

University of Groningen

Proposing and empirically validating change impact analysis metrics

Arvanitou, Elvira Maria

IMPORTANT NOTE: You are advised to consult the publisher's version (publisher's PDF) if you wish to cite from it. Please check the document version below.

Document Version

Publisher's PDF, also known as Version of record

Publication date:
2018

[Link to publication in University of Groningen/UMCG research database](#)

Citation for published version (APA):

Arvanitou, E. M. (2018). *Proposing and empirically validating change impact analysis metrics*. [Thesis fully internal (DIV), University of Groningen]. University of Groningen.

Copyright

Other than for strictly personal use, it is not permitted to download or to forward/distribute the text or part of it without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license (like Creative Commons).

The publication may also be distributed here under the terms of Article 25fa of the Dutch Copyright Act, indicated by the "Taverne" license. More information can be found on the University of Groningen website: <https://www.rug.nl/library/open-access/self-archiving-pure/taverne-amendment>.

Take-down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

Downloaded from the University of Groningen/UMCG research database (Pure): <http://www.rug.nl/research/portal>. For technical reasons the number of authors shown on this cover page is limited to 10 maximum.

Chapter 4 – A Metric for Class Instability

Context: *Change impact analysis investigates the negative consequence of system changes, i.e., the propagation of changes to other parts of the system (also known as the ripple effect). Identifying modules of the system that will be affected by the ripple effect is an important activity, before and after the application of any change.*

Goal: *However, in the literature, there is only a limited set of studies that investigate the probability of a random change occurring in one class, to propagate to another. In this paper we discuss and evaluate the Ripple Effect Measure (in short REM), a metric that can be used to assess the aforementioned probability.*

Method: *To evaluate the capacity of REM as an assessor of the probability of a class to change due to the ripple effect, we: (a) mathematically validate it against established metric properties (e.g., non-negativity, monotonicity, etc.), proposed by Briand et al., and (b) empirically investigate its validity as an assessor of class proneness to the ripple effect, based on the 1061-1998 IEEE Standard on Software Measurement (e.g., correlation, predictive power, etc.). To apply the empirical validation process, we conducted a holistic multiple-case study on java open-source classes.*

Results: *The results of REM validation (both mathematical and empirical) suggest that REM is a theoretically sound measure that is the most valid assessor of the probability of a class to change due to the ripple effect, compared to other existing metrics.*

4.1 Motivation

Change impact analysis is related to assessing and identifying the consequences, caused by modifications in one part of a system, on other parts of the same system, known as the *ripple effect* (Horowitz and Williamson, 1986). Studying and quantifying the ripple effect can provide benefits both before and after the application of a change:

- **before the actual application of changes:** for program comprehension and effort estimation (Bohner, 1996; Hassaine et al., 2011); and
- **after changes have been applied:** for test case prioritization and the identification of relationships among software components (Rovegard et al., 2008).

Despite these benefits, to the best of our knowledge, there is only a limited set of metrics for assessing whether a class is prone to the ripple effect (i.e., its probability to change due to a modification in another class of the system). When exploring the proneness of a class to the ripple effect, there are two aspects that need to be assessed: (a) the probability of the source class to undergo changes, i.e., the event that will trigger the ripple effect phenomenon, and (b) the dependencies that will potentially propagate the change to dependent classes. Although the former can only be estimated by analyzing the source code change history, the latter can be estimated by structural analysis, which will be the focus of this study. In order to come up with potential metrics for structurally assessing class proneness to the ripple effect, we need to consider: a) the number of dependencies, and b) a parameter named *propagation factor*, i.e., the probability of a change to propagate from one class to the other through a dependency (Tsantalis et al., 2005). The only type of metrics that might be able to assess these two aspects are coupling metrics, since class coupling is defined *as the degree to which one class is connected to other classes of the system*. However, the ability of coupling metrics to assess¹⁵ if a class would change due to the ripple effect has not been empirically validated.

By inspecting the definitions of the most popular coupling metrics (see Chapter 4.6), we observed that none of them can capture both the number of dependencies of a class and their propagation factor. Thus, in this paper we describe a structural coupling metric, namely Ripple Effect Measure (REM)¹⁶, and validate its capacity as an assessor of the probability of a class to change due to the ripple effect. The calculation of REM is based on the identification of

¹⁵ The term *assess* is used as defined in the *1061 IEEE Standard for Software Quality Metrics* (1998), i.e., as the ability to substitute, track the changes of, and predict the levels of a quality factor; and to discriminate (IEEE-1061, 1998) between low and high values of the quality factor.

¹⁶ REM has been introduced by Ampatzoglou et al. (2015), as part of a case study on the instability of pattern-participating classes. In this paper, we discuss its foundations to enable the self-containment of this manuscript, but our main focus is on its validation.

efferent class dependencies, the quantification of the propagation factor for each one of them, and the aggregation from the dependency level to the class level. The validation process of REM is two-fold:

- a theoretical validation aiming at mathematically proving that REM holds basic properties of software measurement (e.g., non-negativity, normalization, etc.) (Briand et al., 1999);
- an empirical one, aiming at investigating the validity (IEEE-1061, 1998) of REM as an assessor of a class proneness to the ripple effect, by comparing it to existing coupling metrics.

The rest of the paper is organized as follows: In Chapter 4.2 we discuss *related work*, whereas in Chapter 4.3 we present the *REM* and its *calculation* process. Chapter 4.4 discusses the *metric validation processes* that will be used during REM validation. Chapter 4.5 presents the *outcome of the theoretical validation* of REM. Chapter 4.6 presents the *design* and the *results of the empirical validation* of REM. Chapter 4.7 *discusses the main findings* of the validation; finally, Chapters 4.8 and 4.9 present *threats to validity* and *conclude* the paper, respectively.

