



Universiteit  
Leiden  
The Netherlands

## **An online corpus of UML Design Models : construction and empirical studies**

Karasneh, B.H.A.

### **Citation**

Karasneh, B. H. A. (2016, July 7). *An online corpus of UML Design Models : construction and empirical studies*. Retrieved from <https://hdl.handle.net/1887/41339>

Version: Not Applicable (or Unknown)

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/41339>

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/41339> holds various files of this Leiden University dissertation.

**Author:** Karasneh, B.H.A.

**Title:** An online corpus of UML Design Models : construction and empirical studies

**Issue Date:** 2016-07-07

# Quality Assessment of UML Class Diagrams

*In this chapter, we present an experiment conducted for comparing how experts and students assess the quality of class diagrams. Six quality attributes were addressed: Understandability, Layout, Extensibility, Modifiability, Completeness and Correctness. From this study, we aim to find out how well students are capable of evaluating the quality of UML designs. One particular scenario that we have in mind that where students perform grading of peer-produced UML software design as part of a Software Engineering course. Moreover, we aim to learn which features experts and students use for assessing the quality attributes of class diagrams. The study reveals that experts and students' assessment of the six quality attributes differs significantly. However, a qualitative analysis of experts and students' feedback suggests that students use similar features as experts use for assessing the quality of diagrams. Hence, peer-feedback from students can be useful in educational settings.*

This chapter is based on the following publications:

- Bilal Karasneh, Dave Stikkolorum, Enrique Larios, Michel R.V. Chaudron. **Quality Assessment of UML Class Diagrams - A Study Comparing Experts and Students-**. *MoDELS2015, Ottawa, Canada*. 2015.
- Bilal Karasneh, Michel R. V. Chaudron. **Online Img2UML Repository: An Online Repository for UML Models**. *In Proceedings of the 3rd International Workshop on Experiences and Empirical Studies in Software Modeling (EESS-MOD@MODELS 2013), pages 61-66, Miami, USA*. 2013.

Quality is a multidimensional concept, and in practice people make different interpretations of the same concept. Nowadays UML [34] is the de-facto standard for modeling software systems. UML offers a rich set of symbols for describing software. Modern software designs contain many abstraction levels, and designing them is an iterative process [95]. The collection of design documents is an important part of the system documentation which will be used and maintained for a long time by a development organization. In the software engineering class, students should understand the importance of software models and their design process.

Software design is considered as a difficult task in comparison with programming for many students. One reason is that current Integrated Development Environments (IDEs) help students to improve the quality of their code, for example using code metrics such as a maintainability index and a cyclomatic complexity. On the other hand, current UML CASE tools do not give any hints to improve models, except some layout algorithms and syntax. Although there are some proposed tools that give students some feedback about their design, these tools still suffer from many limitations, such as availability and connectivity [96][97].

During programming courses, students are taught about the quality of the source code (including for example, naming and layout conventions and API design guidelines). In software engineering courses, students are taught to understand basic modeling concepts and modeling notations. For instance, what UML diagrams are, when to use a class diagram, how to create a sequence diagram, what are elements of use case diagrams. Many teachers focus on teaching students the proper use of syntactical elements in creating UML diagrams. In both programming and modeling, the completeness and correctness are key attributes of the quality of a solution. However, there are no specific rules or guidelines for assessing quality attributes of designs. This leaves students to self-learning on how to make a good design.

For proper learning of modeling and designing, students need to get feedback from their teachers, or from peers to evaluate their design. For example: this class should have more operations, this class name should be changed, this operation should have more parameters. One way of providing feedback could be to use a method for assessing the quality of UML models. Unfortunately, currently no such method exists.

In this study, we explore the use of ISO standards for software quality as a basis for reviewing UML models. Software product quality models such as ISO/IEC 25010 [11] have categories of quality characteristics, and each characteristic is composed of a set of sub-characteristics. One difficulty with this standard is that there are many ways of interpreting every characteristic.

In this chapter, we want to study the ability of students to evaluate their designs and other students designs. Also, we want to study whether the evaluations of students are consistent with those of experts. In addition to the quantitative analysis, we do a qualitative comparison of the feedback provided by students and experts.

