

Generic Program Monitoring by Trace Analysis

Erwan Jahier , Mireille Ducassé

N°4323

Novembre 2001

————— THÈME 2 —————



*Rapport
de recherche*



Generic Program Monitoring by Trace Analysis

Erwan Jahier* , Mireille Ducassé†

Thème 2 — Génie logiciel
et calcul symbolique
Projet Lande

Rapport de recherche n° 4323 — Novembre 2001 — 40 pages

Abstract: Program execution monitoring consists of checking whole executions for given properties in order to collect global run-time information. Monitoring is very useful to maintain programs. However, application developers face the following dilemma: either they use existing tools which never exactly fit their needs, or they invest a lot of effort to implement monitoring code. In this report we argue that, when an event-oriented tracer exists, the compiler developers can enable the application developers to easily code their own, relevant, monitors. We propose a high-level operator, called `foldt`, which operates on execution traces. One of the key advantages of our approach is that it allows a clean separation of concerns; the definition of monitors is totally distinct from both the user source code and the language compiler. We give a number of applications of the `foldt` operator to compute monitors for Mercury program executions: execution profiles, graphical abstract views, and two test coverage measurements. Each example is implemented by a few simple lines of Mercury.

Key-words: dynamic program analysis, monitoring, trace analysis, test coverage, logic and functional programming, Mercury programming language

(Résumé : `tsvp`)

* IRISA/IFSIC ; email : Erwan.Jahier@irisa.fr

† IRISA/INSA ; email : Mireille.Ducasse@irisa.fr

Monitoring générique de programmes par analyse de trace

Résumé : Le monitoring d'exécutions de programmes consiste à vérifier des propriétés sur des exécutions complètes afin de collecter des informations globales sur le comportement des programmes. Le monitoring est très utile pour maintenir les programmes. Cependant, les développeurs d'applications font face au dilemme suivant : soit ils utilisent les outils existants qui ne remplissent jamais totalement leurs attentes, soit ils investissent énormément d'efforts pour implémenter des moniteurs adéquats. Dans ce rapport, nous argumentons que, dès lors qu'un traceur orienté événements existe, le développeur du compilateur peut permettre aux développeurs d'applications de coder facilement leurs propres moniteurs. Nous proposons un opérateur de haut niveau, appelé `foldt`, qui opère sur des traces d'exécution. Un des avantages clés de notre approche est qu'elle permet une nette séparation des préoccupations ; la définition des moniteurs est totalement distincte à la fois du code source des applications et du compilateur du langage de programmation. Nous donnons un certain nombre d'exemples de moniteurs utilisant l'opérateur `foldt` pour analyser des exécutions de programmes Mercury : profils d'exécutions, vues abstraites graphiques, ainsi que deux mesures de couverture de test. Chaque exemple est implémenté par quelques lignes simples de Mercury.

Mots-clé : analyse dynamique de programmes, monitoring, analyse de traces, couverture de test, programmation logique et fonctionnelle, langage de programmation Mercury

1 Introduction

Program maintenance and trace analysis. Several studies (e.g., [Hat97]) show that maintenance is the most expensive phase of software development: the initial development represents only 20 % of the cost, whereas error fixing and addition of new features after the first release represent, each, 40 % of the cost. Thus, 80 % of the cost is due to the maintenance phase.

A key issue of maintenance is program understanding. In order to fix logical errors, programmers have to analyze their program symptoms and understand how these symptoms have been produced. In order to fix performance errors, programmers have to understand where the time is spent in the programs [Ger00]. In order to add new functionality, programmers have to understand how the new parts will interact with the existing ones.

Program analysis tools help programmers understand programs. For example, type checkers [Pfe92] help understand data inconsistencies. Slicing tools [GL91, Tip95] help understand dependencies among parts of a program. Tracers give insights into program executions [EB88].

Some program analysis tools automatically analyze program execution traces. They can give very precise insights of program (mis)behavior. We have shown how such trace analyzers can help users debug their programs. In our automated debuggers, a trace query mechanism helps users check properties of parts of traced executions in order to understand misbehavior [Duc99a, Duc99b, JD99a].

In this article, we show that trace analysis can be pushed toward monitoring to further help understand program behavior. Program execution monitoring consists of collecting global run-time information relative to whole program executions. For example, some monitors gather statistics which help detect heavily used parts that need to be optimized. Other monitors build graphs (e.g., control flow graphs, dynamic call graphs, proof trees) that give a global understanding of the execution.

Execution monitoring. Monitors are trace analyzers which differ from debuggers. Monitoring is mostly a “batch” activity whereas debugging is mostly an interactive activity. In monitoring, a *set* of properties is specified beforehand; the whole execution is checked; and the global collected information is displayed. In debugging, the end-user is central to the process; he specifies on the fly the very next property to be checked; each query is induced by the user current understanding of the situation at the very moment it is typed in. Monitoring is therefore less versatile than debugging. The properties specified for monitoring have a much longer life time, they are meant to be used over several executions.

It is, nevertheless, impossible to foresee all the properties that programmers may want to check on executions. One intrinsic reason is that these properties are often driven by the application domain. Therefore monitoring systems must provide some genericity.

Existing approach to implement monitors. Unfortunately, monitors are generally implemented by ad hoc instrumentation. This instrumentation requires a significant pro-

gramming effort. When done at a low level, for example by modifying the runtime system, it requires a deep knowledge of the system to instrument. The language compiler implementors are the most qualified persons to implement low-level instrumentation. However, the monitored information, as already mentioned, is often application-dependent, hence application programmers and end-users know better what has to be monitored, but it is almost impossible for them to instrument compilers at a low level.

An alternative to low-level instrumentation is source-level instrumentation. For example, run-time behavior information can be extracted by source-to-source transformation for ML [TA95, KHC91] or Prolog [DN00]). Such instrumentation, although simpler than low-level compiler instrumentation, can still be too complex for most programmers. Furthermore, for certain new declarative programming languages like Mercury [SHC96], they may even be impossible. Indeed, in Mercury, the declarative semantics is simple and well defined, but the operational semantics is implicit and complex. For example, the compiler reorders the goals according to its needs. Furthermore, input and output can be made only in deterministic predicates. This complicates code instrumentation.

Thus, ad hoc instrumentation is tedious at a low level and it may be impossible at a high level. On the other hand, the difficult task of instrumenting the code to extract run-time information has, in general, already been achieved to provide a debugger. Debuggers, which help users locate faults in programs are based on tracers. These tracers generate execution traces which provide a precise and faithful image of the operational semantics of the traced language. These traces often contain sufficient information to base monitors upon them.

Our Proposal. In this article, we propose a high-level operator built on top of an execution tracer. The proposed monitoring operator, called `foldt`, is a `fold` which operates on the execution trace. One of the key advantages of our approach is that it allows a clean separation of concerns; the definition of the monitors is totally distinct from both the user source code and the language compiler.

We have implemented `foldt` on top of the Mercury trace. We give a number of applications of the `foldt` operator to compute various monitors: execution profiles, graphical abstract views, and test coverage measurements. Each example is implemented by a few lines of Mercury which can be written by any Mercury programmer. These applications show that the Mercury trace, indeed, contains enough information to build a wide variety of interesting monitors. Detailed measurements show that, under some requirements, `foldt` can have acceptable performance for executions of several millions of execution events. Therefore, our operator, lays the foundation for a generic and powerful monitoring environment.

Note that we have implemented the `foldt` operator on top of Mercury mostly for historical reasons. We acknowledge that some of the monitors were particularly easy to write thanks to the power of Mercury libraries, in particular the set library (e.g., Figure 10). Nevertheless, `foldt` could be implemented for any system

- with an event-oriented tracer, and such that
- this tracer can be modified to integrate `foldt`.

An event oriented *trace* is a sequence of events. An *event* is a tuple of event attributes. An *event attribute* is an elementary piece of information that can be extracted from the current state of the program execution. A trace can be seen as a sequence of tuples of a database ordered by time. Many tracers are event-oriented: for example, Prolog tracers based on Byrd box model [Byr80], tracers for C such as Dalek [OCH90] and Coca [Duc99a], the Egadt tracer for Pascal [FAS94], the Esa tracer for Ada [HS96], and the Ebba tracer for distributed systems [Bat95].

Plan. In Section 2, we introduce the `foldt` operator and describe its current implementation on top of the Mercury tracer. In Section 3, we illustrate the genericity of `foldt` with various kinds of monitors. All the examples are presented at a level of detail that does not presuppose any knowledge of Mercury. Section 4 discusses performance issues of `foldt`. Section 5 compares our contribution with related work. A thorough description of the Mercury trace can be found in Appendix A. Appendix B lists a Mercury program solving the n queens problem, which is used at various places in the article as an input for our monitors.

2 A high-level trace processing operator: `foldt`

In this section, we first define the `foldt` operator over a general trace in a language-independent manner. We describe an implementation of this operator for Mercury program executions, and then present its current user interface.

2.1 Language independent `foldt` definition

A trace is a list of events; analyzing a trace therefore requires to process such a list. The standard functional programming operator `fold` encapsulates a simple pattern of recursion for processing lists. It takes as input arguments a function, a list, and an initial value of an accumulator; it outputs the final value of the accumulator; this final value is obtained by successively applying the function to the current value of the accumulator and each element of the list. As demonstrated by Hutton [Hut99], `fold` has a great expressive power for processing lists. Therefore, we propose a `fold`-like operator to process execution traces; we call this operator `foldt`.

Before defining `foldt`, we define the notions of event and trace for sequential executions.

Definition 1 (*Execution event, Event attributes, Execution trace*)

An execution event is an element of the Cartesian product $\mathbb{E} = A_1 \times \dots \times A_n$, where A_i for $i \in \{1, \dots, n\}$ are arbitrary sets called event attributes. An execution trace is a (finite or infinite) sequence of execution events; the set of all execution traces is denoted by \mathbb{T} . We note $|t|$ the size of the finite sequence of events of a finite trace $t \in \mathbb{T}$ and $|t| = \infty$ the size of infinite traces.

The following definition of `foldt` is a predicative definition of a `fold` operating on a finite number of events of a (possibly infinite) trace. The set of predicates over $\tau_1 \times \dots \times \tau_n$ is denoted by $\text{pred}(\tau_1, \dots, \tau_n)$.