4.2 Related Work

In this chapter, we present studies that are related to the quantification of the ripple effect. As indirect related work, we present studies that attempt to quantify stability, i.e., the ability of the software to remain unchanged, regardless of the changes, occurring to other parts of the system. The connection between stability and ripple effect has been discussed from the earliest papers on ripple effect by Yau and Collofello (1978; 1980; 1985).

Specifically, in the early '80s Yau and Collofello proposed some measures for *design* and *source code stability*. Both measures were considering the probability of an actual change to occur, the complexity of the changed module, the scope of the used variables, and the relationships between modules (Yau and Collofello, 1980; Yau and Collofello, 1985). The metrics proposed by Yau and Collofello are similar to REM, in the sense that they both take into account the dependencies between modules, and the global variables that are part of the global interface of a module. The main points of differentiation of these metrics, compared to REM, is: (a) the inclusion of the probability of the change to occur in the class that emits the ripple effect – leading to a metric that is not

structural, (b) the inclusion of a complexity metric that is used for quantifying the effort needed to apply the change – which is irrelevant to the probability of a change to propagate, but only related to the effort that will be needed to apply the change.

In a more recent study (2007), Black, proposed an approach for calculating a *complexity-weighted total variable definition propagation for a module*, based on the model proposed by Yau and Collofello. The approach calculates complexity metrics, coupling metrics, and control flow metrics, and their combination provides the proposed ripple effect measure (Black, 2008). The main points of similarity and differentiation are the same, as those of Yau and Collofello. An *empirical study* conducted by Elsih and Rine (2003), investigated the *ability of existing metrics* to assess the design stability of classes. The results suggested that coupling metrics (CBO and RFC) are the optimum assessors of class stability, i.e., the reciprocal of the ripple effect. However, the used stability measure was not the actual one, but an estimation based on dependencies and attribute sharing. Finally, in a research effort on a different direction, by Alshayeb and Li (2005), the authors propose a system *design instability measure* that quantifies the actual changes that have been performed from one version of the system to the other. Nevertheless, this approach is not comparable to REM, in the sense that this study is an after-the-fact analysis, whereas REM is an estimator of the proneness of a class to the change due to the ripple effect.

4.3 Ripple Effect Measure

Dependency analysis is at the core of algorithms that explore the ripple effect caused by class change (see e.g., (Black, 2008; Conejero et al., 2012; Dantas et al., 2012; Yau and Collofello, 1980; Yau and Collofello, 1985), in the sense that changes propagate, across system classes, through their dependencies. Such change propagations (i.e., the most common ripple effect) (Haney, 1972), are the result of certain types of changes in one class (e.g., a change in the method signature—i.e., method name, types of parameters and return type—that is invoked inside another method) that potentially emit changes to other classes¹⁷. Such types of changes vary across different types of dependencies.

¹⁷ We note that this work, does not handle conceptual dependencies among classes, as they would be enforced by class contracts. Our work, aims at providing a purely structural metric for assessing the probability of a class to change due to the ripple effect.

According to van Vliet (van Vliet, 2008), there are three types of class dependencies, namely: generalization, containment, and association. Next, we describe the abovementioned types of dependencies, along with the types of changes that can propagate through them:

- **Generalization** is used to represent “is-a” relationships. In a generalization, there are three possible reasons for change propagation: (a) super method invocation (use of super), (b) access of protected fields, and (c) override or implementation of abstract methods by the subclass.
- **Containment** is used to represent “has-a” or “part-whole” relationships. In a containment relationship, changes can propagate due to method calls of the container class to the public interface of the containee class.
- **Association** is used to represent relationships due to the declaration of a local variable inside a method, or due to the use of a class object as a parameter/return type in a method of another class. In an association relationship, changes can propagate due to method calls of a class to the public interface of another class.

We note that the aforementioned way that changes are propagated through class dependencies are described for designs that follow basic object-oriented design principles, i.e. encapsulation (classes do not hold public attributes). In cases that classes hold public attributes, these public attributes are also considered as a reason for change propagation, in the sense that they belong to the class public interface.

Let us consider the design of Figure 4.3.a, in which two relations exists (A1 depends on A, and A1 depends on B). Are these relationships equally strong? Is it equally probable for a change occurring in A and a change occurring in B to ripple into A1? In order to answer such questions, a metric that quantifies the probability of a ripple effect between classes is necessary. To this end, we defined the Ripple Effect Measure (REM), a metric that quantifies the probability of a random change occurring in the public interface of a *source class* (A or B) to be propagated to a *dependent class* (A1) that uses it, by assuming that all elements of the source classes (i.e., attributes and methods) have the same probability to change¹⁸. To calculate REM, we consider:

¹⁸ As a static measure REM discards changes that occur in other parts of the source code (private parts or method bodies), in the sense that changes in the public interface affect only the class that changes (other

- The count of all members of the source classes (A or B) that are accessed by the dependent class (A1), which (if changed) will emit one or more changes to the dependent class (A1); and
- The count of all members of the public interface of source classes (A or B).

The ratio of the two aforementioned counts is an estimate of the probability that a random change in the public interface of source class will occur in a member that will emit this change to the dependent class. In other words, as the number of the members of the source class that emit changes to another dependent class, approaches the total number of members that can change in the source class, it becomes more probable for changes to propagate from the source class to the dependent class. Thus, and by taking into account the three types of dependencies and the way changes propagate across them, REM for a dependency (i.e., the propagation factor) between the dependent and the source class, can be calculated as follows:

$$\text{REM}_{\text{dependency}} = \frac{\text{NDMC} + \text{NOP} + \text{NPrA}}{\text{NOM} + \text{NA}} \quad (1)$$

NDMC: Number of distinct method calls from the dependent class to the source class (plus, super class method invocations for the case of generalization)

NOP: Number of polymorphic methods in the source class (valid only for generalization)

NPrA: Number of protected attributes in the source class (valid only for generalization or friend classes)

NOM: Number of methods in the source class

NA: Number of attributes in the source class

Based on the example of Figure 4.3.a, there are three changes which would ripple from superclass A to subclass A1:

- **change in the signature of a1.** Changing the signature of a1(), would lead to a compile error in the corresponding invocation (see body of method a3).
- **change in the signature of a2.** Changing the signature of a2() would lead to a compile error in the place where a2() is overridden, since a2() is

classes do not have access to it) and co-change of method bodies can only be tracked by historical analysis (for more details on the differences between class change proneness and instability see (Ampatzoglou et al., 2015)).

declared as abstract in the superclass.

