

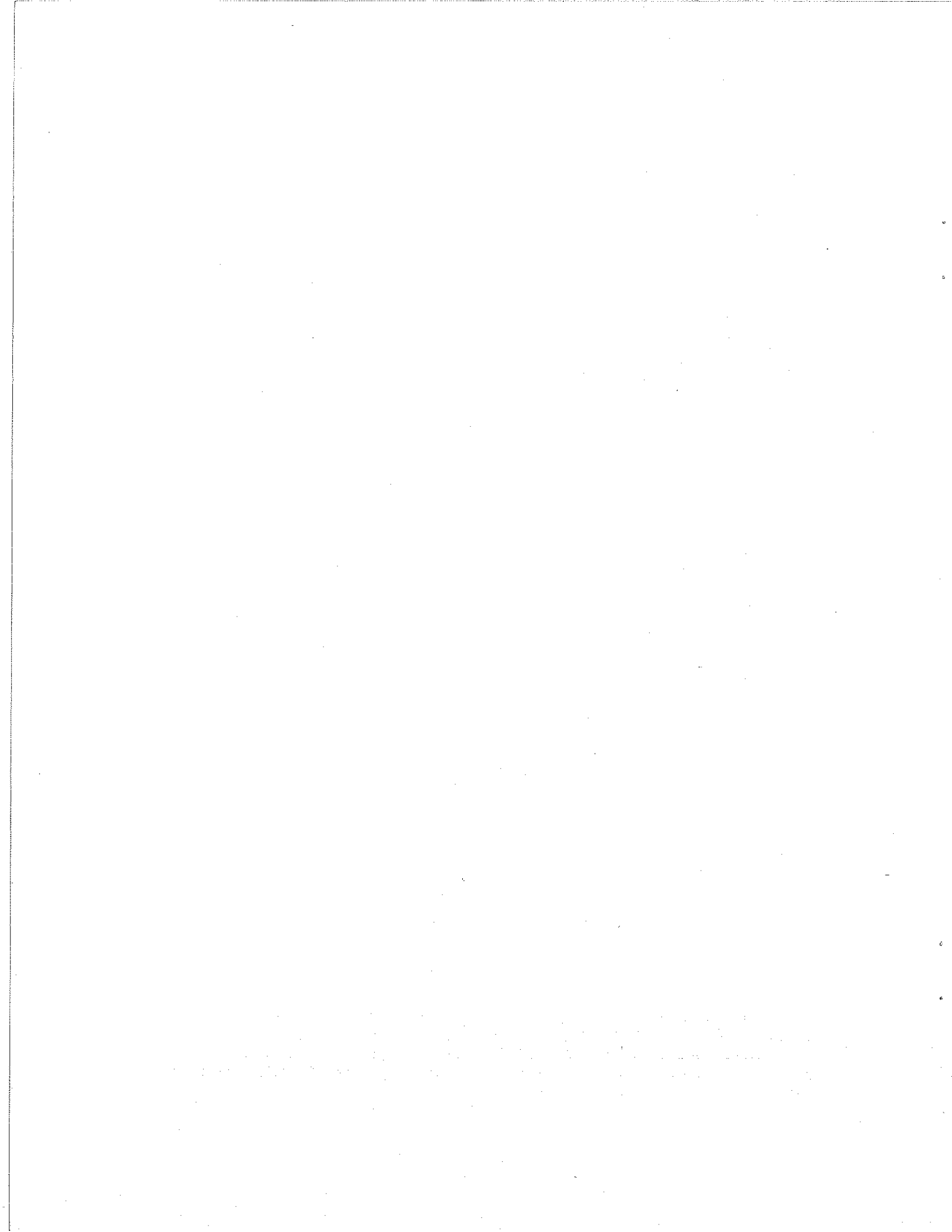
TR-1575

Nov. 20, 1985

EXPERIMENTATION
IN
SOFTWARE ENGINEERING

Victor R. Basili,
Richard W. Selby, Jr.,
David H. Hutchens

Research supported in part by the Air Force Office of Scientific Research Contract AFOSR-F49620-80C-001 and the National Aeronautics and Space Administration Grant NSG-5123 to the University of Maryland. Computer support provided by the Computer Science Center at the University of Maryland.



Experimentation in Software Engineering

Victor R. Basill¹, Richard W. Selby, Jr.²,
and David H. Hutchens³

¹ Department of Computer Science, University of Maryland, College Park, MD 20742
(301) 454-2002

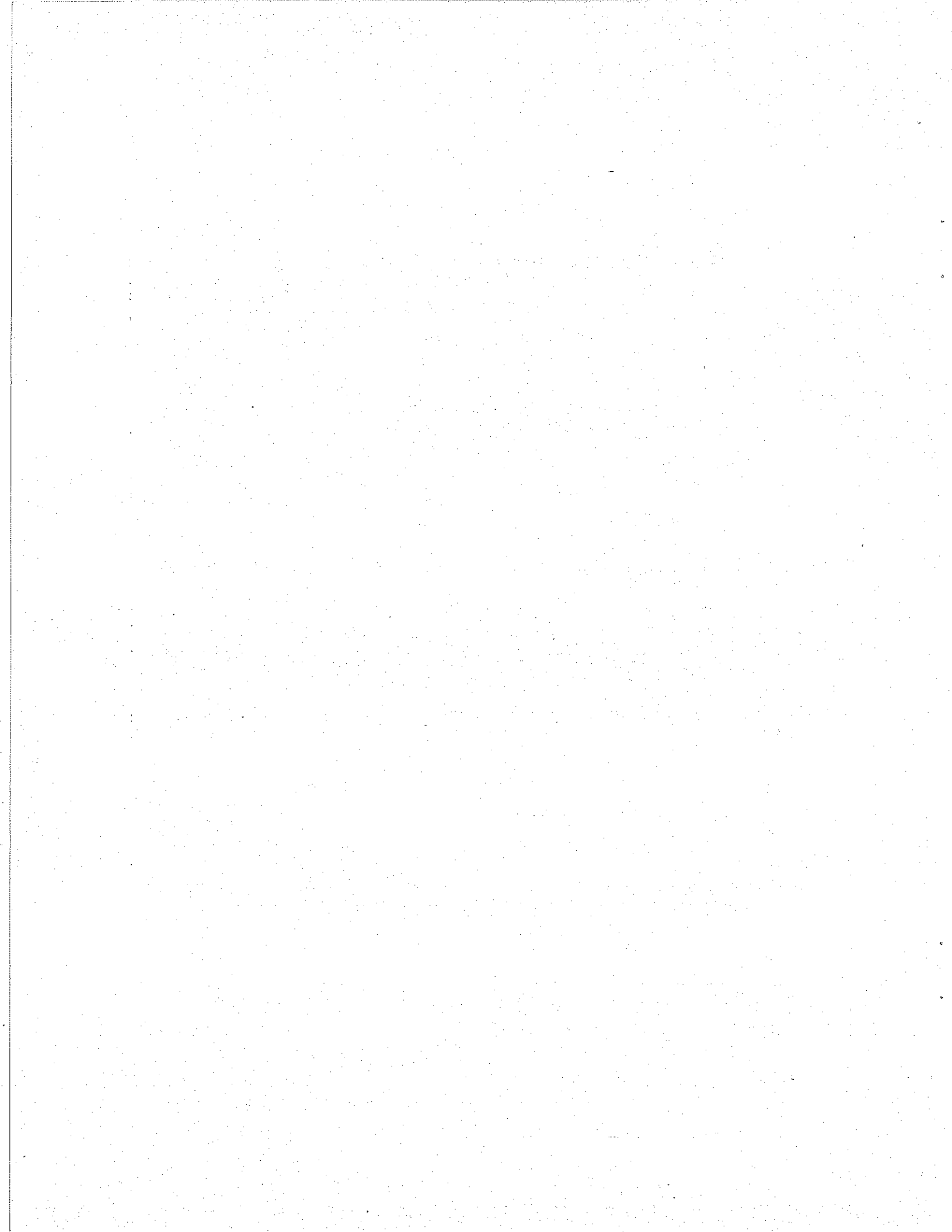
² Department of Information and Computer Science, University of California, Irvine, CA
92717 (714) 856-7403; was with the Department of Computer Science, University
of Maryland, College Park, MD 20742

³ Department of Computer Science, Clemson University, Clemson, SC 29631 (803) 654-
4464

KEYWORDS:

software technology measurement and evaluation, data collection and analysis, soft-
ware metrics, controlled experiment, experimental design, empirical study

Research supported in part by the Air Force Office of Scientific Research Contract AFOSR-F49620-80-C-001 and the National Aeronautics and Space Administration Grant NSG-5123 to the University of Maryland. Computer support provided in part by the Computer Science Center at the University of Maryland.



ABSTRACT

Experimentation in software engineering supports the advancement of the field through an iterative learning process. In this paper we present a framework for analyzing most of the experimental work performed in software engineering over the past several years. We describe a variety of experiments in the framework and discuss their contribution to the software engineering discipline. Some useful recommendations for the application of the experimental process in software engineering are included.