Definition 2 (*foldt operator*)

A `foldt` monitor of type $\tau \times \tau'$ is a 3-tuple : $(\text{initialize}, \text{collect}, \text{post_process}) \in \text{pred}(\tau) \times \text{pred}(\mathbb{E}, \tau, \tau) \times \text{pred}(\tau, \tau')$ such that: $\forall t = (e_i)_{i>0} \in \mathbb{T}$, either

$$(1) |t| < \infty \wedge (\exists!(V_0, \dots, V_n) \in \tau^{n+1}. \\ (\text{initialize}(V_0) \wedge \bigwedge_{i=1}^n \text{collect}(e_i, V_{i-1}, V_i) \wedge \text{post_process}(V_n, \text{Res})))$$

$$(2) \exists!n < |t|, \exists!(V_0, \dots, V_n) \in \tau^{n+1}, \forall x \in \tau. \\ (\text{initialize}(V_0) \wedge \bigwedge_{i=1}^n \text{collect}(e_i, V_{i-1}, V_i) \wedge \text{post_process}(V_n, \text{Res}) \\ \wedge \neg \text{collect}(e_{n+1}, V_n, x))$$

Res is called the result of the monitor $(\text{initialize}, \text{collect}, \text{post_process})$ on trace t .

Operationally, an accumulator of type τ is used to gather the collected information. It is first initialized (V_0). The predicate `collect` is then applied to each event of the trace in turn, updating the accumulator along the way (V_i). There are two ways to stop this process: (1) the folding process stops when the end of the execution is reached if the trace is finite ($|t| < \infty$); (2) if `collect` fails before the end of the execution is reached ($\forall x \in \tau. (\neg \text{collect}(e_{n+1}, V_n, x))$). The last value of the accumulator (V_n) is processed by `post_process`, and put in the result (*Res*) of type τ' ($\text{post_process}(V_n, \text{Res})$).

Note that this definition holds for finite and infinite traces (thanks to the second case of Definition 2). This is convenient to analyze programs that run permanently. The ability to end the `foldt` process before the end of the execution is also convenient to analyze executions part by part as explained in Section 2.3.3. A further interesting property, useful to execute several monitors in a single program execution, is the possibility to simultaneously apply several `fold` on the same list using a tuple of `fold` [Bir87]; in other words:

$$\text{foldt}(i_1, c_1, p_1) \times \dots \times \text{foldt}(i_n, c_n, p_n) = \text{foldt}(i_1 \times \dots \times i_n, c_1 \times \dots \times c_n, p_1 \times \dots \times p_n) \\ \text{where:} \\ \forall a_1, \dots, a_n \in \tau_1 \times \dots \times \tau_n, \\ i_1 \times \dots \times i_n(a_1, \dots, a_n) \Leftrightarrow i_1(a_1) \wedge \dots \wedge i_n(a_n), \\ \forall e \in \mathbb{E}, \forall a_1, \dots, a_n \in \tau_1 \times \dots \times \tau_n, \forall a'_1, \dots, a'_n \in \tau'_1 \times \dots \times \tau'_n, \\ c_1 \times \dots \times c_n(e, a_1, \dots, a_n, a'_1, \dots, a'_n) \Leftrightarrow c_1(e, a_1, a'_1) \wedge \dots \wedge c_n(e, a_n, a'_n), \\ \forall a_1, \dots, a_n \in \tau_1 \times \dots \times \tau_n, \\ p_1 \times \dots \times p_n(a_1, \dots, a_n, a'_1, \dots, a'_n) \Leftrightarrow p_1(a_1, a'_1) \wedge \dots \wedge p_n(a_n, a'_n).$$

2.2 An implementation of `foldt` for Mercury

We prototyped an implementation of `foldt` for the Mercury programming language. After a brief presentation of Mercury and its trace system, we describe our `foldt` implementation.

2.2.1 Mercury and its trace

Mercury [SHC96] is a logic and functional programming language. The principal differences with Prolog are as follows. Mercury supports functions and high order terms. Mercury programs are free from side-effects; even input and output are managed in a declarative way. Mercury strong type, mode and determinism system allows a lot of errors to be caught at compile time, and a lot of optimizations to be done.

The trace generated by the Mercury compiler [SH99] is adapted from Byrd box model [Byr80]. Its attributes are the event number, the call number, the execution depth, the event type (or port), the determinism, the procedure (defined by a module name, a name, an arity and a mode number), the live arguments, the live non-argument variables, and the goal path. A detailed description of these attributes together with an example of event is given in appendix A.

2.2.2 The foldt implementation

An obvious and simple way to implement `foldt` would be to store the whole trace into a single list, and then to apply a `fold` to it. This naive implementation is highly inefficient, both in time and in space. It requires to create and process a list of possibly millions of events. Most of the time, creating such a list is simply not feasible because of disk space limitations. With the current Mercury trace system, several millions of events are generated each second, each event requiring several bytes. To implement realistic monitors, runtime information needs to be collected and analyzed simultaneously (on the fly), **without explicitly creating the trace.**

In order to achieve analysis on the fly, we have implemented `foldt` by modifying the Mercury trace system [SH99]. The Mercury trace system works as follows: when a program is compiled with tracing switched on, the generated C code¹ is instrumented with calls to the tracer (via the C function `trace`). Before the first event (resp. after the last one), a call to an initialization C function `trace_init` (resp. to a finalization C function `trace_final`) is inserted.

When the trace system is entered through either one of the functions `trace`, `trace_init`, or `trace_final`, the very first thing it does is to look at an environment variable that tells whether the Mercury program has been invoked from a shell, from the standard Mercury debugger (`mdb`), or from an other debugger (e.g., Morphine [JD99a]). We have added a new possible value for that environment variable which indicates whether the program has been invoked by `foldt`. In that case, the `trace_init` function dynamically links the Mercury program under execution with the object file that contains the object code of `collect`, `initialize`, and `post_process`. Dynamically linking the program to its monitor is very convenient because neither the program nor the monitor need to be recompiled.

Once the monitor object file has been linked with the program, the C function `trace_init` can call the procedure `initialize` to set the value of a global variable `accumulator_variable`

¹Currently, the only Mercury back-end that has a tracer is the one that relies on a C compiler to produce its executable code.

```

1      % 1 - Define the type of the accumulator:
2      :- type accumulator_type == < A Mercury type >.
3
4      % 2 - Initialize the accumulator:
5      initialize(Accumulator) :-
6          < Mercury goals which initialize the accumulator >.
7
8      % 3 - Update the accumulator:
9      collect(Event, AccumulatorIn, AccumulatorOut) :-
10         < Mercury goals which update the accumulator >.
11
12     % 4 - Optionally, post-process the last value of the accumulator:
13     :- type collected_type == < A Mercury type >.
14
15     post_process(Accumulator, FoldtResult) :-
16         < Mercury goals which post-process the accumulator >.

```

Figure 1: What the user needs to define to use `foldt`

(of type τ). At each event, the C function `trace` calls the procedure `collect` which updates `accumulator_variable`. If `collect` fails or if the last event is reached, the C function `trace_final` calls the procedure `post_process` with `accumulator_variable` and returns the new value of this accumulator (now of type τ').

2.3 The current user interface of `foldt` for Mercury

In this Section, we first describe what the user needs to do in order to define a monitor with `foldt`. Then, we show how this monitor can be invoked.

2.3.1 Defining monitors

We chose Mercury to be the language in which users define the `foldt` monitors to monitor Mercury programs. As a matter of fact, it could have been any other language that has an interface with C, since the trace system of Mercury is written in C. The choice of Mercury, however, is quite natural; people who want to monitor Mercury programs are likely to be Mercury programmers.

The items users need to implement in order to define a `foldt` monitor are given in Figure 1. Lines preceded by ‘%’ are comments. First of all, since Mercury is a typed language, one first needs to define the type of the accumulator variable `accumulator_type` (line 2). Then, one needs to define `initialize` which gives the initial value of the accumulator, and `collect` which updates the accumulator at each event (line 9). Optionally,

```

1  :- import_module int.
2  :- type accumulator_type == int.
3  initialize(0).
4  collect(Event, C0, C) :-
5    if port(Event) = call then C = C0+1 else C = C0.

```

Figure 2: `count_call`, a monitor that counts the number of calls using `foldt`

one can also define the `post_process` predicate which processes the last value of the accumulator. `post_process` takes as input a variable of type `accumulator_type` (τ) and outputs a variable of type `collected_type` (τ'). If `collected_type` is not the same type as `accumulator_type`, then one needs to provide its definition too (line 13). Types and modes of predicates `initialize`, `collect` and `post_process` should be consistent with the following Mercury declarations:

```

:- pred initialize(accumulator_type::out) is det.
:- pred collect(event::in,accumulator_type::in,accumulator_type::out)
    is semidet.
:- pred post_process(accumulator_type::in,collected_type::out) is det.

```

These declarations state that `initialize` is a deterministic predicate (`is det`), namely it succeeds exactly once, and it outputs a variable of type `accumulator_type`; `collect` is a semi-deterministic predicate, namely it succeeds at most once, and it takes as input an event and an accumulator. If `collect` fails, the monitoring process stops at the current event. This can be very useful, for example to stop the monitoring process before the end of the execution if the collecting data is too large, or to collect data part by part (e.g., collecting the information by slices of 10000 events). This also allows `foldt` to operate over non-terminating executions.

The type `event` is a structure made of all the event attributes. To access these attributes, one can use specific functions whose declarations are of the form:

“`:- func <attribute_name>(event::in) = <attribute_type>::out.`”, which means that each such function takes an event and outputs the event attribute corresponding to its name. For example, the function call `depth(Event)` returns the depth of `Event`. The complete list of attribute names is given in Appendix A.

Figure 2 shows an example of monitor that counts the number of predicate invocations (calls) that occur during a program execution. We first import library module `int` (line 1) to be able to manipulate integers (line 2). Predicate `initialize` initializes the accumulator to ‘0’ (line 3). Then, for every execution event, `collect` increments the counter if the event port is `call`, and leaves it unchanged otherwise (line 5). Since `collect` can never fail here, the calls to `collect` proceed until the last event of the execution is reached.

Note that those five lines of code constitute *all the necessary lines* for this monitor to be run. For the sake of conciseness, in the following figures containing monitors, we sometimes

```

[morphine]: run_mercury(queens), foldt(count_call, Result).
                                     A 5 queens solution is [1, 3, 5, 2, 4]
      Last event of queens is reached
      Result = 146      More? (;)
[morphine]:

```

Figure 3: Invoking `foldt` monitor of Figure 2 from an interpreter

forget the module importation directives as well as the type of the accumulator when the context makes them clear.

2.3.2 Invoking `foldt`

Currently, `foldt` can be invoked from a Prolog query loop interpreter. We could not use Mercury for that purpose because there is no Mercury interpreter yet.

We have implemented a Prolog predicate named `run_mercury`, which takes a Mercury program as argument, and which forks a process in which this Mercury program runs in corouting with the Prolog process. When the first event of the Mercury program is reached, the hand is given to the Prolog process which waits for a `foldt` query. `foldt` takes as argument the name of the file where the monitor is defined, and it binds its second argument with the result of the monitor.

A possible session for invoking the monitor of Figure 2 is given in Figure 3. At the right-hand side of the `[morphine]:` prompt, there are the characters typed in by a user. The line in italic is output by the Mercury program; all the other lines are output by the Prolog process. We can therefore see that the program `queens` (which solves the 5 queens problem, cf Appendix B) produces 146 procedure calls.