To study this, we asked students to perform a modeling task, and then to evaluate their models and to evaluate models of other students in terms of six quality attributes:

Understandability, Layout, Extensibility, Modifiability, Completeness, and Correctness. Then we asked five experts to evaluate students models in terms of the same quality attributes.

In this chapter we address the following research questions:

- Can we trust students assessments and for the quality of UML class diagrams? Why?
- What kind of features that experts and students focus on when they measure the quality of UML class diagrams?

The aim of this study is to empirically investigate whether students evaluations are different from experts evaluation, and what are the differences and similarities between students' feedback and experts' feedback. The differences and similarities are useful to assess if the feedback of students can be useful for improving the quality of the design.

## 7.1 Related Work

We follow the general guidelines for experimental design and analysis from [98][99]. Tichy [100] shows that there are good reasons for conducting experiments with students, for testing experimental design and initial hypotheses, or for educational purposes. Depending on the actual experiment, students may also be representative of inexperienced professionals [101].

Boustedt [102] did an empirical study on how students understand class diagram using phenomenographic investigation. He found that the purpose of class diagrams and various elements of the UML notation were understood in a varied way. He recommended that teachers should put more effort in assessing skills in proper usage of the basic symbols and models, and students should have opportunities to practice collaborative design. Our experiment is different from [102] as we ask the students to evaluate class diagram directly in terms of six quality attributes, and we ask them to give feedback and explanation about their evaluation.

Ali et. al. [103] presented the UML class diagram assessor (UCDA) that evaluates class diagrams automatically based on their structure, correctness and language used. The aim of the proposed assessor is to guide students to represent class diagram correctly. The results of our experiment are useful for the kind of assessors in [103] because from the information collected we know which kind of feedback experts and students use for describing violations of modeling conventions and/or models improvement.

Hoggart et. al. [104] found that students understand UML design in classroom settings but find it hard to apply in exercises and tasks. They proposed a tool that gives students feedback about their diagrams in comparison with a model answer proposed by the student's teacher. Generating feedback based on model answers is a bit difficult

because in modeling there is typically more than one solution. Because it is difficult to find a sufficient number of experts, we decide to explore whether the feedback from peer students can help improve model quality.

Kaneda et. al. [105] show that class diagrams reflect the cognitive structure of English based cognitive linguistic. They found that there is impedance some mismatch for understanding of class diagram by students who are not native English speakers. In our experiment, our students are a mix of various nationalities. Before admission, our students have to pass an English (TOEFL) test, and therefore we consider that their English language skills will not influence our study too much.

Aguilera et. al. [106] show that names of elements in UML diagrams have a strong influence on their understandability. They proposed guidelines for naming various kinds of elements in UML.

Selic [107] shows that understandability is the most important characteristics of models. In our study, we show the features that experts and students focus on for assessing understandability of models.

## 7.2 Experiment Design

In this section, we explain our approach, the participants of the experiment and the evaluation form that the participants use for assessing class diagrams.

### 7.2.1 Approach

We conducted an experiment in Leiden University in which both experts and students participated. We gave the students a modeling task and asked them to use the StarUML CASE tool [25] to create their models. Upon completion of the modeling task, they had to upload their models to the UML Repository and evaluate their models based on six quality attributes. Also, they had to mention their background: academia, industry or both, and their experience in UML modeling (less than one year, <1-2>, <2-5> or expert). Subsequently, they had to evaluate other students' design based on the same six quality attributes. We asked students to explain their evaluations through feedback comments for each quality attribute. Students had a trial assignment two weeks before the experiment with another modeling task. This trial is important for the students to be prepared for the experiment. It helped them in getting acquainted with the type of assignment, the tools and thereby limits the learning effects. We also asked the experts to evaluate students' models based on the six quality attributes and to give a feedback of the models and describe their evaluations.

### 7.2.2 Participant

The participants are: five experts and 46 students.

### 7.2.2.1 Experts

Five experts joined this experiment. Each expert has at least five years of experience in UML modeling and software design. Two experts are teachers of software modeling and software engineering for at least three years. One of those two experts also worked in industry. The other three experts are PhD students in the area of software engineering since 2011.

