

From Failure to Proof: The PROB Disprover for B and Event-B

Sebastian Krings, Jens Bendisposto and Michael Leuschel

Institut für Informatik, Universität Düsseldorf**
Universitätsstr. 1, D-40225 Düsseldorf
`{krings,bendisposto,leuschel}@cs.uni-duesseldorf.de`

Abstract. The PROB disprover uses constraint solving to find counter-examples for B proof obligations. As the PROB kernel is now capable of determining whether a search was exhaustive, one can also use the disprover as a prover. In this paper, we explain how PROB has been embedded as a prover into Rodin and Atelier B. Furthermore, we compare PROB with the standard automatic provers and SMT solvers used in Rodin. We demonstrate that constraint solving in general and PROB in particular are able to deal with classes of proof obligations that are not easily discharged by other provers and solvers. As benchmarks we use medium sized specifications such as landing gear systems, a CAN bus specification and a railway system. We also present a new method to check proof obligations for inconsistencies, which has helped uncover various issues in existing (sometimes fully proven) models.

1 Introduction and Motivation

Both the B-method [1] and its successor Event-B [2] are state-based formal methods rooted in set theory. They are used for the formal development of software and systems that are correct by construction. This usually involves formal proofs of different properties of the specification.

In former work [23] we described a disprover based on using PROB's constraint solver to automatically find counter-examples for given proof obligations and thus saving the user from spending time in a futile interactive proof attempt. Say that we have to prove that the goal G is a logical consequence of the hypotheses H_1, \dots, H_n . The PROB disprover then tries to find a solution for the formula $H_1 \wedge \dots \wedge H_n \wedge \neg G$. If it can find a solution, the proof cannot succeed and the solution is a counter-example.

In [23] we already made the observation that in some cases, namely if we neither encounter infinite sets nor deferred sets¹ whose cardinality is unbounded, the absence of a counter-example is actually a proof. We thus suggested as future work to implement an analysis that checks if the absence of a counter-example is

** Parts of this research have been sponsored by the EU funded FP7 project 287563 (ADVANCE) and by an industrial project funded by Alstom.

¹ Deferred sets are sets which are not given upfront by enumerating their elements.

a valid proof. This work has been finalized in the last year: PROB now keeps track of infinite enumeration, in particular the scope in which an infinite enumeration has occurred and whether a solution has been found or not. This enables our technique to detect if the search for a counter-example was exhaustive, i.e., we can now use PROB as a prover. Note that we go beyond the suggested future work of [23]: we allow variables with an infinite domain to occur, as long as they do not have to be enumerated exhaustively. We have also improved the core algorithm of [23] in various ways, by allowing to focus on selected hypotheses and by providing a way to detect inconsistencies in the hypotheses or potential bugs in the disprover. In this paper we have also conducted a thorough empirical evaluation, comparing our constraint-based proof with existing provers for B and Event-B. This study shows that the constraint-based proof fares surprisingly well for a variety of case studies.

2 Constraint-Based Proof Technique

In the following section we describe how PROB can be used as a prover inside Rodin [3] and Atelier B [12]. First, we provide a short introduction to the constraint solving capabilities of PROB in Section 2.1. Further technical details regarding PROB’s kernel can be found in [21,22] or [20]. Following, Section 2.2 will outline how PROB was embedded into Rodin’s proof architecture. Section 2.3 will explain the integration of PROB into Atelier B. Afterwards, in Section 2.4 we will show how PROB can be used to detect inconsistencies in the model.

2.1 PROB’s Constraint Solving Kernel

The PROB constraint solver is based on CLP(FD)-style constraint-propagation [11], i.e., the variables of a B specification are annotated with possible values (e.g., in the form of intervals for integer variables). This information is propagated from one variable to another, e.g., if we know that x is in the range 0..8 and the predicate $x = y + 2$ holds, then y must be in the range $-2..6$. As a last resort, PROB enumerates undetermined variables when no further propagation is possible. While doing so, the solver tracks where and why enumeration occurs. It is able to distinguish between safe and unsafe enumerations, i.e., if all possible values of a variable have to be tried out or if a single solution is sufficient. This is done by observing the context² in which an enumeration occurs. Exhaustive enumeration can then be detected individually for each variable and later be transferred to the whole constraint if possible. Let us look at a few example constraints, where we suppose all free variables to be existentially quantified:

$$- i \in \{1, 2, 1024, 2048\} \wedge i > 2 \quad \vdash? \quad i \bmod 2 = 1$$

Here, we have the two hypotheses $i \in \{1, 2, 1024, 2048\}$ and $i > 2$ and we want to prove that $i \bmod 2 = 1$ is a logical consequence. Hence, we would construct the formula $i \in \{1, 2, 1024, 2048\} \wedge i > 2 \wedge \neg(i \bmod 2 = 1)$ and try

² This includes quantification, negation and arbitrarily nested combinations of them.

to find solutions for i . For this formula, PROB finds two solutions ($i = 1024$ and $i = 2048$) and no infinite enumeration has occurred (PROB has narrowed down the interval of i to 3..2048 before enumeration has started). As such, we can conclude that $G \equiv i \bmod 2 = 1$ is *not* a logical consequence of the hypotheses $H_1 \equiv i \in \{1, 2, 1024, 2048\}$ and $H_2 \equiv i > 2$.