2.3.3 Illustration of the advantage of calling `foldt` from a Prolog query loop

Being able to call `foldt` from a Prolog interpreted loop enables users to write scripts that control several `foldt` invocations. Figures 4 and 5 illustrate this. The monitor of Figure 4 computes the maximal depth for the next 500 events. In the session of Figure 5, a user (via the `[user].` directive) defines the predicate `print_max_depth` that calls the monitor of Figure 4 and prints its result in loop until the end of the execution is reached. This is useful for example for a program that runs out of stack space to check whether this is due to a very deep execution and to know at which events this occurs.

Note that the fact that the monitor is dynamically linked with the monitored program has an interesting side-effect here: one can change the monitor during the `foldt` query resolution (by modifying the file where this monitor is defined). In our example, one could change the interval within which the maximal depth is searched from 500 to 100. The monitor would be (automatically) recompiled, but the `foldt` query would not need to be

```

1  initialize(acc(0, 0)).
2  collect(Event, acc(N0, Depth0), acc(N0+1, max(Depth0, depth(Event)))) :-
3      N0 < 500. % stops after 500 events
4

```

Figure 4: Monitor that computes the maximal execution depth by interval of 500 events

```

[morphine]: [user].
  print_max_depth :-
    foldt(max_depth, acc(_, MaxDepth)),
    print("The maximal depth is "), print(MaxDepth), nl,
    print_max_depth.
^D
[morphine]: run_mercury(qsort), print_max_depth.

  The maximal depth is 54
  The maximal depth is 28
  The maximal depth is 50
                                     [0, 2, 4, 6, 7, 8, ..., 94, 95, 99, 99]
  Last event of qsort is reached
  The maximal depth is 53
[morphine]:

```

Figure 5: A possible session using the monitor of Figure 4

killed and rerun. This can be very helpful to monitor a program that runs permanently; the monitored program is simply suspended while the monitor is recompiled.

As a matter of fact (as the prompt suggests), the Prolog query loop that we use is Morphine [JD99a], an extensible debugger for Mercury “à la Opium” [Duc99b]. The basic idea of Morphine is to build on top of a Prolog query loop a few coroutines primitives connected to the trace system (like `foldt`). Those primitives let one implement all classical debugger commands as efficiently as their hand-written counter-parts; the advantage is, of course, that they let users implement more commands than the usual hard-coded ones, fitting their own needs.

Invoking `foldt` from a debugger has a further advantage; it makes it very easy to call a monitor during a debugging session, and vice versa. Indeed, some monitors are very useful for program runtime behavior understanding, and therefore can be seen as debugging tools.

```

1  :- import_module int, array.
2  :- type accumulator_type == array(int).
3
4  :- mode acc_in  :: array_di.
5  :- mode acc_out :: array_uo.
6
7  initialize(Array) :-
8      init(5, 0, Array).
9
10 collect(Event, Array0, Array) :-
11     Port = port(Event),
12     port_to_int(Port, IntPort),
13     lookup(Array0, IntPort, N),
14     set(Array0, IntPort, N+1, Array).
15
16 :- pred port_to_int(trace_port_type::in, int::out) is det.
17 port_to_int(Port, Number) :-
18     ( if   Port = call then Number = 0
19       else if Port = exit then Number = 1
20       else if Port = redo then Number = 2
21       else if Port = fail then Number = 3
22       else Number = 4 ).

```

Figure 6: A monitor that counts the number of events at each port

3 Applications

In this section, we describe various execution monitors that can be implemented with `foldt`. We first give monitors which compute three different execution profiles: number of events at each port, number of goal invocations at each depth, and sets of solutions. Then, we describe monitors that produce two types of execution graphs: dynamic control flow graph and dynamic call graph. Finally, we introduce two test coverage criteria for logic programs, and we give the monitors that measure them.

3.1 Execution profiles

3.1.1 Counting the number of events at each port

In Figure 2, we have given a monitor that counts the number of goal invocations. Figure 6 shows how to extend this monitor to count the number of events at each port. We need 5 counters that we store in an array. The default mode of the second and third argument of `collect`, respectively equal to `in` and `out`, can be overridden; here, we override them with the values `array_di` and `array_uo` (lines 4 and 5). Modes `array_di` and `array_uo` are special modes that allow arrays to be destructively updated. Predicate `initialize` creates

```

1  % Module importation and accumulator type (array of int) omitted
2
3  initialize(Acc) :-
4      init(32, 0, Acc).
5
6  collect(Event, Acc0, Acc) :-
7      ( if port(Event) = call then
8          Depth = depth(Event),
9          ( if semidet_lookup(Acc0, Depth, N) then
10             set(Acc0, Depth, N+1, Acc)
11         else
12             size(Acc0, Size),
13             resize(Acc0, Size*2, 0, Acc1),
14             set(Acc1, Depth, 1, Acc)
15         )
16     else
17         Acc = Acc0
18     ).

```

Figure 7: A monitor that counts the number of calls at each depth

an array `Array` of size 5 with each element initialized to 0 (line 8). Predicate `collect` extracts the port from the current event (line 11) and converts it into an integer (line 12)². This integer is used as an index to get (`lookup/3`) and set (`set/4`) array values. The goal `lookup(Array0, IntPort, N)` returns in `N` the `IntPortth` element of `Array0`. The goal `set(Array0, IntPort, N+1, Array)` sets the value `N+1` in the `IntPortth` element of `Array0` and returns the resulting array in `Array`.

3.1.2 Counting the number of calls at each depth

Figure 7 implements a monitor that counts the number of calls at each depth. Predicate `initialize` creates an array of size 32 with each element initialized to 0 (line 4). At call events (line 7), predicate `collect` extracts the depth from the current event (line 8) and increments the corresponding counter (lines 10 and 14). Whenever the upper bound of the array is reached, i.e., whenever `semidet_lookup/4` fails (line 9), the size of the array is doubled (lines 13).

3.1.3 Collecting solutions

The monitor of Figure 8 collects the solutions produced during the execution. We define the type `solution` as a pair containing a procedure and a list of arguments (line 1). The

²As a matter of fact, there are more ports than the ones handled by `port_to_int/2` in Figure 6 (cf Appendix A); we ignore them here for the sake of conciseness.

```

1  :- type solution ---> proc_name/arguments.
2  :- type accumulator_type == list(solution).
3
4  initialize([]).
5
6  collect(Event, AccIn, AccOut) :-
7      ( if
8          port(Event) = exit,
9          Solution = proc_name(Event)/arguments(Event),
10         not(member(Solution, AccIn))
11     then
12         AccOut = [Solution | AccIn]
13     else
14         AccOut = AccIn
15     ).

```

Figure 8: A monitor that collects all the solutions

collected variable is a list of `solutions` (line 2), which is initialized to the empty list (line 4). If the current port is `exit` (line 8) and if the current solution has not been already produced (lines 9,10), then the current solution is added to the list of already collected solutions (line 12).

3.2 Graphical abstract views

Other execution abstract views, that are widely used and very useful in terms of program understanding, are views given in terms of graphs. In the following, we show how to implement monitors that generate graphical abstractions of program executions such as control flow graphs and dynamic call graphs. We illustrate the use of these monitors by applying them to the 5 queens program given in Appendix B. This 100 line program generates 698 events for a board of 5×5 . In this article, we use the graph drawing tool `dot` [KN91]. More elaborated visualization tools such as in [SDBP98] would be desirable, especially for large executions. This is, however, beyond the scope of this article.

3.2.1 Dynamic control flow graphs

We define the *dynamic control flow graph* of a logic program execution as the directed graph where nodes are predicates of the program, and arcs indicate that the program control flow went from the origin to the destination node. The dynamic control flow graph of the 5 queens program is given in Figure 9. We can see, for example, that, during the program execution, the control moves from predicate `main/2` to predicate `data/1`, from predicate `data/1` to predicate `data/1` and predicate `queen/2`.

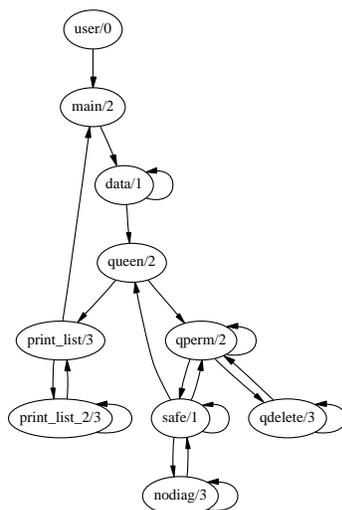


Figure 9: The dynamic control flow graph of 5 queens

```

1  :- type predicate ---> proc_name/arity.
2  :- type arc ---> arc(predicate, predicate).
3  :- type graph == set(arc).
4  :- type accumulator_type ---> collected_type(predicate, graph).
5
6  initialize(collected_type("user"/0, set__init)).
7
8  collect(Event, Acc0, Acc) :-
9      Port = port(Event),
10     ( if (Port = call ; Port = exit ; Port = fail ; Port = redo) then
11         Acc0 = collected_type(PreviousPred, Graph0),
12         CurrentPred = proc_name(Event) / proc_arity(Event),
13         Arc = arc(PreviousPred, CurrentPred),
14         set__insert(Graph0, Arc, Graph),
15         Acc = collected_type(CurrentPred, Graph)
16     else
17         % other events
18         Acc = Acc0
19     ).

```

Figure 10: A monitor that calculates dynamic control flow graphs

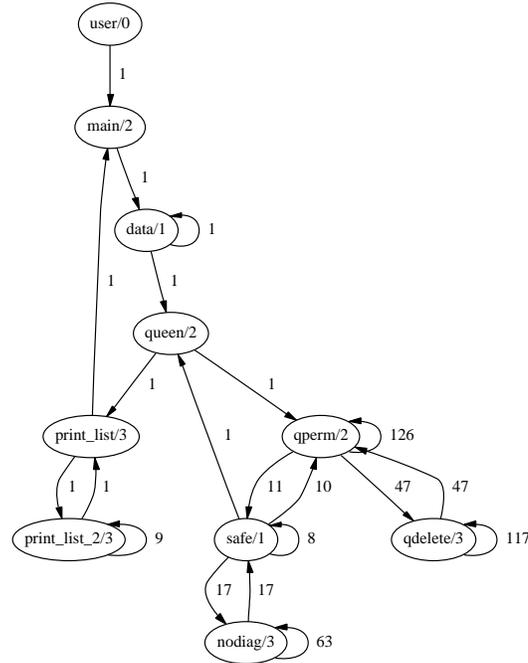


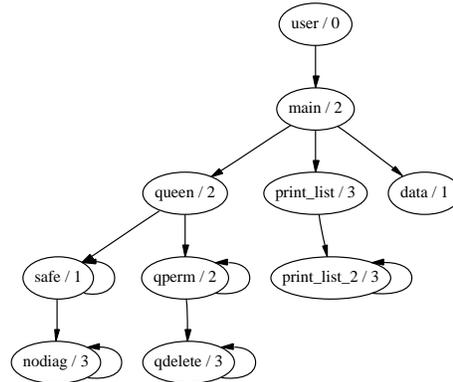
Figure 11: The dynamic **control flow graph** of 5 queens annotated with counters

An implementation of a monitor that produces such a graph is given in Figure 10. Graphs are encoded by a set of arcs, and arcs are terms composed of two predicates (lines 1 to 3). The collecting variable is composed of a predicate and a graph (line 4); the predicate is used to remember the previous node. The collecting variable is initialized with the fake predicate `user/0`, and the empty graph (line 6). At `call`, `exit`, `redo`, and `fail` events (line 10), we insert in the graph an arc from the previous predicate to the current one (lines 11 to 14).