### 7.2.2.2 Students

46 master students of the ICT in business M.Sc. program<sup>1</sup> in Leiden University participated in the experiment during their course on software engineering. All of them have less than one year experience in UML modeling. Some of them have some (mostly short) background in industry, but most of them have an academic background (just finished their B.Sc. degree).

### 7.2.3 Evaluation Form

The form for evaluating class diagrams was implemented through an online system. This system showed a form that contains:

1. The number of models that were evaluated by the participant (out of 46 models).
2. An image of a student's class diagram. The image is created by some other students and has not been evaluated by the participant before.
3. A list of radio-buttons for entering assessments for 6 quality attributes. Each quality attribute can be rated on a scale ranging from 1 to 8:
  - For Understandability, Extensibility and Modifiability: (1) is difficult, (8) is easy.
  - For Layout: (1) is complex, (8) is simple.
  - For Completeness: (1) is not complete, (8) is complete.
  - For Correctness: (1) is not correct, (8) is correct.
4. A comment box. For each quality attribute participants can submit details about their evaluation using a text box. We perform a qualitative analysis of the comments provided by experts and students to figure out which features they focus on when they assess the quality of a model.
5. A submit button. Stores the assessment and navigates participants to evaluate another design.

---

<sup>1</sup>This is a degree in the Science faculty of the University of Leiden. This degree is a mix of topics from Information System, Software Engineering and Management and Business Administration.

### 7.2.4 Modeling Assignment

The modeling assignment was about a library system. The modeling assignment, evaluation form, post-questionnaire and students designs are available in the supplemental materials of the experiments [108].

## 7.3 Comparing Model Evaluation

For comparing the evaluation between experts and students, we use a Multivariate General Linear Model (MGLM). This model is used because it considers multiple dependent variables and multiple independent variables. We also use bootstrapping [109], which is a method that approximates the sampling distribution of the sample mean. In our experiment, the dependent variables are the six quality attributes, and the independent variables are the assessors (experts and students scores). We use IBM SPSS [110] as statistics tool.

### 7.3.1 Experts Evaluation and Students *self* Evaluation

Each class diagram was evaluated by at least two experts and one student (each student evaluated his/her model). In the upload form, students were asked to evaluate their models. For making this comparison, we do resampling (bootstrapping) of 1000 times of size 121 for the experts evaluation and independently the same resampling time of size 46 for the students evaluation.

### 7.3.2 Experts Evaluation and Students *peer* Evaluation

From the set of evaluations, we leave out seven class diagrams because the variation (standard deviation) of student's evaluation is high. This leaves a total of 39 class diagrams. We have 95 model-evaluations from experts per each quality attribute. Moreover, each student evaluated at least class diagrams. For each quality attribute, we have 435 evaluations in total from students.

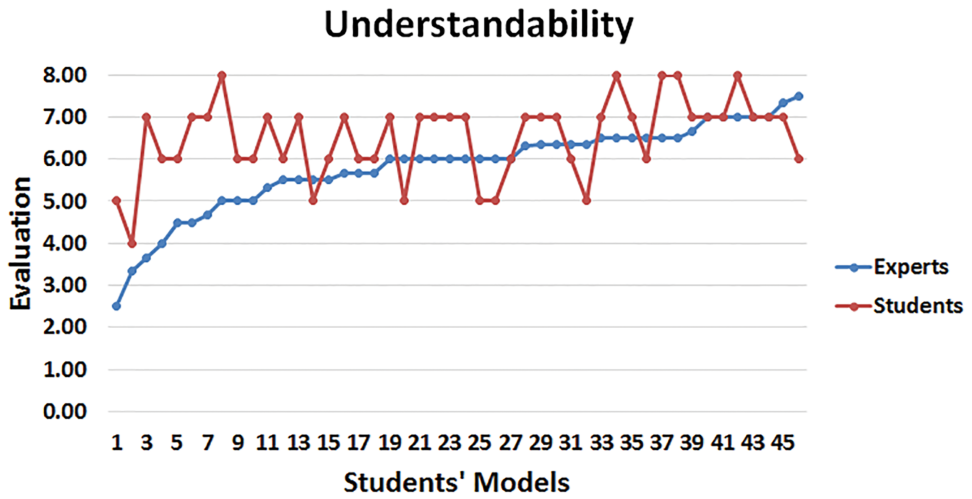
## 7.4 Results and Analysis

We did the quantitative analysis for experts and students evaluation, and qualitative analysis for their feedback and comments.

