

Putting Software Testing Terminology to the Test

Samuel J. Crawford*, Spencer Smith*, Jacques Carette*

*Department of Computing and Software

McMaster University

Hamilton, Canada

{crawfs1, smiths, carette}@mcmaster.ca

Abstract—Despite the prevalence and importance of software testing, it lacks a standardized and consistent taxonomy. This hinders precise communication, leading to discrepancies across the literature (even within individual documents!) and to potential misunderstandings when planning and performing testing. In this paper, we systematically explore the current state of software testing terminology. We 1) identify established standards and prominent testing resources, 2) capture relevant testing terms from these sources, along with their definitions and relationships—both explicit and implicit—and 3) construct graphs to visualize and analyze these data. This process uncovers 557 test approaches and 76 software qualities that may imply additional related test approaches. We also build a tool for generating graphs that illustrate relations between test approaches and track discrepancies captured by this tool and manually through the research process. This reveals 259 discrepancies, including nine terms used as synonyms to two (or more) disjoint test approaches and 15 pairs of test approaches that may either be synonyms or have a parent-child relationship. This also highlights notable confusion surrounding functional, operational acceptance, recovery, and scalability testing. Our findings make clear the urgent need for improved testing terminology so that the discussion, analysis and implementation of various test approaches can be more coherent. We provide some preliminary advice on how to achieve this standardization.

Index Terms—Software testing, terminology, taxonomy, literature review, test approaches

I. Introduction

As with all fields of science and technology, software development should be approached systematically and rigorously. Reference [1] claims that “to be successful, development of software systems requires an engineering approach” that is “characterized by a practical, orderly, and measured development of software” [p. 3]. When a NATO study group decided to hold a conference to discuss “the problems of software” in 1968, they chose the phrase “software engineering” to “imply[] the need for software manufacture to be based on the types of theoretical foundations and practical disciplines, [sic] that are traditional in the established branches of engineering” [2, p. 13]. “The term was not in general use at that time”, but conferences such as this “played a major role in gaining general acceptance ... for the term” [3]. While one of the goals of the conference was to “discuss possible techniques, methods and developments which might lead to the[]

solution” to these problems [2, p. 14], the format of the conference itself was difficult to document. Two competing classifications of the report emerged: “one following from the normal sequence of steps in the development of a software product” and “the other related to aspects like communication, documentation, management, [etc.]” [p. 10]. Perhaps more surprisingly, “to retain the spirit and liveliness of the conference, ... points of major disagreement have been left wide open, and ... no attempt ... [was] made to arrive at a consensus or majority view” [p. 11]!

Perhaps unsurprisingly, there are still concepts in software engineering without consensus, and many of them can be found in the subdomain of software testing. Reference [4] gives the example of complete testing, which may require the tester to discover “every bug in the product”, exhaust the time allocated to the testing phase, or simply implement every test previously agreed upon [p. 7]. Having a clear definition of “complete testing” would reduce the chance for miscommunication and, ultimately, the tester getting “blamed for not doing ... [their] job” [p. 7]. Because software testing uses “a substantial percentage of a software development budget (in the range of 30 to 50%)”, which is increasingly true “with the growing complexity of software systems” [1, p. 438], this is crucial to the efficiency of software development. Even more foundationally, if software engineering holds code to high standards of clarity, consistency, and robustness, the same should apply to its supporting literature!

Unfortunately, a search for a systematic, rigorous, and complete taxonomy for software testing revealed that the existing ones are inadequate and mostly focus on the high-level testing process rather than the testing approaches themselves:

- Tebes et al. [5] focus on parts of the testing process (e.g., test goal, test plan, testing role, testable entity) and how they relate to one another,
- Souza et al. [6] prioritize organizing test approaches over defining them,
- Firesmith [7] similarly defines relations between test approaches but not the approaches themselves, and
- Unterkalmsteiner et al. [8] focus on the “information linkage or transfer” [p. A:6] between requirements engineering and software testing and “do[] not aim at providing a systematic and exhaustive state-of-the-art survey of [either domain]” [p. A:2].

In addition to these taxonomies, many standards documents (see Section III-A1) and terminology collections (see Section III-A2) define testing terminology, albeit with their own issues.

For example, a common point of discussion in the field of software is the distinction between terms for when software does not work correctly. We find the following four to be most prevalent:

- Error: “a human action that produces an incorrect result” [9, p. 128], [10, p. 399].
- Fault: “an incorrect step, process, or data definition in a computer program” [9, p. 140] inserted when a developer makes an error [9, pp. 128, 140], [10, pp. 399–400], [11, p. 12-3].
- Failure: the inability of a system “to perform a required function or ... within previously specified limits” [9, p. 139], [12, p. 7] that is “externally visible” [12, p. 7] and caused by a fault [10, p. 400], [11, p. 12-3].
- Defect: “an imperfection or deficiency in a project component where that component does not meet its requirements or specifications and needs to be either repaired or replaced” [9, p. 96].

This distinction is sometimes important, but not always [13, p. 4-3]. The term “fault” is “overloaded with too many meanings, as engineers and others use the word to refer to all different types of anomalies” [11, p. 12-3], and “defect” may be used as a “generic term that can refer to either a fault (cause) or a failure (effect)” [9, p. 96], [14, p. 124]. Software testers may even choose to ignore these nuances completely! [15, pp. 13–14] “just call[s] it what it is and get[s] on with it”, abandoning these four terms, “problem”, “incident”, “anomaly”, “variance”, “inconsistency”, “feature” (!), and “a list of unmentionable terms” in favour of “bug”; after all, “there’s no reason to dice words”!

But why are minor differences between terms like these even important? Our previous list of terms “error”, “fault”, “failure”, and “defect” are used to describe many test approaches, including:

- 1) Defect-based testing
- 2) Error forcing
- 3) Error guessing
- 4) Error tolerance testing
- 5) Error-based testing
- 6) Error-oriented testing
- 7) Failure tolerance testing
- 8) Fault injection testing
- 9) Fault seeding
- 10) Fault sensitivity testing
- 11) Fault tolerance testing
- 12) Fault tree analysis
- 13) Fault-based testing

When considering which approaches to use or when actually using them, the meanings of these four terms

inform what their related approaches accomplish and how to they are performed. For example, the tester needs to know what a “fault” is to perform fault injection testing; otherwise, what would they inject? Information such as this is critical to the testing team, and should therefore be standardized.

Discrepancies such as these that can lead to miscommunications—such as that previously mentioned by [4, p. 7]—are prominent in the literature. ISO/IEC and IEEE categorize experience-based testing as both a test design technique and a test practice on the same page—twice [16, Fig. 2, p. 34]! The structure of tours can be defined as either quite general [16, p. 34] or “organized around a special focus” [17]. Load testing is performed with loads “between anticipated conditions of low, typical, and peak usage” [16, p. 5] or loads that are as large as possible [15, p. 86]. Alpha testing is performed by “users within the organization developing the software” [14, p. 17], “a small, selected group of potential users” [11, p. 5-8], or “roles outside the development organization” conducted “in the developer’s test environment” [17]. It is clear that there is a notable gap in the literature, one which we attempt to describe and fill. While the creation of a complete taxonomy is unreasonable, especially considering the pace at which the field of software changes, we can make progress towards this goal that others can extend and update as new test approaches emerge.

This document describes this process, as well as its results, in more detail. We first define the scope of what kinds of software testing are of interest (Section II) and examine the existing literature (Section III). Despite the amount of well understood and organized knowledge, there are still many discrepancies in the literature, either within the same source or between various sources (Section IV). This reinforces the need for a proper taxonomy! We provide some potential solutions covering some of these discrepancies (Section V).

II. Scope

Since our motivation is restricted to testing code, only this component of Verification and Validation (V&V) is considered. However, some test approaches are used for testing things other than code, and some approaches can be used for both! In these cases, only the subsections of these approaches focused on code are considered. For example, reliability testing and maintainability testing can start without code by “measur[ing] structural attributes of representations of the software” [18, p. 18], but only reliability and maintainability testing performed on code itself is in scope of this research. Therefore, some practices are excluded from consideration either in part or in full; hardware testing (Section II-A) and the V&V of other artifacts (Section II-B) are completely out of scope, as well as relevant areas of other testing approaches that are otherwise in scope. Static testing can be performed on code, so while it isn’t relevant to the original motivation of

this work, it is a useful component of software testing and is therefore included at this level of analysis (Section II-C).

A. Hardware Testing

While testing the software run on or in control of hardware is in scope, testing performed on the hardware itself is out of scope. The following are some examples of hardware testing approaches:

- Ergonomics testing and proximity-based testing (see [17]) are out of scope, since they are used for testing hardware.
- Similarly, EManations SECurity (EMSEC) testing [19], [20, p. 95], which deals with the “security risk” of “information leakage via electromagnetic emanation” [20, p. 95], is also out of scope.
- Orthogonal Array Testing (OAT) can be used when testing software [21] (in scope) but can also be used for hardware [22, pp. 471-472], such as “processors ... made from pre-built and pre-tested hardware components” [p. 471] (out of scope). A subset of OAT called “Taguchi’s Orthogonal Array Testing (TOAT)” is used for “experimental design problems in manufacturing” [23, p. 1573] or “product and manufacturing process design” [24, p. 44] and is thus also out of scope.