– $i \in \{1, 2, 1024, 2048\} \wedge i > 2 \vdash? i \bmod 2 = 0$

For the opposite of the goal, i.e., $i \bmod 2 \neq 1$ or equivalently $i \bmod 2 = 0$, we construct the formula $i \in \{1, 2, 1024, 2048\} \wedge i > 2 \wedge \neg(i \bmod 2 = 0)$. In this case PROB finds no solution and no infinite enumeration has occurred. As such, we have proven that $i \bmod 2 = 0$ follows logically from $i \in \{1, 2, 1024, 2048\} \wedge i > 2$.

– $i > 20 \vdash? i \bmod 2 = 1$

If we want to prove that $(i \bmod 2 = 1)$ is a logical consequence of $i > 20$, we construct the formula $i > 20 \wedge \neg(i \bmod 2 = 1)$. PROB finds a solution ($i = 22$), but infinite enumeration has occurred in the sense that the possible values of i lie in the interval 22.. ∞ . However, in this context this is not an issue, as a solution has been found. As such, we can conclude that $i \bmod 2 = 1$ is not a logical consequence of $i > 20$.

– $i > 20 \vdash? (i \bmod 2 = 0 \vee i \bmod 1001 \neq 800)$

Finally, if we want to prove that $(i \bmod 2 = 0 \vee i \bmod 1001 \neq 800)$ is a logical consequence of $i > 20$, we get the formula $i > 20 \wedge \neg(i \bmod 2 = 0 \vee i \bmod 1001 \neq 800)$. Here PROB finds no solution, but an “enumeration warning” is produced. Indeed, the constraint solver has narrowed down the possible solutions for i to the interval 801.. ∞ , but with the default search settings no solution has been found. Here, we cannot conclude that $i \bmod 2 = 0 \vee i \bmod 1001 \neq 800$ is a logical consequence of $i > 20$. Indeed, $i = 1801$ is a counter-example.³

2.2 Integration into Rodin for Event-B

When working on a proof obligation, Rodin keeps track of two sets of hypotheses: the set of *all* available hypothesis for the target goal and a *user-selected* subset. The idea is to be able to reduce the search space of the automatic provers by excluding irrelevant hypotheses. In the case of the PROB prover we could, for instance, get rid of hypotheses that are irrelevant for the proof but contain variables over infinite domains, deferred sets or complicated constraints.

This approach cannot lead to false positives, because limiting the number of available hypothesis cannot render a formerly unprovable sequent provable. However, disproving while omitting hypotheses can lead to false negatives if the hypotheses are too weak for a proof. For instance, say the goal G is $i \bmod 2 = 1$ and the hypotheses are $i \in \{1, 2, 3\}$ (H_1) and $i \neq 2$ (H_2). PROB will not find a counter-example for $H_1 \wedge H_2 \wedge \neg G$ but it will find a (false) counter-example for $H_1 \wedge \neg G$.

Figure 1 outlines how the disprover proceeds in more detail:

³ Which PROB can find if we enlarge the default search space, e.g., by adding $i < 10000$ as additional constraint.

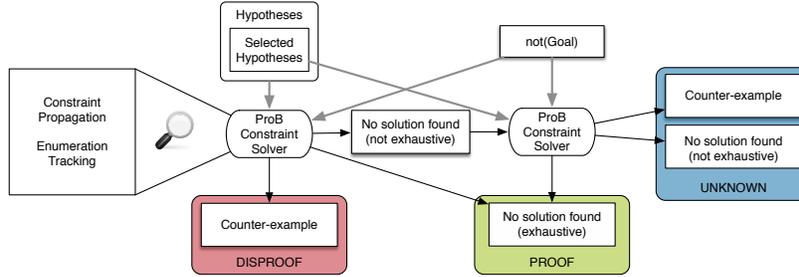


Fig. 1. Disproving Algorithm

1. We first try to solve the predicate $H_1 \wedge \dots \wedge H_m \wedge \neg G$, i.e., the negated goal together with *all* available hypotheses. If we find a solution, we report the proof obligation as *unprovable* and insert the counter-example inside the Rodin proof tree. If no counter-example is found and search was exhaustive, the initial sequent is *proven*, because no counter-example *exists*.
2. If the constraint solver is unable to prove or disprove the predicate in step 1, we reduce the number of hypotheses to the user-selected hypotheses and again look for a counter-example. The three possible outcomes are:
 - A contradiction is detected with the reduced set of hypotheses. This is still a valid *proof*, as removing hypotheses can only introduce further counter-examples but not remove them.
 - If we find a solution, we report a *possible* counter-example, but leave the proof obligation status as *unknown*. However, we do not interfere with the ongoing proof effort, as the proof obligation might still be provable using all hypotheses.
 - Otherwise we return without a result (status is *unknown*).