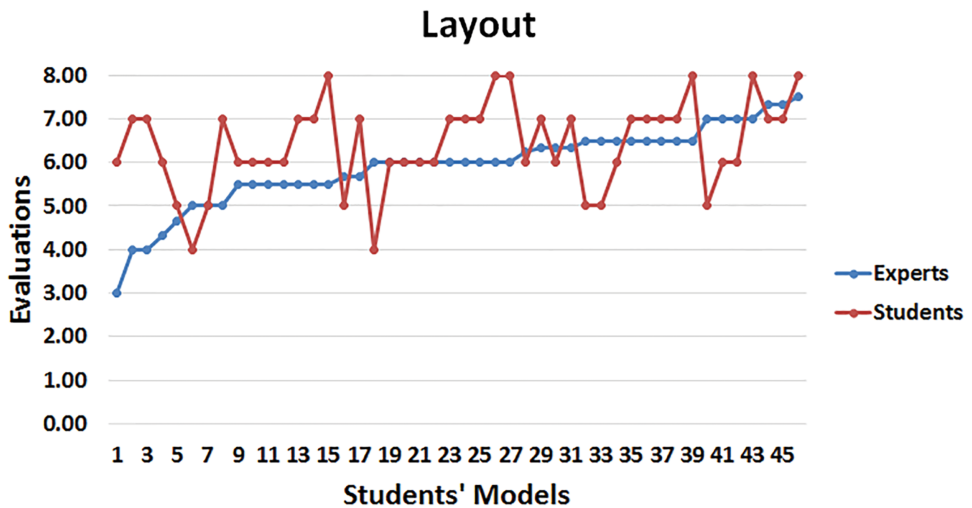
### 7.4.1 Quantitative Analysis

Figures 7.1 and 7.2 show the average evaluations of experts, and students' self-evaluation for understandability and layout respectively. The evaluation is sorted