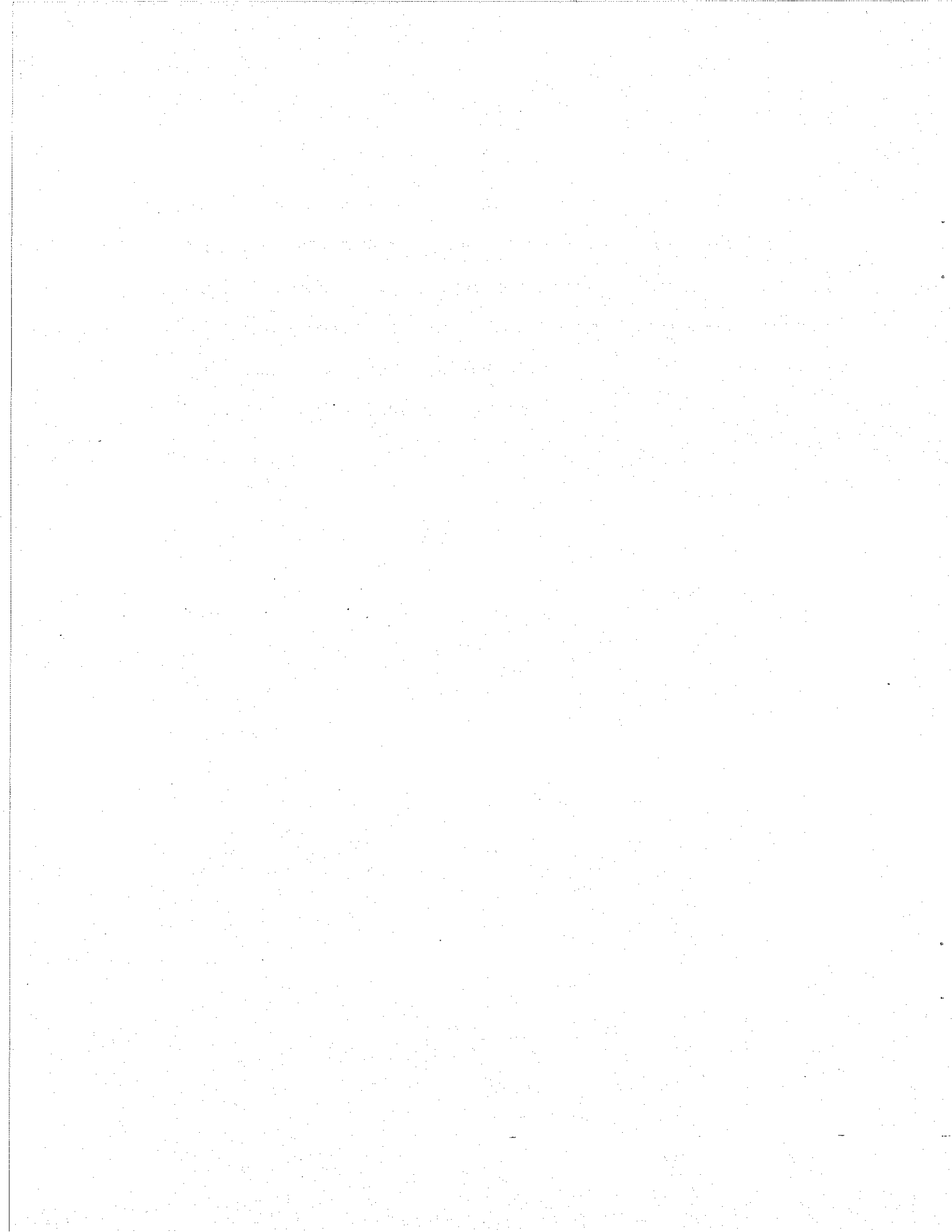
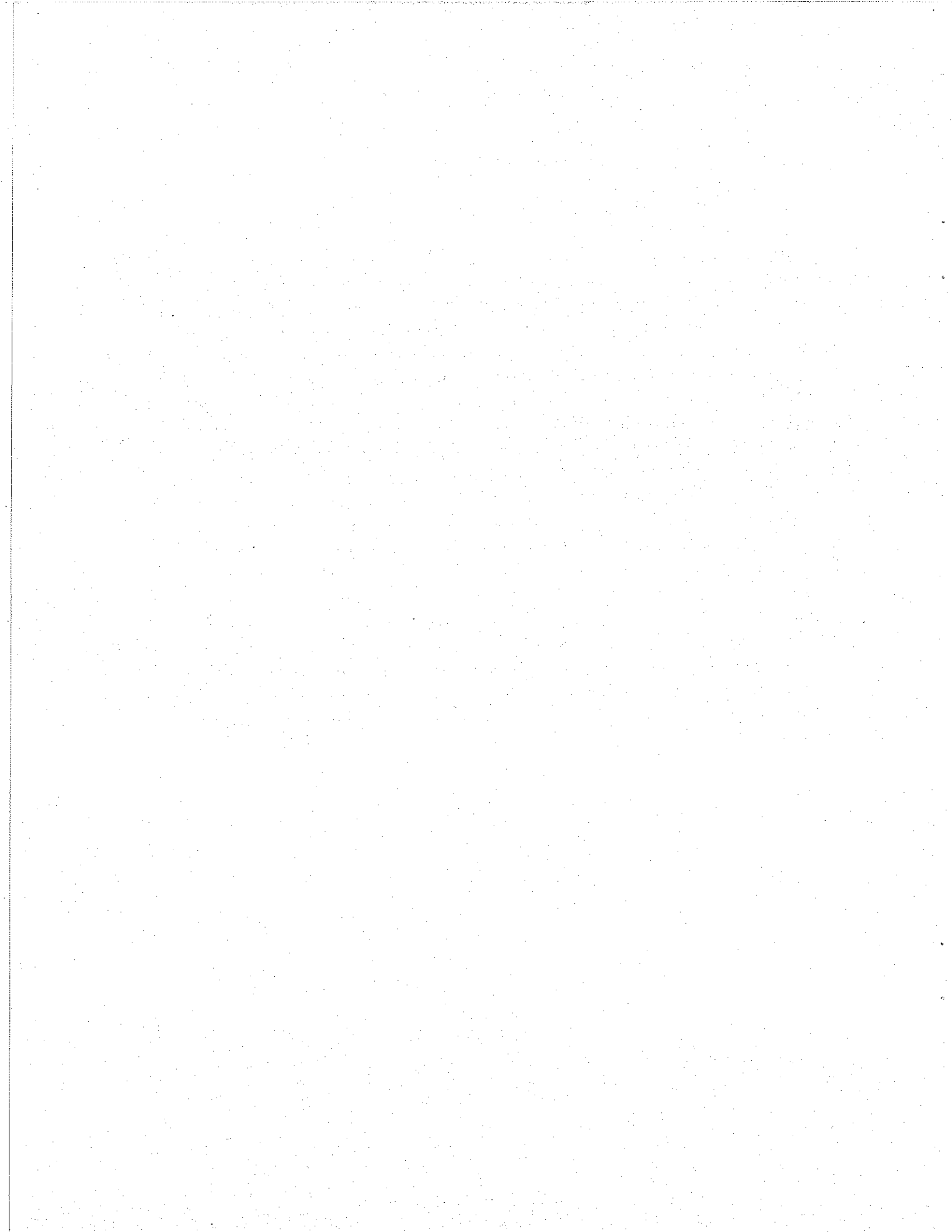


Table of Contents

1 Introduction	1
2 Objectives	1
3 Experimentation Framework	2
3.1 Experiment Definition	2
3.2 Experiment Planning	4
3.3 Experiment Operation	5
3.4 Experiment Interpretation	6
4 Classification of Analyses	6
4.1 Blocked Subject-Project Studies	7
4.2 Replicated Project Studies	11
4.3 Multi-Project Variation Studies	16
4.4 Single Project Studies	19
5 Problem Areas in Experimentation	22
5.1 Experimentation Overall	23
5.2 Experiment Definition	23
5.3 Experiment Planning	23
5.4 Experiment Operation	24
5.5 Experiment Interpretation	24
6 Conclusion	25
7 References	26



1. Introduction

As any area matures, there is the need to understand its components and their relationships. An experimental process provides a basis for the needed advancement in knowledge and understanding. Since software engineering is in its adolescence, it is certainly a candidate for the experimental method of analysis. Experimentation is performed in order to help us better evaluate, predict, understand, control, and improve the software development process and product.

Experimentation in software engineering, as with any other experimental procedure, involves an iteration of a hypothesize and test process. Models of the software process or product are built, hypotheses about these models are tested, and the information learned is used to refine the old hypotheses or develop new ones. In an area like software engineering, this approach takes on special importance because we greatly need to improve our knowledge of how software is developed, the effect of various technologies, and what areas most need improvement. There is a great deal to be learned and intuition is not always the best teacher.

In this paper we lay out a framework for analyzing most of the experimental work that has been performed in software engineering over the past several years. We then discuss a variety of these experiments, their results, and the impact they have had on our knowledge of the software engineering discipline.

2. Objectives

There are three overall goals for this work. The first objective is to describe a framework for experimentation in software engineering. The framework for experimen-

tation is intended to help structure the experimental process and to provide a classification scheme for understanding and evaluating experimental studies. The second objective is to classify and discuss a variety of experiments from the literature according to the framework. The description of several software engineering studies is intended to provide an overview of the knowledge resulting from experimental work, a summary of current research directions, and a basis for learning from past experience with experimentation. The third objective is to identify problem areas and lessons learned in experimentation in software engineering. The presentation of problem areas and lessons learned is intended to focus attention on general trends in the field and to provide the experimenter with useful recommendations for performing future studies. The following three sections address these goals.

