

© 2012 Shin Hwei Tan

@TCOMMENT: TESTING JAVADOC COMMENTS TO DETECT COMMENT-CODE  
INCONSISTENCIES

BY

SHIN HWEI TAN

THESIS

Submitted in partial fulfillment of the requirements  
for the degree of Master of Science in Computer Science  
in the Graduate College of the  
University of Illinois at Urbana-Champaign, 2012

Urbana, Illinois

Adviser:

Associate Professor Darko Marinov  
Assistant Professor Lin Tan, University of Waterloo

# Abstract

Code comments are important artifacts in software. Javadoc comments are widely used for API specifications in Java. API developers write Javadoc comments, and API users often read these comments to understand the API they use, e.g., an API user can read a Javadoc comment for a method instead of reading the method body of the method. An inconsistency between the Javadoc comment and body for a method indicates either a fault in the body or, effectively, a fault in the comment that can mislead the method callers to introduce faults in their code.

This thesis presents a novel approach, called @TCOMMENT, for testing Javadoc comments, specifically for testing method properties about null values and related exceptions. Our approach consists of two components. The first component takes as input source files for a Java project and automatically analyzes the English text in Javadoc comments to infer a set of likely properties for a method in the files. The second component generates random tests for these methods, checks the inferred properties, and reports inconsistencies. We evaluated @TCOMMENT on seven open-source projects and found 28 inconsistencies between Javadoc comments and method bodies. We reported all inconsistencies, and 12 have already been confirmed and fixed by the developers.

*To my family and friends, for their love and support.*

# Acknowledgments

I would like to express my deepest gratitude to:

- My parents and my siblings for their love and support.
- My supportive friends who have helped me through my Master's studies, including Wei Chieh Wong, Sue Yen Tay, Hui Ean Teh, Harshitha Menon and Nikhil Jain.
- My adviser, Prof. Darko Marinov for exposing me to academic research and for sharing his knowledge, advice, support.
- My co-adviser, Prof. Lin Tan for sharing her knowledge, experience and advice.
- Prof. Gary T. Leavens for his collaborative support that resulted in this thesis.
- Many great professors under whom I have had the privilege to learn and grow as a student, including Prof. Darko Marinov, Prof. Madhu Parthasarathy, Prof. Sarita V. Adve, Prof. Nitin Vaidya, Prof ChengXiang Zhai, and others.
- James Tonyan and Aditya Dasgupta for help in labeling properties inferred by @TCOMMENT and comment-code inconsistencies reported by @TCOMMENT.
- My fellow colleagues who have exposed me to other interesting research and domains of study, including Mathew Alan Kirn, Damion Mitchell, Jurand Noguec, Brett Daniel, Vilas Jagannath, Milos Gligoric, Qingzhou Luo, Yu Lin, Adrian Nistor, Steven Lauterburg, Rajesh Karmani, and others.

# Table of Contents

<b>List of Tables</b> . . . . .	<b>vi</b>
<b>List of Figures</b> . . . . .	<b>vii</b>
<b>List of Abbreviations</b> . . . . .	<b>viii</b>
<b>Chapter 1 Introduction</b> . . . . .	<b>1</b>
1.1 Contributions . . . . .	2
1.2 Organization . . . . .	4
<b>Chapter 2 Background</b> . . . . .	<b>5</b>
2.1 Javadoc Comments . . . . .	5
2.2 Comment-Code Inconsistency . . . . .	7
2.3 Randoop . . . . .	8
<b>Chapter 3 Examples</b> . . . . .	<b>11</b>
3.1 Any Exception . . . . .	11
3.2 Specific Exception and One Null Parameter . . . . .	13
3.3 Specific Exception and Two Null Parameters . . . . .	15
3.4 Expected Exception . . . . .	15
<b>Chapter 4 @tComment Design</b> . . . . .	<b>17</b>
4.1 Inferring Properties from Comments . . . . .	17
4.2 Heuristics . . . . .	20
4.3 Checking Properties in Test Generation . . . . .	21
<b>Chapter 5 Evaluation</b> . . . . .	<b>27</b>
5.1 Experimental Setup . . . . .	27
5.2 Comment-Code Inconsistency Detection . . . . .	28
5.3 Detailed Comment-Code Inconsistency Detection Results . . . . .	29
5.4 Comment Analysis Results . . . . .	32
5.5 Sensitivity of @Randoop Options . . . . .	35
<b>Chapter 6 Related Work</b> . . . . .	<b>36</b>
<b>Chapter 7 Conclusions and Future Work</b> . . . . .	<b>39</b>
<b>References</b> . . . . .	<b>40</b>

# List of Tables

2.1	Subject projects and their number of Javadoc comments per method . . . .	6
4.1	Properties to be extracted . . . . .	18
4.2	Categories of sequences that @Randoop classifies based on matches for parameters with null values. . . . .	23
5.1	Subject projects and their description . . . . .	27
5.2	Subject projects and their basic statistics . . . . .	28
5.3	Overall results for the default configuration of our @Randoop (nullRatio=0.6, timeLimit=3600s) . . . . .	30
5.4	Comment analysis results . . . . .	34

# List of Figures

2.1	Example Javadoc comment . . . . .	6
2.2	Example comment that led us to find a fault in code . . . . .	8
2.3	Example comment in which we found a problem . . . . .	8
2.4	Randoop test-generation algorithm . . . . .	9
2.5	Example contract-violating test generated by Randoop . . . . .	10
2.6	Example regression test generated by Randoop . . . . .	10
3.1	Example test generated by @Randoop. . . . .	12
3.2	Two more example tests generated by @Randoop. . . . .	14
3.3	Example test that Randoop would generate. . . . .	16
3.4	Null-related Javadoc comment helps identify false alarms that Randoop would generate . . . . .	16
4.1	Integration of @Randoop checking of @TCOMMENT-inferred properties into test generation . . . . .	22



# List of Abbreviations

API	Application Programming Interface
LOC	Lines of Codes
NLP	Natural Language Processing
RANDLOOP	RANDom tester for Object-Oriented Programs

# Chapter 1

## Introduction

Source code comments are important artifacts in software and have been around for as long as code has been around. While comments are removed by compilers and not executed, they aid in many software engineering tasks such as code comprehension, reuse, or maintenance. Comments can be broadly categorized into those that appear in the body of a method to describe its inner working and those that appear in the header of a method to describe its specification [29]. Java has standardized the writing of API specifications as Javadoc comments with tags such as `@param` to describe method parameters and `@throws` to describe what exceptions the method could throw. API developers write Javadoc comments to describe their classes and methods. API users often read these comments to understand the code, e.g., an API user can read a Javadoc comment for a method instead of reading the body of the method.

A comment-code inconsistency between the Javadoc comment for a method and the code of that method's body is highly indicative of a fault. First, it can be the case that the comment is correct (in that it properly specifies what the code should do) but the method body has a fault (in that it improperly implements the specification). Second, it can be the case that the method body is correct (in that it properly implements the intended specification) but the comment is incorrect (in that it does not properly describe the intended specification). While the second case does not by itself have an executable fault, it can mislead the users of the method to introduce faults in their code [40].

Because comment-code inconsistencies are indicative of faults, it is important to check for such inconsistencies. However, automating such checking is challenging because it requires

automated understanding of the natural-language text in the comments. While natural-language processing (NLP) techniques have made much progress in the recent decades [30], they are still challenged by ambiguities that are inherent in understanding general text. For example, consider the Javadoc snippet “@param chrono Chronology to use, null means default”, which describes a certain method parameter `chrono` that is an object of type `Chronology`. The part “null means default” is hard to understand because it could specify that `null` is treated in some “default” manner (e.g., throwing a `NullPointerException`) or that `null` is used to represent some default value of `Chronology`.

The only currently viable solution for automated understanding of the natural-language text in the comments is to build *domain-specific analyses*. Tan et al. [40, 41] pioneered automated checking of comment-code inconsistencies based on NLP analysis. Their `iComment` [40] and `aComment` [41] projects focus on systems code written in C/C++ and analyze comments in the domains of locking protocols (e.g., “the caller should grab the lock”), function calls (e.g., “function `f` must be called only from function `g`”), and interrupts (e.g., “this function must be called with interrupts disabled”). Their tools extract rules from such comments and use *static analysis* to check source code against these rules to detect comment-code inconsistencies.

## 1.1 Contributions

This thesis presents a novel approach, called `@TCOMMENT`<sup>1</sup>, for *testing* comment-code inconsistencies in Javadoc comments and Java methods. Specifically, the thesis makes the following contributions.

