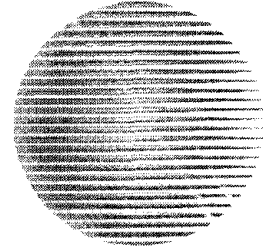


01 **Research**

02
03 **Assessing the capability of**
04 **internal metrics as early**
05 **indicators of maintenance effort**
06 **through experimentation**
07
08
09



10
11 Marcela Genero Bocco^{1,*}, Daniel L. Moody^{2,3} and Mario Piattini¹
12

13 ¹*Department of Computer Science, Universidad de Castilla-La Mancha, Ciudad Real, Spain*

14 ²*Department of Computer Science, University of Iceland, Reykjavik, Iceland*

15 ³*Department of Cybernetics, Czech Technical University (CVUT), Prague, Czech Republic*
16
17

18
19
20 **SUMMARY**

21
22 **The complexity of software artifacts is widely believed to be an important determinant of maintenance**
23 **effort. This paper conducts an experimental analysis of the impact of complexity on the maintenance of the**
24 **Unified Modeling Language (UML) class diagrams. This represents an analysis of the effect of an internal**
25 **quality attribute on an external quality attribute. A range of complexity metrics are proposed based on an**
26 **ontological analysis of the UML language and previous research. The relative influence of these metrics on**
27 **maintenance effort is then evaluated using a laboratory experiment. A within-subjects design was used, with**
28 **subjects required to modify a range of UML class diagrams with different levels of complexity. Only two of**
29 **the metrics emerged as significant determinants of maintenance effort: number of methods and number of**
30 **associations. Together these explain around 28% of the variation in maintenance effort. While these findings**
31 **are encouraging, further research is necessary to explore the ability of these metrics to predict maintenance**
32 **effort. Copyright © 2005 John Wiley & Sons, Ltd.**

33 **KEY WORDS:** complexity; Unified Modeling Language (UML); class diagram; object-oriented (OO) modeling;
34 experimental research; software quality
35
36
37
38

39
40 *Correspondence to: Marcela Genero Bocco, Paseo de la Universidad, N^o 4, 13071 Ciudad Real, Spain.

41 †E-mail: marcela.genero@uclm.es

42 Contract/grant sponsor: Consejería de Ciencia y Tecnología de la Junta de Comunidades de Castilla-La Mancha; contract/grant
43 number: PCC-03-003-1

44 Contract/grant sponsor: Dirección General de Investigación del Ministerio de Ciencia y Tecnología; contract/grant number:
TIC2003-07804-C05-03



01 **1. INTRODUCTION**

02
03 **1.1. Maintainability of information systems**

04
05 The adaptability of a system to changes in its environment is considered to be one of a system's
06 most important characteristics [1,2]. In practice, maintenance represents the major source of work
07 for software organizations [3]. It is estimated that 61% of the professional life of programmers is
08 devoted to maintenance [4]. Also, maintenance costs represent two thirds of the total lifetime costs of
09 information systems in commercial mainframe environments [5]. For this reason, *maintainability* is an
10 important issue in the design of information systems. However in practice, it is often given relatively
11 little attention in the interests of completing development projects 'on time, on budget'. The result
12 is that information systems often have substantially reduced operational lifetimes, and are sometimes
13 obsolete before they are delivered [6,7].

14 Maintainability is generally accepted to be an important aspect of software quality, and is one of the
15 six software quality characteristics defined in the international standard for software quality: ISO/IEC
16 9126 [8]. Maintainability is defined as 'the capability of the software to be modified' and consists of
17 four subcharacteristics:

- 18 • *analyzability*, the capability of the software to be diagnosed for deficiencies or causes of failures
- 19 in the software, or for the parts to be modified to be identified;
- 20 • *changeability*, the capability of the software product to enable a specified modification to be
- 21 implemented;
- 22 • *stability*, the capability of the software to minimize unexpected effects from modifications of the
- 23 software; and
- 24 • *restability*, the capability of the software to enable modified software to be validated.

25
26 Maintainability is an external quality characteristic, as it relates to the external behavior of
27 the product [9]. However, internal quality characteristics (characteristics of the product which are
28 not visible to the user) are the key to satisfying external quality characteristics. Internal quality
29 characteristics define the 'hows' by which external quality characteristics or 'whats' will be met [10].
30 Understanding the relationship between internal and external quality characteristics is essential for
31 developing software which satisfies user needs [11].

32
33 **1.2. The effect of complexity on maintenance effort**

34
35 Software complexity has long been recognized as a determinant of maintenance effort [12–20].
36 There are several reasons why complexity affects maintenance effort.

- 37 • Complex systems are more difficult to understand, and understanding a system is a prerequisite
- 38 for making changes to it [21,22]. Empirical studies show that humans have a strictly limited
- 39 capacity for processing information [23–25]. Once the level of complexity exceeds this limit, a
- 40 state of *information overload* ensues and comprehension degrades rapidly [26–29]. Both field
- 41 and experimental studies show that information overload results in reduced performance on
- 42 cognitive tasks [24,29]. This relates to the *analyzability* subcharacteristic in ISO/IEC 9126 [8].
- 43 • More complex systems will tend to have higher *entropy* because of the number of interacting
- 44 components, which increases their volatility and susceptibility to change [30,31]. This increases



01 the number of changes required, thus contributing to maintenance effort. This relates to the
 02 *stability* subcharacteristic in ISO/IEC 9126 [8].

- 03 • As systems increase in complexity, the impact of changes becomes more complex and
 04 unpredictable, which increases the effort required to implement changes [32–34]. Changes to
 05 one part of the system are more likely to affect many other parts of the system, often referred
 06 to as the ‘ripple effect’ [2,32,35]. This relates to the *changeability* subcharacteristic in ISO/IEC
 07 9126 [8].

