

1
2
3 **An Environment for Conducting**
4 **Families of Software Engineering**
5 **Experiments**
6

7
8
9 LORIN HOCHSTEIN

10 *University of Nebraska*
11

12
13 TAIGA NAKAMURA

14 *University of Maryland*
15

16
17 FORREST SHULL

18 *Fraunhofer Center Maryland*
19

20
21 NICO ZAZWORKA

22 *University of Maryland*
23

24
25 VICTOR R. BASILI

26 *University of Maryland, Fraunhofer Center, Maryland*
27

28
29 MARVIN V. ZELKOWITZ

30 *University of Maryland, Fraunhofer Center, Maryland*
31

32
33
34
35 **Abstract**

36 The classroom is a valuable resource for conducting software engineering
37 experiments. However, coordinating a family of experiments in classroom
38 environments presents a number of challenges to researchers. Understanding
39 how to run such experiments, developing procedures to collect accurate data,
40 and collecting data that is consistent across multiple studies are major problems.

1	This paper describes an environment, the Experiment Manager that simplifies	1
2	the process of collecting, managing, and sanitizing data from classroom experi-	2
3	ments, while minimizing disruption to natural subject behavior. We have success-	3
4	fully used this environment to study the impact of parallel programming languages	4
5	in the high-performance computing domain on programmer productivity at	5
6	multiple universities across the United States.	6
7		7
8		8
9	1. Introduction	176
10	1.1. Collecting Accurate Data	177
11	1.2. Classroom Studies	178
12	2. Classroom as Software Engineering Lab	179
13	3. The Experiment Manager Framework	182
14	3.1. Instrumentation Package	183
15	3.2. Experiment Manager Roles	184
16	3.3. Data Collection	186
17	4. Current Status	190
18	4.1. Experiment Manager Effectiveness	191
19	4.2. Experiment Manager Evolution	192
20	4.3. Supported Analyses	193
21	4.4. Evaluation	197
22	5. Related Work	197
23	6. Conclusions	198
24		24
25		25
26		26
27		27

1. Introduction

Scientific research advances by the creation of new theories and methods, followed by an experimental paradigm to either validate those theories and methods or to offer alternative models, which may be more appropriate and accurate. Computer science is not different from other sciences, and the field has been moving to adopt such experimental approaches [25]. However, in much of computer science, and in software engineering in particular, the experimental model poses difficulties possibly unique among the sciences. Software engineering is concerned about the appropriate models applicable to the development of large software systems. As such it involves the study of numerous programmers and other professionals over long periods of time. Thus, much of this research involves human behavior and in many ways is similar to research in psychology or the social sciences.

1 The major experimental approach accepted in scientific research is the replicated 1
2 study. However, by being expensive to produce, software is not amenable to such 2
3 studies. While a typical medical clinical trial may involve hundreds of subjects 3
4 testing a drug or a new treatment, even one duplication of a software development is 4
5 beyond the resources of most organizations. Although this approach is commonly 5
6 used in various fields of research involving humans, such as clinical study in 6
7 medical science, conducting many studies is still difficult and expensive, which is 7
8 often a major obstacle for good software engineering research. 8

9 The problems with empirical studies in software engineering can be classified by 9
10 two major problems: cost of such studies and accuracy of the data. 10
11

12 1.1 Collecting Accurate Data 12

13 1.1.1 *Costs of Software Engineering Studies* 13

14 15 Because of the cost of developing a piece of software, typically case studies of a 15
16 single development are monitored, and after many such studies, general trends can 16
17 be observed. As an example of this, we will look at the NASA Goddard Space Flight 17
18 Center (GSFC) Software Engineering Laboratory (SEL), which from 1976 to 2001 18
19 conducted many such studies [6]. The experiences of the SEL are illustrative of the 19
20 problems encountered in data collection. The data was collected, beginning in 1976, 20
21 at GSFC from NASA developers and the main development contractor, Computer 21
22 Sciences Corporation (CSC). The data was then manually reviewed at GSFC before 22
23 being sent to the University of Maryland for entry into the project measures database 23
24 using a UNIX-based Ingres system. 24

25 The naïve simplicity in which data was collected broke down by 1978 and a more 25
26 rigorous set of processes was instituted. This could not be a part-time activity by 26
27 faculty using undergraduate employees. In addition, the university researchers 27
28 wanted a considerable amount of data, and soon realized that the GSFC program- 28
29 ming staff did not have the time to comply with their requests. They had to 29
30 compromise on the amount of data desired versus the amount of data that realisti- 30
31 cally could be collected. Data, which was collected on forms filled out by the 31
32 programming staff, were shortened to allow for more complete collection. 32

33 The data collection process for the 20 projects then under study became more 33
34 rigorous with this five-step approach: 34
35

- 36 1. Programmers and managers completed forms. 36
- 37 2. Forms were initially verified at CSC. 37
- 38 3. Forms were encoded for entry at GSFC. 38
- 39 4. Encoded data checked by validation program at GSFC. 39
- 40 5. Encoded data revalidated and entered into database at University (after several 40
years, CSC took over total management of the database).

1 But, to obtain contractor cooperation, a 10% overhead cost to projects was 1
2 allocated for data collection and processing activities. Eventually the overhead 2
3 cost of collecting data was reduced, but the total cost of collecting, processing, 3
4 and analyzing data continued to remain between 5% and 10%. However, on a \$500K 4
5 project, this still amounted to almost \$50K just for data collection – an amount that 5
6 few organizations are willing to invest. While the SEL believed that the payoff in 6
7 improved software development justified this cost, it is beyond the scope of this 7
8 chapter to prove that. Suffice it to say that most organizations consider the additional 8
9 costs of data collection as unjustified expenses. 9
10

11 **1.1.2 Accuracy in Collected Data** 11

12
13 Most data collected on software development projects can be generally classified 13
14 as self-reported data – the technical staff fill out effort reports on hours worked, 14
15 change reports on defects found and fixed, etc. The care in which such data is 15
16 reported and collected greatly affects the accuracy of the process. Unfortunately, the 16
17 process is not very accurate. Self-reported measures can vary over time, due to 17
18 history or maturation effects [7], and the accuracy of such measures varies across 18
19 individuals. This is a particular problem when the subjects have more interest in 19
20 completing the task than complying with the protocols of the study. Basili et al. [2] 20
21 evaluated Software Science metrics against self-reported effort data collected from 21
22 GSFC software projects. There was very little correlation between this effort and 22
23 metrics known to predict effort, and there was concern that poor self-reported data 23
24 was distorting the results. Perry et al. [18] analyzed previous data from project 24
25 notebooks and free-form programmer diaries which were originally kept for per- 25
26 sonal use. They found that the free-form diaries were too inconsistent across subjects 26
27 and sometimes lacked sufficient resolution. 27
28

29 **1.2 Classroom Studies** 29

30
31 The overhead in collecting accurate detailed data at the SEL was too high to 31
32 maintain a data collection process. This same result has been found in other data 32
33 collection studies. Instead of large replications, running studies in classrooms using 33
34 students as subjects have become a favorite approach for trying out new develop- 34
35 ment techniques [20]. Even though a conclusion drawn from student subjects cannot 35
36 always be generalized to other environments, such experiments aid in developing 36
37 approaches usable elsewhere. However, conducting larger-scale software engineer- 37
38 ing research in classroom environments can be quite complex and time consuming 38
39 without proper tool support. Proper tool support is a requirement if we want to 39
40 improve on the poor quality of self-reported data. 40

1 A single environment (e.g., a single class) is often insufficient for obtaining 1
2 significant results. Therefore, multiple replications of a given experiment must be 2
3 carried out by different faculty at different universities in order to provide sufficient 3
4 data to make conclusions. This means, each experiment must be handled in a similar 4
5 manner to allow for combining partial results. The complexity of providing consistent 5
6 data across various experimental protocols has been overwhelming so far. 6

7 In this chapter, we describe an environment we are developing to simplify the 7
8 process of conducting software engineering experiments that involve development 8
9 effort and workflow and ensure consistency in data collection across experiments in 9
10 classroom environments. We have used this environment to carry out research to 10
11 identify the effect of parallel programming languages on novice programmer pro- 11
12 ductivity in the domain of high-performance computing (e.g., MPI [12], OpenMP 12
13 [11], UPC [8], and Matlab*P [10]). Although there are often issues regarding the 13
14 external validity of students as subjects, this is not a major concern here since we are 14
15 explicitly interested in studying student programmers. 15

16 This work was carried out in multiple universities across the United States in 16
17 courses where the professors were not software engineering researchers and were, 17
18 therefore, not experienced with conducting experiments that involved human subjects. 18
19

20 **2. Classroom as Software Engineering Lab** 20

21
22 The classroom is an appealing environment for conducting software engineering 22
23 experiments, for several reasons: 23
24

- 25 • Most researchers are located at universities. Being close to your subjects is 25
26 often necessary to obtain accurate results. 26
- 27 • Training can be integrated into the course. No extra effort is then required by 27
28 the subjects since there is the assumption that the training is a valuable 28
29 academic addition to the classroom syllabus. 29
- 30 • Required tasks can be integrated into the course. 30
- 31 • All subjects are performing identical programming tasks, which are not gener- 31
32 ally true in industry. This provides an easy source for replicated experiments. 32
33