2.3 Integration into Atelier B for classical B

The integration of PROB into Atelier B is closer to the original implementation of the disprover explained in [23].⁴ Within Atelier B a proof obligation is translated into a B machine, where all hypotheses are put into the properties clause and the assertions clause contains an implication of the form $SelHyp \Rightarrow Goal$. Here, $Goal$ is the proof goal, and $SelHyp$ are the selected hypotheses. The latter are empty if `prob(0)` is called from Atelier B and contain all hypotheses H_1 which have a variable in common with $Goal$ if `prob(1)` is called. When `prob(2)` is called, Atelier B recursively adds all further hypotheses which have variables in common with H_1 . The selection algorithm is the same that is used for the other Atelier B provers (e.g., `pp(0)`, `pp(1)`, `pp(2)`). It is also possible to specify a

⁴ This work was conducted in a joint project with ClearSy (Lilian Burdy, Etienne Prun) and funded by Alstom (Fernando Mejia).

time-out \mathfrak{t} in milliseconds: `prob(n| \mathfrak{t})`. Once the machine is constructed, Atelier B calls the command line version of PROB, which tries to find a counter-example to $SelHyp \Rightarrow Goal$ and writes the result to an intermediate file. The possible result values are very similar to above:

- no counter-example exists: the proof obligation is proven,
- no counter-example found (with reason being either time-out, deferred sets used or enumeration warning): the proof obligation status is unknown
- counter-example found: the proof obligation status is still unknown, but not provable from the selected hypotheses.

2.4 Inconsistency Detection

After the algorithms outlined in Sections 2.2 and 2.3 return a proof, a second phase can be triggered as outlined in Figure 2: We try to find a proof for the negation of the goal. This time, we send $H_1 \wedge \dots \wedge H_m \wedge G$ to the constraint solver. The result allows us to decide, whether the goal predicate G played a role in the original proof. If the negated goal can be proven as well, we detected a contradiction in the hypotheses. Contradicting hypotheses might occur due to an error in the model, in particular if they are detected at the root of the proof tree.⁵ Hence, the user should be notified if they occur in a successful proof.

If contradicting hypotheses or disproven obligations have been found, PROB can afterwards compute the unsat core in order to provide smaller counter-examples and ease understanding of shortcomings in the underlying model. This helped us to identify the cause of several bugs in the Stuttgart 21 model and in one of the published landing gear case studies (see Section 3.2).

Furthermore, this two-phase analysis can be used to detect bugs in PROB: if the search for a counter-example fails to explore certain cases, it might be independent of the goal. Hence, we can detect if PROB correctly spots contradictions using crafted sequents. In fact, we did detect an error in a prototypical optimisation (common-subexpression elimination), which we did not use in this paper. We could even go further and apply other provers to the unsat core generated by PROB in order to validate a proof effort by a second toolchain.

3 Empirical Evaluation and Comparison

In this section, we compare PROB to several other provers available for the Rodin platform [3], i.e., Rodin’s automatic tactic and the SMT plug-in [14,15].

Our evaluation leads us to the following conclusions:

- In many cases PROB can discharge proof obligations that cannot be discharged by other provers. Each additional obligation that is discharged actually saves time and money.
- None of the provers can be replaced by the others.

⁵ Deeper within a proof, contradicting hypotheses can occur “naturally”.

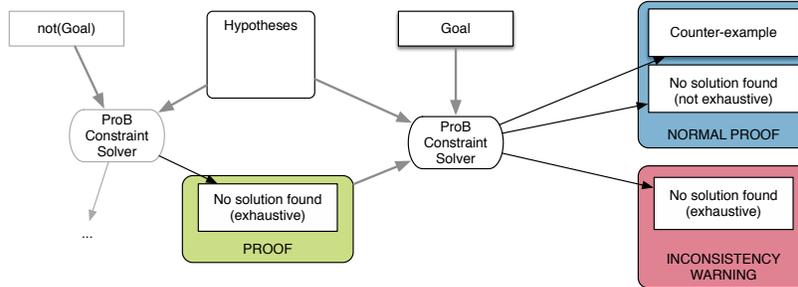


Fig. 2. Inconsistency Detection

- The performance of a prover is influenced by the surrounding tactic, including other provers. While the influence of a tactic on PROB is only marginal, it is quite strong for other provers.

3.1 Experimental Setup

For our experiments, we have used Rodin 3.1, version 2.1.0 of the Atelier B provers plugin and version 1.2.1 of the SMT plugin, with the bundled version 2.4.1 of CVC3 and the bundled development version of veriT. We have used a timeout of 5 seconds for each SMT solver, run in succession. PROB was used in version 1.5.1-beta1, connected through the disprover plugin version 3.0.8. Again, a timeout of 5 second was used for each constraint solving attempt with a maximum of two attempts per proof obligation (see Figure 1). We used a global timeout of 25 seconds for a whole tactic.

All benchmarks were run on a MacBook Pro featuring a 2.6 GHz i7 CPU and 8 GB of RAM. We did not run proof attempts in parallel to avoid issues due to hyper-threading or scheduling. We developed an evaluation plugin⁶ for the Rodin platform that applies the user- or pre-defined proof tactics to selected proof obligations.

We used the following combined tactics as they represent closely what can be utilized by end-users:

- The automatic tactic that comes with Rodin. It applies a number of rewriting rules and decision procedures to the proof tree. For instance, it checks if the goal is included in the set of hypotheses and thus discharged. The automatic tactic is applied until a fixpoint is reached or the process times out. This is the “Default Auto Tactic” of Rodin where the calls to PP and ML have been removed.
- In a second step, we used this tactic in its original state, i.e., with the PP and ML provers from Atelier B enabled.

⁶ See <https://github.com/wysiib/ProverEvaluationPlugin> for sources and instructions.