**New Domain:** We focus `@TCOMMENT` on a *new domain* in comment analysis, specifically method properties for *null values and related exceptions* in Java libraries/frameworks.

---

<sup>1</sup>The name `@TCOMMENT` follows the convention used by prior work (i.e., `iComment` [40] and `aComment` [41]), where the ‘@’ sign represents block tags in Javadoc comments and the word ‘t’ represents testing. The pronunciation of the word “`@TCOMMENT`” is the same as for “at-comment”.

This domain was not studied in the previous work on detecting comment-code inconsistencies, but our inspection of several Java projects showed this domain to be important and widely represented in well-documented Java code. Detecting comment-code inconsistencies in this new domain has its unique challenges that require new solutions, as discussed below.

**Dynamic Analysis:** @TCOMMENT uses a *dynamic analysis* to check comment-code inconsistencies, unlike previous work that used static analysis. Specifically, our @TCOMMENT implementation builds on the Randoop tool [34] for random test generation of Java code. Randoop randomly explores method sequences of the code under test, checks if execution of these sequences violates a set of default *contracts* such as throwing an uncaught exception [34], and generates as tests those sequences that violate some constraint. We modify Randoop to check @TCOMMENT-inferred properties during random test generation and to report violations that correspond to comment-code inconsistencies. We refer to our modified Randoop as *@Randoop*.

We chose dynamic analysis to address the following challenges imposed by the new domain. First, even widely used tools for static checking of Java code, such as FindBugs [23], can have a large number of false alarms when checking properties related to `null` values and exceptions if these properties are available for only some parts of the code, which is the case when inferring properties from Javadoc comments that are not available for all methods. Second, we focus on Java libraries and frameworks, which have few calls to their own methods from their own code bases. Therefore, an analysis cannot focus on callers to these methods to see what parameters they pass in. Instead, a dynamic approach such as Randoop, which generates call sequences and parameters to test library methods, is particularly beneficial.

**Improved Testing:** @Randoop allows us not only to *detect comment-code inconsistencies* but also to *improve test generation* in Randoop. For detecting inconsistencies, @Randoop generates tests that Randoop would not necessarily generate otherwise because these tests need not violate the default contracts that Randoop checks. For improving test gen-

eration, @Randoop identifies some tests as likely false alarms in Randoop and ranks them lower so that developers can focus on the true faults. A false alarm refers to a test that Randoop generates but that does *not* find fault in the code under test, e.g., a test that causes a method to throw exception, but the exception is expected according to the Javadoc comment for the method.

**Evaluation:** We applied @TCOMMENT on seven open-source Java projects (Apache Commons Collections, GlazedLists, JFreeChart, Joda Time, Apache Lucene, Apache Log4j, and Apache Xalan) that have well-developed, well-documented, and well-tested code. We found 28 methods with inconsistencies between Javadoc comments and method bodies in these projects. We reported all inconsistencies, and 12 of them have been already confirmed and fixed by the developers (some by fixing the code and some by fixing the comment), while the rest await confirmation by the developers. @TCOMMENT automatically inferred 2479 properties regarding null values and their related exceptions from Javadoc comments, with a high accuracy of 97–100%. The high accuracy was achieved without using NLP techniques, largely due to the Javadoc null-related comments being well structured with few paraphrases and variants.

## 1.2 Organization

The rest of this thesis is organized as follows. Chapter 3 shows example inconsistencies found with @TCOMMENT. Chapter 4 describes in detail how @TCOMMENT infers properties from Javadoc comments and uses our modified Randoop to check comment-code inconsistencies. Chapter 5 presents our evaluation of @TCOMMENT. Chapter 6 reviews related work, and Chapter 7 concludes the thesis.

# Chapter 2

## Background

This chapter presents the background information necessary for understanding the concepts and techniques introduced later in this thesis. We begin by describing a special type of API documentation, Javadoc comments. We then explain the operation of a test generation tool, Randoop.

### 2.1 Javadoc Comments

Code comments play an important role in software development. As code comments are more flexible and easier to understand than code, programmers often write comments along with code. Programmers write code comments to document the usage of code segments, to express their assumptions and requirement, and to ease future code maintenance.

Java developers follow a standardized way of writing code comments. These code comments are delimited by `/**...*/` (i.e., they start with a double ‘\*’) and can be applied to Java classes, fields, constructors, and methods. They are referred to as Javadoc comments since they are parsed by the Javadoc tool [8]. The Javadoc tool generates well-formatted API documentation from Javadoc comments. The basic structure of a Javadoc comment consists of two parts. The first part is a free-form text that gives an overview description of the corresponding class, field, constructor, or method. The second part consists of more structured block tags that allow generating well-formatted API. The work in this thesis focuses on comments for methods and constructors.

Figure 2.1 shows an example Javadoc comment for a method called `chainedIterator`.

```

/**
 * Gets an iterator that iterates through an array of {@link
 * Iterator}s one after another.
 *
 * @param iterators the iterators to use, not null or empty or
 *                 contain nulls
 * @return a combination iterator over the iterators
 * @throws NullPointerException if iterators array is null or
 *                 contains a null
 */
public static Iterator chainedIterator(Iterator [] iterators)

```

Figure 2.1: Example Javadoc comment

This Javadoc comment has a summary (i.e., “Gets an iterator ... another.”) and the tags `@param`, `@return`, and `@throws`. The summary also contains an inline tag `@link` that links to the documentation of the `Iterator` class. Each `@param` tag is followed by the parameter name (i.e., `iterators`) and the description of the parameter, whereas the `@return` tag is followed by the description of the return value. (If a method has several parameters, then the order of the parameter names for the `@param` tags should follow the order of the parameters in the parameter list.) Similarly, each `@throws` tag has the name of the expected exception (i.e., `NullPointerException`) and the description of the condition under which the exception will be thrown.

Project	# Methods	# Javadoc Comments	Percentage = $\frac{\# \text{Javadoc Comments}}{\# \text{Methods}}$
Collections	3,874	2,434	63%
GlazedLists	2,753	1,741	63%
JFreeChart	6,205	6,186	100%
JodaTime	3,887	2,917	75%
Log4j	2,115	958	45%
Lucene	5,222	2,205	42%
Xalan	5,404	3,229	60%
Total	29,460	19,670	64%

Table 2.1: Subject projects and their number of Javadoc comments per method

Table 2.1 shows the subject projects that will be used in our evaluation (Chapter 5) and the number of Javadoc comments per method in these projects. For example, `JFreeChart` has almost 100% of methods with Javadoc comments. The total number of Javadoc comments (19,670) and the high percentage (42–100%) of methods with Javadoc comments indicate that Javadoc comments are commonly used by API developers to document methods in Java libraries such as those used in our evaluation.

## 2.2 Comment-Code Inconsistency

Code comments are difficult to analyze automatically since many of their parts are written in a natural language (e.g., English). The variation across programmers in terms of style of writing and natural languages used also makes comments harder to parse accurately by a machine. Because code comments are not executed, they could be even more error prone than code.

In general, there are two types of comments. One type of comments explains the usage of a code segment, while the other type states programmers' expectation and constraints [40]. For example, the comment “Gets an iterator that iterates through two Iterators one after another.” in Figure 2.1 belongs to the first type, while the description in the `@throws` tag “throws `NullPointerException` if either iterator is null” belongs to the second type. Since the first type of comments only gives a high-level overview of the code segment, it is potentially less likely to be inconsistent with the code than the second type.

There are several reasons for the cause of inconsistencies between code and comments [40]. One reason could be that software evolves due to changes in requirement or modification to facilitate bug fixes. Another reason could be that mistakes are made by developers due to careless programming or misunderstanding of the code usage.

An inconsistency between code and comment could either indicate a fault in the source code or a bug in the comment. Figure 2.2 shows an example inconsistency between the code



```

/** ...
 * @param anchor the anchor (<code>null</code> not permitted).
 */
public void setRotationAnchor(TextAnchor anchor)

```

Figure 2.2: Example comment that led us to find a fault in code

```

/** ...
 * @param map the map to use to transform the objects
 * @return the transformer
 * @throws IllegalArgumentException if the map is null
 */
public static <I, O> Transformer<I, O>
    mapTransformer(Map<? super I, ? extends O> map)

```

Figure 2.3: Example comment in which we found a problem

and the Javadoc comment that was due to a fault in the code. The comment states that null is not a valid input, but the code executes normally without throwing any exception. We further confirmed this bug by looking at the methods with similar Javadoc comments within the same class which throw `IllegalArgumentException` when null is passed in. This fault has been reported and fixed by the developers. Figure 2.3 shows an example inconsistency between the code and its comments that was due to a bug in the comment. The comment states that `IllegalArgumentException` will be thrown if the parameter `map` is null, whereas the code returns a `NULL_INSTANCE` object when null is passed in.

