



The abc Group

---

---

## **Making Trace Monitors Feasible**

abc Technical Report No. abc-2007-1

Pavel Avgustinov, Julian Tibble, Oege de Moor  
Programming Tools Group, University of Oxford, UK

April 6, 2007

---

---

**a s p e c t b e n c h . o r g**

# Contents

<b>1</b>	<b>INTRODUCTION</b>	<b>3</b>
<b>2</b>	<b>TRACE MONITORING</b>	<b>4</b>
<b>3</b>	<b>BENCHMARKS</b>	<b>6</b>
<b>4</b>	<b>MEMORY LEAK ELIMINATION</b>	<b>9</b>
<b>5</b>	<b>INDEXING</b>	<b>12</b>
5.1	Choosing a Partition . . . . .	14
5.2	A Data-Structure for Partitioning . . . . .	15
5.3	Updating Indexed Disjuncts . . . . .	16
5.4	Evaluation . . . . .	18
<b>6</b>	<b>RELATED WORK</b>	<b>19</b>
<b>7</b>	<b>CONCLUSIONS</b>	<b>21</b>

## List of Figures

1	An automaton with states labelled by constraints . . . . .	6
2	An annotated automaton for safe-enumeration. . . . .	13
3	The calculations performed when using simple sets to store disjuncts. . . . .	13
4	The new constraints after a NEXT event has occurred. . . . .	13
5	The safe-enumeration constraints, as stored using indexing. . . . .	15
6	Negative update tree for the safe-enumeration example. . . . .	17
7	Positive update tree for the safe-enumeration example. . . . .	17

## Abstract

A *trace monitor* observes an execution trace at runtime; when it recognises a specified sequence of events, the monitor runs extra code. In the aspect-oriented programming community, the idea originated as a generalisation of the advice-trigger mechanism: instead of matching on single events (joinpoints), one matches on a sequence of events. The runtime verification community has been investigating similar mechanisms for a number of years, specifying the event patterns in terms of temporal logic, and applying the monitors to hardware and software.

In recent years trace monitors have been adapted for use with mainstream object-oriented languages. In this setting, a crucial feature is to allow the programmer to quantify over groups of related objects when expressing the sequence of events to match. While many language proposals exist for allowing such features, until now no implementation had scalable performance: execution on all but very simple examples was infeasible.

This paper rectifies that situation, by identifying two optimisations for generating *feasible* trace monitors from declarative specifications of the relevant event pattern. We restrict ourselves to optimisations that do not have a significant impact on compile-time: they only analyse the event pattern, and not the monitored code itself.

The first optimisation is an important improvement over an earlier proposal in [1] to avoid space leaks. The second optimisation is a form of indexing for partial matches. Such indexing needs to be very carefully designed to avoid introducing new space leaks, and the resulting data structure is highly non-trivial.

## 1 INTRODUCTION

*Trace monitors* observe the current execution trace, and execute some extra code when the trace matches a given pattern. Many runtime verification concerns can be expressed as trace monitors very naturally, simply by picking out violating traces.

A trace monitor is usually specified declaratively in two parts: firstly, a pattern describing which traces should match, and secondly, an action that should be executed when a program trace matches. The actual implementation of the trace monitor is automatically generated from its specification.

There is a very large amount of previous research on this topic, *e.g.* [1, 6, 7, 9–13, 16, 18, 19, 22, 23, 25, 27]. These studies range from applications in medical image generation through business rules to theoretical investigations of the underlying calculus. The way the patterns are specified varies, and temporal logic, regular expressions and context-free languages have all been considered.

One theme shines through all of these previous works: trace monitors are an attractive, useful notion, worthy of integration into a mainstream programming language. This has not happened, however, because it turns out to be very difficult to generate efficient code when the trace monitor is phrased as a declarative specification.

The challenge is particularly severe when the specifications trace the behaviour of a group of several objects simultaneously. For example, when checking that a lock is always acquired and released in the same method invocation, we need to track the locked resource, as well as the method invocation in question. Similarly, when checking that a collection is not modified while iteration is in progress, we need to track the state of the collection as well as its iterators. For this reason, numerous of the above proposals (in particular [1, 9, 11, 18, 22]) have a facility for *binding* multiple variables during the matching process. Many more systems do not allow free variables in the trace patterns; in our experience, this feature is indispensable for many real-world examples, and therefore the focus of this work is optimising systems *with* free variables.

In our experiments, however, none of these systems (including our own in [1]) managed to generate *feasible* trace monitors. The generated monitors are adequate as a proof of concept, but they cannot be used in practice on substantial systems.

For some applications, it may be possible to use the substantial body of related work on *static* type-state verification, *e.g.* [15]. However, that invariably entails interprocedural analysis of the observed code. This is costly, and the assumption that all the code is available before load-time cannot be satisfied in practice.

Furthermore, before such costly techniques are explored, we must first determine how far one can get with cheaper techniques.

**Contribution** This paper shows, for the first time, how to generate *feasible* trace monitors from declarative specifications that use free variables. Furthermore, our techniques rely only on analysis of the trace specification, not on costly whole-program analyses of the monitored code. Any monitoring system, regardless of the chosen specification formalism, must implement the techniques presented here to achieve feasibility.

Specifically, the detailed contributions are as follows:

**Benchmarks:**

- We present the first benchmark set of substantial, realistic applications of runtime monitoring.
- Each of these monitors has been coded by hand (in the programming language AspectJ), and also in the specification formalism of [1].
- This set of benchmarks (which is publicly available) thus provides the first solid basis for experimental evaluation of trace monitoring features.

**Leak detection and prevention:**

- We demonstrate that space leaks are a show-stopping bottleneck when naively generating trace monitors.
- In [1] an analysis was briefly sketched for eliminating that problem. We identify a crucial flaw in that analysis, and show how it can be remedied via the novel notion of *persistent weak references*.
- Unlike [1], we then proceed to carefully evaluate the effects of the optimisation.

**Indexing of partial solution sets:**

- We exhibit another show-stopping performance problem, namely the need to update the set of partial solutions whenever a relevant event occurs.
- We propose an automatic technique for choosing an appropriate index structure on the set of partial solutions, which is purely based on the monitor specification.
- Furthermore, we present an indexing data structure that does not introduce new space leaks, and thus combines well with the above leak prevention technique. Again this involves very careful and subtle use of weak references: naive indexing would worsen space leakage.
- Finally, we provide a set of algorithms to update such indexing data structures (using a notion of *tree patches*), and present an argument for why these algorithms are correct.

## 2 TRACE MONITORING

We first outline the basic strategy for generating executable code from trace monitors, via a specific example. Variations of the same code generation strategy can be found in any trace monitoring system, but our primary focus is on tracematches [1] — the system that this work evaluates and optimises.

As previously mentioned, a trace monitor is a combination of a pattern and an action to run when the program trace matches that pattern. In a tracematch, the pattern is defined in three parts: a set of variables, an alphabet of symbols, and a regular expression. The symbols are defined in terms of the variables, and the regular expression is over the alphabet of symbols.

A common runtime verification property is that of safe-enumeration:

After an enumeration is created, the data-source upon which it is based may not be modified while the enumeration is in use — that is, until the last call to its `nextElement()` method.

To check this property with a tracematch, two variables are required. One of them will range over collections that might be iterated; in the Java 1.2 API these are instances of the class `Vector`, and so we will use the identifier  $v$ . The second variable,  $e$ , ranges over enumerations. Symbols are defined using AspectJ pointcuts — a language for intercepting runtime events [20] — but for clarity we give an informal definition of the three required symbols.

CREATE		an enumeration $e$ is created from a Vector $v$
UPDATE		the Vector $v$ is modified
NEXT		the <code>nextElement()</code> method is called on $e$

Finally, the regular expression which specifies a violation of the safe-enumeration property is ‘CREATE NEXT\* UPDATE+ NEXT’.

Let us now consider the following question: What does it mean for such a regular expression (which includes free variables) to match the program trace? Consider the following execution history:

Symbol	$v$	$e$	$v = v_2$	$v = v_1$
			$e = e_3$	$e = e_1$
CREATE	$v_1$	$e_1$		CREATE
CREATE	$v_1$	$e_2$		
NEXT		$e_2$		
NEXT		$e_1$		NEXT
UPDATE	$v_1$			UPDATE
CREATE	$v_2$	$e_3$	CREATE	
NEXT		$e_3$	NEXT	
NEXT		$e_1$		NEXT
UPDATE	$v_2$		UPDATE	

(\*)

The left side of the table above shows a sequence of symbol-matching events from a program trace. For each event, the symbol that it corresponds to is shown, together with the objects bound to  $v$  and  $e$ .

