

Who watches the watchers: Validating the ProB Validation Tool*

Jens Bendisposto

Sebastian Krings

Michael Leuschel

Institut für Informatik, Heinrich-Heine Universität Düsseldorf
Universitätsstr. 1, D-40225 Düsseldorf

{bendisposto, krings, leuschel}@cs.uni-duesseldorf.de

Over the years, PROB has moved from a tool that *complemented* proving, to a development environment that is now sometimes used *instead of* proving for applications, such as exhaustive model checking or data validation. This has led to much more stringent requirements on the integrity of PROB. In this paper we present a summary of our validation efforts for PROB, in particular within the context of the norm EN 50128 and safety critical applications in the railway domain.

1 Introduction

Over the years, PROB has moved from being a tool that *complemented* proving, to being a tool that is now sometimes used *instead of* proving for certain applications. For example, PROB is sometimes used to exhaustively model check B specifications which are not easily amenable to proving (e.g., where it is not easy to find inductive invariants). In 2009 PROB started to be used for data validation within Siemens [7, 8], replacing custom proof procedures. In this application, complicated properties had to be validated on concrete data values for safety critical railway applications. As this turned out to be highly successful, PROB is now being used for data validation by a variety of other companies (Alstom, ClearSy, Systemrel) and for a variety of industrial end users (e.g., RATP).

All this means that we can now longer rely on the integrity of the B provers in case a bug in PROB prevents the detection of an error. As such, validation of the PROB tool became very important and even mandatory. Indeed, according to the railway norm EN 50128 [4] [Sect. 3.1], PROB is now being used as a tool of class T2, which mandates a certain amount of documentation and validation. Our long term goal is to enable PROB to also be used as a tool of class T3, i.e., moving from a tool of class T2 that “*supports the test or verification of the design or executable code, where errors in the tool can fail to reveal defects but cannot directly create errors in the executable software*” [4] to a tool that “*generates outputs which can directly or indirectly contribute to the executable code (including data) of the safety related system*” [4]. In addition to the validation efforts described in this paper, T3 would require a secondary toolchain as described in “Checking Computations of Formal Method Tools - A Secondary Toolchain for PROB”.¹ In this paper we outline the various validation efforts that have been conducted in that light, and present a summary of our experiences.

*This research is being carried out as part of the EU funded FP7 research project 287563 (Advance).

¹To appear in proceedings of F-IDE 2014.

2 Architecture

2.1 Source Code

The kernel of PROB is developed in Prolog, but many other programming languages come into play. More precisely, the source code of PROB contains >45,000 lines of Prolog, >15,000 lines of Tcl/Tk, >5,000 lines of C (for LTL and symmetry reduction), 1,951 lines of SableCC grammar along with >10,000 lines of Java for the parser (which are expanded by SableCC into 91,000 lines of Java). In addition, there are >5,000 lines of Haskell code for the CSP parser and about 70,000 lines of Java code for the Rodin [2] plugin. These are the statistics for version 1.3.7-beta9 of PROB.

We use Git for our source code repositories, managing a history that traces back to 2005. The repository has been ported to different SCMs multiple times, starting from an initial CVS repository, to SVN and to GIT. In addition to our main Git repositories we manage several SVN repositories for examples, test cases and documentation. These have not been moved to Git for practical reasons.

We have two versions of the PROB binary:

- the ProB Tcl/Tk version, which has a GUI and
- the probcli binary. It does not depend on Tcl/Tk and can be used from the command-line, from within Makefiles and via a socket interface which listens to requests. The latter is used for integrating PROB into the Rodin Eclipse platform.

Both binaries are used for industrial applications such as data validation; the Tcl/Tk version has the additional benefit of providing a predicate debugger and graphical visualization of formulas.

2.2 Important Components

A graphical representation of the PROB architecture can be found in Figure 1. The most important components are the following:

- The PROB kernel deals with B's datatypes (integers, booleans, pairs, sets, ...) and the various operations on them (arithmetic, set operations, ...).
- The B interpreter interprets the structure of a B machine to compute the initial states, the enabled operations, the effect of performing an operation, etc. The B interpreter calls the PROB kernel (via the `kernel_mappings` interface file).
- The B Type Checker reads in the parsed abstract syntax tree (AST), type checks it and performs certain pre-processing.
- the B Parser reads in a series of B files and produces an abstract syntax tree (AST) in Prolog form (written into a file ending with ".prob").

Any error in those components is likely to result in a potential error during data validation or model checking, and as such these components need to be tested very thoroughly.

2.2.1 ProB kernel

The kernel supports the underlying data structures and operators of the B language: arithmetic, sets, relations, functions and sequences. As such, the specification for the ProB kernel is the reference book by Abrial [1], in particular Chapters 2 and 3. As additional reference, the Atelier-B handbook [5] has been used, for those operators not covered in [1].