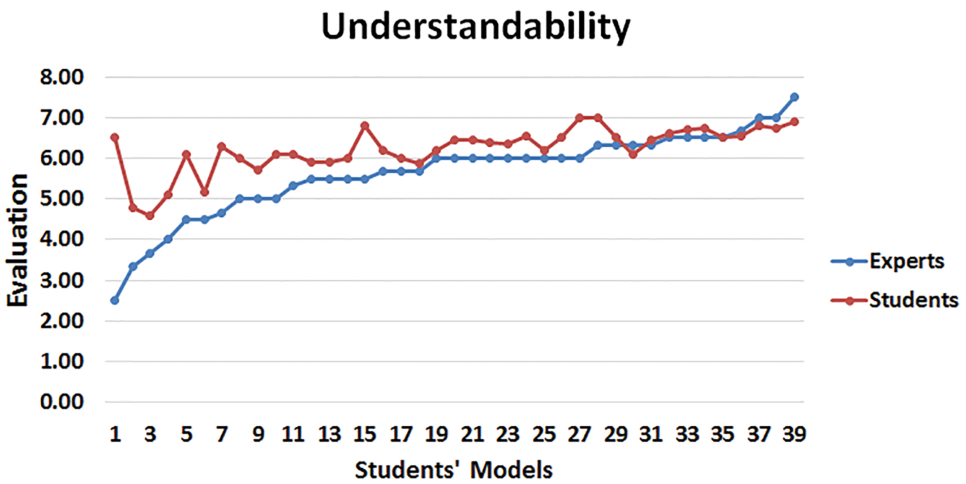


**Figure 7.1:** Experts and students evaluation (self-evaluation) for Understandability

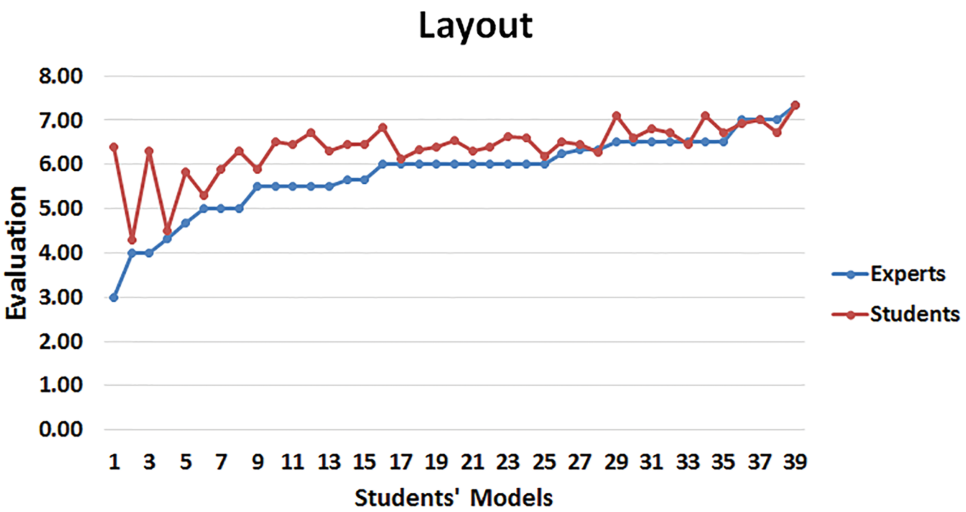


**Figure 7.2:** Experts and students evaluation (self-evaluation) for Layout

in ascending order based on experts' evaluation. Figures 7.1 and 7.2 show that experts and students differ in the most cases (high and low quality diagrams). Figures 7.3 and 7.4 show the average evaluations of experts and students' peer-evaluation for understandability and layout respectively. The evaluation is sorted in ascending order based on experts' evaluation. From Figures 7.3 and 7.4, students assessment is mostly



**Figure 7.3:** Experts and students evaluation (peer-evaluation) for Understandability



**Figure 7.4:** Experts and students evaluation (peer-evaluation) for Layout

higher than the experts' for understandability and layout, and sometimes they are close.

Table 7.1 shows the results of MGLM relating experts with students. Table 7.1 shows that there is a significant difference between expert evaluation and students self-evaluation. In addition, Table 7.1 also show that there is significant difference

**Table 7.1:** Results of Multivariate General Linear Model

Dependent Variable	Assessors		95% Sig.
<b>Understandability</b>	Experts	<i>Students Self-Evaluation</i>	0.00
		<i>Students Peer-Evaluation</i>	0.021
<b>Layout</b>	Experts	<i>Students Self-Evaluation</i>	0.003
		<i>Students Peer-Evaluation</i>	0.006
<b>Extensibility</b>	Experts	<i>Students Self-Evaluation</i>	0.003
		<i>Students Peer-Evaluation</i>	0.009
<b>Modifiability</b>	Experts	<i>Students Self-Evaluation</i>	0.001
		<i>Students Peer-Evaluation</i>	0.011
<b>Completeness</b>	Experts	<i>Students Self-Evaluation</i>	0.001
		<i>Students Peer-Evaluation</i>	0.053
<b>Correctness</b>	Experts	<i>Students Self-Evaluation</i>	0.00
		<i>Students Peer-Evaluation</i>	0.004

**Table 7.2:** Description of Experts and Students Evaluation

Dependent Variables	Experts	Self-Evaluation	Peer-Evaluation
	Mean	Mean	Mean
<b>Understandability</b>	5.71	6.51	6.22
<b>Layout</b>	5.74	6.4	6.38
<b>Extensibility</b>	5.59	6.18	6.14
<b>Modifiability</b>	5.52	6.2	6.06
<b>Completeness</b>	5.77	6.53	6.3
<b>Correctness</b>	5.07	6.31	5.86

between experts evaluations and students peer-evaluations. Table 7.2 shows the description of experts and students evaluations. From Table 7.2, all means of experts evaluations are less than the means of students evaluations in both cases (self/peer evaluation). For analyzing the evaluations, we show the correlation metrics between experts and students peer-evaluations in Table 7.3. In Table 7.3, it is possible to see many high correlations between quality attributes. First, the correlation between experts evaluation, understandability has a high correlation with all quality attributes. Second, on the student side, understandability also has a high correlation with most other quality attributes. We notice that the correlation between understandability and layout for students evaluations is higher than in experts evaluations. Third, regarding the correlation between experts' and students evaluations, the highest correlation is between experts understandability and students understandability. The second highest correlation is between experts and students evaluations for layout.