Note that in our definition of dynamic control flow graph, the number of times each arc is traversed is not given. Even if the control goes between two nodes several times, only one arc is represented. One can imagine a variant where, for example, arcs are labeled by a counter; one just needs to use multi-sets instead of sets. The result of such a variant applied to the 5 queens program is displayed Figure 11.

3.2.2 Dynamic call graphs

A *static call graph* of a program is a graph where the nodes are labeled by the predicates of the program, and where arcs between nodes indicate potential predicate calls. We define the

Figure 12: The dynamic **call graph** of 5 queens

dynamic call graph of a logic program execution as the sub-graph of the (static) call graph composed of the arcs and nodes that have actually been traversed during the execution. For example, in Figure 12, we can see that predicate `main/2` calls predicates `data/1`, `queen/2`, and `print_list/2`. In this particular example, the static and dynamic call graphs are identical. Note that here, the queens program was linked with a version of the library that has been compiled without trace information. This is the reason why we do not see calls to, e.g. `io_write_string/3`.

An implementation of a monitor that builds the dynamic call graphs is given in Figure 13. In order to define this monitor, we use the same data structures as for the previous one, except that we replace the last traversed predicate by the whole call stack in the collected variable type (line 2). This stack is necessary in order to be able to get the direct ancestor of the current predicate. The set of arcs is initialized to the empty set (lines 4) and the stack is initialized to a stack that contains a fake node `user/0` (line 5). In order to construct the set of arcs, we insert at call events an arc from the previous predicate to the current one (line 12). The call stack is maintained on the fly by the `update_call_stack/4` predicate; the current predicate is pushed onto the stack at `call` and `redo` events (line 22), and popped at `exit`, `fail`, and `exception` events (line 24). The result of the execution of this monitor applied to the 5 queens program is displayed in Figure 12.

3.3 Test coverage

In this section, we define two notions of test coverage for logic programs, and we show how to measure the corresponding coverage rate of Mercury program executions using the `foldt` primitive. The aim here is not to provide the ultimate definition of test coverage for logic programs, but rather to propose two possible definitions, and to show how the corresponding

```

1  % Definition of pred, arc, and graph types omitted (cf previous monitor)
2  :- type accumulator_type ---> ct(stack(predicate), graph).
3
4  initialize(ct(Stack, set__init)) :-
5      stack__push(stack__init, p("user", 0), Stack).
6
7  collect(Event, ct(Stack0, Graph0), Acc) :-
8      Port = port(Event),
9      CurrentPred = p(proc_name(Event), proc_arity(Event)),
10     update_call_stack(Port, CurrentPred, Stack0, Stack),
11     ( if Port = call then
12         PreviousPred = stack__top_det(Stack0),
13         set__insert(Graph0, arc(PreviousPred, CurrentPred), Graph),
14         Acc = ct(Stack, Graph)
15     else
16         Acc = ct(Stack, Graph0) ).
17
18 :- pred update_call_stack(trace_port_type::in, predicate::in,
19     stack(predicate)::in, stack(predicate)::out) is det.
20 update_call_stack(Port, CurrentPred, Stack0, Stack) :-
21     ( if ( Port = call ; Port = redo ) then
22         stack__push(Stack0, CurrentPred, Stack)
23     else ( Port = fail ; Port = exit ; Port = exception ) then
24         stack__pop_det(Stack0, _, Stack)
25     else % other events
26         Stack = Stack0 ).

```

Figure 13: A monitor that computes dynamic call graphs

coverage rate measurements can be quickly prototyped. As a consequence, the proposed monitors cannot pretend to be optimal neither in functionality, nor in implementation.

3.3.1 Test coverage and logic programs

The aim of test coverage is to assess the quality of a test suite. In particular, it helps to decide whether it is necessary to generate more test cases or not. For a given coverage criterion, one can decide to stop testing when a certain percentage of coverage is reached.

The usual criterion used for imperative languages are *instruction* and *branch*³ criteria [Bei90]. The *instruction coverage rate* achieved by a test suite is the percentage of instructions that have been executed. The *branch coverage rate* achieved by a test suite is the percentage of branches that have been traversed during its execution.

One of the weaknesses of instruction and branch coverage is due to Boolean expressions. The problem occurs when a Boolean expression is composed by more than one atomic instruction: it may be that a test suite covers each value of the whole condition without covering all values of each atomic part of the condition. For example, consider the condition ‘A or B’ and a test suite where the two cases ‘A = true, B = false’ and ‘A = false, B = false’ are covered. In that case, every branch and every instruction is exercised, and nevertheless, B never succeeded. If B is erroneous, even 100% instruction and branch coverage will miss it. In logic programming, this problem is crucial because every instruction is a goal, i.e., a boolean expression.

3.3.2 Predicate coverage

In order to address the above problem, we need a coverage criterion that checks that each single predicate defined in the tested program succeeds and fails a given number of times. But we do not want to expect every predicate to fail because some, like printing predicates, are intrinsically deterministic. Therefore, we want a criterion that allows the test designer to specify how many times a predicate should succeed and fail. Therefore we define a *predicate criterion* as a pair composed of a predicate and a list of `exit` and `fail` ports that represent respectively successes and failures. In the case of Mercury, we can take advantage of the determinism declaration to automatically determine if a predicate should succeed and fail. Here is an example of predicate criterion that can be automatically defined according to the determinism declaration of each predicate:

‘det’ predicates:	1 success
‘semidet’ predicates:	1 success, 1 failure
‘multi’ predicates:	2 successes
‘nondet’ predicates:	2 successes, 1 failure

³where a branch is an arm of an alternative; for example, an `if-then-else` produces two branches.

```

1  :- type pred_crit ---> pc(proc_name, list(trace_port_type)).
2  :- type accumulator_type == list(pred_crit).
3
4  initialize([                pc("main", [exit]),
5  pc("data", [exit]),         pc("queen", [exit,exit,fail]),
6  pc("qperm", [exit,exit,fail]), pc("qdelete", [exit,exit,fail]),
7  pc("safe", [exit,fail]),     pc("nodiag", [exit,fail]),
8  pc("print_list", [exit]),    pc("print_list_2", [exit])  ]).
9
10 collect(Event, CSL0, CSL) :-
11     solutions(update_pred_list(port(Event), proc_name(Event), CSL0), Sol),
12     if Sol = [CLS1|_] then CSL = CSL1 else CSL = CSL0.
13
14 :- pred update_pred_list(trace_port_type::in, string::in, string::in,
15     accumulator_type::in, accumulator_type::out) is nondet.
16 update_pred_list(Port, ProcName, CSL0, CSL) :-
17     ( Port = exit ; Port = fail ),
18     list__delete(CSL0, pc(ProcName, Crit), CSL1),
19     list__delete(Crit, Port, NewCrit),
20     ( if NewCrit = [] then CSL = CSL1
21     else list__insert(pc(ProcName, NewCrit), CSL1, CSL) ).

```

Figure 14: A monitor that measures the predicate coverage rate of the `n queens` program

Then, we define the *predicate coverage rate of a logic program test suite* as the percentage of program predicate criteria that are covered during the execution of the suite. To compute that rate, one just needs to look at `exit` and `fail` events to see which criterion is covered.

Figure 14 shows a `foldt` monitor that measures the predicate coverage rate of the `queens` program. The type `predicate_criterion` is represented by a pair containing a procedure name and a list of ports (line 2). The accumulator will contain at each step the list of `predicate_criterion` not yet covered. The initial list of predicate criteria to be covered (lines 5 to 9) has been generated automatically by parsing the source file⁴. The call to `update_pred_list` succeeds if and only if the current event corresponds to an uncovered `pred_crit`; in that case `update_pred_list` outputs in its fourth argument the new value of the accumulator where the current `pred_crit` has been removed (lines 22 and 23). Otherwise, the accumulator remains unchanged (line 14).

3.3.3 Call site coverage

The previous coverage criterion only checks that at least one call of each predicate is covered. The problem is that 100% predicate coverage does not imply 100% instruction nor branch

⁴cf `extras/morphine/source/generate_pred_cov.m` available in the Mercury distribution

```

1  :- type call_site_crit ---> csc(proc_name,line_number,list(trace_port_type)).
2  :- type accumulator_type == list(call_site_crit).
3
4  initialize([
5      csc("queen", 16, [exit,fail]),      csc("print_list", 17, [exit]),
6      csc("write_string", 19, [exit]),    csc("qperm", 43, [exit,fail]),
7      csc("safe", 44, [exit,fail]),      csc("qdelete", 48, [exit,fail]),
8      csc("qperm", 50, [exit,fail]),     csc("qdelete", 54, [exit,fail]),
9      csc("nodiag", 58, [exit,fail]),    csc("safe", 59, [exit,fail]),
10     csc("nodiag", 73, [exit,fail]),    csc("write_string", 82, [exit]),
11     csc("write_string", 84, [exit]),    csc("print_list_2", 85, [exit]),
12     csc("write_string", 86, [exit]),    csc("write_int", 94, [exit]),
13     csc("write_string", 100, [exit]),   csc("print_list_2", 101, [exit])  ]).
14
15 collect(Event, CSL0, CSL) :-
16     solutions(update_call_site_list(port(Event), proc_name(Event),
17         line_number(Event), CSL0), Sol),
18     if Sol = [CSL1|_] then CSL = CSL1 else CSL = CSL0.
19
20 :- pred update_call_site_list(trace_port_type::in, string::in, string::in,
21     int::in, accumulator_type::in, accumulator_type::out) is nondet.
22 update_call_site_list(Port, ProcName, Ln, CSL0, CSL) :-
23     ( Port = exit ; Port = fail ),
24     list_delete(CSL0, csc(ProcName, Ln, Crit), CSL1),
25     list_delete(Crit, Port, NewCrit),
26     ( if NewCrit = [] then CSL = CSL1
27         else list__insert(csc(ProcName, Ln, NewCrit), CSL1, CSL) ).