- **change in the name of att2.** Changing the name of attribute att2, would lead to a compile error in the body of a3.

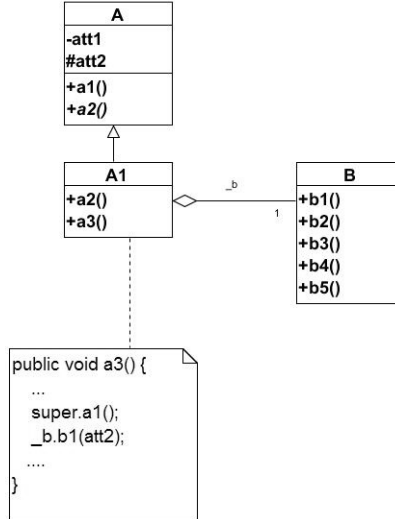


Figure 4.3.a: Coupling Intensity Example

Thus, REM, for the dependency of A₁ on A, can be calculated as follows:

$$REM_{A_1(A)}^{19} = \frac{NDMC+NOP+NPrA}{NOM+NA} = \frac{1+1+1}{2+2} = 0.75$$

Similarly, there is only one change which would propagate from containee class (B) to the container class (A1), i.e. change in the signature of b1. Changing the signature of b1, would lead to a compile error in the corresponding invocation. We note that if any other method of B is changed, the change will not propagate to A1, as long as the method is not called in the implementation of A1 (in the example no other method of B is invoked from A1, except b1). Thus, REM, for the dependency of A₁ on B, can be calculated as:

$$REM_{A_1(B)} = \frac{NDMC+NOP+NPrA}{NOM+NA} = \frac{1+0+0}{5+0} = 0.20$$

Therefore, based on REM, for the example of Figure 4.3.a, we can claim that a change occurring in A is more probable to propagate in A1, than a change occurring in B.

¹⁹ This refers to the REM of class A1, due to its dependency to A

Until this point, the calculation of REM is performed at the class dependency level, i.e., from a class to one other class. In order to aggregate REM to a higher level of granularity, i.e., at class level, one should take into account the dependencies of a class on all other classes. In the example of Figure 4.3.a, class A1 will have to change if a change propagates from either class B or A or both. Therefore, for one class having several dependencies, we propose the use of the *joint probability of all events* (i.e., change in any dependency) (Aho and Ullman, 1995), for aggregating the score of REM from dependency to class level. For example, in the case of Figure 4.3.a, the joint probability is calculated as follows:

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A) \cdot P(B) \\ &= 0.75 + 0.20 - 0.75 \cdot 0.20 = 0.8 \end{aligned}$$

P(A ∪ B): REM_{A1(A, B)}

P(A): REM_{A1(A)}

P(B): REM_{A1(B)}

4.4 Validation Process

According to Briand et al. (1999), metric validation should be performed as a two-step process, i.e., providing a theoretical and empirical validation. The *theoretical validation* aims at mathematically proving that a metric satisfies certain criteria, whereas the *empirical validation* aims at investigating the accuracy with which a specific metric quantifies the corresponding quality factor. The validation of the REM metric will be performed using both types of validation. Concerning the theoretical validation of REM, we will use the properties for coupling metrics proposed by Briand et al. (1999)²⁰. The used properties are:

- **Normalization and Non-Negativity:** The metric score should have a pre-defined lower and upper bound, and not be negative (Briand et al., 1999);
- **Null Value and Maximum Value:** The metric score should be null or take the maximum value, only under very specific circumstances (Briand et al., 1999);

²⁰ We acknowledge the existence of various other metric properties (e.g., the ones proposed by Weyuker (1988)), we preferred to use the ones proposed by Briand et al., because they are coupling-specific.

- **Monotonicity:** The metric score for a system S should be higher than the metric score for a system S' (which is identical to S , except from one class, i.e., C and C' respectively), if the metric for class C is higher than the metric score for class C' (Briand et al., 1999); and
- **Merging of Classes:** If a class C' is the union of two classes $C1$ and $C2$, then the metric score for C' should be less or equal to the sum of the scores of $C1$ and $C2$ (Briand et al., 1999). For the special case of unconnected classes $C1$ and $C2$, the property suggests that the merged class should have a metric score that is equal to the sum of the individual scores.

Regarding the empirical validation of REM we will use the properties described in the *1061 IEEE Standard for Software Quality Metrics* (2009), for comparing REM to existing coupling metrics (for details on these coupling metrics see Chapter 4.6). In this Standard, six metric validation criteria are introduced, accompanied by the statistical test that shall be used for evaluating every criterion, as follows:

- **Correlation** assesses the association between a quality factor and the metric under study. The criterion is desirable, to ensure that the use of the metric can substitute the quality factor. The criterion is quantified by using a correlation coefficient (IEEE-1061, 1998).
- **Consistency** assesses if the metric under study is consistently correlated with the quality factor, by using their ranks. The criterion is desirable to ensure that the metric under study can accurately rank, by quality, a set of products or processes. The criterion is quantified by using the coefficient of rank correlation (IEEE-1061, 1998).
- **Tracking** assesses if the metric under study is capable of tracking changes in product quality over their lifecycle. The criterion is quantified by using coefficient of rank correlation for a set of project versions (IEEE-1061, 1998).
- **Predictability** assesses the accuracy with which the metric under study is able to predict the levels of the quality factor. The criterion is quantified through the standard estimation error for a regression model using as predictor the metric under study (IEEE-1061, 1998).
- **Discriminative Power** assesses if the metric under study is capable of

separating groups of high-quality and low-quality components. Although the criterion is proposed to be quantified through a contingency table (IEEE-1061, 1998), the proposed quantification was not applicable for coupling metrics, because they cannot be recoded to categorical variables, without setting arbitrary thresholds. Therefore, we use an equivalent test for assessing discriminative power, i.e., independent t-test (Field, 2013).

