

Experiment Management Support for Performance Tuning

Karen L. Karavanic and Barton P. Miller
 {karavan,bart}@cs.wisc.edu

Computer Sciences Department
 University of Wisconsin
 Madison, WI 53706-1685

1 INTRODUCTION

“An experimental science is supposed to do experiments that find generalities. It’s not just supposed to tally up a long list of individual cases and their unique life histories. That’s butterfly collecting.”

- Richard C. Lewontin [1]

The development of a high-performance parallel system or application is an evolutionary process. It may begin with models or simulations, followed by an initial implementation of the program. The code is then incrementally modified to tune its performance and continues to evolve throughout the applications’s life span. At each step, the key question for developers is: *how and how much did the performance change?* This question arises comparing an implementation to models or simulations; considering versions of an implementation that use a different algorithm, communication or numeric library, or language; studying code behavior by varying number or type of processors, type of network, type of processes, input data set or work load, or scheduling algorithm; and in benchmarking or regression testing. Despite the broad utility of this type of comparison, no existing performance tool provides the necessary functionality to answer it; even state of the art research tools such as Paradyn[2] and Pablo[3] focus instead on measuring the performance of a single program execution.

We describe an infrastructure for answering this question at all stages of the life of an application. We view each program run, simulation result, or program model as an *experiment*, and provide this functionality in an *Experiment Management* system. Our project has three parts: (1) a representation for the space of executions, (2) techniques for quantitatively and automatically comparing two or more executions, and (3) enhanced performance diagnosis abilities based on historic performance data. In this paper we present initial results on the first two parts. The measure of success of this project is that an activity that was complex and cumbersome to do manually, we can automate.

The first part is a concise representation for the set of executions collected over the life of an application. We store information about each experiment in a *Program Event*, which enumerates the components of the code executed and the execution environment, and stores the performance data collected. The possible combinations of code and execution environment form the multi-dimensional *Program Space*, with one dimension for each axis of variation and one point for each Program Event. We enable exploration of this space with a simple naming mechanism, a selection and query facility, and a set of interactive visualizations. Queries on a Program Space may be made both on the contents of the performance data and on the metadata that describes the multi-dimensional program space. A graphical representation of the Program Space serves as the user interface to the Experiment Management system.

The second part of the project is to develop techniques for automating comparison between experiments. Performance tuning across multiple executions must answer the deceptively simple question: what changed in this run of the program? We have developed techniques for determining the “difference” between two or more program runs,



automatically describing both the structural differences (differences in program execution structure and resources used), and the performance variation (how were the resources used and how did this change from one run to the next). We can apply our technique to compare an actual execution with a predicted or desired performance measure for the application, and to compare distinct time intervals of a single program execution. Uses for this include performance tuning efforts, automated scalability studies, resource allocation for metacomputing [4], performance model validation studies, and dynamic execution models where processes are created, destroyed, migrated [5], communication patterns and use of distributed shared memory may be optimized [6,9], or data values or code may be changed by steering [7,8]. The difference information is not necessarily a simple measure such as total execution time, but may be a more complex measure derived from details of the program structure, an analytical performance prediction, an actual previous execution of the code, a set of performance thresholds that the application is required to meet or exceed, or an incomplete set of data from selected intervals of an execution.

The third part of this research is to investigate the use of the predicted, summary, and historical data contained in the Program Events and Program Space for performance diagnosis. We are exploring novel opportunities for exploiting this collection of data to focus data gathering and analysis efforts to the critical sections of a large application, and for isolating spurious effects from interesting performance variations. Details of this are outside of the scope of this paper.

2 PROGRAM EVENTS AND AUTOMATED COMPARISON

The *Program Event* is a representation of a program run or execution. Associated with each Program Event is a collection of performance data. We can automatically compare the changes between two (or more) Program Events with our structural difference operator.

2.1 The Program Event

We represent a program as a collection of discrete program resources. Possible resources include the program code, application processes, machine nodes, synchronization points, data structures, and files. Each group of resources provides a unique view of the application. We organize the program resources into classes according to the aspect of the application they represent, and structure each class as a tree, called a *resource hierarchy*.