08 1.3. Objectives of this research

09 Most previous research on maintainability [16,17,36–38] has focused on the maintainability of the
 10 final system (i.e., code-level maintenance). Our focus is different in the sense that we are interested in
 11 evaluating maintainability from the early stages of software development. For that reason, this paper
 12 investigates the relationship between complexity (an internal quality characteristic) and maintainability
 13 (an external quality characteristic). The specific focus of this research is on the Unified Modeling
 14 Language (UML) *static structure models*, commonly called *class diagrams* [39]. We evaluate the
 15 effect of a set of metrics for measuring complexity of UML class diagrams [40,41] on maintenance
 16 effort through a controlled experiment.

17 We performed a previous controlled experiment [42], pursuing a similar objective. In that
 18 experiment, as in this, the independent variable was UML class diagram complexity. In the previous
 19 experiment the dependent variables were three maintainability subcharacteristics (understandability,
 20 analyzability and modifiability) measured by means of subjective ratings on a seven-point scale.
 21 Even though the results obtained in the previous experiment (through correlation analysis) reflect that
 22 the metrics proposed were related to class diagram maintainability, we are aware that the dependent
 23 measures were subjective and therefore less reliable. Therefore, we decided to carry out another
 24 experiment measuring the dependent variable in a more objective way, which is the experiment
 25 presented in this paper.

26 The outline of the paper is as follows:

- 27 • Section 2 defines a set of candidate metrics for measuring complexity of UML class diagrams;
- 28 • Section 3 describes the experimental design used to evaluate the metrics;
- 29 • Section 4 presents the results of the experiment and evaluates their validity; and
- 30 • Section 5 summarizes the findings and discusses the main contributions of the research.

31 2. METRICS FOR MEASURING COMPLEXITY OF UML CLASS DIAGRAMS

32 This section defines a set of metrics for measuring complexity of UML class diagrams.

33 2.1. UML class diagrams

34 UML defines an industry standard approach to modeling object-oriented (OO) systems [39]. So far
 35 there has been very little research into the maintainability of UML models. However, UML is usually
 36 used in a *model-driven* development approach, so changes to systems should be driven by changes to



01 the underlying models. Class diagrams define the architecture of an OO system—the objects and their
 02 behavior—and represent the starting point for all changes to an OO system.

03
 04 **2.2. Metrics for UML class diagrams**

05
 06 A set of candidate metrics was developed for measuring the complexity of UML class diagrams.
 07 These have been published in [41,42], but are summarized here to provide context for understanding
 08 the experiment. The metrics were developed based on:

- 09
 10
 11
- an analysis of the ontological structure of UML class diagrams;
 - a review of the OO metrics literature.

12
 13 *2.2.1. Ontological analysis*

14
 15 According to general systems theory [21,31], the complexity of a system is determined by the *number*
 16 *and variety of elements and the number and variety of relationships between them*. In developing
 17 measures of complexity for UML class diagrams, a natural starting point therefore is to define the types
 18 of elements and relationships such diagrams may contain: the *ontological structure* of the notation [2].
 19 At the requirements level, a UML class diagram may consist of the following constructs [39].

- 20
 21
 22
 23
 24
 25
 26
 27
 28
 29
 30
- Classes, which consist of:
 - attributes (with their name, type and visibility scope);
 - methods (with their name, list of parameters, return type, visibility and scope).
 - Relationships, which may be of four types:
 - association;
 - aggregation: this includes both composition relationships and part-whole relationships as defined in UML 1.4;
 - generalization;
 - dependency.

31
 32 Potential measures of UML class diagram complexity include counts of each construct. This results
 33 in seven (simple) measures of complexity.

- 34
 35
 36
 37
 38
 39
 40
 41
 42
 43
 44
- M1. Number of classes: the number of classes in a class diagram. This metric corresponds to OA1 in [43], DSC (Design Size in Classes) in [44] and C_s in [45]. This is the most common measure of class diagram complexity used in the literature.
 - M2. Number of methods: the total number of methods defined in a class diagram, including those defined at class and instance level, but not including inherited methods (this would lead to double counting). This is a modification of the WMC (Weighted Method per Class) metric [46], in which: (a) the weightings of all methods are 1; and (b) the value is totaled across all classes (as WMC is a class-level metric). Weightings are not applicable at the requirements stage because the code for methods is not available. This also corresponds to the NOM (Number of Methods) metric in [17,44] and the combination of NIM (Number of Instance Methods) and NCM (Number of Class Methods) in [47] (although all of these metrics are defined at the class level).





- 01 M3. Number of attributes: the total number of attributes defined in a class diagram, including those
 02 defined at class and instance level, but not including inherited attributes or attributes defined
 03 within methods. This is a generalization of the NOA (Number of Attributes) metric [45] to the
 04 class diagram level.
- 05 M4. Number of associations: the number of association relationships in a class diagram. This is a
 06 generalization of the NAS (Number of Associations) metric [48] to the class diagram level.
- 07 M5. Number of aggregations: the number of aggregation relationships in a class diagram.
- 08 M6. Number of dependencies: the number of dependency relationships in a class diagram.
- 09 M7. Number of generalizations: the number of generalization relationships in a class diagram (each
 10 parent-child pair in a generalization relationship).

