

Investigating Reading Techniques for Object-Oriented Framework Learning

Forrest Shull, Filippo Lanubile, and Victor R. Basili, *Fellow, IEEE*

Abstract--The empirical study described in this paper addresses software reading for construction: how application developers obtain an understanding of a software artifact for use in new system development. This study focuses on the processes that developers would engage in when learning and using object-oriented frameworks. We analyzed 15 student software development projects using both qualitative and quantitative methods to gain insight into what processes occurred during framework usage. The contribution of the study is not to test predefined hypotheses but to generate well-supported hypotheses for further investigation. The main hypotheses we produce are that example-based techniques are well suited to use by beginning learners while hierarchy-based techniques are not because of a larger learning curve. Other more specific hypotheses are proposed and discussed.

Index Terms--object-oriented frameworks, software reading, empirical study

I. INTRODUCTION

In almost any software development environment, the various work documents used (e.g. requirements documents, code and design plans) require continual review and modification throughout the lifecycle. This is due to the central role such documents play in many software engineering tasks (e.g. verification and validation, maintenance, evolution, and reuse). Software reading, i.e., the individual analysis of textual software work products aimed at achieving whatever degree of understanding is needed to accomplish a particular task, is thus a key technical activity for software development. A number of studies have examined reading techniques applied to a variety of software engineering tasks, such as detecting defects in requirements [3, 4, 41, 44, 49], assessing user interfaces [52], and reading Object-Oriented code [27].

The reading techniques presented in this paper are classified under Reading for Construction, which is aimed at answering the question: Given an existing artifact, how do I understand how to use it as part of my new system? Reading for construction is important for comprehending what a system does, what capabilities exist and do not exist; it helps us abstract the important information in the system.

It is useful for maintenance as well as for building new systems from reusable components and architectures [4].

We chose to focus on the understanding of *object-oriented frameworks* as the artifact to be used in system development. An object-oriented framework (hereinafter, simply framework) is "a reusable design of all or part of a system that is represented by a set of abstract classes and the way their instances interact" [22, 24]. From the perspective of the application developer, it is a skeleton application that can be customized to produce specific applications in some domain [14, 24]. Some of the best-known frameworks support the development of graphical user interfaces (e.g. MacApp, ET++, Interviews, MFC and AWT). Frameworks are also spreading into other domains such as communication software [38], manufacturing systems [8, 37], and banking applications [5].

The choice to focus on frameworks was motivated primarily by two reasons:

1. Frameworks are a promising means of reuse. Although class libraries are often touted as an effective means of building new systems more cheaply or reliably, these libraries provide only functionality at a low level. This forces the developer to provide the interconnections both between classes from the library and between the library classes and the system being developed. Greater benefits, such as faster application development [28], are expected from reusable, domain specific frameworks that usefully encapsulate these interconnections themselves.
2. Frameworks have associated learning problems that affect their usefulness. The effort required to learn enough about the framework to begin coding is very high, especially for novices [32, 46]. Developing an application by using a framework is closer to maintaining an existing application than to developing a new application from scratch: in framework-based development, the static and dynamic structures must first be understood and then adapted to the specific requirements of the application. As in maintenance, for a developer unfamiliar with the system to obtain this understanding is a non-trivial task. Little work has yet been done on minimizing this learning curve.

Recognizing that one study cannot address issues for all types of frameworks, this paper concentrates on *white-box* frameworks. Frameworks of this type are tailored by deriving new classes through inheritance, and by writing

Forrest Shull and Victor Basili are with the Computer Science Department, University of Maryland, College Park, MD 20742. E-mail: {fshull, basili}@cs.umd.edu.

Filippo Lanubile is with the Dipartimento di Informatica, University of Bari, Via Orabona 4, 70126 Bari, Italy. E-mail: lanubile@di.uniba.it.

application-specific methods. Black-box frameworks, in contrast, provide a set of application-specific components that are plugged together through polymorphic composition [22]. It has been suggested that frameworks evolve towards black-box as the system design for their application domain becomes better understood [22]. Conclusions about white-box frameworks are of interest primarily for two reasons:

1. White-box frameworks are in wide use [14].
2. Not every framework inevitably evolves into a black-box framework. For example, some frameworks are retired from use before evolving to black-box stage; other frameworks are in application domains that are not understood to a sufficient degree to support black-box frameworks. [33]

II. RESEARCH QUESTIONS

Since we approached this study from the viewpoint of software reading, our primary focus was on the processes developers would engage in, as they attempted to discover enough information about the framework to be able to use it effectively. We reasoned that the best approach would be to observe a number of different approaches or techniques and their effects in practice. From this information, we hoped to determine what kinds of strategies could be used, and for which situations they were likely to be particularly well- or ill-suited. Ultimately we hoped to gain an understanding of the deeper principles involved in framework usage by studying the interaction between the different techniques and the specific task undertaken.

Since our study took place in the context of a classroom assignment, we felt it necessary to give our students a starting point for using frameworks. In the absence of any empirical evidence or general agreement in the literature on the best way to teach developers how to use a framework, we selected two promising approaches. Each approach was the basis for a set of guidelines that was taught to half of the class, so that the strengths and weaknesses of the approaches could be compared. (We discuss the guidelines themselves and the process of their creation in section V.) At the same time, we did not want to prevent the students from using work practices that they already knew, or discovered during the project, to be effective. Therefore, we allowed the students to modify these guidelines as desired. Our intention was to study the work practices of the students to determine when our guidelines were used, what other work practices were used, and how effective they were for particular tasks. We did not wish to constrain our subjects in any way to an artificial procedure, but to study and understand what they felt the most suitable approaches were for the problem. Our main research questions can be phrased as:

Can strategies for learning frameworks be identified?

What are their characteristics?

As it turned out, one of the learning approaches was viewed by subjects as too cumbersome for the environment of this

study, and was not used by any subject for the duration of the study. Although aspects of this approach were incorporated into the work practices of the subjects (these practices are described in some detail in section VII.A), a straightforward quantitative comparison of the results of using each approach was not possible. Instead, most of the results presented in this paper come from a quantitative and qualitative analysis of the student experiences with learning and using the framework over the course of the semester. This type of information is useful for giving us a deeper understanding of what is important in learning to use frameworks.

III. RELATED WORK

A survey of the literature on frameworks shows that relatively little has been written on *using* frameworks (as opposed to building or designing them). Most of the work on using and learning frameworks tends to concentrate on strategies for framework designers to use in documenting their work. The primary weaknesses of this approach are that, first, the results are only applicable to frameworks for which the prescribed documentation has been constructed (that is, they do not directly contribute to general guidelines that would help developers in approaching any framework) and secondly, that usually very little empirical evidence is presented in order to demonstrate that the prescribed method is as effective as claimed. We present some of the main areas of framework documentation here and discuss representative papers for each. We would like to reiterate that we in no way see our study as competitive with these works. Our aim is not to suggest that these other approaches are right or wrong, or to present an alternative approach which we argue to be superior. Rather, we aim to provide a more low-level indication of what sources of information or types of activities are important in framework use - information which may be used to help identify weaknesses in higher-level techniques and focus them on aspects of the framework which are most important.

1. **Patterns and recipes:** Beck and Johnson [6, 23] advocate the use of “patterns” (interlocking descriptions of problem/solution pairs, similar to object-oriented design patterns [17] or cookbook recipes, e.g. [1]) to describe frameworks. Each pattern describes a functionality supported by the framework, demonstrates how to implement the functionality, and discusses the impact of the implementation on the system. This seems a promising approach, because it seems capable of both showing the developer only as much detail as he or she needs for the current task and directing the developer’s attention to only the most relevant portions of the framework. Patterns have in fact been used as the sole form of documentation for the HotDraw framework. However, the only evidence presented as to their effectiveness is an informal study in which subjects were asked to learn HotDraw using patterns and provide feedback [23]. This study seems to have been very successful at its primary goal of

helping the patterns' authors debug their work, but does not provide much detail as to how the learning process was influenced. Thus the presentation leaves unanswered questions as to whether the observed effectiveness was a function of the specific project undertaken or would be true in any environment.

A related approach is the use of "hooks," which are meant to be similar to Beck and Johnson's patterns, although more structured, more uniform, and less narrative in style [16]. Like patterns, each hook provides only the information necessary to solve a specific, focused problem. They are produced by the framework developer to illustrate how the framework is intended to be used.

Johnson states [23] that it "would probably be worthwhile to try out the patterns in a controlled setting where it would be possible to watch how people use the patterns and what aspects of [the framework] are hard to learn." We feel that our study examines framework usage in exactly this way, although as we did not want to make *a priori* assumptions about the effectiveness of one method of documentation over another we did not work in an environment documented using patterns.

2. **Formal and/or searchable specifications of behavior:** Another tactic has been to formalize descriptions of the behavior of framework components, which then allows the creation of a search mechanism for finding useful components given a query. One such example is the prototype framework browser constructed at the University of Quebec [30], which is especially promising in that it concentrates on finding a general solution which can be applied to any existing framework, regardless of the level of documentation supplied. However, Gangopadhyay and Mitra point out two major stumbling blocks for query-based learning of frameworks: first, searching for and reusing one component at a time does not allow the potentially subtle connections between components to be understood, and second, it is a very difficult problem to match a query which has been specified in a way meaningful to the developer with the description of the framework components [18].
3. **Architectural approaches:** Gangopadhyay and Mitra recommend instead a top-down approach to learning frameworks, by which they mean a concentration on the framework architecture rather than on individual components [18]. They recommend the development of exemplars, executable visual models that consist of instances of concrete framework classes along with explicit representations of their collaborations. An exemplar should contain at least one concrete subclass for each abstract class in the framework. This approach might prove difficult to use for frameworks that do not conform to good design style issues, such as having only a few abstract classes, and using abstract

classes to implement important sites for customization in the framework. It is also unclear how helpful the exemplar approach would be in cases in which the developer wants to make a modification the framework designer has not anticipated (i.e. a modification at a location that is not represented as an abstract class in the framework).