3. Experimentation Framework

The framework of experimentation, summarized in Figure 1, consists of four categories corresponding to phases of the experimentation process: I) definition, II) planning, III) operation, and IV) interpretation. The following sections discuss each of these four phases.

3.1. Experiment Definition

The first phase of the experimental process is the study definition phase. The study definition phase contains six parts: A) motivation, B) object, C) purpose, D) perspective, E) domain, and F) scope. Most study definitions contain each of the six parts; an example definition appears in Figure 2.

There can be several motivations, objects, purposes, or perspectives in an experimental study. For example, the motivation of a study may be to understand, assess, or improve the effect of a certain technology. The "object of study" is the primary entity examined in a study. A study may examine the final software product, a development process (e.g., inspection process, change process), a model (e.g., software reliability model), etc. The purpose of a study may be to characterize the change in a system over time, to evaluate the effectiveness of testing processes, to predict system development cost by using a cost model, to motivate¹ the validity of a theory by analyzing empirical evidence, etc. In experimental studies that examine "software quality," the interpretation usually includes correctness if it is from the perspective of a developer or reliability if it is from the perspective of a customer. Studies that examine metrics for a given project type from the perspective of the project manager may interest certain project managers, while corporate managers may only be interested if the metrics apply across several project types.

Two important domains that are considered in experimental studies of software are 1) the individual programmers or programming teams (the "teams") and 2) the programs or projects (the "projects"). "Teams" are (possibly single-person) groups that work separately, and "projects" are separate programs or problems on which teams work. Teams may be characterized by experience, size, organization, etc., and projects may be characterized by size, complexity, application, etc. A general classification of the scopes

¹ For clarification, the usage of the word "motivate" as a study purpose is distinct from the study "motivation."

of experimental studies can be obtained by examining the sizes of these two domains considered (see Figure 3). Blocked subject-project studies examine one or more objects across a set of teams and a set of projects. Replicated project studies examine object(s) across a set of teams and a single project, while multi-project variation studies examine object(s) across a single team and a set of projects. Single project studies examine object(s) on a single team and a single project. As the representativeness of the samples examined and the scope of examination increase, the wider-reaching a study's conclusions become.

3.2. Experiment Planning

The second phase of the experimental process is the study planning phase. The following sections discuss aspects of the experiment planning phase: A) design, B) criteria, and C) measurement.

The design of an experiment couples the study scope with analytical methods and indicates the domain samples to be examined. Fractional factorial or randomized block designs usually apply in blocked subject-project studies, while completely randomized or incomplete block designs usually apply in multi-project and replicated project studies [33, 40]. Multivariate analysis methods, including correlation, factor analysis, and regression [75, 80, 89], generally may be used across all experimental scopes. Statistical models may be formulated and customized as appropriate [89]. Non-parametric methods should be planned when only limited data may be available or distributional assumptions may not be met [99]. Sampling techniques [41] may be used to select representative programmers and programs/projects to examine.

Different motivations, objects, purposes, perspectives, domains, and scopes require the examination of different criteria. Criteria that tend to be direct reflections of cost/quality include cost [111, 106, 86, 4, 28], errors/changes [49, 14, 109, 2, 81, 19], reliability, [42, 64, 56, 70, 69, 76, 77, 95], and correctness [51 61, 68]. Criteria that tend to be indirect reflections of cost/quality include data coupling [62, 48, 102, 78], information visibility [85, 83, 55], programmer understanding [98, 100, 107, 110], execution coverage [103, 21, 24], and size/complexity [17, 59, 71].

The concrete manifestations of the cost/quality aspects examined in the experiment are captured through measurement. Paradigms assist in the metric definition process: the goal-question-metric paradigm [20, 22, 25, 93] and the factor-criteria-metric paradigm [39, 72]. Once appropriate metrics have been defined, they may be validated to show that they capture what is intended [12, 18, 44, 50, 106, 113]. The data collection process includes developing automated collection schemes [15] and designing and testing data collection forms [22, 10]. The required data may include both objective and subjective data and different levels of measurement: nominal (or classificatory), ordinal (or ranking), interval, or ratio [99].

3.3. Experiment Operation

The third phase of the experimental process is the study operation phase. The operation of the experiment consists of A) preparation, B) execution, and C) analysis. Before conducting the actual experiment, preparation may include a pilot study to confirm the experimental scenario, help organize experimental factors (e.g., subject expertise), or inoculate the subjects [44, 43, 63, 24, 110, 73]. Experimenters collect and

validate the defined data during the execution of the study [18, 109]. The analysis of the data may include a combination of quantitative and qualitative methods [30]. The preliminary screening of the data, probably using plots and histograms, usually proceeds the formal data analysis. The process of analyzing the data requires the investigation of any underlying assumptions (e.g., distributional) before the application of the statistical models and tests.

3.4. Experiment Interpretation

The fourth phase of the experimental process is the study interpretation phase. The interpretation of the experiment consists of A) interpretation context, B) extrapolation, and C) impact. The results of the data analysis from a study are interpreted in a broadening series of contexts. These contexts of interpretation are the statistical framework in which the result is derived, the purpose of the particular study, and the knowledge in the field of research [15]. The representativeness of the sampling analyzed in a study qualifies the extrapolation of the results to other environments [20]. Several follow-up activities contribute to the impact of a study: presenting/publishing the results for feedback, replicating the experiment [33, 40], and actually applying the results by modifying methods for software development, maintenance, management, and research.

4. Classification of Analyses

Several investigators have published studies in the four general scopes of examination: blocked subject-project, replicated project, multi-project variation, or single project. The following sections cite studies from each of these categories. Note that sur-

veys on experimental methodology in empirical studies include [35, 96, 74]. Each of the sections first discusses one experiment in moderate depth, using italicized keywords from the framework for experimentation, and then chronologically presents an overview of several others in the category.

4.1. Blocked Subject-Project Studies

With a *motivation* to improve and better understand unit testing, [24] conducted a study whose *purpose* was to characterize and evaluate the processes (i.e., *objects*) of code reading, functional testing, and structural testing from the *perspective* of the developer. The testing processes were examined in a blocked subject-project *scope*, where 74 student through professional programmers (from the programmer *domain*) tested four unit-size programs (from the program *domain*) in a replicated fractional factorial *design*. Objective *measurement* of the testing processes was in several *criteria* areas: fault detection effectiveness, fault detection cost, and classes of faults detected. Experiment *preparation* included a pilot study [63], *execution* incorporated both manual and automated monitoring of testing activity, and *analysis* used analysis of variance methods [33, 90]. The major results (in the *interpretation context* of the study purpose) included 1) with the professionals, code reading detected more software faults and had a higher fault detection rate than did the other methods; 2) with the professionals, functional testing detected more faults than did structural testing, but they were not different in fault detection rate; 3) with the students, the three techniques were not different in performance, except that structural testing detected fewer faults than did the others in one study phase; and 4) overall, code reading detected more interface faults and functional

testing detected more control faults than did the other methods. A major result (in the *interpretation context* of the field of research) is that the study suggests that non-execution based fault detection, as in code reading, is at least as effective as on-line methods. The particular programmers and programs sampled qualify the *extrapolation* of the results. The *impact* of the study is an advancement in the understanding of effective software testing methods.

In order to understand program debugging, [57] evaluated several related factors, including effect of debugging aids, effect of fault type, and effect of particular program debugged from the perspective of the developer and maintainer. Thirty experienced programmers independently debugged one of four one-page programs that contained a single fault from one of three classes. The major results of these studies were 1) debugging is much faster if the programmer has had previous experience with the program, 2) assignment bugs were harder to find than other kinds, and 3) debugging aids did not seem to help programmers debug faster. Consistent results were obtained when the study was conducted on ten additional experienced programmers [58]. These results and the identification of possible "principles" of debugging contribute to the understanding of debugging methodology.

In order to improve experimental methodology and its application, [110] evaluated programmers' ability to understand and modify a program from the perspective of the developer and modifier. Various measures of programmer understanding were calculated, in a series of factorial design experiments, on groups of 16 - 48 university students performing tasks on two small programs. The study emphasized the need for well-structured and well-documented programs, and provided valuable testimony on and

worked toward a suitable experimentation methodology.

In order to assess the impact of language features on the programming process, [53] characterized the relationship of language features to software reliability from the perspective of the developer. Based on an analysis of the deficiencies in a programming language, nine different features were modified to produce a new version. Fifty-one advanced students were divided into two groups and asked to complete implementations of two small but sophisticated programs (75-200 line) in the original language and its modified version. The redesigned features in the two languages were contrasted in program fault frequency, type, and persistence. The experiment identified several language-design decisions that significantly affected reliability, which contributes to the understanding of language design for reliable software.

In order to understand the unit testing process better, [60] evaluated a reading technique and functional and "selective" testing (a composite approach) from the perspective of the developer. Thirty-nine university students applied the techniques to three unit-size programs in a Latin square design. Functional and "selective" testing were equally effective and both superior to the reading technique, which contributed to our understanding of testing methodology.