11 Together these metrics enable measurement of 'the number and variety of elements and the number
 12 and variety of relationships between them' in a UML class diagram, as required by our definition of
 13 complexity.
 14

15 2.2.2. Review of OO metrics literature

16
 17 A thorough review of the OO metrics literature was conducted in order to identify additional
 18 metrics for measuring the complexity of UML class diagrams at the requirements stage. Most
 19 research on OO software measurement has focused on the measurement of code and detailed design
 20 artifacts [9,45,49,50]. However, there are now a significant number of proposals for metrics which can
 21 be applied at the analysis level [17,43-48,51-55]. Of these, most are defined at the class level, but
 22 some can be generalized to the class diagram level. A number of proposals evaluate the complexity of
 23 the UML notation itself (*meta-level* complexity), so are not applicable in this research [53-55]. Based
 24 on these previous proposals, we identified four additional metrics.
 25

- 26 M8. Number of generalization hierarchies (NGenH): the number of distinct generalization hierarchies
 27 in a class diagram. This corresponds to the OA2 metric in [43] and the NOH metric in [44].
- 28 M9. Number of aggregation hierarchies (NAggH): the number of distinct aggregation hierarchies in
 29 a class diagram. This is a generalization of M8 to aggregation hierarchies.
- 30 M10. Maximum depth of inheritance (MaxDIT): the maximum DIT value for all classes in a class
 31 diagram. The DIT value for a class in a generalization hierarchy is defined as the length of the
 32 longest path from the class to the root of the hierarchy. This is a generalization of the DIT metric
 33 in [45] to the class diagram level.
- 34 M11. Maximum height of aggregation (MaxHAgg): the maximum HAgg value for all classes in a
 35 class diagram. The HAgg value for a class in an aggregation hierarchy is defined as the length
 36 of the longest path from the class to the leaves. This is a generalization of M10 to aggregation
 37 hierarchies.
 38

39 All of the metrics (M1-M11) satisfy the ISO/IEC 9126 required properties for measures used for
 40 comparison, as they are objective, empirical and repeatable [8]. They are also simple, which Fenton and
 41 Neil [56] argue is a desirable property for software metrics in general. Another issue is that the proposed
 42 metrics could also allow designers and conceptual modelers to make better decisions when they are
 43 modeling class diagrams. In particular, they allow selection among several class diagram alternatives
 44 with equivalent semantic content. Moreover, they could be useful for predicting, for example, the

01 maintenance effort needed for modifying UML class diagrams. This potential need is to be investigated
02 empirically and is the main objective of the current study reported here.

03
04 2.2.3. *Predicted impact of the metrics on maintenance effort*
05

06 The metrics M1–M7 (which represent counts of individual UML constructs) are expected to increase
07 maintenance effort, as adding a construct of any type will increase overall complexity of the class
08 diagram. However, some types of elements are likely to affect complexity more than others e.g., adding
09 a new class is likely to increase complexity more than adding a new attribute. Thus, we would expect
10 overall complexity of a class diagram to be a *weighted sum* of these measures.

- 11 • We predict that the number of classes (M1) will have the greatest effect on maintenance effort.
12 The reason for this is that the number of classes corresponds to the number of distinct concepts or
13 'things' in the model. Based on theories of human cognition, increasing the number of concepts
14 will increase the effort required to understand the model, which in turn will increase maintenance
15 effort [23–25,57,58]. This is also consistent with software estimation models, in which the
16 number of entities (which correspond to classes in an OO environment) is the most influential
17 determinant of software size [59–61].
- 18 • We predict that the metrics corresponding to the number of relationships (M4–M7) will be the
19 next most influential. According to *cognitive load theory*, cognitive load is increased by the
20 level of interactivity among elements (the number of relationships among concepts) [62,63].
21 Increasing cognitive load reduces understanding and is therefore likely to increase maintenance
22 effort.
- 23 • The metrics M2 and M3 provide measures of the internal complexity of classes. One would
24 expect a class diagram with high internal complexity (many methods and attributes per class) to
25 be more difficult to understand than a model with the same number of classes and relationships
26 (external or system complexity) but lower internal complexity. However, this is expected to have
27 a weaker impact than external complexity.
28

29 Because generalization hierarchies and aggregation hierarchies define subsystems within the class
30 diagram, we predict that the metrics M8 and M9 will *reduce* complexity and therefore reduce
31 maintenance effort. Given two models with the same number of classes, we would expect one with
32 more generalization and/or aggregation hierarchies to be easier to understand, because hierarchical
33 structures facilitate 'chunking' and thus reduce load on working memory [23,24,64].

34 Finally, we predict that M10 and M11 will increase maintenance effort. The rationale for
35 this is that human classification hierarchies tend to be limited in depth, usually to around three
36 levels [27]. Thus, very deep hierarchies may overtax powers of human information processing [65,66].
37 The literature has shown conflicting results for the effects of inheritance hierarchies on maintenance.
38 For example, Basili *et al.* [67] found that the larger the DIT value, the greater the probability of fault
39 detection, Daly *et al.* [68] found that systems with three levels of inheritance had lower maintenance
40 time than those with no inheritance and Cartwright [69] found that three levels of inheritance reduced
41 maintenance time. On the other hand, Unger and Prechelt [70] found that inheritance depth had no
42 effect on maintenance effort, Harrison *et al.* [71] found that systems with no inheritance are easier to
43 modify than systems with three or five levels, while Poels and Dedene [72] found that extensive use
44 of inheritance leads to models that are more difficult to modify.



01 In summary, we predict that the metrics will influence maintenance effort in ascending order of
02 influence:

- 03 • M1 (corresponding to the number of elements or concepts);
- 04 • M4, M5, M6, M7 (corresponding to the number of relationships or interactivity);
- 05 • M2, M3 (corresponding to the internal structure of elements);
- 06 • M10, M11 (corresponding to depth of hierarchies);
- 07 • M8, M9 (expected to reduce complexity and therefore maintenance effort).

08
09 We are conscious that these predictions must be empirically evaluated to extract final conclusions.

12 3. RESEARCH METHODOLOGY

14 3.1. Research objectives

15
16 The purpose in generating the initial set of metrics was to be as complete as possible—to try to measure
17 all possible aspects of UML class diagram complexity. However, some of these metrics are likely
18 to be more influential than others in determining maintenance effort. The purpose of this research is to
19 empirically evaluate the relative influence of all of these metrics on maintenance effort and use this to
20 build a prediction model. The validity of each of the predictions stated in Section 2.2.3 (related to the
21 order of influence) is out of the scope of the paper and will be pending for further research.