4. **Tutorials:** Other work has focused on tutorials created for users to follow which will presumably guide users through the most important points of the framework. (Two examples are [47] for Unidraw and [15] for ET++.) An interesting example of work in this area is Rosson *et al.*'s tutorial for learning Smalltalk [34] which applies Minimalist instruction techniques [9] and seems to corroborate the benefits that may result from a well-designed tutorial course (claiming to allow new users to develop code for interactive applications after only four hours). Like Johnson, Rosson undertakes some testing which is aimed not at testing hypotheses but at helping to debug the documentation. However, no study has been undertaken to examine the breadth of knowledge achieved, although this becomes an exceptionally pertinent question when the goal is radical decreases in learning time for a constant breadth of knowledge.

An important weakness which is shared by all of these approaches is that they assume that the framework developer will be able to anticipate future uses of the framework adequately to provide enough patterns (or exemplars, or tutorial lessons) in sufficient detail. An alternate research approach would be to avoid making any such assumptions about framework usage in order to undertake an empirical study of how developers go about performing the necessary tasks. Such an approach is not new, and has in fact proven useful in understanding how developers perform related tasks such as understanding code [48] or performing maintenance [43]. Similar methods can be used to study the process of learning frameworks, since white-box framework understanding is a specialized kind of program understanding. A framework can be thought of as a set of object classes that collaborate to carry out a set of responsibilities; the developer needs to gain an understanding of what the various classes are and the functionality they provide. Since classes cannot be reused in isolation, it is also necessary to understand how these classes interact with each other. Understanding the dependencies between the components is a difficult task and the source of many complaints about framework complexity [24].

Example applications play a key role in the documentation of frameworks, by showing what the framework is good for and pointing out features that the framework provides. However, examples do not explicitly show how these features are provided [23]. The problem remains that it is difficult to understand the interactions between objects using source code.

Other authors [11] who have applied an empirical approach to studying framework usage in industrial environments agree that, in most cases, framework customization will be more complex than just making modifications at a limited number of predefined spots. Documentation that assumes this is possible will be too constraining and will not provide support for many realistic development problems, which far from requiring isolated changes may sometimes even require changes to the underlying framework architecture.

Our study belongs in this category of empirical study of practical framework use. It is similar in type to the study undertaken by Schneider and Repenning [39], which draws conclusions about the process of software development with frameworks from fifty application-building efforts supervised by the authors. The large number of projects followed allowed the authors to examine both successful and unsuccessful projects, and their observation of (and sometime participation in) the process allowed them to both characterize the usual process and to identify some conditions that contribute to situations where the process breaks down and leads to unsuccessful projects. Our results complement and in some cases extend the results from the Schneider and Repenning study, and we consequently discuss them in greater detail later in section VIII.E.

IV. DEFINITIONS

We include the following definitions to clarify and illustrate some terminology that we use often in this discussion. It is our hope that these definitions will help to make clear our model of framework usage and to keep certain concepts distinct throughout the discussion. For example, we should first be careful to differentiate the framework developer (the developer(s) who designed and implemented the framework) from the application developer (the developer(s) who design and implement a new system using the framework to provide certain key functionality), who is usually referred to as simply the “developer” in this discussion.

example application: an application which has been constructed using the framework. Such examples may be created by the framework developer to illustrate how to produce some functionality using the components provided by the framework, or may be a “real-life” application in whose development process the framework happened to be used.

functionality supported by framework: a particular functionality is either provided by an example application, or there is a one-to-one mapping between the functionality and a framework component at some level of granularity (e.g. subsystem, class, method, ...). An example might be

the behavior of radio buttons (primitive GUI objects would be likely to be supported by a class in a GUI framework) or a linked list (most frameworks provide classes that encapsulate reusable abstract data types).

functionality provided by examples: an example application exhibits dynamic behavior which corresponds to the required functionality. The functionality may be implemented in one component or a combination of components which contain some code specific to the example (i.e. some but not all of the functionality may be inherited from the framework).

functionality supported by object model: the required functionality can be implemented as part of the system represented by the object model without too much change to the object model. Obviously the phrase “too much change” is unacceptably vague, but we do not yet have a useful way of characterizing the concept of when adjustment to an object model becomes excessive. Although there are some high-level studies of software architecture underway, we look upon the effort to provide guidance as to how much adjustment is possible in a given situation as an important and promising area of future research.

V. INITIAL DESIGN OF FRAMEWORK LEARNING TECHNIQUES

In the absence of a definitive approach to framework learning, we turned to the literature to identify useful approaches. Our first step was to identify helpful models of the framework, that is, helpful ways of thinking about the framework that would highlight the truly important features and could be used for finding particular functionality. The most common description of a framework uses a class hierarchy to describe the functionality supported by the framework and an object model to describe how the dynamic behavior is implemented. Most of the common descriptions of a framework in the literature (e.g. [28, 46]) present a model of the framework similar to this one. To teach subjects how to use this model, we created a set of guidelines that could be used to gain an understanding of the class hierarchy. The guidelines help developers understand the functionality provided by the framework by concentrating on abstract classes in the hierarchy. The developer is guided first through the broad classes of functionality, then through deeper and deeper levels of concrete classes to find the most specific instantiation. We refer to this procedure as the *Hierarchy-Based* (HB) procedure, to emphasize that the underlying model of the framework is the class hierarchy.

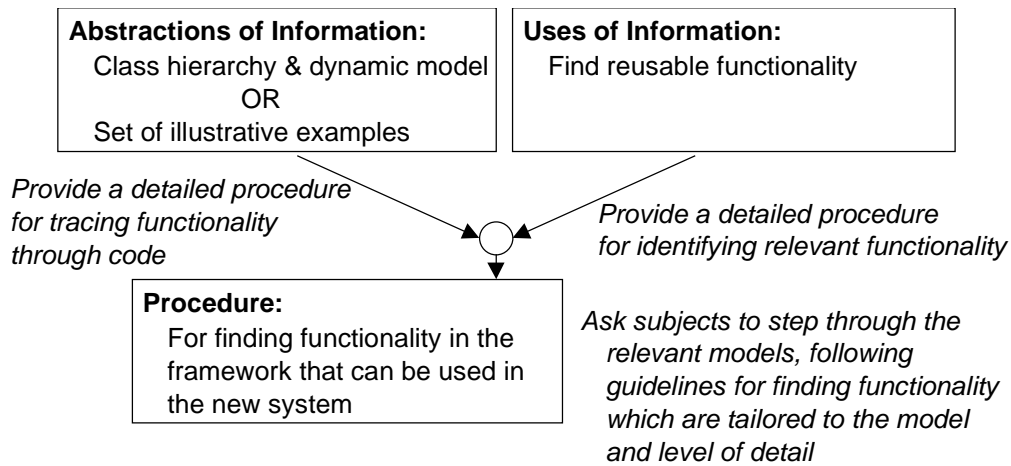


Figure 1: Producing focused, tailored procedures

As an alternative model, we decided to look at the framework through a set of example applications which, taken together, were meant to illustrate the range of functionality and behavior provided by the framework. Although a detailed examination of learning frameworks by means of examples has not been undertaken, learning by example also seemed a promising approach. Sets of example applications have been used to document some frameworks (the framework we used came with such a set) and the approach has been recommended for similar types of activities, such as learning effective techniques for problem solving [10], or learning how to write programs in a new programming language [26, 34]. It has also been argued that learning by examples is well-suited for “domains where multiple organizational principles and irregularities in interaction exist” [7], which may be a fair assessment of the large hierarchy of classes in a framework.

The framework we used in this study came with a set of examples at varying levels of complexity that was constructed to demonstrate the important concepts of the framework. To help subjects use these examples to learn the framework we created a set of guidelines that would guide exploration through the example set, to particular examples, to particular objects in the implementation, to particular lines of code in the object. This procedure is referred to as the *Example-Based* (EB) procedure.

Once we had identified suitable models, we constructed detailed guidelines for identifying functionality that would be relevant to the system being developed. To do this we concentrated on identifying similarities between the classes that were specified by subjects in their original object models, and the classes in the framework. These guidelines were then also tailored to the models, and integrated with the procedures for understanding the framework (see figure 1). The final guidelines were intended to be step-by-step procedures that could be taught to the students and used to find functionality in the framework.

VI. DESCRIPTION OF THE STUDY

To undertake an exploratory analysis into framework usage, we ran a study as part of a software engineering course at the University of Maryland. Our class of 43 upper-level undergraduates and graduate students was divided into 15 two- and three-person teams. Teams were chosen randomly and then examined to make certain that each team met certain minimum requirements (e.g. no more than one person on the team with low C++ experience) for the class. Each team was asked to develop an application during the course of the semester, going through all stages of the software lifecycle (interpreting customer requirements into object and dynamic models, then implementing the system based on these models). The application to be developed was one that would allow a user to edit OMT-notation diagrams [36]. That is, the user had to be able to graphically represent the classes and different types of relations between them of a system, to be able to perform some operations (e.g. moving, resizing) directly on these representations, and to be able to enter descriptive attributes (class operations, object names, multiplicity of relations, etc.) that would be displayed according to the notational standards. The project was to be built on top of the ET++ framework [50], which assists the development of GUI-based applications. ET++ provides a hierarchy of over 200 classes that provide windowing functionality such as event handling, menu bars, dialog boxes, and the like. ET++ was considered an attractive choice for this study because it possesses an “outstanding design” and is available from the public domain [32]. The design of ET++ is of sufficient quality that it was the source of seventeen of the design patterns in [17].