On the right side of the table, we show two substitutions of objects for tracematch variables. Each substitution defines a projection of the original sequence to a string over the symbol alphabet, found by including the symbols from lines on the left of the table that have compatible variable bindings.

The tracematch action is triggered if there is some substitution for which the regular expression matches a suffix of the projected string. For example, in the table above, the substitution  $v = v_1 \wedge e = e_1$  results in the projected string `CREATE NEXT UPDATE NEXT`, which matches the regular expression. The action is triggered *at the event* which completed the match (*i.e.* the line marked (\*), not the line after, even though the projected string still matches).

When a tracematch is compiled, the regular expression is translated into a finite state automaton. Such automata are well understood, and can be easily constructed from regular expressions. Variants also exist for other formalisms such as context-free grammars and linear temporal logic. In particular, automata are much better suited to *online* matching, which is why they are used here instead of keeping a regular expression representation.

As we saw above, matches with different variable bindings should be independent of each other and may be interleaved. This means that, conceptually at least, a separate finite automaton is required for every possible binding — an unbounded number. However, it is possible to simulate this collection of automata by using a single automaton and labelling its states with constraints. An example of such a labelled automaton for the safe-enumeration property is shown in Figure 1. A constraint  $C$ , labelling a state  $i$ , is interpreted as follows: “there is an automaton in state  $i$  for each variable binding that satisfies  $C$ ”. Constraints are boolean expressions, consisting of  $(x = v)$ ,  $(x \neq v)$ ,  $\wedge$ ,  $\vee$ , and  $\neg$  (for tracematch variables  $x$  and objects  $v$ ), and are stored in disjunctive normal form.

For the interested reader, the actual source for this monitor is shown below, but the focus of this paper is code generation, not language design.

```

1 tracematch(Vector v, Enumeration e) {
2   sym create_enum after returning(e) :
```

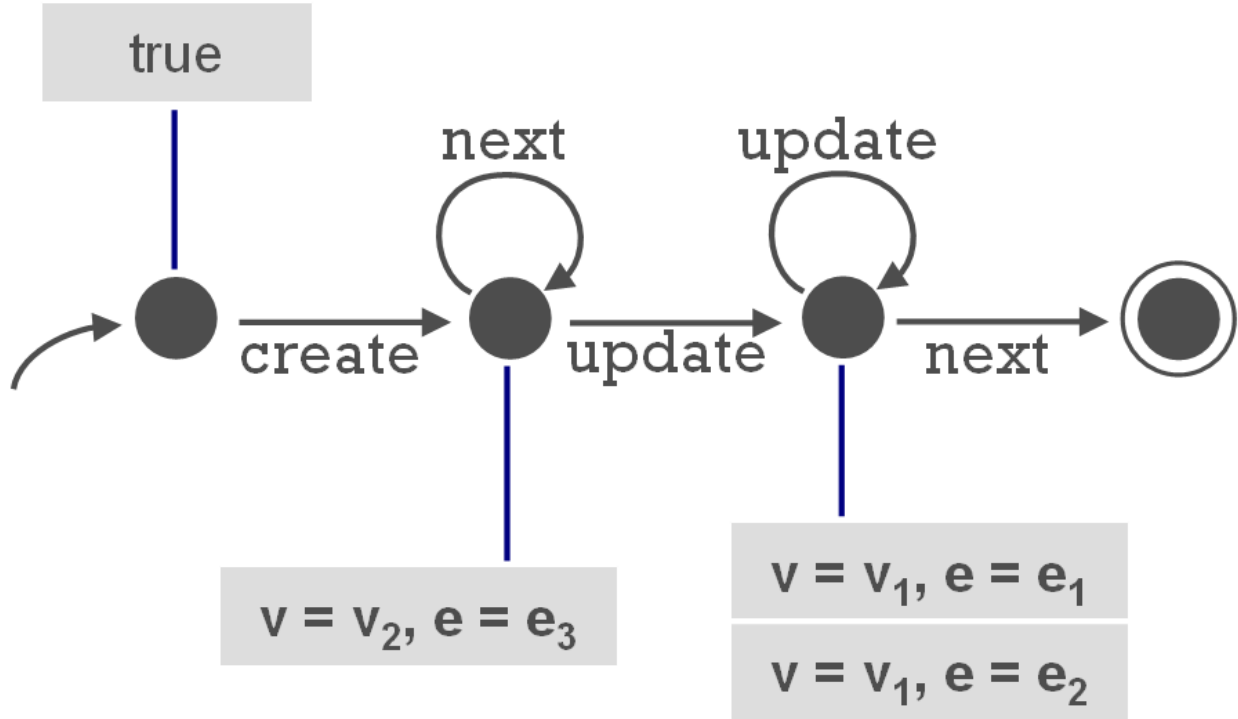


Figure 1: An automaton with states labelled by constraints

```

3     call(Enumeration+.new(..) && args(v);
4     sym call_next before :
5     call(Object Enumeration.nextElement()) && target(e);
6     sym update_source after :
7     vector_update() && target(v);
8
9     create_enum call_next* update_source+ call_next
10
11    {
12        throw new ConcurrentModificationException();
13    }
14 }

```

Lines 2–7 declare the individual event patterns, Line 9 contains the regular pattern that is matched against the current trace, and Lines 11–13 constitute the extra code that is executed upon a successful match. A detailed discussion of the pertinent language design decisions can be found in [1].

### 3 BENCHMARKS

We now present our benchmark collection. In selecting benchmarks, our main inspiration was the literature on runtime verification, and also on trace-based aspect-oriented programming. We scrupulously included base programs where the overheads are likely to be substantial. An alternative, taken in [8], is to apply a single tracematch to many unrelated base programs. This allows research into cheap strategies for eliminating inapplicable monitors (indeed [8] demonstrates such a strategy). However, it is not appropriate in this context because it would wrongly give the impression that overheads are low, simply because the benchmarks may

not heavily engage in the monitored events. Also, we provide a ‘gold standard’ for each benchmark, consisting of a hand-coded best possible solution.

To our knowledge, this is the first collection of trace monitoring benchmarks especially constructed to cover a wide variety of properties, while highlighting potential overheads. It is publicly available for others to use in comparative experiments. Indeed JavaMOP has recently adopted some of our examples for their own system [14], confirming that these benchmarks are not tailored just for tracematches: they set a standard for all runtime verification systems.

Below we first informally describe the benchmarks, each consisting of a (monitor, base program) pair. Then we define the gold standard to compare against, namely highly hand-optimised AspectJ implementations. Finally we take our initial measurements, using the naive implementation of tracematches to determine where the main performance problems are.

## Monitors and base programs

**SAFEENUM** This is the example mentioned in the previous section: we would like to throw an exception when a collection is modified while an enumeration over that collection is in progress. We apply it to JHotDraw [17], a Java graphics package that allows the user to animate a drawing’s components in a visually appealing manner.

The animation loop contains a call to `Thread.sleep` to slow down movement; we removed that to truly observe the overheads incurred by trace monitoring, and also minimised the window to factor out the cost of display operations. The effect of these measures is that enumerating the figure elements at each frame dominates the cost of the animation loop, and one would thus expect overheads to be very large indeed. JHotdraw’s use of enumerations is actually unsafe because one can edit the drawing while animation is in progress, and when that is done, our monitor catches the violation.

**NULLTRACK** This is a debugging concern. We aim to provide a trace monitor that can help track down the cause for a `NullPointerException`. We do this by intercepting all field writes where a field  $f$  is set to **null**, followed by no intervening assignment of a non-**null** value until we see a read of  $f$  followed by a null pointer exception — the trace monitor then reports all fields that were set to **null**, and where that happened. (It is necessary to report all fields that were set to **null** and not just a single one, because the exception may be caused by a method returning one of these fields, and then dispatching on the result.)

Note that this is intrinsically a very expensive monitor, because it involves a large number of instrumentation points: all field reads and writes, and all method calls. We applied it to `CertRevSim`, a discrete event simulator used to evaluate the performance of various certificate revocation schemes.

**HASHCODE** Another common runtime verification concern, this trace monitor checks the property that whenever an object is stored in a data structure indexed by its hash code (*e.g.* a `HashSet`), and subsequently a membership test is executed, the object’s hash code hasn’t changed in the meantime (if it has, we could get false negatives). This is expensive to check at runtime for a different reason: if the base program makes heavy use of hash-based data structures, the trace monitor will have to keep track of very many partial matches, one for each object stored.

We applied this to two different base programs: Firstly, `AProVE` [2], a termination prover for term-rewriting systems which we chose for its heavy use of hash sets, and also because it is precisely the type of complex, large application where this tracematch can be helpful. We also applied `HASHCODE` to `Weka` [28], a machine learning library that makes extensive use of `HashMaps`, but which is substantially smaller than `AProVE`.