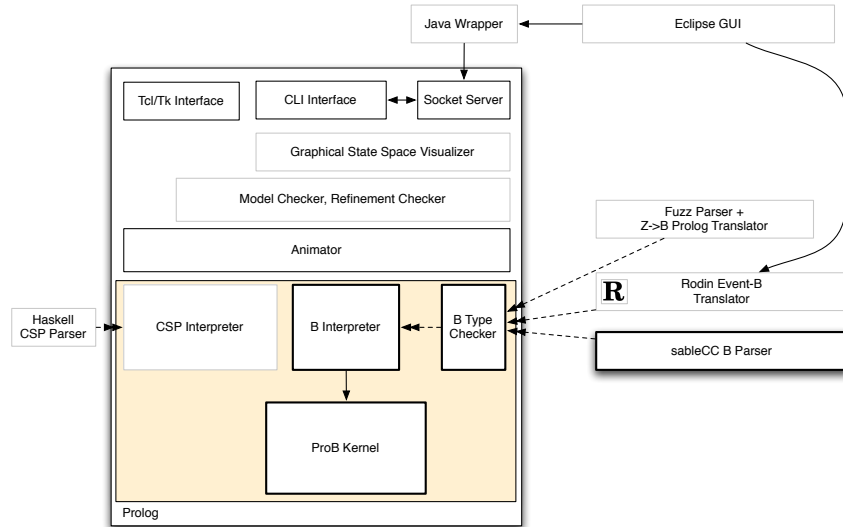


Figure 1: Architecture of ProB with highlighting of critical components (for data validation)

Many of the laws from [1] are tested below, in various forms, e.g., instances are used as unit tests or more systematically, in Section 3.4.2 we ensure that ProB can find no counter-example to those laws. For example, on page 124 we can find the law $\text{union}(SS) = \{x \mid \exists y. (y \in SS \wedge x \in y)\}$, which is added as one of the laws in the machine checking set laws. The particular example from that page — $u = \{\{0,5,2,4\}, \{2,4,5\}, \{2,1,7,5\}\}$ with $\text{union}(u) = \{0,1,2,4,5,7\}$ — is added as a regression test to the machine checking explicit computations of laws.

Several mistakes in the B-Book were found in this process; e.g., the first definition of `override` on page 80 uses p instead of r . In the second definition, it should be (a,b) instead of (x,y) . On page 102, the property involving $\text{ran}(p;q) = \text{ran}(p)$ is false, it should be $\text{ran}(q)$.

2.2.2 B Interpreter

The B interpreter evaluates B predicates, expressions² and substitutions. As such, again the book by Abrial [1] is the main specifications.

In particular, Chapter 1.2 is used for predicates, and Chapter 4 is used for substitutions. As additional reference, the Atelier-B handbook [5] has been used, for those aspects not covered in [1].

For substitutions, the definitions in [1] are in form of rules for the weakest pre-condition. As PROB uses an interpretive approach, these definitions do not apply directly, and PROB is validated in other ways.

2.2.3 Parser, Type Checker and Syntax Tree Manipulations

For the parser, the syntax of Atelier-B [5] is the reference specification. The type checking basics are specified in [1]. For the visibility rules, the tables in [5] have been used as reference.

²See page 31 of [1] for the difference.

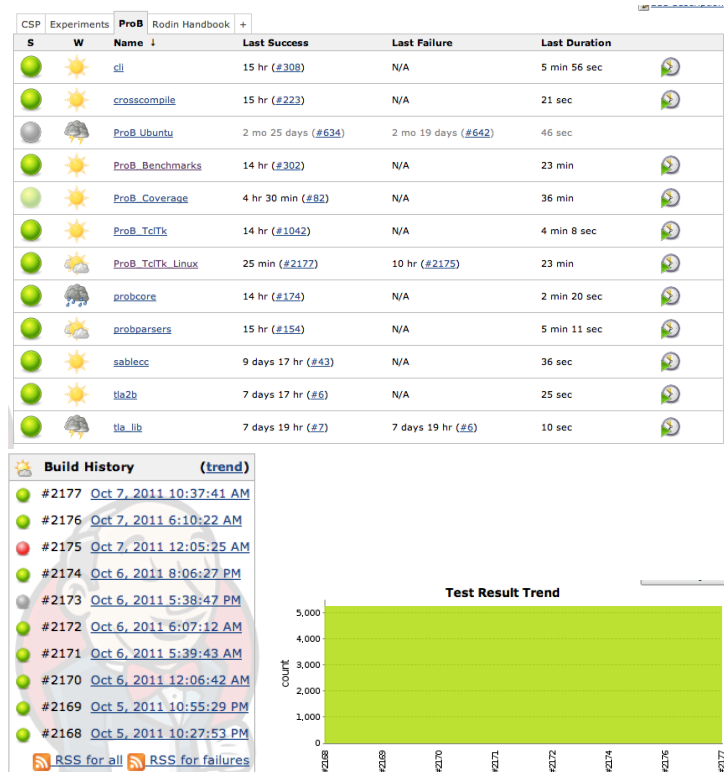


Figure 2: Different web pages of our Jenkins server

3 Building and Testing through Continuous Integration

PROB contains unit tests, integration and regression tests as well as self-model check tests for mathematical laws (see below). All of these tests are run automatically on our continuous integration platform “Jenkins”.³ The tests are run on all platforms for which we ship PROB. When a test fails, an email is sent automatically to the PROB development team. The status of the tests can be monitored on the web pages of our Jenkins instance (see Figure 2):

