



HOW TO READ PAPERS



Prerequisite for researchers:
You really should like to read.

I'm not talking about reading
papers. I mean just reading.
Anything.



I'd rather have a book, but in a pinch I will settle for a set of Water Pik instructions. Once, long ago, I bested a desperate bout of insomnia by studying the only piece of written material in my apartment that I had not already read at least twice: my roommate's 1974 Toyota Corolla manual.

- Ex Libris: Confessions of a Common Reader
("서재 결혼시키기"), Anne Faldman

ANATOMY OF RESEARCH PAPERS

- X Title
- X Authors
- X Abstract
- X Introduction
- X Backgrounds (or Related Work)
- X Technique X
- X Research Questions
- X Experimental Setup
- X Results
- X Discussions / Future Work
- X Threats to Validity
- X Related Work (or Backgrounds)
- X Conclusion
- X Acknowledgement
- X References
- X Appendix
- X Online Supplements

READING VS. WRITING

- X Much of what we discuss today also work as advice for writing (i.e., write what where)
- X One helpful meta-reading is to focus on writing skills rather than contents – hard, but worthy skill to learn
 - Make note of words, phrases, even context-specific references



TITLE

“A rose by any other name would smell as sweet”

– Romeo & Juliet, Act 2 Scene 1

But really? 😊



TITLE

- X What the paper is about, obviously
- X The “granularity”
 - “On [something really big]”
 - “Towards [something really difficult]”
 - ”Evaluating the impact of X-zation of Y-ified Z on [something widely studied]”



AUTHORS

“If my doctor told me I had only six minutes to live, I would not brood. I’d type a little faster.”

– Isaac Asimov



AUTHORS

- X Knowing individual authors do help you reading the paper
 - Research trajectory
 - Area of expertise
 - Expected quality (ideally)
 - Social landscape of a field

- X Remembering the mapping between names and papers will help you when you actually run into them at conferences 😊



ABSTRACT

“Yes.”

– J. K. Gardner & L. Knopoff



ABSTRACT

- X A short summary of the entire paper
- X Usually 250–300 words
- X Some venues require structured abstract
 - Introduction / Methods / Results / Conclusion (IMRaD)
- X You SHOULD get some idea about what the paper is really about

JOURNAL ARTICLE

Is the sequence of earthquakes in Southern California, with aftershocks removed, Poissonian? 🛒

J. K. Gardner; L. Knopoff

Bulletin of the Seismological Society of America (1974) 64 (5): 1363-1367.

Published: October 01, 1974 Article history ▼

“ Cite 🔗 Share ▼ 🛠 Tools ▼

Abstract

Yes.



INTRODUCTION

“You never get a second chance to make the first impression.”
– anonymous



INTRODUCTION

- X A slightly longer summary of the entire paper (abstract expanded)
- X Ordering can be different, but usually contains the following parts:
 - Description of the world/status-quo
 - What the problem is
 - Why the problem is important (i.e., motivation)
 - What other people have tried before
 - Why they fall short
 - Description of the current paper
 - What “we” do for the problem, differently from others (i.e., justification)
 - A brief summary of technical contributions
 - Organisation of the remaining sections

Diversity-Aware Mutation Adequacy Criterion for Improving Fault Detection Capability

Donghwan Shin
School of Computing
KAIST
Daejeon, Republic of Korea
Email: donghwan@se.kaist.ac.kr

Shin Yoo
School of Computing
KAIST
Daejeon, Republic of Korea
Email: shin.yoo@kaist.ac.kr

Doo-Hwan Bae
School of Computing
KAIST
Daejeon, Republic of Korea
Email: bae@se.kaist.ac.kr

Abstract—Many existing testing techniques adopt diversity as an important criterion for the selection and prioritization of tests. However, mutation adequacy has been content with simply maximizing the number of mutants that have been killed. We propose a novel mutation adequacy criterion that considers the diversity in the relationship between tests and mutants, as well as whether mutants are killed. Intuitively, the proposed criterion is based on the notion that mutants can be distinguished by the sets of tests that kill them. A test suite is deemed adequate by our criterion if the test suite distinguishes all mutants in terms of their kill patterns. Our hypothesis is that, simply by using a stronger adequacy criterion, it is possible to improve fault detection capabilities of mutation-adequate test suites. The empirical evaluation selects tests for real world applications using the proposed mutation adequacy criterion to test our hypothesis. The results show that, for real world faults, test suites adequate to our criterion can increase the fault detection success rate by up to 76.8 percentage points compared to test suites adequate to the traditional criterion.