Definition 2.1. A *Program Event*, E , is a forest composed of zero or more unique resource hierarchies:

$$E = R_0 \cup R_1 \cup \dots \cup R_n, \quad n \geq 0. \quad \square$$

A resource hierarchy is a collection of related program resources. The root node of each resource hierarchy represents the entire program execution and is labeled with the name of the entire resource hierarchy. Each descendant of the root node represents a particular program resource within that view. As we move down from the root node, each level of the hierarchy represents a finer-grained description of the program. For example, a code hierarchy might have one level for nodes that represent modules, below that a level with function nodes, and below that a level for loops or basic block nodes.

In the Code resource hierarchy of Figure 2, the root level (level 0) is the program-level view of an application, which represents the cumulative behavior of the whole program. Level 1 is the module-level view of the source code,

and level 2 is its function-level view. Each level of a resource hierarchy is a set of resources, and each level above the leaf level is a partition of the set of nodes in the next lower level. For example, the module level of Figure 2, which contains `testutil.C`, `main.C`, and `vect.C`, is a partition of the set of all leaf nodes. We specify a particular level of a hierarchy using a superscript notation: R_0^1 refers to level one (the children of the root node) of resource hierarchy zero. The Code hierarchy in Figure 2 shows the set of functions partitioned into modules:

$$\begin{aligned} R_0^0 &= \{\text{Code}\} & R_0^1 &= \{\text{testutil.C}, \text{main.C}, \text{vect.C}\} \\ R_0^2 &= \{\text{printstatus}, \text{verifyA}, \text{verifyB}, \text{main}, \text{vect}::\text{addel}, \text{vect}::\text{findel}, \text{vect}::\text{print}\}. \end{aligned}$$

Definition 2.2. A resource hierarchy R is a tree of the form

$$R = (r, T)$$

where r is a resource and T is the set of all children of r in the resource hierarchy:

$$T = \{(r_0, T_0), (r_1, T_1), \dots, (r_m, T_m)\}. \quad \square$$

Figure 2 shows a sample Program Event for a parallel application called *Tester*. The Program Event is the set $\{\text{Code}, \text{Machine}, \text{Process}\}$. The Code hierarchy contains nodes that represent the program's modules and functions; the Machine hierarchy contains one node for each CPU on which *Tester* executed; and the Process hierarchy contains one node for each process.

A resource is a representation of a logical or physical component of a program execution. A single resource might be used to represent a particular aspect of the program or the environment in which it executes: a process, a function, a CPU, or a variable. In Figure 2, the leaf node labeled “verifyA” represents the resource $\langle \text{Code/testutil.C/verifyA} \rangle$. The semantic meaning attached to a particular resource is relevant only to the programmer and does not affect the model functionality; however, a program execution component must be uniquely represented by a single resource. Each internal node of a resource hierarchy tree represents a set of one or more resources; for example, $\langle \text{Code/testutil.C} \rangle$ is a single resource that represents the aggregation of the set $\{\text{printstatus}, \text{verifyA}, \text{verifyB}\}$.

A *resource name* is formed by concatenating the labels along the unique path within the resource hierarchy from the root to the node representing the resource. For example, the resource name that represents function `verifyA` (shaded) in Figure 2 is $\langle \text{Code/testutil.C/verifyA} \rangle$.

2.2 Representing Performance Information in a Program Event

For a particular performance measurement, we may wish to specify certain part or parts of a program. For example, we may be interested in measuring CPU time as the average for all executions, as the total for one entire execution, or as the total for a single function. The *focus* constrains our view of the program to a selected part. Selecting the root node of a hierarchy represents the unconstrained view, the whole program. Selecting any other node narrows the view to include only those leaf nodes that are descendents of the selected node. For example, the shaded nodes in Figure 2 represent the constraint: functions `verifyA` and `verifyB` of process *Tester:2* running on any CPU.

Definition 2.3. A *focus* F is formed by selecting a subset of resource nodes from each of the resource hierarchies R :

$$\{R_1^{j_1}, R_2^{j_2}, R_3^{j_3}, \dots, R_n^{j_n}\}$$

where jx is a selection from level x of resource hierarchy R_x and n is the total number of resource hierarchies. All of the nodes selected in a hierarchy must be from the same level. \square

The Program Event provides an intuitive naming scheme for program components called *resource normal form*. We convert the selected set to *resource normal form* by concatenating the selections from each resource hierarchy. For example, the shaded selection of Figure 2 is represented as:

`< /Code/testutil.C/(verifyA,verifyB), /Machine, /Process >`.