- The SMT plugin [14,15] applies two different SMT solvers (veriT [10] and CVC3 [7]) to the original goal. We used the default SMT tactic that calls PP and ML as well.
- Finally, we add PROB to the tactic as well. It is applied to the goal before the other provers.

In addition we benchmarked the provers alone, i.e. without tactics. This gives us a better picture of the individual power of each prover.

- PP and ML from Atelier B together,
- SMT plugin on its own, using both veriT and CVC3, and
- ProB alone.

We used the following models for our benchmarks:

- Answers to the ABZ-2014 landing gear case study [9]. Beside our own version [18], we also used the three models by Su and Abrial [26], a model by André, Attiogbé and Lanoix [4], as well as a model by Mammam and Laleau [24].
- A model of the Stuttgart 21 Railway station interlocking by Wiegard, derived from chapter 17 of [2] with added timing and performance modeling.
- A model of a controller area network (CAN) bus developed by Colley.
- A formal development of a graph coloring algorithm by Andriamiarina and Méry. The graphs to be colored are finite, but unbounded and not fixed in the model.
- A model of a pacemaker by Méry and Singh [25].

The models were selected so as to cover a variety of use cases. The landing gear model [18] contains mainly enumerated sets; hence we suspected PROB to perform well. We included several other versions of the case study to investigate how modelling style influenced prover performance. On the other end of the spectrum, the graph coloring model uses only deferred sets. Hence, we expected PROB not to perform well, as finite enumeration is not possible. The other models were expected to lie in between those extremes. We do not claim that our selection is representative. Indeed, we could have selected more models using (mostly) deferred sets; but this would have just confirmed that PROB's prover is disabled for proof obligations involving deferred sets.

For raw data and additional visualizations see <http://www.stups.hhu.de/ProB/index.php5/Sefm2015>. Rodin is available on <http://www.event-b.org>. The provers are available from update sites included in Rodin.⁷

3.2 Results

The benchmark results for the tactics can be found in Tables 1 and 2 and Figures 3 to 5, while the results for the provers alone are in Table 3 and part (b) of Figure 3. Table 1 shows the total number of proof obligations discharged, as well

⁷ For a standalone version of PROB see <http://www.prob2.de>.

Table 1. Benchmark results: Discharged Event-B Proof Obligations

Model	# POs	Tactic alone	+ML/PP	+ML/PP+SMT	+ML/PP+SMT+PROB
Landing Gear System 1, Su, et. al.	2328	2022	2190	2303	2306
Landing Gear System 2, Su, et. al.	1188	817	915	1169	1173
Landing Gear System 3, Su, et. al.	341	134	152	205	262
CAN Bus, Colley	534	289	398	403	388
Graph Coloring, Andriamiarina, et. al.	254	122	166	170	169
Landing Gear System, Hansen, et. al.	74	64	65	67	74
Landing Gear System, Mammar, et. al.	433	218	297	381	397
Landing Gear System, Andre, et. al.	619	180	214	319	450
Pacemaker, Neeraj Kumar Singh	370	258	354	364	369
Stuttgart 21 interlocking, Wiegard	202	37	33	97	147

Table 2. Benchmark results: Event-B Average Runtimes (in seconds / po)

Model	Tactic alone	+ML/PP	+ML/PP+SMT	+ML/PP+SMT+PROB
Landing Gear System 1, Su, et. al.	0.23	0.35	0.3	0.55
Landing Gear System 2, Su, et. al.	0.34	0.64	0.74	0.79
Landing Gear System 3, Su, et. al.	8.29	9.71	11.08	6.01
CAN Bus, Colley	5.29	5.93	6.03	7.13
Graph Coloring, Andriamiarina, et. al.	1.48	2.56	7.44	8.04
Landing Gear System, Hansen, et. al.	0	2.1	2.7	0.2
Landing Gear System, Mammar, et. al.	1.68	2.02	2.05	2.39
Landing Gear System, Andre, et. al.	11.64	11.89	11.92	7.01
Pacemaker, Neeraj Kumar Singh	0	0.1	0.04	0.4
Stuttgart 21 interlocking, Wiegard	11.7	13.26	13.2	9.84

as the percentage of proof obligations discharged using ML/PP together with SMT and in the last column the percentage discharged by using these two proof tactics together with the PROB disprover. Each Venn diagram shows how many proof obligations are discharged by which prover. Table 2 shows the runtimes of the different provers for all proof obligations and for discharged proof obligations individually. Note that for the Stuttgart 21 model and the Andre et al. model, PROB found several unprovable proof obligations, i.e., errors in the model as can be seen in Table 3. E.g., for Stuttgart 21 PROB found a counter-example for two proof obligations, while it found five counter-examples in the landing gear model. This is very useful feedback to the developer of the model, and the initial purpose of the PROB disprover.

The diagram in Figure 3 shows the gain of using PROB in addition to the other decision procedures. Compared to the SMT Tactic, adding PROB leads to an additional 304 (238+1+11+54) proof obligations being discharged. However, due to the time consumption by PROB, 47 (35+7+5) proof obligations cannot be discharged anymore. With a higher time-out, these could again be proven. The second diagram in Figure 3 shows how the individual provers alone contribute: Each of them has a set of proof obligations that cannot be solved by any of the others (192 for ML/PP, 43 for SMT and 2000 for PROB).