34 In addition to the results that are obtained directly by these studies, such experi- 34
35 ments are also useful for piloting experimental designs and protocols which can later 35
36 be applied to industry subjects, an approach which has been used successfully 36
37 elsewhere (e.g., [3, 4, 5]). 37

38 While there are threats to validity of such studies by using students as subjects as 38
39 proxies for professional programmers (e.g., the student environment may not be 39
40 representative of the ones faced by professional programmers), there are additional 40

1 complexities that are specific to research in this type of environment. We encountered each of these issues when conducting research on the effects of parallel programming model on effort in high-performance computing [14]:

- 2 1. *Complexity*: Conducting an experiment in a classroom environment is a complex process that requires many different activities (e.g., planning the experimental design, identifying appropriate artifacts and treatments, enrolling students, providing for data collection, checking for process compliance, sanitizing data for privacy, and analyzing data). Each such activity identifies multiple points of failure, thus requiring a large effort to organize and run multiple studies. If the study is done at multiple universities in collaboration with other professors, these professors may have no experience in organizing and conducting such experiments.
- 3 2. *Research versus pedagogy*: When the experiment is integrated into a course, the experimentalist must take care to balance research and pedagogy [9]. Studies must have minimal interference with the course. If the students in one class are divided up into treatment groups and the task is part of an assignment, then care must be taken to ensure that the assignment is of equivalent difficulty across groups. Students who consent to participate must not have any advantage or disadvantage over students who do not consent to participate, which limits additional overhead required by the experiment. In fact, each university's Institutional Review Board (IRB), required in all United State universities performing experiments with human subjects, insists that participation (or nonparticipation) must have no effect on the student's grade in the course.
- 4 3. *Consistent replication across classes*: To build empirical knowledge with confidence, researchers replicate studies in different environments. If studies are to be replicated in different classes, then care must be taken to ensure that the artifacts and data collection protocols are consistent. This can be quite challenging because professors have their own style of giving assignments. Common projects across multiple locations often differ in crucial ways making meta-analysis of the combined results impossible [16].
- 5 4. *Participation overhead for professors*: In our experience, many professors are quite willing to integrate software engineering studies into their classroom environment. However, for professors who are unfamiliar with experimental protocols, the more effort required of them to conduct a study, the less likely it will be a success. In addition, collaborating professors who are not empirical researchers may not have the resources or the inclination to monitor the quality of captured data to evaluate process conformance. Therefore, empirical researchers must try to minimize any additional effort required to run an empirical study in the course while ensuring that data is being captured correctly.

- 1 The required IRB approval, when attempted for the first time, seems like a 1
2 formidable task. Help in understanding IRB approval would greatly aid the 2
3 ability of conducting such research experiments. 3
- 4 5. *Participation overhead for students:* An advantage of integrating a study into a 4
5 classroom environment is that the students are already required to perform the 5
6 assigned task as part of the course, so the additional effort involved in 6
7 participating in the study is much lower than if subjects were recruited from 7
8 elsewhere. However, while the additional overhead is low, it is not zero. The 8
9 motivation to conform to the data collection process is, in general, much lower 9
10 than the motivation to perform the task, because process conformance cannot 10
11 be graded. In addition, the study should not subvert the educational goals of 11
12 the course. Putting the experiment in the context of the course syllabus is 12
13 never easy. 13
- 14 This can be particularly problematic when trying to collect process data 14
15 from subjects (e.g., effort, activities, and defects), especially for assignments 15
16 that take several weeks. (e.g., we saw a reduction in process conformance over 16
17 time when subjects had to fill out effort logs over the course of multiple 17
18 assignments). 18
- 19 6. *Automatic data collection of software process:* To reduce subject overhead and 19
20 increase data accuracy, it is possible to collect data automatically from the 20
21 programmer's environment. Capturing data at the right level of granularity is 21
22 difficult. All user-generated events can be captured (keyboard and mouse 22
23 events), but this produces an enormous volume of data that may not abstract 23
24 to useful information. Allowing this raw data to be used can create privacy 24
25 issues, such as revealing account names, with the ability to then determine how 25
26 long specific users took to build a product or how many defects they made. 26
- 27 All development activities taking place within a particular development 27
28 environment (e.g., Eclipse) simplifies the task of data collection, and tools 28
29 exist to support such cases (e.g., Marmoset [21]). However, in many 29
30 domains, development will involve a wide range of tools and possibly 30
31 even multiple machines. For example, in the domain of high-performance 31
32 computing, preliminary programs may be compiled on a home PC, final 32
33 programs are developed on the university multiprocessor, and are ulti- 33
34 mately run on remote supercomputers at a distant datacenter. Programmers 34
35 typically use a wide variety of tools, including editors, compilers, build 35
36 tools, debuggers, profilers, job submission systems, and even web browsers 36
37 for viewing documentation. 37
- 38 7. *Data management:* Conducting multiple studies generates an enormous vol- 38
39 ume of heterogeneous data. Along with automatically collected data and 39
40 manually-reported data, additional data includes versions of the programs, 40

1 pre- and post-questionnaires, and various quality outcome measures (e.g., 1
2 grades, code performance, and defects). Because of privacy issues, and to 2
3 conform to IRB regulations, all data must be stored with appropriate access 3
4 controls, and any exported data must be appropriately sanitized. Managing this 4
5 data manually is labor-intensive and error-prone, especially when conducting 5
6 studies at multiple sites. 6

3. The Experiment Manager Framework

7
8
9
10 We evolved the Experiment Manager framework (Fig. 1) to mitigate the complex- 10
11 ities described in the previous section. The framework is an integrated set of tools 11
12 to support software engineering experiments in HPC classroom environments. While 12
13 aspects of the framework have been studied by others, the integration of all features 13
14 allows for a uniform environment that has been used in over 25 classroom studies 14
15 over the past 4 years. The framework supports the following. 15
16
17

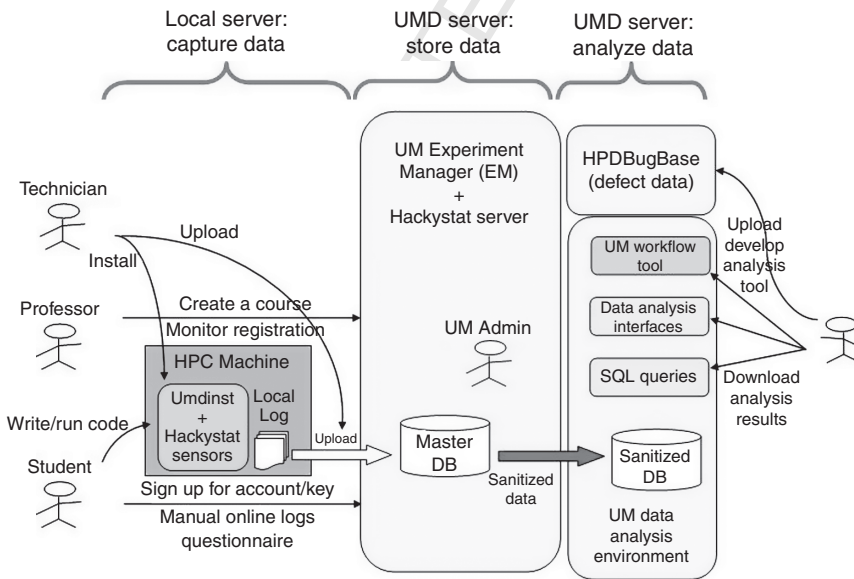


FIG. 1. Experiment Manager structure.

- 1 1. *Minimal disruption of the typical programming process*: Study participants solve 1
2 programming tasks under investigation using their typical work habits, spreading 2
3 out programming tasks over several days. The only additional activity required is 3
4 filling out some online forms. Since we do not require them to complete the task in 4
5 an alien environment or work for a fixed, uninterrupted length of time, we 5
6 minimize any negative impact on pedagogy or subject overhead. 6
- 7 2. *Consistent instruments and artifacts*: Use of the framework ensures that the 7
8 same type of data will be collected and the same type of problems will be 8
9 solved, which increases confidence in meta-analysis across studies at different 9
10 universities. 10
- 11 3. *Centralized data repository with web interface*: The framework provides a 11
12 simple, consistent interface to the experimental data for experimentalists, sub- 12
13 jects, and collaborating professors. This reduces overhead for all stakeholders 13
14 and ensures that data is consistently collected across studies. 14
- 15 4. *Sanitization of sensitive data*: The framework provides external researcher 15
16 with access to the data sets that have been stripped of any information that 16
17 could identify subjects, to preserve anonymity and comply with the protocols 17
18 of human subject research as set out by IRBs at American universities. 18
19
20
21
22

3.1 Instrumentation Package

23 Our instrumentation package, called *UMDinst*, supports automatic collection of 23
24 software process data in a Unix-based, command-line development environment, 24
25 which is commonly used in high-performance computing. The package is designed 25
26 to be installed in a master account on the local server and then enabled in the accounts 26
27 of each subject by executing a set up script, or be installed in the account of individual 27
28 subjects. The appropriate installation mode depends on the need of a specific experi- 28
29 ment and the configuration of the machine to be instrumented. In either case, the 29
30 package can be used without the intervention of system administrators. 30