Before implementation began on the project, the class was randomly divided into two groups and each was taught only one of the two framework models and its corresponding guidelines for use. During the implementation, we then monitored the activities undertaken by the students as much as possible in order to understand if our learning techniques were used, what other learning approaches were applied, and which of these were effective. To do this, we asked the

students to provide records of the activities they undertook so that we could monitor their progress and the effectiveness of the techniques, and augmented this with interviews after the project was completed. (These data collection mechanisms are discussed below.) Although we felt our learning techniques made good starting points, we did not constrain students to use the techniques exactly as they were taught. We wanted to leave the students the flexibility to modify their development methods if the techniques were not well suited to a particular part of the implementation, or if they had personal techniques that would work better for them.

Since the analysis was carried out both for individuals and the teams of which they were part, we were able to treat the study as an embedded case study [51]. Over the course of the semester, we used a number of different methods to collect a wide variety of data, each of which we discuss briefly below. Most of our collection methods are mentioned by Singer and Lethbridge in their discussion of the pros and cons of various methods for studying maintenance activities [43], and we respond to some of their comments where appropriate. We hope this study provides an additional illustration of their conclusion that, in order to obtain an accurate picture of the work involved, a variety of methods must be used at different points of the development cycle in order to balance out the advantages and disadvantages of each method.

1. Questionnaires were used at the beginning (to report previous programming experience) and end (to report effort spent during the last week of implementation and the level of completion for each functional requirement for the project) of the semester. Although the information reported on the beginning questionnaires could not be verified, the end questionnaires were verified against the executables submitted for each team. The unit of analysis for the beginning questionnaires was the individual student, while the end questionnaires were filled out for the entire team. Both were mandatory although self-reported, and did not impact the students' grades.
2. Exam grades were recorded for certain questions on the midterm that could be used by us to gauge the students' level of understanding of framework concepts. These grades were recorded for the individual students, were assigned by us after evaluating the students' responses for the level of understanding exhibited, and were mandatory (as they constituted part of the students' grades for the course).
3. Progress reports were to be submitted by each team for each week of the implementation phase. They consisted of an estimate of the number of hours worked by the team for the week in implementing the project, and a list of which functional requirements had been begun and which completed. As the students were told that the progress reports had no bearing on their grades, many teams opted to submit them only sporadically or not at all. (In some ways, these reports were similar to Singer and Lethbridge's idea of logbooks, which allow the developer to record information at certain times throughout the development process. Singer and Lethbridge concentrate on the dangers of making the report too time-consuming, but we have noticed a quite opposite phenomenon: if the experimenter makes the report too minimal, the developer may assume that the information to be collected cannot be truly important and thus make completing the report a very low priority.)
4. Problem reports were requests for clarification or for help with ET++ that the students submitted (via email) to the course instructors. A record was kept of the general subject of each request, and by which team it had been submitted. In this way we hoped to maintain a record of the kinds of difficulties encountered by teams during the course of the project. Problem reports were obviously not mandatory and had no effect on student grades, but were a resource that the students knew could be made use of at their discretion. (Singer and Lethbridge focus on the inaccuracies of retrospective reports, but our problem reports were actually an excellent way to get an accurate picture of where teams were having problems at the time they were having them - which may be, admittedly, unique to the classroom environment.)
5. Implementation score was assigned by us to each team at the end of the semester. Projects were graded by assessing how well the submitted system met each of the original functional requirements (on a 6 point scale based upon the suggested scale for reporting run-time defects in the NASA Software Engineering Laboratory [45]: "required functionality missing", "program stops when functionality invoked", "functionality cannot be used", "functionality can only partly be used", "minor or cosmetic deviation", "functionality works well"). The score for each functional requirement was then weighted by the subjective importance of the requirement (assigned by us) and used to compute an implementation score that reflects the usefulness and reliability of the delivered system. The weights were chosen in such a way that if each functionality worked well, an implementation score of 100 would be obtained. Scores less than 100 provided a rough indication of what percentage of system functionality had been implemented. (Because extra credit was awarded in rare instances that functionality beyond what was required was implemented, it was also possible for implementation scores to be slightly greater than 100.)
6. Final reports were collected from each team at the end of the semester. These reports consisted of documentation for the submitted system (object models and use cases) as well as records of the activities undertaken while implementing the project (object models of examples that had been studied, lists of classes that had been examined for some functionalities). Additionally, in-class presentations were given by each team in which they could present interesting details of the functionality available in their system, their experiences and difficulties with the techniques they used, and/or their general approach to

implementation. The completeness of the final reports counted toward each team's grade, although their conformance to any particular technique did not.

7. Self-assessments were mandatory ratings in which each student was asked to rate the effectiveness of each member of his or her team (including him- or herself) as well as the team performance as a whole. Partly this was to detect if every team member had done their share of the work, and partly it was to ask students to think about what they had done rightly and wrongly during the course of the implementation. Although it was mandatory that each student return a self-assessment, they did not count directly toward the student grade (although in some cases, evidence from the self-assessments and the interviews led to individual grades being slightly adjusted).
8. Interviews were mandatory "debriefing" sessions at the end of the semester. Each team would come as a group to the course instructors, to be asked questions about what kinds of activities they did during the course of the semester, which of these they found particularly useful or useless, and what parts of the project were easiest and hardest. A set of base questions was asked in every interview, although additional questions were conducted in a dynamic manner. That is, the course of the interview was directed in new directions by us as unforeseen but interesting themes were raised.

Table 1 summarizes the types of data we collected along with the collection methods we used. In the interest of space, we do not present the actual data in this paper, but maintain a website that contains them in a form suitable for downloading [42].

VII. ANALYSIS

We then analyzed this mix of qualitative and quantitative data to gain some insight into what was going on within each team. By comparing and contrasting teams, we began to see implications that addressed our research questions. Since there has not yet been a large amount of work spent on understanding this area of framework use, our focus was on using this information to look for tentative but reasonable hypotheses and not on testing known hypotheses. The process of building theories from empirical research has been first proposed in the social science literature [13, 19] but it is also followed in the software engineering discipline [40].

A. Development Processes

The analysis approach we used was primarily a mix of qualitative and quantitative, in order to understand in detail the development strategies our subjects undertook. Our first step was to get an overview of what development processes teams had used. (By "development processes" we mean how the team had been organized, what techniques they had used to understand the framework and implement the functionality, whether they based their implementation on an example or started from scratch, and what tools they had used to support their techniques.) To

this end, we performed a qualitative analysis of the explanations given by members of the teams during the interviews and final reports, and on the self-assessments. We first focused on understanding what went on within each of the teams during the implementation of the project. We identified important concepts by classifying the subjects' comments under progressively more abstract themes, then looked for themes that might be related to one another (an example is given below). Once we felt we had a good understanding of what teams did, we made comparisons across groups to begin to hypothesize what the relevant variables were in general. This allowed us to look for variations in team effectiveness that might be the result of differences in those key variables, as well as to rule out confounding factors.

In order to illustrate this analysis technique better, let us consider a small example. While trying to categorize the types of remarks students made during the final interviews, we noticed a lot of comments (some made spontaneously, some with prompting by the interviewers) concerning how teams spent most of their time during the implementation phase of the project. We grouped some of these remarks into a general category that deals with what kinds of activities students found useful in implementation; combined with other categories (which deal with, for example, what kinds of tools students found useful or what kinds of examples were helpful to examine) we began to get a better idea of what techniques students developed to help them in implementation. To get an accurate picture of what teams did over the entire course of the semester, however, we needed to look at other categories of remarks, concerning for example which of the two initial procedures (HB or EB) the team was taught, what their experiences were with the technique, and what parts of the technique were found not to be useful and were discarded. By making such abstractions from the students' comments, we built an understanding of what each separate team did.

With this understanding of the processes at work within teams, we compared experiences across teams to try to identify themes that emerge. For example, we noticed that most teams who began by modifying an example tended to do better than teams who began implementing from scratch. Our next step of the analysis was to test these provisional hypotheses, again by making comparisons across teams to find possible refuting evidence or confounding factors. For example, suppose Team X started their implementation from an example but turned in a very poor implementation in the end. Does this refute our provisional hypothesis? Perhaps it signifies that the trend we thought we had observed was simply a fluke, or perhaps we may notice a confounding factor - say, Team X was also very poorly organized - that may account for the seemingly anomalous results and needs to be taken into account as part of the analysis.

Aspect of Interest	Measures	Form of Data	Unit of Analysis	Collection Methods
Development Processes	Techniques used	Qualitative	team	interviews, final reports
	Tools used	Qualitative	team	interviews
	Team organization	Qualitative	team	interviews, self-assessments
	Starting point for implementation	Qualitative	team	interviews, final reports
	Difficulties encountered with technique	Qualitative	team	problem reports, self assessments, final reports
Product	Degree of implementation for each functionality	Quantitative	team	implementation score, final reports
Other Factors Influencing Effectiveness	Effort	Quantitative	team	progress reports, questionnaires
	Level of understanding of technique taught	Quantitative	individual	exam grades
	Previous experience	Quantitative	individual	questionnaires

Table 1: Types of measurements and means for collecting.

Category	Number of Teams	Number Originally Taught EB, HB	Description
EB	5	5, 0	Students in this category used the EB technique as it was taught, following the guidelines as closely as they could.
EB/HB	5	0, 5	This is a hybrid approach that focuses on using examples to identify classes important in the implementation of a particular functionality. The main difference from the EB technique is that, in the hybrid technique, the student does not always begin tracing the functionality through the code, but may instead use the example to suggest important classes and then return to the framework hierarchy to focus on learning related classes.
ad hoc EB	4	2, 2	This was an ad hoc approach that emphasized the importance of learning the framework via examples, but ignored the detailed guidelines given. The primary difference between these techniques and the EB category is that ad hoc EB techniques are missing a consistent mechanism for selecting examples and tracing code through them.
EB/scratch	1	1, 0	The team used the EB technique to identify basic structure and useful classes, but implemented the functionality mostly from scratch.

Table 2: Description of development processes observed in the study.