In order to improve and better understand the maintenance process, [43] conducted two experiments to evaluate factors that influence two aspects of software maintenance, program understanding and modification, from the perspective of the developer and maintainer. Thirty-six junior through advanced professional programmers in each experiment examined three classes of small (36 - 57 source line) programs in a factorial design. The factors examined include control flow complexity, variable name mnemoni-

city, type of modification, degree of commenting, and the relationship of programmer performance to various complexity metrics. In [44] they continued the investigation of how software characteristics relate to psychological complexity, and presented a third experiment to evaluate the ability of 54 professional programmers to detect program bugs in three programs in a factorial design. The series of experiments showed that software science [59] and cyclomatic complexity [71] measures are related to the difficulty experienced by programmers in locating errors in code.

In order to improve and better understand program debugging, [108] evaluated the theory that "programmers use 'slicing' (stripping away a program's statements that do not influence a given variable at a given statement) when debugging" from the perspective of the developer, maintainer, and researcher. Twenty-one university graduate students and programming staff debugged a fault in three unit-size (75 - 150 source line) programs in a non-parametric design. The study results supported the slicing theory, that is, programmers during debugging routinely partitioned programs into a coherent, discontinuous piece (or slice). The results advance the understanding of software debugging methodology.

In order to improve design techniques, [87] evaluated flowcharts and program design languages (PDL) from the perspective of the developer. Twenty-two graduate students designed two small (approximately 1000 source line) projects, one using flowcharts and the other using PDL. Overall, the results suggested that design performance and designer-programmer communication were better for projects using PDL.

In order to validate a theory of programming knowledge, [101] conducted two studies, using 139 novices and 41 professional programmers, to evaluate programmer behavior from the perspective of the researcher. The theory was that programming knowledge contained programming plans (generic program fragments representing common actions sequences) and rules of programming discourse (conventions used in composing plans into programs). The results support the existence and use of such plans and rules by both novice and advanced programmers.

Other blocked subject-project studies include [82, 112].

4.2. Replicated Project Studies

With a *motivation* to assess and better understand team software development methodologies, [15] conducted a study whose *purpose* was to characterize and evaluate the development processes (i.e., *objects*) of a a) disciplined-methodology team approach, b) ad hoc team approach, and c) ad hoc individual approach from the *perspective* of the developer and project manager. The development processes were examined in a replicated project *scope*, in which advanced university students comprising seven three-person teams, six three-person teams, and six individuals (from the programmer *domain*) used the approaches, respectively. They separately developed a small (600 - 2200 line) compiler (from the program *domain*) in a non-parametric *design*. Objective *measurement* of the development approaches was in several *criteria* areas: number of changes, number of program runs, program data usage, program data coupling/binding, static program size/complexity metrics, language usage, and modularity. Experiment *preparation* included presentation of relevant material [68, 7, 34], *execution* included automated

monitoring of on-line development activity and *analysis* used non-parametric comparison methods. The major results (in the *interpretation context* of the study purpose) included 1) the methodological discipline was a key influence on the general efficiency of the software development process; 2) the disciplined team methodology significantly reduced the costs of software development as reflected in program runs and changes; and 3) the examination of the effect of the development approaches was accomplished by the use of quantitative, objective, unobtrusive, and automatable process and product metrics. A major result (in the *interpretation context* of the field of research) is that the study supports the belief that incorporating discipline in software development reflects positively on both the development process and final product. The particular programmers and program sampled qualify the *extrapolation* of the results. The *impact* of the study is an advancement in the understanding of software development methodologies and their evaluation.

In order to improve the design and implementation processes, [84] evaluated system modularity from the perspective of the developer. Twenty university undergraduates each developed one of four different types of implementations for one of five different small modules. Then each of the modules were combined with others to form several versions of the whole system. The major results suggested that minor effort was required in assembling the systems and that major system changes can be confined to small, well-defined subsystems. The results support the ideas on formal specifications and modularity discussed in [83, 85] and advance the understanding of design methodology.

In order to assess the impact of static typing of programming languages in the development process, [54] evaluated the use of a statically typed language (having integers and strings) and a "typeless" language (e.g., arbitrary subscripting of memory) from the perspective of the developer. Thirty-eight students programmed a small (48 - 297 source line) problem in both languages, with half doing it in each order. The two languages were compared in the resulting program faults, the number of runs containing faults, and the relation of subject experience to fault proneness. The major result was that the use of a statically typed language can increase programming reliability, which assists in the design and use of programming languages.

In order to improve program composition, comprehension, debugging, and modification, [98] evaluated the use of detailed flowcharts in these tasks from the perspective of the developer, maintainer, modifier, and researcher. Groups of 53 - 70 novice through intermediate subjects, in a series of five experiments, performed various tasks using small programs. No significant differences were found between groups that used and those that did not use flowcharts, questioning the merit of using detailed flowcharts.

In order to improve and better understand the unit testing process, [79] evaluated the techniques of three-person walk-throughs, functional testing, and a control group from the perspective of the developer. Fifty-nine junior through advanced professional programmers applied the techniques to test a small (100 source line) but nontrivial program. The techniques were not different in the number of faults they detected, all pairings of techniques were superior to single techniques, and code reviews were less cost-effective than the others. These results assist in the selection of appropriate software

testing techniques.

In order to validate a particular metric family, [17] evaluated the ability of a proposed metric family to explain differences in system development methodologies and system changes from the perspective of the developer, project manager, and researcher. The metrics were applied to 19 versions of a small (600 - 2200) compiler, which were developed by teams of advanced university students using three different development approaches (see the first study [15] described in this section). The major results included 1) the metrics were able to differentiate among projects developed with different development methodologies; and 2) the differences among individuals had a large effect on the relationships between the metrics and aspects of system development. These results suggest insights into the formulation and appropriate use of software metrics.

In order to improve the understanding of why software errors occur, [65] characterized programmer misconceptions, cognitive strategies, and their manifestations as bugs in programs from the perspective of the developer and researcher. Two-hundred-four novice programmers separately attempted implementations of an elementary program. The results supported the programmers' intended use of "programming plans" [100] and revealed that most people preferred a read-process strategy over a process-read strategy. The results advance the understanding of how individuals write programs, why they sometimes make errors, and what programming language constructs should be available.

In order to understand the effect of coding conventions on program comprehensibility, [73] conducted a study to evaluate the relationship between indentation levels and program comprehension from the perspective of the developer. Eighty-six novice through professional subjects answered questions about one of seven program variations

with different level and type of indentation. The major result was that an indentation level of two or four spaces was preferred over zero or six.

In order to improve software development approaches, [29] characterized and evaluated the prototyping and specifying development approaches from the perspective of the developer, project manager, and user. Seven two- and three-person teams, consisting of university graduate students, developed versions of the same application software system (2000 - 4000 line); four teams used a requirement/design specifying approach and three teams used a prototyping approach. The systems developed by prototyping were smaller, required less development effort, and were easier to use. The systems developed by specifying had more coherent designs, more complete functionality, and software that was easier to integrate. These results contribute to the understanding of the merits and appropriateness of software development approaches.

In order to validate the theoretical model for N-version programming [66], [67, 3] conducted a study to evaluate the effectiveness of N-version programming for reliability from the perspective of the customer and user. N-version programming uses a high-level driver to connect several separately designed versions of the same system, the systems "vote" on the correct solution, and the solution provided by the majority of the systems is output. Twenty-seven graduate students were asked to independently design an 800 source line system. The factors examined included individual system reliability, total N-version system reliability, and classes of faults that occurred in systems simultaneously. The major result was that the assumption of independence of the faults in programs is not justified, and therefore, the reliability of the combined "voting" system may not be as high as given by the model.

In order to improve and better understand software development approaches, [94] characterized and evaluated the Cleanroom development approach [47, 46], in which software is developed without execution (i.e., completely off-line), from the perspective of the developer, project manager, and customer. Fifteen three-person teams of advanced university students separately developed a small system (800 - 2300 source line); ten teams used Cleanroom and five teams used a traditional development approach in a non-parametric design. The major results included 1) most developers using the Cleanroom approach were able to build systems without program execution; and 2) the Cleanroom teams' products met system requirements more completely and succeeded on more operational test cases than did those developed with a traditional approach. The results suggest the feasibility of complete off-line development, as in Cleanroom, and advance the understanding of software development methodology.

Other replicated project studies include [37, 5, 63].

4.3. Multi-Project Variation Studies

With a *motivation* to improve the understanding of resource usage during software development, [4] conducted a study whose *purpose* was to predict development cost by using a particular model (i.e., *object*) and to evaluate it from the *perspective* of the project manager, corporate manager, and researcher. The particular model generation method was examined in a multi-project *scope*, with baseline data from 18 large (2500 - 100,000 source line) software projects in the NASA S.E.L. production environment (from the program *domain*), in which teams contained from two to ten programmers (from the programmer *domain*) [10, 11, 38, 91]. The study *design* incorporated multivariate

methods to parameterize the model. Objective and subjective *measurement* of the projects was based on 21 *criteria*² in three areas: methodology, complexity, and personnel experience. Study *preparation* included preliminary work [52], *execution* included an established set of data collection forms [10], and *analysis* used forward multivariate regression methods. The major results (in the *interpretation context* of the study purpose) included 1) the estimation of software development resource usage improved by considering a set of both base-line and customization factors; 2) the application in the NASA environment of the proposed model generation method, which considers both types of factors, produced a resource usage estimate for a future project within one standard deviation of the actual; and 3) the confirmation of the NASA S.E.L. formula that the cost per line of reusing code is 20% of that of developing new code. A major result (in the *interpretation context* of the field of research) is that the study highlights the difference of each software development environment, which influences the use of resource estimation models. The particular programming environment and projects sampled qualify the *extrapolation* of the results. The *impact* of the study is an advancement in the understanding of estimating software development resource expenditure.