- **Reliability** assesses if the metric under study can fulfill all five aforementioned validation criteria, over multiple systems. This criterion can be assessed by replicating the previously discussed tests (for each of the aforementioned criteria) to various software systems (IEEE-1061, 1998).

4.5 Theoretical Validation

Based on the definition of REM (see Chapter 4.3), we validate it against all coupling properties in the next paragraphs.

4.5.1 Normalization and Non-Negativity

REM at both dependency and class levels is normalized and non-negative. Concerning normalization, the **lower bound of REM is zero** (since the lower bound of the numerator is zero – i.e., classes that *do not call any method* from other classes, and *do not override polymorphic methods*, and *do not access any protected attributes* of a superclass); whereas the **upper bound of REM is 1** (since $NMDC + NOP \leq NOM$ and $NPrA \leq NA$). Regarding non-negativity, **REM is always assigned to a positive value**, because it involves only additions and divisions of positives numbers.

4.5.2 Null Value and Maximum Value

REM at both dependency and class levels is defined for all possible cases of dependencies and classes, and has a specific maximum value. REM would not be defined (assigned a null value) only in the case that the denominator would be zero, i.e., classes depending on other **classes with no attributes and no methods**. However, such classes are not expected to exist, in the sense that they would not offer any functionality to the system. As mentioned in Chapter 4.5.1, the **maximum value of REM is 1**.

4.5.3 Monotonicity

In the original definition of monotonicity, two levels of granularity are involved, i.e., the system level and the class level, where the values from class

level are aggregated to system level. In the case of REM, the two levels of granularity are different, i.e., the REM at dependency level, is aggregated to class level as discussed in Chapter 4.3. Therefore, in order to explore the monotonicity of REM, we apply the definition, by mapping class to system (original definition) and dependency to class (original definition).

For simplicity, assume a class C_1 with two dependencies (D_{11} and D_{12}), and a class C_2 with one dependency having exactly the same REM as the first dependency of C_1 and a second dependency that is different (D_{21} and D_{22} , where $REM_{D_{11}}=REM_{D_{21}}$ and $REM_{D_{12}} > REM_{D_{22}}$). Based on the aforementioned assumptions and the definition of monotonicity, in order for REM to be monotonic, we need to prove that $REM_{C_1(D_{11},D_{12})} > REM_{C_2(D_{21},D_{22})}$. We mathematically prove the previous relationship, as follows:

$$REM_{C_1(D_{11},D_{12})} > REM_{C_2(D_{21},D_{22})} \quad \Rightarrow$$

$$\cancel{REM}_{D_{11}} + REM_{D_{12}} - \cancel{REM}_{D_{11}} * \cancel{REM}_{D_{12}} >$$

$$\cancel{REM}_{D_{21}} + REM_{D_{22}} - \cancel{REM}_{D_{21}} * \cancel{REM}_{D_{22}} \quad \Rightarrow$$

$$REM_{D_{12}} * (1 - \cancel{REM}_{D_{11}}) > REM_{D_{22}} * (1 - \cancel{REM}_{D_{21}}) \quad \Rightarrow$$

$$REM_{D_{12}} > REM_{D_{22}} ,$$

which holds, based on the original assumption ($REM_{D_{12}} > REM_{D_{22}}$). Therefore, REM is a **monotonic coupling** metric. A proof of monotonicity for classes with n dependencies, could be provided, but it is omitted due to space limitations.

4.5.4 Merging of Classes

For similar reasons to monotonicity, we will argue that REM holds for the merging of classes property, through an example with two dependencies. Let class C_1 be a class with one dependency (D_1) and class C_2 be another class with a different dependency (D_2). The merged class C' will have two dependencies (D_1 and D_2). In order for the ‘merging of classes’ property to hold for REM, we need to prove that $REM_{C'(D_1,D_2)} \leq REM_{C_1(D_1)} + REM_{C_2(D_2)}$. The mathematical proof of the previous relationship is as follows²¹:

²¹ We note that since class C_1 has only one dependency (i.e., D_1), $REM_{C_1} = REM_{D_1}$

$$REM_{C'(D1,D2)} \leq REM_{C1(D1)} + REM_{C2(D2)} \quad \Rightarrow$$

$$\begin{aligned} REM_{D1} + REM_{D2} - REM_{D1} * REM_{D2} \\ \leq REM_{D1} + REM_{D2} \quad \Rightarrow \end{aligned}$$

$$-REM_{D1} * REM_{D2} \leq 0,$$

which holds, because both REM_{D1} and REM_{D2} are non-negative numbers (based on the explanations provided in Chapter V.A). Thus, REM is able to handle **merging of classes** as implied by the corresponding property. We note that the property of the special case of *merging unconnected classes* is not satisfied by the REM, due to its probabilistic nature.

4.6 Empirical Validation

In order to empirically investigate the validity of REM to assess if a class will change due to the ripple effect, we performed a case study on two open source projects. We compare REM to existing coupling metrics with respect to all the criteria described in Chapter 4.4 (IEEE-1061, 1998). The coupling metrics that have been used as control variables are:

- **Coupling Between Objects (CBO)**: Number of classes to which a class is coupled (Chidamber and Kemerer, 1994).
- **Response For a Class (RFC)**: Number of local methods, plus the number of methods called by local methods in the class (Chidamber and Kemerer, 1994).
- **Message Passing Coupling (MPC)**: Number of send statements defined in the class (Li and Henry, 1993).
- **Data Abstraction Coupling (DAC)**: Number of abstract data types defined in the class (Li and Henry, 1993).
- **Measure of Aggregation (MOA)**: Number of data declarations of a user defined type (Bansiya and Davies, 2002).