Although this is not a common method of analysis in computer science, it is a recommended approach for social sciences and other fields that require the analysis of human behavior [13, 29]. It is well suited for our purposes here because our variables of interest are heavily influenced by human behavior and because we are not attempting to prove hypotheses about framework usage, but rather to begin formulating hypotheses about this process, about which we currently know little.

We found that teams used development processes that fell into 1 of 4 categories (Table 2). While no team used the HB technique for the entire semester (although these guidelines were partly incorporated into some of the other techniques that were used), the EB technique did enjoy consistent use, both as taught and in combination with other techniques. It can be noted that even teams who were taught HB and not exposed to EB tended to reinvent a technique similar to EB on their own. (Some teams were even contrite about this. “We didn’t realize at the time that this was the technique taught to the other part of the class, but it seemed the natural thing to do.”)

B. Potentially Confounding Factors

We also undertook quantitative analyses to gauge the effects of potentially confounding factors in the study. Due to the small sample size and the exploratory nature of this study, we used an α -level of 0.20 for the statistical tests reported in this paper, which is higher than standard levels. Although not common, this α -level has been used in similar hypothesis-building studies, e.g. [2]. We realize that statistical tests at this significance level do not provide strong evidence of a relationship, but instead see their contribution as helping detect patterns in the data that can be specifically tested in future studies. The relatively high α -level makes it more likely that we err in the direction of finding false correlations than of inadvertently missing significant relationships. However, the low number of data points in this study means that the second case is still a possibility.

We noticed that teams that had been taught the HB technique experienced many early difficulties, particularly in tracking flow of control in the framework, and in finding and correctly parameterizing reusable classes from the hierarchy. We wanted to test if the progress of these teams had been adversely affected by the reading technique. We therefore undertook an analysis of the number of person-hours spent by different groups in order to understand if the amount of effort a team spent on implementation was largely dependent upon the technique they had been taught. The test for differences in the two groups could not be conclusive, as only six teams reported their overall effort data. The analysis did, however, show no significant difference between the average amount of effort spent on the project over the course of the semester by teams who had been taught each of the different techniques (t-value of 0.5916 and p-value of 0.5859). Student remarks from the interviews tended to support this. Most teams who had

been taught the HB technique reported that they switched their development approach usually after trying to apply HB for the first 1 to 2 weeks. Regardless of the technique a team had been taught, the heaviest investments of effort for almost all teams came toward the end of the implementation phase.

We also wished to examine whether students actually had a different initial understanding of the framework based on the models and procedures we had taught them. We gauged their initial understanding by means of two questions that appeared on their midterm, after they had been taught one or the other of the framework models, but before they had actually had any experience using the ET++ framework. The first question was intended to measure how well the students grasped the concept of the framework hierarchy of classes (Table 3 shows the distribution of grades on this question by procedure taught). The second question measured understanding of the model of interaction (Table 4). We tested whether response rates were independent of the procedure taught by using Wilcoxon rank sum tests. The results of neither test were significant ($Z = -0.6933$, $p = 0.4881$ and $Z = 0.6343$, $p = 0.5259$ for the first and second questions, respectively) showing that there is no difference in how well the questions are answered with respect to the technique taught. This shows that, although taught different procedures for using the framework, neither group of subjects started out at a disadvantage to the other in terms of their understanding of the framework itself.

		Grade Achieved			
		A	B	C	D
Technique Taught	EB	5	5	5	6
	HB	9	3	4	6

Table 3: Distribution of grades on exam question dealing with the framework class hierarchy.

		Grade Achieved			
		A	B	C	D
Technique Taught	EB	10	3	2	6
	HB	6	7	4	5

Table 4: Distribution of grades on exam question dealing with the framework model of interaction.

A final concern in all studies of this type has to do with the experience of the subjects. We were concerned that the effectiveness of our teams might have more to do with the level of experience the team members had with implementing similar projects than with any of the variables under study in our experiment. We removed one team (which had had extreme organizational difficulties) from our analysis as an extreme outlier (as defined by [31]). We then used the Pearson correlation coefficient [21] to measure the strength of the linear relationship between experience and implementation score (with correlation values close to 1 or -1 representing an exact linear relationship and values close to zero representing no linear

relationship). We found no correlation between the total amount of experience of a team with programming in an academic environment and its effectiveness at the implementation of the project (Pearson's correlation coefficient of -0.0015), but did uncover a correlation between total experience programming in industry and effectiveness at implementation (Pearson's correlation coefficient of 0.55 , Figure 2). However, an R^2 value of 0.31 for the model means that industrial experience accounted for only 31% of the observed variation in the implementation score. This low correlation implies that the previous experience of our subjects is not sufficient to explain the observed results, and thus that the way in which this implementation was undertaken has affected the quality of implementation.

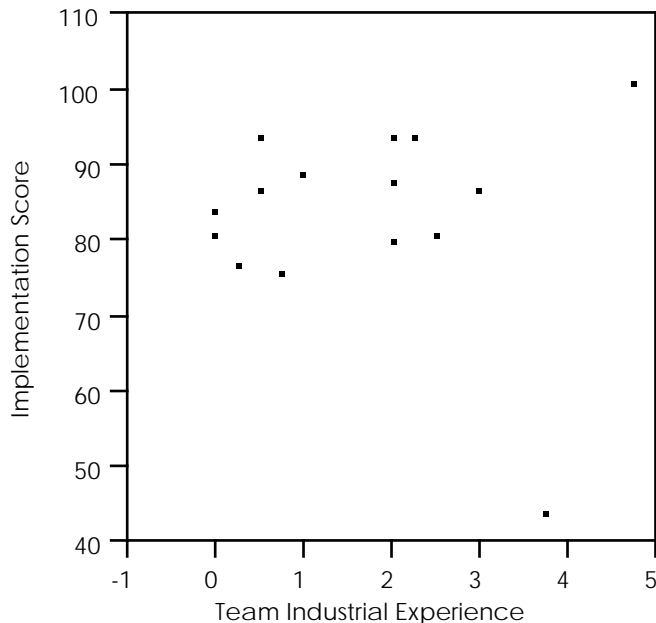


Figure 2: The total industrial experience of a team (in years) and its correlation with its effectiveness at implementation of the project (the team with implementation score of 44 has been removed from analysis as an outlier).

VIII. RESULTS

In this section we present the hypotheses that result from our study of framework usage. Along with each we present the relevant indications from our study, so that the reader may judge the strength of the evidence which supports these hypotheses.

A. Hypotheses About Our Models of the Framework

In order to understand the process of learning a framework by means of examples, we can generalize not only from the EB procedure itself but from its example-based derivatives as well. From Table 1, it can be seen that all students who were taught EB ended up using this technique (or a close derivative) throughout the implementation phase. Perhaps

more importantly, all students who were taught the other technique ended up employing a less rigorous example-based technique on their own. It seemed, therefore, that not only our EB technique, but example-based learning in general, was a natural way to approach learning such a complicated new system. This leads us to formulate:

HYPOTHESIS 1: Example-based techniques are well-suited to use by beginning learners.

In contrast, subjects tried to use the HB procedure but eventually abandoned it. Qualitative analysis was necessary to understand why this happened. What was wrong with the HB procedure that made it not useful in this environment? We analyzed the student remarks from the problem reports, self-assessments, and final reports to see if characteristic difficulties had been reported. A common theme (remarked upon by half of the teams who had been taught HB) was that the technique gave subjects no idea which piece of functionality provided the best starting place for implementation, or where in the massive framework hierarchy to begin looking for such functionality.

Teams also registered complaints about the time-consuming nature of the HB technique - especially compared to an example-based approach where implementation can begin much more rapidly, hierarchy-focused approaches seem to require a much larger investment of effort before any payoff is observed. One team pointed out that they had explicitly compared their progress against other (as it happened, Example-Based) teams: "We talked to other groups, and they seemed to be getting done faster with examples. So after the first week we started going to examples, too."

Despite these difficulties, students reported that they felt that the HB technique would have been very effective if they had had both sufficient documentation to support it and more time to use it. "The Hierarchy-Based procedure would be helpful if you have the *time*," said one group in the final interviews, "but on a tight schedule it doesn't help at all." Another opinion was expressed by the group who said, "It's the technique I normally use anyway - and it would have been especially good here when the examples are not enough for implementing the functionality." There seemed to be a consensus that it would have allowed them to escape from the limitations of the example-based approach and engage in greater customization of the resulting system, but simply wasn't effective in the current environment. Five teams were able to create effective strategies that were hybrids of the hierarchy-based and example-based methods (EB/HB). However, the lack of guidance as to how to get started, and the time required to learn the necessary information to use it effectively, meant that no development teams used it exclusively for a significant portion of the implementation phase. (By no means was this a completely negative development, as we now have more detail on techniques that minimize that crucial learning curve.)

HYPOTHESIS 2: A hierarchy-focused technique is not well-suited to use by beginners under a tight schedule.

B. Practical Implications for Using an Example-Based Procedure

The analysis in the last section should not be taken to imply that our EB procedure was without problems. We also undertook a qualitative analysis of subject satisfaction with the Example-Based procedure in order to understand where it could be strengthened for future use. We found that there were also characteristic problems with EB that were encountered by beginners.

Although the teams using this technique usually managed to get the functionality working in the end, almost all (4 out of 5) of the groups in the study who used our EB technique reported difficulties in finding the necessary functionality within the example set. The problems with finding it in the first place seemed especially acute when the functionality needed was a very small part of a much larger example (e.g., a characteristic way of displaying items onscreen, or the dialog boxes discussed under “key functionalities”, below). This indicates a problem with our EB technique – further guidance is required to assist developers in finding and extracting small pieces of functionality embedded within larger examples. As there is currently little guidance available in this regard, more work will have to be done in this area to enable effective example-based techniques.