Performance data is represented as a collection of tuples of the form $P[m, f, t, r]$:

- m : the name of the *metric*, which is a measurable execution characteristic, such as CPU time.
- f : the *focus* for this performance data, for example CPU time for focus `< /Code/Main.C, /Machine, /Process >`,
- t : the *time interval*, specifies when during an execution the data was collected, and
- r : a *performance result* (r), which may be a simple scalar or more complex object.

Because data is uniquely identified using a focus, we say the performance data is stored in *resource normal form*.

2.3 The Structural Difference Operator

When comparing the performance of two program executions, a natural first question is, how different were the code and environments used in the two tests, and where did they differ? If we perform two test runs using identical code running on identical, dedicated platforms, every resource hierarchy of the first execution has an identical counterpart in the second execution; performance data may be meaningfully compared for every focus which is valid for one of the individual executions. The comparison becomes more complex if we consider cases in which either the code or the run time environment (or both) differ between our two test runs. We need to determine the common set of valid resources to determine the new set of valid foci for the program.

Given two Program Events, \mathcal{E}_1 and \mathcal{E}_2 (see Figure 4), we can calculate their structural difference. To perform a structural comparison of two program executions, we compare the two sets of resource hierarchies in a top-down manner. The *Structural Difference Operator*, \oplus , takes two Program Events as operands and yields a Program Event Group (PEG). The result contains all resources from the original two Program Events, so the structural difference operation works as a hierarchical set union operation. This operator can be applied iteratively to build up a single set of resource hierarchies that characterizes any number of distinct program runs.

Figure 1 shows an application of the structural difference operator to two program events \mathcal{E}_1 and \mathcal{E}_2 . The top set of resource hierarchies describes event \mathcal{E}_1 , and the middle set describes \mathcal{E}_2 . The bottom of the figure shows the result. Resources common to both executions are outlined with mixed dash/dot; `<Code>` and `<Code/Module2>` are examples of such resources. Resources unique to \mathcal{E}_1 or \mathcal{E}_2 are outlined with dot or solid boxes respectively.

2.4 Moving Beyond the Single Execution Model

A collection of two or more Program Events is stored in a Program Event Group (PEG). The interface to the PEG is centered around a single set of resource hierarchies, constructed using the Structural Difference Operator, which represents the components of all of its Program Events. We developed two aggregation operators of particular use in viewing PEG performance data: the *list* aggregation operator, when applied to a hierarchy that contains more than one

selected node, yields a performance result which is a list of values rather than a single value; and the *cluster* aggregation operator groups the requested performance data by value – each group contains values that were within a specified Δ of each other (represented by the average of the group). We have defined several new metrics particular to the multi-execution PEG. The *Performance Distance* metric takes performance data from two executions and measures how much the performance differed between the two runs. The *Performance Difference* operator returns a list of all locations [f, t] (focus, time interval) for which the discrete distance metric yields an answer of true for a specified metric. The *Discrete Distance* metric is a binary function that indicates whether two performance results differ by more than a specified interval.

3 THE PROGRAM SPACE

The collection of Program Events is stored in a multi-dimensional Program Space. Users can navigate this space, visualize its form and contents, and make queries on the structure and contents. *A more complete description of these features is given in the full paper.*

4 CASE STUDIES

We have implemented a preliminary prototype of an experiment management performance tool. To test our design we have examined data collected while tuning several parallel programs. We provide examples of using our prototype to (1) compare implementations based on alternate communication libraries, (2) evaluate performance as a program evolves through versions, and (3) track data for a scalability analysis. In each case, the use of an experiment management system simplifies and speeds the programmer's task. *Complete study details appear in full paper.*

4.1 Comparing Alternate Implementations: Porting a PVM Application to MPI

In this study we compared two versions of a parallel message-passing FFT code called ns, which was ported from the PVM to the MPI message passing libraries. The application solves the Navier-Stokes equation in three dimensions. A scientist porting an application wants directed feedback about the resulting changes in performance to the application, and hopefully some idea of the cause of any performance degradation. In Figure 3 we show the resource hierarchies that resulted from applying the Structural Difference Operator to two Program Events, a 4-node run of nspvm (PVM version) and a 4-node run of nsmpif (MPI version), both on the IBM SP-2. This display provides a quick way to see what differed in both the code and the environment between the two runs. The left in Figure 3 shows that Code, Machine, Memory, Process, and SyncObject resources appear in both Program Events. In the Code hierarchy, three modules (dfft.c, ns3d.c, and p3d.c) appear in both runs. At the leaf, we can see four procedures (strip, and strip1-3) appeared only in Program Event 1. The Message hierarchy on the right shows the changes in message tags: tags 0_1, 0_3, 0_5, and 0_2 represent MPI message tags, and the rest represent the message tags for the PVM version. By selecting any of these resources, we can display the performance data for the Program Event(s) that include that resource.