B. V&V of Other Artifacts

The only testing of a software artifact produced by the software life cycle that is in scope is testing of the software itself, as demonstrated by the following examples:

- Design reviews and documentation reviews are out of scope, as they focus on the V&V of design [14, pp. 132] and documentation [14, pp. 144], respectively.
- Error seeding is the “process of intentionally adding known faults¹ to those already in a computer program”, done to both “monitor[] the rate of detection and removal”, which is a part of V&V of the V&V itself (out of scope), “and estimat[e] the number of faults remaining” [14, p. 165], which helps verify the actual code (in scope).
- Fault injection testing, where “faults are artificially introduced¹ into the SUT [System Under Test]”, can be used to evaluate the effectiveness of a test suite [11, p. 5-18], which is a part of V&V of the V&V itself (out of scope), or “to test the robustness of the system in the event of internal and external failures” [16, p. 42], which helps verify the actual code (in scope).
- “Mutation [t]esting was originally conceived as a technique to evaluate test suites in which a mutant is

¹While error seeding and fault injection testing both introduce faults as part of testing, they do so with different goals: to “estimat[e] the number of faults remaining” [14, p. 165] and “test the robustness of the system” [16, p. 42], respectively. Therefore, these approaches are not considered synonyms, and the lack of this relation in the literature is not included in Section IV-B2 as a synonym discrepancy.

a slightly modified version of the SUT” [11, p. 5-15], which is in the realm of V&V of the V&V itself (out of scope). However, it “can also be categorized as a structure-based technique” and can be used to assist fuzz and metamorphic testing [11, p. 5-15] (in scope).

C. Static Testing

Sometimes, static testing is excluded from software testing [1, p. 439], [7, p. 13], [25, p. 222], restricting “testing” to mean dynamic validation [11, p. 5-1] or verification “in which a system or component is executed” [14, p. 427]. However, “terminology is not uniform among different communities, and some use the term ‘testing’ to refer to static techniques² as well” [11, p. 5-2]. This is done by [16, p. 17] and [26, pp. 8-9]; the authors of the former even explicitly exclude static testing in another document [14, p. 440]!

III. Methodology

At a high level, our methodology follows the following steps:

- 1) Identify authoritative sources (Section III-A)
- 2) Identify software testing terminology from each source, focusing on test approaches and software qualities
- 3) For each test approach, record its: (Section III-C)
 - a) Name
 - b) Category³ (Section III-B1)
 - c) Definition
 - d) Synonyms (Section III-B2)
 - e) Parents (Section III-B3)
 - f) Other relevant notes (e.g., prerequisites, uncertainties, and other sources to investigate)
- 4) Repeat steps 1 to 3 for any missing or unclear terminology (Section III-D)
- 5) Analyze these data for discrepancies
 - a) Record discrepancies as they arise during data collection
 - b) Generate relation graphs (Section III-E)
- 6) Report results of discrepancy analysis (Section IV)
- 7) Seek to resolve these discrepancies (Section V)

A. Sources

As there is no single authoritative source on software testing terminology, we need to look at many sources to observe how this terminology is used in practice. Since we are particularly interested in software engineering, we start from the vocabulary document for systems and software engineering [14] and two versions of the Guide to the SoftWare Engineering Body Of Knowledge (SWE-BOK Guide)—the newest one [13] and one submitted

²Not formally defined, but distinct from the notion of “test technique” described in Table I.

³There may be more than one category given for a single test approach which is indicative of a discrepancy (see Section IV-B1).



Fig. 1: Summary of how many sources comprise each source tier.

for public review⁴ [11]—as suggested by Dr. Carrette. To gather further sources, we then use a version of “snowball sampling”, which “is commonly used to locate hidden populations ... [via] referrals from initially sampled respondents to other persons” [27]. We apply this concept to “referrals” between sources. For example, [17] cites [28] as the original source for its definition of “scalability” (see Section V-B); we verified this by looking at this original source. We similarly “snowball” on terminology itself; when a term requires more investigation (e.g., its definition is missing or unclear), we perform a miniature literature review on this subset to “fill in” this missing information (see Section III-D). If these additional sources provide more information and are “trustworthy”, we may then investigate them in their entirety (as opposed to just the original subset of interest). We define a source to be “trustworthy” if it:

- 1) has gone through a peer-review process,
- 2) is written by numerous, well-respected authors,
- 3) is informed by many sources, and
- 4) is accepted and used in the field of software.

For ease of discussion and analysis, we group the complete set of sources into “tiers” based on their format, method of publication, and this metric of “trustworthiness”. We therefore create the following tiers, given in order of descending trustworthiness: established standards (Section III-A1), terminology collections (Section III-A2), textbooks (Section III-A3), and papers and other documents (Section III-A4). Note that most sources used to “fill in” missing information are papers. A summary of how many sources comprise each tier is given in Figure 1.

1) **Established Standards:** These are documents written for the field of software engineering by reputable standards bodies, namely ISO, the International Electrotechnical Commission (IEC), and IEEE. Their purpose is to “encourage the use of systems and software engineering standards” and “collect and standardize terminology” by “provid[ing] definitions that are rigorous, uncomplicated, and understandable by all concerned” [14, p. viii]. For these reasons, they are the most trustworthy sources. However, this does not imply perfection, as we identify

⁴Reference [11] has been published since we investigated these sources; if time permits, we will revisit this published version.

49 discrepancies within these standards (see Tables II and III)! Only standards for software development and testing are in scope for this research (see Section II). For example, “the purpose of the ISO/IEC/IEEE 29119 series is to define an internationally agreed set of standards for software testing that can be used by any organization when performing any form of software testing” [16, p. vii]. This tier is composed of [9], [12], [14], [16], [29]–[39].

2) **Terminology Collections:** These are collections of software testing terminology built up from multiple sources (such as the established standards outlined in Section III-A1) that are made to be widely applicable. For example, the SWEBOK Guide is “proposed as a suitable foundation for government licensing, for the regulation of software engineers, and for the development of university curricula in software engineering” [4, p. xix]. They are often written by a large organization, such as the International Software Testing Qualifications Board (ISTQB), but not always. We include Firesmith’s taxonomy [7] because it presents relations between many test approaches and Doğan et al.’s literature review [40] because it cites many of sources from which we can “snowball” if desired (see Section III-A). This tier is composed of [7], [11], [13], [17], [40].

3) **Textbooks:** We consider textbooks to be more trustworthy than papers (see Section III-A4) because they are widely used as resources for teaching software engineering and may be used as guides in industry. Although textbooks have smaller sets of authors, they follow a formal review process before publication. Textbooks used at McMaster University [1], [10], [15] served as the original (albeit ad hoc and arbitrary) starting point of this research, and we investigate other books as they arise. For example, [17] cites [28] as the original source for its definition of “scalability” (see Section V-B); we verified this by looking at this original source. This tier is composed of [1], [4], [10], [15], [28], [41], [42].

4) **Papers and Other Documents:** The remaining documents all have much smaller sets of authors and are much less widespread than those in higher source tiers. While most documents are journal articles and conference papers, the following document types are also present. Some of these are less than academic, but show how terms are used in practice; we include them in this source tier for brevity:

- Report [26], [43], [44]
- Thesis [45]
- Website [46], [47]
- Booklet [48]
- ChatGPT [49] (with its claims supported by [50])

The full set of sources that comprise this tier is [6], [21]–[24], [26], [43]–[73].

B. Terminology

This research is intended to describe the current state of software testing literature. To reduce potential bias, we

TABLE I: Categories of testing given by ISO/IEC and IEEE.

Term	Definition	Examples
Test Approach	A “high-level test implementation choice” that includes “test level, test type, test technique, test practice and ... static testing” [16, p. 10] and is used to “pick the particular test case values” [14, p. 465]	black or white box, minimum and maximum boundary value testing [14, p. 465]
Test Level ^a	A stage of testing “typically associated with the achievement of particular objectives and used to treat particular risks”, each performed in sequence [16, p. 12], [29, p. 6] with their “own documentation and resources” [14, p. 469]	unit/component testing, integration testing, system testing, acceptance testing [14, p. 467], [16, p. 12], [29, p. 6]
Test Practice	A “conceptual framework that can be applied to ... [a] test process to facilitate testing” [14, p. 471], [16, p. 14]	scripted testing, exploratory testing, automated testing [16, p. 20]
Test Technique ^b	A “procedure used to create or select a test model, identify test coverage items, and derive corresponding test cases” [16, p. 11] (similar in [14, p. 467]) that “generate evidence that test item requirements have been met or that defects are present in a test item” [29, p. vii]	equivalence partitioning, boundary value analysis, branch testing [16, p. 11]
Test Type	“Testing that is focused on specific quality characteristics” [14, p. 473], [16, p. 15], [29, p. 7]	security testing, usability testing, performance testing [14, p. 473], [16, p. 15]

^a Also called “test phase” or “test stage” (see relevant synonym discrepancies in Section IV-B2).