31 *UMDinst* package instrument programs that are involved in the software devel- 31
32 opment process by replacing each command that invokes a tool (e.g., compiler) with 32
33 a script that first collects the desired information and then calls the original tool. It is 33
34 used for instrumenting compilers, although it is also designed to support job 34
35 schedulers (common in high-performance computing environments), debuggers, 35
36 and profilers. For each compile, the following data is captured: 36
37

- 38 • a time stamp when the command is executed 38
- 39 • contents of the source file that were compiled 39
- 40 • contents of local header files referenced in the source file 40

- 1 • the command used to invoke the compiler 1
- 2 • the return code of the compiler 2
- 3 • the time to compile 3

4
5 The UMDinst package includes Hackystat sensors [15] to instrument supported 5
6 editors such as Emacs and vi, and to capture shell commands and time stamps. The 6
7 collected data is used in studies to estimate total effort as well as to infer develop- 7
8 ment activities (e.g., debugging, parallelizing, and tuning). Hackystat is a system 8
9 developed by Johnson at the University of Hawaii that captures low-level events 9
10 from a set of instrumented tools. Thus, while UMDinst captures data at the com- 10
11 mand-line level, Hackystat captures time stamps and events from editors and related 11
12 tools that have been instrumented. The pair of tools provides a complete history of 12
13 user interaction in developing a program. 13

16 3.1.1 *Web Portal* 16

17
18 The heart of the Experiment Manager framework is the web portal, which serves 18
19 as a front-end to the database that contains all of the raw data, along with metadata 19
20 about individual experiments. Multiple stakeholders use the web interface: experi- 20
21 menters, subjects, and data analysts. For example, experimenters specify treatments 21
22 (in our case, parallel programming models), assignment problem, participation rate, 22
23 and grades. They also upload data captured automatically from UMDinst. Subjects 23
24 fill in questionnaires, and report on process data such as time worked on different 24
25 activities and defects. Analysts would export data of interest, such as total effort [13] 25
26 for hypothesis testing, or a stream of time stamped events for workflow modeling. 26

29 3.2 Experiment Manager Roles 29

30
31 We have divided the functionality of the Experiment Manager into four roles. For 31
32 each role, we developed several use cases that describe its functionality, thus 32
33 simplifying the design of the software. 33

- 34 1. *Technician*: The technician sets up the environment on the local server, usually 34
35 not at the University of Maryland. This will be someone at a university with 35
36 access to the machine the students use for the class. Often it is the Teaching 36
37 Assistant in the course the software will be used in. The tasks for the technician 37
38 are to install UMDinst so that students can use it. At the end of the semester, 38
39 the technician also sends the collected data to the University of Maryland 39
40 server in case it was collected locally. 40

1 2. *Professor*: A database provides the professor with sample IRB questionnaires 1
2 for submittal. This cannot be fully automated since each university has its own 2
3 guidelines for submitting the IRB approval form. But experience with many 3
4 universities over the last 4 years allows us to help in answering the most 4
5 common questions on these forms. 5

6 The instructor first registers each class with the Experiment Manager to set 6
7 up a classroom experiment. For each such class, the professor can assign 7
8 several programming projects from our collected database of assignments or 8
9 assign one of his own. During the semester, the system allows the professor to 9
10 see if students have completed their assignments, but does not allow access 10
11 to any of the collected data until the grades for the assignment are completed. 11
12 In reality, the Teaching Assistant may be the person to actually perform this 12
13 task, but conceptually is acting in the role of the professor. 13

14 3. *Student*: A student who takes part in the experiment provides data on HPC 14
15 development. This requires the student to: 15

16 1. Register with the Experiment Manager by filling out a background ques- 16
17 tionnaire on courses taken and experiences in computer science in general 17
18 and HPC programming in particular. Although this registration process can 18
19 take up to 15 min, it is required only once during the semester. 19

20 2. Run the script to set up the wrappers for the commands that edit, compile 20
21 and run programs. Once an assignment is underway, the data collection 21
22 process is mostly automatic and data is collected mostly painlessly. 22

23 4. *Analyst*: An analyst accesses the collected data for evaluating some hypothesis 23
24 about HPC development. At the present time, the analysis tools are relatively 24
25 simple. Analysts can see the total effort and defects made by each student and 25
26 collect workflow data. 26

27 Many tools exist in prototypes to support the various types of studies. For 27
28 example, the HPC community is developing concepts of what productivity means 28
29 in the HPC environment [26] and we have been looking at developing workflow 29
30 models (e.g., how much time is spent in various activities, such as developing code, 30
31 testing, parallelizing the code, and tuning the code for better performance). 31
32 To support the study of workflows and productivity, we have developed a tool 32
33 to allow the experimenter to apply various heuristics to the base data to see if we can 33
34 automatically deduce the developer's programming activity, for example, testing 34
35 and debugging versus development. This tool takes raw data, collected from online 35
36 tools such as Hackystat and manual logs generated by students, and we have been 36
37 developing algorithms for automatically inferring the workflow cycle [23]. Results 37
38 of this work are described in Section 4. Current activities are looking at extending 38
39 these tools. 39
40 40

3.3 Data Collection

The actual data collection activity was designed to present minimal complexity to the student (Figs. 2 and 3). Within the Experiment Manager, the student has two options. If data was not collected automatically, the student can enter a set of activities, with the times each activity started and ended (Fig. 3) (e.g., self-reported data, which we discussed earlier to be less reliable). However, the effort tool simplifies the process greatly. If the student clicks to start the tool (small oval near bottom of Fig. 3), then a small window opens on the top left corner of the screen (large oval in upper left in Fig. 3). Each time the student starts a different activity, the student only needs to pull down the menu in the effort tool and set the new activity type (Fig. 2). The time between clicks is recorded as the time of the previous activity. Thus, while the data is not totally automatic, we believe we have minimized the overhead of collecting such data.

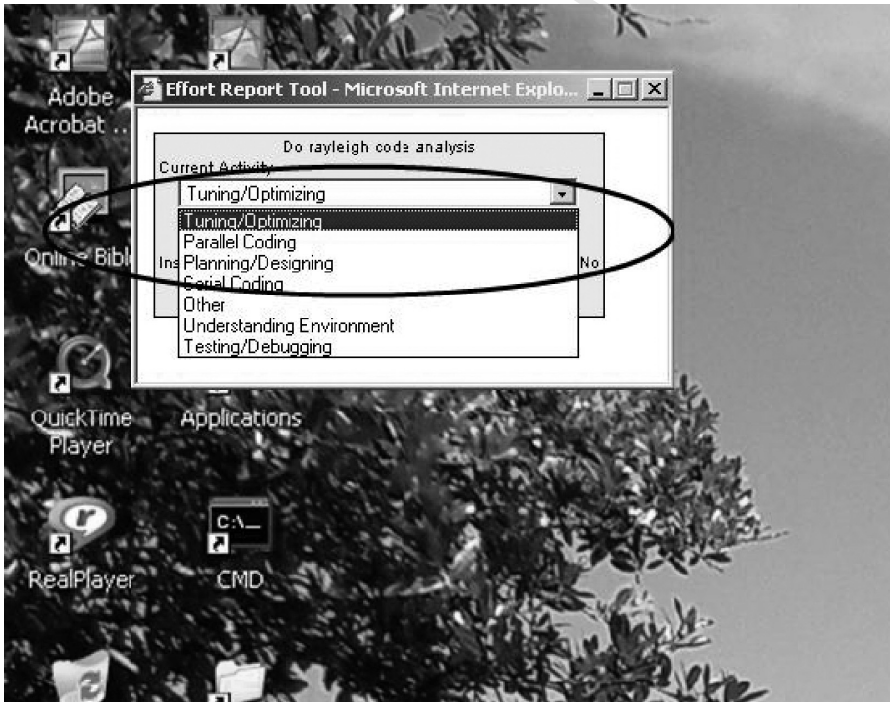


FIG. 2. Effort capture tool.

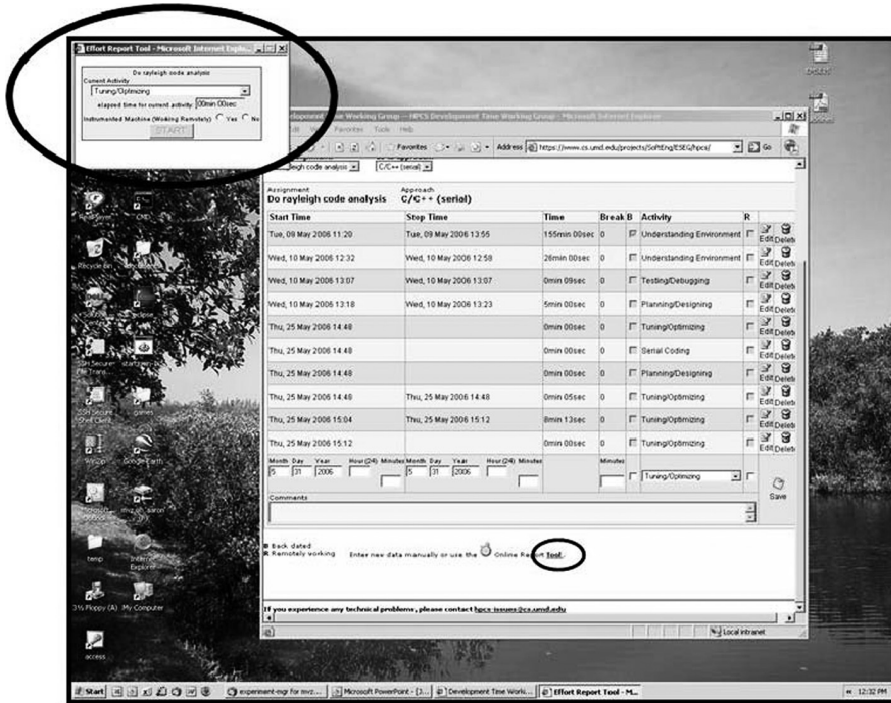


FIG. 3. Effort collection screen.