4.2 An Evolving Application: Performance Tuning a Shared Memory Application

A protein-folding application, called Fold4, was recently developed in our Chemical Engineering Department. The tuning effort was reported in detail elsewhere [10]. We analyzed the program versions and performance data from

this study and were able to automate the identification of changes from one version to the next.

We ran three versions of Fold4, taken from different steps in the performance tuning study. The researchers conducting the study had located several problems. Version 1 was a simple port to the Wisconsin Cluster of Workstations (COW) from the SGI PowerChallenge. A problem was identified with a serial portion of the code that consumed 40% of the execution time on 8 nodes. Version 2 changed to data partitioning to try to relieve the bottleneck. A problem was identified with false sharing of data blocks. Version 3 padded and aligned data to improve the cache behavior. We present selected results here to demonstrate the benefit of the experiment management approach in navigating the large space of resources and data involved in a complete performance tuning study.

To consider the changes from version 1 to version 2, we merged these two Program Events (by applying the Structural Difference Operator) and generated a Program Event Group (PEG). The PEG display distinguishes foci valid for both Program Events. *Display omitted; see Figure 3.* The Memory hierarchy shows that a data change of some kind occurred, since there are more data structures in the memory hierarchy for Version 2. Next we applied the performance difference operator with metric memoryBlockingTime. The perfDiff operator pairwise compares performance data available for both Program Events. The results are presented in the perfDiff display (shown in Figure 5). This display shows that memory blocking behavior differed overall in the two runs; further it differed in each process, and that those differences were localized to five data structures (GM, GM->part0-3).

4.3 A Scalability Study

In this study we compared 4 executions of a parallel global circulation message-passing code called Ocean. We ran the same code on 2, 4, 8, and 16 nodes of the Wisconsin COW, and used the merge, cluster, and performance difference features to explore the behavior of the application while scaling.

5 RELATED WORK

There has been extensive work on parallel profiling tools [11,12,3,2,13], and a more limited amount of work on tools that allow displays of performance data from multiple runs [14]. In the physical and life sciences, there has been work to automate and store experimental data [15]. In the debugging area, work has been done to find bugs by comparing the execution of a new program version to an old one by comparing program output and user-selected variables [16]. Our representation of a Program Space and automated comparison of program runs represents a new direction. *A complete discussion of this related work is included in the full paper.*

6 CONCLUSIONS

Taking an experiment management view of performance tuning allows us to describe an application's behavior throughout its lifetime and across a variety of environments and code ports. Our multi-execution view provides much valuable feedback for the scientist developing or maintaining a parallel code. We are currently working to expand the prototype described here, with a richer set of data types, a web-based interface to service developers engaged in geographically distributed cooperative program development and use, and additional performance displays. We are investigating the use of an object oriented database for data storage, and exploring techniques to improve Paradyn's on-the-fly diagnosis capabilities using feedback from our experiment management system.