```

Figure 15: A monitor that measures the call site coverage rate of the n queens program

coverage. To ensure 100% instruction and branch coverage, we need a criterion that ensures that every *predicate invocation* in the program succeeds and fails. Hence we need a definition attached to call sites (or goals) and not only to predicates. We define the *call site criterion* as being a tuple composed of a predicate, a line number, and a list of `exit` and `fail`. The *call site coverage rate of a logic program test suite* is therefore the percentage of program call site criteria that are covered during the execution of the test suite.

A monitor that measures call site coverage rate of Mercury program executions is given in Figure 15. This monitor is very similar to the previous one. The difference is that now we have line numbers in the criteria. Here again, the list of call sites to cover has been generated automatically⁵.

⁵cf the Mercury program `extras/morphine/source/generate_pred_cov.m`

4 Experimentations

In the previous section we have shown the flexibility and power of the `foldt` primitive. The aim of this section is to assess the performance of the current `foldt` implementation. When executing a monitor, some time is spent in the normal program execution (T_{prog}), some extra time is spent in the trace system of Mercury (Δ_{trace}), some extra time is spent in the interface between the tracer and the `foldt` mechanism⁶ (Δ_{inte}), some extra time is spent in the basic `foldt` mechanism (Δ_{foldt}), and some extra time is spent in the monitor itself ($\Delta_{\langle monitor \rangle}$). Hence, if we call T the execution time of a monitored program, we have:

$$T = T_{prog} + \Delta_{trace} + \Delta_{inte} + \Delta_{foldt} + \Delta_{\langle monitor \rangle}$$

In the following, we measure: T_{prog} , $T_{trace} = T_{prog} + \Delta_{trace}$, $T_{inte} = T_{prog} + \Delta_{trace} + \Delta_{inte}$, $T_{foldt} = T_{prog} + \Delta_{trace} + \Delta_{inte} + \Delta_{foldt}$ and $T_{\langle monitor \rangle} = T_{prog} + \Delta_{trace} + \Delta_{inte} + \Delta_{foldt} + \Delta_{\langle monitor \rangle}$ for the different monitors defined in the article.

We compare T_{trace} , T_{inte} , and T_{foldt} against T_{prog} . We will therefore compute the following ratios:

$$R_t = T_{trace}/T_{prog}, \quad R_i = T_{inte}/T_{prog} \quad \text{and} \quad R_f = T_{foldt}/T_{prog}$$

We also give the relative cost of each of the monitors:

$$\Delta_{\langle monitor \rangle} = T_{\langle monitor \rangle} \Leftrightarrow T_{foldt}$$

4.1 Methodology

Hardware and software. The measurements given in the following show the results of experiments run on a DELL inspiron 7500, with a 433 MHz Celeron, 192 Mb of RAM, running under the Linux 2.2.14 operating system. The machine was very lightly loaded; no X server, no network, simply the basic operating system and a Prolog process in a console to run the measurement scripts. The Prolog system is Eclipse 4.1 [Ecl99]. The Mercury compiler is a stable snapshot of 14 June 2001⁷. The results are consistent with experiments run on a SUN Sparc Enterprise 250 (2 x UltraSPARC-II, 296MHz, 512 Mb of RAM) running Solaris 2.7 (figures not given here).

Time measuring command. In order to measure the program execution times, we use the `benchmark_det` predicate of the `benchmarking.m` Mercury standard library. This predicate repeats the body of a program any given number of times. This is very important for small programs, as the startup cost very often dominates the execution cost. In the following experiments, each program is re-executed until it runs at least for 20 seconds. Each experiment has been done several times, and the deviation was smaller than 1 %.

⁶The Mercury predicate `collect` is called from the Mercury tracer which is written in C.

⁷The last official release is numbered 0.10.1

Monitored programs. The monitored programs are the Mercury benchmark suite⁸, composed of programs adapted from the Prolog benchmark suite of [RD92]. In order to have a wider range of execution sizes, we also measure `n queens` for $n=10,11$, as well as `mastermind`, a 1100 lines Mercury program which solves a mastermind game⁹.

Mercury compilation grades. In the following, the compilation grade g_{nt} refers to Mercury modules compiled with the command `mmc --grade asm_fast.gc.picreg`, which means that no trace event is generated. It is the grade used to measure the plain execution time of programs (T_{prog}).

The compilation grade g_t refers to Mercury modules compiled with the command `mmc --grade asm_fast.gc.picreg trace --deep -trace-optimized`, which means that all events related to all the predicates of the module, except library predicates, are generated. This grade is used in the following to measure the time spent in the basic trace system (T_{trace}), the time spent in the interface between the basic tracer and the `foldt` mechanism (T_{inte}), the time spent in the basic `foldt` mechanism (T_{foldt}) and the time spent in the monitors ($T_{<monitor>}$).

Measuring T_{prog} . If the programs are compiled in grade g_{nt} , then their execution does not produce any trace and their measured execution duration ($T_{measured}$) is exactly T_{prog} .

$$T_{prog} = T_{measured} \quad \text{if compilation grade is } g_{nt}$$

Measuring T_{trace} . If the programs are compiled in grade g_t , then their execution calls the Mercury tracer. In order to measure the cost of the basic tracer, we have to ensure that the tracer is systematically called at each event, but that it does not do anything else than entering and exiting the top-level switch of the trace system.

$$T_{trace} = T_{measured} \quad \text{if compilation grade is } g_t \text{ and} \\ \Delta_{inte} = 0, \quad \Delta_{foldt} = 0, \quad \text{and} \quad \Delta_{<monitor>} = 0$$

In order to ensure that $\Delta_{inte} = 0$, $\Delta_{foldt} = 0$ and $\Delta_{<monitor>} = 0$, we use the `continue` command of the Mercury tracer without specifying any break-point. Indeed, with that command, at each event, the trace system is entered; if the event does not correspond to one of the specified breakpoints, the normal execution is resumed. In our measurements, as no break-point is specified, the whole execution is traversed, and nothing is executed but the basic tracing mechanism.

⁸The source code of this benchmark suite can be found on the Mercury ftp site ftp://www.mercury.cs.mu.oz.au/pub/mercury/mercury-tests-*.tar.gz

⁹The full source code of the Mercury mastermind program can be found at the url <http://www.irisa.fr/lande/jahier/>

Measuring T_{inte} . In order to measure the cost of the interface between the basic tracer and the `foldt` mechanism, we have to ensure that the tracer is systematically called at each event, that it enters and exits the top-level switch of the trace system, that it prepares the context to call the `collect` predicate defined for the monitor, but that it does not compute anything else, in particular it should not retrieve any event attribute.

$$T_{inte} = T_{measured} \quad \text{if compilation grade is } g_t \text{ and } \Delta_{foldt} = 0, \text{ and } \Delta_{\langle monitor \rangle} = 0$$

In order to ensure that $\Delta_{foldt} = 0$, we have implemented a degenerate `foldt` such that no event attribute is computed (we have replaced these computations by void values). In order to ensure that $\Delta_{\langle monitor \rangle} = 0$, we use a monitor that does not compute anything (`collect(_E, A, A).`).

Measuring T_{foldt} . In order to measure the cost of the basic `foldt` mechanism, we have to ensure that `collect` is called at each event for a monitor that computes nothing.

$$T_{foldt} = T_{measured} \quad \text{if compilation grade is } g_t \text{ and } \Delta_{\langle monitor \rangle} = 0$$

In order to make sure that $\Delta_{\langle monitor \rangle} = 0$, we call `foldt` with a trivial monitor, that does not compute anything (`collect(_E, A, A).`).

Two of the current attributes are very costly to retrieve: the live arguments and the line number. The live arguments can be very large data structures. The line number corresponds to the line where the goal is called and not where the predicate is defined. It is dynamically retrieved. Many interesting monitors can be run without these attributes. Indeed, for the monitors we propose in this article, only one monitor uses the live arguments and one monitor uses the line number. Monitors that do not use these costly attributes can disable them. As a consequence, the measurements of T_{foldt} is made with these two attributes disabled.

Measuring $\Delta_{\langle monitor \rangle}$. In order to measure the cost of a given monitor, we simply subtract from its execution time the execution time of the empty monitor

$$\Delta_{\langle monitor \rangle} = T_{measured} \Leftrightarrow T_{foldt} \quad \text{for each monitor.}$$

If the line number or the live arguments are needed, they are explicitly retrieved. Their retrieval cost therefore appears in the related $\Delta_{\langle monitor \rangle}$.

4.2 Resulting tables

In this subsection we introduce two tables of measurements. The discussion of their contents is in the following section.

Table 1 illustrates the cost of the basic tracer, the cost of the interface between the tracer and `foldt`, as well as the cost of `foldt` on the benchmark programs described in the previous sections. The first column contains the names of the monitored programs. The second column contains the numbers of execution events generated by the program

program	events	T_{prog} in <i>ms</i>	T_{trace} in <i>ms</i>	R_t	T_{inte} in <i>ms</i>	R_i	T_{foldt} in <i>ms</i>	R_f	R_f^*
queens-5	698	0.03	0.24	7.5	2	61.5	2.07	63	9.5
query	935	0.09	0.28	3	1.53	16.5	1.58	17	3.5
deriv	1,540	0.05	0.11	2.5	0.64	14	0.66	14.5	3
qsort	1,564	0.1	0.48	5	4.03	42	4.16	43.5	6.5
nrev	1,619	0.14	0.54	4	4.44	32.5	4.58	33.5	5
primes	2,192	0.21	0.8	4	6.23	30	6.51	31	5.5
cqueens	3,789	0.14	1.26	9.5	11.11	81	11.39	82	11.5
crypt	4,602	0.72	1.8	3	11.16	16	11.54	16.5	3.5
poly	79,070	6.44	29	4.5	226.2	35.5	233.4	36.5	6
tak	190,831	3.88	57.1	15	553.7	142	572.4	148	20
queens-10	4,289,986	257	1530	6	12760	50.5	13200	51.5	8
mastermind	9,106,510	3630	6500	2	30490	8.5	31520	9	2.5
queens-11	32,384,320	2103	12030	6	97010	46.5	100190	48	7.5

Table 1: Cost of the `foldt` mechanisms on benchmarks. $T_{trace} = T_{prog} + \Delta_{trace}$, $T_{inte} = T_{prog} + \Delta_{trace} + \Delta_{inte}$, $T_{foldt} = T_{prog} + \Delta_{trace} + \Delta_{inte} + \Delta_{foldt}$, $R_t = T_{trace}/T_{prog}$, $R_i = T_{inte}/T_{prog}$, $R_f = T_{foldt}/T_{prog}$, $R_f^* = (T_{foldt} \ominus \Delta_{inte})/T_{prog}$

executions compiled in grade \mathbf{g}_t (all events are generated, except events relative to library predicates). The programs are sorted wrt this number of events, from `queens-5`, 698 events, to `queens-11`, more than 32 millions of events. The third column contains the execution times of the programs compiled without any trace information (T_{prog}). The fourth column contains the execution times of the programs compiled in grade \mathbf{g}_t and run under the control of the tracer without tracing anything ($T_{trace} = T_{prog} + \Delta_{trace}$). The fifth column contains the overhead factor of the basic trace mechanism ($R_t = T_{trace}/T_{prog}$). The sixth column contains the execution times of the programs compiled in grade \mathbf{g}_t and run under the control of the tracer, and where the degenerate `foldt` is called with an empty monitor ($T_{inte} = T_{prog} + \Delta_{trace} + \Delta_{inte}$). The seventh column contains the overhead factor of the trace and the interface between the tracer and `foldt` ($R_i = T_{inte}/T_{prog}$). The eighth column contains the execution time of the programs compiled in grade \mathbf{g}_t and run under the control of `foldt` with an empty monitor ($T_{foldt} = T_{prog} + \Delta_{trace} + \Delta_{inte} + \Delta_{foldt}$). The ninth column contains the overhead factor of the trace, the interface and the basic `foldt` mechanism ($R_f = T_{foldt}/T_{prog}$). The tenth column contains the overhead factor of the trace and the basic `foldt` mechanism, with the interface cost subtracted ($R_f^* = (T_{foldt} \ominus \Delta_{inte})/T_{prog}$). The time measurements have been rounded off two digits after the dot. The ratios have been rounded up to the nearest half.