This study also provides indications that characteristics of the example set can influence the performance of an example-based approach as well. One-fourth of all the teams in our study had trouble making use of the examples because the example set provided did not conform to a consistent organization or structure. Some examples were based on Model-View-Controller interaction [20] while others were not constrained by any such separation of functionality, and different examples seemed to achieve the same functionality by using different classes from the framework hierarchy. As others [35] have pointed out, learning how to implement functionality from existing applications is difficult because the rationales for design choices, which explain why the finished implementation looks the way it does, are usually not included in the documentation. When attempting to reuse functionality from existing applications, developers are implicitly asked to reconstruct the choices that led to the finished implementations they are studying. This situation can actually be made worse in a framework-based environment, where effective reuse requires the developer to understand the rationales behind a number of applications, not just one. This is a problem that will have to be addressed by any example-based technique.

HYPOTHESIS 3: The effectiveness of a technique for adapting framework-based applications depends on breadth of functionality and other characteristics of the existing applications

This hypothesis is very high-level, and further work is required to understand better the effect of characteristics of the example set. However, hypotheses 4 through 6 in the next section are concrete applications that may demonstrate the usefulness of this hypothesis.

C. Hypotheses About the Level of Specificity in the Procedures

This study provides us with an excellent opportunity to understand whether the procedure we created was at a useful level of specificity. Because some teams followed the procedure exactly while others followed it only to a certain extent, we can distinguish between the following two types of example-based processes:

- **Strictly adopted:** The EB technique is considered “strictly adopted” when followed in a step-by-step fashion. This procedure focuses entirely on guiding the developer to find useful functionality in the example applications, which can then be tailored to the current system.
- **Ad hoc adopted:** The EB/HB, ad hoc EB, and EB/scratch techniques are considered to fall in this general category. Subjects who used these techniques are still considered to have adopted the basic approach behind the EB technique (viz. learning from examples). However, the subjects have augmented this basic approach to a greater or lesser degree with other techniques that were found to be effective.

We can compare the experiences of strictly adopted to *ad hoc* adopted teams to understand if our technique, at its current level of detail, is useful.

To perform this analysis, we focused on certain *key functionalities*, that is, certain requirements for which there was a large degree of variation between teams in terms of the quality of the implementation. Recall that we had graded each required piece of functionality for each project on a 6-point scale. To select key functionalities we looked for functionalities for which there was enough variation among all 15 of the teams on the rating scale (regardless of what type of technique they used) that teams could be easily divided into at least two groups based on their score. We then analyzed whether the level of detail at which the procedure was followed had any correlation with whether the team was able to implement a particular functionality in a more complete or sophisticated way.

To back up this qualitative analysis, we attempted to use statistical tests to verify the correlation. Since both of the variables in this analysis (technique followed and result from implementation) are on a nominal scale (i.e. expressed as categories) we can organize data into contingency tables where each dimension corresponds to a variable and numbers represent frequencies. We want to test whether the proportion of teams in each type of implementation varies due to the type of technique that was used. Specifically, the null hypothesis is that the proportion of teams who achieved each type of implementation is independent of the technique that was used to perform the

implementation. In order to test this difference we apply the chi-square test of probability¹. Due to our small sample sizes and the exploratory nature of this study, we again used an α -level of 0.20. We also present the product moment correlation coefficient, r , as a measure of the effect size [25]. (An r -value of 0 would show no correlation between the variables, whereas a value of 1 shows a perfect correlation.) We realize that these tests do not provide strong statistical evidence of any relationship, but instead see their contribution as helping detect patterns in the data that can be specifically tested in future studies.

We identified 4 such key functionalities: links, dialog boxes, deletion, and multiple views. They illustrate 3 types of situations that may arise when functionality is being sought in examples:

1. **The examples don't provide all of the functionality desired.** Key functionality 1 (links) fits into this category. OMT notation requires that links be drawn between classes in the model which interact in some way with each other. The specifications for our OMT editor stated certain other specifications for handling links, such as how they should be entered into the diagram, and that they should automatically update when their associated classes are moved. Provided with the example set was the "er" entity-relation diagram editor which provided similar functionality that allowed two linked objects to be represented by means of a line that connected their centers. Although useful as a starting point, this implementation was not sophisticated enough for the project, because the same two classes in an OMT diagram may be connected by multiple links. If these are all represented by lines drawn from center to center, every link between these classes will overlap. There is no example provided with ET++ that shows functionality that explicitly addresses this concern.

About half of the teams implemented the link functionality as in the er example. Of the teams that implemented a more sophisticated version that allowed links to be uniquely represented (i.e. multiple links between two classes do not overlap), there were a number of solutions: randomly displacing links by a small offset, recording the point where the user drew the link across the boundary of the class, and breaking the line down into a series of line segments which may then be individually positioned.

Almost all (4/5) of the teams who used the EB technique implemented the less sophisticated version of the functionality found in the er example. By comparison, less than half (4/10) of the teams who used a modified version of EB turned in the less sophisticated implementation (Table 5). A chi-square test of independence was undertaken to test whether

¹ We base our use of the chi-square test, rather than the adjusted chi-square test, on [12], which argues that the adjusted test "tends to be overly conservative."

the level of sophistication was dependent on the type of technique used, although we recognize that the small number of data points involved can lead to some inaccuracies in the results [31]. The test resulted in a p-value of 0.143 ($\chi^2 = 2.143$), which is statistically significant at the selected α -level. An r -value of 0.38 confirms that this shows a moderate correlation between level of sophistication and type of technique [21]. From this example we hypothesize that:

HYPOTHESIS 4: A detailed Example-Based procedure can cause developers to not go beyond the functionality that is to be found in the example set.

		Technique Followed	
		strictly adapt.	<i>ad hoc</i> adapt.
Result of Implementation	Sophisticated	1	6
	Simple	4	4

Table 5: Number of teams achieving sophisticated versus simple implementation of links, whether using strictly or *ad hoc* adopted techniques.

2. **The functionality was completely contained in (perhaps multiple) examples.** Key functionalities 2 (dialog boxes) and 3 (deletion) provide evidence that the EB technique performed about the same as the variant techniques in this case.

For dialog boxes, examples existed which showed how to create dialog boxes containing graphical devices (e.g. text fields, radio buttons) and how to use them to display and store information. The difficulty was that this functionality was spread piecemeal over multiple examples and students had a hard time finding and integrating all of the functionality they needed. About half of the class (7/15) managed to get the dialog box functionality working correctly and interfaced with the rest of the system (Table 6). All techniques were distributed roughly equally between the teams who did and did not get this functionality working correctly. The chi-square test here yielded a p-value of 0.714 ($\chi^2 = 0.134$), for which the related r -value is 0.10. This confirms that response levels are very likely effectively equal between the two categories.

For deletion, there was at least one example that clearly contained functionality to delete classes and links from a diagram. All teams were able to implement this functionality. Getting the functionality to support the ability to undo or redo a deletion was apparently more challenging, however, although the examples covered this as well. Partly, this may have been due to students simply forgetting to implement this part of the functionality since it was not explicitly mentioned in the requirements (although adequate testing should have revealed the need to interface correctly with this standard ET++ functionality!).

Teams basically implemented deletion in one of three ways:

1. The ability to undo or redo was integrated fully with deletion.
2. The ability to undo or redo was not implemented for deletion, but deletion was implemented in such a way that a later call to undo or redo would not cause a core dump (this *minimally* satisfied the requirements for the project, in letter if not in spirit!).
3. The ability to undo or redo was not implemented at all for deletion; deletion could be used but a later use of undo or redo could cause a core dump if the program tried to manipulate an object that had been deleted.

Again, the different techniques for learning frameworks seem pretty equally distributed among these three categories (Table 7). The chi-square test for this functionality yielded a p-value of 1 ($\chi^2 = 0.000$), with an associated *r*-value of 0, which is not significant and indicates that response rates are exactly equal regardless of the type of technique used.

From these two examples we hypothesize

HYPOTHESIS 5: When the functionality sought is contained in the example set, Example-Based techniques will perform about the same, regardless of the level of detail provided.

		Technique Followed	
		strictly adapt.	<i>ad hoc</i> adapt.
Result of Implementation	Fully correct	2	5
	Incomplete	3	5

Table 6: Number of teams who achieved fully correct versus incomplete implementations of dialog boxes, whether using strictly or *ad hoc* adopted techniques.

		Technique Followed	
		strictly adapt.	<i>ad hoc</i> adapt.
Result of Implementation	Impl. 1	1	2
	Impl. 2	2	4
	Impl. 3	2	4

Table 7: Number of teams achieving each of the three levels of implementation for deletion, whether using strictly or *ad hoc* adopted techniques.

3. **The examples provide a more sophisticated implementation than is required.** Key functionality 4 (views) fits here. The requirements for the OMT editor stated that the program must provide multiple views of the currently opened document. Examples existed which satisfied the project’s requirements about views (that they update automatically, be independently scrollable, and allow resizing).

However, there were also examples that gave an even more sophisticated implementation that allowed views to be dynamically added and deleted.

All but three teams chose the more sophisticated implementation. Two of these three teams turning in less functionality used the EB technique (Table 8). The chi-square test resulted in a p-value of 0.171 ($\chi^2 = 1.875$), which is significant at the selected α -level. An *r*-value of 0.35 shows a moderate correlation between the variables and allows us to hypothesize:

HYPOTHESIS 6: When the example set contains functionality beyond what is required for the system, a sufficiently detailed Example-Based procedure can help focus developers on just what is necessary.

		Technique Followed	
		strictly adapt.	<i>ad hoc</i> adapt.
Result of Implementation	Sophisticated	3	9
	Simple	2	1

Table 8: Number of teams achieving sophisticated versus simple implementation of views, whether using strictly or *ad hoc* adopted techniques.