23 3.2. Experimental design

24
25 The experimental design is summarized in Figure 1. A *within-subjects* design was used, in which each
26 participant was given nine UML class diagrams of different levels of complexity. Participants were
27 required to make a number of prescribed changes to each model and to record how long it took them
28 to complete the task.

30 3.3. Participants

31
32 There were 30 participants in the experiment: 20 final-year Computer Science students and 10 Software
33 Engineering academics at the University of Castilla-La Mancha. All participants had prior experience
34 using UML to develop software. The student participants had all completed at least two software
35 engineering courses in UML. All participants were also given an intensive ‘refresher course’ in UML
36 class diagrams prior to the experiment.

38 3.4. Materials

39
40 Participants were given a guide summarizing UML notation and nine UML class diagrams representing
41 different application domains. The experimental models had different levels of complexity, covering
42 a broad range of metric values (see Table I). Each experimental model included a brief description of
43 the underlying problem domain and an associated maintenance task. An example of one of the class
44 diagrams used (CD5) and maintenance task is included in Appendix A (all experimental materials

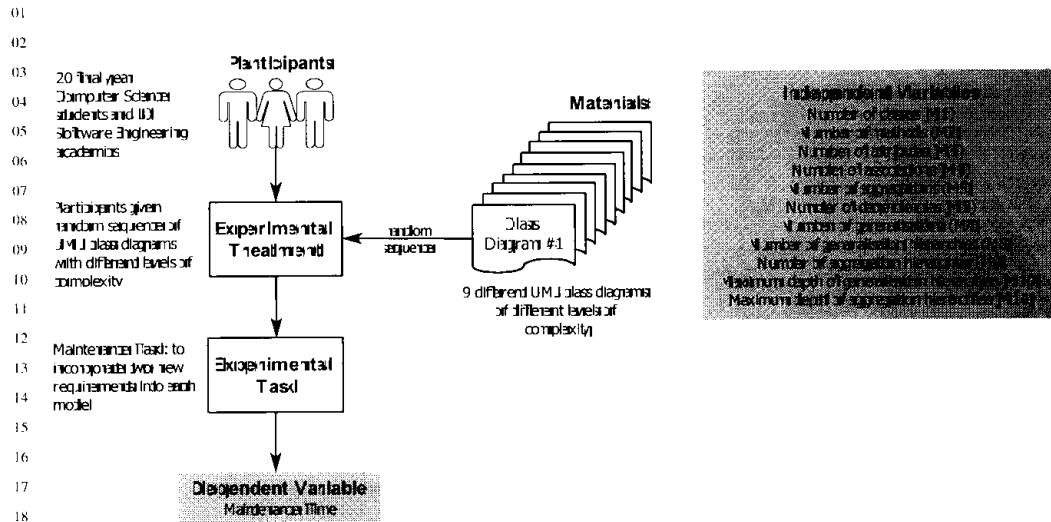


Figure 1. Summary of experimental design.

Table I. Complexity metrics for the experimental models.

22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

	M1	M2	M3	M4	M5	M6	M7	M8	M9	M10	M11
CD1	7	11	22	1	0	0	5	1	0	2	0
CD2	8	12	31	1	6	0	1	1	1	1	2
CD3	3	17	24	2	0	0	0	0	0	0	0
CD4	10	12	21	15	3	0	0	0	2	0	1
CD5	9	19	29	3	3	0	3	1	3	2	1
CD6	7	16	7	6	0	0	0	0	0	0	0
CD7	23	33	66	4	5	2	16	3	2	3	3
CD8	20	30	65	6	5	0	14	3	4	3	2
CD9	23	65	80	20	3	2	3	1	3	2	3
AVG	12.22	23.89	38.33	6.44	2.78	0.44	4.67	1.11	1.67	1.44	1.33
STD DEV	7.63	17.3	25.27	6.65	2.33	0.88	6.12	1.17	1.5	1.24	1.22

38 are available upon request to the corresponding author). The experimental models were given to
39 participants in random order, to control for learning effects.

40 3.5. Experimental task

41 Each experimental model had an associated experimental task, consisting of two new requirements to
42 be incorporated into the model. The activities that the subjects had to perform fit into the category
43
44



01 of *enhancive maintenance*, according to the taxonomy of software maintenance activities proposed
 02 by Chapin *et al.* [73], as they involved adding new functionality to the system. The tasks defined for
 03 each class diagram were all of equivalent complexity (in terms of number of diagram elements added,
 04 modified or deleted), so the only source of variation in effort to perform the tasks should be model
 05 complexity. An example of one of the experimental tasks is included in Appendix A. Each subject
 06 had to modify the class diagrams according to the specifications and to record the time taken for each
 07 task—they were asked to write down their start time (to the nearest second) prior to beginning the task,
 08 and then to write down their finish time at the end. Each subject was allowed an elapsed calendar time
 09 of one week to carry out the experimental task.

11 3.6. Independent variable

12
 13 The independent variable was the complexity of the experimental models. This was measured by
 14 the metrics previously defined—these represent empirical indicators of the underlying theoretical
 15 construct. While there is only one *latent* independent variable (complexity), there are 11 *observed*
 16 independent variables (M1–M11). The *levels* of the independent variable are embodied in the different
 17 experimental class diagrams: each model represents a different combination of metric values (Table I).

19 3.7. Dependent variable

20
 21 The dependent variable in this experiment was maintenance effort, which was defined as ‘the effort
 22 required by a participant to make the prescribed changes to an experimental model’. This was measured
 23 by the time taken to perform the experimental task, expressed in seconds. This includes the time
 24 to comprehend the class diagram, to analyze what changes were required and to complete them.
 25 The dependent variable thus incorporates the subcharacteristics of *analyzability* and *changeability*
 26 defined in ISO/IEC 9126 [8]. Time taken is an objective measure of the effort required to perform the
 27 experimental task, and is one of the most commonly used measures of effort in performing cognitive
 28 tasks [27].