Except for the graph coloring algorithm PROB performs surprisingly well. The graph coloring algorithm uses unbounded sets, meaning that some of the proof obligations cannot be proven using constraint solving and enumeration.

Table 3. Results of running Provers alone (without pre-processing by Rodin)

Model	# POs	ML/PP	SMT	PROB	
				prove	disprove
Landing Gear System 1, Su, et. al.	2328	1396	1477	2311	0
Landing Gear System 2, Su, et. al.	1188	341	567	1176	0
Landing Gear System 3, Su, et. al.	341	99	146	290	0
CAN Bus, Colley	534	481	282	276	0
Graph Coloring, Andriamarina, et. al.	254	90	97	0	0
Landing Gear System, Hansen, et. al.	74	70	59	74	0
Landing Gear System, Mammar, et. al.	433	227	257	400	0
Landing Gear System, Andre, et. al.	619	189	268	567	5
Pacemaker, Neeraj Kumar Singh	370	356	224	354	0
Stuttgart 21 interlocking, Wiegard	202	51	44	125	2

As can be seen in Table 1, adding PROB improves the results of automatic proving for all other models. In some cases, such as the landing gears, the improvement is substantial (cf., Figure 4). The reason for the rather big improvement is that these models only use enumerated sets, booleans and integers as base types. In these cases PROB can produce elaborate case distinctions, combined with constraint solving to narrow down the search space. This type of proof is not supported by the classical provers ML and PP. Generally, the proof obligations that pose problems to PROB are certain well-definedness proof obligations. For instance, function application requires to prove that the parameter is in the domain of the function. Usually this leads to expensive enumeration of the possible parameter values.

For some of the models, using PROB slows down the prove process. As shown in Table 2 PROB’s runtime is above average for some proof obligations, while it considerably speeds up other proof attempts. We suspect that this is due to the multiple constraint solver calls PROB performs on different sets of hypotheses

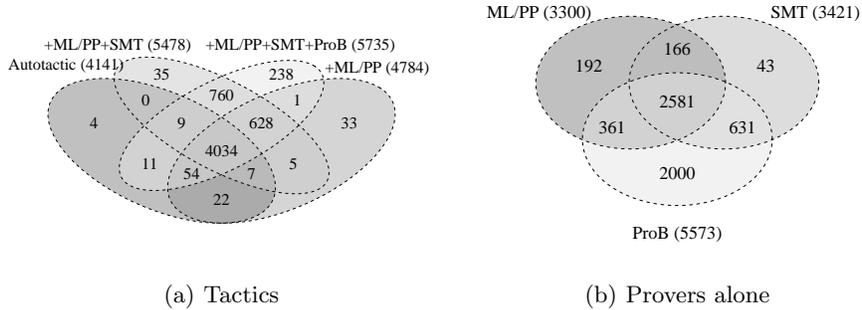


Fig. 3. Visualization of the full benchmark results

Table 4. Performance of provers on different kinds of proof obligations

Kind of PO	# POs	ML/PP	SMT	PROB
Feasibility of non-det. action	59	53 (89.8 %)	40 (67.8 %)	44 (74.6 %)
Guard strengthening	300	27 (9.0 %)	13 (4.3 %)	258 (86.0 %)
Invariant preservation	4938	2877 (58.3 %)	3111 (63.0 %)	4488 (90.9 %)
Action simulation	153	118 (77.1 %)	108 (70.6 %)	134 (87.6 %)
Theorem	97	13 (13.4 %)	29 (29.9 %)	66 (68.0 %)
Well definedness	779	200 (25.7 %)	109 (14.0 %)	570 (73.2 %)
	6326	3288 (52.0 %)	3410 (53.9 %)	5560 (87.9 %)

as shown in Figure 1. Also, PROB is looking for proofs and counter-examples. This often means that PROB will continue the computation, even after it has realized that no proof is possible (in the hope of finding a counter-example).

It is also interesting to note that, on their own, the ML and PP provers do not fare quite so well as in Table 1: they require pre-processing and tactic support to be fully effective: See Table 3 containing the results without any pre-processing.

All models except the Landing Gear System by Mammar, et. al. show the same behavior: The rate of discharged proof obligations drops significantly if Rodin’s default tactics are not applied. Adding SMT solvers or PROB does not replace the tactics either.

In contrast, the model by Mammar, et. al. shows the opposite behavior: without pre-processing, more proof obligations can be discharged. This is probably due to the timeouts leaving less time for the actual prover, if we include a pre-processing phase. In future, we want to examine whether better pre-processing can improve the performance of the PROB disprover.

The same effect can be observed in Table 4. Here, the performance of the provers on different kinds of proof obligations is given. For most kinds, PROB does perform quite well when compared to ML/PP and the SMT solvers, especially for guard strengthening proofs, theorem proofs and well-definedness proofs. For feasibility and finiteness proof obligations, on the other hand, PROB fares less well.

Unexpected Performance of SMT To our surprise, the SMT solvers did not perform as well as we expected when compared to PROB. For certain kinds like guard strengthening or initialization in Table 4, the SMT solvers prove less proof obligations than ML/PP or PROB. We suspect that this is due to the translation from Event-B to SMT-LIB:

- The λ -based approach [14,15] does not support sets of sets. Thus, a whole class of proof obligations cannot be solved by it. Therefore, the SMT plugin uses the second approach presented in [14,15] as the default:
- The *ppTrans* approach [19] translates set theory to predicate calculus. The resulting SMT-LIB problem is then enriched by the predicate calculus version of certain set-theoretic axioms.