I. INTRODUCTION

One fundamental limitation of software testing is the fact that, to validate the behavior of the Program Under Test (PUT), we can only ever sample a very small number of test inputs out of the vast input space. Almost all existing testing techniques, at some level, attempt to answer the following question: *How does one sample a finite number of test inputs to cover a wide a range of program behaviors as possible?*

The concept of diversity has received much attention while answering the above question. For example, Adaptive Random testing (ART) [1] seeks to create diversity of randomly sampled test inputs by choosing an input that is as different from those already sampled as possible. Clustering-based selection and prioritization [2], [3] assumes that a diverse set of test inputs would explore and validate a wider range of program behaviors. Diversity in test input has been studied as a test adequacy criterion for black box testing of web applications [4]. Information theoretic measures of diversity are also been studied as a test selection criterion [5], [6].

In contrast, diversity has received little attention in mutation testing. Mutation adequacy remains essentially as a simple count of the number of mutants killed. Many of existing works focus either on reducing the cost of mutation testing (i.e., do fewer, do smarter, and do faster as summarized by Bluff and Utch [7]) or analyzing equivalent mutants [8], [9].

Intuitively, mutants semantically equivalent to an original program relatively very few attention has been paid to improve its fault detection capability of the mutation adequacy criterion.

The existing mutation adequacy as a count of killed mutants does not cater for diversity. Consequently, despite its potential correlation to the fault detection capability, many diverse mutants are generated but not killed. Suppose a pathological case in which a single test can kill all generated mutants. In terms of the kill information, it means that the single test does not capture the diversity of the mutants enough, while the traditional mutation adequacy simply determines the single test as adequate. Such a case calls for a richer adequacy criterion.

This paper proposes a novel adequacy criterion called distinguishing mutation adequacy criterion, which include the notion of diversity. The proposed metric is based on previous work on theoretical framework for mutation testing [10]. At the core of the new criterion lies the idea that mutants can be distinguished from each other by the set of tests that kill them. Our mutation adequacy criterion aims not only to kill, but also to distinguish as many mutants as possible. The aforementioned pathological case of a single test killing all mutants will perform poorly under our new criterion.

By aiming to distinguish the maximum number of mutant the proposed adequacy criterion can select more diverse set of test cases. Suppose there exist two mutants. Test t_1 kill both, while t_2 and t_3 kill different one each. Under the existing criterion, the set $\{t_1\}$ is deemed adequate, whereas the proposed criterion will choose the set $\{t_2, t_3\}$ instead. We hypothesize that this, more diverse set of tests will show high-fault detection capability.

The hypothesis on fault detection capability is validated by an empirical study. We use the Defect4J [11] data to study real world faults in non-critical systems. The control group consists of sets of test suites selected based on the traditional mutation adequacy (i.e., ones that kill all mutants) while the treatment group consists of sets of test suites that can collectively kill and distinguish all mutants. Both groups are evaluated for their fault detection capabilities by executing them against faulty and fixed versions of programs collected from the real world. The results show that our novel mutation

adequacy criterion shows either equal or higher fault detection capabilities for all studied subject faults.

The technical contributions of this paper are as follows:

- This paper introduces a novel mutation adequacy criterion called distinguishing mutation criterion. A test suite is mutation adequate with respect to this criterion if all considered mutants have unique sets of tests that kill them, i.e. can be distinguished by their kill patterns.
- The proposed adequacy criterion is empirically evaluate using real faults in the Defect4J repository using random testing. The results show that the new adequacy criterion shows at least equal or higher fault detection capability than the traditional mutation adequacy criterion. The increase in the fault detection success rate is up to 76.8 percentage points.