The tool automatically computes elapsed time between events and saves the data in the database. If the student stops for a period of time (e.g., goes to lunch and surfs the web), there is a *stop* button on the tool. Upon returning, the user simply clicks on *start* to resume timing.

For most HPC development, the student simply has to:

1. Log into Experiment Manager to go to effort page (Fig. 4), then click on effort tool (Fig. 3).
2. Develop program as usual.
3. Each time a new activity starts, click on the new activity in the effort tool (Fig. 2).
4. If any errors are found, the student records that defect by invoking the defect tool (Fig. 4) to explain the defect on a separate page (Fig. 5).

Only steps 3 and 4 involve any separate activity for participating in these experiments, and such activity is minimal.

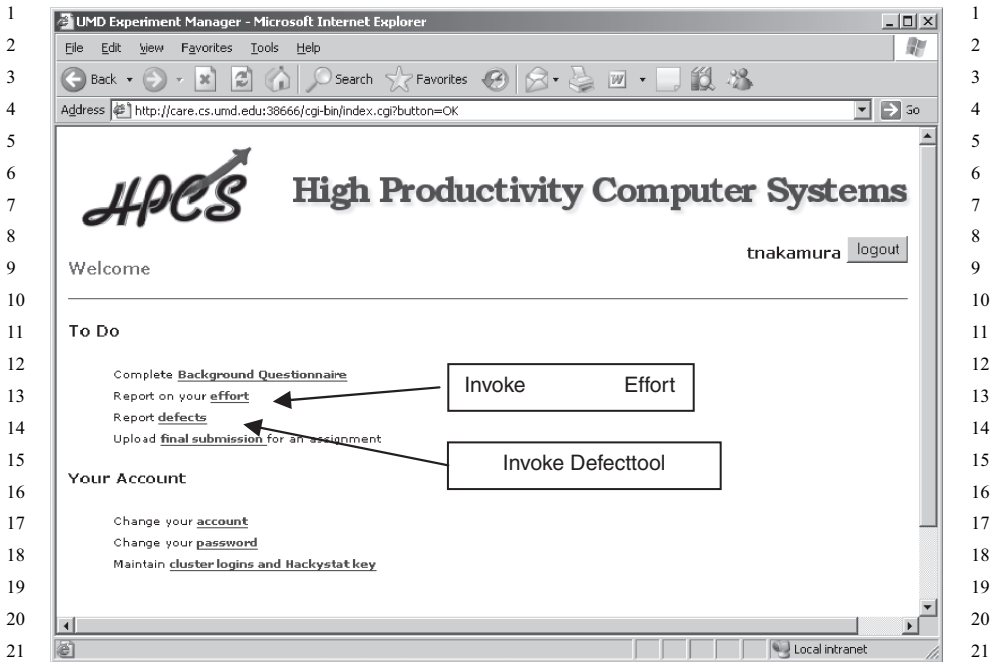


FIG. 4. Student view of Experiment Manager.

3.3.1 Data Sanitization

While personal data collected by the experiments must be kept private, we would like to provide as much data as possible to the community as part of our analysis. The sanitization process exports 'safe' data into a database that can be made accessible to other researchers, running on a separate machine.

The sanitization process is briefly described in Fig. 6. Each data object we obtained in an experiment is classified as one of:

1. *Prohibited*: Data contains personal data we cannot reveal (e.g., name or other personal identifiers).
2. *Clean*: Data we can reveal (e.g., effort data for development of another clean object).
3. *Modified*: Data we can modify to make it clean (e.g., removing all personal identification in the source program).

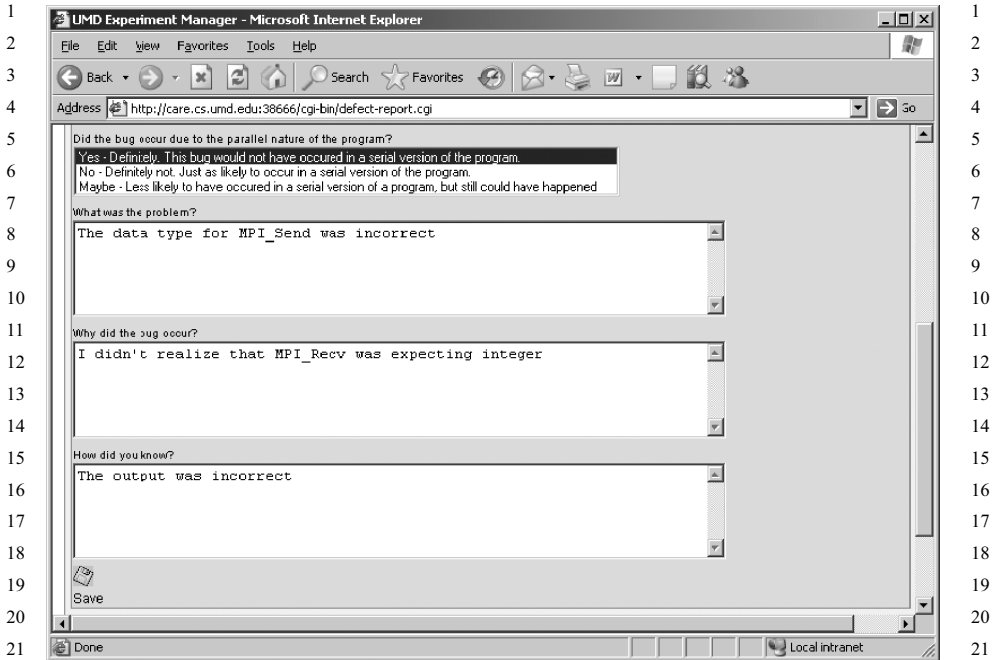


FIG. 5. Defect reporting tool.

Clean data can be moved to the analysis server and modified data can also be moved. Only prohibited data cannot be exported to others desiring to look at our collected database. Our sanitization process on data consists of the following four functions:

1. *Normalization* – Normalize the time stamps for each class on a common basis. By making each time stamp relative to 0 from the beginning of that experiment, information about in which semester it was collected (and hence from which school the data was collected) is hidden.
2. *Discretization* – Since grades are considered private data, we define a mapping table that maps grades on a small set, such as {good,bad}. Converting other interval or ratio data into less specific ordinal sets, while it loses granularity, it helps to present anonymity.
3. *Summarization* – With some of the universities, we can give out source code if we remove all personal identifiers of the students who wrote the code. But in

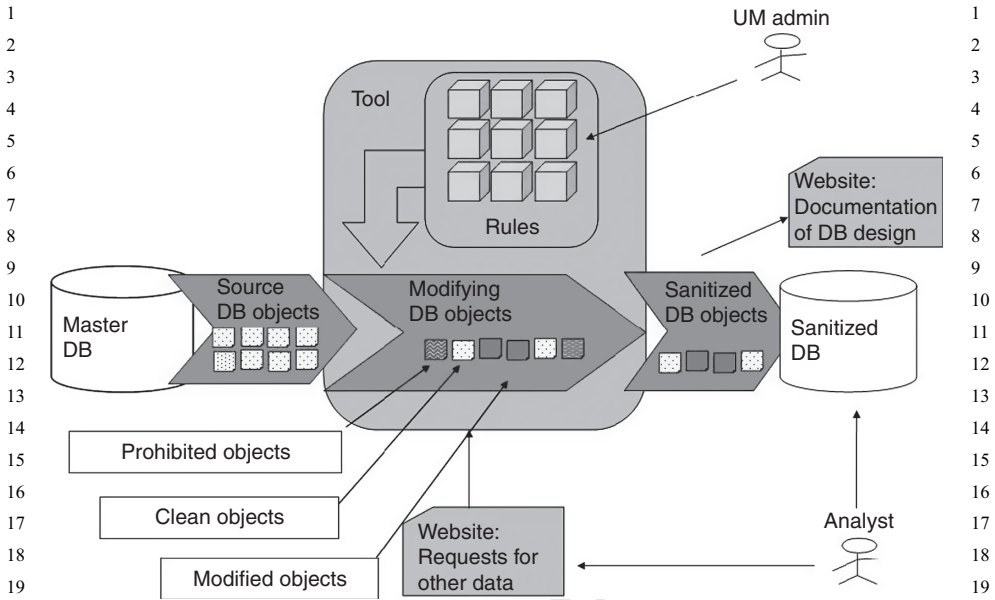


FIG. 6. Basic sanitization process.

some cases, we are prohibited from doing even that. If we cannot give out source code, we can collect patterns and counts of their occurrence in the source code. For example, we can count lines of code, or provide analyses of 'diffs' of successive versions of code.

4. *Anonymization* – We can create hash values for dates, school names, and other personal identifiers.