<http://cobra.cs.uni-duesseldorf.de/jenkins/>.

Over the years many tests have been developed. They can be classified into unit tests for the PROB kernel, parser tests, type checker tests, and system tests exercising the whole PROB tooling infrastructure.

3.1 Unit Tests for Kernel

PROB contains over a 1,000 manually entered unit tests at the Prolog level. Most of them check the proper functioning of the various PROB kernel predicates operating on B’s data structures. For example, the Prolog implementation of the union operator is exercised to check that among others $\{1\} \cup \{2\}$ evaluates to $\{1,2\}$.

There are two situation in which we manually add tests. The first is to add unit tests following the test-driven development process. The second is to add tests aiming to increase code coverage as we will

³See [http://en.wikipedia.org/wiki/Jenkins_\(software\)](http://en.wikipedia.org/wiki/Jenkins_(software)).

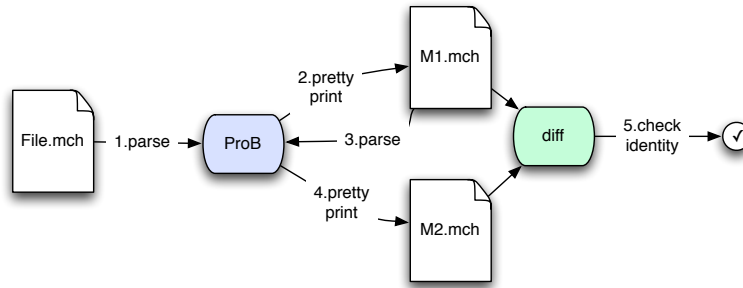


Figure 3: Overview of the Parser Validation Process

explain in Section 4.

In addition, we have developed an automatic *unit test generator*, which tests the PROB kernel predicates for many different scenarios and different set representations. For example, starting from the initial Prolog call `union([int(1)], [int(2)], [int(1), int(2)])`, the test generator will derive 1358 individual unit tests. The various tests will use different set representations (AVL-tree or symbolic representations), will swap the order of the first two arguments as `union` is declared commutative, and will check various orderings in which the sets are instantiated (e.g., it could be that first the result of the union is known, then the second argument). The latter point is particularly important for the PROB kernel, which relies on co-routines. Indeed, we have to check that the kernel predicates behave correctly no matter in which order (partial) information is propagated.

3.2 Validation of the parser

The PROB parser is written in Java and has been developed using SableCC. We use about 1800 JUnit based tests to check various parts of the parser for errors. While most of them verify the correct behavior on small examples or single expressions, some employ full scale B machines. The unit tests are run independently on both Windows and Linux. Several of these tests are automatically adopted to different newline characters and encodings.

Additionally, we execute our parser on a large number of our regression test machines and pretty print the internal representation. We then parse the internal representation and pretty print it again, verifying (with `diff`) that we get exactly the same result (see Figure 3). This type of validation can easily be applied to a large number of B machines, and will detect if the parser omits, reorders or modifies expressions, provided that the pretty printer does not compensate errors of the parser. On the downside, the validation will only detect those errors in machines generated by the pretty printer, which may prevent us from catching errors which only appear in non-pretty printed machines, e.g., when parentheses in expressions are set incorrectly.

3.3 Validation of the type checker

For the moment we also read in a large number of our regression test machines and pretty print the internal representation, this time with explicit typing information inserted. We now run this automatically generated file through the Atelier B parser and type checker (see Figure 4). With this we test whether the typing information inferred by our tool is compatible with the Atelier B type checker. (Of course,

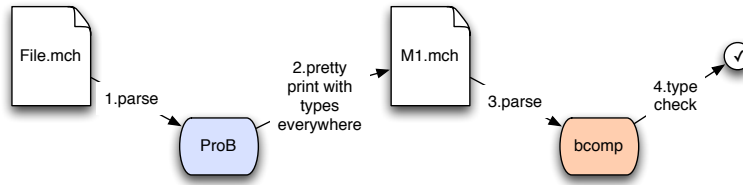


Figure 4: Overview of the Type Checker Validation Process

we cannot use this approach in cases where our type checker detects a type error.) Also, as the pretty printer only uses the minimal number of parentheses, we also ensure to some extent that our parser is compatible with the Atelier B parser (see below). Again, this validation can easily be applied to a large number of B machines. More importantly, it can be systematically applied to the particular B machines that PROB validates: provided the parser and pretty printer are correct, this gives us a guarantee that the typing information for those machines is correct. The latest version of PROB Tcl/Tk has a command to cross check the typing of the internal representation with Atelier B in this manner.