7 REFERENCES

- [1] D. L. Wheeler. Top population geneticist delivers scathing critique of field. *The Chronicle of Higher Education*, February 14, 1997.
- [2] B. P. Miller, M. D. Callaghan, J. M. Cargille, J. K. Hollingsworth, R. B. Irvin, K. L. Karavanic, K. Kunchithapadam, and T. Newhall. The Paradyn parallel performance measurement tool. *IEEE Computer*, 28(11):37–46, November 1995.
- [3] D. Reed *et al.* Scalable performance analysis: The Pablo performance analysis environment. I. C. Press, editor, *Proc. Scalable Parallel Libraries Conference*, pages 135–142, Los Alamitos CA, 1993.
- [4] I. Foster and C. Kesselman. Globus: A metacomputing infrastructure toolkit. *International Journal of Super-computer Applications*, 1997. to appear.
- [5] J. Dongarra, G. Geist, R. Manchek, and V. Sunderam. Integrated PVM framework supports heterogeneous network computing. *Computers in Physics*, 7:166–74, 1993.
- [6] S. K. Reinhardt, J. R. Larus, and D. A. Wood. Tempest and typhoon: User-level shared memory. *21st Annual International Symposium on Computer Architecture*, April 1994.
- [7] W. Gu, G. Eisenhauer, E. Kraemer, J. Stasko, J. Vetter, and K. Schwan. Falcon: On-line monitoring and steering of large-scale parallel programs. *Symposium on the Frontiers of Massively Parallel Computation*, McLean, Virginia, February 1995.
- [8] K. Kunchithapadam and B. P. Miller. Integrating a Debugger and a Performance Tool for Steering, in **Debugging and Performance Tools for Parallel Computing Systems**. IEEE Computer Society Press, 1996. ed. by M.L Simmons, A.H. Hayes, J.S. Brown, and D.A. Reed.
- [9] S. Fink, S. Kohn, and S. Baden. Flexible communication mechanisms for dynamic structured applications. *IR-REGULAR '96*, Santa Barbara, CA, August 1996.
- [10] Z. Xu, J. R. Larus, and B. P. Miller. Shared-memory performance profiling. *6th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, Las Vegas, Nevada, June 1997.
- [11] J. Yan, S. Sarukhai, and P. Mehra. Performance measurement, visualization and modeling of parallel and distributed programs using the AIMS toolkit. *Software – Practice and Experience (SPE)*, 25(4):429–461, April 1995.
- [12] W. Williams, T. Hoel, and D. Pase. The MPP Apprentice performance tool: Delivering the performance of the Cray T3D. In *Programming Environments for Massively Parallel Distributed Systems*, Monte Verita, 1994.
- [13] M. Heath and J. Etheridge. Visualizing the performance of parallel programs. *IEEE Software*, 8(5):29–39, September 1991.
- [14] K. L. Karavanic, J. Myllymaki, M. Livny, and B. P. Miller. Integrated visualization of parallel program performance data. *Parallel Computing*, to appear.
- [15] Y. Ioannidis and M. Livny. Conceptual schemas: Multi-faceted tools for desktop scientific experiment management. *International Journal of Intelligent and Cooperative Information Systems*, 1(3):451–474, December 1992.
- [16] R. Sasic and D. Abramson. Guard: A relative debugger. *Software Practice and Experience*, to appear.
- [17] J. Kohn and W. Williams. ATExpert. *Journal of Parallel and Distributed Computing*, 18:205–222, 1993.
- [18] J. K. Hollingsworth and B. P. Miller. Dynamic control of performance monitoring on large scale parallel systems. *International Conference on Supercomputing*, Tokyo, July 1993.

