noted that while this is not its strongest point, OCaml is an acceptable scripting language. In other words, when OCaml is chosen as the main language for a new project, the project may be saved the introduction of additional dependencies towards various dedicated scripting languages down the road. For instance, the HTML pages of the Frama-C web site are processed with the `yamllpp` preprocessor, which is written in OCaml. For comparison, in its 15 years of development, Caveat had at one point accumulated dependencies towards Perl, Python, Bash and Zsh (in addition to C++, and in addition to OCaml, used for journalizing and re-playing). Some of these dependencies have since been removed, by replacing some tools with OCaml equivalents.

3.5 Module system

OCaml’s module system (Leroy 1996) has direct advantages: it creates separate namespaces and, when the modules are in separate files and interfaces have been defined, fully type-checked separate compilation. It is easy to underestimate the importance of these features in the management of a big project because they make the compiler transparent, but when they are missing, their absence is unpleasantly noticeable. We discuss these, and the (theoretically more interesting) functor system *per se*.

Separate compilation With OCaml, in bytecode, separate compilation has the same meaning as everywhere: compilation is parallelizable and only modified files need to be recompiled, with a quick final link phase. With native compilation, all the ancestors of the modified modules in the dependency graph must be recompiled, and the compilation of two files with a parenthood relationship can not be parallelized. Depending on the structure of an OCaml

project, recompilation after an incremental change in a low-level module may sometimes feel longish, but in truth, it is much faster to recompile Frama-C with `ocamlOpt` than to recompile Caveat with `g++`.

The existence of two OCaml compilers, one with blazingly fast compilation in general, the other with acceptable recompilation time and producing reasonably fast code, allows very short modify-recompile-test cycles. Again, it is easy to take short recompilation times for granted but with other languages, when a software project grows in size, this can sometimes be lost, and sorely missed.

The OCaml compiler tries very hard not to get in the way between a programmer and his program, and it does not force the programmer to write interfaces. However, if the interface `m.mli` is missing for module `M`, the compiled interface is generated from `m.ml`. This means that any change to `m.ml` changes the compiled interface and forces the recompilation of every module that uses `M`, even in bytecode. In a large project, modules should always have interfaces, if only for the sake of separate compilation. OCaml has an option to generate automatically the interface `m.mli` that exports everything from `M`.

Separate namespaces for compilation units Orthogonally to separate compilation, but as importantly for big projects, OCaml’s module system provides separate namespaces. Better yet, the compiler option `-pack` allows to group several compilation units into a namespace. As a consequence, compilation units that have been put into different packs may safely use the same name for a type, variable or module. For instance the types `tree` in files both called `m.ml` in packs `lib1` and `lib2` are seen as `Lib1.M.tree` and `Lib2.M.tree`.

This feature is very useful for libraries because libraries may use very common filenames (`util.ml`) with the guarantee that there will not be a clash at link-time for users of this library (on condition that the pack name itself is unique).

In Frama-C, plug-ins are independent from each other: each plug-in only interfaces with the Frama-C kernel, and does not see the implementation details of other plug-ins. In order to implement this separation, the Frama-C system automatically packs each plug-in. Thus, two different plug-ins may use files with identical names and still be linked together within Frama-C.

Interfaces and functors The possibility to write functors (modules that are parameterized by other modules or functors), introduced before objects (at the time of Caml Special Light), has proved a workable, and completely statically checked, alternative to object-oriented programming. We use OCaml objects only when interfacing with existing OCaml code that uses objects (`CIL` and `lablgtk`), and use functors for the rest.

Some very structural idioms seem destined to be expressed with objects (or, for that matter, class types): equality, pretty-printing, hashconsing or marshaling functions³. Most of our data structure are complicated enough that automatically produced pretty-printers or equality functions would not fit the bill. Consequently, it is in our case neither more nor less tedious to write modules (and interfaces) that sport, in addition to a type `t`, functions such as `pretty: Format.formatter -> t -> unit` and `equal: t -> t -> bool`, and for unmarshaling values of a hash-consed type, `rehash: t -> t`. But, speaking of equality, it should be noted on the other hand that OCaml’s polymorphic comparison functions (including `=`, `>=` and even `==`) are dangerous pitfalls. The type-checker does not complain when they are applied wrongly instead of, for instance, `equal` above.

³ In the presence of hashconsing, not only do you have to write your own unmarshaling functions, but they are extremely tricky to get right

The module system allows to define and encapsulate the definitions of data structures, and in particular to give a purely functional interface to a sophisticated data structure that uses mutable values internally for optimization. This is the widely misunderstood nuance between purity and persistency. In fact, some positive reports on the industrial use of Haskell (Nanavati 2008) resonate deeply with our own programming experience, except that we attribute to OCaml’s module system the advantages attributed there to Haskell’s purity.

3.6 Labels and optional arguments

OCaml allows to use labels for function arguments. This feature does not make anything possible that was not already, but in practice, labels provide a concise way to remove the risk of confusion when a function takes several arguments of the same type with no obvious normal order between them. The only language that we know of with a feature vaguely similar to OCaml’s labels is Objective C’s infix notation for function calls.