4. Current Status

Our Experiment Manager framework currently contains data from 25 classroom experiments conducted at various universities in the United States (Fig. 7). While some of the experiments preceded the Experiment Manager (and motivated its development) and their data was imported into the system; perhaps half the experiments used parts or all of the system.

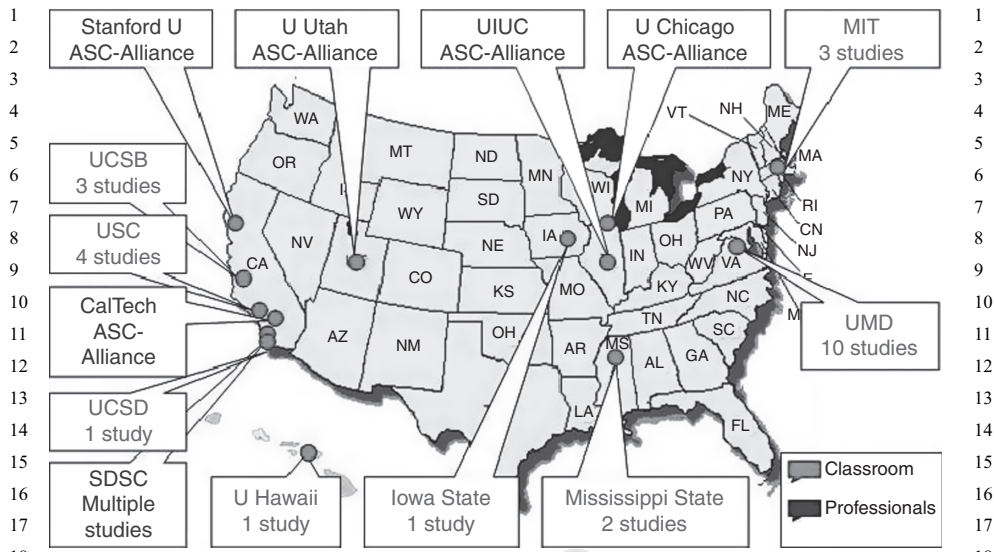


FIG. 7. Completed studies.

4.1 Experiment Manager Effectiveness

Before using the experiment manager, we discovered many discrepancies between successive classroom projects which prevented merging the results. Some of these were:

1. *It was not often obvious for how many processors the final programs were written.* Since a major goal of HPC programming is to divide an algorithm to run on multiple processors, this speedup (i.e., relative decrease in execution time by using multiple processors) is a critical measure of performance. Without knowing the initial goals for each student assignment, it was unclear how to measure performance goals for each class.
2. *Related to the previous problem, the projects all had differing requirements for final program complexity.* (e.g., the number of replicated cells needing to be computed). How big a grid (e.g., number of replicated cells) were required in which to compute an answer and measure performance? This affected student programming goals.
3. *Grading requirements differed.* Was performance on an HPC machine important? Sometimes just getting a valid solution mattered. Maximum speedup, or the decrease in execution time of an HPC machine over a serial implementation, was sometimes the major goal.

1 By using our collected database of potential assignments, as well as our checklist 1
2 of project attributes, this problem has lessened across multiple classes recently, 2
3 allowing for the combination of results across different universities. 3

4 The IRB process seems like a formidable roadblock the first time any new 4
5 professor encounters it. Often, in contacting faculty at a new university we would 5
6 lose a semester's activity simply because the IRB approval process was too onerous 6
7 the first time it was attempted. With our experience of IRB issues, and our collection 7
8 of IRB forms required by various universities, this no longer is a major problem. 8

9 A related problem was the installation of software on the host computer for the 9
10 collection of data. Again, this often meant the delay by a semester since the 10
11 installation process was too complex. This was a major driving force to host much 11
12 of this software as a web server at the University of Maryland, with a relatively 12
13 simple UMDinst package that needed to be installed at each university's site. 13

14 The effort tool (pictured earlier as Figs. 2 and 3) also solved some of our data 14
15 collection problems. We can collect effort data by three ways (Hackystat at the level 15
16 of editor and shell event time stamps, manual data via programmer filled-in forms, 16
17 and compiler time stamps via UMDinst). All give different results [13]. The use of 17
18 the effort tool greatly eases the data collection problem, which we believe increases 18
19 the reliability of such data. 19

20 Most of our results, so far, are anecdotal. But we have been able to address new 20
21 universities and additional classes in a more methodical manner at present and 21
22 believe the Experiment Manager software is a major part of this improvement. 22
23

24 4.2 Experiment Manager Evolution 24

25 The system is continuing to evolve. Current efforts focus on the following tasks: 25
26
27

- 28 • We are evolving the user interface to the Experiment Manager web-based tool. The 28
29 goal is to minimize the workload of various stakeholders (i.e., roles) for setting up 29
30 the experiment environment, registering with the system, entering the data, and 30
31 conducting an analysis. We would like to develop small native applications 31
32 that provide a more integrated interface to the operating system, making it less 32
33 disruptive to users. 33
34
- 35 • We want to continue our experimentation with organizations that are often behind 35
36 firewalls. Although we are currently studying professionals in an open environ- 36
37 ment, we want to use the Experiment Manager in this environment. Although the 37
38 UMDinst instrumentation package can be set up in secure environments, the 38
39 collected data cannot be directly uploaded to the University of Maryland servers. 39
40 We have planned extensions to the Experiment Manager architecture to better 40

- 1 support the experimentations with these organizations. Working in these 1
2 environments is necessary to see how professionals compare to the students. 2
- 3 • We will continue to evolve our analysis tools. For example, our prototype 3
4 experience bases for evolving hypotheses and high end computing defects (e.g., 4
5 www.hpccbugbase.org) will continue to evolve both in content and usability. 5
 - 6 • We will evolve problem-specific harnesses that automatically capture informa- 6
7 tion about correctness and performance of intermediate versions of the code 7
8 during development to ensure that the quality of the solutions (specifically, 8
9 correctness and performance) is measured consistently across all subjects. 9
- 10

11 This also requires us to evolve our experience bases to generate performance 11
12 measures for each program submitted in order to have a consistent performance and 12
13 speedup measure for use in our workflow and time to solution studies. 13

14 Our long-range goal is to allow the community access to our collected data. This 14
15 requires additional work on a sanitized database that removes personal indicators 15
16 and fulfills legal privacy requirements for use of such data. 16

17

18 4.3 Supported Analyses 18

19

20 The Experiment Manager was designed to ease data analysis, in order to support 20
21 the investigation of a range of different research questions. Some of these analyses 21
22 are focused in detail on a single developer being studied, while others aggregate data 22
23 over several classes, allowing us to look across experimental data sets to discover 23
24 influencing factors on effective HPCS development. 24

25

26 4.3.1 Views of a Single Subject 26

27

28 One view of a subject's work patterns is provided directly by the instrumentation. 28
29 We refer to this view as the *physical level view* since it objectively reports incontro- 29
30 vertible occurrences at the operating system level, such as the time stamp of each 30
31 compilation. 31

32 Figure 8 shows such a physical view for a 9-h segment of work done by a given 32
33 subject. The *x*-axis represents time and each dot on the graph represents a compile 33
34 event. Although we have the tools to measure physical activities with a high degree 34
35 of accuracy, this type of analysis does not yield much insight. For example, although 35
36 we know how often the compiler was invoked, we do not know why: We cannot 36
37 distinguish compilations which add new functionality from compilations which 37
38 correct problems or defects that could have been avoided. This is an important 38
39 distinction to make, if we want to know the cause of unnecessary rework so it can be 39
40 avoided. 40

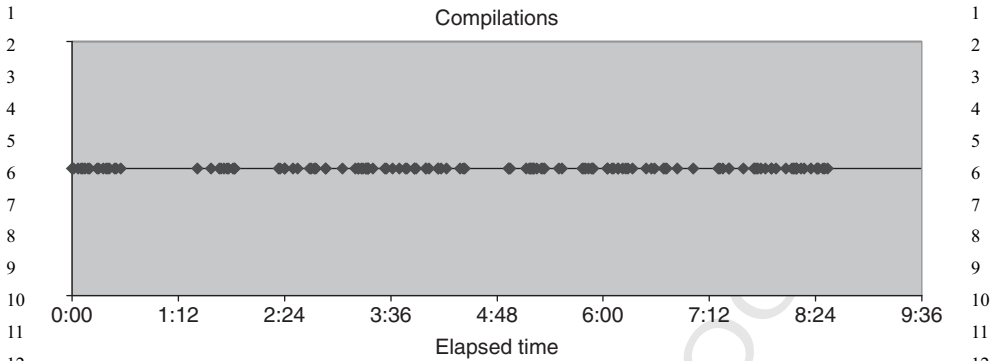


FIG. 8. Compilation events for one subject.