8 FIGURES

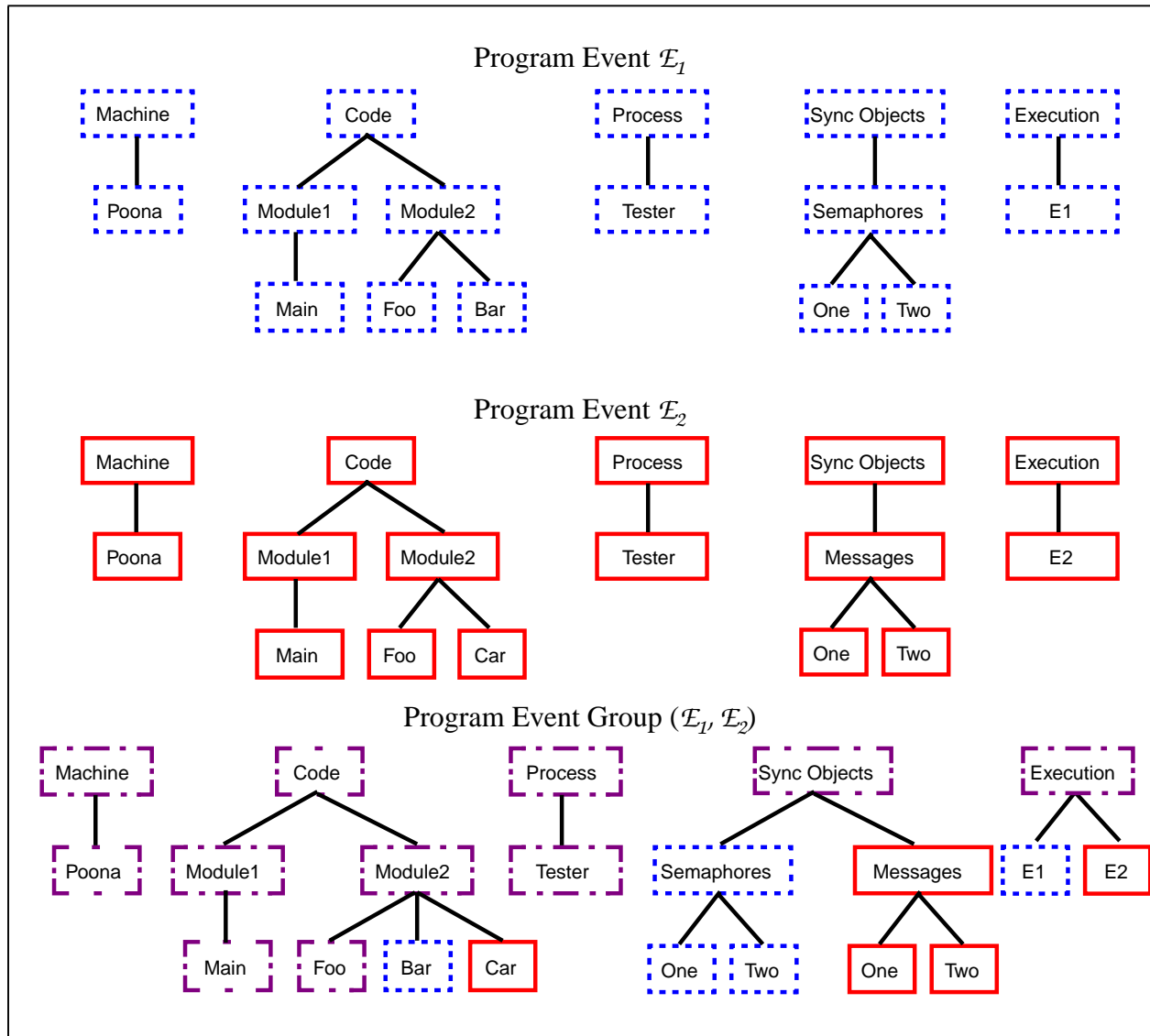


Figure 1: An Example of the Structural Difference Operator

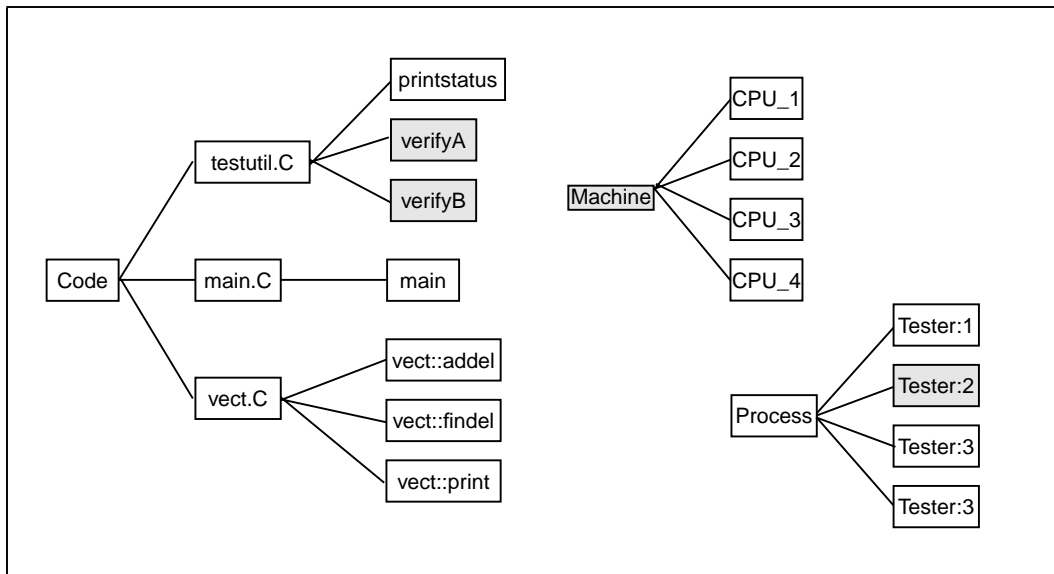


Figure 2: Program Event for Program *Tester*, with Three Resource Hierarchies

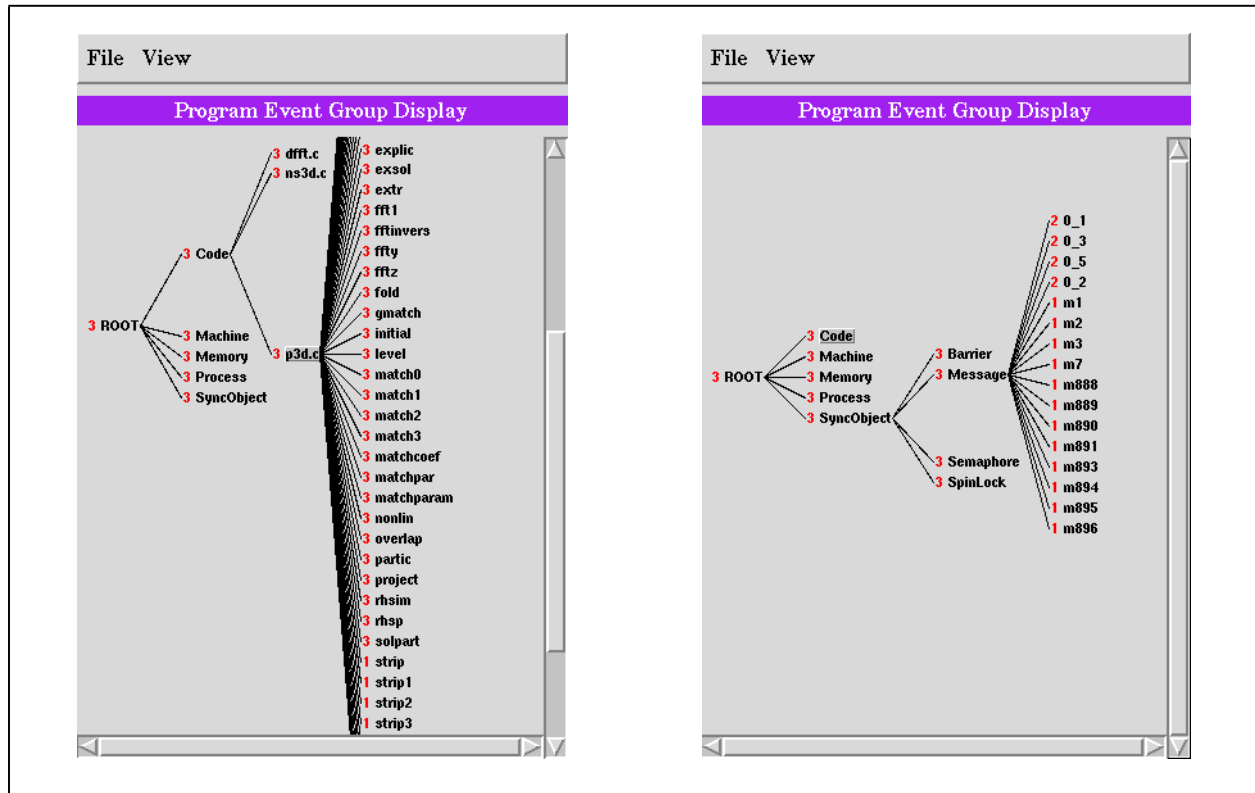


Figure 3: Resource Hierarchies for Program Event Group: *nspvm*, *nsmipif*. The PEG Display allows the developer to navigate resource hierarchies and quickly see what differed structurally between two (or more) program runs. The display is organized like the resource hierarchies from Figure 2, with the addition of an integer run identifier (RID). Each run is labeled with a value 1, 2, 4, 8, etc. Resources that appeared in only one run, are labeled with the RID of that run (1 or 2 in this example). Resources that appear in more run are labeled with the sum of the RIDs; the resources labeled with 3 appeared both runs 1 and 2.

```

[1]  $S \leftarrow \{ \}$ 
[2]  $\forall (r_i, T_i) \in S_1$ 
[3]   if  $\exists (r_j, T_j) \in S_2$ , s.t.  $\text{match}(r_i, r_j)$  then
[4]      $S \leftarrow S \cup \{ r_i \oplus r_j, T_i \oplus T_j \}$ 
[5]      $S_2 \leftarrow S_2 - (r_j, T_j)$ 
[6]   else  $S \leftarrow S \cup \{ (r_i, T_i) \}$ 
[7]  $S \leftarrow S \cup S_2$ 