**OBSERVER** The observer pattern is a popular example in the aspect-oriented programming community. Whenever a subject changes its state, all registered observers must be notified; with AOP (and with trace monitoring) it is possible to achieve this without the subject explicitly making such notifications.

The base program here is `AJHotDraw` [24], an aspect-oriented rewrite of `JHotDraw`, which uses the observer pattern for its display updates. Again, this benchmark has been chosen to fairly measure



worst-case overheads: in AJHotDraw, each subject has precisely one observer; there is thus considerable cost involved in using a data structure that caters for multiple observers, without knowledge of the one-to-one correspondence. Again, we removed all delays from the animation loop and minimised the window to factor out display operation overheads.

**DBPOOLING** This is an example proposed by Laddad in his AspectJ text book [21]. The idea is that in a particular database application, creating and establishing connections is the most expensive operation, so, whenever possible, we want to pool existing connections and reuse them, rather than creating new ones.

The base program in this case is artificial — a slight modification of the example in Laddad’s book, simply connecting to the database multiple times and performing database operations. The AspectJ version we used is Laddad’s. Note that for this example, we expect the trace monitor to *improve* performance rather than hinder it, since it would prevent unnecessary connections from being established.

**LUINMETH** Here we are concerned with checking a stylistic rule. Whenever a lock is acquired, it should be released before the enclosing method returns. Checking this rule with a trace monitor was proposed as a motivating example by the authors of PQL [22]. It is particularly interesting because it requires matching method entry/exit pairs at runtime, which is in general a context-free language problem. However, making judicious use of variable bindings, it turns out that it is possible to express this using weaker formalisms; we encoded it in tracematches, which only allow regular expressions as patterns.

The base program we chose here is Jigsaw [26], the W3C’s leading edge web server platform. It makes frequent use of locking and unlocking, so there are plenty of points of interest. In fact, Jigsaw violates the style rule being checked, and our monitor catches the violations.

**REWEAVE** The *abc* compiler makes use of an optimisation phase termed *reweaving*, during which the effects of weaving aspect code are undone and weaving is repeated using more precise analysis results obtained on the previously woven code [4]. Of course the correctness of this relies on properly undoing the effects of the first weaving step, so that the reweaving process can start from a clean state. In terms of trace monitoring, if a field is written to during the execution of the *weave()* method, is *not* written to during unweaving, and is read during reweaving, then it is very likely that it was not reset properly, and an error should be reported.

We applied this to the *abc* compiler itself, and used the instrumented version to compile a small program to obtain our numbers.

ID	MONITOR	BASE	KSLOC	NONE	ASPECTJ	NOOPT
1	SAFEENUM	JHOTDRAW	9.5	5.0s	5.3s	>90M
2	NULLTRACK	CERTREVSIM	1.4	0.17s	0.47s	47.27s
3	HASHCODE	APROVE	438.7	345.0s	485.8s	>90M
4	HASHCODE	WEKA	9.9	2.8s	2.9s	47.7
5	OBSERVER	AJHOTDRAW	21.1	6.5s	7.0s	41.3s
6	DBPOOLING	ARTIFICIAL	<0.1	70s	3.4s	3.7s
7	LUINMETH	JIGSAW	100.9	13.6s	18s	21.9s
8	REWEAVE	ABC	51.2	4.5s	5.4s	27.9s

Table I: Run times in seconds.

**Gold standard AspectJ implementations** We implemented the trace monitoring concern for each of these benchmarks as a tracematch, and also created a hand-coded, hand-optimised set of plain AspectJ aspects that manually achieve the goal. We did this in order to have a ‘gold standard’ for each benchmark, the best possible implementation an optimising compiler might aim for.

In all cases, the AspectJ code is much longer and more complex than the tracematch specification. For example, it often makes nifty use of weak references, and sometimes it makes additional assumptions not

available to the compiler (for instance that all collections being enumerated have been created in user code and not in libraries). In short, the AspectJ version is what an expert programmer would do for the problem in hand.

Discussing the entire set of hand-coded solutions is beyond the scope of this paper, but we may highlight just one example in detail: `SAFEENUM`. We subclass `Vector` to `MyVector`. `MyVector` keeps a version number as a field, which gets incremented upon each update operation. That increment is implemented with a piece of advice; with another piece of advice, all constructor calls on `Vector` are replaced by constructor calls on `MyVector`. It is for the replacement of those constructor calls that we need to be sure all instances of `Vector` are created within `JHotDraw` and not in some library code. An automatic test for verifying this property would require costly interprocedural analysis of the monitored code. Each implementation of `Enumeration` also has such a version number, copied from the underlying vector upon creation; this is introduced by the aspect as an intertype declaration on the `Enumeration` interface. Each time an enumeration step is taken, a piece of advice compares the version number of the enumeration with that of the underlying vector: when they are not equal, an exception is thrown. All pieces of advice need to be declared as **synchronized**.

**Measurements** The initial performance measurements are given in Table I. For this set of measurements, we disabled most tracematch optimisations, simply allowing the implementation to perform its code generation. In particular, no attention is paid to potential space leaks or to organising partial match sets; this is the approach that seems to be prevalent in the field.

Let us examine some of the trends exhibited by these numbers in detail.

We can see that the tracematch performance is very close to that of AspectJ in the case of `DBPOOLING`; this can be explained by the fact that the benchmark does expensive database processing that dominates the monitoring overheads. More surprisingly, performance is very good in the case of `LUINMETH`, which is applied to a significantly larger base program. Here, the explanation is in the tracematch pattern. Recall that it is intended to find occurrences of a lock being acquired and not released by the time the acquiring method returns. When a method that acquired a lock does return, therefore, it is clear that none of the associated partial matches will need further updates, and they are invalidated. There is no build-up of ‘live’ partial matches over the course of the benchmark.

Performance looks less promising for some of the other benchmarks. `REWEAVE` shows a slowdown of more than 5 times, `HASHCODE/WEKA` more than 16 times, `NULLTRACK` — 100 times. `SAFEENUM` and `HASHCODE/APROVE` are completely infeasible for our generated trace monitors. Even though we let each of them run for several hours, they did not reach anywhere near the end of the respective computations. The huge slowdown was especially visible with `SAFEENUM`, applied to a `JHotDraw` animation: the amount of animation steps per second dropped very rapidly until it was taking more than 10 seconds per frame, and that time was still increasing.

These numbers alone should be sufficient to motivate the need for further optimisations.

## 4 MEMORY LEAK ELIMINATION

In [1], Allan *et al.* briefly sketch an analysis and code generation strategy to avoid introducing memory leaks into the system. They give one example to show its effectiveness (a version of `SAFEENUM`), but there is no proper experimental validation. Crucially, there is a subtle but important flaw in their proposal, which we shall correct below by introducing the novel notion of *persistent weak references*.

The overall aim is to enable garbage collection of partial matches that are guaranteed not to reach a final state in the automaton. Roughly, that can be achieved when ‘completing the match’ would require an extra event on an object that is already garbage-collected itself; but because that object has expired, the extra event cannot occur. We first describe the analysis required to detect that situation.

Crucial to our strategy is the concept of a *weak reference*. Normally, a reference to a runtime object prevents that object from being reclaimed by the garbage collector. Weak references are present in most garbage-collected languages; they allow the programmer to refer to an object without preventing its destruc-

tion. In Java, this takes the form of the special class `java.lang.ref.WeakRef`, whose constructor takes the referent object as an argument. `WeakRef` provides a method called `get()` — calling it will return the referent, if it still exists, and `null` otherwise.

**Categorising references** The analysis works on the finite state automaton generated from the regular expression of the tracematch. For each non-initial non-final state in that automaton, the free variables of the tracematch are divided into three categories:

**collectableWeakRefs** Variables that are bound on *every* path from the current state to a final state.

**weakRefs** Variables that are not used in the tracematch body and are not in the above set.

**strongRefs** Variables that are not in the above two sets.

As an example, consider the `SAFEENUM` tracematch presented above. It only has two non-initial non-final states (*cf.* Figure 1). From the first of these, we need to take both an `UPDATE` and a `NEXT` transition to reach the final state; `UPDATE` binds the `v` variable and `NEXT` binds `e`, so both of these are **collectableWeakRefs**. On the second state, we only need to see a `NEXT` to get to the final state, so the only **collectableWeakRefs** variable is `e`. Since `v` isn't used in the tracematch body, it is a **weakRef**.