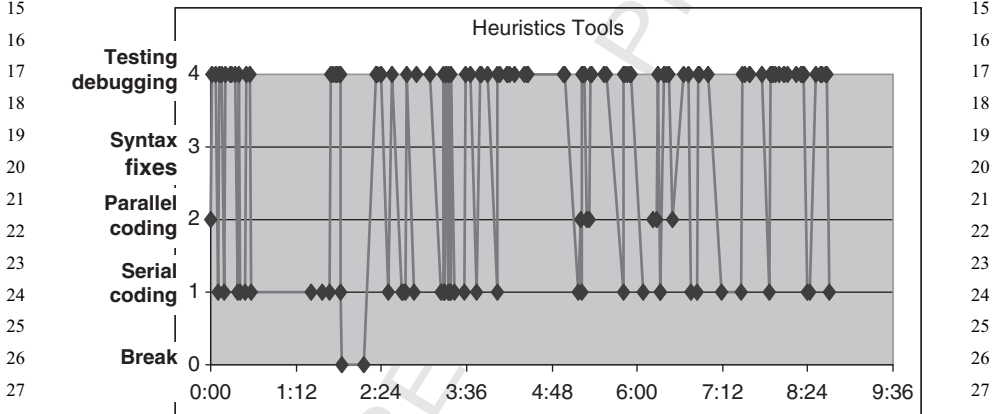


FIG. 9. Types of development activities for one subject.

A second approach is to use the logged compile times along with a snapshot of the code at each such compile and then apply a set of heuristics to guess at the semantics of the activities requiring that compilation. We have built such a tool that allows us to use various algorithms to manipulate these heuristics. The purpose of determining the semantic activities is to build baselines for predicting and evaluating the impact of new tools and languages.

Figure 9 uses the same data as Fig. 8 to illustrate this view. By evaluating the changes in the code for each compile, we can infer an activity being performed by the programmer. Again, each dot represents one compile but in this case the activity preceding each compile has been classified as one of:

- 1 • *Serial coding*: The developer is primarily focused on adding functionality 1
2 through serial code. This is inferred since most of the changes since the 2
3 previous compiler was in new code being added. 3
- 4 • *Parallel coding*: The developer is adding code to take advantage of multiple 4
5 processors, not just adding function calls to the parallel library. We decided to 5
6 separate out this activity since the amount of effort spent in this activity is 6
7 indicative of how difficult it is to take advantage of the parallel architecture for 7
8 solving the problem. This is inferred since parallel execution calls (such as to 8
9 the MPI library) were added to the program. 9
- 10 • *Syntax fixes*: The developer is fixing errors from a previous compile. We can 10
11 determine this since the previous compile failed, and the source program is 11
12 changed with no intervening execution. 12
- 13 • *Testing and debugging*: The developer is focused on finding and fixing a 13
14 problem, not adding new functionality. This activity can be identified via 14
15 some typical and recognizable testing strategies, such as when a high percentage 15
16 of the code added before a compile were output statements (so that variable 16
17 values can be checked at runtime); creating/modifying test data files instead of 17
18 the main code block; or removing test code back out of the system at the end of a 18
19 debugging session. Our hypothesis is that effort spent on these activities can 19
20 come from misunderstanding of the problem or the proposed solution and so 20
21 could be addressed with more sophisticated aides given to developers. 21
22

23 Such a view helps us understand better the approach used by the subject, and how 23
24 much of his/her time was spent on rework as opposed to adding new functionality. 24
25 The duration data associated with each activity also helps us identify interesting 25
26 events during the development: For example, when a large amount of time is spent 26
27 debugging, analysts can focus on events preceding the debugging activity to under- 27
28 stand what type of bug entered the system and how it was detected. This type of 28
29 information can be used to understand how hard or easy it is for the developer to 29
30 program in a given environment, and allows us to reason about what could be 30
31 changed to improve the situation. 31
32

34 4.3.2 *Validation of Workflow Heuristics* 34

35 The heuristics to date have been developed by having researchers examine the 35
36 collected data in detail (e.g., examining successive changes in source code versions). 36
37 However, the accuracy of these heuristics is not generally known. We have devel- 37
38 oped a tool that provides a programmer with information about the inferred activities 38
39 in real time. Using the tool, the developer then provides feedback about whether the 39
40

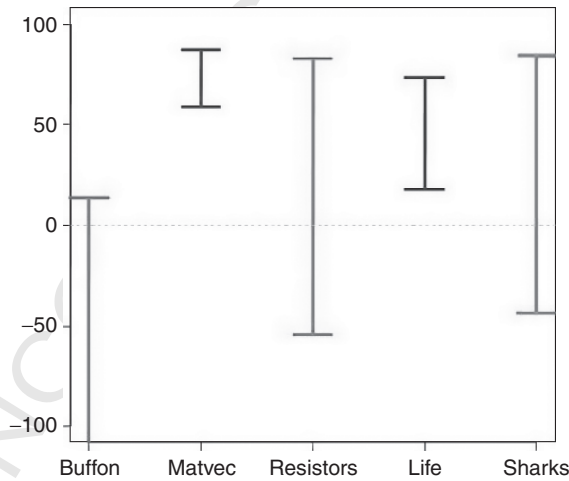
1 heuristic has correctly classified the activity. This allows us to evaluate how well the 1
2 heuristics agree with the programmer’s belief about the current development 2
3 activity. 3
4 4

5 4.3.3 Views of Multiple Subjects Across 5 6 Several Classes 6 7 7

8 From the same data, total effort can be calculated for each developer and 8
9 examined across problems and across classes to understand the range of variation 9
10 and whether causal factors can be related to changes in other measures of outcome, 10
11 such as the performance of the code produced. 11

12 We note that for our experimental paradigm to support effective analyses, subjects 12
13 in different classes who tackle the same problem using the same HPC approach should 13
14 exhibit similar results regarding effort and the performance achieved. Moreover, we 14
15 must be able to find measurable differences between subjects who, for example, 15
16 applied different approaches to the same problem. In a previous paper [19], we 16
17 presented some initial analyses of the data showing that both conditions hold. 17

18 Such analyses have been instrumental in developing an understanding of the 18
19 effectiveness of different HPC approaches in different contexts. For example, 19
20 Fig. 10 shows a comparison of effort data for two HPC approaches, ‘OpenMP’ 20
21 21