```

Figure 4: Algorithm to find the Structural Difference of two Program Events, $S_1 \oplus S_2$

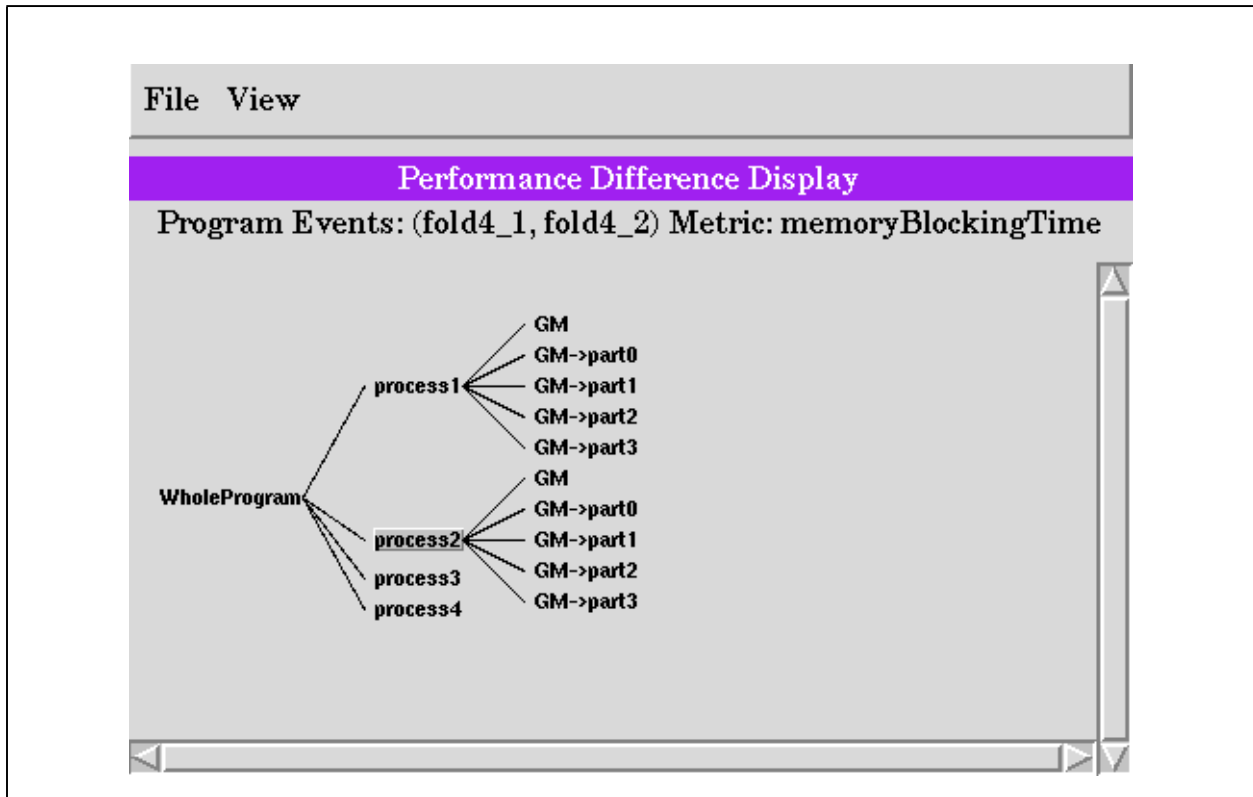


Figure 5: Results of the Performance Difference Operator for metric **MemoryBlockingTime.** The nodes shown represent resource combinations for which there is a performance difference detected between the Program Events. Starting at the left, the node **WholeProgram** means that there is some performance change between the two runs. The process nodes indicate that performance changed in some way for each of the four processes. The next level details individual shared data structures: **GM** is a shared index structure, and **GM->part0**, **GM->part1**, **GM->part2**, and **GM->part3** are the shared data structures listed in the index. The data structures listed are those common to both runs for which **memoryBlockingTime** changed. This snapshot was taken after selecting two nodes, **process1** and **process2**, to see a detailed display of their children. A visualization of the performance data for any node, showing the plots of the values for each run, can be launched by selecting the node on this display.