The practical consequences of this analysis may be that the level of detail that is appropriate in an Example-Based procedure is strongly dependent on the breadth of the example set provided. We hypothesize that the more detailed a procedure is, the more it focuses the developer on using only functionality provided by the examples. This may be because developers become too reliant on the examples and do not understand the system at a sufficient level of detail to implement effectively from scratch when necessary. Alternatively, it may be that integrating functionality written from scratch becomes more and more difficult when more and more of the system is taken from examples that someone else has written.

Of course, there are other benefits to providing additional detail in such procedures. Among the most important may be packaging experience to guide new developers. For example, in our study we noticed that half of the teams who used a variant of the EB procedure wasted time and effort during the course of the implementation phase by having to re-implement some functionality that they had implemented previously in a short-sighted way. Since only 1 of the 5 teams who used our EB procedure exactly reported the same difficulty, we feel that a definite benefit of the procedure was that it helped guide our subjects to focus more on features which were important to the implementation as a whole, rather than just the portion they happened to be working on at a particular time.

D. Hypotheses About Beginning Implementation

Because we did not constrain the students as to how they went about implementing the functionality, we could also study whether the way in which the implementation was begun had an effect on the rest of the implementation process. Some teams decided to start from scratch, beginning with no functionality on top of the framework and slowly growing the code as they determined how to implement specific functionalities. Other teams chose to start with an example which seemed to contain some of the functionality they would need in the finished system. They then modified the functionality that was there and added whatever else was required to produce the finished system. We wanted to see if one approach tended to get more functionality working more reliably (as measured by the team's implementation score).

All teams who started from an example chose the same one, a simple entity-relation diagram editor (known as "er"). It was similar to the OMT editor to be developed, but much simpler: as specified for the OMT editor, the ER diagram editor allowed simple shapes to be added to an editable document, and allowed these shapes to be selected, moved, and associated with one another. However, more sophisticated functionality, such as entry of attributes via dialog boxes and the maintenance of separate regions within a particular shape, was lacking.

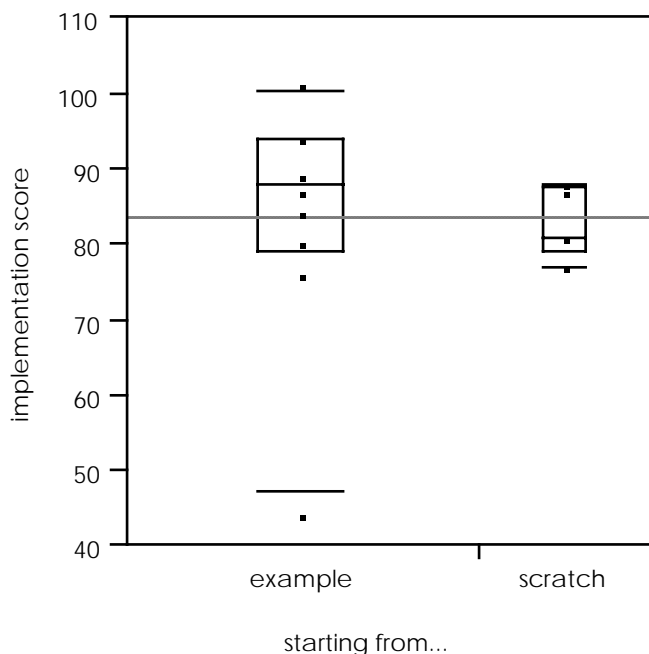


Figure 3: Team implementation scores for teams starting from scratch versus teams starting by modifying the "er" example. 90%, 75%, 50%, 25%, and 10% quantiles are shown. The team with implementation score of 44 has been removed from analysis as an outlier.

We used a t-test to determine whether teams starting from the "er" example tended to perform significantly better than teams starting from scratch (Figure 3). One point, representing a team which experienced severe organizational difficulties which were primarily responsible for a very low implementation score, must be removed from this analysis as an extreme outlier (according to the definition given in [31]).

The test yielded a p-value of 0.15 ($t = 1.538$), which is significant at the 0.20-level and provides some evidence that teams who started by modifying an example tended to be more effective than those starting from scratch. This indicates that, overall, the benefits of relying on an existing example as a starting point (which may include being able to exploit an existing file structure and to model new classes on similar ones which already exist in the example) outweigh the negatives (the extra work involved in identifying relevant functionality and removing irrelevant code).

HYPOTHESIS 7: For implementing a set of requirements in a framework-based environment, if a suitable example application can be found, then adapting that application is a more effective strategy than starting from scratch.

It was especially noteworthy that every one of the teams who started from an example used the "er" editor as a starting point. There was, in fact, a second example application which could conceivably have served as a starting point also. This second important example application was a small drawing editor that seemed to provide most of the functionality needed for the OMT editor, as well as much extraneous functionality. It is very interesting that no teams chose the drawing editor as a starting point. Of the teams who commented on this decision during the final interviews, most seemed to agree with the reasoning that the drawing editor just had too much "extra functionality and complicated code that we do not need", and that the large amount of extraneous functionality would be too confusing to enable the code responsible for the functionalities of interest to be easily separated out. If we had been able to compare the performance of teams who had started from each of the examples, we might have been able to draw useful conclusions about what attributes of an example make it a better or worse starting point for implementation. This is a promising direction for future work.

E. Hypotheses About the Importance of the Object Model

From the qualitative analysis of interviews and self-assessments, there were some general indications about the importance of the object model:

- 6 of our development teams mentioned – without being prompted by us – the importance of their object models as guides to implementation. 3 of these teams reported that they had been able to stay fairly close to their original object model of the system during the course of implementation. All of these teams ranked in the

top half of the class with regards to implementation score. The remaining 3 teams were reporting problems; they had strayed from their original model during implementation. It seems that this inability to follow the model had some negative effects, as all 3 were ranked in the bottom half of the class. (Since 5 of these 6 teams had all received the same grade on the original model, it seems unlikely that the variation in performance could have been caused by factors outside the implementation phase, such as the quality of the model itself.)

- 2 of the 3 teams who were not able to follow their object model also reported difficulties due to the inconsistent nature of the examples.

Throughout this study we have observed the central importance of the object models in development. Although we assume that object models are important and necessary parts of any software development effort, their interaction with framework development seems even more noteworthy. When using frameworks, there is the important question of whether to modify the object model of the system, so as to exploit a piece of functionality offered by the framework that might not exactly fit the original plan, or to keep the object model “as is”, even if it makes implementing the application on top of the framework harder.

The general indications from the interviews and self-assessments leads us to:

HYPOTHESIS 8: Overuse of framework functionality can lead to negative effects (e.g., rework) which might overcome the positive ones (e.g., short cycle time).

It seems that one cause of trouble may have been that these teams did not understand the examples well enough to adapt the functionality illustrated by the examples to the system being developed. Teams who implemented incorrect functionality may have been simply too willing to make modifications to their planned system to accommodate the examples more easily.

Schneider and Repenning [39] present an interesting study of framework use that comes to similar conclusions. What we describe as straying from the original object model, they identify as projects for which “the application design process had been driven by features of the framework,” a condition which is described as being present in some development efforts which “ended up with really messy designs, or were cancelled.” As we noticed through observation of our teams, they identify the likelihood that a team will have to reimplement some functionality as one of the major negative consequences arising from this situation: “Premature design decisions made during the feature-driven phase can corrupt application system architecture or require abandonment of much work.” They further point out that the developer’s temptation is to deal with the easier problem of adapting functionality provided by the framework first, leaving more difficult functionality (which

is not guaranteed to fit nicely into the framework-based development) for later and making breakdowns all the more devastating when they occur.

Some of the successes of our EB technique may also have come from its treatment of the object model. We noticed in our analysis of problems that a much higher percentage of teams using modified versions of EB reported wasted effort due to incorrect implementations than teams using EB. We attribute this to the EB technique’s focus on the object model as a guide for implementation, which may not have been so evident in ad hoc variants.

These results create an interesting tension with our experiences with example-based techniques, which seem to imply that reuse of framework functionality is always a positive thing. Taken together, our conclusions indicate that, within proper bounds, exploiting functionality from the framework and example set is the most helpful direction - otherwise, it can be very bad indeed. Schneider [39] draws a similar conclusion: although overuse of framework functionality can lead to negative effects, as described above, exploitation of the low-level framework features is a sensible trend that presumably pays off, when used within bounds. Discerning just where these bounds may lie is a crucial question that awaits further study; for now, we conclude only that development difficulties are liable to appear unless the object model of the system is well supported by the framework and its example set (i.e., the framework permits an implementation of the system that does not require major deviations from the object model).

IX. ANSWERS TO RESEARCH QUESTIONS

In this section, we relate our observations to our specific research questions for the study.

A. Can strategies for learning frameworks be identified?

Hypotheses 1 and 2 address this question.

The example-based technique was identified as an effective strategy in our environment. While we cannot say in all cases that an example-based learning approach would be superior to one based on the class hierarchy and model of interaction, the indication of this study was that for novice users, the examples were a more effective way to learn. Within the category of example-based techniques, we further differentiate “strictly adopted” from “ad hoc adopted.” Although these techniques share many common features, our analysis of their use in practice discovered certain characteristic differences. This leads us to classify them as separate strategies.

Although the hierarchy-based approach cannot be deemed effective in our environment, as it was abandoned and not used to implement any of the semester projects, we cannot assume that a hierarchy-based approach is always inferior.

The most important environmental condition appeared to be the subjects' familiarity with the particular framework being used. Several of our subjects had recognized benefits of the HB technique but were unable to apply it due to their lack of familiarity. They expressed this in comments during the interviews such as: "The HB procedure was more similar to what I normally do, but..." or "I found the examples limiting in some ways and thought the HB procedure would address this problem, but..." It is possible that, if our subjects had had more experience with the framework, the HB procedure would have proven better suited to their needs. Indications are that hierarchy-based procedure required more experience with the framework to be used effectively.