22 22
23 23
24 24
25 25
26 26
27 27
28 28
29 29
30 30
31 31
32 32
33 33
34 34
35 35
36 36
37 37
38 38
39 39
40 40
41 41
42 42
43 43
44 44
45 45
46 46
47 47
48 48
49 49
50 50
51 51
52 52
53 53
54 54
55 55
56 56
57 57
58 58
59 59
60 60
61 61
62 62
63 63
64 64
65 65
66 66
67 67
68 68
69 69
70 70
71 71
72 72
73 73
74 74
75 75
76 76
77 77
78 78
79 79
80 80
81 81
82 82
83 83
84 84
85 85
86 86
87 87
88 88
89 89
90 90
91 91
92 92
93 93
94 94
95 95
96 96
97 97
98 98
99 99
100 100
101 101
102 102
103 103
104 104
105 105
106 106
107 107
108 108
109 109
110 110
111 111
112 112
113 113
114 114
115 115
116 116
117 117
118 118
119 119
120 120
121 121
122 122
123 123
124 124
125 125
126 126
127 127
128 128
129 129
130 130
131 131
132 132
133 133
134 134
135 135
136 136
137 137
138 138
139 139
140 140
141 141
142 142
143 143
144 144
145 145
146 146
147 147
148 148
149 149
150 150
151 151
152 152
153 153
154 154
155 155
156 156
157 157
158 158
159 159
160 160
161 161
162 162
163 163
164 164
165 165
166 166
167 167
168 168
169 169
170 170
171 171
172 172
173 173
174 174
175 175
176 176
177 177
178 178
179 179
180 180
181 181
182 182
183 183
184 184
185 185
186 186
187 187
188 188
189 189
190 190
191 191
192 192
193 193
194 194
195 195
196 196
197 197
198 198
199 199
200 200
201 201
202 202
203 203
204 204
205 205
206 206
207 207
208 208
209 209
210 210
211 211
212 212
213 213
214 214
215 215
216 216
217 217
218 218
219 219
220 220
221 221
222 222
223 223
224 224
225 225
226 226
227 227
228 228
229 229
230 230
231 231
232 232
233 233
234 234
235 235
236 236
237 237
238 238
239 239
240 240
241 241
242 242
243 243
244 244
245 245
246 246
247 247
248 248
249 249
250 250
251 251
252 252
253 253
254 254
255 255
256 256
257 257
258 258
259 259
260 260
261 261
262 262
263 263
264 264
265 265
266 266
267 267
268 268
269 269
270 270
271 271
272 272
273 273
274 274
275 275
276 276
277 277
278 278
279 279
280 280
281 281
282 282
283 283
284 284
285 285
286 286
287 287
288 288
289 289
290 290
291 291
292 292
293 293
294 294
295 295
296 296
297 297
298 298
299 299
300 300
301 301
302 302
303 303
304 304
305 305
306 306
307 307
308 308
309 309
310 310
311 311
312 312
313 313
314 314
315 315
316 316
317 317
318 318
319 319
320 320
321 321
322 322
323 323
324 324
325 325
326 326
327 327
328 328
329 329
330 330
331 331
332 332
333 333
334 334
335 335
336 336
337 337
338 338
339 339
340 340
341 341
342 342
343 343
344 344
345 345
346 346
347 347
348 348
349 349
350 350
351 351
352 352
353 353
354 354
355 355
356 356
357 357
358 358
359 359
360 360
361 361
362 362
363 363
364 364
365 365
366 366
367 367
368 368
369 369
370 370
371 371
372 372
373 373
374 374
375 375
376 376
377 377
378 378
379 379
380 380
381 381
382 382
383 383
384 384
385 385
386 386
387 387
388 388
389 389
390 390
391 391
392 392
393 393
394 394
395 395
396 396
397 397
398 398
399 399
400 400
401 401
402 402
403 403
404 404
405 405
406 406
407 407
408 408
409 409
410 410
411 411
412 412
413 413
414 414
415 415
416 416
417 417
418 418
419 419
420 420
421 421
422 422
423 423
424 424
425 425
426 426
427 427
428 428
429 429
430 430
431 431
432 432
433 433
434 434
435 435
436 436
437 437
438 438
439 439
440 440
441 441
442 442
443 443
444 444
445 445
446 446
447 447
448 448
449 449
450 450
451 451
452 452
453 453
454 454
455 455
456 456
457 457
458 458
459 459
460 460
461 461
462 462
463 463
464 464
465 465
466 466
467 467
468 468
469 469
470 470
471 471
472 472
473 473
474 474
475 475
476 476
477 477
478 478
479 479
480 480
481 481
482 482
483 483
484 484
485 485
486 486
487 487
488 488
489 489
490 490
491 491
492 492
493 493
494 494
495 495
496 496
497 497
498 498
499 499
500 500
501 501
502 502
503 503
504 504
505 505
506 506
507 507
508 508
509 509
510 510
511 511
512 512
513 513
514 514
515 515
516 516
517 517
518 518
519 519
520 520
521 521
522 522
523 523
524 524
525 525
526 526
527 527
528 528
529 529
530 530
531 531
532 532
533 533
534 534
535 535
536 536
537 537
538 538
539 539
540 540
541 541
542 542
543 543
544 544
545 545
546 546
547 547
548 548
549 549
550 550
551 551
552 552
553 553
554 554
555 555
556 556
557 557
558 558
559 559
560 560
561 561
562 562
563 563
564 564
565 565
566 566
567 567
568 568
569 569
570 570
571 571
572 572
573 573
574 574
575 575
576 576
577 577
578 578
579 579
580 580
581 581
582 582
583 583
584 584
585 585
586 586
587 587
588 588
589 589
590 590
591 591
592 592
593 593
594 594
595 595
596 596
597 597
598 598
599 599
600 600
601 601
602 602
603 603
604 604
605 605
606 606
607 607
608 608
609 609
610 610
611 611
612 612
613 613
614 614
615 615
616 616
617 617
618 618
619 619
620 620
621 621
622 622
623 623
624 624
625 625
626 626
627 627
628 628
629 629
630 630
631 631
632 632
633 633
634 634
635 635
636 636
637 637
638 638
639 639
640 640
641 641
642 642
643 643
644 644
645 645
646 646
647 647
648 648
649 649
650 650
651 651
652 652
653 653
654 654
655 655
656 656
657 657
658 658
659 659
660 660
661 661
662 662
663 663
664 664
665 665
666 666
667 667
668 668
669 669
670 670
671 671
672 672
673 673
674 674
675 675
676 676
677 677
678 678
679 679
680 680
681 681
682 682
683 683
684 684
685 685
686 686
687 687
688 688
689 689
690 690
691 691
692 692
693 693
694 694
695 695
696 696
697 697
698 698
699 699
700 700
701 701
702 702
703 703
704 704
705 705
706 706
707 707
708 708
709 709
710 710
711 711
712 712
713 713
714 714
715 715
716 716
717 717
718 718
719 719
720 720
721 721
722 722
723 723
724 724
725 725
726 726
727 727
728 728
729 729
730 730
731 731
732 732
733 733
734 734
735 735
736 736
737 737
738 738
739 739
740 740
741 741
742 742
743 743
744 744
745 745
746 746
747 747
748 748
749 749
750 750
751 751
752 752
753 753
754 754
755 755
756 756
757 757
758 758
759 759
760 760
761 761
762 762
763 763
764 764
765 765
766 766
767 767
768 768
769 769
770 770
771 771
772 772
773 773
774 774
775 775
776 776
777 777
778 778
779 779
780 780
781 781
782 782
783 783
784 784
785 785
786 786
787 787
788 788
789 789
790 790
791 791
792 792
793 793
794 794
795 795
796 796
797 797
798 798
799 799
800 800
801 801
802 802
803 803
804 804
805 805
806 806
807 807
808 808
809 809
810 810
811 811
812 812
813 813
814 814
815 815
816 816
817 817
818 818
819 819
820 820
821 821
822 822
823 823
824 824
825 825
826 826
827 827
828 828
829 829
830 830
831 831
832 832
833 833
834 834
835 835
836 836
837 837
838 838
839 839
840 840
841 841
842 842
843 843
844 844
845 845
846 846
847 847
848 848
849 849
850 850
851 851
852 852
853 853
854 854
855 855
856 856
857 857
858 858
859 859
860 860
861 861
862 862
863 863
864 864
865 865
866 866
867 867
868 868
869 869
870 870
871 871
872 872
873 873
874 874
875 875
876 876
877 877
878 878
879 879
880 880
881 881
882 882
883 883
884 884
885 885
886 886
887 887
888 888
889 889
890 890
891 891
892 892
893 893
894 894
895 895
896 896
897 897
898 898
899 899
900 900
901 901
902 902
903 903
904 904
905 905
906 906
907 907
908 908
909 909
910 910
911 911
912 912
913 913
914 914
915 915
916 916
917 917
918 918
919 919
920 920
921 921
922 922
923 923
924 924
925 925
926 926
927 927
928 928
929 929
930 930
931 931
932 932
933 933
934 934
935 935
936 936
937 937
938 938
939 939
940 940
941 941
942 942
943 943
944 944
945 945
946 946
947 947
948 948
949 949
950 950
951 951
952 952
953 953
954 954
955 955
956 956
957 957
958 958
959 959
960 960
961 961
962 962
963 963
964 964
965 965
966 966
967 967
968 968
969 969
970 970
971 971
972 972
973 973
974 974
975 975
976 976
977 977
978 978
979 979
980 980
981 981
982 982
983 983
984 984
985 985
986 986
987 987
988 988
989 989
990 990
991 991
992 992
993 993
994 994
995 995
996 996
997 997
998 998
999 999
1000 1000

1 and ‘MPI.’ The percentage of effort saved by using OpenMP instead of MPI is shown 1
2 for each of five parallel programming problems (‘Buffon,’ ‘Matvec,’ etc.) represent- 2
3 ing different classes of parallel programs. Thus, 50 on the y-axis represents 50% less 3
4 effort for OpenMP; -50% would indicate that OpenMP required 50% more effort. 4
5 The height of each bar represents the range of values from across an entire dataset of 5
6 subjects. As can be seen, in two cases OpenMP yielded better results than MPI as all 6
7 subjects required less effort; for two other cases although there were some who 7
8 required less effort for MPI, the majority of data points indicated an effort savings 8
9 associated with OpenMP. In only one case, for the Buffon problem, did MPI appear 9
10 to give most subjects a savings in effort. As we gather more datasets using the 10
11 Experiment Manager tool suite, we will continue this type of analysis to understand 11
12 what other problems are in the set for which MPI requires less effort, and what it 12
13 is about these situations that sets them apart from the ones where OpenMP was the 13
14 less effort-intensive approach. (Interested readers can find a description of these 14
15 approaches and programming problems in other publications [19].) 15
16

17 4.4 Evaluation 17

18
19 We have been performing classroom experiments in the HPC domain since early 19
20 2003. While we have not performed a careful controlled experiment of its effective- 20
21 ness, we have observed anecdotally that the Experiment Manager avoids many of 21
22 the problems others (including ourselves) have observed in running experiments. 22
23 Many of these problems have already been reported in this paper. We have been able 23
24 to run the same experiment across multiple classes in multiple universities and 24
25 combine the results. Data is collected reliably and consistently across multiple 25
26 development platforms. We have been able to obtain data to install into our database 26
27 effortlessly without the need for students to perform any post-development activity. 27
28 Faculty, who are not experimental researchers, have been able to run their own 28
29 experiments with only minimal help from us. And finally, the ability to sanitize data 29
30 allows us to provide copies of datasets to others wanting to perform their own 30
31 analysis without running into IRB and privacy restrictions. 31
32

33 5. Related Work 34

35
36 There are various other projects that either support software engineering experi- 37
38 ments, or support automatic data collection during development, but not both. 38
39 The SESE system [1] has many similarities: it is web-based and supports features 39
40 such as managing subjects, supporting multiple roles, administering questionnaires, 40

1 capturing time spent during the experiment, collection of work products, and 1
2 monitoring of subject activity. By comparison, Experiment Manager supports addi- 2
3 tional data capture (e.g., intermediate source files and defects) and data analysis 3
4 (e.g., sanitization and workflow analysis). 4