Errors Detected During validation of the type checker we found several issues. (We did not detect errors in the type checker, but rather differences between our parser and the Atelier B parser bcomp.) For example, PROB accepts identifiers with single letter without warning and PROB accepts enumerated set elements that are used as identifiers while bcomp does not. We also detected that bcomp reports a lexical error (“illegal token |–”) if the vertical bar (|) of a lambda abstraction is followed directly by the minus sign. Most importantly, however, we found errors in the priority table of the english Atelier B “B Language Reference Manual 1.8.6” [5]. All in all there are 26 errors in the English reference manuals, upon which our pretty printer and BParser was based. These issues have been fixed both in our parser and in the handbook.

3.4 System Tests

These tests exercise the entire PROB tool chain, including the parser, the type checker, the PROB interpreter and the PROB kernel.

3.4.1 Regression Tests

PROB contains over 1100 regression tests which are made up of B models along with saved animation traces. These models are loaded, the saved animation traces are replayed, and the models are also run through the model checker. These tests have turned out to be extremely valuable in ensuring that a bug once fixed remains fixed. They are also very effective at uncovering errors in arbitrary parts of the system (e.g., the parser, type checker, the interpreter, the PROB kernel, etc.). The various features of PROB are also checked here: there are tests to check that PROB’s visualization features work as expected, there are tests for PROB’s test case generation algorithms, there are tests for the various constraint-based checks that PROB can perform, etc.

```

law1 == (dom(ff ∪ gg) = dom(ff) ∪ dom(gg));
law2 == (ran(ff ∪ gg) = ran(ff) ∪ ran(gg));
law3 == (dom(ff ∩ gg) ⊆ dom(ff) ∩ dom(gg));
law4 == (ran(ff ∩ gg) ⊆ ran(ff) ∩ ran(gg));
....

```

Figure 5: A small selection of the laws about B functions

3.4.2 Self-Model Check with Mathematical Laws

With this approach we use PROB’s model checker to check itself, in particular the PROB kernel and the B interpreter. The idea is to formulate a wide variety of mathematical laws and then use the model checker to ensure that no counterexample to these laws can be found.

Concretely, PROB now checks itself for over 500 mathematical laws. In version 1.3.3 there were laws for booleans (39 laws), arithmetic laws (40 laws), laws for sets (81 laws), relations (189 laws), functions (73 laws) and sequences (61 laws), as well as some specific laws about integer ranges (24 laws) and the various basic integer sets (7 laws). The number of laws is continuously being expanded. Figure 5 contains some of these laws about function.⁴

This technique can be viewed as using the model checker to generate a very large number of tests from the mathematical laws. The self-model check has been very effective at uncovering errors in the PROB kernel and interpreter. Furthermore, the self-model checking tests rely on every component of the PROB execution environment working perfectly; otherwise a violation to a mathematical law could be found. In addition to the PROB main code, the parser, type checker, Prolog compiler, hardware and operating system all have to work perfectly. Indeed, we have identified a bug in our parser (FIN was treated the same as FIN1) using the self-model check. Furthermore, we have even uncovered two bugs in the underlying SICStus Prolog compiler using self-model check:

- The Prolog `findall` did sometimes drop a list constructor, meaning that instead of `[[[]]]` it sometimes returned `[]`. In terms of B, this meant that instead of $\{\emptyset\}$ we received the empty set \emptyset . This violated some of our mathematical laws about sets, notably the law $\text{POW1}(SS) = \text{POW}(SS) - \{\{\}\}$. This bug was reported to SICS and it was fixed in SICStus Prolog 4.0.2.
- A bug in the AVL library (notably in the predicate `avl_max` computing the maximum element of an AVL-tree) was found and reported to SICS. The bug was fixed in SICStus Prolog 4.0.5.

Note that these problems would not have been detected by validating or proving the code of PROB correct. It was essential to test the actual code of PROB along with the entire execution environment. The model checker together with the mathematical laws enabled this testing to be performed very effectively (see also Example 1 below).

We also applied some of these tests to other tools. For example, we found various bugs in the first implementation of the JEB animator [9], and the TLC [10] model checker when applied to TLA translations of our B laws.