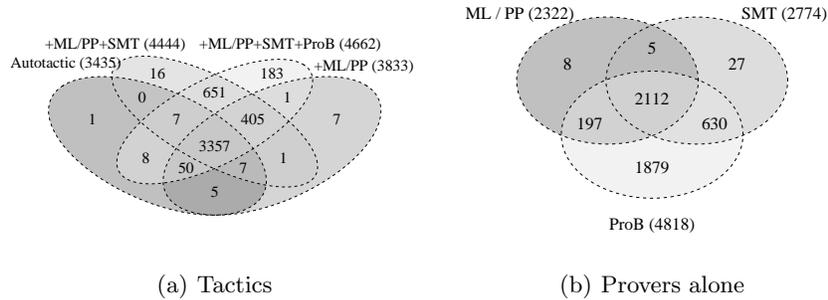


Fig. 4. Visualization of the benchmark results. Part 1: Landing gear systems

Newer releases of SMT-Solvers like CVC4 [5] support finite sets natively as an extension to the SMT-LIB language [27]. Thus, certain classes of proof obligations could be passed to the SMT-Solvers directly instead using one of the approaches mentioned above. We assume that this would increase the number of proof obligations that could be discharged successfully. In summary, while the SMT plugin has been very successful, we recommend critically examining the current SMT-LIB translation and believe there is scope for considerable improvement by using an alternate translation.

Inconsistency in Hypothesis Detection The inconsistency detection of Section 2.4 found also various contradictions in the theorems (at lower refinement levels) of the Stuttgart 21 model. It also highlighted an issue in the first development of the ABZ landing gear from [26]. The PROB disprover was flagging, e.g., the proof obligation `treat_hndl_up_112/inv1/INV` in the machine LPN4 as containing a contradiction in the hypothesis. The PROB unsat core algorithm found out the following root cause:

```
close_EV = FALSE & open_EV = FALSE & door = op2c1 &
((open_EV = FALSE & close_EV = FALSE) => door = c1)
```

The seen context LPNC0 contains the axiom `partition(D, {c1}, {c12op}, {op2c1}, {op})`.⁸ The first line comes from the guard of the event `treat_hndl_up`, the second line is the invariant `inv1` from LPN4. In other words, the disprover has detected that this event can never be executed given the invariant. A similar issue was detected for several other events.⁹

When (not) to use the PROB Disprover In summary, we present the following insights on when to use the PROB disprover (+) and when not to (-):

⁸ For technical reasons this axiom is not yet included in the unsat core; partition axioms are never removed from the core by the current algorithm.

⁹ In LPN4 of [26]: `treat_hndl_up_122`, `treat_hndl_up_132`, `treat_hndl_dn_112`, `treat_hndl_dn_122`, `treat_hndl_dn_132`.

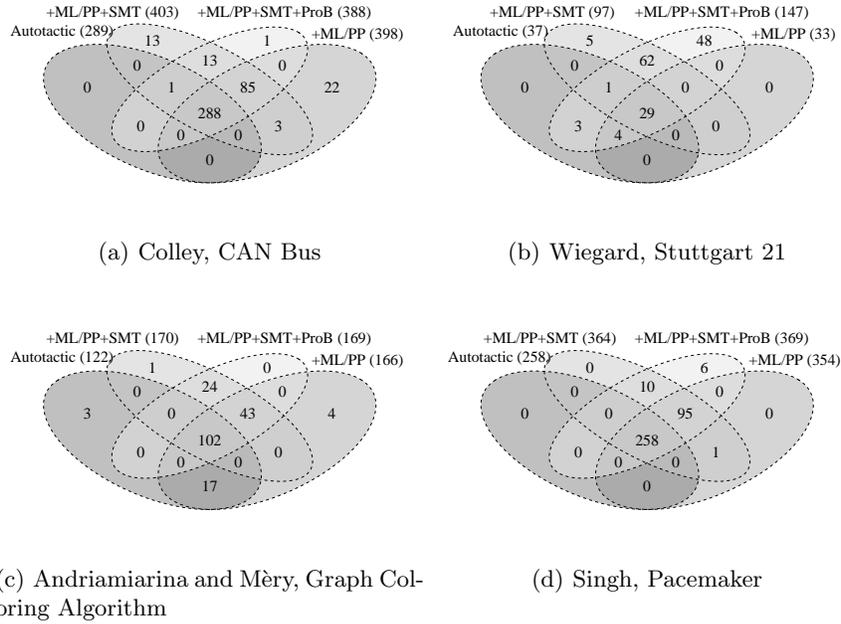


Fig. 5. Visualization of the benchmark results. Part 2: Miscellaneous models

- + Used solely as a disprover, PROB can prevent futile interactive proof attempts. This is always worthwhile.
- + The inconsistency detection is very useful for finding subtle modelling errors.
- + On models such as the ABZ landing gear models (Figure 4), which rely heavily on enumerated sets, booleans and/or bounded integers as base types, PROB performs very well.
- + The Stuttgart 21 model shows that explicit data, e.g., track layouts or time tables, can often be used effectively by PROB. Often, this results in a proof by an elaborate case distinction.
- + PROB performs reasonably well on unbounded intervals, when interval reasoning can be applied. This occurs for example in the lower refinement levels of the ABZ case study models or the pacemaker model.
- As soon as the proof goal references deferred sets (e.g., in the graph coloring model), no proof can be done by construction of the disprover (see Figure 1).
- When unbounded datastructures are used, PROB cannot exhaustively enumerate cases and is much less powerful. This happens for example in the CAN bus model, that represents a buffer as an unbounded partial function from \mathbb{N} to \mathbb{Z} .