## 2.3 Randoop

Randoop [34] is a test-generation tool that creates a set of test sequences by randomly choosing methods and their inputs. It takes as input a set of classes and user-specified options, such as time limit and null ratio, and produces as output a set of unit tests that either reveal a fault in the classes under test or capture the current behavior of the classes.

```

GenerateSequences(classes, contracts, timeLimit)
  errorSeqs ← {} // Their execution violates a contract.
  nonErrorSeqs ← {} // Their execution violates no contract.
  while timeLimit not reached do
    // Create new sequence
    m(T1 . . . Tk) ← randomPublicMethod(classes)
    ⟨seqs, val⟩ ← randomSeqsAndVals(nonErrorSeqs, T1 . . . Tk)
    newSeq ← extend(m, seqs, vals)
    // Discard duplicates
    if newSeq ∈ nonErrorSeqs ∪ errorSeqs then
      continue
    end if

    // Execute new sequence and check contracts.
    violated ← execute(newSeq, contracts)
    // Classify new sequence and outputs.
    if violated then
      errorSeqs ← errorSeqs ∪ {newSeq}
    else
      nonErrorSeqs ← nonErrorSeqs ∪ {newSeq}
    end if
  end while
  return ⟨ nonErrorSeqs, errorSeqs ⟩

```

Figure 2.4: Randoop test-generation algorithm

During the test generation, Randoop first randomly select a method call from the classes under test. It then generates a sequence to construct an object (or select null as input with some probability) for each parameter (of non-primitive type) of the selected method. Figure 2.4 shows the pseudo-code for Randoop test-generation algorithm [34].

As random choices ignore program semantics, they may produce redundant and illegal test sequences. After the generation of the test sequences, Randoop executes these sequences and checks them against a set of contracts. The result of the execution is used to determine if the sequence is contract-violating, redundant, or useful in generating more sequences. This technique of gathering feedback during execution is called *feedback-directed random testing*.

Randoop can output two types of test suites. Test suites of the first type contain *contract-*

```

// In automatically generated test class:
void testViolation() {
    java.util.Iterator [] var1 = null;
    java.util.Iterator var2 =
        org.apache.commons.collections.IteratorUtils.
            chainedIterator(var1);
}

```

Figure 2.5: Example contract-violating test generated by Randoop

```

// In automatically generated test class:
void testRegression() {
    java.util.Collection var0 = null;
    java.lang.Object var2 = null;
    org.apache.commons.collections.Predicate var3 =
        org.apache.commons.collections.PredicateUtils.
            equalPredicate((java.lang.Object)(byte)0);
    org.apache.commons.collections.CollectionUtils.filter(var0, var3);
    //Regression assertion(captures the current behavior of the code)
    assertNotNull(var3);
}

```

Figure 2.6: Example regression test generated by Randoop

*violating tests* that have some violation of a property that a class, an object, or a method is expected to preserve. Test suites of the second type contain *regression tests* that capture the actual behavior of the current implementation.

Figure 2.5 shows an example contract-violating test generated by Randoop for the method `chainedIterator` from Figure 2.1. The test fails due to `NullPointerException` being thrown during execution. Figure 2.6 shows an example regression test generated by Randoop. The assertion in the test captures the current state (i.e., not null) of the `Predicate` object.

# Chapter 3

## Examples

We illustrate how `@TCOMMENT` can be used by showing three examples of comment-code inconsistencies that we found with `@TCOMMENT` in two projects, and one example of a Randoop false alarm identified by `@TCOMMENT` where a method throws an exception but it is expected according to the relevant comment. The first three examples show progressively more complex cases: (1) inferring that *some* exception should be thrown when a method parameter is `null`; (2) inferring what *type* of exception should be thrown when *one* method parameter is `null`; and (3) inferring what *type* of exception should be thrown when *two* method parameters are `null`. The last example shows a Randoop false alarm.

### 3.1 Any Exception

Consider first the `JodaTime` project [10], a widely used Java library for representing dates and times. `JodaTime` provides several classes that support multiple calendar systems. `JodaTime` code is fairly well commented with Javadoc comments (about 75% of `JodaTime` methods have Javadoc comments as shown in Table 2.1). We ran `@TCOMMENT` to first infer properties for the methods in the `JodaTime` code and then to check these properties with our `@Randoop`. For each test that `@Randoop` generates, it marks whether the test, when executed, violated some `@TCOMMENT`-inferred property or a default Randoop contract.

Figure 3.1 shows an example test that violates an `@TCOMMENT`-inferred property. This test creates a `MutablePeriod` object `var2` and invokes several methods on it. The executions of `setSeconds` and `setValue` methods finish normally, but for `setPeriod`, `@TCOMMENT`

```

// In automatically generated test class:
void test1() {
    org.joda.time.MutablePeriod var2 =
        new org.joda.time.MutablePeriod(1L, 0L);
    var2.setSeconds(0);
    var2.setValue(0, 0);
    org.joda.time.Chronology var9 = null;
    try {
        var2.setPeriod(0L, var9);
        fail("Expected some exception when chrono==null");
    } catch (Exception expected) {}
}

// In a class under test:
/** ...
 * @param duration the duration, in milliseconds
 * @param chrono the chronology to use, not null
 * @throws ArithmeticException if the set exceeds
 *         the capacity of the period
 */
void setPeriod(long duration, Chronology chrono)

```

Figure 3.1: Example test generated by @Randoop. Method under test and its comment.

reports that there is a likely comment-code inconsistency: the parameter `var9` is `null`, but the method execution throws no exception, which disagrees with the corresponding comment indicating that some exception should be likely thrown. Note that this test *passes* if some exception is thrown and fails otherwise.

Figure 3.1 also shows the relevant parts of the `setPeriod` method. It has two parameters, and the Javadoc comment provides a description for each of them. As discussed in Section 2.1, a typical Javadoc comment has the main, free-flow text (for brevity omitted in our examples) and specific *tags/clauses* such as `@param`, `@throws`, `@return`, etc. We call the entire block of text before a method *one comment* with several *comment tags*. Figure 3.1 shows one Javadoc comment with two `@param` tags and one `@throws` tag.

The key part here is “not null” for the `chrono` parameter. `@TCOMMENT` infers the

property that whenever `chrono` is `null`, the method should throw *some* exception (although it does not know which exception type it should be because the tag for `ArithmeticException` is not related to `null`). `@Randoop` finds that the shown test violates this property. Note that it may not be a comment-code inconsistency; the inference could have been wrong, e.g., “`not null`” could represent the method precondition—such that if `chrono` is `null`, the method could do anything and is not required to throw an exception—or “`not null`” could be a part of a larger phrase, say, “`not a problem to be null`”—such that the method must not throw an exception.

In this case, our inspection showed that `@TCOMMENT` performed a correct inference and detected a real comment-code inconsistency. In fact, `@TCOMMENT` also found a similar inconsistency for another overloaded `setPeriod` method. We reported both inconsistencies in the `JodaTime` bug database [11], and `JodaTime` developers fixed them by changing comments. It is important to note that `Randoop` would have *not* generated this example test because it does not throw an exception. More precisely, `Randoop` internally produces the method sequence but would not output it to the user as a possibly fault-revealing test. Indeed, `@Randoop` generates the test precisely because it does not throw any exception when some exception is expected according to the comment.

## 3.2 Specific Exception and One Null Parameter

Consider next the `Apache Commons Collections` project (which we will call `Collections` for short) [1], a popular library for representing collections of objects. Figure 3.2 shows two example tests, each of which violates an `@TCOMMENT`-inferred property, and the corresponding method declarations and their comments.