In order to be as inclusive as possible, we selected metrics from three different metric suites (Chidamber and Kemerer (1994), Li and Henry (1993), and QMOOD (Bansiya and Davies, 2002), which are well-known and tool-supported. Also, the aforementioned list of metrics includes both code- and design-level coupling metrics.

4.6.1 Case study Design

The case study has been designed and reported according to the guidelines of Runeson et al. (2012). Therefore, in this Chapter we present: (a) the goal of the case study and the derived research questions, (b) the description of cases and units of analysis, (c) the case selection criteria, (d) the data collection procedure, and the (e) data analysis process.

Research Objectives and Research Questions: The aim of this case study, expressed through a GQM formulation, is: *to analyze coupling metrics for the purpose of evaluation with respect to their validity to assess if a class will change due to the ripple effect, from the point of view of software engineers in the context of change impact analysis.*

Based on this goal we set two research questions: the first, on the first five validity criteria (i.e., correlation, consistency, tracking, predictability and discriminative power), in which no distinction between projects is necessary (i.e., all projects are handled as one dataset); and the second on reliability. Reliability is investigated through a separate research question, because it entails performing the analysis for the aforementioned five validity criteria separately for each project (i.e., each project is considered as a different dataset, and those datasets are cross-checked in order to assess metrics' reliability). Therefore, we ask the following research questions:

- RQ1:** Can REM assess if a class will change due to the ripple effect, based on the criteria of the 1061-1998 IEEE Std., compared to other existing metrics?
- RQ1.1:** How does REM compare to the other coupling metrics, w.r.t. correlation?
- RQ1.2:** How does REM compare to the other coupling metrics, w.r.t. consistency?
- RQ1.3:** How does REM compare to the other coupling metrics, w.r.t. tracking?
- RQ1.4:** How does REM compare to the other coupling metrics, w.r.t. their predictive power?
- RQ1.5:** How does REM compare to the other coupling metrics, w.r.t. their discriminative power?

RQ₂: How does REM compare to the other coupling metrics, w.r.t. their reliability?

Case and Units of Analysis: This study is a holistic multiple-case study, i.e. it studies multiple cases where each case is comprised of a single unit of analysis. Specifically the subjects of the study are open source projects, where classes are cases and at the same time units of analysis.

Case Selection: As subjects for our study, we selected to use two open source projects, i.e., JFlex (versions 1.4 and earlier - in total 14 versions) and JMoI (versions 0.9 and earlier - in total 10 versions). The main motivation for selecting these projects was the intention to reuse an existing dataset, which has been developed and used for a research effort with similar goals (Tsantalis et al., 2005). Specifically, Tsantalis et al. (2005) have explored the specific versions of these projects, to identify classes that have changed from the projects' previous versions. These changes have been manually inspected, so as to check whether they are due to reasons other than the ripple effect (main effect—referred as internal changes in the original publication (Tsantalis et al., 2005)) or due to the ripple effect (referred as propagated changes in the original paper (Tsantalis et al., 2005)). Consequently, we were able to reuse the extracted data that concerned changes due to ripple effect. All classes of these systems have been used as cases for this study. Therefore, our study was performed on 150 java classes. We clarify that for assessing correlation, consistency, predictability and discriminative power, we used only the last version of JFlex and JMoI as cases: according to the used IEEE standard, system evolution is not considered for these four validation criteria. We used the last versions of the projects, so as to examine the largest (in terms of lines of code and number of classes) versions of the systems. On the other hand, regarding tracking and reliability we used all examined versions of both systems.

Data Collection: For each case (i.e., class), we primarily recorded ten variables, as follows:

- **Demographics:** 3 variables (i.e., project, version, class name).
- **Evaluated metrics:** 6 variables (i.e., REM, CBO, RFC, MPC, DAC, and MOA). These variables are going to be used as the independent variables for testing correlation, predictability and discriminative power.

- **Observed ripple effect:** For any transition between two successive versions of a class, this variable indicates if a class has changed due to the ripple effect (value equals 1), or whether it remained unaffected by the ripple effect (value equals 0). This variable is going to be used as the dependent variable in all tests.

The coupling metrics have been calculated by using two tools:

- **REM** has been calculated by modifying the tool of Tsantalis et al. (2005). The tool in its original version was created in order to calculate class change proneness, by using a constant value for the propagation factor. In the updated version (Ampatzoglou et al., 2015), which is freely available for download in the web²², REM has substituted the propagation factor, so as to increase the realism of the calculated change probability. From this tool, we use the REM calculation, and not the complete instability calculation.
- **the rest of the coupling metrics**, have been calculated using the *Percerons Client* tool²³. Percerons is an online platform (Ampatzoglou et al., 2013a) created to facilitate empirical research in software engineering, by providing, among others, source code quality assessment (Ampatzoglou et al., 2013b). The platform has been used for similar reasons in (Alhusain et al., 2013; Ampatzoglou et al., 2013b; Griffith and C. Izurieta, 2014).

Data Analysis: For investigating both research questions, we used the statistical analysis that is advocated by the 1061-1998 IEEE Standard (1998) (see Chapter 4.4), as summarized in Table 4.6.1.a. We note that a coupling metric is expected to have a proportional relationship to the probability of a class to change due to the ripple effect: the larger the number of dependencies of a class and the stronger the dependencies, the more probable for this class to receive changes, from other classes.