4 Discussion and Conclusion

One motivation for the experiments conducted in this paper was the empirical evaluation of our constraint solver, more precisely its capability to detect inconsistencies (a successful proof with the disprover requires finding a contradiction without enumerating unbounded variables; see Figure 1). Finding inconsistencies is important for many other features of PROB, e.g., detecting disabled events during animation. Furthermore, it is useful for constraint-based validation, such as deadlock checking [17], where it avoids the constraint solver exploring infeasible alternatives. In the context of model-based testing, it enables PROB to detect uncoverable alternatives, and not spend time trying to find test cases for them.

An important issue is the soundness of the PROB disprover. In [8] we have presented the various measures we are taking to validate PROB’s results in general. In addition, we have developed an SMT-LIB [6] importer for PROB and have applied our disprover to a large number of SMT-LIB benchmarks, checking that no “false theorems” are proven. For this paper, we have also double checked many of the proof obligations which were only provable by PROB, to ensure that they are indeed provable. As the Venn diagrams in Figures 4 and 5 show, a large number of proof obligations can be proven by two or even three different provers. As the three provers rely on completely different technologies and have been developed by independent teams, we can have a very high confidence that those proof obligations are indeed provable.

We have demonstrated that constraint-based proof in general, and PROB in particular, is capable of discharging proof obligations that currently cannot be proven using Rodin’s auto tactic and the SMT solvers. Our prover typically deals well with a different kind of proof obligations than the other provers, and is thus an orthogonal extension rather than a replacement. Rodin’s auto tactic performs well in the realm of set theoretic constructs and relational expressions, some of which cannot be easily represented in the SMT syntax. SMT on the other hand performs well on arithmetic expressions, where the auto tactics often fail. ProB finally covers predicates over enumerated sets, explicit data and explicit computations and has a good support for integer arithmetic over finite domains.

However, for models which make heavy use of deferred sets, such as the graph coloring algorithm model, the PROB disprover can currently mainly play its role as disprover. More precisely, for any proof obligation which involves deferred sets and where no precise value of the cardinality of the deferred set is known, the disprover can only return either a counter-example or the result “unknown”.

In future, we plan to improve the treatment of deferred sets in PROB, and to have the constraint solver determine the cardinalities of those sets while solving.

We also plan to conduct experimental evaluations for PROB within Atelier B, and compare with efforts such as [13] or the BWare project [16]. First results on industrial case studies within Alstom are already very promising.

We think that the PROB Disprover is a valuable extension to the existing set of provers, because it can increase the number of proof obligations that are automatically discharged, thus saving time and money. Overall, the outcome

of the empirical evaluation was a positive surprise, as PROB's main domain of application is finding concrete counter-examples, not discharging proof obligations. In particular, the fact that the number of discharged proof obligations (5573 in Figure 3 (b)), for the models under consideration, is better than that of the two SMT solvers of the SMT plugin (3421 in Figure 3 (b)) was completely unexpected.

Acknowledgements We would like to thank the various developers for giving us access to their Event-B models, and for discussions and feedback: Jean-Raymond Abrial, Andre, Attiogbe, John Colley, Régine Laleau, Luis-Fernando Mejia, Lanoix, Amel Mammam, Dominique Méry, Neeraj Kumar Singh, Wen Su.

References

1. J.-R. Abrial. *The B-Book*. Cambridge University Press, 1996.
2. J.-R. Abrial. *Modeling in Event-B: System and Software Engineering*. Cambridge University Press, 2010.
3. J.-R. Abrial, M. Butler, and S. Hallerstede. An open extensible tool environment for Event-B. In Z. Liu and J. He, editors, *Proceedings ICFEM'06*, LNCS 4260, pages 588–605. Springer-Verlag, 2006.
4. André, Attiogbé, and Lanoix. Modelling and Analysing the Landing Gear System: a Solution with Event-B/Rodin. <http://www.lina.sciences.univ-nantes.fr/aelos/software/LGS-ABZ2014/index.php>. Solution to ABZ-2014, Accessed: 2014-03-17.
5. C. Barrett, C. L. Conway, M. Deters, L. Hadarean, D. Jovanović, T. King, A. Reynolds, and C. Tinelli. Cvc4. In *Proceedings of the 23rd International Conference on Computer Aided Verification, CAV'11*, pages 171–177, Berlin, Heidelberg, 2011. Springer-Verlag.
6. C. Barrett, A. Stump, and C. Tinelli. The SMT-LIB Standard: Version 2.0. In A. Gupta and D. Kroening, editors, *Proceedings of the 8th International Workshop on Satisfiability Modulo Theories*, 2010.
7. C. Barrett and C. Tinelli. CVC3. In W. Damm and H. Hermanns, editors, *Proceedings of the 19th International Conference on Computer Aided Verification*, volume 4590 of *Lecture Notes in Computer Science*, pages 298–302. Springer-Verlag, July 2007. Berlin, Germany.
8. J. Bendispoto, S. Krings, and M. Leuschel. Who watches the watchers: Validating the prob validation tool. In *Proceedings of the 1st Workshop on Formal-IDE, EPTCS XYZ*, 2014. Electronic Proceedings in Theoretical Computer Science, 2014.
9. F. Boniol and V. Wiels. The landing gear system case study. In F. Boniol, V. Wiels, Y. Ait Ameer, and K.-D. Schewe, editors, *ABZ 2014: The Landing Gear Case Study*, volume 433 of *Communications in Computer and Information Science*, pages 1–18. Springer International Publishing, 2014.
10. T. Bouton, D. C. B. de Oliveira, D. Déharbe, and P. Fontaine. verit: an open, trustable and efficient smt-solver. In R. A. Schmidt, editor, *Proc. Conference on Automated Deduction (CADE)*, Lecture Notes in Computer Science, pages 151–156. Springer-Verlag, 2009.
11. M. Carlsson and G. Ottosson. An Open-Ended Finite Domain Constraint Solver. In H. G. Glaser, P. H. Hartel, and H. Kuchen, editors, *Proceedings PLILP'97*, LNCS 1292, pages 191–206. Springer-Verlag, 1997.