4 Coverage

The validation techniques explained above are complemented by code coverage analysis techniques. In particular, we try to ensure that the unit tests (Section 3.1) and the self-model checks (Section 3.4.2)

⁴Throughout the paper, we replaced B-style notation with mathematical notation for better comprehensibility.

above cover all predicates and clauses of the PROB kernel.⁵ We count a clause as covered if it has been executed successfully by at least one test case. A predicate is regarded as covered if at least one of the clauses belonging to it is covered.

Initially, we have developed our own code coverage tool for SICStus Prolog. The tool uses Prolog's term expansion facility to keep a record of which program points are covered. As of version 4.2, SICStus Prolog now provides code coverage profiling, along with visualizations in Eclipse and Emacs. The coverage is slightly less fine-grained than our original implementation, but it is much faster (as the profiling is directly supported by the Prolog abstract machine) and supports all of Prolog's features out of the box. We have thus decided to switch to SICStus' code coverage profiler, but have implemented various backends to output the information in HTML, LaTeX, or other formats (see <http://nightly.cobra.cs.uni-duesseldorf.de/coverage/>). One of the backends features a colorized version of our source code, highlighting the coverage status. With this backend, it is possible to inspect the source code itself, finding out exactly which parts of the code are covered. This can be used to identify poorly tested source code and helps in adding missing test cases.

The code coverage statistics are also systematically computed by our continuous integration server (see Section 3). As of version 1.3.5 of PROB, we now also keep track of clauses which should not be covered as they relate to internal errors:

- the `error_manager` module was extended to distinguish internal and normal errors. A normal error is, for example, a division by zero in a B machine. An internal error would be a typing error or unexpected failure within a piece of PROB Prolog code. Normal errors can be covered by tests; internal errors must not be coverable.
- a Prolog term expander keeps track of clauses and execution paths which relate to internal errors, and excludes those from the coverage statistics.

So far, the code coverage analysis has been very useful in extending our unit tests and mathematical laws. As an example, 96.5 % of the clauses of the `kernel_mappings` module are now covered by the unit tests and the mathematical laws. The only uncovered clauses relate to the Z compaction operator, and one error condition that was not triggered by the tests. In summary, the code coverage has helped us write better tests and has allowed us to uncover a few undetected errors in the kernel.

However, even fully covered code might still contain errors. Executing all parts of the source code does not necessarily mean that all paths were covered. There might, for instance, be bugs that only occur if certain combinations of branches are executed. Thus we can not rely solely on the coverage reports to judge the effectiveness of our test suite. We experimented with artificially introduced bugs to evaluate how certain classes of programming errors can be found. The following section will give a small overview.

4.1 Judging Effectiveness of Tests

In Section 3 we already described the effectiveness of the current testing procedure at uncovering various errors in PROB and in the underlying compiler. We now check the effectiveness for artificially introduced errors.

The tests in this subsection were conducted on PROB version 1.3.4-rc1, svn revision 9543.

⁵We do not count the regression tests towards coverage, as there we basically compare against previous versions of PROB, which could in theory already contain a serious bug.

Error	Unit Tests	Regression	Laws
Replaced \cap by set difference	Passed	Passed	Failed
Subtle error in integer multiplication	Failed	Passed	Failed
Interchanged argument in not partial function test	Passed	Passed	Failed
Skip values when enumerating a variable	Failed	Passed	Passed

As you can see, all errors were detected and the mathematical laws checks detected most of the errors. However, as checking the mathematical laws is based on searching for counter-examples, there is a certain class of errors that can not be found. If due to a bug the search space is pruned, mathematical law checking will obviously not be able to find more counter-examples withing the pruned area. Therefore the mathematical law checks do not make the unit tests obsolete.

In the future, we plan to conduct a more extensive study about the effectiveness of our tests.

5 Performance and Codespeed

We run several benchmarks for PROB on a daily basis. The results are stored using the codespeed web application⁶ which allows us to compare selected revisions on selected benchmarks. Furthermore, an unusual increase or decrease in performance is reported right on the front page. Strictly speaking this does not add value to our test infrastructure, as the benchmarks do not report error or success. However, in several occasions the existence of a bug in PROB lead to strange speed-ups or slowdowns, for instance, if a bug causes PROB to check every possible value of an existential quantification instead of stopping after finding a solution. Therefore, monitoring the overall speed of our applications is a useful addition to testing.

In order to identify the particular revision causing unexpected benchmark results, we have found `git bisect` to be useful. It allows us to mark certain revisions as good (i.e. the last benchmarked revision that returned the expected results) and certain other revisions as bad (i.e. the revision showing strange behavior). Git then performs an interactive binary search, checking out revisions in between until the one responsible has been found.

The process of comparing certain revisions of PROB has been partly automated by several scripts that control checkout, benchmarking and result reporting.

6 Static Analysis

Some bugs are hard to detect by testing but easy to detect with static analysis. For instance, we had some calls to error handling predicates that were not used correctly, e.g., they were called with the wrong arity. In most cases these calls happened in source code locations that are unreachable unless a previous bug occurs in the underlying execution engine. There were also bugs in the reachable parts of ProB that were not well tested, i.e., those parts that are of a more experimental nature.