For the `synchronizedMap` method, `@TCOMMENT` *correctly* infers that the method should throw `IllegalArgumentException` when the parameter `map` is `null`; while this is explicit in the `@throws` tag, note that the `@param map` tag could be contradicting by allowing any

```

// In automatically generated test class:
void test2() {
    java.util.Map var0 = null;
    try {
        java.util.Map var1 = org.apache.commons.collections.MapUtils.
            synchronizedMap(var0);
    } catch (IllegalArgumentException expected) {return;}
    fail("Expected IllegalArgumentException," +
        "got NullPointerException");
}
void test3() {
    java.util.Collection var0 = null;
    java.util.Iterator [] var1 = new java.util.Iterator []{};
    try {
        org.apache.commons.collections.CollectionUtils.
            addAll(var0, (java.lang.Object []) var1);
    } catch (NullPointerException expected) {return;}
    fail("Expected NullPointerException when collection==null");
}

// In classes under test:
/** ...
 * @param map the map to synchronize, must not be null
 * @return a synchronized map backed by the given map
 * @throws IllegalArgumentException if the map is null
 */
static Map synchronizedMap(Map map)

/** ...
 * @param collection the collection to add to, must not be null
 * @param elements the array of elements to add, must not be null
 * @throws NullPointerException if the collection or array is null
 */
static void addAll(Collection collection, Object [] elements)

```

Figure 3.2: Two more example tests generated by @Randoop. Methods under test and their comments.

behavior when `map` is `null`. Inferring a specific type of expected exception is unlike in the previous example when `@TCOMMENT` could only infer that *some* exception should be

thrown. Indeed, inferring the type in this case is important because, when `map` is `null`, the method does throw an exception but of a different type. In this case, even the original Randoop generates `test2` because it throws an exception. However, Randoop also generates dozens of others tests that are not fault-revealing, so this comment-code inconsistency would be proverbially “the needle in a haystack” among the tests generated by Randoop. In contrast, `@Randoop` prominently highlights the inconsistency. We reported this comment-code inconsistency [2], and the `Collections` developers fixed it by removing the `@throws` part of the comment.

### 3.3 Specific Exception and Two Null Parameters

For the `addAll` method in the same project, `@TCOMMENT` *correctly* infers that the method should throw `NullPointerException` when either the parameter `collection` is `null` *or* the parameter `elements` is `null`. This is similar to the previous case where the specific exception type is inferred, but in this case two method parameters are involved. The inconsistency that `@TCOMMENT` finds is, in fact, related to the situation where only one parameter—`collection`—is `null` (while the array is empty), but the exception is not thrown as expected. We also reported this comment-code inconsistency [3], and it is under consideration.

### 3.4 Expected Exception

For all examples presented so far, an exception is expected (according to the Javadoc comments), but the method under test either does not throw an exception or throws an exception of a different type. We next discuss an example where an exception is thrown by the method under test, but it is expected as indicated by the relevant comment (Figure 3.4). This example illustrates a case where Randoop would generate a false alarm, but `@TCOMMENT` can help remove this false alarm automatically by lowering its ranking



```

// In automatically generated test class:
public void test4() {
    ca.odell.glazedlists.DebugList var0 =
        new ca.odell.glazedlists.DebugList();
    ca.odell.glazedlists.DebugList var1 =
        new ca.odell.glazedlists.DebugList();
    boolean var2 = var0.retainAll((java.util.Collection)var1);
    boolean var4 = var1.add((java.lang.Object)(byte)0);
    java.lang.Object var5 = null;
    int var6 = var1.indexOf(var5);
    java.lang.Object [] var7 = null;
    java.lang.Object [] var8 = var1.toArray(var7);
}

```

Figure 3.3: Example test that Randoop would generate.

```

// In a class under test:
/** ...
 * @param array the array into which the ... purpose.
 * @return an array containing the elements of this list.
 * @throws ArrayStoreException if the runtime type ... list.
 * @throws NullPointerException if the specified array is
 *         <tt>null</tt>.
 */
public <T> T[] toArray(T[] array)

```

Figure 3.4: Null-related Javadoc comment helps identify false alarms that Randoop would generate

through analyzing the comment. `@TCOMMENT` infers that the method `toArray` should throw a `NullPointerException` if `null` is passed to the `array` parameter. Randoop generates a test (Figure 3.3) for `toArray` when `array` is `null`, because the execution of this test indeed throws a `NullPointerException`, and Randoop reports as potentially fault-revealing all tests that throw uncaught exceptions during execution. In contrast, `@TCOMMENT` marks that the exception is expected according to the comment, which lowers the ranking of this false alarm to improve the testing accuracy.

# Chapter 4

## @tComment Design

Our @TCOMMENT approach consists of two components. The first component takes as input the source code for a Java project, automatically analyzes the English text in the Javadoc comments in the project, and outputs a set of inferred likely properties for a method. The second component takes as input the same code and inferred properties, generates random tests for the methods in the code, checks the inferred properties, and reports inconsistencies.

Similar to prior work [40, 41], we build a domain-specific comment analysis, due to the difficulty of inferring arbitrary properties from general comments. In particular, we focus on null-pointer related comments because null-pointer dereferences are common memory bugs [18], and a large percentage of Javadoc comments (24.2% in the seven projects we evaluated) contain the keyword `null`.

### 4.1 Inferring Properties from Comments

Our goal is to infer from Javadoc comments null-related properties about method parameters. For a parameter that is of non-primitive type and can take a `null` value, @TCOMMENT infers one of these four kinds of properties: (1) *Null Normal*: if the parameter is `null`, the method should execute normally (and throw no exception); (2) *Null Any Exception*: if the parameter is `null`, the method should throw *some* exception; (3) *Null Specific Exception*: if the parameter is `null`, the method should throw a *specific* type of exception; or (4) *Null Unknown*: if the parameter is `null`, we do not know the expected behavior of the method.

Properties	Meaning	Comment Example	Notation
<i>Null Normal</i>	If the parameter is null, the method should execute normally (and throw no exception).	@param predicate the predicate to use, may be null	predicate==null => normal
<i>Null Any Exception</i>	If the parameter is null, the method should throw <i>some</i> exception.	@param collection the collection to add to, must not be null	collection==null => exception
<i>Null Specific Exception</i>	If the parameter is null, the method should throw a <i>specific</i> type of exception.	@throws IllegalArgumentException if the id is null	id==null => IllegalArgumentException
<i>Null Unknown</i>	We do not know the expected behavior of the method when the parameter is null.	@param array the array over which to iterate	array==null => unknown

Table 4.1: Properties to be extracted. The comment examples are real comments from the projects used in our evaluation.

Table 4.1 shows four examples of comment tags and their corresponding inferred properties. For example, `@TCOMMENT` infers from the second tag, “`@param collection the collection to add to, must not be null`”, that if the method parameter `collection` is null, then the method is expected to throw some exception, represented as *collection == null => exception*. Based on our experience with the null-related Javadoc comments, we found that this interpretation—the method should throw an exception—often matches developers’ intention, and thus we adopted it for `@TCOMMENT`. The comment-code inconsistencies that we reported and developers confirmed, as well as the low false-positive rate of our detection for reported comment-code inconsistencies, confirm our decision. However, note that we focus on *library projects*, where the methods need not trust their callers. The interpretation may differ for applications with more of a design-by-contract mentality where callers were trusted more. As discussed earlier in the introduction, this example tag could have another interpretation, describing a precondition such that passing null value for the parameter `collection` allows the method to do anything, not necessarily throw an exception.

Our `@TCOMMENT` implementation leverages the Javadoc doclet [13] to parse Javadoc comments. For example, consider the above tag for `collection`; Figure 3.2 shows the comment for this tag and the corresponding method declaration. The Javadoc doclet parses this tag and outputs the name of the method parameter (`collection`), its type (`java.util.Collection`), the corresponding method’s full signature (`org.apache.commons.collections.CollectionUtils#addAll(Collection collection, Object[] elements)`), and the free-form comment text (“`the collection to add to, must not be null`”). The method parameter, its type, and the full method signature are used later by `@Randoop`, the test-generation component of `@TCOMMENT`, to check the generated tests.

`@TCOMMENT` first extracts all Javadoc `@param` (for parameters of non-primitive types) and `@throws` tags that contain a non-empty free-form text, since a Javadoc tag with an empty free-form text does not describe a *Null Normal*, *Null Any Exception*, or *Null Specific Exception*

property. Then, `@TCOMMENT` infers *Null Normal* and *Null Any Exception* properties from the `@param` comment tags and *Null Specific Exception* properties from the `@throws` comment tags. It assigns *Null Unknown* to a method parameter if neither its `@param` tag nor `@throws` tag describes any other property. In this thesis, however, we do not count *Null Unknown* toward the 2479 inferred properties, since one cannot test against these *Null Unknown* properties.