**Exploiting the categorisation** For variables in the first category, it is sufficient to keep weak references (*i.e.* references that do not prevent garbage collection) to the bound values, since we are guaranteed to bind them again before reaching a final state, and could keep a strong reference then (if necessary). Moreover, if one of these weak references expires, then we can discard the entire partial match, since it cannot possibly reach a final state — any path to a final state would have to bind the expired runtime value, which is impossible. This observation might well improve the memory behaviour of trace monitors, since it could reduce the number of live partial matches.

Allan *et al.* claim that for **weakRef** variables, we also only need to keep weak references. The reason is that even if the runtime object expires, it would not actually be used, and so keeping a strong reference would unnecessarily prevent its garbage collection. A reference to it is only kept for matching purposes. Note, however, that discarding partial matches when such a variable expires is not justified, since by definition we can reach a final state without necessarily binding it again. It turns out that this is an oversimplified view that can lead to not all matches being successfully completed; we will explain this in a moment.

Finally, variables that are not necessarily re-bound on every path to a final state and are used in the tracematch body must be kept alive; hence we need to keep strong references — such variables form the **strongRef** category.

Of course, there are certain tracematches which inherently *do* introduce space leaks, and this categorisation of free variables allows the compiler to issue a warning to that effect: if there exists a non-initial non-final state for which **collectableWeakRefs** is empty, then partial matches in that state could conceivably accumulate to an unbounded number without ever being discarded, and a warning should be emitted. Such warnings are very helpful in practice, because it is easy to forget about performance when writing declarative monitor specifications.

**Persistent weak references** To see why the original proposed treatment of **weakRefs** is not sound, consider the following simple example: Suppose we have a tracematch with two symbols, `A` and `B`, and that `A` binds a tracematch variable `x`. The pattern is `AB`, and the tracematch body doesn't use `x`, so that it is a **weakRef**. Imagine at some point during program execution, we have the following constraint on `x`:

$$(x = v_1) \vee (x = v_2)$$

and that then both `v1` and `v2` expire. Weak references to expired objects return `null`, and so now the constraint becomes

$$(x = ?) \vee (x = ?)$$

that is, we cannot tell the two disjuncts apart any more. Thus, when we see another  $B$  event, the tracematch body would be run once instead of twice.

This small example shows that we need to treat weak references that do not invalidate their entire partial match specially: We need to be able to tell them apart even after they expire. We propose the concept of *persistent weak references*, as explained below, to address this issue.

The defining characteristic of a persistent weak reference should be that after its referent expires, calling `get()` returns not `null`, but some object that uniquely identifies the original referent. Moreover, all persistent weak references to the same object should return the same value after it has been garbage-collected.

It is not immediately obvious how one could achieve such behaviour, but our work on *indexing* (cf. Section 5) suggests an approach that works: Make use of *collectable key identity maps*. We proceed by defining `PersistentWeakRef`, a subclass of the standard `WeakRef` class which has no publicly visible constructors. Rather, it provides a static public method `getRefFor(Object o)` that can be used to create new references. This method maintains a static identity map  $m$  from runtime objects to associated instances of `PersistentWeakRef`; when called, it first checks if  $m$  already contains its parameter, and if so simply returns the associated value. Otherwise, it constructs a new `PersistentWeakRef`, records the correspondence in  $m$  and returns it. Effectively this ensures that only one persistent weak reference object is ever constructed for each runtime value. The map  $m$  has special handling for its keys: They are stored as weak references (so as not to prevent their garbage collection), and moreover when they expire, the associated key/value pairs are discarded. In this way, the memory used by the `PersistentWeakRef` class is proportional to the number of *live* referents, that is, it doesn't introduce any memory leaks itself.

Finally, we need to define the behaviour of the `get()` method on our new class. This proceeds as follows: First of all, dispatch to the superclass. If the result is non-`null`, the object is still alive and we can simply return it. If the result is `null`, we can return `this` — that is, the `PersistentWeakRef` instance. This satisfies our two requirements above, as we can still tell apart weak references to expired objects, and references to the same runtime object will return the same value when it expires.

**Measurements** Let us now examine the effects of this optimisation on our set of benchmarks. Table II gives the relevant figures. We see that performance has improved significantly in some cases: REWEAVE (8) runs 3 times faster than before, HASHCODE/WEKA (4) also 3 times faster, the overhead of NULLTRACK (2) is almost halved.

ID	ASPECTJ	LEAKELIM	NOOPT
1	5.3s	133.3	>90M
2	0.47s	25.6s	47.27s
3	485.8s	>90M	>90M
4	2.9s	15s	47.7
5	7.0s	39.7s	41.3s
6	3.4s	3.7s	3.7s
7	18s	20.9s	21.9s
8	5.4s	9.2s	27.9s

Table II: Run times in seconds, including numbers after leak elimination

Particularly pleasing is the fact that SAFEENUM (1) is now feasible, and it is worth examining the situation there a bit more closely. As stated earlier, the base program is JHotDraw, a Java figure editor. For the benchmark, its animation routines were modified by removing all delays, so that the figure elements are moved around the screen as fast as possible. Animation is implemented by enumerating the figure elements at each step, and moving each of them slightly; in total, 100000 animation steps are performed.

There is a crucial difference between this trace monitor and the one in our LUINMETH benchmark, which proved to have very low overheads: SAFEENUM checks that the `next()` method is never called on an enumeration after the underlying collection has been updated; in particular, there is no single event after

which we can be sure that a specific partial match will never lead to a successful match, so the number of potential matches grows unboundedly with base program execution. In LUI<sub>IN</sub>METH, we could discard partial matches when returning from the associated method.

Thus, SAFEENUM is plagued by terrible memory performance, since it has to keep a steadily growing number of objects in memory; also, all of these must be updated after every relevant event, and this explains the huge slowdown we described in the previous section.

Consider now the effects of the space leak elimination on this. The enumeration object will be classified as one of the **collectableWeakRefs**; thus, it can still be garbage-collected, and when it is, all associated partial matches are discarded. Now, each enumeration expires when the associated animation step is completed — the following animation step creates a new one. So whenever we complete an animation step, we drop all associated matching state, which leads to the benchmark’s becoming feasible. Moreover, rather than having unboundedly increasing memory behaviour, it exhibits practically constant memory usage very slightly above that of the uninstrumented program, as described by Allan *et al.*

So, it seems that space leak elimination is indispensable for many applications. We have observed significant speed-ups after enabling the optimisation, particularly for “open-ended” trace monitors (*i.e.* those for which it is not possible to rule out a match’s completion before program termination). Most notably, many *liveness* or *safety* properties, which are popular in the runtime verification community and assert that some good condition always holds or some bad sequence of events never occurs, are of this type. Without observing object garbage collection and invalidating partial matches based on that, such trace monitors would have to keep track of ever-growing sets of potential matches.

In our experience, the analysis of Allan *et al.* succeeds in eliminating space leaks, and correctly emits warnings when there could be a leak. Moreover, the early cleanup of invalidated matching state leads to unexpected performance gains in many situations.

However, tracematch performance is still rather worse than that of equivalent AspectJ instrumentation. For SAFEENUM, the AspectJ version implements the usual technique for implementing safe iterators by putting logical time stamps on collections and iterators - it seems unlikely that this idea can be automatically synthesised from the specification.

The HASHCODE/APROVE (3) benchmark remains infeasible. Close examination of its behaviour reveals that space leaks have been eliminated, but all time is spent iterating over a large set of partial matches. AProVE makes heavy use of hash sets, and each object stored in a hash set is potentially the source of a match completion. We will address this in the next section.

## 5 INDEXING

Recall that the basic implementation of trace monitors (as explained in Section 2) consists of a finite state machine where the states have been labelled with constraints; and that constraints are boolean combinations of variable bindings to objects.

The performance numbers shown so far were taken with a simple implementation that represents such constraints naively as sets of disjuncts. As we will see, large numbers of stored disjuncts will likely be irrelevant to any single update, but with a simple set representation every single one must be iterated over for each update. The overheads this causes make trace-monitoring infeasible for large classes of programs and monitors. This section details a data-structure for partitioning disjuncts, and algorithms for updating such structures, that avoid processing irrelevant disjuncts. The methods shown preserve matching behaviour and extend the techniques from Section 4 so that no new memory leaks are introduced.

To see what it means for a disjunct to be *irrelevant* we must summarise parts of the previous work on the semantics of tracematch matching [1]. A tracematch symbol is modelled as a function from events to constraints.

$$\text{symbol} = \text{event} \rightarrow \text{constraint}$$

For example, if a symbol  $a$  does not match the event  $e$  then  $a(e) = \text{false}$ , but if  $a$  does match the event  $e$  just in the case that the tracematch variable  $x$  is bound to the object  $o$ , then  $a(e) = (x = o)$ .