^b Also called “test design technique” [16, p. 11], [17].

do not invent or add our own classifications or kinds of relations. Instead, the notions of test approach categories (Section III-B1), synonyms (Section III-B2), and parent-child relations (Section III-B3) presented here arose naturally from the literature. We define them here for clarity since we use them throughout this paper, even though they are “results” of our research. While most information is presented explicitly in the sources we investigate, some appears more implicitly. This is a useful distinction to make, as implicit claims carry less weight than explicit ones. We call this property “rigidity” and define it in Section III-B4.

1) Approach Categories: While there are many ways to categorize software testing approaches, perhaps the most widely used is the one given by ISO/IEC and IEEE [16], where a test approach can be categorized as a test level, test type, test technique, or test practice (see Table I). The categories of “level” and “type” are particularly common; for example, six non-IEEE sources also give unit testing, integration testing, system testing, and acceptance testing as examples of test levels [1, pp. 443–445], [17], [11, pp. 5–6 to 5–7], [26, pp. 9, 13], [42, pp. 807–808], [66, p. 218], although they may use a different term for “test level” (see Table I). Because of their widespread use and their usefulness in dividing the domain of software testing into more manageable subsets, these categories are used for now. These four subcategories of test approaches can be loosely described by what they specify as follows:

- Level: What code is tested
- Practice: How the test is structured and executed
- Technique: How to derive inputs and/or outputs
- Type: Which software quality is evaluated

For example, boundary value analysis is a test technique since its inputs are “the boundaries of equivalence partitions” [16, p. 2], [29, p. 1]. Similarly, acceptance testing is a test level since its goal is to “enable a user, customer, or other authorized entity to determine whether to accept a

system or component” [14, p. 5], which requires the system or component to be developed and ready for testing.

While the vast majority of identified test approaches can be categorized in this way, we also note the potential significance of an “artifact” category, since some terms could refer to the application of a test approach and/or the resulting document(s). Because of this, a test approach being categorized as a category from Table I and an artifact is not a discrepancy (see Section IV-B1). Excluding this “artifact” category, the categories given in Table I seem to be orthogonal. For example, “a test type can be performed at a single test level or across several test levels” [16, p. 15], [29, p. 7], and “Keyword-Driven Testing [sic] can be applied at all testing levels ... and for various types of testing” [30, p. 4]. This means that a specific test approach can be derived by combining multiple test approaches from different categories; for example, formal reviews are a combination of formal testing and reviews.

One important side effect of the particularity of these terms is that they can be “overloaded”; for example, someone could reasonably yet imprecisely use any of these four categories as a synonym for “approach”. Even the prompt in [49] was imprecise, asking for the “type of software testing that focuses on looking for bugs where others have already been found.” Interestingly, ChatGPT later “corrected” this by calling defect-based testing an approach! Because of this, careful consideration needs to be given to discrepancies of this nature. For example, [43, p. 45] defines interface testing as “an integration test type that is concerned with testing ... interfaces”, but since it does not define “test type”, this may not have special significance.

2) Synonym Relations: The same approach often has many names. For example, specification-based testing is also called:

- 1) Black-Box Testing [16, p. 9], [17], [14, p. 431], [11, p. 5–10], [29, p. 8], [10, p. 399], [70, p. 344]

- 2) Closed-Box Testing [16, p. 9], [14, p. 431]
- 3) Functional Testing⁵ [14, p. 196], [10, p. 399], [43, p. 44] (implied by [14, p. 431], [29, p. 129])
- 4) Domain Testing [11, p. 5-10]
- 5) Specification-oriented Testing [1, p. 440, Fig. 12.2]
- 6) Input Domain-Based Testing (implied by [13, pp. 4-7 to 4-8])

These synonyms are the same as synonyms in natural language; while they may emphasize different aspects or express mild variations, their core meaning is nevertheless the same. Throughout our work, we use the terms “specification-based testing” and “structure-based testing” to articulate the source of the information for designing test cases, but a team or project also using grey-box testing may prefer the terms “black-box” and “white-box testing” for consistency. Thus, synonyms are not inherently problematic, although they can be (see Section IV-B2).

Synonym relations are often given explicitly in the literature. For example, [16, p. 9] lists “black-box testing” and “closed box testing” beneath the glossary entry for “specification-based testing”, meaning they are synonyms. “Black-box testing” is likewise given under “functional testing” in [14, p. 196], meaning it is also a synonym for “specification-based testing” through transitivity. However, these relations can also be less “rigid” (see Section III-B4); “functional testing” is listed in a cf. footnote to the glossary entry for “specification-based testing” [14, p. 431], which supports the previous claim but would not necessarily indicate a synonym relation on its own.

Similarly, [11, p. 5-10] says “specification-based techniques ... [are] sometimes also called domain testing techniques” in the SWEBOK Guide V4, from which the synonym of “domain testing” follows logically. However, its predecessor V3 only implies the more specific “input domain-based testing” as a synonym. The section on test techniques says “the classification of testing techniques presented here is based on how tests are generated: from the software engineer’s intuition and experience, the specifications, the code structure ...” [13, p. 4-7], and the first three subsections on the following page are “Based on the Software Engineer’s Intuition and Experience”, “Input Domain-Based Techniques”, and “Code-Based Techniques” [p. 4-8]. The order of the introductory list lines up with these sections, implying that “input domain-based techniques” are “generated[] from ... the specifications” (i.e., that input domain-based testing is the same as specification-based testing). Furthermore, the examples of input domain-based techniques given—equivalence partitioning, pairwise testing, boundary-value analysis, and random testing—are all given as children⁶

⁵This may be an outlier; see Section IV-C1.

⁶Pairwise testing is given as a child of combinatorial testing, which is itself a child of specification-based testing, by [11, pp. 5-11 to 5-12] and [29, Fig. 2], making it a “grandchild” of specification-based testing according to these sources.

of specification-based testing [16], [17], [29, Fig. 2]; even V4 agrees with this [11, pp. 5-11 to 5-12]!

3) Parent-Child Relations: Many test approaches are multi-faceted and can be “specialized” into others; for example, there are many subtypes of performance-related testing, such as load testing and stress testing (see Section V-C). These “specializations” will be referred to as “children” or “subapproaches” of the multi-faceted “parent”. This nomenclature also extends to other categories (such as “subtype”; see Section III-B1 and Table I) and software qualities (“subquality”). There are many reasons two approaches may have a parent-child relation, such as:

- 1) The parent approach is part of a mutually exclusive set. It is often trivial to classify a test approach as a child of one of a set of other, mutually exclusive test approaches. For example, ISO/IEC and IEEE say that “testing can take two forms: static and dynamic” [16, p. 17] and provide examples of subapproaches of static and dynamic testing [Fig. 1]. Likewise, Gerrard says “tests can be automated or manual” [26, p. 13] and gives subapproaches of automated and manual testing [Tab. 2], [44, Tab. 1].
- 2) One is “stronger than” or “subsumes” the other. When comparing adequacy criteria that “specif[y] requirements for testing” [10, p. 402], “criterion X is stronger than criterion Y if, for all programs P and all test sets T, X-adequacy implies Y-adequacy” [p. 432]. While this relation only “compares the thoroughness of test techniques, not their ability to detect faults” [p. 434], it is sufficient to consider one a child of the other.
- 4) Rigidity: A consequence of the use of natural language and the lack of standardization is the considerable degree of nuance that can get lost when referring to information sources. While most information is presented explicitly in the sources we investigate, some appears more implicitly. This is a useful distinction to make, as implicit claims carry less weight than explicit ones. We call this property “rigidity” and capture it when citing sources. This allows us to provide a more complete picture of the state of the literature; for example, we can view implicit discrepancies separately in Tables II and III, since additional context may rectify them. The following non-mutually exclusive reasons for information to be considered “implicit” emerged, and the given keywords are used to identify them (see the relevant source code):

- 1) The information is implied. The implicit categorizations of “test type” by [7, pp. 53–58] (see Table IV) are an example of this. The given test approaches are not explicitly called “test types”, as the term is used more loosely to refer to different kinds of testing—what should be called “test approaches” as per Table I. However, this set of test approaches are “based on the associated quality characteristic and its associated quality attributes” [p. 53], implying that

they are test types. Cases such as this are indicated by a question mark or one of the following keywords: “implied”, “inferred”, or “likely”.