The rest of the paper is organized as follows. Section I presents a formal definition of the existing mutation adequacy criterion. Section III introduces the distinguishing mutation adequacy criterion using the new formal notations. Section IV describes the design of our empirical evaluation, the results of which are presented and analysed in Section V. Section VI discusses related work, and finally, Section VII concludes.

II. BACKGROUND

A. Formal Model of Mutation Testing

To formally represent mutation adequacy criteria considered in this paper, we summarize the essential elements of the formal framework for the mutation-based testing methods. Detailed descriptions for the formal framework are presented in [10].

Let P be a set of programs which includes the program under test. In mutation testing, there are three essential programs in P : an original program $p_o \in P$, a mutant $m \in M \subseteq P$ generated from p_o , and a correct program $p_c \in P$ which represents the true requirements¹ about p_o . For a test $t \in T$ for P , if the behaviors of p_o and p_c are different, it is said that t detects a fault in p_o . Similarly, if the behaviors of p_o and m are different for t , it is said that t kills m . Note that the notion of behavioral difference is an abstract concept. It is formalized by a testing factor, called a test differentiator, which is defined as follows:

Definition 1: A test differentiator $d: T \times P \times P \rightarrow \{0, 1\}$ is a function,² such that

$$d(t, p_o, p_c) = \begin{cases} 1 & (\text{true}), & \text{if } p_o \text{ is different with } p_c \text{ for } t \\ 0 & (\text{false}), & \text{otherwise} \end{cases}$$

for all tests $t \in T$ and programs $p_o, p_c \in P$.

By definition, a test differentiator concisely represents whether the behaviors of $p_o \in P$ and $p_c \in P$ are different for t .

¹While p_c is not a real program, this is not a serious assumption, because we only require the behavior of p_c for a given set of tests. In practice, a human may play the role of p_c , acting as a hidden oracle.

²This function-style definition is replaceable by a predicate-style definition, such as $d \subseteq T \times P \times P$.

We make no attempt to incorporate any specific definition of program differences. The specific definition of differences can only be decided in context. For example, while 0.3333 is different with 1/3 in the strict sense, 0.3333 will be regarded as the same as 1/3 in some cases. To keep things general, we consider a set of test differentiators D that includes all possible test differentiators for P .

A test differentiator, or simply differentiator, can formally describe the notion of differences in mutation testing. For example, when t detects a fault in p_o , it is clearly formalized as follows³:

$$d(t, p_o, p_o) = 1$$

On the other hand, when t kills a mutant m , it is also clearly formalized as follows:

$$d(t, p_o, m) = 1$$

Note that p_o , p_c , and m are general entities, and largely separated from any specifics such as programming languages or mutation methods.

B. Mutation Adequacy Criterion

Since mutation testing was first proposed in the 1970s, it has been widely studied in the aspects of both theory and practice, and a mutation adequacy criterion has played the key role in the studies of mutation testing. A mutation adequacy criterion is a predicate that determines the adequacy of a test suite using mutants. It is said that a test suite is mutation-adequate when the test suite kills all of the generated mutants. Using a differentiator, it is clearly and concisely formalized as follows:

$$\forall m \in M, \exists t \in T, S, d(t, p_o, m) = 1. \quad (1)$$

In other words, a test suite T is mutation-adequacy if all mutants $m \in M$ are killed by at least one test $t \in T$.

Equation (1) is general enough to consider various mutation testing approaches. For example, there is a spectrum of mutation approaches from a strong mutation [12] to a weak mutation [13], depending on which d is used. In a strong mutation analysis, a test t kills a mutant m when the output of m differs from the output of the original program p_o for t . In a weak mutation analysis, t kills m when the internal states of m and p_o are different for t . In the rest of this paper, we refer (1) as the traditional mutation adequacy criterion in comparison to the new mutation adequacy proposed in Section III-B.