²² <http://iwi.eldoc.ub.rug.nl/root/2014/ClassInstability/>

²³ <http://www.percerons.com>

Table 4.6.1.a: Measure Validation Results

| Criterion | Test | Variables |
|----------------------|---|---|
| Correlation | Point bi-serial correlation | Coupling Metrics Observed ripple effect (<i>last version of the projects</i>) |
| Consistency | Rank bi-serial correlation | Coupling Metrics Observed ripple effect (<i>last version of the projects</i>) |
| Tracking | Rank bi-serial correlation | Coupling Metrics Observed ripple effect (<i>across all versions</i>) |
| Predictability | Logistic Regression | Independent: Coupling Metrics Dependent: Observed ripple effect (<i>last version of the projects</i>) |
| Discriminative Power | Independent Sample t-test | Testing: Coupling Metrics Grouping: Observed ripple effect (<i>last version of the projects</i>) |
| Reliability | all the aforementioned tests (<i>seperately for each project –across all versions</i>) | |

For presenting the results on *Correlation* and *Consistency*, we use the correlation coefficients (coeff.) and the levels of statistical significance (sig.). The value of the correlation coefficient denotes the degree to which the value of the observed ripple effect is in analogy to the value of the predictor. We note that since the dependent variable (i.e., Observed ripple effect) is binary, we used the ***point bi-serial correlation*** and ***rank bi-serial correlation*** for assessing correlation and consistency, respectively. To represent the *Tracking* property of the evaluated metrics, we report on the consistency for multiple project versions, through the mean correlation coefficient and the percentage of versions, in which the correlation was statistically significant.

For reporting on *Predictability*, with a regression model, we present the level of statistical significance of the effect (sig.) of the independent variable on the dependent (how important is the predictor in the model), and the accuracy of the model (i.e., the total correctness value—percentage of correctly classified cases). While investigating predictability, we produced a separate logistic regression model for each predictor (univariate analysis), because our intention was not to investigate the cumulative predictive power of all coupling metrics, but of each metric individually. Similarly to correlation/consistency, we

performed a logistic regression, since the dependent variable is binary, therefore the classification threshold was set to 0.5.

For presenting the *Discriminative Power*, we investigate whether the two groups (actually propagated changes, and not propagated changes) differ with respect to the corresponding coupling metric score. For reporting on the independent sample t-tests, we present the mean difference (diff.) between the values of the testing variable of the compared groups, and the level of statistical significance (sig.). We note that in order for a metric to adequately discriminate groups of cases, the significance value should be less than 0.05.

Finally, for reporting on the *Reliability* of coupling metrics while assessing if a class will change due to the ripple effect, we present the results of all the aforementioned tests, separately for the two explored OSS projects. The extent to which the results on the two projects are in agreement (i.e., Are they statistically significant for both? Is the same metric the most valid assessor of class proneness to the ripple effect for both projects?) represents the reliability of the considered metric.

4.6.2 Results

In this Chapter, we present the results of the empirical validation of REM, organized by research question. We first present the outcome of comparing REM to the other coupling metrics, w.r.t. their correlation, consistency, tracking, predictive and discriminative power as assessors of a class proneness to the ripple effect. Subsequently we present the results of comparing the reliability of REM to the reliability of the other coupling metrics.

RQ1: Correlation, Consistency, Tracking, Predictability and Discriminative Power

Table 4.6.2.a presents the outcome of the statistical analysis related to RQ₁. Each row of Table 4.6.2.a corresponds to one validity criterion, whereas each column to one metric. To enable the easy reading of Table 4.6.2.a, we use visual aids to help the reader in focusing on the most important observations. Specifically, we denote statistically significant relationships at the 0.01 level with grey cell shading, and statistically significant relationships at the 0.05 level with light grey cell shading (when applicable). Finally, the optimal

assessor of class proneness to the ripple effect, for each criterion, is denoted with bold fonts.

To answer RQ₁ we can suggest that *REM* and *CBO* are the only coupling metrics that can provide a statistically significant assessment (albeit weak) whether a class will change due to the ripple effect, w.r.t. the validity criteria described in IEEE 1061 Standard. The results can be summarized as follows:

- **Correlation, Predictability and Discriminative Power:** REM is the most valid metric while assessing if a class will change due to the ripple effect, followed by CBO. The results on REM validity are statistically significant at the 1% level, whereas for CBO at the 5% level. The result for other metrics are not statistically significant. Concerning the discriminative power of metrics, we note that, although the absolute value of the difference offered by CBO is higher than the one offered by REM, this difference is due to the range of values of the two metrics²⁴. Therefore, we consider REM as the most valid assessor of class proneness to the ripple effect, based on its statistical significance.
- **Consistency:** REM is the most consistent assessor of class proneness to the ripple effect, followed by CBO and RFC. The results on the REM are statistically significant at the 1% level. CBO and RFC at the 5% level.
- **Tracking:** REM achieves the best results with respect to tracking class proneness to the ripple effect, followed by CBO. The results on REM are statistically significant for 26% of the examined cases.

Table 4.6.2.a: Measure Validation Results

| | | REM | RFC | MPC | DAC | CBO | MOA |
|-------------|--------------|-------------|------------|------------|------------|------------|------------|
| Correlation | <i>coef.</i> | 0.25 | 0.11 | 0.06 | 0.07 | 0.18 | 0.13 |
| | <i>sig.</i> | 0.00 | 0.18 | 0.46 | 0.34 | 0.02 | 0.09 |
| Consistency | <i>coef.</i> | 0.26 | 0.15 | 0.11 | 0.11 | 0.18 | 0.01 |
| | <i>sig.</i> | 0.00 | 0.05 | 0.16 | 0.18 | 0.02 | 0.97 |
| Tracking | <i>coef.</i> | 0.21 | 0.16 | 0.18 | 0.14 | 0.19 | 0.17 |
| | <i>pct.</i> | 26% | 17% | 21% | 13% | 21% | 17% |

²⁴ Range of REM is [0, 1], whereas range of CBO is [0, +∞).

| | | REM | RFC | MPC | DAC | CBO | MOA |
|----------------------|-------------|--------------|-------|-------|-------|-------|-------|
| Predictability | <i>acc.</i> | 62.3% | 54.3% | 54.3% | 55.0% | 57.6% | 55.0% |
| | <i>sig.</i> | 0.01 | 0.99 | 0.92 | 0.99 | 0.03 | 0.67 |
| Discriminative Power | <i>diff</i> | -0.16 | -5.21 | -2.57 | -0.02 | -1.64 | -0.08 |
| | <i>sig</i> | 0.01 | 0.49 | 0.93 | 0.27 | 0.02 | 0.68 |

RQ2: Reliability

Regarding RQ₂, we performed all the aforementioned tests separately for each one of the analyzed OSS projects, and compared them. The results for each project are summarized in Tables 4.6.2.b and 4.6.2.c. The notation used in Tables 4.6.2.b and 4.6.2.c is the same as that of Table 4.6.2.a. Therefore, the shading of the cell corresponds to the level of statistical significance, whereas the bold fonts indicate the most valid assessor of class proneness to the ripple effect concerning a certain criterion.