Prolog is not statically typed and therefore we cannot rely on a typechecker to identify incorrect calls to predicates, however it is very easy to implement an analysis for these sort of problems in Prolog by augmenting the interpreter with a term expander. Term expansion is very similar to macro expansion in Lisp. A term expander hooks into the loading process of a Prolog file and allows modification of a term while it is being loaded. In our case, we did not modify the terms, but extract information from it. Hence, the expanded code is equal to the code used in production. This allows us to relate

⁶<https://github.com/tobami/codespeed/>

all detected bugs directly to our code. For instance, if we analyze the term `:- dynamic foo/1` in a module `bar`, we store the information that a predicate `foo` of arity 1 exists in the module `bar` and that it is a dynamic predicate. We also extract information about predicate and fact definitions, mode declarations, multifile and meta annotations, exported and imported predicates, blocking declarations for co-routines and most importantly calls. Strictly speaking this is not traditional static analysis, because Prolog directives such as `use_module` are executed while loading the Prolog program. But the program is not run either, so we think it is appropriate to classify the approach as a static analysis. To run the analyzer, we first load our term expander and then the entry point module of ProB. During loading the module, the term expander inspects all terms that are being loaded and stores the extracted information in a Prolog database. The analyzer automatically detects calls where no definition is known at compile time. These calls are potential bugs. Using the tool, we were able to find 17 missing import statements, 9 cases where predicates are used with the wrong arity, 5 cases where the wrong predicate was called and 5 missing dynamic declarations. The missing imports did not actually cause wrong behavior at runtime because of a Sicstus Prolog implementation detail. As long as a module is loaded (even if it was loaded by any other module), a predicate can be called using a fully qualified form `module:predicate(...)` at runtime. This is rather fragile because if for some reasons the module is not loaded the calls will break. Each module should declare its dependencies explicitly. The calls to predicates with wrong arity were also not critical because they were in dead code (if the Prolog engine works properly).

Three cases where the wrong predicate was called were actual bugs in some experimental features of ProB. One feature was removed, the other two were fixed. The other two cases were dead code. The missing dynamic declarations are not dangerous in the current version of Sicstus Prolog, but there is no guarantee that this won't cause actual bugs in future versions. The problems we found using the tool are summarized in the following table:

Problem Type	Fragile code/Bad practice	Bug	Dead code
Import missing	17	0	0
Wrong arity	0	0	9
Wrong predicate	0	3	2
Missing dynamic delaration	5	0	0

Because Prolog is a dynamic language, finding dead code is harder than in a static language such as Java. The compiler cannot know if a call is produced dynamically at runtime. This can easily lead to dead code, if we change our code. For instance, if we replace a call to predicate `p` by a call to predicate `q`, is `p` now dead code or is it used somewhere else? Using our analysis we get a list of predicates that are not being called in a static way. This clearly produces false positives, but it is possible to manually inspect the candidates. We think that we can further improve the analysis by detecting typical dynamic usage pattern such as using a predicate in a higher order predicate, e.g. `maplist`.

In addition to the static analysis, we are currently working on a tool that can be used to inspect the ProB source code. The tool should contain a code browser that allows to directly navigate between caller and callee, search for predicates, filter according to specific criteria (e.g., parts of the name). This tool should also incorporate the documentation. We want to display certain properties of the source code (e.g., percentage of predicates that are documented or percentage of predicates with a mode declaration) to encourage the developers to increase the quality of the documentation and code.

Our custom tools are accompanied by Spider, an Eclipse based IDE specifically developed by Sicstus for Prolog development. It features several static analyses, including calling mode and checks for potential non-determinacy and potential non-termination. In addition, there is an option to check for

common mistakes and bad practice.

7 Dynamic Runtime Analysis

In this section we examine various ways we ensure consistency of PROB's output at runtime.

7.1 Double chain

One way to dramatically increase the confidence in PROB's output is to double check the result with a separately developed tool. Thus far, this has been done with the following three tools:

1. PredicateB. Together with the company ClearSy we have performed Data Validation tasks for Alstom [6]. Here, the Java based tool PredicateB was used to double check PROB's output. Quite often, PredicateB took considerably more time than PROB to check the properties. Some properties needed to be rewritten to be suitable for PredicateB (whose constraint solving capabilities are much weaker). Overall, this approach was very successful and led to the development of the DTVT (Data Table Validation Tool) which integrates PROB and PredicateB.
2. Ovado. In this work [3], the company Systemel is using the tool Ovado as primary tool, and uses PROB to double check Ovado's output.
3. PyB. Finally, in recent work we are developing an independent tool written in Python to check B predicates. The idea is to feed PROB's output into PyB (in a more integrated fashion, e.g., passing witnesses for existential predicates) and to let PyB check the output. An independent companion paper about PyB has been submitted to this workshop.