In order to assess, manage, and improve multi-project environments, [28, 26, 106, 13, 36, 18, 62, 109, 97, 105] characterized, evaluated, and/or predicted the effect of several factors from the perspective of the developer, modifier, project manager, and corporate manager. All the studies examined moderate to large projects from produc-

² Twenty-one factors were selected after examining a total of 82 factors that possibly contributed to project resource expenditure, including 36 from [106] and 16 from [28].

tion environments. The relationships investigated were among various factors, including structured programming, personnel background, development process and product constraints, project complexity, human and computer resource consumption, error-prone software identification, error/change distributions, data coupling/binding, project duration, staff size, degree of management control, and productivity. These studies have provided increased project visibility, greater understanding of classes of factors sensitive to project performance, awareness of the need for project measurement, and efforts for standardization of definitions. Analysis has begun on incorporating project variation information into a management tool [18, 23].

In order to improve and better understand the software maintenance process, [104] conducted an experiment to evaluate the relationship between the rate of maintenance repair and various product and process metrics from the perspective of the developer, user, and the project manager. A total of 447 small (up to 600 statements) commercial and clerical Cobol programs from one Australian organization and two U.S. organizations were analyzed. The product and process metrics included program complexity, programming style, programmer quality, and number of system releases. The major results were 1) in the Australian organization, program complexity and programming style significantly affected the maintenance repair rate; and 2) in the U.S. organizations, the number of times a system was released significantly affected the maintenance repair rate.

In order to improve the software maintenance process, [1] evaluated operational faults from the perspective of the user, customer, project manager, and corporate manager. The fault history for nine large production products (e.g., operating system

releases or their major components) was empirically modeled. He developed an approach for estimating whether and under what circumstances preventively fixing faults in operational software in the field was appropriate. Preventively fixing faults consists of installing fixes to faults that have yet to be discovered by particular users, but have been discovered by the vendor or other users. The major result is that for the typical user, corrective service is a reasonable way of dealing with most faults after the code has been in use for a fairly long period of time, while preventively fixing high-rate faults is advantageous during the time immediately following release.

In order to assess the effectiveness of the testing process, [31] evaluated estimations of the number of residual faults in a system from the perspective of the customer, developer, and project manager. The study was based on fault data collected from three large (2000 - 6000 module) systems developed in the Hughes-Fullerton environment. The study partitioned the faults based on severity and analyzed the differences in estimates of remaining faults according to stage of testing. Insights were gained into relationships between fault detection rates and residual faults.

4.4. Single Project Studies

With a *motivation* to improve software development methodology, [8] conducted a study whose *purpose* was to characterize the process (i.e., *object*) of iterative enhancement in conjunction with a top-down, stepwise refinement development approach from the *perspective* of the developer. The development process was examined in a single project *scope*, where the authors, two experienced individuals (from the programmer *domain*), built a 17,000 line compiler (from the program *domain*). The study *design* in-

incorporated descriptive methods to capture system evolution. Objective *measurement* of the system was in several *criteria* areas: size, modularity, local/global data usage, and data binding/coupling [62, 102]. Study *preparation* included language design [9], *execution* incorporated static analysis of system snapshots, and *analysis* used descriptive statistics. The results (in the *interpretation context* of the statistical framework) included 1) the percentage of global variables decreased over time while the percentage of actual vs. possible data couplings across modules increased, suggesting the usage of global data became more appropriate over time; and 2) the number of procedures and functions rose over time while the number of statements per procedure or function decreased, suggesting increased modularity. The major result of the study (in the *interpretation context* of the study purpose) was that the iterative enhancement technique encouraged the development of a software product that had several generally desirable aspects of system structure. A major result (in the *interpretation context* of the field of research) is that the study demonstrates the feasibility of iterative enhancement. The particular programming team and project examined qualify the *extrapolation* of the results. The *impact* of the study is an advancement in the understanding of software development approaches.

In order to improve, better understand, and manage the software development process, [6] evaluated the effect of applying chief programming teams and structured programming in system development from the perspective of the user, developer, project manager, and corporate manager. The large (83,000 line) system, known as "The New York Times Project," and was developed by a team of professionals organized as a chief programmer team, using structured code, top down design, walk-throughs, and program

libraries. Several benefits were identified, including reduced development time and cost, reduced time in system integration, and reduced fault detection in acceptance testing and field use. The results of the study demonstrated the feasibility of the chief programmer team concept and the accompanying methodologies in a production environment.

In order to improve their development environments through increased understanding, [49, 14, 2, 81, 19] each conducted single project studies to characterize the errors and changes made during a development project. They examined the development of a moderate to large software project, done by a multi-person team, in a production environment. They analyzed the frequency and distribution of errors during development and their relationship with several factors, including module size, software complexity, developer experience, method of detection and isolation, effort for isolation and correction, phase of entrance into the system and observance, reuse of existing design and code, and role of the requirements document. Such analyses have produced fault categorization schemes and have been useful in understanding and improving a development environment.

In order to improve design methodology, [55, 27] examined a ground-support system written in Ada³ to characterize the use of Ada packages from the perspective of the developer. Four professional programmers developed a project of 10,000 source lines of code. Factors such as how package use affected the ease of system modification and

³ Ada is a trademark of the Department of Defense.

how to measure module change resistance were identified, as well as how these observations related to aspects of the development and training. The major results were 1) several measures of Ada programs were developed, and 2) there was a indication that a lot of training will be necessary if we are to expect the facilities of Ada to be properly used.

In order to assess and improve software testing methodology, [21, 88] characterized and evaluated the relationship between system acceptance tests and operational usage from the perspective of the developer, project manager, customer, and researcher. The execution coverage of functionally generated acceptance test cases and a sample of operational usage cases was monitored for a medium-size (10,000 line) software system developed in a production environment. The results calculated that 64% of the program statements were executed during system operation and that the acceptance test cases corresponded reasonably well to the operational usage. The results give insights into the relationships among structural coverage, fault detection, system testing, and system usage.

5. Problem Areas in Experimentation

The following sections identify several problem areas of experimentation in software engineering. These areas may serve as guidelines in the performance of future studies. After mentioning some overall observations, cautions in each of the areas of experiment definition, planning, operation, and interpretation are discussed.

5.1. Experimentation Overall

There appears to be no "universal model" or "silver bullet" in software engineering. There are an enormous number of factors that differ across environments, in terms of desired cost/quality goals, methodology, experience, problem domain, constraints, etc. [106, 26, 4, 13, 28]. This results in every software development/maintenance/... environment being different. Another area of wide variation is the many-to-one differential in human performance [17, 45, 24]. The particular individuals examined in an empirical study can make an enormous difference. Among other considerations, these variations suggest that metrics need to be validated for a particular environment and a particular person to show that they capture what is intended [17, 18]. Thus, experimental studies should consider the potentially vast differences among environments and people.

5.2. Experiment Definition

In the definition of the purpose for the experiment, the formulation of intuitive problems into precisely stated goals is a nontrivial task [20, 22]. Defining the purpose of a study often requires the articulation of what is meant by "software quality." The many interpretations and perceptions of quality [32, 39, 72] highlight the need for considering whose perspective of quality is being examined. Thus, a precise specification of the problem to be investigated is a major step toward its solution.

5.3. Experiment Planning

Experimental planning should have a horizon beyond a first experiment. Controlled studies may be used to focus on the effect of certain factors, while their results may be confirmed in replications [92, 98, 101, 110, 57, 58, 44, 43, 24] and/or larger case

studies [4, 15]. When designing studies, consider that a combination of factors may be effective as a "critical mass," even though the particular factors may be ineffective when treated in isolation [15, 105]. Note that formal designs and the resulting statistical robustness are desirable, but we should not be driven exclusively by the achievement of statistical significance. Common sense must be maintained, which allows us, for example, to experiment just to help develop hypotheses [19, 109]. Thus, the experimental planning process should include a series of experiments for exploration, verification, and application.

5.4. Experiment Operation

The collection of the required data constitutes the primary result of the study operation phase. The data must be carefully defined, validated, and communicated to ensure its consistent interpretation by all persons associated with the experiment: subjects under observation, experimenters, and literature audience [18]. There have been papers in the literature that do not define their data well enough to enable a comparison of results across many projects and environments. We have often contacted the experimenter to discover that we are measuring different things. Thus, the experimenter should be cautious about the definition, validation, and communication of data, since they play a fundamental role in the experimental process.

