Study of Novice Programmers using Eclipse and Gild

Peter C. Rigby CHISEL Group Dept. of Computer Science University of Victoria

pcr@uvic.ca

ABSTRACT

In this paper we discuss a pilot user study that compares the use of two integrated development environments (IDEs), Eclipse and Gild, by novice programmers. Gild is a perspective for Eclipse that is intended to be more suitable for first-year students who are learning how to program in Java. This study focuses on qualitative and quantitative measures; the quantitative measures include: efficiency, effectiveness, satisfaction and understanding. Two statistically significant results are obtained from the satisfaction measure, in particular: the frustration level and the overall level of satisfaction. The mean differences for the remaining measures indicate that Gild was more suitable for novices than Eclipse. Qualitative analysis yields suggestions for improvement for both interfaces and also identifies areas of success.

Keywords

Novice, Integrated Development Environment, user study

1. INTRODUCTION

The Gild (Groupware-enabled Integrated Learning and Development) environment, a perspective of the Eclipse IDE, was created to be a more appropriate development environment for novice users [4]. Gild consists of a teaching perspective and a student perspective. The teaching perspective allows instructors to create course units that can be easily imported by students. This perspective also provides basic assignment grading support. For students, the Gild environment reduces and simplifies the menu and toolbar contents, debugger, code editor, task/todo list, and the import and export process; Gild also provides an integrated Web browser, enforces clear steps for the save/build/run process, and allows students to run fragments of their code without a main class through JPages. Additional novice-friendly error descriptions are offered for common compilation errors.

Since Gild's deployment in January 2003, it has been used in Canada (University of Victoria), in the USA (Berkeley, Virginia Tech and the University of Washington), in Germany (Eberhard Karls Universität Tübingen), and possibly other institutions that have not Suzanne Thompson CHISEL Group Dept. of Computer Science University of Victoria zazam@uvic.ca

contacted us. At the University of Victoria it has been successfully used in first-year introductory programming courses. Feedback from students has been elicited in the form of questionnaires, interviews, and an ethnographic study; however, a quantitative comparison between Gild and Eclipse had yet to be completed. If it can be empirically shown that Gild is indeed more useful for novice programmers, this would give weight to the claims made on its behalf and potentially lead to greater acceptance. We have conducted an initial experiment that compares Gild and Eclipse from a novice programmer's perspective based on efficiency, effectiveness, satisfaction, and understanding.

2. HYPOTHESIS

It is hypothesized that Gild will be more efficient, effective, and satisfactory when compared to Eclipse by novice programmers. It is also expected that Gild will lead to a better understanding of the programming environment.

3. SUBJECTS

A total of six subjects voluntarily participated in this study. Subjects were all students enrolled in CSC 110, an introductory Java programming course, at the University of Victoria. Students were recruited from CSC 110 because they had a basic understanding of Java, but few had any experience with advanced IDEs. Students used TextPad for their coursework, which provides basic editing functionality. All six subjects were enrolled in their first semester of their first year at the University of Victoria, and all had just recently graduated from high school. Each subject was taking CSC 110 because it was required by their programs; four subjects were computer science students, while two were physics students. The subjects all had at least two months of Java programming experience and were beginning to explore IDEs.

4. DESIGN

For this experiment a repeated measures design was used, that is, subjects used both Eclipse and Gild, cancelling individual skill differences. With a small sample size this design provides more statistical power than an independent samples design; in an independent samples design differences in individual skill would have a greater impact on the results. Subjects were equally and randomly assigned Eclipse or Gild as their initial IDE.

Separate installations of Gild and Eclipse were placed on the same computer. This was done because the Gild plug-in alters the Eclipse environment (e.g., new buttons on the task bar.) Camtasia (with audio recording) was used to record the user study session. There were only the two principal researchers in the room with the subject during the user study. Subjects were asked to complete several tasks from a modified assignment that they had completed at the beginning of the current semester. It was hoped that basing the code on a class assignment would reduce the learning time of the code itself, thus lessening the impact of individual code understanding abilities; this focused the experiment on the IDE. There were three tasks in a set to solve for each IDE: a compilation problem, a logical problem, and an extra coding assignment. The tasks were always presented in the same order: set one tasks first, set two tasks second. Although the sets of tasks were different from each other, they had the same level of difficulty.