III. EXTENDING MUTATION ADEQUACY CONSIDERING DIVERSITY OF MUTANTS

A. Limitation of Traditional Mutation Adequacy

To use the limitation of the traditional mutation adequacy criterion, we provide a working example with four mutants and three tests in Figure 1. Each of the values represents whether a test kills a mutant. For example, $d(t_1, p_o, m_1) = 1$ which means that t_1 kills m_1 .

³In experiments, when the correct version of a program for a fault is known in advance, the correct version can be used as p_c . In this case, the corresponding faulty version should be used as p_o , so that the difference between p_o and p_c implies the fault.

PART 1: STATE OF THE WORLD

One fundamental limitation of software testing is the fact that, to validate the behavior of the Program Under Test (PUT), we can only ever sample a very small number of test inputs out of the vast input space. Almost all existing testing techniques are, at some level, attempts to answer the following question: *how does one sample a finite number of test inputs to cover as wide a range of program behaviours as possible?*

The concept of diversity has received much attention while answering the above question. For example, Adaptive Random Testing (ART) [1] seeks to increase diversity of randomly sampled test inputs by choosing an input that is as different from those already sampled as possible. Clustering-based test selection and prioritization [2], [3] assumes that a diverse set of test inputs would explore and validate a wider range of program behaviors. Diversity in test output has been studied as a test adequacy criterion for black box testing of web applications [4]. Information theoretic measures of diversity has also been studied as a test selection criterion [5], [6].

- X State of the fact
- X What has been happening in other field

PART 2: DESCRIPTION OF THE PROBLEM

In contrast, diversity has received little attention in relation to mutation testing. Mutation adequacy remains essentially as a simple count of the number of killed mutants. Many of existing work (i.e., mutants semantically equivalent to an original program). (i.e. Relatively very few attention has been paid to improve the Off fault detection capability of the mutation adequacy criterion itself.

The existing mutation adequacy as a count of killed mutants does not cater for diversity. Consequently, despite its potential correlation to the fault detection capability, many diverse mutants are generated but wasted. Suppose a pathological case in which a single test can kill all generated mutants. In terms of the kill information, it means that the single test does not capture the diversity of the mutants enough, while the traditional mutation adequacy simply determines the single test as adequate. Such a case calls for a richer adequacy criterion in mutation testing.

- X Compared to what has been going on in other fields, ours hasn't put much effort in X
- X Existing Y fails to do X because...

PART 3: WHAT WE PROPOSE

This paper proposes a novel adequacy criterion called distinguishing mutation adequacy criterion, which includes the notion of diversity. The proposed metric is based on our previous work on theoretical framework for mutation testing [10]. At the core of the new criterion lies the idea that mutants can be *distinguished* from each other by the set of tests that kill them. Our mutation adequacy criterion aims not only to kill, but also to distinguish as many mutants as possible. The aforementioned pathological case of a single test killing all mutants will perform poorly under our new criterion.

By aiming to distinguish the maximum number of mutants, the proposed adequacy criterion can select more diverse set of test cases. Suppose there exist two mutants. Test t_1 kills both, while t_2 and t_3 kill different one each. Under the existing criterion, the set $\{t_1\}$ is deemed adequate, whereas the proposed criterion will choose the set $\{t_2, t_3\}$ instead. We hypothesize that this, more diverse set of tests will show higher fault detection capability.

- X High level description of what this paper proposes
- X Try to imagine what the low-level, technical details would be

PART 4: HOW WE ARE GOING TO EVALUATE

The hypothesis on fault detection capability is validated by an empirical study. We use the `Defects4J` [11] data set to study real world faults in non-trivial systems. The control group consists of sets of test suites selected based on the traditional mutation adequacy (i.e., ones that kill all mutants), while the treatment group consists of sets of test suites that can collectively kill and distinguish all mutants. Both groups are evaluated for their fault detection capabilities by executing them against faulty and fixed versions of programs collected from the real world. The results show that our novel mutation

- X Peek into the actual evaluation
- X Allows you to set your expectation – what evidence will I see?