5 PLUM, back in 1976, was one of the first systems to automatically collect 5
6 development data [24]. It, along with Hackystat [15], Ginger2 [22], Marmoset 6
7 [21], and Mylyn [17] are examples of systems which are designed to collect data 7
8 during the development process, but do not have data management facilities that are 8
9 specifically oriented towards running multiple experiments. Hackystat, which we 9
10 are using in the Experiment Manager, can collect data from several different types of 10
11 applications (e.g., vi, Emacs, Eclipse, JUnit, and Microsoft Word) via sensors. It was 11
12 originally designed for project monitoring rather than running experiments. We have 12
13 adopted the use of some of the Hackystat sensors into our data collection system. 13
14 Ginger2 is an environment for collecting an enormous amount of low-level detail 14
15 during software development, including eye-tracking and skin resistance. Marmoset 15
16 is an Eclipse-specific system which captures source code snapshots at each compile, 16
17 and is designed for computer science education research. Mylyn (originally called 17
18 Mylar) is also an Eclipse-specific system. Mylyn provides support for task-focused 18
19 code development and includes a framework for capturing and reporting on 19
20 information about Eclipse usage. 20
21 21

22 **6. Conclusions** 22

23 23
24 24
25 The classroom provides an excellent opportunity for conducting software engi- 25
26 neering experiments, but the complexities inherent in this environment makes such 26
27 research difficult to perform across multiple classes and at multiple sites. The 27
28 Experiment Manager framework supports the end-to-end process of conducting 28
29 software engineering experiments in the classroom environment. This allows 29
30 many others to run such experiments on their own in a way that allows for the 30
31 appropriate controls of the experiment so that results across classes and organization 31
32 at geographically diverse locations can be compared. The Experiment Manager 32
33 significantly reduces the effort on behalf of the experimentalists who are managing 33
34 the family of studies, and on the subjects themselves, by applying heuristics to infer 34
35 programmer activities. 35

36 We have successfully applied the Experiment Manager framework and with each 36
37 application are learning and improving the interface, simplifying the use by students, 37
38 making its use of value in shrinking the overall problem solving process by students, 38
39 for example, various forms of harnesses, the support for analysis, in order to get a 39
40 thorough understanding of the HPC development model. 40

ACKNOWLEDGMENTS

This research was supported in part by Department of Energy contract DE-FG02-04ER25633 and Air Force grant FA8750-05-1-0100 to the University of Maryland. Several students worked on the Experiment Manager including Patrick R. Borek, Thiago Escudeiro Craveiro, and Martin Voelp.

REFERENCES

- [1] Arisholm E., Sjoberg D. I. K., Carelius G. J., and Lindsjom Y., September 2002. A web-based support environment for software engineering experiments. *Nordic Journal of Computing*, **9**(3): 231–247.
- [2] Basili V. R., Selby W. R., and Phillips T.-Y., November 1983. Metric analysis and data validation across Fortran projects. *IEEE Transactions on Software Engineering*, *SE-9*, **6**: 652–663.
- [3] Basili V., and Green S., July 1994. Software process evolution at the SEL. *IEEE Software*, **11**(4): 58–66.
- [4] Basili V., July 1997. Evolving and packaging reading technologies. *Journal of Systems and Software*, **38**(1): 3–12.
- [5] Basili V., Shull F., and Lanubile F., July 1999. Building knowledge through families of experiments. *IEEE Transactions on Software Engineering*, **25**(4): 456–473.
- [6] Basili V., McGarry F., Pajerski R., and Zelkowitz M., May 2002. Lessons learned from 25 years of process improvement: The rise and fall of the NASA Software Engineering Laboratory. In *IEEE Computer Society and ACM International Conference on Software Engineering*, pp. 69–79. Orlando, FL.
- [7] Campbell D. T., and Stanley J. C., 1963. *Experimental and Quasi-Experimental Designs for Research*. Houghton-Mifflin.
- [8] Carlson W., Culler D., Yellick K., Brooks E., and Warren K., 1999. Introduction to UPC and language specification (CCS-TR-99-157). Technical report, Center for Computing Sciences.
- [9] Carver J., Jaccheri L., Morasca S., and Shull F., 2003. Issues in using students in empirical studies in software engineering education. In *International Symposium on Software Metrics*, pp. 239–249. Sydney, Australia.
- [10] Choy R., and Edelman A., 2003. *MATLAB*P 2.0: A unified parallel MATLAB*. Singapore MIT Alliance Symposium.
- [11] Dagum L., and Memon R., January 1998. OpenMP: An industry-standard API for shared-memory programming. *IEEE Computational Science & Engineering*, **5**(1): 46–55.
- [12] Dongarra J. J., Otta S. W., Snir M., and Walker D., July 1996. A message passing standard for MPP and workstations. *Communications of the ACM*, **39**(7): 84–90.
- [13] Hochstein L., Basili V., Zelkowitz M., Hollingsworth J., and Carver J., September 2005. Combining self-reported and automatic data to improve effort measurement. In *Joint 10th European Software Engineering Conference and 13th ACM SIGSOFT Symposium on the Foundations of Software Engineering*, pp. 356–365. Lisbon, Portugal.
- [14] Hochstein L., Carver J., Shull F., Asgari S., Basili V., Hollingsworth J. K., and Zelkowitz M., November 2005. HPC Programmer Productivity: A Case Study of Novice HPC Programmers, Supercomputing 2005. Seattle, WA.
- [15] Johnson P. M., Kou H., Agustin J. M., Zhang Q., Kagawa A., and Yamashita T., August 2004. Practical automated process and product metric collection and analysis in a classroom setting: Lessons learned from Hackstat-UH. In *International Symposium on Empirical Software Engineering*, Los Angeles, California.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

Au2

Au3

1 [16] Miller J., September 2000. Applying meta-analytical procedures to software engineering experi- 1
2 ments. *Journal of Systems and Software*, **54**(1): 29–39. 2
3 [17] Murphy G. C., Kersten M., and Findlater L., July/August 2006. How are Java Software Developers 3
4 using the eclipse IDE? *IEEE Software*, **23**(4): 76–83. 4
5 [18] Perry D. E., Staudenmayer N. A., and Votta L. G., 1995. Understanding and improving time usage in 5 [Au4]
6 software development. Volume 5 of *Trends in Software: Software Process* John Wiley & Sons. 6
7 [19] Shull F., Carver J., Hochstein L., and Basili V., 2005. Empirical study design in the area of high 7
8 performance computing (HPC). In *International Symposium on Empirical Software Engineering*,
9 Noosa Heads, Australia. 8
10 [20] Sjoberg D., Hannay J., Hansen O., Kampenes V., Karahasanovic A., Liborg N., and Rekdal A. C.,
11 September 2005. A survey of controlled experiments in software engineering. *IEEE Transactions on*
12 *Software Engineering*, **31**(9): 733–753. 10
13 [21] Spacco J., Strecker J., Hovemeyer D., and Pugh W., 2005. Software repository mining with 11
14 Marmoset: an automated programming project snapshot and testing system. In *Proceedings of the*
15 *International Workshop on Mining Software Repositories*, pp. 1–5. St. Louis, Missouri. 13
16 [22] Torii K., Mastumoto K., Nakakoji K., Takada Y., Takada S., and Shima K., July 1999. Ginger2: An 14
17 environment for computer-aided empirical software engineering. *IEEE Transactions on Software*
18 *Engineering*, **25**(4). 1 [Au5]
19 [23] Voelp M., August 2006. Diploma Thesis, Computer Science. University of Applied Sciences, 16
20 Mannheim, Germany. 17
21 [24] Zelkowitz M. V., October 1976. Automatic program analysis and evaluation. In *International*
22 *Conference on Software Engineering*, pp. 158–163. San Francisco, CA. 18
23 [25] Zelkowitz M. V., and Wallace D., May 1998. Experimental models for validating computer 19
24 technology. *IEEE Computer*, **31**(5): 23–31. 20
25 [26] Zelkowitz M. V., Basili V., Asgari S., Hochstein L., Hollingsworth J., and Nakamura T., September 21
26 2005. Productivity measures for high performance computers. In *International Symposium on*
27 *Software Metrics*, Como, Italy. 22
28 23
29 24
30 25
31 26
32 27
33 28
34 29
35 30
36 31
37 32
38 33
39 34
40 35
36
37
38
39
40

Author Query Form

Book Series: Advances in Computers, 74
Chapter 05

Dear Author,

During the preparation of your manuscript for typesetting some questions have arisen. These are listed below. Please check your typeset proof carefully and mark any corrections in the margin of the proof or compile them as a separate list. This form should then be returned with your marked proof/list of corrections to Elsevier Science.

Disk use

In some instances we may be unable to process the electronic file of your article and/or artwork. In that case we have, for efficiency reasons, proceeded by using the hard copy of your manuscript. If this is the case the reasons are indicated below:

- Disk damaged Incompatible file format LaTeX file for non-LaTeX journal
- Virus infected Discrepancies between electronic file and (peer-reviewed, therefore definitive) hard copy.
- Other:

We have proceeded as follows:

- Manuscript scanned Manuscript keyed in Artwork scanned
- Files only partly used (parts processed differently:.....)

Bibliography

If discrepancies were noted between the literature list and the text references, the following may apply:

- The references listed below were noted in the text but appear to be missing from your literature list. Please complete the list or remove the references from the text.
- Uncited references: This section comprises references which occur in the reference list but not in the body of the text. Please position each reference in the text or, alternatively, delete it. Any reference not dealt with will be retained in this section.

Query Refs.	Details Required	Author's response
AU1	Please check the running title for correctness.	
AU2	Please provide publisher's location in Ref. [7].	
AU3	Please provide publisher's name in Ref. [14].	
AU4	Please provide publisher's location in Ref. [18].	
AU5	Please provide page range in Ref. [22].	