Table 2 illustrates the relative cost of each of the monitors presented earlier in this article. For each monitor and each benchmark program, we give the execution time where the cost of the basic `foldt` mechanism related to the benchmark program has been subtracted ($\Delta_{\langle monitor \rangle} = T_{\langle monitor \rangle} \ominus T_{foldt}$). The time measurements in *ms* have been rounded off one digit after the dot. The time measurements in seconds have been rounded off to the closest integer. Two results are noted ε because they were below $10^{-1}ms$.

program	Δ_{call} in <i>ms</i>	Δ_{stat} in <i>ms</i>	Δ_{histo} in <i>ms</i>	Δ_{sol} in <i>ms</i>	Δ_{cfg} in <i>ms</i>	Δ_{dcg} in <i>ms</i>	Δ_{pred_cov} in <i>ms</i>	Δ_{call_cov} in <i>ms</i>
queens-5	1.3	0.1	1.0	2	2.0	3.7	4.7	25.0
query	0.9	ε	0.7	1.4	1.0	1.1	2.7	8.6
deriv	0.4	ε	0.3	0.6	0.5	0.5	2.0	47.4
qsort	2.6	0.1	2.0	3.9	4.2	6.1	7.5	28.4
nrev	2.9	0.1	2.1	4.4	4.0	4.2	8.8	27.5
primes	4.0	0.2	3.0	5.8	6.2	7.7	14.0	45.5
cqueens	7.3	0.3	5.2	9.6	9.9	15.4	28.0	118.7
crypt	6.9	0.3	5.0	11.2	11.7	19.8	42.8	234.3
poly	148.6	6.3	107.2	217.2	340.2	603.9	979.6	$11 \cdot 10^3$
tak	355.8	13.2	264.0	479.8	521.7	782.8	774.2	$2 \cdot 10^3$
queens-10	$8 \cdot 10^3$	410.0	$6 \cdot 10^3$	$13 \cdot 10^3$	$12 \cdot 10^3$	$23 \cdot 10^3$	$38 \cdot 10^3$	$168 \cdot 10^3$
mastermind	$25 \cdot 10^3$	300.0	$10 \cdot 10^3$	$37 \cdot 10^3$	$36 \cdot 10^3$	$31 \cdot 10^3$	$437 \cdot 10^3$	$5947 \cdot 10^3$
queens-11	$63 \cdot 10^3$	$3 \cdot 10^3$	$45 \cdot 10^3$	$101 \cdot 10^3$	$87 \cdot 10^3$	$179 \cdot 10^3$	$296 \cdot 10^3$	$1326 \cdot 10^3$

Table 2: Relative cost of the monitors defined in this article. $\Delta_{\langle monitor \rangle} = T_{\langle monitor \rangle} \Leftrightarrow T_{foldt}$

Δ_{call} gives the relative execution time of the monitor which counts the number of goal invocations, see Section 2.3.1. Δ_{stat} gives the relative execution time of the monitor which counts the number of events at each port, see Section 3.1.1. Δ_{histo} gives the relative execution time of the monitor which counts the number of calls at each depth, see Section 3.1.2. Δ_{sol} gives the relative execution time of the monitor which collects all solutions, see Section 3.1.3. Δ_{cfg} gives the relative execution time of the monitor which builds the dynamic control flow graph, see Section 3.2.1. Δ_{dcg} gives the relative execution time of the monitor which builds the dynamic call graph, see Section 3.2.2. Δ_{pred_cov} gives the relative execution time of the monitor which computes the predicate coverage ratio, see Section 3.3.2. Δ_{call_cov} gives the relative execution time of the monitor which computes the call site coverage ratio, see Section 3.3.3.

4.3 Discussion

In this section, we discuss the resulting ratios of Table 1 and the relative costs of the monitors of Table 2.

Two extremes: tak and mastermind. The `tak` program has ratios much higher than the other programs. The tracer overhead is already 15, the interface overhead is 142 and the `foldt` overheads are 148 and 20. This program is actually a single predicate four times recursive. It already broke the stacks of the reference tracer used in [DN00]. This code is an extreme case to test compiler optimization capabilities. As many optimizations, such as last call optimization, cannot be applied to traced code, the better the compiler is, the worse debugger and monitor ratios look. Program `Tak` is very untypical of casual programs. The ratios related to `tak` are not taken into account in the averages given below.

On the other hand `mastermind` has very low ratios. The tracer overhead is 2, the interface overhead is 8.5 and the `foldt` overheads are 9 and 2.5. The `mastermind` program uses a lot of library predicates which are not traced and monitored in detail. This is very typical of casual programs. It is very encouraging that the more realistic program has the best ratios. However, the other benchmarks do not use untraced libraries, in order to be fair, the ratios related to `mastermind` are not taken into account in the averages given below.

In the following, averages are, thus, computed without the figures related to `tak` and `mastermind`.

The ratios are not correlated with the number of events. Program `queens-5`, which has only 698 events, has very bad ratios, whereas `crypt`, almost five thousand events, and `poly`, almost eighty thousand events, have better ratios than the average. The same program, `n queens`, run for $n=5,10$ and 11, always has comparable ratios. This seems to indicate that the overheads depend mostly of the monitored program and is somewhat constant for a given program.

Overhead of the tracer. The average of the tracer overhead is 5. It is a very good ratio. Prolog tracers can easily have an overhead over 20 [DN00]. A low ratio for the tracer is, of course, a good starting point to build efficient generic monitors.

Overhead of the interface between the tracer and `foldt`. The average of the interface overhead is 39. This is very high and it is the main source of inefficiency of our current implementation. It illustrates how crucial the implementation of this interface is for efficient generic monitoring.

The problems comes from the fact that the monitor programs have to be integrated in the tracer. In our particular case, the Mercury predicate `collect` is called from the Mercury tracer which is written in C. In order to call Mercury code from C with the current (low-level back-end of the) Mercury compiler, machine registers need to be saved and restored. Since the `collect` predicate is called several million times, this has a noticeable influence on the performance. A way to remove this overhead could be to use the new MLDS back-end of the Mercury compiler, which generates high-level C code that does not use machine registers; unfortunately, the trace system is not supported for the MLDS back-end at the time of this writing. Once the MLDS Mercury back-end is available, calling the `collect` predicate will actually be compiled as a simple C procedure call from within C code. The overhead of the interface between the tracer and `foldt` should thus become negligible.

One important lesson learned from these measurements is a follows. Whether the monitors and the tracer are implemented in the same programming language or not, the integration of the compiled monitors should not cost more than a procedure call. The monitors must therefore be compiled into the same target language as the tracer. Furthermore, no run-time verifications should be made. The monitors should therefore have no side-effect on the traced execution. It should thus be statically checked that monitors only update their own (fresh) variables.

Overhead of foldt. The average of the `foldt` overhead is 40. Most of it is coming from the overhead of the interface discussed above. Assuming that the above fix could be done and that the interface overhead indeed becomes negligible, we have computed an ideal ratio: $R_f^* = (T_{foldt} \Leftrightarrow \Delta_{inte}) / T_{prog}$. The average of the overhead of `foldt` without the interface cost is 6.5. As the average of the tracer overhead is 5, we can say that the `foldt` overhead is acceptable.

In the context of Mercury, this is especially true because the developers of Mercury claim that Mercury programs executed in trace mode are faster than the equivalent Prolog programs executed in optimized mode in the faster Prolog systems [SH99].

For example, T_{fold} of `mastermind` is only approximatively half of a second, even with the bad interface overhead. And if we subtract Δ_{inte} from T_{fold} in the worse case, ie for `queens-11`, the resulting time is only 15 s.

Unused event attributes. As already mentioned, some attributes can be very costly to compute. When they are not needed it should be possible to disable them. In the current implementation of `foldt`, this is already the case for the list of live arguments and the line numbers. The `foldt` overhead has been measured without the cost of the live arguments and the line numbers. Some measurements, not reported here, showed that this has an impact on the performance.

Granularity of the instrumentation. In order to measure the worst case, we have made the Mercury tracer systematically generate all the possible events. Not all the monitors need such a fine grained instrumentation. For example, the monitor that counts the number of events at each depth, only `call` events are necessary. When manually implementing monitors, people only instrument the code where it is necessary. As a matter of fact, the Mercury tracer enables users to specify what type of events, if any, should be generated for a given module (the only restriction is that, if some events are generated for a predicate, `call` events must be present). Hence, programmers can already take advantage of this possibility to optimize their monitor. As further work, we plan to automate this optimization, namely, to automatically generate the appropriate compiling option for a given monitor.

Relative costs of monitors (Table 2). One can notice that surprisingly, the `count_call` monitor is more time demanding than the `statistics` one. The reason is probably that `count_call` collects an integer whereas `statistics` collects an array, which is a Mercury data structure that can be destructively updated. At the time of writing, the current implementation of the Mercury compiler cannot destructively update general data types yet (although work is in progress). This optimization would have a major performance impact on most of the monitors presented here (namely, all the monitors excepted the ones that are collecting arrays).

The live argument attribute is used by only one monitor, to compute all the solutions; Δ_{sol} takes into account the extra time needed to retrieve the live arguments. One monitor

uses the line number attribute, to compute call site coverage; Δ_{call_cov} takes into account the extra time needed to retrieve the line numbers.