PART 5 & 6: CONTRIBUTIONS AND ORGANISATION

The technical contributions of this paper are as follows:

- This paper introduces a novel mutation adequacy criterion called *distinguishing mutation criterion*. A test suite is mutation adequate with respect to this criterion if all considered mutants have unique sets of tests that kill them, i.e. can be *distinguished* by their kill patterns.
- The proposed adequacy criterion is empirically evaluated using real faults in the Defects4J repository and random testing. The results show that the new adequacy

- X Point-by-point, declarative summary of what the paper actually contributes
- X A brief ToC...

The rest of the paper is organised as follows. Section II presents a formal definition of the existing mutation adequacy criterion. Section III introduces the distinguishing mutants adequacy criterion using the same formal notations. Section IV describes the design of our empirical evaluation, the results of which are presented and analysed in Section V. Section VI discusses related work, and finally Section VII concludes.

5.

BACKGROUND / RELATED WORK

“Knowledge is power.”

– Francis Bacon

(Bonus: definitely read this AFTER CLASS–

https://www.reddit.com/r/AskReddit/comments/dxosj/what_word_or_phrase_did_you_totally_misunderstand/c13pbyc/)



BACKGROUND / RELATED WORK

- X Where you list existing work that are related to yours
- X Depending on context, it may be better to:
 - Explain existing material before describing yours: have a “Background” section
 - Describe your technique first, and compare it to existing work before: have a “Related Work” section
 - Can have both in rare cases
- X Summarise existing work briefly
- X ***Position your work***
 - highlight differences between them and yours



RESEARCH QUESTIONS / EXPERIMENTAL SETUP

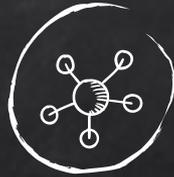
“Judge a man by his questions rather than his answers.”

– Voltaire



RESEARCH QUESTIONS

- X The “solution” in a research question is usually presented in the form of question and answer.
- X You can learn how the proposed solution is being evaluated by reading RQs
 - Effectiveness: how good is the new technique?
 - Efficiency: how cheap/fast is the new technique?



EXPERIMENTAL SETUP

- X For domains that involves empirical evaluation (i.e., your “solution” is a “technique” therefore you apply it to a set of subjects and evaluate how well it works), you need to specify:
 - What your subjects are, and why/how they were chosen
 - How your technique is implemented, using which external dependencies
 - In which environment the experiments have been executed
 - What the parameters to your algorithm are
- X This helps both replication and understanding of results



THREATS TO VALIDITY

“The validity of internet quotes are getting sketchy these days.”
– Abraham Lincoln



THREATS TO VALIDITY

- X This is where you honestly discuss any factors that may threaten the validity of your claims
 - Internal validity: have you evaluated your solution correctly?
 - E.g., implementational correctness
 - External validity: how well does your claim generalise?
 - E.g., any potential bias in the selection of subjects
 - Construct validity: are you measuring what you are supposed to measure?
 - E.g., the use of correct evaluation metrics
- X Use extra scepticism

READING STRATEGY

IT IS OKAY TO FEEL OVERWHELMED

- X Reading research papers is difficult. By definition.
- X Keep in mind *why* you are reading this paper (when applicable, ask the person who told you to read it)
 - Comfort zone vs. extension of your knowledge
 - Theoretical vs. empirical

MULTI-PASS READING

- X Almost all guidelines recommend multi-pass reading
 - Skim: focus on title, abstract, figures, and conclusion – identify what it is about
 - Read: read again, using the common structure as your map – form your own questions
 - Interpret: deep-dive into the text, tables, and figures – find answers to your own questions
 - Summarise: add the paper to your own map of research landscape

CRITICAL/CREATIVE READING

X Read Critically

- Do not be afraid to question the contents – be sceptical and suspicious
- Is the problem right? Is the argument justifiable? Are the data not biased?

X Read Creatively

- Criticising something is easy: building up is harder
- Can this be done for other problems? What would I do differently, or subsequently?