- 2) The information is not universal. Reference [14, p. 372] defines “regression testing” as “testing required to determine that a change to a system component has not adversely affected functionality, reliability or performance and has not introduced additional defects”. While reliability testing, for example, is not always a subset of regression testing (since it may be performed in other ways), it can be accomplished by regression testing, so there is sometimes a parent-child relation (defined in Section III-B3) between them. Cases such as this are indicated by one of the following keywords: “can be”, “should be”, “ideally”, “usually”, “most”, “likely”, “often”, or “if”.
- 3) The information is conditional. As a more specific case of information not being universal, sometimes prerequisites must be satisfied for information to apply. For example, branch condition combination testing is equivalent to (and is therefore a synonym of) exhaustive testing if “each subcondition is viewed as a single input” [1, p. 464]. Likewise, statement testing can be used for (and is therefore a child of) unit testing if there are “less than 5000 lines of code” [p. 481]. Cases such as this are indicated by the keyword “can be” or “if”.
- 4) The information is dubious. This happens when there is reason to doubt the information provided. If a source claims one thing that is not true, related claims lose credibility. For example, the incorrect claim that “white-box testing”, “grey-box testing”, and “black-box testing” are synonyms for “module testing”, “integration testing”, and “system testing”, respectively, casts doubt on the claim that “red-box testing” is a synonym for “acceptance testing” [73, p. 18]. Doubts such as this can also originate from other sources. Reference [43, p. 48] gives “user scenario testing” as a synonym of “use case testing”, even though “an actor [in use case testing] can be ... another system” [29, p. 20], which does not fit as well with the label “user scenario testing”. However, since a system can be seen as a “user” of the test item, this synonym relation is treated as implicit instead of as an outright discrepancy. Cases such as this are indicated by a question mark or one of the following keywords: “inferred”, “should be”, “ideally”, “likely”, “if”, or “although”.

Discrepancies based on implicit information are themselves implicit. These are automatically detected when generating graphs and analyzing discrepancies (see Section III-E) by looking for the indicators of uncertainty mentioned above (see the relevant source code). These are used when creating the glossaries to capture varying degrees of nuance, such as when a test approach “can be”

a child of another or is a synonym of another “most of the time” but not always.

C. Procedure

We track terminology used in the literature by building glossaries. The one most central to our research is our test approach glossary, where we give each test approach its own row to record its name and any given categories (see Section III-B1), synonyms (see Section III-B2), parents (see Section III-B3), and definitions. If no category is given, the “approach” category is assigned (with no accompanying citation) as a “catch-all” category. All other fields may be left blank, but a lack of definition indicates that the approach should be investigated further to see if its inclusion is meaningful (see Section III-D). Any additional information from other sources is added to or merged with the existing information in our glossary where appropriate. This includes the generic “approach” category being replaced with a more specific one, an additional synonym being mentioned, or another source describing an already-documented parent-child relation. If any new information contradicts existing information (or otherwise indicates something is wrong), this is investigated and documented (see Section IV), which may be done in a separate document and/or in the glossary itself. Sometimes, new information does not conflict with existing information, in which case the clearest and most concise version is kept, or they are merged to paint a more complete picture. Finally, we record any other notes, such as questions, prerequisites, and other resources to investigate.

We use similar procedures to track software qualities and supplementary terminology (either shared by multiple approaches or too complicated to explain inline) in separate glossaries with a similar format. The name, definition, and synonym(s) of all terms are tracked, as well as any precedence for a related test type for a given software quality. We use heuristics to guide this process for all three glossaries to increase confidence that all terms are identified, paying special attention to the following when investigating a new source:

- glossaries and lists of terms,
- testing-related terms (e.g., terms containing “test(ing)”, “validation”, or “verification”),
- terms that had emerged as part of already-discovered testing approaches, especially those that were ambiguous or prompted further discussion (e.g., terms containing “performance”, “recovery”, “component”, “bottom-up”, or “configuration”), and
- terms that implied testing approaches.

We apply these heuristics to most investigated sources, especially established standards (see Section III-A1), in their entirety. Some sources, however, are only partially investigated, such as those chosen for a specific area of interest or based on a test approach that was determined

to be out-of-scope. These include the following sources as described in Section III-D: [23], [24], [38], [39], [74]–[77].

During the first pass of data collection, all software-testing-focused terms are included. Some of them are less applicable to test case automation or too broad, so they will be omitted during future analysis.

D. Undefined Terms

The search process led to some testing approaches being mentioned without definition; [16] and [7] in particular introduced many. Once the standards in Section III-A1 had been exhausted, we devised a strategy to look for sources that explicitly define these terms, consistent with our snowballing approach. This uncovers new approaches, both in and out of scope (such as EManations SECurity (EMSEC) testing and aspects of orthogonal array testing; see Section II).

The following terms (and their respective related terms) were explored in the following sources, bringing the number of testing approaches from 473 to 557 and the number of undefined terms from 174 to 190 (the assumption can be made that about 81% of added terms also included a definition):

- Assertion Checking: [55], [57], [67]
- Loop Testing⁷: [59], [61], [62], [69]
- EMSEC Testing: [19], [20]
- Asynchronous Testing: [64]
- Performance(-related) Testing: [68]
- Web Application Testing: [40], [43]
 - HTML Testing: [44], [58], [73]
 - Document Object Model (DOM) Testing: [51]
- Sandwich Testing: [71], [72]
- Orthogonal Array Testing⁸: [21], [22]
- Backup Testing⁹: [45]

E. Tools

To better visualize how test approaches relate to each other, we develop a tool to automatically generate graphs of these relations. All parent-child relations are graphed, since they are guaranteed to be visually meaningful. Synonym relations, however, are either excluded from or included in graphs as follows. For each synonym pair, at least one term will have its own row (or else it would not appear in the glossary at all), so the following cases are possible:

1. (Excluded)

Only one synonym has its own row. This is a “typical” synonym relation (see Section III-B2) where the terms are interchangeable. The synonym could be included as an alternate name inside the node of its partner, but this would unnecessarily clutter the graphs.

⁷References [38] and [39] were used as reference for terms but not fully investigated, [75] and [76] were added as potentially in scope, and [74] and [77] were added as out-of-scope examples.

⁸References [23] and [24] were added as out-of-scope examples.

⁹See Section IV-D.

2. (Included)

Both synonyms have their own row in the glossary. This may indicate that the synonym relation is incorrect, since separate rows in the glossary define separate approaches (with their own definitions, nuances, etc.).

3. (Included)

Two synonym pairs share a synonym without its own row. This is a transitive extension to the previous case. If two distinct approaches share a synonym, that implies that they are synonyms themselves, resulting in the same possibility of the relation being incorrect.

Since these graphs tend to be large, it is useful to focus on specific subsets of them. These can be generated from a given subset of approaches, such as those in a selected approach category (see Section III-B1) or those pertaining to recovery or scalability; the latter are shown in Figures 2a and 3a, respectively. By specifying sets of approaches and relations to add or remove, these generated graphs can then be updated in accordance with our recommendations; applying those given in Sections V-A, V-B, and V-C results in the updated graphs in Figures 2b, 3b, and 4, respectively. Any added approaches or relations are colored **orange**.

IV. Discrepancies

After gathering all these data¹⁰, we find many discrepancies. To better understand and analyze them, we group them by their syntax and their semantics. Syntactic discrepancies (Section IV-A) describe how a discrepancy manifests, such as information that is wrong (a “mistake”; Section IV-A1) or missing (an “omission”; Section IV-A2). On the other hand, semantic discrepancies (Section IV-B) describe the knowledge domain in which a discrepancy manifests, such as a discrepancy between synonyms (Section IV-B2) or parent-child relations (Section IV-B3). As an example, the structure of tours can be defined as either quite general [16, p. 34] or “organized around a special focus” [17]. This is a case of contradictory definitions, so it appears with both the contradictions (Section IV-A3) and the definition discrepancies (Section IV-B4). Within these sections, “less significant” discrepancies are omitted for brevity, and those remaining are then sorted based on their source tier (see Section III-A).

A summary of how many discrepancies there are by syntax and by semantics is shown in Tables II and III, respectively, where a given row corresponds to the number of discrepancies either within that source tier (see Section III-A) and/or with a “more trusted” one (i.e., a previous row in the table). The numbers of (Exp)licit and (Imp)licit (see Section III-B4) discrepancies are also presented in these tables. Since each discrepancy is grouped by syntax and by semantics, the totals per source and

¹⁰Available in ApproachGlossary.csv, QualityGlossary.csv, and SuppGlossary.csv at <https://github.com/samm82/TestGen-Thesis>.

TABLE II: Breakdown of identified Syntactic Discrepancies by Source Tier.

Source Tier	Mistakes		Omissions		Contradictions		Ambiguities		Overlaps		Redundancies ^a		Total
	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	
Established Standards	7	1	2	0	17	10	4	0	8	0	0	0	49
Terminology Collections	11	0	1	0	32	17	14	3	5	1	2	0	86
Textbooks	6	0	1	0	37	4	5	0	1	0	0	0	54
Papers and Others	8	1	4	0	21	19	9	3	2	1	2	0	70
Total	32	2	8	0	107	50	32	6	16	2	4	0	259

^aSection omitted for brevity.

TABLE III: Breakdown of identified Semantic Discrepancies by Source Tier.