From Tables 4.6.2.b and 4.6.2.c, we can observe that REM is the only coupling metric that produces reliable results for all required tests (i.e., similar between the two software systems). Specifically, REM is the optimal assessor of class proneness to the ripple effect regarding *correlation*, *tracking*, *predictive* and *discriminative power*, for both projects. Concerning *consistency*, REM is the most valid assessor of class proneness to the ripple effect for JMol, but the 4th for JFlex. However, for both projects, the results on the consistency of REM are statistically significant, and therefore reliable. An interesting result that we can observe from Tables 4.6.2.b and 4.6.2.c is that CBO, i.e., the second most valid assessor of class proneness to the ripple effect according to Table 4.6.2.a, suffers from reliability issues. Specifically, CBO is not reliable w.r.t any criterion, since it provides statistically significant results, at most, only for one out of the two projects (i.e., consistency and tracking for JFlex and correlation for JMol).

Table 4.6.2.b: Measure Validation Results (JFlex)

| | | REM | RFC | MPC | DAC ²⁵ | CBO | MOA |
|----------|--------------|--------------|-------------|-------|-------------------|-------------|-------|
| Correlat | <i>coef.</i> | 0.37 | 0.24 | 0.00 | N/A | 0.27 | 0.18 |
| | <i>sig.</i> | 0.02 | 0.15 | 1.00 | N/A | 0.10 | 0.28 |
| Consist | <i>coef.</i> | 0.37 | 0.41 | 0.39 | N/A | 0.45 | 0.18 |
| | <i>sig.</i> | 0.02 | 0.01 | 0.02 | N/A | 0.00 | 0.30 |
| Trackin | <i>coef.</i> | 0.24 | 0.18 | 0.18 | N/A | 0.24 | 0.17 |
| | <i>pct.</i> | 28% | 14% | 21% | N/A | 28% | 7% |
| Predicta | <i>acc.</i> | 67.6% | 58.9% | 54.9% | N/A | 62.2% | 64.9% |
| | <i>sig.</i> | 0.03 | 0.16 | 0.99 | N/A | 0.08 | 0.29 |
| Descri | <i>diff</i> | -0.26 | -20.57 | -0.02 | N/A | -2.74 | -0.46 |
| | <i>sig.</i> | 0.02 | 0.15 | 1.00 | N/A | 0.10 | 0.28 |

Table 4.6.2.c: Measure Validation Results (JMol)

| | | REM | RFC | MPC | DAC | CBO | MOA |
|------------|--------------|--------------|-------|-------|-------|-------|-------|
| Correlati | <i>coef.</i> | 0.24 | 0.09 | 0.08 | 0.07 | 0.20 | 0.15 |
| | <i>sig.</i> | 0.00 | 0.33 | 0.37 | 0.45 | 0.03 | 0.10 |
| Consiste | <i>coef.</i> | 0.24 | 0.09 | 0.05 | 0.11 | 0.15 | 0.00 |
| | <i>sig.</i> | 0.01 | 0.36 | 0.56 | 0.25 | 0.11 | 0.98 |
| Tracking | <i>coef.</i> | 0.19 | 0.14 | 0.15 | 0.14 | 0.17 | 0.16 |
| | <i>pct.</i> | 22% | 11% | 11% | 11% | 22% | 22% |
| Predictiab | <i>acc.</i> | 60.5% | 45.0% | 45.0% | 51.8% | 52.9% | 50.9% |
| | <i>sig.</i> | 0.03 | 0.97 | 0.87 | 0.98 | 0.08 | 0.87 |
| Discrimina | <i>diff</i> | -0.13 | -0.38 | -5.12 | -0.02 | -1.40 | 0.03 |
| | <i>sig.</i> | 0.04 | 0.97 | 0.87 | 0.31 | 0.07 | 0.87 |

²⁵ We were unable to calculate the validity of DAC as an assessor of class proneness to the ripple effect in JFlex, due to the limited use of inheritance in the project: no class had any data abstraction coupling.

4.7 Discussion

In this Chapter, the outcomes of this study are discussed from two different perspectives. First, we provide possible interpretations of the obtained results; and second, we present possible implications to researchers and practitioners.

Interpretation of the Results: The results of this study suggest that REM is a theoretically sound coupling metric and that exceeds all other metrics on the considered validation criteria, followed by CBO. REM outperforms other coupling metrics, mostly because it is specifically designed for assessing the probability of a class to change due to the ripple effect, and not coupling in general. To this end, REM has been designed so as to resolve the limitations of the other metrics and combine their strengths, as follows:

- REM takes into account **the number of dependencies**, similarly to CBO, DAC, and MOA. However, none of the aforementioned metrics considers the strength of these dependencies.
- REM is a **normalized measure** that is able to quantify the **strength of dependencies**. For similar metrics, which quantify the strength of coupling and at the same time are proportional to the number of dependencies, i.e. RFC and MPC, it is not clear to what extent their value is related to the two components they mix. For example, MPC is equal for two classes, for which the first has two dependencies and every dependency is used three times, and the second has six dependencies and each dependency is used once.
- REM is the only metric that takes into account the changes that can be propagated by using **protected attributes**, either through hierarchies or friend methods.