30 3.8. Summary of theoretical model

31
 32 Figure 2 summarizes the theoretical model tested by this experiment.

- 34 • Circles indicate latent variables (theoretical constructs) while rectangles indicate observed or
 35 measured variables (empirical indicators).
- 36 • The solid lines indicate *causal relationships* hypothesized between latent variables. The direction
 37 of the arrows indicates direction of causation.
- 38 • The dotted lines indicate *measurement relationships* between latent variables and their empirical
 39 indicators. The direction of the arrows indicates whether they are reflective or formative. The \ominus
 40 symbol indicates an inverse indicator (where the measure is inversely related to the underlying
 41 construct).

42
 43 As shown in Figure 2, we hypothesize a causal relationship between complexity of UML class
 44 diagrams and maintenance effort. The metrics defined in Section 2 are all *formative* measures of

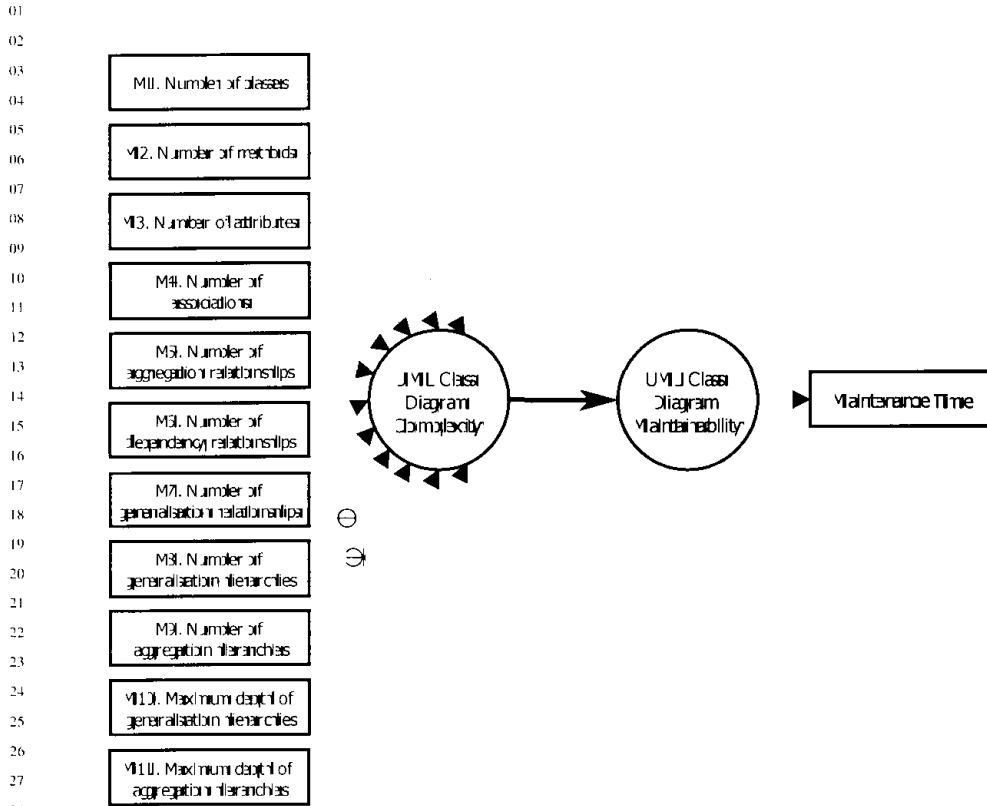


Figure 2. The *a priori* theoretical model.

32 complexity, while maintenance time is a reflective measure of effort. The purpose of this research
33 is to analyze the effects of different aspects or dimensions of complexity (represented by the different
34 metrics) on maintenance effort.
35

36 3.9. Hypotheses

37
38 We want to investigate whether the complexity metrics proposed can be used as predictors of
39 maintenance effort, so we formulate the following null and alternative hypotheses.
40

- 41 • H_0 : There will be no relationship between the complexity metrics (M1–M11) and maintenance
42 time.
- 43 • H_1 : There will be positive relationships between (M1–M7, M10, M11) and maintenance time
44 and negative relationships between (M8, M9) and maintenance time.



01 For simplicity, we express these as single hypothesis, but in reality they incorporate 11 distinct null
02 and alternative hypotheses as shown in Figure 2.

03
04

05 4. RESULTS

06

07 This section begins with the data collection and continues with the analysis the empirical data and the
08 interpretation of the analysis results.

09

10 4.1. Data collection

11

12 Data were collected by calculating the value of the metrics for each experimental data model and
13 entering start and finish times for each experimental task. Metrics were calculated automatically
14 using a tool, while maintenance time was calculated from start and finish times using a spreadsheet.
15 The experiment resulted in a total of 270 data points (30 participants \times 9 experimental tasks) for the
16 dependent variable, which formed the basis for the data analysis.

17

18 4.2. Statistical significance versus practical meaningfulness

19

20 In interpreting empirical results, it is important to distinguish between statistical significance and
21 practical meaningfulness of results [74,75]. *Statistical significance* measures whether the observed
22 relationship between variables is likely to be a 'real' relationship or if it is likely to be the result of
23 random variation. Statistical significance is measured by alpha (α), which represents the probability
24 that the result could have occurred by chance.

25 *Effect size* has been suggested by a number of researchers as an index of practical
26 meaningfulness [74,76–78]. In studies involving relationships between variables, effect size is most
27 commonly measured by the standardized regression coefficient (β). In terms of deciding whether an
28 effect is 'meaningful' or not, the most frequently used guidelines are those proposed by Cohen [74]:

29

- 30 • $|\beta| < 0.1$, not meaningful;
- 31 • $|\beta| \geq 0.1$, small;
- 32 • $|\beta| \geq 0.3$, medium;
- 33 • $|\beta| \geq 0.5$, large.