5. PROCEDURE

Subjects were asked to set aside 90 minutes for this user study. To begin, a background (entry) questionnaire was asked by one of the researchers. The subject was shown a simple program previously created, and was walked though the steps that needed to be completed for the code to run. After this, the subject was shown how to use the debugger. None of the subjects had any previous experience with debuggers. The subject was shown how to set breakpoints, what the different debugger views were, and how to step through the code. The subject was shown how to run the program again. This was necessary since in Eclipse the subject needed to know that the run configuration did not have to be created again; rather, the run button could be used. Users were not shown how to create new files, the problem view, or told anything about how Eclipse and Gild informed users of errors. The subject was directed to a list of the three tasks that were written up as an HTML page. The subject was given 25 minutes to try to complete the tasks for the first IDE. After 25 minutes had passed, the subject completed the understanding questionnaire; subjects could check their answers with the IDE. The subject was then asked to complete a satisfaction questionnaire.

The subject was provided with a tour of the second IDE through the use of the same simple program. They were given 25 minutes to complete the second set of tasks. At the end of 25 minutes they were again asked to complete the understanding questionnaire and a satisfaction questionnaire.

6. MEASURES AND VARIABLES

ISO 9241-11 [1] defines usability as "The effectiveness, efficiency, and satisfaction with which specified users achieve specified goals in particular environments." These three measures are often the focus of HCI-related experiments and are used in this experiment. Frøkjær *et al.* [3] found that there was at most a weak correlation between these three measures. Therefore, this experiment evaluated these three dimensions independently. This experiment was also intended to measure the subject's understanding of the IDE and the programming process.

6.1 Efficiency

ISO [1] defines efficiency as "the resources expended in relation to the accuracy and completeness of goals achieved." Task completion time was the only resource measured in this experiment. For each IDE we considered three tasks: fixing the compiler error, fixing a logical error, and an extra coding assignment.

6.2 Effectiveness

ISO [1] defines effectiveness as "the accuracy and completeness with which specified users can achieve specified goals in particular environments". Completeness was measured based on the subject's correct or incorrect completion of each given task. The granularity was to the quarter of a mark since some subjects came close to the correct answers within varying degrees. A good measure of accuracy was difficult to establish.

Accuracy was measured by the number of wrong turns (i.e. incorrect clicks) taken by the subject. Counting the number of incorrect clicks can be used to capture errors since the subjects did not use keyboard shortcuts (besides copy and paste) when working within the IDE. The question that determined if an action was an error was: Does the user's click help him or her accomplish the current task? If the answer was no, then that click was considered to be an error. Scores for accuracy were completed by independent trained evaluators.

Things that are not considered to be errors include coding errors, the long way of doing a task, and any hovering and reading. Examples of errors include: selecting the incorrect menu for a desired action, running an application when a compilation error is present, and the repeat undoing and redoing of an action. Note that once a wrong path has been taken the subject has a new goal: to get out of this wrong path. For example, if they run the code when there is an error they will be asked if they wish to run it anyways; if they select "cancel" they have left the incorrect path, but if they select "ok" then they have performed another error.

Questions and hints are also worth noting, and are considered to be similar to errors, but are counted independently. It was not suggested that subjects ask the investigators questions; although questions would be answered and hints were given when a subject became exasperated with the current task.

6.3 Satisfaction

ISO [1] defines satisfaction as "the comfort and acceptability of the work system to its users and other people affected by its use." Satisfaction was measured by a five point Likert scale. (Usually "Agree" to "Disagree".) Having the satisfaction questionnaire following the completion of the understanding questionnaire may have led to some bias due to the subject's ability to answer the understanding questions.

6.4 Understanding

Understanding the programming environment and the programming process is essential to the success of a novice programmer. Efficiency, effectiveness and satisfaction fail to measure understanding. Bohlen *et al.* [2] use the measure of achievement instead of efficiency in their study of end-users learning new software. Achievement is measured through test, practicum, and assignment scores. Although it would be ideal to measure achievement as per Bohlen *et al.*, it was not feasible to do so in this study. The study used a related measure, the subject's understanding of the IDE. The questionnaire consisted of simple true/false questions that were assigned to each subject twice throughout the experiment: after using the first IDE, and after using the second IDE. The qualitative results will be discussed first followed by the quantitative results.