Source Tier	Categories		Synonyms		Parents		Definitions		Terminology		Citations		Total
	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	Exp	Imp	
Established Standards	10	7	4	2	6	0	13	2	5	0	0	0	49
Terminology Collections	10	14	9	3	9	2	20	0	13	2	4	0	86
Textbooks	2	0	14	0	10	3	17	1	7	0	0	0	54
Papers and Others	10	12	13	9	11	1	5	0	7	2	0	0	70
Total	32	33	40	14	36	6	55	3	32	4	4	0	259

grand totals in these tables are equal. However, the numbers of discrepancies listed in Sections IV-A and IV-B are not equal, as those automatically uncovered based on semantics are only listed in the corresponding category section for clarity; they still contribute to the counts in Table II!

Moreover, certain “subsets” of testing revealed many interconnected discrepancies. These are given in their respective sections as a “third view” to keep related information together, but still count towards the syntaxes and semantics of discrepancies listed above, causing a further mismatch between the counts in Tables II and III and the counts in Sections IV-A and IV-B. The “problem” subsets of testing include Functional Testing, Recovery Testing, Scalability Testing, and Compatibility Testing.

A. Syntactic Discrepancies

The following sections list observed discrepancies grouped by how the discrepancy manifests. These include Mistakes, Omissions, Contradictions, Ambiguities, Overlaps, and Redundancies¹¹.

1) Mistakes: The following are cases where information is incorrect; this includes cases Terminology included that should not have been, untrue claims about Citations, and simple typos:

- Since errors are distinct from defects/faults [11, p. 12-3], [9, pp. 128, 140], [10, pp. 399–400], error guessing should instead be called “defect guessing” if it is based on a “checklist of potential defects” [29, p. 29] or “fault guessing” if it is a “fault-based technique” [13, p. 4-9] that “anticipate[s] the most plausible faults in each SUT” [11, p. 5-13]. One (or both) of these proposed terms may be useful

in tandem with “error guessing”, which would focus on errors as traditionally defined; this would be a subapproach of error-based testing (implied by [10, p. 399]).

- Similarly, “fault seeding” is not a synonym of “error seeding” as claimed by [14, p. 165] and [10, p. 427]. The term “error seeding”, used by [7, p. 34], [14, p. 165], and [10, p. 427], should be abandoned in favour of “fault seeding”, as it is defined as the “process of intentionally adding known faults to those already in a computer program ... [to] estimat[e] the number of faults remaining” [14, p. 165] based on the ratio between the number of new faults and the number of introduced faults that were discovered [10, p. 427].
- Since keyword-driven testing can be used for automated or manual testing [30, pp. 4, 6], the claim that “test cases can be either manual test cases or keyword test cases” [p. 6] is incorrect.
- Reference [17] classifies ML model testing as a test level, which it defines as “a specific instantiation of a test process”: a vague definition that does not match the one in Table I.
- The terms “acceleration tolerance testing” and “acoustic tolerance testing” seem to only refer to software testing in [7, p. 56]; elsewhere, they seem to refer to testing the acoustic tolerance of rats [78] or the acceleration tolerance of astronauts [79, p. 11], aviators [80, pp. 27, 42], or catalysts [81, p. 1463], which don’t exactly seem relevant...
- The differences between the terms “error”, “failure”, “fault”, “defect” are significant and meaningful [9, pp. 128, 139–140], [10, pp. 399–400], [11, p. 12-3], but [15, pp. 13–14] “just call[s] it what it is and get[s]

¹¹Section omitted for brevity.

on with it”, abandoning these four terms, “problem”, “incident”, “anomaly”, “variance”, “inconsistency”, “feature” (!), and “a list of unmentionable terms” in favour of “bug”; after all, “there’s no reason to dice words”!

- Peters and Pedrycz claim that “structural testing subsumes white box testing” but they seem to describe the same thing; it says “structure tests are aimed at exercising the internal logic of a software system” and “in white box testing ..., using detailed knowledge of code, one creates a battery of tests in such a way that they exercise all components of the code (say, statements, branches, paths)” on the same page [1, p. 447]!
- Reference [43, p. 46] says that the goal of negative testing is “showing that a component or system does not work” which is not true; if robustness is an important quality for the system, then testing the system “in a way for which it was not intended to be used” [17] (i.e., negative testing) is one way to help test this!

2) Omissions: The following are cases where information (usually Definitions) should have been included but was not:

- Integration testing, system testing, and system integration testing are all listed as “common test levels” [16, p. 12], [29, p. 6], but no definitions are given for the latter two, making it unclear what “system integration testing” is; it is a combination of the two? somewhere on the spectrum between them? It is listed as a child of integration testing by [17] and of system testing by [7, p. 23].
- Similarly, component testing, integration testing, and component integration testing are all listed in [14], but “component integration testing” is only defined as “testing of groups of related components” [14, p. 82]; it is a combination of the two? somewhere on the spectrum between them? As above, it is listed as a child of integration testing by [17].
- Reference [43, p. 42] says “See boundary value analysis,” for the glossary entry of “boundary value testing” but does not provide this definition.

3) Contradictions: The following are cases where multiple sources of information (sometimes within the same document!) disagree; note that cases where all sources of information are incorrect are considered contradictions and not Mistakes, since this would require analysis that has not been performed yet:

- Regression testing and retesting are sometimes given as two distinct approaches [16, p. 8], [7, p. 34], but sometimes regression testing is defined as a form of “selective retesting” [14, p. 372], [11, pp. 5-8, 6-5, 7-5 to 7-6], [52, p. 3]. Moreover, the two possible variations of regression testing given by [10, p. 411] are “retest-all” and “selective retest”, which

is possibly the source of the above misconception. This discrepancy creates a cyclic relation between regression testing and selective retesting.

- A component is an “entity with discrete structure ... within a system considered at a particular level of analysis” [34] and “the terms module, component, and unit [sic] are often used interchangeably or defined to be subelements of one another in different ways depending upon the context” with no standardized relationship [14, p. 82]. For example, [17] defines them as synonyms while [53, p. 107] says “components differ from classical modules for being re-used in different contexts independently of their development”. Additionally, since components are structurally, functionally, or logically discrete [14, p. 419] and “can be tested in isolation” [17], “unit/component/module testing” could refer to the testing of both a module and a specific function in a module, introducing a further level of ambiguity.
- Performance testing and security testing are given as subtypes of reliability testing by [33], but these are all listed separately by [7, p. 53].
- Similarly, random testing is a subtechnique of specification-based testing [16, pp. 7, 22], [17], [11, p. 5-12], [29, pp. 5, 20, Fig. 2] but is listed separately by [7, p. 46].
- Path testing “aims to execute all entry-to-exit control flow paths in a SUT’s control flow graph” [11, p. 5-13] (similar in [15, p. 119]), but [14, p. 316] adds that it can also be “designed to execute ... selected paths.”
- The structure of tours can be defined as either quite general [16, p. 34] or “organized around a special focus” [17].
- Alpha testing is performed by “users within the organization developing the software” [14, p. 17], “a small, selected group of potential users” [11, p. 5-8], or “roles outside the development organization” conducted “in the developer’s test environment” [17].
- “Use case testing” is given as a synonym of “scenario testing” by [17] but listed separately by [16, Fig. 2] and described as a “common form of scenario testing” in [29, p. 20]. This implies that use case testing may instead be a child of user scenario testing (see Table V).
- The terms “test level” and “test stage” are given as synonyms ([17]; implied by [26, p. 9]), but [11, p. 5-6] says “[test] levels can be distinguished based on the object of testing, the target, or on the purpose or objective” and calls the former “test stages”, giving the term a child relation (see Section III-B3) to “test level” instead. However, the examples listed—unit testing, integration testing, system testing, and acceptance testing [11, pp. 5-6 to 5-7]—are commonly categorized as “test levels” (see Section III-B1).
- While [15, p. 120] implies that condition testing is a subtechnique of path testing, [10, Fig. 13.17] says

that multiple condition coverage (which seems to be a synonym of condition coverage [p. 422]) does not subsume and is not subsumed by path coverage.

- Load testing is performed with loads “between anticipated conditions of low, typical, and peak usage” [16, p. 5] or loads that are as large as possible [15, p. 86].
- State testing requires that “all states in the state model ... [are] ‘visited’” in [29, p. 19] which is only one of its possible criteria in [15, pp. 82-83].
- Reference [14, p. 456] says system testing is “conducted on a complete, integrated system” (which [1, Tab. 12.3] and [10, p. 439] agree with), while [15, p. 109] says it can also be done on “at least a major portion” of the product.
- “Walkthroughs” and “structured walkthroughs” are given as synonyms by [17] but [1, p. 484] implies that they are different, saying a more structured walkthrough may have specific roles.
- Reference [15, p. 92] says that reviews are “the process[es] under which static white-box testing is performed” but correctness proofs are given as another example by [10, pp. 418-419].
- Reference [43, p. 46] says “negative testing is related to the testers’ attitude rather than a specific test approach or test design technique”; while [29] seems to support this idea of negative testing being at a “higher” level than other approaches, it also implies that it is a test technique [pp. 10, 14].