34

35 4.3. Construct validity

36

37 The metrics defined in Section 2.2 are proposed as measures of a single underlying construct: UML
38 class diagram complexity. Construct validity is one of the most important constructs in measurement
39 and refers to the level of confidence that empirical indicators reflect the construct they purport to
40 measure [75]. Factor analysis is the preferred technique among researchers for evaluating construct
41 validity. However in this case, the sample size was too small, so *inter-item correlation analysis* was
42 used instead. Out of 55 inter-item correlations between all metrics, 49 were significant at the 0.01
43 level, two more were significant at the 0.05 level and four were non-significant. Convergent validity
44 (measured as the average correlation of each metric with the other metrics) for nine of the 11 metrics



01 was large (> 0.5) and medium for the remaining two ($r > 0.3$). This provides strong evidence that they
 02 are all measuring the same underlying construct: each of the metrics measures a different dimension of
 03 the construct.

04 The construct validity of the metrics used for the independent variable is guaranteed by the
 05 DISTANCE framework [79] which was used for their theoretical validation [40].

07 4.4. Regression analysis

08
 09 Multivariate regression analysis is a technique for analyzing causal relationships between continuous
 10 variables. All of the independent and dependent variables in this experiment are continuous variables so
 11 regression analysis is an appropriate technique for analyzing the relationships between them. The result
 12 of regression analysis is a prediction model relating the independent and dependent variables.

13 *Stepwise regression analysis* was used to evaluate the relative influence of the metrics in determining
 14 maintenance effort. This is a statistical technique used in exploratory research to discover the
 15 relationship between variables when a range of different causal factors have been hypothesized [75].
 16 This is an appropriate technique to use here as the research is exploratory. Stepwise regression analysis
 17 may be carried out in the following two ways [80].

- 19 • Top down: start by including the most highly correlated variable in the regression equation and
 20 then adding the variable with the highest partial correlation until there is no significant change
 21 in r^2 .
- 22 • Bottom up: start by including all variables in the regression equation and then removing the least
 23 influential variable in each stage.

24
 25 In this paper, we use top down analysis, but both methods lead to the same results.

27 4.4.1. Regression analysis stage 1

28
 29 Bivariate correlation analysis was first carried out between all observed independent variables (the
 30 complexity metrics) and the dependent variable (maintenance time). All correlations were found to
 31 be statistically significant except for number of dependencies (M6). The number of methods (M2)
 32 metric emerged as the most influential variable, with a correlation coefficient of 0.499. Regression
 33 analysis was therefore carried out using M2 as the sole independent variable and maintenance time as
 34 the dependent variable. The results of the analysis were:

- 36 • $r^2 = 0.249$;
- 37 • $p = 0.000$;
- 38 • $\beta = 0.499$.

39
 40 The r^2 statistic shows that M2 explains almost 25% of the variation in maintenance time.
 41 The standardized regression coefficient is 0.499, which represents a medium effect [74]. The regression
 42 equation which results is

$$43 \text{Maintenance Time} = 2.07 + (0.05 \times \text{Number of Methods}) \quad (1)$$



01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44

Table II. Separate effects of metrics on maintenance time.

Metric	Significance (p)	Effect size (β)
Number of methods (M2)	0.000	0.388
Number of associations (M4)	0.002	0.217

4.4.2. Regression analysis stage 2

In the second stage of the analysis, the metric with the highest partial correlation (while controlling for M2) was entered into the regression model. Partial correlation analysis revealed that number of associations (M4) was the only metric which had a significant partial correlation with Maintenance Time. Most of the remaining metrics were negatively correlated with the dependent variable after controlling for M2. M4 was therefore entered into the regression model. The results of the analysis were:

- $r^2 = 0.277$;
- $p = 0.000$;
- $\beta = 0.533$.

This shows that number of methods (M2) and number of associations (M4) together explain around 28% of the variation in maintenance effort, which is a significant increase in the predictive power of the model. The standardized regression coefficient is 0.533, which represents a large effect [74]. The regression equation which results is

$$\text{Maintenance Time} = 1.94 + (0.04 \times \text{Number of Methods}) + (0.09 \times \text{Number of Associations}) \quad (2)$$

The coefficients in the regression equation measure the 'raw' effect of each independent variable on the dependent variable. The regression coefficient for M2 is 0.04, which means that each additional method increases maintenance time by 2.4 seconds or 1%. The regression coefficient for M4 is 0.09, which means that each additional association increases maintenance time by 5.4 seconds or 2.27%.

The separate effects of each independent variable on Maintenance Time are shown in Table II. Both metrics were found to be significant predictors of maintenance time ($\alpha < 0.01$), with M2 having a medium effect and M4 having a small effect.

While each association has a larger incremental effect on maintenance time than each method (as measured by the regression coefficient), the much larger number of methods (on average, there were about six times as many methods as associations in the experimental models) means that they have a much greater effect on maintenance time overall (as measured by the effect size). While the regression coefficient for M2 is less than half the size of that for M4, the standardized regression coefficient (β) is almost twice the size.

4.4.3. Regression analysis stage 3

In the third stage of analysis, partial correlation analysis was carried out on the remaining metrics to see if there was a significant correlation with maintenance time after controlling for M2 and M4.



01 None of the remaining metrics had significant effects on maintenance time, nor were they even close
02 (p values were all greater than 0.6). This provides strong evidence that the stepwise regression analysis
03 is complete, and that the regression model defined in (2) is fully specified.