Figure 2: An annotated automaton for safe-enumeration.

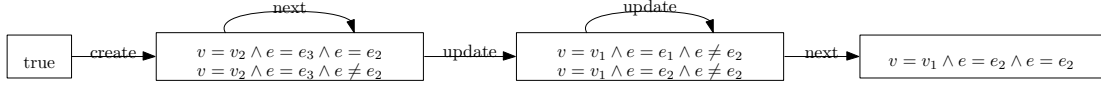


Figure 3: The calculations performed when using simple sets to store disjuncts.

The set of all symbols declared by a tracematch is written  $A$ . We write  $j \xrightarrow{a} i$  to mean there is a transition in the tracematch automaton from state  $j$  to state  $i$  that is labelled with the symbol  $a$ . For each state  $i$ ,  $label_i$  denotes the constraint labelling it. When an event  $e$  occurs, the new label at each state  $i$ , written  $label'_i$ , is

$$label'_i \stackrel{\text{def}}{=} \left( \bigvee_{j \xrightarrow{a} i} (label_j \wedge a(e)) \right) \vee \left( label_i \wedge \bigwedge_{a \in A} \neg a(e) \right) \quad (1)$$

The first line of this equation says that if there is a partial match in state  $j$ , and the variable bindings for that partial match are compatible with  $a(e)$ , then that partial match can transition to state  $i$ . These are called positive updates. The second line states that some transition *must* be taken for each partial match, unless no symbol can be matched to  $e$  that would result in compatible bindings. These are called negative updates.

To illustrate, consider the safe-enumeration monitor from Section 2, together with the variable bindings shown in Figure 2. Suppose that a NEXT event occurred with the variable binding ( $e = e_2$ ). The calculations that should be performed to obtain the new constraints, in accordance with Equation 1, are shown in Figure 3. Indeed, when using a simple-set implementation, these calculations must be performed. However over half of them are redundant: note that three out of the five new disjuncts, when simplified, are either *false* or unchanged — that is, a disjunct that is identical to one previously labelling the same state. These are the *irrelevant* disjuncts. The simplified version of the constraints is shown in Figure 4.

In general, suppose that the constraint labelling some state  $j$  has a disjunct  $d$ , which contains the equality

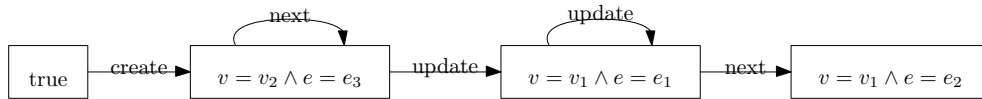


Figure 4: The new constraints after a NEXT event has occurred.

$(x = o_1)$ . If an event  $e$  occurs, and there is an  $a$  transition from  $j$  to  $i$ , we can see from the positive updates in Equation 1 that  $d \wedge a(e)$  will be calculated as part of the new constraint labelling state  $i$ . Suppose, however, that the constraint generated by matching the event to the symbol  $a$  contains  $(x = o_2)$  where  $o_1 \neq o_2$ . It is guaranteed that  $d \wedge a(e) \equiv \text{false}$ , because it contains two contradictory constraints on  $x$ . The disjunct  $d$  is therefore irrelevant to  $a$  at the event  $e$ .

A similar situation is found when calculating negative updates. Suppose that state  $i$  is labelled with a disjunct  $d = (x = o_1) \wedge d'$ , and the same event  $e$  occurs such that  $a(e) = (x = o_2) \wedge c$  (for some predicate  $c$ ). Equation 1 shows that computing the negative updates for  $i$  will involve calculating  $d \wedge \neg a(e)$ :

$$\begin{aligned}
d \wedge \neg a(e) &\equiv d \wedge \neg((x = o_2) \wedge c) \\
&\equiv d \wedge ((x \neq o_2) \vee \neg c) \\
&\equiv (d \wedge (x \neq o_2)) \vee (d \wedge \neg c) \\
&\equiv ((x = o_1) \wedge d' \wedge (x \neq o_2)) \vee (d \wedge \neg c) \\
&\equiv ((x = o_1) \wedge d') \vee (d \wedge \neg c) \\
&\equiv d \vee (d \wedge \neg c) \\
&\equiv d
\end{aligned}$$

In this case,  $d$  is also irrelevant for negative updates — not because it is falsified, but because  $d$  is *unchanged* after the update and continues to label state  $i$ .

The goal of indexing is to partition the disjuncts stored at each state so that as many irrelevant disjuncts are ignored as possible for each update.

## 5.1 Choosing a Partition

The tracematch implementation automatically chooses, for each state, a set of variables with which to partition the disjuncts stored at that state. For illustration, consider a state  $i$ , where only the variables  $x$  and  $y$  are guaranteed to be bound. This state has three outgoing transitions labelled  $a$ ,  $b$ , and  $c$ . The symbol  $a$  binds  $x$  and  $y$ ,  $b$  binds  $x$  and  $z$ , and  $c$  binds just  $z$ . What variables should be used to partition disjuncts labelling  $i$ ?

Firstly, a variable can only be used to partition disjuncts at a state if it is guaranteed to be bound at that state and is also bound by an outgoing transition. If this is not the case, then the definition of irrelevance shown above does not apply. There are therefore some transitions which cannot benefit from disjunct partitioning; the  $c$ -transition on state  $i$  is such a transition because it only binds  $z$ , and  $z$  is not guaranteed to be bound at state  $i$ .

Transitions that cannot benefit from indexing are ignored when choosing a partition for a state. Whether or not partitioning is used at state  $i$ , updates for  $c$  would proceed the same way: should an event corresponding to  $c$  occur, then all the disjuncts labelling  $i$  would be iterated over. The cost of that iteration is the same (modulo constant factors) whatever partition is used.

For each state, the set of variables to partition on is found by considering all transitions that can benefit from partitioning, and intersecting the sets of potential partition-variables from each of them. In the case of state  $i$ , the set of potential partition-variables for  $a$  is  $\{x, y\}$ , and for  $b$  it is  $\{x\}$  ( $c$  is not considered because it cannot benefit from partitioning). Therefore, taking the intersection of these sets, the disjuncts at this state would be partitioned by their binding for  $x$ .

It is possible that this method results in no partition at all because there are two or more mutually exclusive sets that could be partitioned on. The safe-enumeration monitor we have been considering is an example of this: the UPDATE event only binds the vector  $v$ , whilst the NEXT event only binds the enumeration  $e$ . Indeed, in general, there may be some examples where it is most performant to not partition the disjuncts at all. However, it is likely that the programmer will be able to judge which symbols are likely to match the most often. For this reason, symbols may be marked as ‘frequent’ in a tracematch. If no partition can be chosen by the method described above then the process is repeated for just the ‘frequent’ symbols.

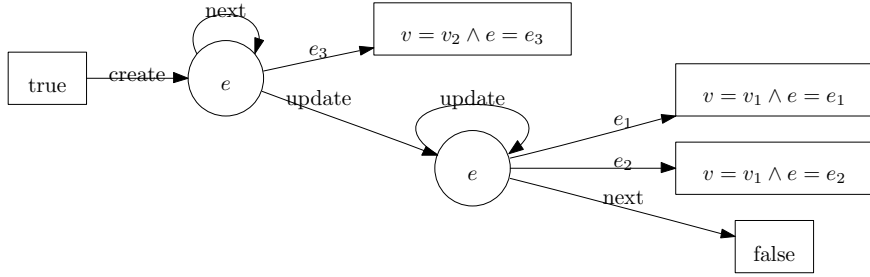


Figure 5: The safe-enumeration constraints, as stored using indexing.

In the case of safe-enumeration, we marked the NEXT symbol as frequent, which meant that the variable  $e$  was chosen to index on.

## 5.2 A Data-Structure for Partitioning

Partitions are represented by the tracematch implementation as trees with the following properties:

- each non-leaf node is labelled with a tracematch variable;
- every node at the same level is labelled with the same variable;
- each edge is labelled with an object;
- the leaves of the tree are sets of disjuncts; and
- if a node is labelled by a variable  $v$ , and an edge from that node is labelled by an object  $o$ , then every disjunct in the branch underneath that edge contains  $(v = o)$ .

For safe-enumeration, the same constraints that appeared in Figure 2 are stored using indexing as shown in Figure 5.

At runtime, the variable labelling each non-leaf node is implicit, and the nodes are represented as maps from objects to child nodes. Writing the implementation of these map-objects requires careful effort in order not to break the optimisations of Section 4, because each map’s keys are objects in the monitored program and the map must keep references to them.