7.2 Run Time Checking

The Prolog code contains a monitoring module which — when turned on — will check pre- and post-conditions of certain predicate calls and also detect unexpected failures. Many kernel predicates also check for unexpected arguments. All of this overcomes the fact that Prolog has no static typing to some extent. So far we have not systematically relied on this feature; but we could turn the run time checking on during property validation. The only drawback is that the runtimes will be increased.

7.3 Positive and Negative Evaluation

As already mentioned, all properties and assertions were checked twice, both positively and negatively (see Figure 6). Indeed, PROB has two Prolog predicates to evaluate B predicates: one positive version which will succeed and enumerate solutions if the predicate is true and a negative version, which will succeed if the predicate is false and then enumerate solutions to the negation of the predicate. The reason for the existence of these two Prolog predicates is that Prolog's built-in negation is generally unsound and cannot be used to enumerate solutions in case of failure. For example, given the Prolog rule $p(\text{eq}(X, Y)) :- X=Y$. the query $\text{not}(p(\text{eq}(X, 0)))$ would fail, i.e., one could erroneously conclude that there is no value X which is different from 0. In PROB, we have a dedicated predicate for the negation, which will suspend until it can determine the result correctly. For the above example, one could write $\text{not_p}(\text{eq}(X, Y)) :- \text{dif}(X, Y)$.

With these two predicates we can uncover undefined predicates: if for a given B predicate both the positive and negative Prolog predicates fail then the formula is undefined. For example, the property

P	$\neg P$	Conclusion
true	true	bug in PROB
true	false	P is true
false	true	P is false
false	false	P is not well-defined
true	timeout	P is (probably) true
false	timeout	P is false or undefined
timeout	true	P is (probably) false
timeout	false	P is true or undefined
timeout	timeout	P is unknown

Figure 6: Positive and negative evaluation of predicates

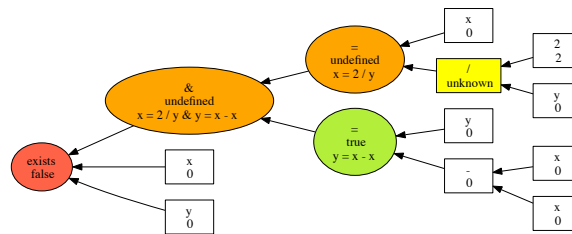


Figure 7: Visualising an undefined property

$x = 2/y \ \& \ y = x - x$ over the constants x and y would be detected as being undefined, and would be visualized by our graphical formula viewer as in Figure 7 (yellow and orange parts are undefined).

In the context of validation, this approach has another advantage: for a formula to be classified as true the positive Prolog predicate must succeed *and* the negative Prolog predicate must fail, introducing a certain amount of redundancy (admittedly with common error modes). In fact, if both the positive and negative Prolog predicates would succeed for a particular B predicate then a bug in PROB would have been uncovered. If both fail, then either the B predicate is undefined or we have a bug in PROB.⁷

This validation aspect can detect errors in the predicate evaluation parts of PROB, i.e., the treatment of the Boolean connectives \vee , \wedge , \Rightarrow , \neg , \Leftrightarrow , quantification \forall , \exists , and the various predicate operators such as \in , \notin , $=$, \neq , $<$, ... This redundancy can not detect bugs inside expressions (e.g., $+$, $-$, ...) or substitutions (but the other validation aspects mentioned above can).

Example 1. We want to illustrate the error detection of double evaluation by a simple example. Take the following very simple B machine, with two assertions which should obviously be true.

```

MACHINE DoubleEvaluationTest
SETS ID={aa , bb}
CONSTANTS iv
PROPERTIES iv ∈ ID & iv ≠ bb
ASSERTIONS iv ∈ {aa}; iv ∉ {bb}

```

⁷Typically, PROB would also generate error messages for well-definedness problems. Currently, there is however no guarantee that all well-definedness problems will result in an error message; the only certain way to detect them is to do both positive and negative evaluation.

END

Say that we introduce an error in the part of the source code that deals with testing membership of elements in a singleton set. An element is a member of a singleton set if it is equal to the single element. We replaced this test by its negation, and thus introduced an error into the PROB kernel. When analyzing the assertions, we now get the following result, which tells us that there is a bug in the ProB kernel, because the assertion stating that `iv` is not a member of the set `{bb}` is both true and false at the same time:

```
iv ∈ {aa}
== unknown

iv ∉ {bb}
== both_true_false
```

Also note that the self-model checking described in Section 3.4.2 immediately provides a counter example ($SS=\{e11\}$, $TT=\emptyset$) for several of the mathematical laws for sets, for example:

```
{xx | xx ∈ SS ∨ xx ∈ TT } = SS ∪ TT
== false
```

8 Future Work and Discussion

In the longer term we want to further increase the confidence in PROB's output so that it can eventually be used as a tool of class T3 within the railway norm EN 50128.