7. QUALITATIVE RESULTS

In this paper the most interesting qualitative results are discussed. Due to the nature of the assignments given, that is, the program would block waiting for input, there were problems that came to light that had not been encountered before. It should be noted that the assignments included a Keyboard.java file that was also supplied to the subjects in their class assignment; this file provided a simple way to get input from the keyboard. Subjects were not required to know how the keyboard class was implemented; the investigators did not alter the usage of this file in any way.

7.1 Gild-specific results

Although Gild has a relatively simple interface, there were some problems encountered by subjects. When code is executing (running) or being debugged in Gild, the background of the code editor pane changes from a white colour to a light yellow colour to indicate the change in the state of the IDE. Some subjects may have believed that the yellow background colour was only associated with debugging. After subjects set a breakpoint and went to debug the code, the program would eventually block, waiting for input. Gild always implicitly sets a breakpoint on the first line of code to be executed, so subjects would step through their code quickly, trying to get to their first breakpoint. Several subjects did not understand why the code was not reaching their first breakpoint, which was almost always located below the first block for input. A possible solution for this may be to flash the console window when the application is waiting for input.

Related to subjects' understanding of compilation errors, was their use of the "extra help" for compiler errors that is provided to users. Although subjects were told that there was a compilation error in the program, few initially compiled the program. A hint was given for subjects to look at the "Problem View" window, but since few subjects compiled the code before going to look at the problem view, no compiler error messages were present. Further to this, few subjects took the time to fully read the possible solutions for the particular compiler errors.

Another problem encountered by two Gild users was that when they searched in the help for a solution for their particular problem, they ended up searching the help for all of Eclipse. The results returned referred to far more complicated examples then the subjects could understand, although one subject did find the answer they were looking for by chance. It is recommended that the help menu in Gild makes the help for Gild more easily accessible and more prominent then the help for Eclipse; searching for help should only search within the help for Gild.

7.2 Eclipse-specific results

Some subjects were quite apprehensive about configuring Eclipse to run their Java application, despite a demonstration by the investigators on how to do this. No subjects ran into major troubles with this process, though many kept running their programs though the Run menu instead of using the run or debug icon shortcuts.

A major difficulty subjects had with Eclipse was the use of the debugger. Like in Gild, subjects would set a breakpoint, but the application would block for input before switching to the debug perspective. What was happening was that the application would run, and then when a breakpoint was encountered Eclipse would prompt subjects to switch to the debug perspective. Subjects would then, in confusion, either debug the program again (a new thread of execution) or switch to the debug perspective. The two options were often combined by subjects in interesting and creative ways. Few subjects noticed that the program was blocking. If the subject chose to switch over to the debug perspective, the "Console" window would not show that the program was blocking - this would only be shown in the Java perspective. This led to subjects running another debug thread within the debug perspective, and often

double-clicking on breakpoints in an effort to get Eclipse to run the code from the breakpoint.

There was also an additional problem only encountered in Eclipse whereby "ghost projects", i.e., the Hello World project that was used to demonstrate how to use Eclipse, would surface occasionally despite being closed. In one instance there was an old terminated thread belonging to the Hello World application that surfaced in the debug perspective. It is recommended that Eclipse remove all references to closed projects.

7.3 Results for Gild and Eclipse

A major incorrect step that many subjects did in both Eclipse and Gild was to try to run code that they knew contained a compilation error. Few subjects wanted to take the time to read through the program and understand what was wrong; perhaps this was due to the time constraint and the number of tasks that they felt they needed to accomplish. The investigators do not believe that the subjects felt any more rushed than when usually doing an assignment. It is recommended that Gild prevent code from being executed when there is a visible compilation error. Also, subjects made a habit of following the suggestions provided by the compiler error messages literally. For example, the compiler told subjects to delete a token when in fact a bracket needed to be added.