CRITICAL/CREATIVE READING

- X Read Between/Beyond Lines
 - Notice what is NOT being said
 - Learn to distinguish nuances
- X Reading as a group
 - Will bring out points that you missed
 - Discussion is not a fight to win or lose
 - Reading critically does not mean you have to judge the paper

"IT HAS LONG BEEN KNOWN"	I didn't look up the original reference.
"A DEFINITE TREND IS EVIDENT"	The data are practically meaningless.
"WHILE IT HAS NOT BEEN POSSIBLE TO PROVIDE DEFINITE ANSWERS TO THE QUESTIONS"	An unsuccessful experiment, but I still hope to get it published.
"THREE OF THE SAMPLES WERE CHOSEN FOR DETAILED STUDY"	The other results didn't make any sense.
"TYPICAL RESULTS ARE SHOWN"	This is the prettiest graph.
"THESE RESULTS WILL BE IN A SUBSEQUENT REPORT"	I might get around to this sometime, if published/funded.
"A CAREFUL ANALYSIS OF OBTAINED DATA"	Three pages of notes were obliterated when I knocked over a glass of beer.
"AFTER ADDITIONAL STUDY BY MY COLLEAGUES"	They didn't understand it, either.
"THANKS ARE DUE TO JOE BLOTZ FOR ASSISTANCE WITH THE EXPERIMENT AND TO CINDY ADAMS FOR VALUABLE DISCUSSIONS"	Mr. Blotz did the work and Ms. Adams explained to me what it meant.
"A HIGHLY SIGNIFICANT AREA FOR EXPLORATORY STUDY"	A totally useless topic selected by my committee.
"IN MY EXPERIENCE"	Once.
"IN CASE AFTER CASE"	Twice.
"IN A SERIES OF CASES"	Three times.
"IT IS BELIEVED THAT"	I think.
"IT IS GENERALLY BELIEVED THAT"	A couple of others think so, too.
"CORRECT WITHIN AN ORDER OF MAGNITUDE"	Wrong.
"ACCORDING TO STATISTICAL ANALYSIS"	Flavor has it.
"IT IS CLEAR THAT MUCH ADDITIONAL WORK WILL BE REQUIRED BEFORE A COMPLETE UNDERSTANDING OF THIS PHENOMENON OCCURS"	I don't understand.
"A STATISTICALLY-ORIENTED PROJECTION OF THE SIGNIFICANCE OF THESE FINDINGS"	A wild guess.
"IT IS HOPED THAT THIS STUDY WILL STIMULATE FURTHER INVESTIGATIONS IN THIS FIELD"	I quit.

BOTH USE AND IGNORE THE STRUCTURE

- X Different papers are written for different model readers in mind
 - Most frequently: for someone with general CS knowledge by little domain expertise
 - Depending on venues: for someone with reasonable domain knowledge on a specific topic (e.g., specialised conference proceedings)
- X Do not read everything sequentially. Jump depending on:
 - Your domain knowledge
 - Your questions

TREAT CITATIONS AS HYPertexts

- X Be prepared to follow references and read other papers
- X People *do* make habitual citations (I am guilty as well), sometimes completely out of context – do not rely on authority, find out on your own

REPRODUCTION

- X Increasingly researchers are required to provide a replication package with a paper (code, data, and anything else required to replicate the study)
- X Make a note of whether such a package is explicitly included in the paper – it will save you time later

READ TO WRITE BETTER

- X Try to expand your vocabulary all the time: sometimes you want the perfect choice of word to get the right nuance
- X Pick up phrases and expressions that are useful
- X Identify fancy LaTeX tricks; maintain a snippet library

NOTE-MAKING

- X Write a short summary of what you read, if possible
- X What the paper is about, which data has been used, what you think are strengths and weaknesses

SUMMARIES

- X There are common structures papers follow – know which contains which
- X Do multi-pass reading, gradually digging into finer details
- X Read both critically and creatively
- X Maintain a note of what you read

MORE READING ABOUT READING

- X How to read a paper – S. Keshav, University of Waterloo ([link](#))
- X How to read a research paper – M. Mitzenmacher, Harvard University ([link](#))
- X How to read a research paper – Grisword, Murphy, and Conati ([link](#))
- X How to seriously read a scientific paper – Science Magazine ([link](#))
 - o How to read a scientific paper ([link](#))