4) Ambiguities: The following are cases where information (usually Definitions or distinctions between Terminology) is unclear:

- The distinctions between development testing [14, p. 136], developmental testing [7, p. 30], and developer testing [7, p. 39], [26, p. 11] are unclear and seem miniscule.
- Reference [17] defines “Machine Learning (ML) model testing” and “ML functional performance” in terms of “ML functional performance criteria”, which is defined in terms of “ML functional performance metrics”, which is defined as “a set of measures that relate to the functional correctness of an ML system”. The use of “performance” (or “correctness”) in these definitions is at best ambiguous and at worst incorrect.
- “Installability testing” is given as a test type [16, p. 22], [29, p. 38], [14, p. 228], while “installation testing” is given as a test level [10, p. 439]. Since “installation testing” is not given as an example of a test level throughout the sources that describe them (see Section III-B1), it is likely that the term “installability testing” with all its related information should be used instead.
- Reference [17] claims that code inspections are related to peer reviews but [15, pp. 94-95] makes them quite distinct.

5) Overlaps: The following are cases where information overlaps, such as nonatomic Definitions and Terminology:

- Reference [16, p. 34] gives the “landmark tour” as an example of “a tour used for exploratory testing”, but they also use the analogy of “a tour guide lead[ing] a tourist through the landmarks of a big city” to describe tours in general. Is the distinction between them the fact that landmark tours are pre-planned and follow a decided-upon sequence [p. 34]?
- ISO/IEC and IEEE say that “test level” and “test phase” are synonyms, both meaning a “specific instantiation of [a] test sub-process” ([14, pp. 469, 470]; [31, p. 9]), but they have other definitions as well. “Test level” can also refer to the scope of a test process; for example, “across the whole organization” or only “to specific projects” [16, p. 24] and “test phase” can also refer to the “period of time in the software life cycle” when testing occurs [14, p. 470], usually after the implementation phase [14, pp. 420, 509], [42, p. 56].
- ISO/IEC and IEEE define “error” as “a human action that produces an incorrect result”, but also as “an incorrect result” itself [9, p. 128]. Since faults are inserted when a developer makes an error [11, p. 12-3], [9, pp. 128, 140], [10, pp. 399-400], this means that they are “incorrect results”, making “error” and “fault” synonyms and the distinction between them less useful.
- Additionally, “error” can also be defined as “the difference between a computed, observed, or measured value or condition and the true, specified, or theoretically correct value or condition” [9, p. 128] (similar in [11, pp. 17-18 to 17-19, 18-7 to 18-8]). While this is a widely used definition, particularly in mathematics, it makes some test approaches ambiguous; for example, back-to-back testing is “testing in which two or more variants of a program are executed with the same inputs, the outputs are compared, and errors are analyzed in case of discrepancies” [9, p. 30] (similar in [17]), which seems to refer to this definition of “error”.
- The SWEBOK Guide V4 defines “privacy testing” as testing that “assess[es] the security and privacy of users’ personal data to prevent local attacks” [11, p. 5-10]; this seems to overlap (both in scope and name) with the definition of “security testing” in [16, p. 7]: testing “conducted to evaluate the degree to which a test item, and associated data and information, [sic] are protected so that” only “authorized persons or systems” can use them as intended.
- “Orthogonal array testing” [11, pp. 5-1, 5-11] and “operational acceptance testing” [7, p. 30] have the same acronym (“OAT”).

B. Semantic Discrepancies

The following sections list observed discrepancies grouped by what area the discrepancy manifests in. These

include Approach Category Discrepancies, Synonym Relation Discrepancies, Parent-Child Relation Discrepancies, Definition Discrepancies, Terminology Discrepancies, and Citation Discrepancies.

1) Approach Category Discrepancies: While the IEEE categorization of testing approaches described in Table I is useful, it is not without its faults. One issue, which is not inherent to the categorization itself, is the fact that it is not used consistently. The most blatant example of this is that ISO/IEC and IEEE [14, p. 286] describe mutation testing as a methodology, even though this is not one of the categories they created! Additionally, the boundaries between approaches within a category may be unclear: “although each technique is defined independently of all others, in practice [sic] some can be used in combination with other techniques” [29, p. 8]. For example, “the test coverage items derived by applying equivalence partitioning can be used to identify the input parameters of test cases derived for scenario testing” [p. 8]. Even the categories themselves are not consistently defined, and some approaches are categorized differently by different sources; these differences are tracked so they can be analyzed more systematically.

- Since keyword-driven testing can be used for automated or manual testing [30, pp. 4, 6], the claim that “test cases can be either manual test cases or keyword test cases” [p. 6] is incorrect.
- Reference [17] classifies ML model testing as a test level, which it defines as “a specific instantiation of a test process”: a vague definition that does not match the one in Table I.
- Reference [43, p. 46] says “negative testing is related to the testers’ attitude rather than a specific test approach or test design technique”; while [29] seems to support this idea of negative testing being at a “higher” level than other approaches, it also implies that it is a test technique [pp. 10, 14].

Some category discrepancies can be detected automatically, such as test approaches with more than one category. These are given in Table IV and include experience-based testing, which is of particular note. ISO/IEC and IEEE categorize experience-based testing as both a test design technique and a test practice on the same page—twice [16, Fig. 2, p. 34]! These authors say “experience-based testing practices like exploratory testing ... are not ... techniques for designing test cases”, although they “can use ... test techniques” [29, p. viii], which they support in [16, p. 33] along with scripted testing. This implies that “experience-based test design techniques” are used by the practice of experience-based testing which is not itself a test technique (and similarly with scripted testing). If this is the case, it blurs the line between “practice” and “technique”, which may explain why experience-based testing is categorized inconsistently in the literature.

Subapproaches of experience-based testing, such as error guessing and exploratory testing, are also categorized am-

biguously, causing confusion on how categories and parent-child relations (see Section III-B3) interact. Reference [16, p. 34] says that a previous standard [29] “describes the experience-based test design technique of error guessing. Other experience-based test practices include (but are not limited to) exploratory testing ..., tours, attacks, and checklist-based testing”. This seems to imply that error guessing is both a technique and a practice, which does not make sense if these categories are orthogonal. These kinds of inconsistencies between parent and child test approach categorizations may indicate that categories are not transitive or that more thought must be given to them.

2) Synonym Relation Discrepancies: As mentioned in Section III-B2, synonyms do not inherently signify a discrepancy. Unfortunately, there are many instances of incorrect or ambiguous synonyms, such as the following:

- A component is an “entity with discrete structure ... within a system considered at a particular level of analysis” [34] and “the terms module, component, and unit [sic] are often used interchangeably or defined to be subelements of one another in different ways depending upon the context” with no standardized relationship [14, p. 82]. For example, [17] defines them as synonyms while [53, p. 107] says “components differ from classical modules for being re-used in different contexts independently of their development”. Additionally, since components are structurally, functionally, or logically discrete [14, p. 419] and “can be tested in isolation” [17], “unit/component/module testing” could refer to the testing of both a module and a specific function in a module, introducing a further level of ambiguity.
- ISO/IEC and IEEE say that “test level” and “test phase” are synonyms, both meaning a “specific instantiation of [a] test sub-process” ([14, pp. 469, 470]; [31, p. 9]), but they have other definitions as well. “Test level” can also refer to the scope of a test process; for example, “across the whole organization” or only “to specific projects” [16, p. 24] and “test phase” can also refer to the “period of time in the software life cycle” when testing occurs [14, p. 470], usually after the implementation phase [14, pp. 420, 509], [42, p. 56].
- “Use case testing” is given as a synonym of “scenario testing” by [17] but listed separately by [16, Fig. 2] and described as a “common form of scenario testing” in [29, p. 20]. This implies that use case testing may instead be a child of user scenario testing (see Table V).
- The terms “test level” and “test stage” are given as synonyms ([17]; implied by [26, p. 9]), but [11, p. 5-6] says “[test] levels can be distinguished based on the object of testing, the target, or on the purpose or objective” and calls the former “test stages”, giving the term a child relation (see Section III-B3) to “test level” instead. However, the examples listed—

TABLE IV: Test approaches with more than one category.

Approach	Category 1	Category 2
Capacity Testing	Technique [29, p. 38]	Type [16, p. 22], [7, p. 53], [31, p. 2]
Checklist-based Testing	Practice [16, p. 34]	Technique [17]
Data-driven Testing	Practice [16, p. 22]	Technique [43, p. 43]
End-to-end Testing	Type [17]	Technique [7, p. 47], [72, pp. 601, 603, 605–606]
Endurance Testing	Technique [29, p. 38]	Type [31, p. 2]
Experience-based Testing	Practice [16, pp. 22, 34], [29, p. viii]	Technique [16, pp. 4, 22], [17], [11, p. 5-13], [7, pp. 46, 50], [29, p. 4]
Exploratory Testing	Practice [16, pp. 20, 22, 34], [29, p. viii]	Technique [16, p. 34], [11, p. 5-14], [7, p. 50]
Load Testing	Technique [29, p. 38]	Type [16, pp. 5, 20, 22], [17], [14, p. 253]
Model-based Testing	Practice [16, p. 22], [29, p. viii]	Technique [43, p. 4]
Mutation Testing	Methodology [14, p. 286]	Technique [11, p. 5-15], [10, pp. 428–429]
Performance Testing	Technique [29, p. 38]	Type [16, pp. 7, 22, 26–27], [29, p. 7]
Stress Testing	Technique [29, p. 38]	Type [16, pp. 9, 22], [14, p. 442]

unit testing, integration testing, system testing, and acceptance testing [11, pp. 5-6 to 5-7]—are commonly categorized as “test levels” (see Section III-B1).