The maps are specialised hash-tables that may keep weak references to the keys and can have extra code that is triggered when a weak reference expires. The behaviour upon a reference expiring differs, depending on the classification of the variable being indexed, as described above. If it is a **collectableWeakRef**, then every disjunct on the sub-tree indexed by that expired weak-reference also must have a collectable weak-reference to that garbage-collected object — it is therefore safe to drop the entire branch when the binding expires. Note that this is a particularly fast way of discarding invalidated constraints: rather than having to iterate over and check each one in turn, all constraints on this state with the same expired binding are dropped at once.

There is a potential problem, however, in that this could introduce a race condition. If an event occurs that does not benefit from the current index variable (as described above), we have to iterate over all key/value pairs at the current level in the indexing map. If the garbage collector invalidates one of the keys during this iteration, it will be removed from the map, and we have violated the Java API requirement of fail-fast iteration. The solution is a custom map implementation that knowingly deviates from the usual iterator contract. We allow *safe* modifications to the map during an iteration, and take care to ensure that dropping a key/value pair due to key expiry is safe in this sense.



In fact, there is one more pitfall of allowing the map implementation to discard branches at undefined times: One cannot rely on an iterator’s `hasNext()` method to give the right result, since all remaining key/value pairs could conceivably be dropped before the call to `next()` even if `hasNext()` returned `true`. Thus, another contract modification is necessary: we allow `next()` to return `null` if there is no further element to iterate.

If the index variable is a **strongRef**, then no further care needs to be taken; we can use a simple identity hash map. Note that this is unlike the standard `HashMap` implementation, which considers keys to be equal subject to their `equals()` method — we really need object identity, due to the `tracematch` semantics.

Finally, let us consider the situation when the indexing variable is a **weakRef**, *i.e.* when we cannot invalidate constraints due to the variable being garbage-collected, but still need to keep a weak reference. Recall that this case proved especially tricky in Section 4. It may seem that indexing does not make sense for such variables. It is, however, still the case that constraints may benefit from indexing, at least while an object is still alive and there are events that bind it. Once it expires, we will never have to explicitly look it up in the map, but we need to keep the associated constraints accessible to iteration.

One approach might be to group together all key/value pairs with an expired key; as stated above, as long as we can iterate over them we can perform updates correctly. However that gives rise to rather unpleasant race conditions. When do we perform this grouping operation? Since a garbage collection can occur at any time, suppose one happened during an iteration of the key/value pairs. By merging the invalidated set into another, we could end up either not iterating it or iterating it twice.

The approach we propose, therefore, is to reuse our work from earlier and use a specialised indexing map that stores its keys in a `PersistentWeakRef` (*cf.* Section 4). Recall that only one such weak reference is constructed for each runtime value, and that once that value expires, calling `get()` returns the weak reference itself. The result is that the indexing map is still fully iterable after the key expires, and key lookups are possible while the key is alive, which is what we aimed to achieve. It is easy to see that this does not result in additional space leaks, since after a key expires the memory overhead for having seen a runtime value is constant and very small, and will be fully eliminated once associated partial matches complete or fail.

### 5.3 Updating Indexed Disjuncts

We follow the code-generation policy previously described for `tracematches` [1] — a method is generated for each `tracematch` symbol. This method is triggered when an event occurs which matches a `tracematch` symbol, given some variable bindings. The method takes the variable bindings and updates the constraints labelling the automaton. Since more than one symbol can match the same event, and updates for every symbol need to know what the constraints were before the event occurred, temporary constraints are used to store the intermediate results.

For example, the pseudo-code for the method which updates the constraints for a symbol *a* is:

```
def update_a(event_bindings):
  for (i,j) in a-transitions:
    label(i).pos = label(i).pos or
                  (label(j).original and event_bindings)

  for i in states:
    label(i).neg = label(i).neg and not(event_bindings)
```

Such a method is run for each symbol that corresponds to an event. Note how the two halves of this method correspond to the two lines of Equation 1, respectively. The only difference is that here the results are imperatively built symbol-by-symbol using temporary variables. After each applicable update method is run, the results are combined:

```
def combine():
  for i in states:
    label(i).original = label(i).pos or label(i).neg
    label(i).pos = false
```

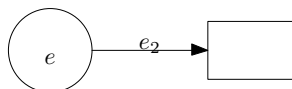


Figure 6: Negative update tree for the safe-enumeration example.

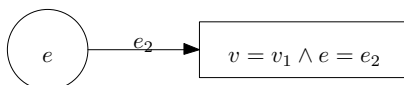


Figure 7: Positive update tree for the safe-enumeration example.

label(i).neg = label(i).original

Note that what is shown here is just pseudo code for clarity; in the actual implementation, specialised code is generated for each operation shown here and the loops are statically unrolled.

When translating this approach to use indexed constraints, one problematic area is the final copy in the combine method. The current constraint labelling the state is copied to a temporary variable which will store the intermediate results of negative updates. However, if this was implemented as a copy-by-reference then later updates would see a modified tree instead of the original, and if it was a copy-by-value, then the whole tree would have to be copied — destroying the performance gains of using indexing at all!

The solution is to build up negative updates starting with an empty tree, but treating it as a ‘patch’ on top of the original. In other words, whenever disjuncts are looked up in the tree for negative updates, if no leaf-node in the tree is found then the original tree is checked instead.

Consider again the safe-enumeration example. This time we will see the same NEXT update for the binding ( $e = e_2$ ) with indexing used — starting with the indexing structure shown in Figure 5.

Firstly, the negative updates must be computed for the two middle states (the initial state always has the constraint *true*). The value of  $e$  from the event is used to look up relevant disjuncts in the trees for each state. For the leftmost of these two states there are no relevant disjuncts, so the negative patch tree for this state is empty. The other state, however, contains the relevant disjunct ( $v = v_1 \wedge e = e_2$ ). Since we are calculating a negative update, we take the conjunction of this disjunct with  $e \neq e_2$ , which is *false*. In a patch tree, *false* is represented by the empty set of disjuncts, as shown in Figure 6.

Next, the positive update trees are computed. Recall that the purpose of positive updates is to propagate disjuncts from state to state along edges labelled with the symbol that matched the current update. Again, the two middle states are considered because they are the only states with outgoing NEXT transitions. Evidently we will find the same single relevant disjunct as when calculating the negative updates, although in general this will be a subset, since not all indexed states will have outgoing transitions for every symbol. This time, the relevant disjunct ( $v = v_1 \wedge e = e_2$ ) is inserted into the positive patch for the final state, as shown in Figure 7.

Finally, the update trees are combined. Recall that the negative update trees are only ‘patches’ — each one contains branches from the original constraint that were *changed* by the negative updates. These changes are folded into the main constraint tree, overwriting entries indexed by the same objects. It is then safe to cleanup any paths to the empty set. Once this is done, the positive changes are folded into the main tree. In contrast, this does *not* involve overwriting because the intent is to find the logical disjunction of the constraint represented by the main tree and the tree represented by the positive update tree. At this stage the results are the same as Figure 4, but in indexed form.

## 5.4 Evaluation

Let us now examine the effect on the benchmark times, as displayed in Table III.

ID	ASPECTJ	FULLOPT	LEAKELIM	NOOPT
1	5.3s	59.1s	133.3	>90M
2	0.47s	1.57s	25.6s	47.27s
3	485.8s	845.0s	>90M	>90M
4	2.9s	3.8s	15s	47.7
5	7.0s	58.7s	39.7s	41.3s
6	3.4s	3.9s	3.7s	3.7s
7	18s	22.4s	20.9s	21.9s
8	5.4s	8.4s	9.2s	27.9s

Table III: Complete benchmark results, with run times in seconds.

Not surprisingly, in some cases indexing *deteriorates* performance, in particular for OBSERVER (5). As we mentioned earlier, in this application each subject has exactly one observer. It follows that the indexing structure only adds to the overhead of accessing that one element. This type of slow-down could be eliminated by introducing indexing in a dynamic fashion, only building the index when the number of disjuncts in a set exceeds a given threshold.

Overall, however, the effect of indexing is hugely beneficial. In the case of NULLTRACK (2), it reduces the execution time from 25.6s to 1.57s. Furthermore, APROVE (3) now becomes feasible to execute, and SAFEENUM becomes just over two times faster.

We conclude that the combination of leak elimination and indexing is a *conditio sine qua non* for the generation of efficient trace monitors.

It is natural to ask whether further improvements are possible. We are currently investigating several possibilities:

**Reducing redundancy** One obvious way to carry our indexing scheme further would be to note that it is redundant to store binding information both labelling the edges of an indexing tree, and on the partial matches at the leaves. We could then further specialise the partial match representation by discarding the redundant information. In the extreme case where every variable appears as an index, we would only have to store a simple counter recording the current automaton state.