Subjects had difficulty in Eclipse and Gild with locating the terminate button. As previously noted, in Eclipse subjects would end up with many threads of execution running, though Gild did not allow this; therefore subjects were forced to find the elusive terminate button in Gild. It is recommended that the terminated button be made more visible in both environments. One subject commented that Gild was less intimidating than Eclipse since all relevant windows were visible at once. The subject felt that the layout was good.

8. QUANTITATIVE RESULTS

The purpose of this study was to provide an experimental comparison of Eclipse and Gild as used by novice programmers. The qualitative results provide insight into how novices use the two tools and suggest possible improvements. Eclipse and Gild were compared quantitatively based on the measures of efficiency, effectiveness, satisfaction, and understanding. The mean values for these measures were in the predicted direction, but none were statistically significant (see Table 1), which was likely due to a small sample size (N = 6). Interestingly, when the individual questions on the satisfaction questionnaire were compared, two statistically significant results were obtained (see Table 2 below). The results appear promising, but a larger study would be required to show that Gild has achieved its goals.

8.1 Efficiency

Efficiency was measured in terms of time required to complete all tasks in the set for an IDE. Since many of the subjects did not complete all of the tasks, a composite measure (time divided by the number of tasks completed) was used to indicate how many tasks were completed in the time the subjects used. A low score indicates that subjects finished many tasks quickly. Gild scored 6.26 units lower than Eclipse (Table 1). This result is not statistically significant (p = 0.184). On average subjects took 4.17 minutes less when using Gild than Eclipse. This trend (p = 0.159) indicated that the simplified, one-perspective Gild interface allowed users to navigate quickly through the IDE; allowing them to complete their

Table 1: Results of a paired t-test comparing efficiency, effectiveness, satisfaction, and understanding for novice users of Eclipse and Gild. *measures composing efficiency, ** measures composing effectiveness

Measures	Mean	Mean	Mean	Sig. (2
	for	for Gild	Diff.	tailed)
	Eclipse			
Time*	24.50	20.33	4.17	0.159
Time / Com-	16.85	10.59	6.26	0.184
pleted*				
Tasks Com-	2.13	2.29	-0.17	0.102
pleted**				
Questions**	3.33	2.00	1.33	0.520
Errors**	26.50	19.67	6.83	0.332
Satisfaction	24.00	27.00	-3.00	0.226
Understanding	5.83	6.83	-1.00	0.111

tasks more quickly than when using the more complex Eclipse IDE. The increased efficiency may also be due to less time lost searching for the correct options within the environment. The additional composite measure that is based on time and number of tasks completed did not reveal any additional information.

8.2 Effectiveness

Three measures were used to assess effectiveness: the number of questions or hints the subject asked or was given, the number of errors the subject made, and the number of tasks the subject completed correctly. No measure achieved statistical significance, though the results were in the predicted direction.

The number of tasks completed by subjects using Eclipse and Gild were close and approached statistical significance (p = 0.102). The mean difference is only -0.17, which is too small to be significant since the assignment of marks was done in increments of 0.25. Ideally the experiment would have only had subjects who had never used an IDE before. However, some subjects were admitted to the sample even thought they had previous experience with IDEs. This skill difference between subjects meant that some subjects were able to solve all the tasks correctly, while other subjects only solved a few tasks correctly; this may explain why statistically significant results were not achieved.

The number of questions asked and hints given to a subject was dependent on the subject's personality. There were some subjects that seemed to lack confidence and would ask many questions, while others would persevere and would sometimes require a hint. Table 1 shows that the mean difference for questions (and hints) asked indicates that on average subjects using Gild asked 1.33 fewer questions than subjects using Eclipse. This result was not statistically significant (p = 0.520) and would require a larger sample size to eliminate between subject differences.

The number of errors was measured by the number of wrong clicks (see criteria in Section 6.2). As was the case with the number of questions asked by subjects, some subjects made many clicks, while other subjects thought more and clicked less. Subjects made on average 6.83 fewer errors when using Gild than when using Eclipse; this result is not statistically significant (p = 0.332). A larger sample size would reduce the impact of subject differences. A more effective change would be to adjust the criteria to look at higher level errors; although, higher level errors would likely be

 Table 2: Results of a paired t-test comparing individual satisfaction questions.