Most of the monitor intrinsic costs are below a few seconds. Some monitors take a few minutes. The worst cases appear for the computation of the call site coverage of the biggest two executions: `mastermind` and `queens-11`. They run respectively for 1 h. 40 min. and 22 min. The monitor which computes the call site coverage could be optimized in order to avoid traversing too long lists. We conjecture that an ad hoc monitor programmed in the same way would take a comparable amount of time.

Conclusion on performance. With a fast tracer, an interface between the tracer and `foldt` reduced to procedure calls, and the possibility to disable the computation of heavy non-necessary attributes, then generic monitoring can be efficient.

5 Related work

Programmable debuggers. We designed 3 programmable debuggers, Opium for Prolog [Duc99b], Coca for C [Duc99a] and Morphine for Mercury [JD99a]. They are based on a Prolog query loop plus a handful of coroutines primitives connected to the trace system. Those primitives allow the Prolog interpreter to communicate with the debugged program. Opium, Coca and Morphine are full debugging programming languages in which all classical debuggers commands can be implemented straightforwardly and efficiently. However, it appears that their set of primitives are not well suited for monitoring. All the monitors implemented with `foldt` can easily be implemented with the debugger set of primitives [JD99b], but the resulting monitors require too many context switches and too much socket traffic between the program and the monitor. With programs of several million execution events, such monitors are four orders of magnitude slower than their counterparts that use `foldt` [Jah00].

Automated development of monitors. Jeffery et al. designed Alamo [TJ98, JZTB98, Jef99], an architecture that aims at easing the development of monitors for C programs. As in our approach, their monitoring architecture is based on event filtering, and monitors can be programmed. Their system performs trace extraction whereas we rely on an already available tracer; this saves us a tedious task which has already been done and optimized. On the other hand, we do not have the full control of the information available in the trace. Note, however, that, so far, we have been able to reconstruct missing information, for example the call stack of Figure 13. Moreover, to avoid code explosion, Jeffery et al. perform part of the event filtering at compilation time. They need to recompile the program each time they want to execute another monitor whereas we only need to dynamically link the monitor to the monitored program. Alamo and the monitored program are running in coroutines, but within the same address space.

Eustace and Srivastava developed Atom [ES95], a system that also aims at easing the implementation of monitors. The difference with Alamo is that monitors are implemented

with procedure calls and global variables which is much more efficient than coroutines. However, the language provided by Atom is less expressive than the Alamo's. Alamo and Atom have influenced the design of `foldt` and we tried to take the best of both: a full and high-level programming language implemented by procedure calls. The advantages of our architecture over [ES95] and [JZTB98] are the following:

- A higher level interface makes the code of our monitors compact, easy to write and read, and therefore to maintain. Of course, this point is difficult to assess. We hope that the numerous and various examples given in Section 3 convince the reader that it is indeed the case. Nonetheless, we see three conjectures explaining why the code is more compact and more elegant. Firstly, users do not have to deal with code instrumentation directives; the instrumentation has already been done. Secondly, we take advantage of the expressive power of `fold`; using high order predicates such as `fold` for processing lists has proven to be concise and far less prone to error than processing lists manually. A third reason is that our monitors are written in Mercury, which is a considerably higher level language than C or C-like languages which are used in [ES95, JZTB98].
- Provided that an event-oriented tracer exists, the `foldt` operator is easy to implement. To implement it, the work done inside the Mercury runtime system, which corresponds to the really technical part, amounts to only 61 modified lines and 292 new lines of (C) code.

Kishon and al. [KH95] use a denotational and operational continuation semantics to formally define monitors for a simple functional programming language. The kind of monitors they define are profilers, debuggers, and statistic collectors. From the operational semantics, a formal description of the monitor, and a program, they derive an instrumented executable file that performs the specified monitoring activity. The semantics of the original programs is preserved. They use partial evaluation to make their monitors reasonably efficient. The main disadvantage with this approach is that they are rebuilding a whole execution system from scratch, without taking advantage of existing compilers. We strongly believe that it is important to have the same execution system for debugging, monitoring and producing the final executable program. As noted by [BHS92], some errors only occur in presence of optimizations, and vice versa; some programs can only be executed in their optimized form because of time and memory constraints; when searching for “hot spots”, it is better to do it as much as possible with the optimized program as many things can be optimized away; and finally, sometimes, the error comes from the optimizations themselves. In our setting we can easily mix traced and non-traced code.

Efficient monitoring. Patil and Fisher [PF97] address the problem of performance monitoring by delegating the monitoring activities to a second processor that they call a shadow processor. Their approach is very efficient; the monitored program is practically not slowed down, but the set of monitoring commands they propose cannot be extended.

We mentioned in the previous section that we could reduce the number of events generated by the tracer. For example, in [BL92], given a static control flow graph, algorithms can place tracing instructions in optimal ways for computing statistics on imperative program executions.

6 Conclusion

In this article we have proposed a generic monitoring framework based on `foldt`, a high-level primitive that allows users to easily specify what they want to monitor. We illustrated it on various examples that demonstrate its genericity and its simplicity of use. We defined two notions of test coverage for logic programs and showed how to prototype coverage rates measurements with our primitive. To our knowledge no definition of test coverage existed for logic programming so far. More testing and monitoring tools are missing from many declarative systems: `foldt` allows some of these tools to be easily defined and implemented. Measurements showed that the performance of the primitive on the above examples can be acceptable for executions of several million trace events.

To sum up the advantages of our framework, we can say that it is:

- Easy to implement: because it is based on an existing tracer (292 new, and 61 modified lines of codes in our current implementation.)
- Efficient: because the trace is not stored (and provided that only procedure calls are made.)
- Flexible and easy to use: as illustrated by the given applications about execution profiles, graphical abstract displays and test coverage.

Acknowledgments

We would like to thank Fergus Henderson for his technical support and his many contributing ideas; Pierre Deransart, Baudouin Le Charlier and Olivier Ridoux for fruitful discussions; Jean-Philippe Pouzol for his comments on earlier versions of this article.

A Mercury execution events

The Mercury trace is an adaptation of Byrd’s box model [Byr80]. In this section, we describe the Mercury execution events that constitute the Mercury execution trace. More information about the Mercury tracer can be found in [SH99]. The different *attributes* provided by the Mercury tracer are:

1. *Chronological event number* (**chrono**¹⁰). Each event has a unique event number according to its rank in the trace. It is a counter of events.
2. *Goal invocation number* or *call number* (**call**). Unlike chronological event number, several events have the same goal invocation number. All events related to a given goal have a unique goal number given at invocation time.
3. *Execution depth* (**depth**). It is the depth of the goal in the proof tree, namely the number of its ancestor goals + 1.
4. *Event type or port* (**port**). We distinguish between *external events* that occur at procedure entries and exits, which are the traditional ports introduced by Byrd [Byr80], and the *internal events* which refers to what is occurring inside a procedure. External events are:

- **call** a new goal is invoked
- **exit** the current goal succeeds
- **fail** the current goal fails
- **redo** another solution for the current goal is asked for on backtracking.
- **exception** the execution raises an exception

Internal events are:

- **disj** the execution is entering a branch of a disjunction
- **switch** the execution is entering a branch of a switch (a *switch* is a disjunction in which each branch unifies a ground variable with a different function symbol. In that case, at most one disjunction provides a solution).
- **if** the execution is entering the condition branch of an if-then-else
- **then** the execution is entering the “then” branch of an if-then-else
- **else** the execution is entering the “else” branch of an if-then-else
- **first** the execution enters a C code fragment for the first time
- **later** the execution re-enters a C code fragment

¹⁰The names of the attribute accessing functions are in bold in between parentheses.

5. *Determinism (det)*. It characterizes the number of potential solutions for a given goal. The determinism markers of Mercury are: **det** for procedures which have exactly 1 solution, **semidet** for those which have 0 or 1 solution, **nondet** for those which have any number of solutions, **multi** for those which have at least 1 solution, **failure** for those that have no solution, and **erroneous** for those which lead to a runtime error.
6. *Procedure (proc)*. It is defined by:
 - a *flag* telling if the procedure is a function or a predicate (**proc_type**)
 - a *definition module* (**def_module**)
 - a *declaration module* (**decl_module**) The declaration module is the module where the user has declared the procedure. The defining module is the module where the procedure is effectively defined from the compiler point of view. They may be different if the procedure has been inlined.
 - a *name* (**name**)
 - an *arity* (**arity**)
 - a *mode number* (**mode_number**).
The mode number is an integer coding the mode of the procedure. When a predicate has only one mode, the mode number of its corresponding procedure is 0. Otherwise, the mode number is the rank in the order of appearance of the mode declaration.
7. *List of live arguments (args)*. A variable is *live* at a given point of the execution if it has been instantiated and if the result of that instantiation is still available in the runtime system. Destructive input (**di** mode), for example, are not kept until the procedure exits.
8. *List of live Argument types (arg_types)*.
9. *List of local live variables (local_vars)*. Some live variables are not arguments of the current procedure.
10. *Goal path (goal_path)*. The goal path indicates in which part of the code the current internal event occurs. **if**, **then** and **else** branches of an *if-then-else* are denoted by **?**, **e** and **t** respectively; *conjuncts*, *disjuncts* and *switches* are denoted by **ci**, **di** and **si**, where **i** is the conjunct (resp. disjunct, switch) number. For example, if an event with goal path [**c3**, **e**, **d1**] is generated, it means that the event occurred in the first branch of a disjunction, which is in the else branch of an if-then-else, which is in the third conjunction of the current goal. External events have an empty goal path.

The event structure is illustrated by Figure 16. The displayed structure is related to an event of the execution of a **qsort** program which sorts the list of integers [3, 1, 2] using a *quick sort* algorithm. The information contained in that structure indicates

chrono	10
call	6
depth	5
port	then
det	det
proc_type	predicate
def_module	qsort
decl_module	qsort
name	partition
arity	4
mode_number	0
arg	[[1, 2], 3, -, -]
arg_types	[list(int), int, -, -]
local_vars	[live_var("H", 1, int), live_var("T", [2], list(int))]
goal_path	[s1, c2, t]

Figure 16: A Mercury trace event

that procedure `qsort:partition/4-0`¹¹ is currently invoked, it is the tenth trace event being generated, the sixth goal being invoked, and it has four ancestors (depth is 5). At this point, only the first two arguments of `partition/4` are instantiated: the first one is bound to the list of integers `[1, 2]` and the second one to the integer `3`; the third and fourth arguments are not live, which is indicated by the atom `'-'`. There are two live local variables: `H`, which is bound to the integer `1`, and `T`, which is bound to the list of integers `[2]`. The goal path tells that this event occurred in the `then` branch (`t`) of the second conjunction (`c2`) of the first `switch` (`s1`) of `partition/4`.