5.5. Experiment Interpretation

The appropriate presentation of results from experiments contributes to their correct interpretation. Experimental results need to be qualified by the particular samples (e.g., programmers, programs) analyzed [20]. The extrapolation of results from a

particular sample must consider the representativeness of the sample to other environments [41, 111, 106, 86, 4, 28]. The visibility of the experimental results in professional forums and the open literature provides valuable feedback and constructive criticism. Thus, the presentation of experimental results should include appropriate qualification and adequate exposure to support their proper interpretation.

6. Conclusion

Experimentation in software engineering supports the advancement of the field through an iterative learning process. The experimental process has begun to be applied in a multiplicity of environments to study a variety of software technology areas. From the studies presented, it is clear that experimentation has proven effective in providing insights and furthering our domain of knowledge about the software process and product. In fact, there is a learning process in the experimentation approach itself, as has been shown in this paper.

We have described a framework for experimentation to provide a structure for presenting previous studies. We also recommend the framework as a mechanism to facilitate the definition, planning, operation, and interpretation of past and future studies. The problem areas discussed are meant to provide some useful recommendations for the application of the experimental process in software engineering. The experimental framework cannot be used in a vacuum; the framework and the lessons learned complement one another and should be used in a synergistic fashion. This work contributes to the understanding and advancement of experimentation in software engineering.

7. References

- [1] E. N. Adams, Optimizing Preventive Service of Software Products, *IBM Journal of Research and Development* **28**, 1, pp. 2-14, Jan. 1984.
- [2] J.-L. Albin and R. Ferreol, Collecte et analyse de mesures de logiciel (Collection and Analysis of Software Data), *Technique et Science Informatiques* **1**, 4, pp. 297-313, 1982. (Rairo ISSN 0752-4072)
- [3] A. Avizienis, P. Gunningberg, J. P. J. Kelly, L. Strigini, P. J. Traverse, K. S. Tso, and U. Voges, The UCLA Dedix System: A Distributed Testbed for Multiple-Version Software, *Digest Fifteenth Int. Sym. Fault-Tolerant Computing*, Ann Arbor, MI, June 19-21, 1985.
- [4] J. W. Bailey and V. R. Basili, A Meta-Model for Software Development Resource Expenditures, *Proc. Fifth Int. Conf. Software Engr.*, San Diego, CA, pp. 107-116, 1981.
- [5] J. W. Bailey, Teaching Ada: A Comparison of Two Approaches, Dept. Com. Sci., Univ. Maryland, College Park, MD, working paper, 1984.
- [6] F. T. Baker, System Quality Through Structured Programming, *AFIPS Proc. 1972 Fall Joint Computer Conf.* **41**, pp. 339-343, 1972.
- [7] V. R. Basili and F. T. Baker, Tutorial of Structured Programming, *Eleventh IEEE COMPCON*, IEEE Cat. No. 75CH1049-6, 1975.
- [8] V. R. Basili and A. J. Turner, Iterative enhancement: a practical technique for software development, *IEEE Transactions on Software Engineering* **SE-1**, 4, Dec. 1975.
- [9] V. R. Basili and A. J. Turner, *SIMPL-T: A Structured Programming Language*, Paladin House Publishers, Geneva, IL, 1976.
- [10] V. R. Basili, M. V. Zelkowitz, F. E. McGarry, R. W. Reiter, Jr., W. F. Truszkowski, and D. L. Weiss, The Software Engineering Laboratory, Software Eng. Lab., NASA/Goddard Space Flight Center, Greenbelt, MD, Rep. SEL-77-001, May 1977.
- [11] V. R. Basili and M. V. Zelkowitz, Analyzing Medium-Scale Software Developments, *Proc. Third Int. Conf. Software Engr.*, Atlanta, GA, pp. 116-123, May 1978.
- [12] V. R. Basili, *Tutorial on Models and Metrics for Software Management and Engineering*, IEEE Computer Society, New York, 1980.
- [13] V. R. Basili and K. Freburger, Programming Measurement and Estimation in the Software Engineering Laboratory, *Journal of Systems and Software* **2**, pp. 47-57, 1981.
- [14] V. R. Basili and D. M. Weiss, Evaluation of a Software Requirements Document By Analysis of Change Data, *Proc. Fifth Int. Conf. Software Engr.*, San Diego, CA, pp. 314-323, March 9-12, 1981.
- [15] V. R. Basili and R. W. Reiter, A Controlled Experiment Quantitatively Comparing Software Development Approaches, *IEEE Trans. Software Engr.* **SE-7**, May 1981.

- [16] V. R. Basili and C. Doerflinger, Monitoring Software Development Through Dynamic Variables, *Proc. COMPSAC*, Chicago, IL, 1983.
- [17] V. R. Basili and D. H. Hutchens, An Empirical Study of a Syntactic Metric Family, *Trans. Software Engr.* SE-9, 6, pp. 664-672, Nov. 1983.
- [18] V. R. Basili, R. W. Selby, Jr., and T. Y. Phillips, Metric Analysis and Data Validation Across FORTRAN Projects, *IEEE Trans. Software Engr.* SE-9, 6, pp. 652-663, Nov. 1983.
- [19] V. R. Basili and B. T. Perricone, Software Errors and Complexity: An Empirical Investigation, *Communications of the ACM* 27, 1, pp. 42-52, Jan. 1984.
- [20] V. R. Basili and R. W. Selby, Jr., Data Collection and Analysis in Software Research and Management, *Proceedings of the American Statistical Association and Biometric Society Joint Statistical Meetings*, Philadelphia, PA, August 13-16, 1984.
- [21] V. R. Basili and J. R. Ramsey, Structural Coverage of Functional Testing, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep. TR-1442, Sept. 1984.
- [22] V. R. Basili and D. M. Weiss, A Methodology for Collecting Valid Software Engineering Data*, *Trans. Software Engr.* SE-10, 6, pp. 728-738, Nov. 1984.
- [23] V. R. Basili and C. L. Ramsey, Arrowsmith-P - A Prototype Expert System for Software Engineering Management, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep., 1985. (submitted to the *Symposium on Expert Systems in Government*, Mclean, VA, Oct. 1985)
- [24] V. R. Basili and R. W. Selby, Jr., Comparing the Effectiveness of Software Testing Strategies, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep., 1985. (submitted to the *IEEE Trans. Software Engr.*)
- [25] V. R. Basili and R. W. Selby, Jr., Four Applications of a Software Data Collection and Analysis Methodology, *Proc. NATO Advanced Study Institute: The Challenge of Advanced Computing Technology to System Design Methods*, Durham, U. K., July 29 - August 10, 1985.
- [26] V. R. Basili and R. W. Selby, Jr., Calculation and Use of an Environment's Characteristic Software Metric Set, *Proc. Eighth Int. Conf. Software Engr.*, London, August 28-30, 1985.
- [27] V. R. Basili, E. E. Katz, N. M. Panlilio-Yap, C. L. Ramsey, and S. Chang, A Quantitative Characterization and Evaluation of a Software Development in Ada, *IEEE Computer*, September 1985.
- [28] B. W. Boehm, *Software Engineering Economics*, Prentice-Hall, Englewood Cliffs, NJ, 1981.
- [29] B. W. Boehm, T. E. Gray, and T. Seewaldt, Prototyping Versus Specifying: A Multiproject Experiment, *IEEE Trans. Software Engr.* SE-10, 3, pp. 290-303, May 1984.
- [30] R. C. Bogdan and S. K. Biklen, *Qualitative Research for Education: An Introduction to Theory and Methods*, Allyn and Bacon, Boston, MA, 1982.
- [31] J. Bowen, Estimation of Residual Faults and Testing Effectiveness, *Seventh Minnowbrook Workshop on Software Performance Evaluation*, Blue Mountain Lake, NY, July 24-27, 1984.