06 4.5. Validity analysis

07 In this section we analyze the internal, external and statistical validity of the experiment.

11 4.5.1. Internal validity

12 The following variables were controlled as part of the experiment.

- 15 • Participant characteristics: use of a within-subjects design minimizes threats to internal validity
16 by equalizing participant characteristics between groups.
- 17 • Task complexity: the experimental tasks for all experimental models were equivalent in
18 complexity.
- 19 • Instrumentation: the same measurement techniques were used for independent and dependent
20 variables for all participants. The risk of measurement error was reduced by calculating metric
21 values automatically.
- 22 • Training: all participants were given the same amount of prior training.
- 23 • Learning effects: experimental models were given to subjects in random order to minimize
24 learning and sequence effects.

25 However, we identified two possible threats to the internal validity of the study.

- 28 • Control of environment: one threat to the validity of the experiment is that it was not conducted
29 under controlled conditions. Participants were allowed to do the experimental tasks in their own
30 time, which means that environmental factors (e.g., time of day, location) were not explicitly
31 controlled. This was done because of the number of separate tasks involved: it would be quite
32 onerous for subjects to complete all nine experimental tasks in a single laboratory session, and
33 may have led to fatigue, which could have adversely affected the results. We argue that because
34 of the within-subjects design, these effects would be randomly distributed across levels of the
35 independent variable, and therefore do not represent a serious threat to the conclusions of the
36 study.
- 37 • Measurement error: another threat to internal validity is the fact that subjects were responsible for
38 recording the time it took to perform the experimental tasks. This increases risk of measurement
39 error for the dependent variable, as subjects may have recorded the time inaccurately. However,
40 again, the within-subjects design means that the error should be randomly distributed across
41 levels of the independent variable, so does not represent a serious threat to the findings.
42 If anything, this would reduce the significance of the regression results by introducing random
43 error, so the results most likely underestimate the strength of the relationship between complexity
44 and maintenance effort.



01 4.5.2. External validity

02
03 In this study, external validity is a much greater issue than internal validity. We identify three possible
04 threats to the external validity of this study.

- 05
06 • Sample population: a clear threat to the generalizability of the findings was the choice of
07 experimental subjects. In general, the population from which one selects subjects for the
08 experiment should be representative of the population to which the researcher wishes to
09 generalize results [81]. In this study, we used a sample population of students and academics,
10 but wish to generalize the results to practice. Most previous experimental studies of conceptual
11 modeling have used Computer Science or Information Systems students as proxies for
12 practitioners, so this study has no greater external validity problems than most (e.g., [64,82–
13 90]). According to the literature, causal relationships are more generalizable across populations
14 than values of variables [75,81]. We can therefore be reasonably confident that the causal
15 relationships found between the complexity metrics and maintenance effort will also hold in
16 the target population, although the level of performance (i.e., time taken to perform the tasks) is
17 likely to be significantly different.
- 18 • Experimental models: all of the class diagrams used in this experiment were created by the
19 researchers specifically for the purpose of the experiment. As a result, they may not be
20 representative of models found in practice in terms of distribution of metric values. In particular,
21 the experimental models were much less complex than models that would be found in practice,
22 where models frequently consist of 100 or more classes. This raises questions about the
23 *scalability* of the results.
- 24 • Experimental task: a final weakness of this experiment was the artificiality of the experimental
25 task. Participants were required to make two rather simple changes to each experimental model—
26 in practice, maintenance changes would be much more complex. It is difficult to realistically
27 simulate the task of maintaining UML models in a laboratory environment and this would be
28 better tested in a real-world maintenance environment.

30 4.5.3. Statistical validity

31
32 There are a number of assumptions underlying the linear regression model that need to be satisfied for
33 the results of regression analysis to be statistically valid. However, in many if not most reports where
34 linear regression analysis is used, these assumptions are not even addressed [75].

- 35
36 • *Linearity*. Inspection of the scatterplots of both independent variables in the final regression
37 model (M2 and M4) against Maintenance Time showed no evidence of curvilinearity.
- 38 • *Random error*. Correlation coefficients were calculated between independent variables and
39 unstandardized residuals in each stage of the regression analysis and found to be 0.000.
- 40 • *Residual normality*. Shapiro–Wilk tests of normality were carried out on unstandardized
41 residuals in each stage of the regression analysis, and they were found to be normally distributed.
- 42 • *Homoscedasticity*. Heteroscedasticity exists when the variance of residuals increases or
43 decreases with values of the independent variable. Inspection of scatterplots showed that the
44 variance of residuals was randomly distributed in all cases.