## 4.2 Heuristics

`@TCOMMENT` uses three relatively simple heuristics to analyze the free-form text. While the heuristics are not perfect, our empirical evaluation shows that they are highly accurate in practice. First, if negation words, such as “not” or “never”, are found *up to three words* before or after the word `null`—e.g., “the collection to add to, must not be null” has “not” two words from `null`—`@TCOMMENT` infers the *Null Any Exception* property. If no negation words are found up to three words around the word `null`—e.g., the first tag in Table 4.1—`@TCOMMENT` infers the *Null Normal* property.

Second, for `@throws` tags—e.g., “@throws `IllegalArgumentException` if the id is null”—the Javadoc doclet parses the comment tag and outputs the specific exception (`IllegalArgumentException`) and the free-form text (“if the id is null”). If the text contains the keyword `null`, `@TCOMMENT` simply splits the text into words and searches each word in the list of all method parameter names generated by the Javadoc doclet. If a valid parameter name is found, e.g., `id`, `@TCOMMENT` infers the property `id == null => IllegalArgumentException` from the comment tag.

Third, if the keyword “or” or “either” is in the `@throws` comment text, e.g., “@throws `NullPointerException` if the collection or array is null” in Figure 3.2, `@TCOMMENT` generates multiple properties, e.g., `collection == null => NullPointerException` and `array == null => NullPointerException`. If both *Null Any Exception* and *Null Specific Exception* properties are inferred for the same method parameter, e.g., `collection`,

@TCOMMENT keeps only the *Null Specific Exception* property.

### 4.3 Checking Properties in Test Generation

After @TCOMMENT infers likely method properties, it uses our modified Randoop, called *@Randoop*, to check these properties using random test generation. Figure 4.1 shows the simplified pseudo-code of the Randoop test-generation algorithm [34], together with our extension for checking @TCOMMENT-inferred properties.

We briefly summarize how Randoop works. It produces random *sequences* of method calls (including constructors) of the code under test. It maintains a set of error sequences (to be output as generated unit tests that are likely fault revealing) and a set of non-error sequences (to be used for creating longer sequences). In a loop, it first randomly selects a method  $m$  whose  $k$  parameters (including the receiver for non-static methods) have types  $T_1 \dots T_k$ . It then selects sequences (previously generated) and values (e.g., “0”, “1L”, or “null”) of appropriate type to use for the method parameters. It concatenates these sequences and adds a new call to  $m$ . It then executes the new sequence to check *contracts* (e.g., no uncaught exception during execution). If there is a violation, it adds the new sequence to the error sequences; otherwise, it adds the new sequence to the non-error sequences. More details of the original Randoop algorithm, including discarding duplicates and filtering extensible sequences, are available elsewhere [34].

```

// inferredProperties is specific to @Randoop
GenerateSequences(classes, contracts, inferredProperties, timeLimit)
  errorSeqs ← {} // These will be generated as unit tests
  // we add a comment-code inconsistency field to sequences
  nonErrorSeqs ← {} // These are used to build longer sequences
  while timeLimit not reached do
    // Create new sequence
    m(T1 . . . Tk) ← randomPublicMethod(classes)
    ⟨seqs, val⟩ ← randomSeqsAndVals(nonErrorSeqs, T1 . . . Tk)
    newSeq ← extend(m, seqs, vals)
    // Execute new sequence and check contracts.
    violated ← execute(newSeq, contracts)
    // Classify new sequence and outputs.
    if violated then
      errorSeqs ← errorSeqs ∪ {newSeq}
    else
      nonErrorSeqs ← nonErrorSeqs ∪ {newSeq}
    end if
    // Execute and check @TCOMMENT-inferred properties.
    match ← execute(newSeq, inferredProperties)
    if match = 'Missing Exception' then
      // Add the new sequence, marked as inconsistency
      errorSeqs ← errorSeqs ∪ { newSeq }
      newSeq.isCommentCodeInconsistency ← highlyLikely
    else if match = 'Different Exception' or
      match = 'Unexpected Exception' then
      // Mark an already added sequence as inconsistency
      newSeq.isCommentCodeInconsistency ← likely
    else if match = 'Unknown Status' then
      // Unknown inconsistency status
      newSeq.isCommentCodeInconsistency ← unknown
    else // match = 'Expected Exception'
      // Mark the sequence as likely consistent
      newSeq.isCommentCodeInconsistency ← unlikely
    end if
  end while
  return ⟨ nonErrorSeqs, errorSeqs ⟩

```

Figure 4.1: Integration of @Randoop checking of @TCOMMENT-inferred properties into test generation

	Comment-Code Inconsistent			Unknown	Comment-Code Consistent
	<i>Missing Exception</i>	<i>Different Exception</i>	<i>Unexpected Exception</i>		
Exception Thrown during execution	No	Yes	Yes	Yes	Yes
Properties that @TCOMMENT inferred about expected exceptions for method parameters with null values	<i>Null Any Exception</i> or <i>Null Specific Exception</i>	(1) at least one is <i>Null Specific Exception</i> & (2) thrown exception is not in the set of all specific exceptions	(1) at least one is <i>Null Normal</i> & (2) there is no <i>Null Specific Exception</i> or <i>Null Any Exception</i>	<i>Null Unknown</i> for all method parameters with null	(1) at least one is <i>Null Specific Exception</i> or <i>Null Any Exception</i> & (2) thrown exception is in the set of expected exceptions

Table 4.2: Categories of sequences that @Randoop classifies based on matches for parameters with null values.



Our `@Randoop` modification follows the similar approach that Randoop performs for checking contracts: `@Randoop` executes the sequence and checks the `@TCOMMENT`-inferred properties for method calls where one or more parameter have `null` values. We distinguish five kinds of matches between method execution (does it throw an exception and of what type) and `@TCOMMENT`-inferred properties (is an exception expected and of what type). Based on the match, `@Randoop` can (1) generate a sequence that Randoop would not generate otherwise, (2) generate the same sequence as Randoop but mark the sequence as a comment-code inconsistency, or (3) generate the same sequence as Randoop but mark the sequence as comment-code consistent.

Before we describe the five possible kinds of matches, we describe how `@Randoop` computes the *set* of expected exceptions. It handles multiple `null` values, which naturally arise for methods with several parameters of non-primitive types, e.g., `addAll` from Figure 3.2 has “`Collection collection`” and “`Object[] elements`”. If only one of these parameters is `null`, `@Randoop` uses only the property inferred for that parameter. If two or more parameters are `null`, `@Randoop` puts in the set all expected exceptions for these parameters, e.g., if we had `collection == null => NullPointerException` and `elements == null => IllegalArgumentException`, then `@Randoop` would assume that either of the two exceptions is expected. If some parameter with `null` value has the *Null Any Exception* property, then all types of exceptions are expected. Finally, `@Randoop` adds to the set exceptions that are not null-related but appear in the `@throws` tags for the Javadoc comment of the method under test. For example, the method `arrayIterator(Object[] array, int start, int end)` has one such tag (“`@throws IndexOutOfBoundsException if array bounds are invalid`”) and one null-related tag (“`@throws NullPointerException if array is null`”); although `IndexOutOfBoundsException` does not correspond to a null input, `@TCOMMENT` always adds it to the set of expected exceptions.

Table 4.2 lists the five kinds of matches, which are discussed below.

**Missing Exception** sequences throw no exception during execution, but the corresponding

inferred properties specify that some exception is expected. These sequences should be generated as test cases that are likely comment-code inconsistencies (although they could be false alarms if the inference obtained incorrect properties from the corresponding comments). To repeat, these test cases would *not* be generated by the original Randoop because they throw no exception.

***Different Exception*** sequences throw an exception that is different from the exception(s) expected according to the inferred properties. These are also likely comment-code inconsistencies. These sequences would be generated by Randoop as potentially fault-revealing test cases, and by inspecting them, the developer might find the inconsistency. However, these inconsistencies would be hard to identify among a large number of other test cases that Randoop generates (in our evaluation in Section 5.2, there are only 4 test cases were *Different Exception* among 1,285 test cases that Randoop would generate). In contrast, @Randoop highlights these test cases.

***Unexpected Exception*** sequences throw an exception whereas @TCOMMENT explicitly expects normal execution with no exception. As for *Different Exception*, Randoop would also generate these sequences as test cases, but @Randoop highlights them due to the inconsistency between code and properties inferred from comments.