**Table 7.3:** *Correlation of Experts and Students peer-Evaluation*

	E_Unders.	E_Layout	E_Extens.	E_Modif.	E_Compl.	E_Correct.	S_Unders.	S_Layout	S_Extens.	S_Modif.	S_Compl.	S_Correct.
E_Unders.	1											
E_Layout	<b>0.74</b>	1										
E_Extens.	<b>0.77</b>	0.63	1									
E_Modif.	<b>0.84</b>	<b>0.71</b>	<b>0.88</b>	1								
E_Compl.	0.67	0.65	<b>0.70</b>	0.68	1							
E_Correct.	<b>0.71</b>	0.57	<b>0.83</b>	<b>0.85</b>	<b>0.74</b>	1						
S_Unders.	<b>0.70</b>	0.62	0.59	0.52	0.44	0.41	1					
S_Layout	0.57	0.67	0.56	0.53	0.47	0.47	<b>0.81</b>	1				
S_Extens.	0.61	0.57	0.62	0.53	0.47	0.45	<b>0.86</b>	<b>0.83</b>	1			
S_Modif.	0.60	0.62	0.57	0.56	0.51	0.37	<b>0.85</b>	<b>0.81</b>	<b>0.88</b>	1		
S_Compl.	0.54	0.37	0.49	0.46	0.58	0.51	0.62	0.56	0.59	0.56	1	
S_Correct.	0.59	0.45	0.58	0.56	0.66	0.61	0.62	0.59	0.60	0.61	<b>0.89</b>	1

### 7.4.2 Qualitative Analysis

We qualitatively analyze experts and students (peer-evaluation) comments/feedback for their evaluation. This analysis is important to see the features that experts and students use for assessing the quality of class diagrams. We discuss the feedback of three of the quality attributes: understandability, layout, and completeness. We choose these quality attributes because from Table 7.3, understandability (0.70) and layout (0.67) are the highest correlated quality attributes between experts and students peer-evaluation. In addition, understandability has the strongest correlation with others quality attributes. We choose completeness because from Table 7.1, experts and students almost differ with the significance of 0.053.

We use NVivo10<sup>2</sup> for qualitatively analyzing experts and students comments. In their comments, they explain how they evaluated models quality attributes, and features they used for their evaluation. Table 7.4 shows 11 features of models that experts and students used for assessing understandability. We observe that: (i) 64% of the features are used by both experts and students. (ii) 27% of the features are used by students but are not used by experts. (iii) 9% are used by experts and not used by students. Table 7.5 shows that experts and students used 12 features for assessing the layout, where: (i) 58% of the features are used by both experts and students. (ii) 9% of the features are used by students, but not used by experts. (iii) 33% used by experts, but not used by students. Table 7.6 shows that experts and students used 10 features for assessing completeness, where: (i) 60% of the features are used by both experts and

<sup>2</sup><http://www.qsrinternational.com/>

**Table 7.4:** *Features that Experts and Students Focus on When They Evaluate Understandability*

Features	Experts	Students
Easy to read	-	X
Completeness	X	X
Extra information	X	X
Complexity	X	X
Correctness	X	X
Data type	-	X
implementation	X	-
Layout	X	X
Class, attributes and operation names	X	X
Relationship names	-	X
Number of classes, operations, and attributes	X	X

**Table 7.5:** *Features that experts and students focus on when they evaluate Layout*

Features	Experts	Students
Classes Hierarchy, alignment	X	X
Classes with similar size	X	-
Complexity	X	X
Number of classes, attributes and operations	-	X
Distance between Classes	X	X
Rectilinear edges and diagonal edges	X	-
Line Style (overlapping, crossing, bend)	X	X
Good Class name	X	-
Neat or chaotic structure	X	X
Easy to read	X	X
Same Layout for same/All relationships	X	-
Extra information	X	X

students. (ii) 30% of the features are used only by students. (iii) 10% are used only by experts. Although we see experts focus on more features in Table 7.5, it may be that students are interested in many of the same features – yet they do not mention them clearly in their feedback.