- 01 27. Eysenck M, Keane M. *Cognitive Psychology: A Student's Handbook*. Lawrence Erlbaum Associates: Mahwah NJ, 2000; 540.
- 02 28. Jacoby J, Speller D, Kohn C. Brand choice as a function of information load. *Journal of Applied Psychology* 1978; **63**:83–91.
- 03 29. Lipowski Z. Sensory and information inputs overload: Behavioural effects. *Comprehensive Psychiatry* 1975; **16**(3):105–124.
- 04 30. Heales J. Factors affecting information system volatility. *Proceedings 21st International Conference on Information Systems (ICIS'2000)*. Ang S, Krcmar H, Orlikowski WJ, Weill P, DeGross JI (eds.). Association for Information Systems: Atlanta GA, 2000; 1–3.
- 05 31. Klir G, Elias D. *Architecture of Systems Problem Solving*. Plenum: New York NY, 2003; 349.
- 06 32. Arthur LJ. *Measuring Programmer Productivity and Software Quality*. Wiley: New York NY, 1985; 292.
- 07 33. Brooks FP. *The Mythical Man-Month: Essays on Software Engineering*. Addison-Wesley: Reading MA, 1975; 195.
- 08 34. Meyer B. *Object Oriented Software Construction*. Prentice-Hall: Nyack NY, 1997; 1296.
- 09 35. Alexander C. *Notes on the Synthesis of Form*. Harvard University Press: Boston MA, 1968; 224.
- 10 36. Oman P, Hagemester J. Metrics for assessing a software system's maintainability. *Proceedings Conference on Software Maintenance*. IEEE Computer Society Press: Los Alamitos CA, 1992; 337–344.
- 11 37. Welker KD, Oman PW. Software maintainability metrics models in practice. *Crosstalk* 1995; **8**(11):19–23, 32.
- 12 38. Welker KD, Oman PW, Atkinson G. Development and application of an automated source code maintainability index. *Journal of Software Maintenance* 1997; **9**(3):127–159.
- 13 39. OMG. *Unified Modeling Language (UML) Specification, Version 1.5*. Object Management Group: Needham MA, 2001; 736.
- 14 40. Genero M. Defining and validating metrics for conceptual models. *Doctoral Dissertation*, University of Castilla-La Mancha, Ciudad Real, Spain, 2002; 437.
- 15 41. Genero M, Piattini M, Calero C. Early measures for UML class diagrams. *L'Objet* 2000; **6**(4):489–515.
- 16 42. Genero M, Olivas J, Piattini M, Romero F. Using metrics to predict OO information systems maintainability. *CAISE 2001 (Lecture Notes in Computer Science, vol. 2068)*. Springer: Berlin, 2001; 388–401.
- 17 43. Marchesi M. OOA metrics for the Unified Modeling Language. *Proceedings 2nd Euromicro Conference on Software Maintenance and Reengineering*. IEEE Computer Society Press: Los Alamitos CA, 1998; 67–73.
- 18 44. Bansiya J, Davis C. A hierarchical model for object-oriented design quality assessment. *IEEE Transactions on Software Engineering* 2002; **28**(1):4–17.
- 19 45. Henderson-Sellers B. *Object-Oriented Metrics—Measures of Complexity*. Prentice-Hall: Englewood Cliffs NJ, 1996; 252.
- 20 46. Chidamber S, Kemerer C. A metrics suite for object oriented design. *IEEE Transactions on Software Engineering* 1994; **20**(6):476–493.
- 21 47. Lorenz M, Kidd J. *Object-Oriented Software Metrics: A Practical Guide*. Prentice-Hall: Englewood Cliffs NJ, 1994; 146.
- 22 48. Harrison R, Counsell S, Nithi R. Coupling metrics for object-oriented design. *Proceedings 5th International Symposium on Software Metrics*. IEEE Computer Society Press: Los Alamitos CA, 1998; 150–157.
- 23 49. Zuse H. *A Framework of Software Measurement*. Walter de Gruyter: Berlin, 1997; 775.
- 24 50. Eitzkorn LH, Bansiya J, Davis C. Design and code complexity metrics for OO classes. *The Journal of Object-Oriented Programming* 1999; **12**(1):335–340.
- 25 51. Briand L, Daly J, Wüst J. A unified framework for coupling measurement in object-oriented systems. *IEEE Transactions on Software Engineering* 1999; **25**(1):91–121.
- 26 52. Brito e Abreu F, Carapuça R. Object-oriented software engineering: Measuring and controlling the development process. *Proceedings 4th International Conference on Software Quality*. American Society for Quality (ASQ): Milwaukee WI, 1994; 212–224.
- 27 53. Siau K, Cao Q. Unified Modeling Language—a complexity analysis. *Journal of Database Management* 2001; **12**(1):26–34.
- 28 54. Siau K, Erickson J, Lee L. Complexity of UML: Theoretical versus practical complexity. *Proceedings 12th Workshop on Information Technology and Systems (WITS'02)*. Association of Information Systems: Barcelona, 2002; 13–18.
- 29 55. Siau K, Tian Y. The complexity of Unified Modeling Language: A GOMS analysis. *Proceedings 14th International Conference on Information Systems (ICIS'01)*. Association for Information Systems (AIS): Atlanta GA, 2001; 443–448.
- 30 56. Fenton N, Neil M. Software metrics: A roadmap. *Future of Software Engineering*, Finkelstein A (ed.). ACM Press: New York NY, 2000; 359–370.
- 31 57. Moody D. Complexity effects on end user understanding of data models: An experimental comparison of large data model representation methods. *Proceedings of the 10th European Conference on Information Systems (ECIS'2002)*. University of Gdańsk: Gdańsk, Poland, 2002; 53–58.
- 32 58. Shanks CG, Moody DL, Nuredini J, Tobin D, Weber RA. Representing things and properties in conceptual modelling: An empirical evaluation. *Proceedings of the 11th European Conference on Information Systems (ECIS 2003)*. Seconda Università degli Studi di Napoli: Naples, 2003; 112–124.

Marcela Genero Bocco

De: "John Wiley & Sons, Ltd." <alistair.smith@sunrise-setting.co.uk>
Para: <marcela.genero@uclm.es>
CC: "SMRedNC2" <NedChapin@acm.org>; "SMR proofs" <smrproofs@wiley.co.uk>; <dm@hi.is>; <mario.piattini@uclm.es>
Enviado: martes, 29 de marzo de 2005 17:04
Adjuntar: smr312.pdf; Offprint order.doc
Asunto: Your proofs of SMR312

Dear Author

Your proofs of SMR312

Please find attached proofs of your article with instructions on how to return your corrections and an editable offprint order form.

Do you want to be alerted when your paper is published online? Simply register for Contents Alerts and automatically receive the table of contents of Journal of Software Maintenance and Evolution: Research and Practice. Just select e-mail alert on the journal home page
<http://www.interscience.wiley.com/jpages/1532-060X>

If you have any queries please contact Sunrise Setting Ltd.

Thank you in advance for your co-operation.

Kind regards,

Tim Williams

SMR Project Manager

Sunrise Setting Ltd
12a Fore Street
St Marychurch
Torquay
Devon
TQ1 4NE
UK

Phone: +44 (0) 1803 322635
Fax: +44 (0) 1803 323565
E-mail: smr@sunrise-setting.co.uk