- The differences between the terms “error”, “failure”, “fault”, “defect” are significant and meaningful [9, pp. 128, 139–140], [10, pp. 399–400], [11, p. 12-3], but [15, pp. 13–14] “just call[s] it what it is and get[s] on with it”, abandoning these four terms, “problem”, “incident”, “anomaly”, “variance”, “inconsistency”, “feature” (!), and “a list of unmentionable terms” in favour of “bug”; after all, “there’s no reason to dice words”!
- Reference [17] claims that code inspections are related to peer reviews but [15, pp. 94–95] makes them quite distinct.
- “Walkthroughs” and “structured walkthroughs” are given as synonyms by [17] but [1, p. 484] implies that they are different, saying a more structured walkthrough may have specific roles.
- Peters and Pedrycz claim that “structural testing subsumes white box testing” but they seem to describe the same thing: it says “structure tests are aimed at exercising the internal logic of a software system” and “in white box testing ..., using detailed knowledge of code, one creates a battery of tests in such a way that they exercise all components of the code (say, statements, branches, paths)” on the same page [1, p. 447]!

There are also cases in which a term is given as a synonym to two (or more) terms that are not synonyms themselves. Sometimes, these terms are synonyms; for example, [17] says “use case testing”, “user scenario testing”, and “scenario testing” are all synonyms (although there may be a slight distinction; see Table V and Section IV-B2). However, this does not always make sense. We identify nine such cases through automatic analysis of the generated graphs. The following three are the most prominent examples:

1) Invalid Testing:

- Error Tolerance Testing [43, p. 45]

- Negative Testing [17] (implied by [29, p. 10])

2) Soak Testing:

- Endurance Testing [29, p. 39]
- Reliability Testing¹² [26, Tab. 2], [44, Tab. 1, p. 26]

3) Link Testing:

- Branch Testing (implied by [29, p. 24])
- Component Integration Testing [43, p. 45]
- Integration Testing (implied by [26, p. 13])

3) Parent-Child Relation Discrepancies: Parent-Child Relations are also not immune to difficulties; for example, performance testing and security testing are given as subtypes of reliability testing by [33], but these are all listed separately by [7, p. 53].

Additionally, some self-referential definitions imply that a test approach is a parent of itself. Since these are by nature self-contained within a given source, these are counted once as explicit discrepancies within their sources in Tables II and III. For example, performance and usability testing are both given as subapproaches of themselves [26, Tab. 2], [44, Tab. 1].

There are also pairs of synonyms where one is described as a subapproach of the other, abusing the meaning of “synonym” and causing confusion. We identify 15 of these pairs through automatic analysis of the generated graphs, with the most prominent given in Table V. Of particular note is the relation between path testing and exhaustive testing. While [10, p. 421] claims that path testing done completely “is equivalent to exhaustively testing the program”¹³, this overlooks the effects of input data [15, p. 121], [29, p. 129], [1, p. 467] and implementation issues [p. 476] on the code’s behaviour. Exhaustive testing requires “all combinations of input values and preconditions ... [to be] tested” [16, p. 4] (similar in [17], [15, p. 121]).

4) Definition Discrepancies: Perhaps the most interesting category for those seeking to understand how to

¹²Endurance testing is given as a child of reliability testing by [7, p. 55], although the terms are not synonyms.

¹³The contradictory definitions of path testing given in Section IV-B4 add another layer of complexity to this claim.

TABLE V: Pairs of test approaches with both parent-child and synonym relations.

“Child”	→	“Parent”	Parent-Child Source(s)	Synonym Source(s)
All Transitions Testing	→	State Transition Testing	[29, p. 19]	[43, p. 15]
Co-existence Testing	→	Compatibility Testing	[16, p. 3], [29, Tab. A.1], [33]	[29, p. 37]
Fault Tolerance Testing	→	Robustness Testing ^a	[7, p. 56]	[17]
Functional Testing	→	Specification-based Testing ^b	[29, p. 38]	[14, p. 196], [10, p. 399], [43, p. 44]
Orthogonal Array Testing	→	Pairwise Testing	[21, p. 1055]	[11, p. 5-11], [22, p. 473]
Path Testing	→	Exhaustive Testing	[1, pp. 466-467, 476]	[10, p. 421]
Performance Testing	→	Performance-related Testing	[16, p. 22], [29, p. 38]	[68, p. 1187]
Static Analysis	→	Static Testing	[16, pp. 9, 17, 25, 28], [17]	[1, p. 438]
Structural Testing	→	Structure-based Testing	[15, pp. 105-121]	[16, p. 9], [17], [14, pp. 443-444]
Use Case Testing	→	Scenario Testing ^c	[29, p. 20]	[17], [43, pp. 47-49]

^aFault tolerance testing may also be a subapproach of reliability testing [14, p. 375], [11, p. 7-10], which is distinct from robustness testing [7, p. 53].

^bSee Section IV-C1.

^cSee Section IV-B2.

apply a given test approach, there are many discrepancies between how test approaches, as well as supporting terms, are defined:

- Reference [16, p. 34] gives the “landmark tour” as an example of “a tour used for exploratory testing”, but they also use the analogy of “a tour guide lead[ing] a tourist through the landmarks of a big city” to describe tours in general. Is the distinction between them the fact that landmark tours are pre-planned and follow a decided-upon sequence [p. 34]?
- Integration testing, system testing, and system integration testing are all listed as “common test levels” [16, p. 12], [29, p. 6], but no definitions are given for the latter two, making it unclear what “system integration testing” is; it is a combination of the two? somewhere on the spectrum between them? It is listed as a child of integration testing by [17] and of system testing by [7, p. 23].
- Similarly, component testing, integration testing, and component integration testing are all listed in [14], but “component integration testing” is only defined as “testing of groups of related components” [14, p. 82]; it is a combination of the two? somewhere on the spectrum between them? As above, it is listed as a child of integration testing by [17].
- ISO/IEC and IEEE define “error” as “a human action that produces an incorrect result”, but also as “an incorrect result” itself [9, p. 128]. Since faults are inserted when a developer makes an error [11, p. 12-3], [9, pp. 128, 140], [10, pp. 399-400], this means that they are “incorrect results”, making “error” and “fault” synonyms and the distinction between them less useful.
- Additionally, “error” can also be defined as “the difference between a computed, observed, or measured value or condition and the true, specified, or theoretically correct value or condition” [9, p. 128] (similar in [11, pp. 17-18 to 17-19, 18-7 to 18-8]). While this is a

widely used definition, particularly in mathematics, it makes some test approaches ambiguous; for example, back-to-back testing is “testing in which two or more variants of a program are executed with the same inputs, the outputs are compared, and errors are analyzed in case of discrepancies” [9, p. 30] (similar in [17]), which seems to refer to this definition of “error”.

- The SWEBOK Guide V4 defines “privacy testing” as testing that “assess[es] the security and privacy of users’ personal data to prevent local attacks” [11, p. 5-10]; this seems to overlap (both in scope and name) with the definition of “security testing” in [16, p. 7]: testing “conducted to evaluate the degree to which a test item, and associated data and information, [sic] are protected so that” only “authorized persons or systems” can use them as intended.
- Path testing “aims to execute all entry-to-exit control flow paths in a SUT’s control flow graph” [11, p. 5-13] (similar in [15, p. 119]), but [14, p. 316] adds that it can also be “designed to execute ... selected paths.”
- The structure of tours can be defined as either quite general [16, p. 34] or “organized around a special focus” [17].
- Alpha testing is performed by “users within the organization developing the software” [14, p. 17], “a small, selected group of potential users” [11, p. 5-8], or “roles outside the development organization” conducted “in the developer’s test environment” [17].
- Reference [17] defines “Machine Learning (ML) model testing” and “ML functional performance” in terms of “ML functional performance criteria”, which is defined in terms of “ML functional performance metrics”, which is defined as “a set of measures that relate to the functional correctness of an ML system”. The use of “performance” (or “correctness”) in these definitions is at best ambiguous and at worst incorrect.
- Load testing is performed with loads “between antici-

pated conditions of low, typical, and peak usage” [16, p. 5] or loads that are as large as possible [15, p. 86].