***Unknown Status*** sequences throw an exception but @TCOMMENT inferred no property to tell if the exception is expected or not. Both Randoop and @Randoop generate these as error sequences. While they may indicate a fault in the code, they do not show an inconsistency between code and comment (unless the inference incorrectly missed inferring some property about exceptions).

***Expected Exception*** sequences throw an exception, but this exception is expected according to the properties inferred by @TCOMMENT from the relevant comments. Hence, @Randoop marks these sequences as consistent. If @TCOMMENT inference is correct for these cases, they are false alarms that Randoop would generate; if the inference is incorrect, @TCOMMENT would increase the time for developers to find the true fault-revealing tests.

Our current implementation (Figure 4.1) of @Randoop *modifies only the checking and not the random selection* performed by Randoop. Randoop randomly selects methods to test and parameter values for the methods, and @Randoop does not perform any additional selection that the original Randoop does not perform. It could be beneficial to additionally bias the selection based on the properties inferred from comments. For example, if it is inferred that an exception should be thrown when a method parameter `p` is `null` (i.e., *Null Specific Exception* or *Null Any Exception*), but Randoop does not select any sequence where `p` is `null`, @Randoop could be extended to (non-randomly) generate such additional sequences to check if the inferred property holds. This extension remains as future work.

# Chapter 5

## Evaluation

### 5.1 Experimental Setup

We evaluate @TCOMMENT on seven open-source Java projects. Table 5.1 and Table 5.2 list information about these projects. We modified Randoop revision 652 to build @Randoop. Randoop provides several options that control random generation, and we consider two options that are the most important for @Randoop: (1) *nullRatio* specifies the frequency that the `null` value is randomly chosen as an input for a method parameter of some non-primitive type (e.g., a `nullRatio` of 0.6 instructs Randoop to use `null` 60% of the time, and it never uses `null` for method receivers because that would directly lead to `NullPointerException`); and (2) *timeLimit* specifies the number of seconds that Randoop should generate tests for one project. All experiments were performed on a machine with a 4-core Intel Xeon 2.67GHz processor and 4GB of main memory, running Linux version 2.6.18, and Java HotSpot 64-Bit Server VM, version 1.6.0\_20.

Project	Source	Description	Version
<code>Collections</code>	[1]	Collection library and utilities	3.2.1
<code>GlazedLists</code>	[7]	List transformations in Java	1.8
<code>JFreeChart</code>	[9]	Chart creator	1.0.13
<code>JodaTime</code>	[10]	Date and time library	1.6.2
<code>Log4j</code>	[4]	Logging service	1.2
<code>Lucene</code>	[5]	Text search engine	2.9.3
<code>Xalan</code>	[6]	XML transformations	2.7.1

Table 5.1: Subject projects and their description

Project	# LOC	# Classes	# Methods
Collections	19,417	274	3,874
GlazedLists	19,203	239	2,753
JFreeChart	51,376	396	6,205
JodaTime	18,428	154	3,887
Log4j	14,425	221	2,115
Lucene	38,051	422	5,222
Xalan	53,642	510	5,404

Table 5.2: Subject projects and their basic statistics

## 5.2 Comment-Code Inconsistency Detection

Table 5.3 shows the overall results of @TCOMMENT with the default values for @Randoop options. (Section 5.5 discusses the sensitivity of the results to the value of these options and the selection of the default @Randoop values.) For each project, we tabulate the number of tests that @Randoop generated based on the five kinds of matches between inferred properties and method executions (as described in Section 4.3). For three kinds of matches that could have comment-code inconsistencies, the cells show the split into True Inconsistencies and False Alarms. The last column also shows the number of @TCOMMENT properties that @Randoop checked during test generation.

In total, @Randoop generated 68 tests with potential comment-code inconsistencies. Note that Randoop would not generate 30 of those (column ‘*Missing Exception*’) where methods execute normally while exceptions are expected as specified by the corresponding comments. @Randoop also generates 4 tests where an exception is thrown but different than specified by comments (column ‘*Different Exception*’) and 34 tests where an exception is thrown but normal execution was expected (column ‘*Unexpected Exception*’). @Randoop generates 1507 tests that throw an exception for cases where no null-related properties were inferred (column ‘*Unknown Status*’). Last but not least, @Randoop identifies 232 tests as throwing some exceptions expected by the comments (column ‘*Expected Exception*’).

## 5.3 Detailed Comment-Code Inconsistency Detection Results

The cells with sums show the split of comment-code inconsistencies reported by @TCOMMENT into those truly inconsistent (summarized in row ‘True Inconsistencies’) and not (row ‘False Alarms’). We inspected all the reports by carefully reading the comments and the code to determine if they are indeed inconsistent or not. A subset of reports was also independently inspected by two more students. @TCOMMENT detected 28 previously unknown comment-code inconsistencies and had 40 false alarms.

Project	Missing Exception = TI + FA	Different Exception = TI + FA	Unexpected Exception = TI + FA	Unknown Status	Expected Exception	Tested Properties
Collections	12 = 12 + 0	4 = 3 + 1	6 = 0 + 6	94	36	115
GlazedLists	0 = 0 + 0	0 = 0 + 0	6 = 1 + 5	151	1	11
JFreeChart	1 = 1 + 0	0 = 0 + 0	2 = 2 + 0	127	6	42
JodaTime	3 = 3 + 0	0 = 0 + 0	13 = 0 + 13	37	3	31
Log4j	1 = 1 + 0	0 = 0 + 0	3 = 0 + 3	186	152	179
Lucene	4 = 0 + 4	0 = 0 + 0	2 = 1 + 1	368	2	12
Xalan	9 = 4 + 5	0 = 0 + 0	2 = 0 + 2	544	32	43
Total	30 = 21 + 9	4 = 3 + 1	34 = 4 + 30	1507	232	433
Total						
True Inconsistencies (TI)	21	3	4	28		
False Alarms (FA)	9	1	30	40		

Table 5.3: Overall results for the default configuration of our @Randoop (nullRatio=0.6, timeLimit=3600s)

The sources of false alarms are incorrectly inferred properties for the method with reported violation itself (11 out of 40), missing properties—‘*Null Unknown*’— (11 out of 40), and incorrect/missing properties for another method in the sequence (18 out of 40). As an example of the first source, @TCOMMENT inferred the property *filter == null => exception* from “@param filter if non-null, used to permit documents to be collected”, because the negation word `non` is next to `null`. However, this comment tag does not imply the parameter `filter` cannot be null. Advanced NLP analysis may be leveraged to analyze this tag correctly. Section 5.4 discusses inference accuracy in detail. The second source of false alarm is when a method has several `null` parameters, at least one parameter missing property (*Null Unknown*) and at least one not missing (either *Null Normal*, *Null Any Exception*, or *Null Specific Exception*). @Randoop reports an inconsistency if the method throws an unexpected exception, even if the parameter with *Null Unknown* caused it, because @Randoop does not identify the cause of exception. For instance, @TCOMMENT inferred the properties *zone == null => normal* and (implicitly) two *Null Unknown* for the parameters from the method `getInstance(DateTimeZone zone, long gregorianCutover, int minDaysInFirstWeek)`. If the execution throws `IllegalArgumentException` when the input `minDaysInFirstWeek` is out of the allowed range, @Randoop still reports this as a potential violation of *Null Normal* for `zone`. The third source results in some null values propagating through fields. For example, @TCOMMENT may not know that a constructor for some class should not have a null parameter. If a test passes null, the value is set to some field. Later on, a method can dereference that field and throw an exception; if the method itself has some null parameters, @Randoop can falsely report an inconsistency (again because it does not identify the real cause of exception).

It is worth pointing out that we set all @TCOMMENT options (e.g., distance of negation words from `null`, treating preconditions as requiring exceptions, ignoring *Null Unknown*, etc.) by looking only at the first six projects. The experiments with `Xalan` were performed with the same options.



## 5.4 Comment Analysis Results

Table 5.4 shows the comment analysis results of @TCOMMENT. Columns ‘param’ and ‘param with null’ show, for parameters of non-primitive types, the total number of @param tags that @TCOMMENT analyzed and the number of @param tags that contain the keyword null, respectively. Similarly, columns ‘throws’ and ‘throws with null’ show the total number of @throws tags that @TCOMMENT and the number of @throws tags that contain the keyword null, respectively. In total, there are 2713 @param and @throws tags that contain the keyword null in the seven evaluated projects. @TCOMMENT inferred 2479 *Null Normal*, *Null Any Exception*, and *Null Specific Exception* properties from these comments. Of these 2479 properties, 433 are tested by @TCOMMENT to automatically detect comment-code inconsistencies and improve testing. As discussed in Section 4.3, it would be beneficial to modify Randoop’s random selection to actively test more of the inferred properties in the future. For all other method parameters of non-primitive types that can take a null value, but for which we cannot infer any of the three null-related properties, we assign them the *Null Unknown* properties.