Some care has to be taken, however — this is only well-defined if any symbol that can occur at the start of a matched trace binds all tracematch variables. In fact, this occurs frequently, and so such an optimisation seems promising. Indeed, JavaMOP has implemented a variant, with some success.

**Bound variable correlation** In many examples, there exists a many-to-one relationship between the objects bound in a trace monitor. For example, every enumeration corresponds to one collection, every observer has one subject, and so on. We can thus improve the implementation of indexing by storing (say) all observers as a field on the subject - that has the added benefit that when the subject is garbage collected, so are its observers.

Implementing this automatically requires some annotations on the specification, however, and a fairly complex analysis of the base program, going well beyond the cheap techniques we have introduced here. In the presence of such an analysis, we could generate code that is much closer to our hand-coded AspectJ gold standard for benchmarks like SAFEENUM, as described above.

**Static analyses** In the introduction we already alluded to techniques that have been devised for static type-state verification *e.g.* [15]. These analyses can be employed to show that a trace can never match at

certain program points, thus avoiding the need for instrumentation — indeed, this is the approach taken by [8].

Surprisingly, their findings showed that such an expensive analysis is rarely justified, and makes only a very small runtime difference in the vast majority of cases. This serves to highlight the importance of cheap, intraprocedural analyses in achieving acceptable performance for a trace monitor.

SYSTEM	PURPOSE		integration	PATTERNS			IMPLEMENTATION					availability	
	fault finding	functionality		variables	exact-match	context-free	semantics	leak busting	indexing	specialisation	static match		
tracematches [1]	±	+	+	+	+	-	+	+	+	+	+	[8]	+
PQL [22]	+	-	-	+	-	+	-	-	-	-	-	+	+
J-LO [23]	+	-	+	+	-	-	+	-	-	-	-	-	+
MOP [9]	+	+	-	±	-	-	-	∓	∓	-	-	-	+
AspectJ [3]	-	+	+	-	-	-	-	-	-	-	-	-	+
tracecuts [27]	±	+	+	-	+	+	-	-	-	-	-	-	-
PTQL [18]	+	-	-	+	-	+	-	-	-	-	-	+	-
HAWK [11]	+	-	-	+	-	+	-	-	-	-	-	-	-
Alpha [7]	±	+	-	+	+	+	-	-	-	-	-	-	+
Arachne [16]	±	+	+	-	+	-	-	-	-	-	-	-	+

Table IV: Systems for trace monitoring.

Certainly, a lot remains to be done to optimise trace monitor performance. It is clear, however, that we have identified the two essential techniques that make trace monitoring feasible at all.

## 6 RELATED WORK

In the introduction, we already indicated that while there is a substantial body of work on trace monitoring, there are not a lot of systems available. As the focus of this paper is efficient implementation, we only review such systems here. Our original intention was to provide a detailed comparative study of the most mature trace monitoring systems; it turned out, however, that many of the systems were not available to the general public, and even with those that were, we frequently ran into basic problems that prevented our experimental evaluation.

Table IV gives an overview, comparing the salient features. The first five systems in the table are all publicly available and allow trace monitors to be applied to Java programs. They are, therefore, broadly comparable — even though in AspectJ event sequences must be matched by hand-coding the monitor. The five systems on the bottom of the table are either not publicly available, or (in the case of Arachne) apply to another programming language, namely C.

The table attempts a comparison with respect to a number of criteria. Firstly, the purpose of the system: Many are geared solely towards runtime verification, whereas others (mostly with a background in aspect-oriented programming) are actually intended to augment the monitored program by running extra code when a matching trace is found, or maybe by replacing an event with new code.

Next, we examine the issue of integration with a programming language. Several of the systems are deeply integrated with AspectJ, but some others (for instance PQL) are stand-alone tools. The advantage of programming language integration is enhanced checking of the specifications at compile-time.

There is considerable variety in the way patterns are specified. Not all systems allow variables to be bound by the matching process: without such binding, it is difficult to write patterns that monitor the

behaviour of a specific set of objects. The ‘exact-match’ column in Table IV refers to the matching process. There are two different styles of semantics: One can either demand that every single event be accounted for by the pattern, or one can allow arbitrary events to occur in between matched statements, as does, for example, PQL. We refer to the former as an ‘exact-match semantics’, and to the latter as a ‘skipping semantics’. The precise implications of this design choice are very interesting, but beyond the scope of this document; we refer the interested reader to [5] for an in-depth discussion.

Finally, a number of systems allow context-free patterns as opposed to merely finite state machines. While we have not considered the implementation of such rich patterns in this paper, it is clear that the same techniques apply there to avoid space leaks, and to index partial solution sets.

The next section of Table IV examines the characteristics of the implementation. Only very few systems have based their implementation on a semantics. For tracematches, a proof of the correspondence between its declarative and operational semantics is presented in [1]. Tracematches pioneered the use of leak prevention and indexing as described in this paper, though these techniques have been picked up by JavaMOP to some extent — see the relevant discussion below. Other systems are not concerned with space leaks, and pay the associated performance penalty. The authors of HAWK kindly agreed to run our SAFEENUM benchmark for us (HAWK is not available for download), but memory leaks proved prohibitive. Our experience with PQL is described below.

Tracematches are also the only system that automatically specialises the generated code to the pattern — again, of course, without using interprocedural analysis. Further drastic improvements in efficiency are possible in some applications when interprocedural analysis *is* employed. The most sophisticated system of this kind is PQL, employing a BDD-based static analysis to rule out instrumentation points at compile-time. Unfortunately, we were not able to get the static analysis to work. A similar optimisation has been tried in the context of tracematches [8]: The findings were not very encouraging, showing that often the static analysis made only a very small difference in runtime.

The final column indicates whether a system can be freely downloaded. Where this was the case and the system could process Java, we tried to express our benchmarks. The performance of J-LO on SAFEENUM was such that we gave up on attempting further experiments (with its author’s blessing). Our experience with AspectJ has been presented in this paper as the hand-optimised gold standard against which other systems should be measured. The time spent coming up with implementations in each case was substantial. The findings with the remaining two systems were more interesting.

**PQL** The Program Query Language (PQL [22]) was proposed as a stand-alone tool to find bugs in Java programs by writing queries over execution traces. A PQL query can be named, can make use of free variables, and picks out events by writing fragments of concrete Java syntax. It employs a *skipping semantics*, that is, it allows any event to occur between matched statements. Due to the fact that the named queries can be (mutually) recursive, PQL can express context-free properties of the trace quite naturally.

PQL does not include any optimisations to avoid space leaks, and indeed when we encoded SAFEENUM, we observed a steep linear growth of memory usage over time. It was impossible to complete the benchmark without providing the JVM with more memory, at the end PQL was using over 500MB of heap space. Still, PQL completed the benchmark in 580 seconds, significantly faster than the naive tracematch — and around 10 times slower than the fully optimised tracematch instrumentation.

Unfortunately, we ran into significant problems when trying our other benchmarks. Several of them (both HASHCODE benchmarks, for example) cannot be expressed due to limitations of the PQL language; the problem with HASHCODE is that PQL cannot bind primitive types like `int` (this fact was confirmed with the authors of PQL). Also, it is impossible to intercept and bind assignments to fields, and so we couldn’t express NULLTRACK or REWEAVE. DBPOOLING and LUINMETH are both expressible, but do not work with the PQL 0.1 or 0.2 implementations for technical reasons.

**JavaMOP** JavaMOP is an implementation of monitor-oriented programming [9]. It provides a framework for so-called logic plugins to generate a trace monitor from their own domain-specific trace pattern language; such a plugin for regular expressions is predefined, making it easy to port tracematch patterns to JavaMOP.

The system generates AspectJ source code, which then needs to be compiled with an AspectJ compiler to produce the instrumented program.

The optimisations and benchmarks proposed in this paper were first disseminated in a technical report, and pleasingly the developers of JavaMOP incorporated some of our work into their system. At the time of writing, two of our benchmark trace monitors had been expressed in their formalism and made available on their web page [14]; also, parts of the optimisations presented in this paper had found their way into the generated code, demonstrating the general applicability of our techniques.