- [32] T. P. Bowen, G. B. Wagle, and J. T. Tsai, Specification of Software Quality Attributes, Rome Air Development Center, Griffiss Air Force Base, NY, Tech. Rep. RADC-TR-85-37 (three volumes), Feb. 1985.
- [33] G. E. P. Box, W. G. Hunter, and J. S. Hunter, *Statistics for Experimenters*, John Wiley & Sons, New York, 1978.
- [34] F. P. Brooks, Jr., *The Mythical Man-Month*, Addison-Wesley Publishing Co., Reading, MA, 1975.
- [35] R. E. Brooks, Studying Programmer Behavior: The Problem of Proper Methodology, *Communications of the ACM* **23**, 4, pp. 207-213, 1980.
- [36] W. D. Brooks, Software Technology Payoff: Some Statistical Evidence, *J. Systems and Software* **2**, pp. 3-9, 1981.
- [37] F. O. Buck, Indicators of Quality Inspections, IBM Systems Products Division, Kingston, NY, Tech. Rep. 21.802, Sept. 1981.
- [38] D. N. Card, F. E. McGarry, J. Page, S. Eslinger, and V. R. Basili, The Software Engineering Laboratory, Software Eng. Lab., NASA/Goddard Space Flight Center, Greenbelt, MD Rep. SEL-81-104, Feb. 1982.
- [39] J. P. Cavano and J. A. McCall, A Framework for the Measurement of Software Quality, *Proc. Software Quality and Assurance Workshop*, San Diego, CA, pp. 133-139, Nov. 1978.
- [40] W. G. Cochran and G. M. Cox, *Experimental Designs*, John Wiley & Sons, New York, 1950.
- [41] W. G. Cochran, *Sampling Techniques*, John Wiley & Sons, Inc., 1953.
- [42] P. A. Currit, M. Dyer, and H. D. Mills, Certifying the Reliability of Software, IBM Corp., Federal Systems Division, 6600 Rockledge Dr., Bethesda, MD, 20817, Tech. Rep., March 1985. (submitted to the *IEEE Trans. Software Engineering*)
- [43] B. Curtis, S. B. Sheppard, P. Milliman, M. A. Borst, and T. Love, Measuring the Psychological Complexity of Software Maintenance Tasks with the Halstead and McCabe Metrics, *IEEE Trans. Software Engr.*, pp. 96-104, March 1979.
- [44] B. Curtis, S. B. Sheppard, and P. M. Milliman, Third Time Charm: Stronger Replication of the Ability of Software Complexity Metrics to Predict Programmer Performance, *Proc. Fourth Int. Conf. Software Engr.*, pp. 356-360, Sept. 1979.
- [45] B. Curtis, Cognitive Science of Programming, *Sixth Minnowbrook Workshop on Software Performance Evaluation*, Blue Mountain Lake, NY, July 19-22, 1983.
- [46] M. Dyer and H. D. Mills, Developing Electronic Systems with Certifiable Reliability, *Proc. NATO Conf.*, Summer, 1982.
- [47] M. Dyer, Cleanroom Software Development Method, IBM Federal Systems Division, Bethesda, MD, October 14, 1982.
- [48] T. Emerson, A Discriminant Metric for Module Cohesion, *Proc. Seventh Intl. Conf. Software Engr.*, Orlando, FL, pp. 294-303, 1984.

- [49] A. Endres, An Analysis of Errors and their Causes in Systems Programs, *IEEE Trans. Software Engr.*, pp. 140-149, June 1975.
- [50] A. R. Feuer and E. B. Fowlkes, Some Results from an Empirical Study of Computer Software, *Proc. Fourth Int. Conf. Software Engr.*, pp. 351-355, 1979.
- [51] R. W. Floyd, Assigning Meaning to Programs, *Am. Math. Soc.* **19**, ed. J. T. Schwartz, Providence, RI, 1967.
- [52] K. Freburger and V. R. Basili, The Software Engineering Laboratory: Relationship Equations, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep. TR-764, May 1979.
- [53] J. D. Gannon and J. J. Horning, The Impact of Language Design on the Production of Reliable Software, *Trans. Software Engr.* **SE-1**, pp. 179-191, 1975.
- [54] J. D. Gannon, An Experimental Evaluation of Data Type Conventions, *Communications of the ACM* **20**, 8, pp. 584-595, 1977.
- [55] J. D. Gannon, E. E. Katz, and V. R. Basili, Characterizing Ada Programs: Packages, *The Measurement of Computer Software Performance*, Los Alamos National Laboratory, Aug. 1983.
- [56] A. L. Goel, Software Reliability and Estimation Techniques, Rome Air Development Center, Griffiss Air Force Base, NY, Rep. RADC-TR-82-263, October 1982.
- [57] J. D. Gould and P. Drongowski, An Exploratory Study of Computer Program Debugging, *Human Factors* **16**, 3, pp. 258-277, 1974.
- [58] J. D. Gould, Some Psychological Evidence on How People Debug Computer Programs, *International Journal of Man-Machine Studies* **7**, pp. 151-182, 1975.
- [59] M. H. Halstead, *Elements of Software Science*, North Holland, New York, 1977.
- [60] W. C. Hetzel, An Experimental Analysis of Program Verification Methods, Ph.D. Thesis, Univ. of North Carolina, Chapel Hill, 1976.
- [61] C. A. R. Hoare, An Axiomatic Basis for Computer Programming, *Communications of the ACM* **12**, 10, pp. 576-583, Oct. 1969.
- [62] D. H. Hutchens and V. R. Basili, System Structure Analysis: Clustering With Data Bindings, *IEEE Trans. Soft. Engr.* **SE-11**, 8, Aug. 1985.
- [63] S-S. V. Hwang, An Empirical Study in Functional Testing, Structural Testing, and Code Reading/Inspection*, Dept. Com. Sci., Univ. of Maryland, College Park, Scholarly Paper 362, Dec. 1981.
- [64] Z. Jelinski and P. B. Moranda, Applications of a Probability-Based Model to a Code Reading Experiment, *Proc. IEEE Symposium on Computer Software Reliability*, New York, pp. 78-81, IEEE, 1973.
- [65] W. L. Johnson, S. Draper, and E. Soloway, An Effective Bug Classification Scheme Must Take the Programmer into Account, *Proc. Workshop High-Level Debugging*, Palo Alto, CA, 1983.
- [66] J. P. J. Kelly, Specification of Fault-Tolerant Multi-Version Software: Experimental Studies of a Design Diversity Approach, UCLA Ph.D. Thesis, 1982.

- [67] J. Knight and N. Leveson, A Large Scale Experiment in N-Version Programming, *Proc. of the Ninth Annual Software Engineering Workshop*, NASA/GSFC, Greenbelt, MD, Nov. 1984.
- [68] R. C. Linger, H. D. Mills, and B. I. Witt, *Structured Programming: Theory and Practice*, Addison-Wesley, Reading, MA, 1979.
- [69] B. Littlewood and J. L. Verrall, A Bayesian Reliability Growth Model for Computer Software, *Applied Statistics* **22**, 3, 1973.
- [70] B. Littlewood, Stochastic Reliability Growth: A Model for Fault Renovation Computer Programs and Hardware Designs, *IEEE Trans. Reliability* **R-30**, 4, Oct. 1981.
- [71] T. J. McCabe, A Complexity Measure, *IEEE Trans. Software Engr.* **SE-2**, 4, pp. 308-320, Dec. 1976.
- [72] J. A. McCall, P. Richards, and G. Walters, Factors in Software Quality, Rome Air Development Center, Griffiss Air Force Base, NY, Tech. Rep. RADC-TR-77-369, Nov. 1977.
- [73] R. J. Miara, J. A. Musselman, J. A. Navarro, and B. Shneiderman, Program Indentation and Comprehensibility, *Communications of the ACM* **26**, 11, pp. 861-867, Nov. 1983.
- [74] T. Moher and G. M. Schneider, Methodology and Experimental Research in Software Engineering, *International Journal of Man-Machine Studies* **16**, 1, pp. 65-87, 1982.
- [75] S. A. Mulaik, *The Foundations of Factor Analysis*, McGraw-Hill, New York, 1972.
- [76] J. D. Musa, A Theory of Software Reliability and Its Application, *IEEE Trans. Software Engr.* **SE-1**, 3, pp. 312-327, 1975.
- [77] J. D. Musa, Software reliability measurement, *Journal of Systems and Software* **1**, 3, pp. 223-241, 1980.
- [78] G. J. Myers, *Composite/Structured Design*, Van Nostrand Reinhold, 1978.
- [79] G. J. Myers, A Controlled Experiment in Program Testing and Code Walkthroughs/Inspections, *Communications of the ACM*, pp. 760-768, Sept. 1978.
- [80] J. Neter and W. Wasserman, *Applied Linear Statistical Models*, Richard D. Irwin, Inc., Homewood, IL, 1974.
- [81] T. J. Ostrand and E. J. Weyuker, Collecting and Categorizing Software Error Data in an Industrial Environment, Dept. Com. Sci., Courant Inst. Math. Sci., New York Univ., NY, Tech. Rep. 47, August 1982 (Revised May 1983).
- [82] D. J. Panzl, Experience with Automatic Program Testing, *Proc. NBS Trends and Applications*, Nat. Bureau Stds., Gaithersburg, MD, pp. 25-28, May, 28 1981.
- [83] D. L. Parnas, On the Criteria to be Used in Decomposing Systems into Modules, *Communications of the ACM* **15**, 12, pp. 1053-1058, 1972.
- [84] D. L. Parnas, Some Conclusions from an Experiment in Software Engineering Techniques, *AFIPS Proc. 1972 Fall Joint Computer Conf.* **41**, pp. 325-329, 1972.