Question	Mean for	Mean for	Mean	Sig. (2
Num.	Eclipse	Gild	Diff.	tailed)
1	3.50	4.17	-0.667	0.102
2	3.00	3.83	-0.833	0.042
3	3.83	3.83	0.000	1.000
4	3.83	3.67	0.17	0.771
5	2.83	3.50	-0.667	0.175
6	3.50	3.83	-0.333	0.661
7	3.50	4.17	-0.667	0.025

more subjective, making it difficult to train independent scorers.

8.3 Satisfaction

The level of satisfaction the subjects experienced while using the tool was measured with a five point Likert scale. The minimum possible score is five and the maximum possible score is 35. Eclipse received a mean score of 24, while Gild received a mean score of 27 (Table 1). Although the difference was not statistically significant (p = 0.226), the mean difference of three points was in the direction predicted. Table 2 presents a paired t-test of each individual satisfaction questions and reveals that questions two and seven are statistically significant. Since the overall questionnaire did not reach statistical significance, but individual questions did, it is useful to discuss each question. The questions are discussed in order of decreasing statistical significance.

8.3.1 Question 7 "What was your overall level of satisfaction using this tool?"

The satisfaction questionnaire was developed to look at different elements of satisfaction. This final question asked in a direct manner how satisfied the subjects were with the tool. Eclipse received a mean score of 3.50 which indicates that subjects were between neutral and moderately satisfied with Eclipse. Gild received a mean score of 4.17 which means that subjects were slightly more than moderately satisfied with Gild. A score of 5 would have implied that the tool was satisfactory. This result is statistically significant (p = 0.025) and implies that novices are more satisfied with Gild than with Eclipse. This is reflected in the qualitative results (see Section 7).

8.3.2 *Question* 2 "What was your frustration level using this tool?"

On average subjects experienced a "moderate" level of frustration using Eclipse (3.00) and a "low" level of frustration using Gild (3.83). This result is statistically significant (p = 0.042) and was noticed by the experimenters (see Section 7). Novices must learn many new complex concepts. If a novice is frustrated by the IDE then learning the concepts taught in class will be more difficult. The level of frustration experienced by the novice must be kept low to make the learning experience fruitful.

8.3.3 Question 1 "I feel that the tool helped me with my tasks."

Table 2 shows a trend (p = 0.102) that subjects found Gild (4.17) more helpful with completing tasks than Eclipse (3.50). The result indicates that subjects probably found Gild more intuitive and simple to use and, therefore, more helpful when completing the assigned first-year programming tasks.

8.3.4 Questions 3, 4, 5 and 6

Questions three through six did not produce significant results or trends and will be discussed in terms of their mean values (p > p)(0.15) see Table 2. Question five (p = 0.175), "I feel intimidated when using the tool." indicated that subjects were more intimidated by Eclipse (2.83) than by Gild (3.50). Question six (p = 0.661), "I would use this tool again." received a mean score between "neutral" and "mildly agree" which indicates that both Gild (3.83) and Eclipse (3.50) should improve their support for novices. Question four (p = 0.771), "I feel comfortable to explore the tool and try new features." revealed that subjects felt relatively comfortable exploring new features, Gild (3.67) and Eclipse (3.83). However, subjects used each IDE for only 25 minutes, which does not leave much exploration time. With such a large p-value the difference been Gild and Eclipse is negligible; the means are not in the predicted direction. Question three (p = 1.000), "I find the debugger helpful." received the same mean score (3.83) for both Eclipse and Gild. However, subjects had never used a debugger and found that any debugger was an improvement over no debugger, a more comparative question would be "I find the debugger easy to use." Qualitatively, subjects had more difficulty with the Eclipse debugger (see Section 7).

8.4 Understanding

How well a subject understood an IDE was measured using a true/false questionnaire. The questionnaire consisted of 8 questions. Table 1 shows that after using Eclipse subjects had a mean score of 5.83 out of 8. After using Gild, subjects had a mean score of 6.83 out of 8. On average, subjects answered 1 more question correctly after using Gild than after using Eclipse. This result approached statistical significance with p = 0.111.

Question two, "The IDE compiles code when it is run", was the question most often answered incorrectly. It was answered incorrectly by five (out of six) subjects using Eclipse, but was answered incorrectly twice by subjects using Gild. In order to improve the novice's understanding, Gild removes many of the shortcuts that are available in Eclipse. For example, in Gild, code must be separately saved and compiled (built) so that the novice understands that these two actions are distinct; in Eclipse code is compiled when it is saved (by default).