Concretely, JavaMOP uses an indexing scheme very similar to that of tracematches — a multi-level tree of maps that keep weak references to their keys. The leaves of the tree are labelled not by partial matches, but, essentially, by state numbers. This is therefore equivalent to the potential optimisation about reducing redundancy that we described in Section 5. Note, however, that using it as the standard code generation strategy rather than as an optimisation when it is applicable restricts the user to patterns in which the first symbol matched must bind all free variables of the monitor, and indeed there is a corresponding check in JavaMOP. Also, our analyses to select a suitable set of indices for each state are not implemented; instead, an indexing tree is maintained for each possible set of indices (*i.e.* each distinct set of variables bound by a symbol). This is likely to exacerbate the problem we described, where indexing actually hampered performance rather than improving it for certain benchmarks.

Tentative experiments suggest that our optimisations have been successfully implemented in the context of JavaMOP. We refrain from providing numbers, however, since there are currently a number of problems with the code generation strategy. Since, to use our terminology, all variables are considered to be **collectableWeakRefs** for the purposes of leak elimination, incorrect behaviour may result on patterns which have variables in other categories. Also, the current implementation cannot deal with the monitor action referring to free variables that were not bound by the final symbol — it actually produces non-compileable code in that case, and this is a very significant restriction as it rules out a large number of our benchmarks.

We intend to work together with the authors of JavaMOP to rectify these problems and hope to present a comprehensive set of numbers for the final version of this paper.

## 7 CONCLUSIONS

This paper demonstrates, for the first time, how feasible trace monitors can be generated from specifications. It thus complements the substantial body of work that argued the desirability of trace monitors as a programming language feature.

This result was obtained through two techniques: the elimination of space leaks, and a sophisticated data structure for organising sets of partial matches. Neither of these techniques requires interprocedural analysis, and they can thus be employed without excessive compile-time costs. Our techniques approach the speed of hand-coded, hand-optimised monitors to within an order of magnitude.

The leak elimination analysis was suggested in [1], but the strategy proposed there contained a crucial flaw that made it unsound. We have shown how to rectify this by introducing the novel notion of *persistent* weak references. Furthermore, we presented a thorough experimental evaluation of the effectiveness of the new solution. The fact that the original flaw went undetected for so long is cause for some concern. At present there are no formal verification techniques available for data structures that make use of weak references. We are currently investigating the development of such verification techniques, using the data structure presented here as a motivating example.

All results of this paper, ranging from our benchmark suite to the optimisations, are applicable to most other trace monitoring systems, and we have thus opened the way for many comparative experiments in future. These are already starting to happen, as witnessed by the recent adoption of some of the techniques presented here by the JavaMOP system.

## References

- [1] Chris Allan, Pavel Avgustinov, Aske Simon Christensen, Laurie Hendren, Sascha Kuzins, Ondřej Lhoták, Oege de Moor, Damien Sereni, Ganesh Sittampalam, and Julian Tibble. Adding Trace Matching with Free Variables to AspectJ. In *Object-Oriented Programming, Systems, Languages and Applications*, pages 345–364. ACM Press, 2005.
- [2] AProVE. Automated Program Verification Environment. <http://aprove.informatik.rwth-aachen.de/>, 2006.
- [3] AspectJ Eclipse Home. The AspectJ home page. <http://eclipse.org/aspectj/>, 2003.
- [4] Pavel Avgustinov, Aske Simon Christensen, Laurie Hendren, Sascha Kuzins, Jennifer Lhoták, Ondřej Lhoták, Oege de Moor, Damien Sereni, Ganesh Sittampalam, and Julian Tibble. Optimising AspectJ. In *Programming Language Design and Implementation (PLDI)*, pages 117–128. ACM Press, 2005.
- [5] Pavel Avgustinov, Oege de Moor, and Julian Tibble. On the semantics of trace monitoring patterns. In *Runtime Verification*, 2007.
- [6] Howard Barringer, Allen Goldberg, Klaus Havelund, and Koushik Sen. Rule-based runtime verification. In *Fifth International Conference on Verification, Model Checking and Abstract Interpretation (VMCAI 04)*, volume 2937, pages 44–57. Lecture Notes in Computer Science, 2003.
- [7] Christoph Bockisch, Mira Mezini, and Klaus Ostermann. Quantifying over dynamic properties of program execution. In *2nd Dynamic Aspects Workshop (DAW05)*, Technical Report 05.01, pages 71–75. Research Institute for Advanced Computer Science, 2005.
- [8] Eric Bodden, Laurie Hendren, and Ondřej Lhoták. A staged static program analysis to improve the performance of runtime monitoring. In *European Conference on Object-Oriented Programming*, 2006.
- [9] Feng Chen and Grigore Roşu. Java-MOP: A monitoring oriented programming environment for Java. In *11th International Conference on Tools and Algorithms for the construction and analysis of systems (TACAS '05)*, volume 3440 of *Lecture Notes in Computer Science*, pages 546–550. Springer Verlag, 2005.
- [10] María Augustina Cibrán and Bart Verheecke. Dynamic business rules for web service composition. In *2nd Dynamic Aspects Workshop (DAW05)*, pages 13–18, 2005.
- [11] Marcelo d’Amorim and Klaus Havelund. Event-based runtime verification of java programs. In *WODA '05: Proceedings of the third international workshop on Dynamic analysis*, pages 1–7. ACM Press, 2005.
- [12] Rémi Douence, Thomas Fritz, Nicolas Lorient, Jean-Marc Menaud, Marc Séguira, and Mario Südholt. An expressive aspect language for system applications with arachne. In *Aspect-Oriented Software Development*, pages 27–38. ACM Press, 2005.
- [13] Rémi Douence, Olivier Motelet, and Mario Südholt. A formal definition of crosscuts. In Akinori Yonezawa and Satoshi Matsuoka, editors, *Reflection 2001*, volume 2192 of *Lecture Notes in Computer Science*, pages 170–186. Springer, 2001.
- [14] Grigore Rosu *et al.* JavaMOP homepage. <http://fsl.cs.uiuc.edu/index.php/JavaMOP>, 2007.
- [15] Stephen Fink, Eran Yahav, Nurit Dor, G. Ramalingam, and Emmanuel Geay. Effective typestate verification in the presence of aliasing. In *ISSTA '06: Proceedings of the 2006 international symposium on Software testing and analysis*, pages 133–144, New York, NY, USA, 2006. ACM Press.
- [16] Thomas Fritz, Marc Séguira, Mario Südholt, Egon Wuchner, and Jean-Marc Menaud. An application of dynamic AOP to medical image generation. In *2nd Dynamic Aspects Workshop (DAW05)*, Technical Report 05.01, pages 5–12. Research Institute for Advanced Computer Science, 2005.
- [17] Erich Gamma. JHotDraw. Available from <http://sourceforge.net/projects/jhotdraw>, 2004.

- [18] Simon Goldsmith, Robert O’Callahan, and Alex Aiken. Relational queries over program traces. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages and Applications*, pages 385–402. ACM Press, 2005.
- [19] Peter Hui and James Riely. Temporal aspects as security automata. In *Foundations of Aspect-Oriented Languages (FOAL 2006), Workshop at AOSD 2006*, Technical Report #06-01, pages 19–28. Iowa State University, 2006.
- [20] Gregor Kiczales, Erik Hilsdale, Jim Hugunin, Mik Kersten, Jeffrey Palm, and William G. Griswold. An overview of AspectJ. In J. Lindskov Knudsen, editor, *European Conference on Object-oriented Programming*, volume 2072 of *Lecture Notes in Computer Science*, pages 327–353. Springer, 2001.
- [21] Ramnivas Laddad. *AspectJ in Action*. Manning, 2003.
- [22] Michael Martin, Benjamin Livshits, and Monica S. Lam. Finding application errors using PQL: a program query language. In *Proceedings of the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages and Applications*, pages 365–383. ACM Press, 2005.
- [23] Volker Stolz and Eric Bodden. Temporal Assertions using AspectJ. In *Electronic Notes in Theoretical Computer Science*, volume 144, pages 109–124, 2006.
- [24] Arie van Deursen, Leon Moonen, and Marius Marin. AJHotDraw. <http://sourceforge.net/projects/ajhotdraw/>, 2006.
- [25] Wim Vanderperren, Davy Suvé, María Augustina Cibrán, and Bruno De Fraine. Stateful aspects in JAsCo. In *Software Composition: 4th International Workshop*, volume 3628 of *Lecture Notes in Computer Science*. Springer, 2005.
- [26] w3c. Jigsaw. <http://www.w3.org/Jigsaw/>, 2006.
- [27] Robert Walker and Kevin Viggers. Implementing protocols via declarative event patterns. In *ACM Sigsoft International Symposium on Foundations of Software Engineering (FSE-12)*, pages 159–169. ACM Press, 2004.
- [28] Ian H. Witten and Eibe Frank. *Data Mining: Practical Machine Learning Tools and Techniques with Java implementations*. Morgan Kaufmann Publishers, 2000.