12. ClearSy. *Atelier B, User and Reference Manuals*. Aix-en-Provence, France, 2009. Available at <http://www.atelierb.eu/>.
13. S. Conchon and M. Iguernelala. Tuning the alt-ergo SMT solver for B proof obligations. In *Abstract State Machines, Alloy, B, TLA, VDM, and Z - 4th International Conference, ABZ 2014, Toulouse, France, June 2-6, 2014. Proceedings*, pages 294–297, 2014.
14. D. Déharbe. Automatic Verification for a Class of Proof Obligations with SMT-Solvers. In *Proceedings ASM 2010*, pages 217–230, 2010.
15. D. Deharbe, P. Fontaine, Y. Guyot, and L. Voisin. SMT solvers for Rodin. In *Proceedings ABZ'2012*, LNCS 7316, pages 194–207. Springer, 2012.
16. D. Delahaye, C. Dubois, C. Marché, and D. Mentré. The bware project: Building a proof platform for the automated verification of B proof obligations. In *Abstract State Machines, Alloy, B, TLA, VDM, and Z - 4th International Conference, ABZ 2014, Toulouse, France, June 2-6, 2014. Proceedings*, pages 290–293, 2014.
17. S. Hallerstede and M. Leuschel. Constraint-based deadlock checking of high-level specifications. *TPLP*, 11(4–5):767–782, 2011.
18. D. Hansen, L. Ladenberger, H. Wiegard, J. Bendisposto, and M. Leuschel. Validation of the abz landing gear system using prob. In F. Boniol, V. Wiels, Y. Ait Ameer, and K.-D. Schewe, editors, *ABZ 2014: The Landing Gear Case Study*, volume 433 of *Communications in Computer and Information Science*, pages 66–79. Springer International Publishing, 2014.
19. M. Konrad and L. Voisin. Translation from set-theory to predicate calculus. Technical report, ETH Zurich, 2011.
20. M. Leuschel, J. Bendisposto, I. Dobrikov, S. Krings, and D. Plagge. From animation to data validation: The prob constraint solver 10 years on. In J.-L. Boulanger, editor, *Formal Methods Applied to Complex Systems: Implementation of the B Method*, chapter Chapter 14, pages 427–446. Wiley ISTE, Hoboken, NJ, 2014.
21. M. Leuschel and M. Butler. ProB: A model checker for B. In A. Keijiro, S. Gnesi, and M. Dino, editors, *FME*, volume 2805 of *Lecture Notes in Computer Science*, pages 855–874. Springer-Verlag, 2003.
22. M. Leuschel and M. Butler. ProB: An automated analysis toolset for the B method. *Software Tools for Technology Transfer (STTT)*, 10(2):185–203, 2008.
23. O. Ligtot, J. Bendisposto, and M. Leuschel. Debugging Event-B Models using the ProB Disprover Plug-in. *Proceedings AFADL'07*, Juni 2007.
24. A. Mammarr and R. Laleau. Modeling a landing gear system in event-b. In F. Boniol, V. Wiels, Y. Ait Ameer, and K.-D. Schewe, editors, *ABZ 2014: The Landing Gear Case Study*, volume 433 of *Communications in Computer and Information Science*, pages 80–94. Springer International Publishing, 2014.
25. D. Méry and N. K. Singh. Formal specification of medical systems by proof-based refinement. *ACM Trans. Embed. Comput. Syst.*, 12(1):15:1–15:25, Jan. 2013.
26. W. Su and J.-R. Abrial. Aircraft landing gear system: Approaches with event-b to the modeling of an industrial system. In F. Boniol, V. Wiels, Y. Ait Ameer, and K.-D. Schewe, editors, *ABZ 2014: The Landing Gear Case Study*, volume 433 of *Communications in Computer and Information Science*, pages 19–35. Springer International Publishing, 2014.
27. G. Weissenbacher, D. Kröning, and P. Rümmer. A proposal for a theory of finite sets, lists, and maps for the smt-lib standard. In *Proceedings of the 7th International Workshop on Satisfiability Modulo Theories, SMT'09*, 2009.