9. RECOMMENDATIONS

In this section, recommendations are given for both IDEs and for conducting future studies. From the qualitative analysis it was noted that Eclipse's debugger and run configurations were particulary confusing for subjects. Although not all subjects were intimidated by these features, their interactions with this functionality did not indicate a full understanding of what was occurring. Gild is designed to simplify Eclipse, but Gild was not free from usability issues. It is recommended that the terminate button be made more obvious in both environments. It is also recommended that developers of Eclipse consider the needs of novice users and programmers. The underlying plug-in architecture of Eclipse lends itself to extension, removing existing functionality is less elegant; this has complicated the development of Gild. Making an IDE more comprehensible to novice and intermediate programmers should be considered as an investment since they are the next professional programmers. The ability to simplify and customize the interface will likely help novice users accept the Eclipse IDE. Ideally, Gild would not exist: Eclipse would be customizable to the point where it could resemble Gild.

The quantitative results are encouraging as two statistically significant results were obtained and the mean differences for all main measures were in the predicted directions. It is likely that with slight modifications to the measures and an larger sample size, more statistically significant results could be obtained. The measure of effectiveness produced the least interesting results (see Table 1). Effectiveness was measured by the number of tasks completed, number of questions asked and hints given, and number of errors made the latter two measures require modification. It is recommended that subjects not be allowed to ask questions or receive hints as certain personality types will be more likely to ask questions than others types. The number of errors was measured using a criteria developed for this experiment that was based on mouse clicks. Since mouse clicks do not capture high-level errors and do not differentiate error types (e.g., on the basis of severity), it is recommended that a criterion for classifying high-level errors be developed and that the criteria for low-level errors be further refined. One of the original goals of the Gild project was to understand how novices use an IDE, not just how novices use Eclipse. If a future study were to be conducted, the measures should be mapped to the individual requirements of Gild with the goal of verifying these reauirements.

10. CONCLUSION

This study provides a first attempt to experimentally compare how usable and useful two IDEs are for novice programmers. The hypothesis, which operationalizes the three standard ISO measures, that Gild will be more effective, efficient, satisfactory, and lead to greater understanding when compared to Eclipse by novice programmers, appears to have merit. The frustration level and the overall level of satisfaction of subjects produced statistically significant results. Additionally, the mean differences of all other measures, time, time divided by tasks completed, tasks completed, questions asked and hints given, number of errors made (as per the criteria developed for this experiment), a satisfaction questionnaire, and an understanding questionnaire, were in the predicted direction. The qualitative results provided insight into novices' mental model and highlighted some interaction design issues. A future study with a larger sample size will be required before this hypothesis can be accepted or rejected.

11. ACKNOWLEDGEMENTS

Our thanks to Margaret-Anne Storey, the CHISEL group, Laura Young, Jody Ryall, and the study participants for their help with this user study.

12. REFERENCES

- [1] Iso 9241-11:1998 ergonomic requirements for office work with visual display terminals (vdts) part 11: Guidance on usability, 1998.
- [2] G. A. Bohlen and T. W. Ferratt. The effect of learning style and method of instruction on the achievement, efficiency and satisfaction of end-users learning computer software. In *SIGCPR '93: Proceedings of the 1993 conference on Computer personnel research*, pages 273–283, New York, NY, USA, 1993. ACM Press.
- [3] E. Frøkjær, M. Hertzum, and K. Hornbæk. Measuring usability: are effectiveness, efficiency, and satisfaction really correlated? In CHI '00: Proceedings of the SIGCHI conference on Human factors in computing systems, pages 345–352, New York, NY, USA, 2000. ACM Press.
- [4] M.-A. Storey, D. Damian, J. Michaud, D. Myers, M. Mindel, D. German, M. Sanseverino, and E. Hargreaves. Improving the usability of eclipse for novice programmers. In *eclipse '03: Proceedings of the 2003 OOPSLA workshop on eclipse technology eXchange*, pages 35–39, New York, NY, USA, 2003. ACM Press.