B The Mercury queens program

```
:- module queens.

:- interface.

:- import_module io.

:- pred main(io__state, io__state).
:- mode main(di, uo) is cc_multi.

:- implementation.
```

¹¹`'-0'` denotes the mode number; here, `0` means that `qsort` was declared with only one mode (namely, `:- mode qsort(in, in, out, out) is det`). If more than one mode is declared, `'-1'` denotes the first mode, `'-2'` the second one, etc.

```

:- import_module list, int.

main -->
  ( { data(Data), queen(Data, Out) } ->
    io__write_string("A 5 queens solution is "), print_list(Out)
    ;
    io__write_string("No solution\n")
  ).

:- pred data(list(int)).
:- mode data(out) is det.

:- pred queen(list(int), list(int)).
:- mode queen(in, out) is nondet.

:- pred qperm(list(T), list(T)).
:- mode qperm(in, out) is nondet.

:- pred qdelete(T, list(T), list(T)).
:- mode qdelete(out, in, out) is nondet.

:- pred safe(list(int)).
:- mode safe(in) is semidet.

:- pred nodiag(int, int, list(int)).
:- mode nodiag(in, in, in) is semidet.

data([1,2,3,4,5]).

queen(Data, Out) :-
  qperm(Data, Out),
  safe(Out).

qperm([], []).
qperm([X|Y], K) :-
  qdelete(U, [X|Y], Z),
  K = [U|V],
  qperm(Z, V).

qdelete(A, [A|L], L).
qdelete(X, [A|Z], [A|R]) :-
  qdelete(X, Z, R).

safe([]).
safe([N|L]) :-
  nodiag(N, 1, L),
  safe(L).

nodiag(_, _, []).
nodiag(B, D, [N|L]) :-
  NmB is N - B,
  BmN is B - N,
  ( D = NmB ->

```

```

        fail
; D = BmN ->
        fail
;
        true
),
D1 is D + 1,
nodiag(B, D1, L).

:- pred print_list(list(int), io__state,
        io__state).
:- mode print_list(in, di, uo) is det.

print_list(Xs) -->
(
    { Xs = [] }
->
    io__write_string("[]\n")
;
    io__write_string("[",
        print_list_2(Xs),
        io__write_string("]\n")
).

:- pred print_list_2(list(int),
        io__state, io__state).
:- mode print_list_2(in, di, uo) is det.

print_list_2([]) --> [].
print_list_2([X|Xs]) -->
    io__write_int(X),
    (
        { Xs = [] }
->
        []
;
        io__write_string(", "),
        print_list_2(Xs)
).

```

References

- [Bat95] Peter C. Bates. Debugging heterogeneous distributed systems using event-based models of behavior. *ACM Transactions on Computer Systems*, 13(1):1–31, February 1995.
- [Bei90] B. Beizer. *Software testing techniques*, volume 2nd ed. Int. Thomson Computer Press, 1990.
- [BHS92] G. Brooks, G.J. Hansen, and S. Simmons. A new approach to debugging optimized code. In *SIGPLAN '92 Conf. on Programming Language Design and Implementation*, pages 1–11, 1992.
- [Bir87] Richard S Bird. An introduction to the theory of lists. In M. Broy, editor, *Logic of Programming and Calculi of Discrete Design*, pages 3–42. Springer-Verlag, 1987.
- [BL92] T. Ball and J.R. Larus. Optimally profiling and tracing programs. In *Principles of Programming Languages*, 1992.
- [Byr80] L. Byrd. Understanding the control flow of Prolog programs. In S.-A. Tärnlund, editor, *Logic Programming Workshop*, 1980.
- [DN00] M. Ducassé and J. Noyé. Tracing Prolog programs by source instrumentation is efficient enough. *Journal of Logic Programming*, 43(2):157–172, May 2000.
- [Duc99a] M. Ducassé. Coca: An automated debugger for C. In *Proc. of the 21st Int. Conf. on Software Engineering*, pages 504–513. ACM Press, May 1999.
- [Duc99b] M. Ducassé. Opium: An extendable trace analyser for Prolog. *Journal of Logic programming*, 39:177–223, 1999. Special issue on Synthesis, Transformation and Analysis of Logic Programs, A. Bossi and Y. Deville (eds).
- [EB88] M. Eisenstadt and M. Brayshaw. The Transparent Prolog Machine TPM: an execution model and graphical debugger for Logic Programming. *Journal of Logic Programming*, 5(4):277–342, 1988.
- [Ecl99] Eclipse. *The ECLiPSe Constraint Logic Programming System, ECLiPSe 4.1 - User Manual*. IC.Parc, 1999. <http://www-icparc.doc.ic.ac.uk/eclipse/>.
- [ES95] A. Eustace and A. Srivastava. ATOM: A flexible interface for building high performance program analysis tools. In *Winter 1995 USENIX Conf.*, 1995.
- [FAS94] P. Fritzson, M. Auguston, and N. Shahmehri. Using assertions in declarative and operational models for automated debugging. *Journal of Systems Software*, 25:223–239, 1994.

- [Ger00] M. Gerndt. Towards automatic performance debugging tools. In M. Ducassé, editor, *Proceedings of the Fourth International Workshop on Automated Debugging (AADEBUG 2000)*, <http://xxx.lanl.gov/html/cs/0010035>, August 2000. Computer research repository.
- [GL91] K.B. Gallagher and J.R. Lyle. Using program slicing in software maintenance. *IEEE Transactions on Software Engineering*, 17(8):751–761, 1991.
- [Hat97] L. Hatton. Does OO sync with the way we think? *IEEE Software*, 15(3):46–54, 1997.
- [HS96] W.E. Howden and G.M. Shi. Linear and structural event sequence analysis. In Steven J. Zeil, editor, *Proc. of the 1996 Int. Symp. on Software Testing and analysis*, pages 98–106. ACM Press, 1996.
- [Hut99] G. Hutton. A tutorial on the universality and expressiveness of fold. *Jour. of Functional Programming*, pages 355–372, 1999.
- [Jah00] E. Jahier. *Analyse dynamique de programmes : mise en oeuvre automatisée d'analyseurs performants et spécification de modèles d'exécution*. PhD thesis, INSA de Rennes, 2000. Partly in English.
- [JD99a] E. Jahier and M. Ducassé. *Morphine 0.2 User and Reference Manuals*, 1999.
- [JD99b] E. Jahier and M. Ducassé. Un traceur d'exécutions de programmes ne sert pas qu'au débogage. In F. Fages, editor, *Actes des Journées francophones de Programmation Logique et par Contraintes*, pages 297–311. Hermès, 1999.
- [Jef99] C. Jeffery. *Program monitoring and visualization*. Springer, 1999. ISBN 0-387-98644-8.
- [JZTB98] C. Jeffery, W. Zhou, K. Templer, and M. Brazell. A lightweight architecture for program execution monitoring. *ACM SIGPLAN Notices*, 33(7):67–74, 1998.
- [KH95] A. Kishon and P. Hudak. Semantics directed program execution monitoring. *Journal of Functional Programming*, 5(4):501–547, 1995.
- [KHC91] A. Kishon, P. Hudak, and C. Consel. Monitoring semantics: a formal framework for specifying, implementing and reasoning about execution monitors. *ACM Sigplan Notices*, 26(6):338–352, 1991.
- [KN91] E. Koutsoufios and S. North. Drawing graphs with *dot*. TR 910904-59113-08TM, AT&T Bell Laboratories, 1991.
- [OCH90] R.A. Olsson, R.H. Crawford, and W.W. Ho. Dalek: A GNU, improved programmable debugger. In USENIX Association, editor, *Summer 1990 USENIX Conf.: 1990, Anaheim, California, USA*, pages 221–232. USENIX Association, 1990.

-
- [PF97] H. Patil and C. Fischer. Low-cost, concurrent checking of pointer and array accesses in C programs. *Software Practice and Experience*, 27(1):87–110, 1997.
- [Pfe92] F. Pfenning, editor. *Types in Logic Programming*. MIT Press, 1992. ISBN 0-262-16131-1.
- [RD92] P. Van Roy and A. M. Despain. High-performance logic programming with the Aquarius Prolog compiler. *Computer*, 25(1):54–68, January 1992.
- [SDBP98] J. Stasko, J. Domingue, M. H. Brown, and B. A. Price, editors. *Software Visualization: Programming as a Multimedia Experience*. MIT Press, 1998.
- [SH99] Z. Somogyi and F. Henderson. The implementation technology of the mercury debugger. In *proceedings of the Tenth Workshop on Logic Programming Environments*, volume 30(4). Elsevier, Electronic Notes in Theoretical Computer Science, 1999. <http://www.elsevier.nl/cas/tree/store/tcs/free/entcs/store/tcs30/cover.sub.sht>.
- [SHC96] Z. Somogyi, F. Henderson, and T. Conway. The execution algorithm of Mercury, an efficient purely declarative logic programming language. *Journal of Logic Programming*, 29:17–64, 1996.
- [TA95] A. Tolmach and A.W. Appel. A debugger for Standard ML. *Journal of Functional Programming*, 5(2):155–200, 1995.
- [Tip95] F. Tip. Generic techniques for source-level debugging and dynamic program slicing. In Peter D. Mosses, Mogens Nielsen, and Michael I. Schwartzbach, editors, *TAPSOFT '95: Theory and Practice of Software Development*, volume 915 of *LNCS*, pages 516–530. Springer-Verlag, 1995.
- [TJ98] K.S. Templer and C.L. Jeffery. A configurable automatic instrumentation tool for ANSI C. In *Proc. Automated Software Engineering Conf.* IEEE Computer Society, 1998.

Contents

1	Introduction	3
2	A high-level trace processing operator: foldt	5
2.1	Language independent foldt definition	5
2.2	An implementation of foldt for Mercury	6
2.2.1	Mercury and its trace	7
2.2.2	The foldt implementation	7
2.3	The current user interface of foldt for Mercury	8
2.3.1	Defining monitors	8
2.3.2	Invoking foldt	10
2.3.3	Illustration of the advantage of calling foldt from a Prolog query loop	10
3	Applications	12
3.1	Execution profiles	12
3.1.1	Counting the number of events at each port	12
3.1.2	Counting the number of calls at each depth	13
3.1.3	Collecting solutions	13
3.2	Graphical abstract views	14
3.2.1	Dynamic control flow graphs	14
3.2.2	Dynamic call graphs	16
3.3	Test coverage	17
3.3.1	Test coverage and logic programs	19
3.3.2	Predicate coverage	19
3.3.3	Call site coverage	20
4	Experimentations	22
4.1	Methodology	22
4.2	Resulting tables	24
4.3	Discussion	26
5	Related work	29
6	Conclusion	31
A	Mercury execution events	32
B	The Mercury queens program	34



Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY
Unité de recherche INRIA Rennes, Irista, Campus universitaire de Beaulieu, 35042 RENNES Cedex
Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN
Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex
Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

Éditeur
INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)
<http://www.inria.fr>
ISSN 0249-6399