These characteristics of REM make it the most valid assessor of the probability of a class to change due to the ripple effect, as demonstrated by the obtained empirical evidence. Also, by comparing the rest of the coupling metrics, we can observe that CBO is the 2nd best assessor of class proneness to the ripple effect. This implies that the number of efferent couplings is indeed an important parameter, when exploring the ripple effect. The fact that CBO is a better assessor than MOA and DAC can be explained, as they capture only a small portion of efferent couplings (i.e., number of aggregations and number of used abstract types respectively). Finally, by contrasting RFC to MPC, we can

conclude that the distinct counting of called methods offered by RFC make it a better assessor than MPC, because regardless of how many times a method is called, a change will be propagated even with only one method call.

Additionally, although REM is the optimal assessor of a class proneness to the ripple effect, the results suggest that the strength of the correlation between them (see Table 4.6.2.a) can only be considered as weak or moderate at best. This is an expected outcome, since in this study REM is validated against the actual class changes due to the ripple effect. However, as explained in Chapter 4.1, the actual probability of change in any class, depends not only on the probability of changes to propagate, but also on the probability of changes to occur in the first place. Nevertheless, the quantification of the actual probability of change in a class, might demand the processing of data that are not always available (e.g., class change history, or subjective identification of design hot-spots, etc.). In such cases, the interested software engineer, would be forced to use a structural metric, among which REM is the optimal predictor of class proneness to the ripple effect, based on the results of this case study.

Implications to Researchers and Practitioners: By taking into account the aforementioned results and discussions, we can provide various implications for researchers and practitioners. Specifically, we:

- Encourage **practitioners** to use REM in their *change impact analysis* activities, in the sense that REM is the optimal available assessor of the probability of a class to change due to the ripple effect. We expect that tool support that automates the calculation process (see Chapter 4.6.1) will ease its adoption. However, the applicability of REM in practice is still in need of further investigation.
- Encourage **researchers** to investigate options for improving the effect size of the metric. For example, by incorporating *the actual change proneness of classes in the calculation of REM*, so as to increase the validity of the proposed metric. Specifically, towards this direction, we plan to replicate this study, by using the tool presented in (Ampatzoglou et al., 2015), feed it with historical data on class change proneness and use all the considered metrics as change propagation factors, and observe: (a) if the results in the fitness of metrics as assessors of class proneness to the ripple effect remain the same, and (b) if the validity of the calculated metrics increases (e.g., strength of correlation).

- Encourage **researchers** to investigate the possibility of *assigning weights to axes through which a class can receive changes* (i.e., protected attributes, called methods, and overridden methods) by a study on class change history, and observe if this can increase the validity of REM.
- Encourage **researchers** to *transform REM* so as to fit *architecture evaluation* purposes, i.e., assess the probability of changes to propagate across components. We believe that such a transformation would be of great interest for the architecture community, which shows increased interest in metrics. Also, such an attempt would increase the benefits of practitioners, in the sense that change impact analysis could scale into larger systems.
- Encourage **researchers** to use *combinations of coupling metrics* to explore the possibility of increasing power of metric for predicting the ripple effect.

4.8 Threats to Validity

In this chapter we present the threats to the validity of our case study²⁶. In this case study, we aimed at exploring if certain coupling metrics are valid assessors of class proneness to the ripple effect. Therefore, possible threats to construct validity (Runeson et al., 2012) deal with the way coupling and the ripple effect are quantified, including both rationale and tool support. On the one hand, the rationale on how the coupling metrics are calculated is not a threat, since their definition is clear and well documented, whereas the used tools have been thoroughly tested, before deployment. On the other hand, in order to assess the probability of a class to change due to the ripple effect, we used the number of actually propagated changes, which is exactly how the ripple effect should be calculated. The calculation of actually propagated changes, has been reused from an already published study (Tsantalis et al., 2005), and their results have been manually validated by the first author.

Additionally, in order to ensure the reliability (Runeson et al., 2012) of this study, we: (a) thoroughly documented the study design in this document (see Chapter 4.6.1), so as to make the study replicable, and (b) all steps of data collection and data analysis have been performed by two researchers in order

²⁶ The mathematical results presented in Chapter 4.5, by definition, cannot suffer from threats to validity.

to ensure the validity of the process. Furthermore, concerning the external validity (Runeson et al., 2012) of our results, we have identified two possible threats:

- *we investigated only two OSS projects.* The low number of subjects is a threat to generalization, in the sense that results on these two projects cannot be generalized to the complete OSS projects population. However, since the units of analysis for this study are classes and not projects, we believe that this threat is mitigated. The only real threat concerns the reliability criterion, as it compares results from different projects and we only compare two OSS projects against each other. For this specific criterion, further investigation is required.
- *we investigated project only written in Java.* Due to tool limitations, we could only analyze projects written in Java. Therefore, the results cannot be generalized in other languages, e.g., C++. This threat becomes, even more important, because C++ projects are expected to also make use of the friend operator, which changes the scope of class attributes.

4.9 Conclusions

In this study, we presented and validated a new coupling metric, named Ripple Effect Measure (REM), which can be used for assessing a class proneness to the ripple effect. The metric is defined at dependency level by calculating the portion of the accessible interface of a class that is used by other classes. After calculating REM for all dependencies, REM can be aggregated to class level, by employing simple probability theory. The validation of REM has been performed in a two-step process: first, we theoretically validated it against well-known coupling metric properties (Briand et al., 1999); and second, we empirically validated against the metric validation criteria defined in the 1061-1998 IEEE Standard for a Software Quality Metrics (1998). To empirically investigate the validity of REM as an assessor of class proneness to the ripple effect, we performed a holistic multiple-case study on over 150 OSS project classes.

The results of the theoretical validation suggested that REM is mathematically sound, in the sense that all coupling metric properties can be proven, based on its definition. Additionally, the results of our case study suggested that REM excels as an assessor of a class proneness to the ripple effect metric compared

to a variety of well-known coupling metrics. Based on these results, implications for researchers and practitioners have been provided.