To evaluate the accuracy of our automatic comment-analysis technique, we randomly sample 100 @param (for parameters of non-primitive types) and @throws tags with non-empty free-form text from each project, and manually read them and the corresponding inferred properties to check if the inferred properties are correct. The *accuracy* is calculated as the number of correctly analyzed tags in a sample over the total number of tags in the sample. Note that the manual inspection is purely for evaluating the accuracy of our comment analysis; @Randoop *directly* uses the automatically inferred properties to detect comment-code inconsistencies and improve testing, and no manual inspection of the inferred properties is required.

In addition, we present the standard precision and recall for *Null Normal*, *Null Any Exception*, and *Null Specific Exception* respectively. For example, the precision for *Null Normal*

is the proportion of identified *Null Normal* properties that indeed are *Null Normal* properties. The recall for *Null Normal* is the proportion of true *Null Normal* properties in our sample that @TCOMMENT identifies.

Our analysis of the free-form comment text achieves a high accuracy of 97–100% (Column ‘Accuracy %’) without using NLP techniques as iComment did [40]. In addition, the precisions and recalls are in high nineties to 100% in most cases. One exception is that the precision for *Null Normal* in `Xalan` is only 50%, where only two *Null Normal* properties are inferred, and one was inferred incorrectly. The general high performance is partially due to the Javadoc API comments being much more structured than the comments in systems code written in C. There is also less variance in paraphrases and sentence structures in the Javadoc comments than in the C/C++ comments in systems code. While the general idea of detecting comment-code inconsistencies through testing should be applicable to C/C++ projects, the comment analysis component may need to leverage more advanced techniques as iComment did [40].

If some null-related properties are described in non-Javadoc style comments, e.g., without using the @param tag, @TCOMMENT would not analyze them. As we do not anticipate many such comments, this thesis focused on properly tagged Javadoc comments.

Project	@param	@throws	@param with null	@throws with null	Properties = $N + A + S$	Precision [%]			Recall [%]			Accuracy [%]
						$N$	$A$	$S$	$N$	$A$	$S$	
Collections	700	431	271	207	$347=81+47+219$	75	92	100	100	100	97	97
GlazedLists	110	55	14	17	$19=14+0+5$	100	100	100	100	100	100	100
JFreeChart	1808	75	902	3	$902=362+537+3$	100	100	100	100	100	100	100
JodaTime	802	726	529	96	$553=445+23+85$	100	75	100	100	100	78	98
Log4j	713	0	488	0	$460=243+217+0$	100	100	100	100	100	100	100
Lucene	498	373	39	15	$67=13+25+29$	100	67	100	80	100	100	99
Xalan	699	110	126	6	$131=33+93+5$	50	100	100	100	100	100	99
Total/Overall	5330	1770	2369	344	$2479=1191+942+346$	98	98	100	99	100	93	99

Table 5.4: Comment analysis results.  $N$  is for *Null Normal*,  $A$  for *Null Any Exception*, and  $S$  for *Null Specific Exception*.

## 5.5 Sensitivity of @Randoop Options

We want to understand how different values for @Randoop options `nullRatio` and `timeLimit` affect our results of comment-code inconsistency detection, which can help us identify good default values for using @TCOMMENT. When time budget allows, users can always run @TCOMMENT with many `nullRatios` and `timeLimits` to try to detect more inconsistencies.

We run @Randoop with 5 `timeLimits`—50sec, 100sec, 200sec, 400sec, and 800sec—and 11 `nullRatios`—from 0.0 to 1.0 in increments of 0.1—on all seven projects and measured the most important metric, the number of *Missing Exception* tests. These are  $5 \times 11 \times 7$ , a total of 385, sets of experiments. Despite the randomness in @Randoop, it identifies more *Missing Exception* tests (thus potentially detects more comment-code inconsistencies) as the `timeLimit` increases for all cases but one combination of the value and the project. We found that when running @Randoop for 800sec, `nullRatio` 0.6 helps @Randoop identify the largest number of *Missing Exception* tests across all seven projects. Therefore, we chose it as the default `nullRatio` value. We found that 0.3, 0.5, 0.7, 0.8, and 0.9 are the next best values for `nullRatios` for these seven projects based on the same metric. Note that 0.0 is clearly not good as it never selects `null`, but also 1.0 is not good as it always selects `null` and thus Randoop cannot “grow” bigger objects with non-`null` parameters.

To further understand the effect of `timeLimits`, we increased the `timeLimits` to up to two hours with `nullRatio` 0.6. We found that the number of the *Missing Exception* tests reaches a plateau at about one-hour mark, which is similar to the fact that the original Randoop reaches a plateau around one-hour mark [33].

To further understand the effect of `nullRatios`, we performed additional experiments with `timeLimit` one hour and `nullRatios` from 0.3 to 0.9. The results show that they produce almost identical numbers for the five kinds of matches, suggesting that if one runs @Randoop for an hour, one can pick any `nullRatio` from 0.3 to 0.9 to obtain similar results.

# Chapter 6

## Related Work

**Automated Software Testing.** Many automated software testing techniques are designed to detect software faults [15, 19, 20, 24, 32, 34, 46], e.g., based on random generation or using specifications. @TCOMMENT leverages an additional source—code comments—and modifies Randoop [34] to detect more faults (in both code and comments), and to potentially identify false alarms generated by Randoop. It is quite conceivable to extend @TCOMMENT to improve other automated testing techniques.

**Detecting Comment-Code Inconsistencies.** iComment [40] and aComment [41] extract rules from comments and check source code against these rules *statically* to detect inconsistencies between comments and code. The differences between iComment/aComment and @TCOMMENT have already been discussed in detail in the introduction, so we only summarize them here: (1) @TCOMMENT leverages a new type of comments, related to null values; (2) @TCOMMENT employs a dynamic approach to check comments during testing; and (3) in addition to finding comment-code inconsistencies, @TCOMMENT could find false alarms generated by test-generation tools such as Randoop.

A recent empirical study [26] examines the correlation between code quality and Javadoc comment-code inconsistencies. It checks only simple issues, e.g., whether the parameter names, return types, and exceptions mentioned in the @param, @return, and @throws tags are consistent with the actual parameter names, return types, and exceptions in the method. Doc Check [12] detects Javadoc errors such as missing and incomplete Javadoc tags. Different from checking for these *style* inconsistencies, @TCOMMENT detects *semantic* comment-code

inconsistencies related to null values and exceptions.

**Empirical Studies of Comments.** Several empirical studies aim to understand the conventional usage of comments, the evolution of comments, and the challenges of automatically understanding comments [25, 29, 44, 45]. None of them automatically analyze comments to detect comment-code inconsistencies or improve automated testing.

**Comment Inference from Source Code.** Several recent projects infer comments for failed test cases [47], exceptions [16], API function cross-references [28], software changes [17], and semantically related code segments [37, 38]. Comments automatically generated by these techniques are more structured than developer-written comments; therefore, it may be easier to leverage such automatically-generated comments for finding inconsistencies. However, it is still beneficial to improve the analysis of developer-written comments because (1) millions of lines of developer-written comments are available in modern software; and (2) these developer-written comments bring in information that is not available in source code [40] (which is also not available in comments inferred from the source code) to help us detect more faults. FailureDoc [47] augments a failed test with debugging clues, which could be extended to help explain why the tests generated by @TCOMMENT fail in order to help developers confirm/fix the faults.

**Analysis of Natural-Language Text for Software.** Various research projects analyze natural-language artifacts such as bug reports [14, 22, 27, 31, 36, 39, 42, 43], API documentation [48], and method names [21] for different purposes such as detecting duplicate bug reports or identifying the appropriate developers to fix bugs. @TCOMMENT analyzes comments written in a natural language to detect comment-code inconsistencies and to improve automated testing. Rubio-González et al. detect error code mismatches between code and manual pages in the Linux kernel by combining static analysis and heuristics [35]. Different from some of these studies [21, 48] that use natural-language processing (NLP) techniques such as part-of-speech tagging and chunking, @TCOMMENT does not use NLP techniques

because our simple comment analysis can already achieve a high accuracy of 97–100%, partially due to the more structured Javadoc comments with less paraphrases and variants (Section 5.4).