B. What are the characteristics of these learning strategies?

Since the subjects of our study only had significant experience with the EB technique, we can report only on the characteristics of example-based strategies. (Our observations on this subject were recorded in hypotheses 3 through 8.) We identified two main types of example-based strategies: strictly adopted and ad hoc adopted, each with its own strengths and weaknesses. The relative effectiveness of each seems to be most strongly determined by how closely the object model of the system to be developed corresponds with the existing applications.

Hypothesis 5 expresses our observation that when the functionality called for by the object model is well-contained in the set of existing applications, just about any example-based technique should be helpful. However, as illustrated by hypothesis 4, a strictly adopted technique can't take the developer far beyond what is provided by the existing applications themselves. In a situation in which the set of applications is sparse and does not contain the necessary functionality, an ad hoc technique may be more appropriate. As hypothesis 6 indicates, if the set of applications is particularly large, then a strict adaptation technique may be most helpful. Despite its weaknesses, such a technique in procedural form was shown to guide the developer toward implementing the object model "as is" and away from "gold-plating," or spending time providing extra features that seem nice but are not necessary.

From the experiences of our beginning learners, we also have evidence about other characteristics that are required by example-based techniques in order to be successful. Hypothesis 3 indicates that future studies need to be undertaken to determine if we can add better guidance for helping developers find functionality in existing example applications. Hypotheses 7 and 8, respectively, may indicate that example-based techniques should guide developers to begin their implementation from an existing application, if a suitable one can be found, and to stay closely to the original object model once implementation has begun. (It is possible that, as developers get more

experience with the framework, it may be possible to synchronize the design of the system more closely to the framework infrastructure from the beginning, thereby minimizing the problem of mismatch between the system design and the framework. Our study did not address this possibility.)

X. THREATS TO VALIDITY

There are three tests which can be considered to evaluate the quality of any empirical study: construct validity, internal validity, and external validity [25].

A. Construct Validity

Construct validity aims to assure that the study correctly measures the concepts of interest. The main problem is that variables never measure only the construct of interest but also other extraneous sources of variation. One tactic to enhance construct validity is triangulation: the use of multiple sources aimed at corroborating the same fact or phenomenon [51, pp.90-94].

In our study we applied data triangulation, by including multiple measures for the same aspect of interest and different collection methods for the same measure (Table 1).

B. Internal Validity

Internal validity aims to establish correct causal relationships between variables as distinguished from spurious relationships. Although a case study cannot have the same internal validity as a controlled experiment, because the investigator has little control over events, there are analysis techniques that can strengthen the internal validity, even for exploratory studies like this.

We made inferences using the qualitative analytic technique described in [13]. It consists of performing a within-case analysis, to gain familiarity with each case and find emerging patterns, followed by cross-case analysis, to look for similarities and differences between cases. Although this is not a common method of analysis in computer science, it is a recommended approach for social sciences and other fields that require the analysis of human behavior [13, 29]. It is well suited for our purposes here because our variables of interest are heavily influenced by human behavior and because we are not attempting to prove hypotheses about framework usage, but rather to begin formulating hypotheses about this process, about which we currently know little.

A specific threat to internal validity might be that we have constructed one of the techniques incorrectly, which would explain the differences in performance. As regards the example-based technique we used, EB, we attempted to minimize the odds of making this mistake by basing the specific technique on the example set which is provided by the framework's authors themselves. For the hierarchy-

focused technique, HB, we based our model of the framework on the most important facets of the framework definition. We feel that the use of the inheritance hierarchy means that our model is complete because all functionality provided by the framework must be encapsulated in one of the classes of the class hierarchy, with the inheritance relations showing how the functionalities provided are related to one another. Thus, while there may be other models more adept at supporting framework learning, we feel confident that our model is adequate to the job.

We also undertook quantitative analyses to test the effects of potentially confounding factors (such as differences in effort spent, understanding achieved, or previous experience) which could be rival explanations to our findings. As pointed out in section VII.B, the relatively low number of data points in this study means that these results should be seen, not as providing a definitive list of the important factors, but of identifying potential factors likely to have some impact on framework learning. It remains for further study to verify this impact, a topic we return to in the next section.

C. External Validity

External validity aims to assure that the findings of the study can be generalized beyond the immediate study. Although generalization can be achieved only through replication in multiple studies, we believe that our findings are relevant for a larger population than this single study.

A first threat to external validity might be that professional developers would have behaved differently than the students that we used as the subjects of the study. Certainly, this is always a danger in studies of this sort. However, in this case we feel that this difference would not be a strongly significant one. Although the level of industrial experience in the class was not high, all students had experience both programming in the language used (C++) and in object-oriented techniques. More importantly, even professional developers would almost certainly have been novices in terms of the use of the ET++ framework, so that the most immediately applicable experience would not have significantly varied in either case.

A second threat to external validity might be that our findings are tied to the framework we used, ET++. Although we cannot completely rule out this threat to validity, ET++ has been thoroughly tested and improved from the initial version and it incorporates seventeen of the design patterns in [17]. From this point of view, we consider ET++ representative of the class of sophisticated white-box frameworks that pose learning problems, which can be major inhibitors against their use.

XI. CONCLUSIONS AND FUTURE RESEARCH

This paper has formulated a set of well-motivated hypotheses concerning white-box frameworks based upon direct observation of white-box framework use in development. As the research community builds up

confidence as to the validity (or falsity) of such hypotheses, a more objective basis can be constructed for activities such as tool support and training for this class of framework. For example, hypotheses 1 and 2 are the result of evidence showing that learning by example (as opposed to gaining familiarity with the framework itself first) is useful for helping beginning learners produce working systems quickly. While the general benefits of example-based approaches may or may not seem intuitive, this study has provided some evidence, based on empiricism rather than intuition, that such approaches can be of use on non-trivial development projects using frameworks. Since an organization that uses frameworks will tend to build up a set of related applications, all based on the same underlying structure, hypotheses 1 and 2 indicate that reuse of components or even whole applications from this set can be especially beneficial in a framework environment. Frameworks that come packaged with a set of example applications can provide an analogous benefit.

Hypotheses 4 through 6 (summarized in hypothesis 3) qualify the benefits of example-based learning techniques by pointing out some of their limitations. A direct implication of these hypotheses is that the emphasis placed on adapting examples should vary according to the relation between the example set and the application to be developed. Hypothesis 7 implies that reusing as much previous functionality as possible (even including whole previous applications) is a useful strategy in this kind of application environment. Finally, hypothesis 8 provides more evidence for the benefits of a general software engineering principle (viz. that a system's design should be closely followed during implementation) by showing that it also applies to development using framework technology.

As we have emphasized, the results of this study are hypotheses for further study, not definitive conclusions. This is partly due to a number of factors (which we presented and discussed in section X) that affect the strength of the conclusions that can be drawn from our observations. However, we hope that this paper has made a useful contribution by identifying an initial set of factors that should be controlled, monitored, or tested in later studies.

One lesson learned from this study is the importance of process conformance; subjects are rarely malicious but are almost always loathe to use processes they are uncomfortable with to achieve a result they know can be reached in another manner. This discomfort can be the result of a steep learning curve for the new technique, or the result of the unsuitability of the technique for the current environment (as with the HB technique in this study). The experimenter must make the important decision of whether:

- To constrain subjects to use a specific process, in order to draw conclusions about that process, or
- To allow subjects the freedom to use processes they feel are useful, while monitoring the processes

undertaken. (This strategy was the one used in our study.)

In the latter case, qualitative data collection and analysis are especially important. Without qualitative analysis in this study, we could only have concluded that the HB technique was unsuitable to the environment; through our analysis of interviews and problem reports, we have some confidence that the contributing factor was low subject experience with the framework.

Another result of the qualitative analysis was the identification of another factor influencing framework learning, namely, the level of specificity at which an example-based technique is followed. This factor, which contributes to our hypotheses 4, 5, 6, and 8, seems nearly impossible to assess without the use of qualitative methods. Because our analysis detected distinct differences between subjects using different levels of specificity (in our terminology, between “strictly adopted” and “ad hoc adopted”) this factor should be assessed in future studies as well.

As shown in section VII.B, a third important factor identified by this study was the previous experience of the subjects. The correlation between subject experience and effectiveness at implementation shows that it is necessary to assess the effect of experience on the use of new techniques. In any case, subject experience can be expected to have an effect on the outcome of software engineering practices and it is necessary to check that it does not overshadow other factors, such as the use of the technique under study.

It is our hope that future studies in this area can use our results as a beginning for verifying, bounding and extending these hypotheses. Certainly, studies in the framework domain have the potential for verifying the practical implications of these hypotheses, such as whether training professional developers to concentrate too much on reuse will result in systems of lower quality (as it did for our student subjects, reflected in hypotheses 4 and 8). Studies can bound the hypotheses by discovering contexts in which the techniques are more or less effective, e.g., EB is effective for beginning framework users but HB is more effective for advanced framework users. But it is to be hoped that studies will be undertaken to test the generalizability of our results to other areas as well. For example, the indications from this study have already proven useful for our study of software reading techniques. In this study, we saw that the level of specificity at which a technique is followed can have a distinct effect on the outcome; we have since run experiments on other reading techniques in which the level of specificity was explicitly varied [41]. These experiments have helped us conclude that specificity is an important variable in software reading research, although its specific effects may vary from case to case.

Other potential areas of generalization exist as well. One example (among many) is that hypotheses 1 and 2,

concerning the effectiveness of example-based learning, need not be limited to framework-based environments. Evaluations of example-based learning in related areas would be of great interest. For example, could these results be taken to imply that an effective way to learn to program would be to first study existing programs (i.e. that the first step in learning to *write* good programs is learning how to *read* them)?

In order to facilitate replication or review of this study we have set up a web site containing as much as possible of our experimental materials and data. This web site may be found at [42].

XII. ACKNOWLEDGEMENTS

Our thanks to Gianluigi Caldiera for his invaluable assistance designing and running this experiment as a part of the course CMSC 435 (Fall 1996) at University of Maryland, College Park. Our thanks also go to the students of CMSC 435 for their cooperation and hard work.