In addition to the development of the double chain PyB, there are several techniques that might help to increase the confidence in our primary validation tool even more. One of these techniques is fuzzing or robustness testing. The idea behind fuzzing would be to generate random input data and feed it to PROB. As mentioned before, test coverage of error cases is hard to archive. Random inputs of data could help us to identify cases in which some part of our tool chain crashes instead of reaching a controlled error state.

We could also further refine the results of the coverage analysis mentioned in Section 4. Instead of just observing which predicates were executed we could strive for coverage of possible execution paths or the evaluation of each branch. However, this would again increase the complexity of the analysis.

Aside from improving the static checks strengthening the run time checks would be another reasonable approach. With the increased usage of our code documentation tools more annotations will be added to the source code. The correctness of several of these annotations, could be enforced at runtime. One example, is the Prolog mode declaration which specifies which arguments are input arguments or output arguments respectively.

9 Conclusion

In summary, the following points increase our confidence in the proper working for PROB in general and for data validation in particular:

1. repeated successful application on case studies (e.g., by Siemens on the lines of San Juan, Paris Line 1, Sao Paulo, Barcelona Line 9, ... [7, 8]), discovering all known problems and in some cases even discovering previously unknown problems in the data properties.

2. extensive regression, unit-tests, automated unit-test generation, and validation via mathematical laws and model checking (see Section 3). These tests have even discovered several errors in the underlying Prolog compiler.
3. coverage analysis of the testing, ensuring that all critical parts are covered by the tests; this has enabled us to detect a few more issues (and has prompted us to develop the automated unit-test generator).
4. validation of the parser and type-checker via pretty-printing and cross-checking with the Atelier-B parser (bcomp) (see Section 3.2). The internal representation of PROB for a particular model can be fed *at runtime* to bcomp, for additional validation that the type inference was properly conducted.
5. double evaluation of predicates (see Section 7.3), to catch undefined predicates as well as errors in the PROB kernel and interpreter relating to predicates.

References

- [1] Jean-Raymond Abrial (1996): *The B-Book*. Cambridge University Press.
- [2] Jean-Raymond Abrial, Michael Butler & Stefan Hallerstede (2006): *An open extensible tool environment for Event-B*. In Zhiming Liu & Jifeng He, editors: *Proceedings ICFEM'06*, LNCS 4260, Springer-Verlag, pp. 588–605. Available at <http://eprints.ecs.soton.ac.uk/12711/>.
- [3] Frédéric Badeau & Marielle Doche-Petit (2012): *Formal Data Validation with Event-B*. CoRR abs/1210.7039. Available at <http://arxiv.org/abs/1210.7039>. Proceedings of DS-Event-B 2012, Kyoto.
- [4] CENELEC (2011): *Railway Applications – Communication, signalling and processing systems – Software for railway control and protection systems*. Technical Report EN50128, European Standard.
- [5] ClearSy (2009): *Atelier B, User and Reference Manuals*. Aix-en-Provence, France. Available at <http://www.atelierb.eu/>.
- [6] Thierry Lecomte, Lilian Burdy & Michael Leuschel (2012): *Formally Checking Large Data Sets in the Railways*. CoRR abs/1210.6815. Available at <http://arxiv.org/abs/1210.6815>. Proceedings of DS-Event-B 2012, Kyoto.
- [7] Michael Leuschel, Jérôme Falampin, Fabian Fritz & Daniel Plagge (2009): *Automated Property Verification for Large Scale B Models*. In A. Cavalcanti & D. Dams, editors: *Proceedings FM 2009*, LNCS 5850, Springer-Verlag, pp. 708–723.
- [8] Michael Leuschel, Jérôme Falampin, Fabian Fritz & Daniel Plagge (2011): *Automated property verification for large scale B models with ProB*. *Formal Asp. Comput.* 23(6), pp. 683–709. Available at <http://dx.doi.org/10.1007/s00165-010-0172-1>.
- [9] Faqing Yang, Jean-Pierre Jacquot & Jeanine Souquières (2012): *The Case for Using Simulation to Validate Event-B Specifications*. In Karl R. P. H. Leung & Pornsiri Muenchaisri, editors: *APSEC, IEEE*, pp. 85–90. Available at <http://dx.doi.org/10.1109/APSEC.2012.66>.
- [10] Yuan Yu, Panagiotis Manolios & Leslie Lamport (1999): *Model Checking TLA⁺ Specifications*. In Laurence Pierre & Thomas Kropf, editors: *Proceedings CHARME'99*, LNCS 1703, Springer-Verlag, pp. 54–66. Available at <http://link.springer.de/link/service/series/0558/bibs/1703/17030054.htm>.