# Chapter 7

## Conclusions and Future Work

An inconsistency between comment and code is highly indicative of program faults. We have presented a novel approach, called @TCOMMENT, for testing consistency of Java method bodies and Javadoc comments properties related to null values and exceptions. Our application of @TCOMMENT on seven open-source projects discovered 28 methods with inconsistencies between Javadoc comments and bodies. We reported all these inconsistencies, and 12 were already confirmed and fixed by the developers.

In the future, @Randoop can be extended to (1) modify the random selection performed by Randoop such that it biases the selection based on the properties inferred by @TCOMMENT; (2) identify some causes of exceptions to reduce the rate of false alarms; and (3) rank the reported inconsistencies. @TCOMMENT can be extended to handle other types of properties and to be integrated with other testing or static analysis tools.



# References

- [1] *Apache Commons Collections*, <http://commons.apache.org/collections/>.
- [2] *Apache Commons Collections Bug Report 384*, <https://issues.apache.org/jira/browse/COLLECTIONS-384>.
- [3] *Apache Commons Collections Bug Report 385*, <https://issues.apache.org/jira/browse/COLLECTIONS-385>.
- [4] *Apache Log4j*, <http://logging.apache.org/log4j/>.
- [5] *Apache Lucene*, <http://lucene.apache.org/>.
- [6] *Apache Xalan*, <http://xml.apache.org/xalan-j/>.
- [7] *Glazed Lists*, <http://www.glazedlists.com/>.
- [8] *Javadoc Tool*, <http://java.sun.com/j2se/javadoc>.
- [9] *JFreeChart*, <http://www.jfree.org/jfreechart/>.
- [10] *Joda Time*, <http://joda-time.sourceforge.net/>.
- [11] *Joda Time Bug Report*, [http://sourceforge.net/tracker/?func=detail&atid=617889&aid=3413869&group\\_id=97367](http://sourceforge.net/tracker/?func=detail&atid=617889&aid=3413869&group_id=97367).
- [12] *Sun Doc Check Doclet*, <http://www.oracle.com/technetwork/java/javase/documentation/index-141437.html>.
- [13] *The Standard Doclet*, <http://download.oracle.com/javase/1,5.0/docs/guide/javadoc/standard-doclet.html>.
- [14] John Anvik, Lyndon Hiew, and Gail C. Murphy, *Who should fix this bug?*, ICSE, 2006.
- [15] Chandrasekhar Boyapati, Sarfraz Khurshid, and Darko Marinov, *Korat: Automated testing based on java predicates*, SIGSOFT Softw. Eng. Notes **27** (2002).
- [16] Raymond P.L. Buse and Westley R. Weimer, *Automatic documentation inference for exceptions*, ISSTA, 2008.

- [17] Raymond P.L. Buse and Westley R. Weimer, *Automatically documenting program changes*, ASE, 2010.
- [18] Maciej Cielecki, Jędrzej Fulara, Krzysztof Jakubczyk, and Lukasz Jancewicz, *Propagation of JML non-null annotations in Java programs*, PPPJ, 2006.
- [19] Ilinca Ciupa, Andreas Leitner, Manuel Oriol, and Bertrand Meyer, *ARTOO: Adaptive random testing for object-oriented software*, ICSE, 2008.
- [20] Valentin Dallmeier, Nikolai Knopp, Christoph Mallon, Sebastian Hack, and Andreas Zeller, *Generating test cases for specification mining*, ISSTA'10.
- [21] Zachary P. Fry, David Shepherd, Emily Hill, Lori Pollock, and K. Vijay-Shanker, *Analysing source code: Looking for useful verb-direct object pairs in all the right places*, IET Software Special Issue on Natural Language in Software Development (2008).
- [22] Michael Gegick, Pete Rotella, and Tao Xie, *Identifying security bug reports via text mining: An industrial case study*, MSR, 2010.
- [23] David Hovemeyer and William Pugh, *Finding more null pointer bugs, but not too many*, PASTE, 2007.
- [24] Hojun Jaygarl, Sunghun Kim, Tao Xie, and Carl K. Chang, *OCAT: Object capture-based automated testing*, ISSTA, 2010.
- [25] Zhen Ming Jiang and Ahmed E. Hassan, *Examining the evolution of code comments in PostgreSQL*, MSR, 2006.
- [26] Ninus Khamis, René Witte, and Juergen Rilling, *Automatic quality assessment of source code comments: the JavadocMiner*, NLDB, 2010.
- [27] Zhenmin Li, Lin Tan, Xuanhui Wang, Shan Lu, Yuanyuan Zhou, and Chengxiang Zhai, *Have things changed now? – An empirical study of bug characteristics in modern open source software*, ASID, 2006.
- [28] Fan Long, Xi Wang, and Yang Cai, *API hyperlinking via structural overlap*, ESEC/FSE, 2009.
- [29] Haroon Malik, Istehad Chowdhury, Hsiao-Ming Tsou, Zhen Ming Jiang, and Ahmed E. Hassan, *Understanding the rationale for updating a function's comment*, ICSM, 2008.
- [30] Christopher D. Manning and Hinrich Schütze, *Foundations of statistical natural language processing*, The MIT Press, 2001.
- [31] Dominique Matter, Adrian Kuhn, and Oscar Nierstrasz, *Assigning bug reports using a vocabulary-based expertise model of developers*, MSR, 2009.
- [32] Carlos Pacheco and Michael D. Ernst, *Eclat: Automatic Generation and Classification of Test Inputs*, ECOOP, 2005.

- [33] Carlos Pacheco, Shuvendu K. Lahiri, and Thomas Ball, *Finding errors in .net with feedback-directed random testing*, ISSTA, 2008.
- [34] Carlos Pacheco, Shuvendu K. Lahiri, Michael D. Ernst, and Thomas Ball, *Feedback-Directed Random Test Generation*, ICSE, 2007.
- [35] Cindy Rubio-González and Ben Liblit, *Expect the unexpected: error code mismatches between documentation and the real world*, PASTE, 2010.
- [36] Per Runeson, Magnus Alexandersson, and Oskar Nyholm, *Detection of duplicate defect reports using natural language processing*, ICSE, 2007.
- [37] Giriprasad Sridhara, Emily Hill, Divya Muppaneni, Lori Pollock, and K. Vijay-Shanker, *Towards automatically generating summary comments for Java methods*, ASE, 2010.
- [38] Giriprasad Sridhara, Lori Pollock, and K. Vijay-Shanker, *Automatically detecting and describing high level actions within methods*, ICSE, 2011.
- [39] Chengnian Sun, David Lo, Xiaoyin Wang, Jing Jiang, and Siau-Cheng Khoo, *A discriminative model approach for accurate duplicate bug report retrieval*, ICSE, 2010.
- [40] Lin Tan, Ding Yuan, Gopal Krishna, and Yuanyuan Zhou, */\* iComment: Bugs or bad comments? \*/*, SOSP, 2007.
- [41] Lin Tan, Yuanyuan Zhou, and Yoann Padioleau, *aComment: Mining annotations from comments and code to detect interrupt-related concurrency bugs*, ICSE, 2011.
- [42] Xiaoyin Wang, Lu Zhang, Tao Xie, John Anvik, and Jiasu Sun, *An approach to detecting duplicate bug reports using natural language and execution information*, ICSE, 2008.
- [43] Jin woo Park, Mu woong Lee, Jinhan Kim, Seung won Hwang, and Sunghun Kim, *CosTriage: A Cost-Aware Triage Algorithm for Bug Reporting Systems*, AAAI, 2011.
- [44] S. N. Woodfield, H. E. Dunsmore, and V. Y. Shen, *The effect of modularization and comments on program comprehension*, ICSE, 1981.
- [45] Annie T. T. Ying, James L. Wright, and Steven Abrams, *Source code that talks: An exploration of Eclipse task comments and their implication to repository mining*, MSR, 2005.
- [46] Sai Zhang, David Saff, Yingyi Bu, and Michael D. Ernst, *Combined static and dynamic automated test generation*, ISSTA, 2011.
- [47] Sai Zhang, Cheng Zhang, and Michael D. Ernst, *Automated documentation inference to explain failed tests*, ASE, 2011.
- [48] Hao Zhong, Lu Zhang, Tao Xie, and Hong Mei, *Inferring resource specifications from natural language API documentation*, ASE, 2009.