- [85] D. L. Parnas, A Technique for Module Specification With Examples, *Communications of the ACM* 15, May 1972.
- [86] L. Putnam, A General Empirical Solution to the Macro Software Sizing and Estimating Problem, *IEEE Trans. Software Engr.* 4, 4, 1978.
- [87] H.R. Ramsey, M.E. Atwood, and J.R. Van Doren, Flowcharts Versus Program Design Languages: An Experimental Comparison, *Communications ACM* 26, 6, pp. 445-449, June 1983.
- [88] J. Ramsey, Structural Coverage of Functional Testing, *Seventh Minnowbrook Workshop on Software Performance Evaluation*, Blue Mountain Lake, NY, July 24-27, 1984.
- [89] Statistical Analysis System (SAS) User's Guide, SAS Institute Inc., Box 8000, Cary, NC, 27511, 1982.
- [90] H. Scheffe, *The Analysis of Variance*, John Wiley & Sons, New York, 1959.
- [91] Annotated Bibliography of Software Engineering Laboratory (SEL) Literature, Software Eng. Lab., NASA/Goddard Space Flight Center, Greenbelt, MD Rep. SEL-82-006, Nov. 1982.
- [92] R. W. Selby, Jr., An Empirical Study Comparing Software Testing Techniques, *Sixth Minnowbrook Workshop on Software Performance Evaluation*, Blue Mountain Lake, NY, July 19-22, 1983.
- [93] R. W. Selby, Jr., Evaluations of Software Technologies: Testing, CLEANROOM, and Metrics, Dept. Com. Sci., Univ. Maryland, College Park, Ph. D. Dissertation, 1985.
- [94] R. W. Selby, Jr., V. R. Basili, and F. T. Baker, CLEANROOM Software Development: An Empirical Evaluation, Dept. Com. Sci., Univ. Maryland, College Park, Tech. Rep. TR-1415, February 1985. (submitted to the *IEEE Trans. Software Engr.*)
- [95] J. G. Shanthikumar, A Statistical Time Dependent Error Occurrence Rate Software Reliability Model with Imperfect Debugging, *Proc. 1981 National Computer Conference*, June 1981.
- [96] B. A. Sheil, The Psychological Study of Programming, *Computing Surveys* 13, pp. 101-120, March 1981.
- [97] V.Y. Shen, T.J. Yu, S.M. Thebaut, and L.R. Paulsen, Identifying Error-Prone Software - An Empirical Study, *IEEE Trans. Soft. Engr.* SE-11, 4, pp. 317-324, April 1985.
- [98] B. Shneiderman, R. E. Mayer, D. McKay, and P. Heller, Experimental Investigations of the Utility of Detailed Flowcharts in Programming, *Communications of the ACM* 20, 6, pp. 373-381, 1977.
- [99] S. Siegel, *Nonparametric Statistics for the Behavioral Sciences*, McGraw-Hill, New York, 1955.
- [100] E. Soloway, K. Ehrlich, J. Bonar, and J. Greenspan, What Do Novices Know About Programming?, in *Directions in Human-Computer Interactions*, ed. A. Badre and B. Shneiderman, Ablex, Inc., 1982.
- [101] E. Soloway and K. Ehrlich, Empirical Studies of Programming Knowledge, *Trans. Software Engr.* SE-10, 5, pp. 595-609, Sept. 1984.

- [102] W. P. Stevens, G. J. Myers, and L. L. Constantine, Structural Design, *IBM Systems Journal* **13**, 2, pp. 115-139, 1974.
- [103] L. G. Stucki, New Directions in Automated Tools for Improving Software Quality, in *Current Trends in Programming Methodology*, ed. R. T. Yeh, Prentice Hall, Englewood Cliffs, NJ, 1977.
- [104] I. Vessey and R. Weber, Some Factors Affecting Program Repair Maintenance: An Empirical Study, *Communications of the ACM* **26**, 2, pp. 128-134, Feb. 1983.
- [105] J. Vosburgh, B. Curtis, R. Wolverson, B. Albert, H. Malec, S. Hoben, and Y. Liu, Productivity Factors and Programming Environments, *Proc. Seventh Int. Conf. Software Engr.*, Orlando, FL, pp. 143-152, 1984.
- [106] C. E. Walston and C. P. Felix, A Method of Programming Measurement and Estimation, *IBM Systems J.* **16**, 1, pp. 54-73, 1977.
- [107] G. Weinberg, *The Psychology of Computer Programming*, Van Nostrand Reinhold Co., 1971.
- [108] M. Weiser, Programmers Use Slices When Debugging, *Communications ACM* **25**, pp. 446-452, July 1982.
- [109] D. M. Weiss and V. R. Basili, Evaluating Software Development by Analysis of Changes: Some Data from the Software Engineering Laboratory, *IEEE Trans. Software Engr.* **SE-11**, 2, pp. 157-168, February 1985.
- [110] L. Weissman, Psychological Complexity of Computer Programs: An Experimental Methodology, *SIGPLAN Notices* **9**, 6, pp. 25 - 36, June 1974.
- [111] R. Wolverson, The Cost of Developing Large Scale Software, *IEEE Trans. Computers* **23**, 6, 1974.
- [112] S. N. Woodfield, H. E. Dunsmore, and V. Y. Shen, The Effect of Modularization and Comments on Program Comprehension, Dept. Com. Sci., Arizona St. Univ., Tempe, AZ, working paper, 1981.
- [113] J. C. Zolnowski and D. B. Simmons, Taking the Measure of Program Complexity, *Proc. National Computer Conference*, pp. 329-336, 1981.

Figure 1. Summary of the framework of experimentation.

I. Definition					
Motivation	Object	Purpose	Perspective	Domain	Scope
Understand Assess Manage Engineer Learn Improve Validate Assure	Product Process Model Metric Theory	Characterize Evaluate Predict Motivate	Developer Modifier Maintainer Project manager Corporate manager Customer User Researcher	Programmer Program/project	Single project Multi-project Replicated project Blocked subject-project
II. Planning					
Design		Criteria		Measurement	
Experimental designs Incomplete block Completely randomized Randomized block Fractional factorial Multivariate analysis Correlation Factor analysis Regression Statistical models Non-parametric Sampling		Direct reflections of cost/quality Cost Errors Changes Reliability Correctness Indirect reflections of cost/quality Data coupling Information visibility Programmer comprehension Execution coverage Size Complexity		Metric definition Goal-question-metric Factor-criteria-metric Metric validation Data collection Automatability Form design and test Objective vs. subjective Level of measurement Nominal/classificatory Ordinal/ranking Interval Ratio	
III. Operation					
Preparation		Execution		Analysis	
Pilot study		Data collection Data validation		Quantitative vs. qualitative Preliminary data analysis Plots and histograms Model assumptions Primary data analysis Model application	
IV. Interpretation					
Interpretation context		Extrapolation		Impact	
Statistical framework Study purpose Field of research		Sample representativeness		Visibility Replication Application	

Figure 2. Study definition example.

Definition element	example
Motivation	To improve the unit testing process,
Purpose	characterize and evaluate
Object	the processes of functional and structural testing
Perspective	from the perspective of the developer
Domain: programmer	as they are applied by experienced programmers
Domain: program	to unit-size software
Scope	in a blocked subject-project study.

Figure 3. Experimental scopes.

#Teams per project	#Projects	
	one	more than one
one	Single project	Multi-project variation
more than one	Replicated project	Blocked subject-project

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE

REPORT DOCUMENTATION PAGE

1a. REPORT SECURITY CLASSIFICATION		1b. RESTRICTIVE MARKINGS		
2a. SECURITY CLASSIFICATION AUTHORITY		3. DISTRIBUTION/AVAILABILITY OF REPORT		
2b. DECLASSIFICATION/DOWNGRADING SCHEDULE				
4. PERFORMING ORGANIZATION REPORT NUMBER(S) Department of Computer Science University of Maryland TR-1575		5. MONITORING ORGANIZATION REPORT NUMBER(S)		
6a. NAME OF PERFORMING ORGANIZATION Department of Computer Sci. University of Maryland	6b. OFFICE SYMBOL (If applicable) Uof MD	7a. NAME OF MONITORING ORGANIZATION		
6c. ADDRESS (City, State and ZIP Code) University of Maryland College Park, Maryland 20742		7b. ADDRESS (City, State and ZIP Code)		
8a. NAME OF FUNDING/SPONSORING ORGANIZATION AFOSR/NASA	8b. OFFICE SYMBOL (If applicable)	9. PROCUREMENT INSTRUMENT IDENTIFICATION NUMBER		
8c. ADDRESS (City, State and ZIP Code)		10. SOURCE OF FUNDING NOS.		
		PROGRAM ELEMENT NO.	PROJECT NO.	TASK NO.
11. TITLE (Include Security Classification) Experimentation in Software Engineering				
12. PERSONAL AUTHOR(S) Victor R. Basili, Richard W. Selby, Jr., and David H. Hutchens				
13a. TYPE OF REPORT Technical/Scientific	13b. TIME COVERED FROM _____ TO _____	14. DATE OF REPORT (Yr., Mo., Day) November 20, 1985	15. PAGE COUNT 32	
16. SUPPLEMENTARY NOTATION				
17. COSATI CODES		18. SUBJECT TERMS (Continue on reverse if necessary and identify by block number)		
FIELD	GROUP			SUB. GR.
19. ABSTRACT (Continue on reverse if necessary and identify by block number) Experimentation in software engineering supports the advancement of the field through an iterative learning process. In this paper we present a framework for analyzing most of the experimental work performed in software engineering over the past several years. We describe a variety of experiments in the framework and discuss their contribution to the software engineering discipline. Some useful recommendations for the application of the experimental process in software engineering are included.				
20. DISTRIBUTION/AVAILABILITY OF ABSTRACT UNCLASSIFIED/UNLIMITED <input checked="" type="checkbox"/> SAME AS RPT. <input type="checkbox"/> DTIC USERS <input type="checkbox"/>		21. ABSTRACT SECURITY CLASSIFICATION Unclassified		
22a. NAME OF RESPONSIBLE INDIVIDUAL Victor R. Basili		22b. TELEPHONE NUMBER (Include Area Code) 301-454-2002	22c. OFFICE SYMBOL	