This work has been supported by UMIACS and NSF grant CCR9706151.

XIII. REFERENCES

- [1] *MacApp Programmer's Guide*. Apple Computer, 1986.
- [2] V. Basili, R. Reiter, Jr, “A Controlled Experiment Quantitatively Comparing Software Development Approaches”, *IEEE Trans. Software Engineering*, vol. SE-7, pp. 299-320, May 1981.
- [3] V. Basili, S. Green, O. Laitenberger, F. Lanubile, F. Shull, S. Soerumgaard, and M. Zelkowitz, “The Empirical Investigation of Perspective-Based Reading”, *Empirical Software Engineering: An International Journal*, vol. 1, no. 2, pp. 133-164, 1996.
- [4] V. Basili, G. Caldiera, F. Lanubile, and F. Shull, “Studies on Reading Techniques”, in *Proc. of the Twenty-First Annual Software Engineering Workshop*, Dec. 1996, pp. 59-65.
- [5] D. Baumer, G. Gryczan, R. Knoll, C. Lilienthal, D. Riehle, and H. Zullighoven, “Framework Development for Large Systems”, *Communications of the ACM*, vol. 40, no. 10, pp.52-59, October 1997.
- [6] K. Beck, R. Johnson, “Patterns Generate Architectures”, in *Proc. ECOOP'94*, 1994.
- [7] D. S. Brandt, “Constructivism: Teaching for Understanding of the Internet”, *CACM*, vol. 40, pp. 112-117, Oct. 1997.
- [8] D. Brugali, G. Menga, and A. Aarsten, “The Framework Life Span”, *Communications of the ACM*, vol. 40, no. 10, pp.65-68, October 1997.
- [9] J. Carroll, *The Nurnberg Funnel: Designing Minimalist Instruction for Practical Computer Skill*. Cambridge, MA: MIT Press, 1990.

- [10] M. Chi, M. Bassok, M. Lewis, P. Reimann, and R. Glaser, "Self-Explanations: How Students Study and Use Examples in Learning to Solve Problems", University of Pittsburgh, Technical Report UPITT/LRDC/ONR/KBC-9, Nov. 1987.
- [11] W. Codenie, K. DeHondt, P. Steyaert, A. Vercammen, "From Custom Applications to Domain-Specific Frameworks", *CACM*, vol. 40, pp. 71-77, Oct. 1997.
- [12] Conover, *Practical Nonparametric Statistics*, 2nd Edition. NY: John Wiley & Sons, 1980.
- [13] K. Eisenhardt, "Building Theories from Case Study Research", *Academy of Management Review*, vol. 14, no. 4, pp. 532-550, 1989.
- [14] M. A. Fayad, and D. C. Schmidt, "Object-Oriented Application Frameworks", *Communications of the ACM*, vol. 40, no. 10, pp.32-38, October 1997.
- [15] C. Frei and H. Schaudt, *ET++ Tutorial: Eine Einführung in das Application Framework*. Software Schule Schweiz, Bern, 1991.
- [16] G. Froehlich, H. Hoover, L. Liu, and P. Sorenson, "Hooking into Object-Oriented Application Frameworks", in *Proc. of the 19th International Conference on Software Engineering*, May 1997, pp. 491-501.
- [17] E. Gamma, R. Helm, R. Johnson, and J. Vlissides, *Design Patterns: Elements of Reusable Object-Oriented Software*. Reading, MA: Addison-Wesley, 1995.
- [18] D. Gangopadhyay, S. Mitra, "Understanding Frameworks by Exploration of Exemplars", in *Proc. of 7th International Workshop on CASE*, July 1995, pp. 90-99.
- [19] H. G. Glaser, A. L. Strauss. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. Hawthorne, NY: Aldine Publishing Company, 1967.
- [20] A. Goldberg, *Smalltalk-80: The Interactive Programming Environment*. Reading, MA: Addison-Wesley, 1984.
- [21] L. Hatcher, E. J. Stepanski, *A Step-by-Step Approach to Using the SAS® System for Univariate and Multivariate Statistics*. Cary, NC: SAS Institute Inc., 1994.
- [22] R.E. Johnson and B. Foote, "Designing Reusable Classes", *Journal of Object-Oriented Programming*, vol.1, no. 5, pp.22-35, June/July 1988.
- [23] R. Johnson, "Documenting Frameworks with Patterns", in *Proc. OOPSLA '92*, October 1992, pp. 63-76.
- [24] R.E. Johnson, "Frameworks = Patterns + Components", *Communications of the ACM*, vol. 40, no. 10, pp.39-42, October 1997.
- [25] C.M.Judd, E.R.Smith, and L.H.Kidder, *Research Methods in Social Relations, sixth edition*. Forth Worth: Holt, Rinehart and Winston, Inc., 1991.
- [26] P. Koltun, L. Deimel Jr., and J. Perry, "Progress Report on the Study of Program Reading", *ACM SIGCSE Bulletin*, vol. 15, pp. 168-176, Feb. 1983.
- [27] O. Laitenberger and C. Atkinson, "Generalizing Perspective-based Inspection to Handle Object-Oriented Development Artifacts", in *Proc. ICSE'99*. (To appear.)
- [28] T. Lewis, L. Rosenstein, W. Pree, A. Weinand, E. Gamma, P. Calder, G. Andert, J. Vlissides, K. Schmucker, *Object Oriented Application Frameworks*. Greenwich: Mannings Publication Co., 1995.
- [29] M. Miles, "Qualitative Data as an Attractive Nuisance: The Problem of Analysis", *Administrative Science Quarterly*, vol. 24, no. 4, pp. 590-601, 1979.
- [30] H. Mili, H. Sahraoui, I. Benyahia, "Representing and Querying Reusable Object Frameworks", in *Proc. of the Symposium on Software Reusability*, May 1997.
- [31] R. Ott. *An Introduction to Statistical Methods and Data Analysis*. Belmont, CA: Duxbury Press, 1993.
- [32] W. Pree, *Design Patterns for Object-Oriented Software Development*. Reading, MA: ACM Press & Addison-Wesley Publishing Co., 1995.
- [33] D. Roberts and R. Johnson, "Patterns for Evolving Frameworks", in *Pattern Languages of Program Design*, R.C. Martin et al. (eds.), Software Patterns Series, Addison Wesley, 1997.
- [34] M. B. Rosson, J. M. Carroll, and R. K. E. Bellamy, "SmallTalk Scaffolding: A Case Study of Minimalist Instruction", in *Proc. CHI '90*, April 1990, pp. 423-429.
- [35] S. Rugaber, S. B. Ornburn, and R. J. LeBlanc, Jr., "Recognizing design decisions in programs", *IEEE Software*, vol. 7, pp. 46-54, Jan. 1990.
- [36] J. Rumbaugh, M. Blaha, W. Premerlani, F. Eddy, and W. Lorensen, *Object-Oriented Modeling and Design*. Englewood Cliffs, NJ: Prentice Hall, 1991.
- [37] H. A. Schmid, "Creating Applications from Components: a Manufacturing Framework Design", *IEEE Software*, vol. 13, no. 6, pp.67-75, November 1996.
- [38] D. C. Schmidt, "Applying Patterns and Frameworks to Develop Object-Oriented Communication Software", in P. Salus, ed., *Handbook of Programming Languages*, vol.1, MacMillan Computer Publishing, 1997.
- [39] K. Schneider, A. Repenning, "Deceived by Ease of Use: Using Paradigmatic Applications to Build Visual Design Environments", in. *Proc. of the Symposium on Designing Interactive Systems*, Aug. 1995.
- [40] C. B. Seaman, V. R. Basili, "An Empirical Study of Communication in Code Inspection", in *Proc. ICSE'97*, May 1997, pp. 96-106.
- [41] F. Shull. *Developing Techniques for Using Software Documents: A Series of Empirical Studies*. Ph.D. thesis, University of Maryland, College Park, December 1998.
- [42] F. Shull, "Reading Techniques for Object-Oriented Frameworks", http://www.cs.umd.edu/projects/SoftEng/ESEG/manual/sbr_package/manual.html.
- [43] J. Singer and T. C. Lethbridge, "Methods for Studying Maintenance Activities", in *Proc. of 1st International Workshop on Empirical Studies of Software Maintenance*, Nov. 1996, pp. 105-110.

- [44] S. Sørungård. *Verification of Process Conformance in Empirical Studies of Software Development*. Ph.D. thesis, Norwegian University of Science and Technology, February 1997.
- [45] Software Engineering Laboratory, *Recommended Approach to Software Development, Revision 3*. National Aeronautics and Space Administration, Software Engineering Laboratory, SEL-81-305, 1992.
- [46] Taligent, Inc., *The Power of Frameworks*. New York: Addison-Wesley, 1995.
- [47] J. Vlissides, *Unidraw Tutorial I: A Simple Drawing Editor*. Stanford University, 1991.
- [48] A. von Mayrhauser and A. M. Vans, "Industrial Experience with an Integrated Code Comprehension Model", *IEEE Software Engineering Journal*, pp. 171-182, Sept. 1995.
- [49] A. Porter, L. Votta Jr., and V. Basili. "Comparing Detection Methods for Software Requirements Inspections: A Replicated Experiment". *IEEE Transactions on Software Engineering*, 21(6): 563-575, June 1995.
- [50] A. Weinand, E. Gamma, and R. Marty, "Design and Implementation of ET++, a Seamless Object-Oriented Application Framework", *Structured Programming*, vol. 10, no. 2, 1989.
- [51] R. Yin, *Case Study Research: Design and Methods*. London: Sage Publications, 1994.
- [52] Z. Zhang, V. Basili, and B. Shneiderman, "Perspective-based Usability Inspection: An Empirical Validation of Efficacy", *International Journal of Empirical Software Engineering*, special issue on Human-Computer Interaction. (To appear.)