Syntax begets style. The “OCaml style” is to write pure functions unless an exception needs to be made because the syntax rewards the use of immutable definitions, as seen in the following:

```
let x = 2 in ... x ...
let x = ref 2 in ... !x ...
```

We argue that using labels rewards consistent naming schemes in a similar fashion. When it is common for an argument to be passed repeatedly as-is from caller to callee without any computations actually happening to it (and in a persistent setting as much of Framac is, this happens with a lot of arguments), the labels syntax rewards the consistent choice of a unique label and eponymous variable name for this argument by a very concise syntax. In this example, `f` is being defined and calls functions `g` and `eval`.

```
let f ~mode ~env x y =
  let context = ... in
  ... g ~mode ~context (x+y) ...
  ... eval ~mode ~context ~env ...
```

If the programmer deviates from this style by using different label names or variable names for `mode`, `context`, or `env`, he receives a gentle slap on the wrist in the form of the awkward `~context:computation_context` syntax. This changes the way of reading labels-enabled OCaml programs, too. The reader can put more trust in the names of variables, without having to look for context all the time. The level of obtrusiveness of the label syntax is exactly the same as with the definition of mutable values, and it is exactly right, too. Good style is encouraged but the system can be circumvented when it needs to.

Optional arguments (a syntax for giving a labeled argument a default value if it is omitted) are convenient when the consequences of the omission (and subsequent use of the default value) are visible and traceable (for instance, to provide a toolkit interface that is both powerful and beginner-friendly). It is in general a bad idea to use an optional argument to add a new mode to an existing function, both because of all the existing calls to this function — that the compiler would be glad to help the programmer inspect if she did not use an optional argument — and because of all the calls to be written in the future where the optional argument will be omitted by accident.

4. Conclusion

We have not yet considered the point of view of the external Framac plug-in developer. We hope to see in the future many useful plug-ins written outside the circle of the initial developers. It is too early to draw conclusions on the consequences of the choice of OCaml as the platform’s language for this goal. Responses so far have ranged from the enthusiastic (“and it’s even written in OCaml”) to the rejection (“[...]drawback that the extensions have to be written in

OCaml[sic]”), with in the middle at least one person who decided to learn OCaml because there was something s/he wanted to do with Framac.

Acknowledgments

We would like to acknowledge the help of our colleagues at ProVal, at INRIA Sophia Antipolis’ projects Everest and now Marelle, and at CEA LIST, in the building of Framac. The feedback of users of Framac within the CAT project, at Fraunhofer FIRST or elsewhere has been great. And special thanks go to the developers of the OCaml system.

References

- Patrick Baudin, Anne Pacalet, Jacques Raguideau, Dominique Schoen, and Nicky Williams. Caveat: a tool for software validation. In *Dependable Systems and Networks, 2002*, pages 537+, 2002.
- Patrick Baudin, Jean-Christophe Filliâtre, Thierry Hubert, Claude Marché, Benjamin Monate, Yannick Moy, and Virgile Prevosto. *ACSL: ANSI C Specification Language (preliminary design V1.4)*, preliminary edition, October 2008. URL <http://frama-c.cea.fr/acsl.html>.
- Sylvain Conchon, Jean-Christophe Filliâtre, and Julien Signoles. Designing a generic graph library using ML functors. In Marco T. Morazán, editor, *Trends in Functional Programming*, volume 8 of *Trends in Functional Programming*, pages 124–140. Intellect, UK/The University of Chicago Press, USA, 2008. ISBN 978-1-84150-196-3.
- Pascal Cuoq. Documentation of Framac’s value analysis plug-in, 2008. URL <http://frama-c.cea.fr/download/frama-c-manual-Lithium-en.pdf>.
- Pascal Cuoq and Damien Doligez. Hashconsing in an incrementally garbage-collected system: a story of weak pointers and hashconsing in OCaml 3.10.2. In *ML ’08: Proceedings of the 2008 ACM SIGPLAN workshop on ML*, pages 13–22, New York, NY, USA, 2008. ACM. ISBN 978-1-60558-062-3.
- David Delmas, Stéphane Duprat, Patrick Baudin, and Benjamin Monate. Proving temporal properties at code level for basic operators of control/command programs. In *4th European Congress on Embedded Real Time Software*, 2008.
- Jean-Christophe Filliâtre and Sylvain Conchon. Type-safe modular hashconsing. In *ML ’06: Proceedings of the 2006 workshop on ML*, pages 12–19, New York, NY, USA, 2006. ACM. ISBN 1-59593-483-9.
- Xavier Leroy. A syntactic theory of type generativity and sharing. *Journal of Functional Programming*, 6:1–32, 1996.
- Yaron Minsky. Caml trading: Experiences in functional programming on Wall Street. In Wouter Swierstra, editor, *The Monad.Reader*, April 2007.
- Benjamin Monate and Julien Signoles. Slicing for security of code. In Peter Lipp, Ahmad-Reza Sadeghi, and Klaus-Michael Koch, editors, *TRUST*, volume 4968 of *Lecture Notes in Computer Science*, pages 133–142. Springer-Verlags, March 2008.
- Ravi Nanavati. Experience report: a pure shirt fits. *SIGPLAN Not.*, 43(9): 347–352, 2008. ISSN 0362-1340.
- George C. Necula, Scott McPeak, Shree P. Rahul, and Westley Weimer. CIL: Intermediate language and tools for analysis and transformation of C programs. In *International Conference on Compiler Construction*, pages 213–228, 2002.
- Famantanantsoa Randimbivololona, Jean Souyris, Patrick Baudin, Anne Pacalet, Jacques Raguideau, and Dominique Schoen. Applying formal proof techniques to avionics software: A pragmatic approach. In *FM ’99: Proceedings of the World Congress on Formal Methods in the Development of Computing Systems-Volume II*, pages 1798–1815, London, UK, 1999. Springer-Verlag. ISBN 3-540-66588-9.
- Morten Rhiger. A foundation for embedded languages. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 25(3):291–315, 2003. ISSN 0164-0925.
- Julien Signoles. Plug-in development guide, 2008. URL <http://frama-c.cea.fr/download/plugin-development-guide.pdf>.