**Table 7.6:** *Features that experts and students focus on when they evaluate Completeness*

Features	Experts	Students
Model Abstraction	-	X
Functionality	X	X
Strange Relationships	X	X
Missing Classes, Attributes, and operations	X	X
Data types	-	X
Multiplicity	X	X
Functions Parameters	-	X
Relationship names	X	X
Requirements	X	X
Model Semantics	X	-

## 7.5 Discussion

From the MGLM results in Table 7.1, we conclude that there is a significant difference between the evaluation of experts and students. The results also show that peer-evaluation of students is closer to the evaluation of experts than self-evaluation (because the mean difference is bigger between experts and self-evaluation than with peer-evaluation for all quality attributes as shown in Table 7.2). We explain this by the different viewpoints in the peer-evaluation. Different points of view may have caused different evaluations that on average became more reliable, or at least better than the self-evaluation.

The qualitative analysis of experts' and students' comments shows that students use most features that experts use for assessing the quality of class diagrams. In Table 7.4, 7.5 and 7.6 we summarize the features that experts and students use for assessing understandability, layout and completeness respectively. In the qualitative analysis, we only take into account the issues that can be clearly identified in the feedback. We notice that feedback from experts is more specific than that from students. For example, some students mentioned they did not like a class, but they did not mention what was the problem with this class: name, size, position, etc. However, this general feedback can still be useful because it can be considered as general hints that direct students to a particular area where they still themselves need to find out what needs to be improved.

We conclude that students largely use similar features for assessing the quality of class diagram as experts use. Hence, their feedback is useful for improving their models. So we expect that if students exchange their feedback about their models, this will be a valuable source of feedback for learning and improving their models. Also, we expect that students can make a better evaluation if they do this in a group because they can then discuss their different viewpoints and improve their evaluation.

In addition, students (peer-evaluation) is not so close to expert evaluation. Students seem reluctant to fail fellow students.

## 7.6 Threats to Validity

In this section, we discuss the threats to validity of our study.

### 7.6.1 Internal Validity

We ensured that students are familiar with class diagrams. The experiment was conducted at the end of the Software Engineering course where they had a trial two weeks before the experiment. The participants did not know the aim of our experiments, nor the measures that we are looking for, in order to avoid their expectations from biasing the results.

### 7.6.2 External Validity

There were 46 students participants in the experiment. To mitigate their representativeness, we only address their experience level with UML, and their background (academic, industry or both). About the modeling task, we chose a system from an application domain that should be familiar to students, which is library system.

## 7.7 Conclusion and Future Work

In this chapter, we presented an experiment that investigates the difference between experts and students in assessing the quality of UML class diagram empirically. We made two comparisons: first between experts and students' self-assessment. Second, between experts and students peer- assessment. We use the Multivariate General Linear Model as a statistical method for making those comparisons. The results show that experts and students (self- assessment) are different in terms of means (95% significance). The students self- assessments are higher than experts assessments in terms of mean for the quality attributes used in the experiment. The results also show that experts and students (peer-evaluation) are different in terms of mean (95% significance). The students peer- assessments are higher than experts assessments in terms of mean for all quality attributes used in the experiment. Analyzing the correlation between experts' assessments and students peer- assessments shows that understandability is the highest correlated quality attribute, and that layout is the second highest. The correlation also shows that understandability is correlated with most of the other quality attributes based on both experts assessments and students assessments.

We did a qualitative analysis of experts' feedback and students feedback in peer-assessments. From this, we observe that students mostly use similar features as experts for their assessments. So we conclude that feedback from students is valuable and can be useful for other students for improving their designs.

In the future, we are planning to replicate the experiment and ask students to assess the quality of class diagram in groups. We believe that having an online community for students where they can exchange their models, and their feedback is very useful for improving modeling education. So we are establishing this community with the collaboration of some experts. From this community, students and experts can upload their models and exchange their feedback.