- State testing requires that “all states in the state model ... [are] ‘visited’” in [29, p. 19] which is only one of its possible criteria in [15, pp. 82-83].
- Reference [14, p. 456] says system testing is “conducted on a complete, integrated system” (which [1, Tab. 12.3] and [10, p. 439] agree with), while [15, p. 109] says it can also be done on “at least a major portion” of the product.
- Reference [15, p. 92] says that reviews are “the process[es] under which static white-box testing is performed” but correctness proofs are given as another example by [10, pp. 418–419].
- Reference [43, p. 46] says that the goal of negative testing is “showing that a component or system does not work” which is not true; if robustness is an important quality for the system, then testing the system “in a way for which it was not intended to be used” [17] (i.e., negative testing) is one way to help test this!
- Reference [43, p. 42] says “See boundary value analysis,” for the glossary entry of “boundary value testing” but does not provide this definition.

5) Terminology Discrepancies: While some discrepancies exist because the definition of a term is wrong, others exist because term’s name or label is wrong! This could be considered a “sister” category of Definition Discrepancies, but these discrepancies seemed different enough to merit their own category. The following are examples of these discrepancies:

- Since errors are distinct from defects/faults [11, p. 12-3], [9, pp. 128, 140], [10, pp. 399–400], error guessing should instead be called “defect guessing” if it is based on a “checklist of potential defects” [29, p. 29] or “fault guessing” if it is a “fault-based technique” [13, p. 4-9] that “anticipate[s] the most plausible faults in each SUT” [11, p. 5-13]. One (or both) of these proposed terms may be useful in tandem with “error guessing”, which would focus on errors as traditionally defined; this would be a subapproach of error-based testing (implied by [10, p. 399]).
- Similarly, “fault seeding” is not a synonym of “error seeding” as claimed by [14, p. 165] and [10, p. 427]. The term “error seeding”, used by [7, p. 34], [14, p. 165], and [10, p. 427], should be abandoned in favour of “fault seeding”, as it is defined as the “process of intentionally adding known faults to those already in a computer program ... [to] estimat[e] the number of faults remaining” [14, p. 165] based on the ratio between the number of new faults and the number of introduced faults that were discovered [10, p. 427].
- The distinctions between development testing [14,

p. 136], developmental testing [7, p. 30], and developer testing [7, p. 39], [26, p. 11] are unclear and seem miniscule.

- The terms “acceleration tolerance testing” and “acoustic tolerance testing” seem to only refer to software testing in [7, p. 56]; elsewhere, they seem to refer to testing the acoustic tolerance of rats [78] or the acceleration tolerance of astronauts [79, p. 11], aviators [80, pp. 27, 42], or catalysts [81, p. 1463], which don’t exactly seem relevant...
- “Orthogonal array testing” [11, pp. 5-1, 5-11] and “operational acceptance testing” [7, p. 30] have the same acronym (“OAT”).
- “Installability testing” is given as a test type [16, p. 22], [29, p. 38], [14, p. 228], while “installation testing” is given as a test level [10, p. 439]. Since “installation testing” is not given as an example of a test level throughout the sources that describe them (see Section III-B1), it is likely that the term “installability testing” with all its related information should be used instead.

6) Citation Discrepancies: Sometimes a document cites another for a piece of information that does not appear! For example, [40, p. 184] claims that [70] defines “prime path coverage”, but it does not.

C. Functional Testing

“Functional testing” is described alongside many other, likely related, terms. This leads to confusion about what distinguishes these terms, as shown by the following five:

1) Specification-based Testing: This is defined as “testing in which the principal test basis is the external inputs and outputs of the test item” [16, p. 9]. This agrees with a definition of “functional testing”: “testing that ... focuses solely on the outputs generated in response to selected inputs and execution conditions” [14, p. 196]. Notably, [14] lists both as synonyms of “black-box testing” [pp. 431, 196, respectively], despite them sometimes being defined separately. For example, the International Software Testing Qualifications Board (ISTQB) defines “specification-based testing” as “testing based on an analysis of the specification of the component or system” and “functional testing” as “testing performed to evaluate if a component or system satisfies functional requirements” [17]. Overall, specification-based testing [16, pp. 2-4, 6-9, 22] is a test design technique used to “derive corresponding test cases” [16, p. 11] from “selected inputs and execution conditions” [14, p. 196].

2) Correctness Testing: Reference [11, p. 5-7] says “test cases can be designed to check that the functional specifications are correctly implemented, which is variously referred to in the literature as conformance testing, correctness testing or functional testing”; this mirrors previous definitions of “functional testing” [16, p. 21], [14, p. 196] but groups it with “correctness testing”. Since

“correctness” is a software quality [14, p. 104], [11, p. 3-13] which is what defines a “test type” [16, p. 15], it seems consistent to label “functional testing” as a “test type” [16, pp. 15, 20, 22], [29, pp. 7, 38, Tab. A.1], [30, p. 4]. However, this conflicts with its categorization as a “technique” if considered a synonym of Specification-based Testing. Additionally, “correctness testing” is listed separately from “functionality testing” by [7, p. 53].

3) Conformance Testing: Testing that ensures “that the functional specifications are correctly implemented”, and can be called “conformance testing” or “functional testing” [11, p. 5-7]. “Conformance testing” is later defined as testing used “to verify that the SUT conforms to standards, rules, specifications, requirements, design, processes, or practices” [11, p. 5-7]. This definition seems to be a superset of testing methods mentioned earlier as the latter includes “standards, rules, requirements, design, processes, ... [and]” practices in addition to specifications!

A complicating factor is that “compliance testing” is also (plausibly) given as a synonym of “conformance testing” [43, p. 43]. However, “conformance testing” can also be defined as testing that evaluates the degree to which “results ... fall within the limits that define acceptable variation for a quality requirement” [14, p. 93], which seems to describe something different.

4) Functional Suitability Testing: Procedure testing is called a “type of functional suitability testing” [16, p. 7] but no definition of that term is given. “Functional suitability” is the “capability of a product to provide functions that meet stated and implied needs of intended users when it is used under specified conditions”, including meeting “the functional specification” [33]. This seems to align with the definition of “functional testing” as related to “black-box/specification-based testing”. “Functional correctness”, a child of “functional suitability”, is the “capability of a product to provide accurate results when used by intended users” [33] and seems to align with the quality/ies that would be tested by “correctness” testing.

5) Functionality Testing: “Functionality” is defined as the “capabilities of the various ... features provided by a product” [14, p. 196] and is said to be a synonym of “functional suitability” [17], although it seems like it should really be a synonym of “functional completeness” based on [33], which would make “functional suitability” a subapproach. Its associated test type is implied to be a subapproach of build verification testing [17] and made distinct from “functional testing” [26, Tab. 2]. “Functionality testing” is listed separately from “correctness testing” by [7, p. 53].

D. Recovery Testing

“Recovery testing” is “testing ... aimed at verifying software restart capabilities after a system crash or other disaster” [11, p. 5-9] including “recover[ing] the data directly affected and re-establish[ing] the desired state of the system” [33] (similar in [11, p. 7-10]) so that the system

“can perform required functions” [14, p. 370]. It is also called “recoverability testing” [43, p. 47] and potentially “restart & recovery (testing)” [26, Fig. 5]. The following terms, along with “recovery testing” itself [16, p. 22] are all classified as test types, and the relations between them can be found in Figure 2a.

- Recoverability Testing: Testing “how well a system or software can recover data during an interruption or failure” [11, p. 7-10] (similar in [33]) and “re-establish the desired state of the system” [33]. Synonym for “recovery testing” in [43, p. 47].
- Disaster/Recovery Testing serves to evaluate if a system can “return to normal operation after a hardware or software failure” [14, p. 140] or if “operation of the test item can be transferred to a different operating site and ... be transferred back again once the failure has been resolved” [29, p. 37]. These two definitions seem to describe different aspects of the system, where the first is intrinsic to the hardware/software and the second might not be.
- Backup and Recovery Testing “measures the degree to which system state can be restored from backup within specified parameters of time, cost, completeness, and accuracy in the event of failure” [31, p. 2]. This may be what is meant by “recovery testing” in the context of performance-related testing and seems to correspond to the definition of “disaster/recovery testing” in [14, p. 140].
- Backup/Recovery Testing: Testing that determines the ability “to restor[e] from back-up memory in the event of failure, without transfer[ing] to a different operating site or back-up system” [29, p. 37]. This seems to correspond to the definition of “disaster/recovery testing” in [29, p. 37]. It is also given as a subtype of “disaster/recovery testing”, even though that tests if “operation of the test item can be transferred to a different operating site” [p. 37]. It also seems to overlap with “backup and recovery testing”, which adds confusion.
- Failover/Recovery Testing: Testing that determines the ability “to mov[e] to a back-up system in the event of failure, without transfer[ing] to a different operating site” [29, p. 37]. This is given as a subtype of “disaster/recovery testing”, even though that tests if “operation of the test item can be transferred to a different operating site” [p. 37].
- Failover Testing: Testing that “validates the SUT’s ability to manage heavy loads or unexpected failure to continue typical operations” [11, p. 5-9] by entering a “backup operational mode in which [these responsibilities] ... are assumed by a secondary system” [17]. While not explicitly related to recovery, “failover/recovery testing” also describes the idea of “failover”, and [7, p. 56] uses the term “failover and recovery testing”, which could be a synonym of both of these

terms.

E. Scalability Testing

There were three ambiguities around the term “scalability testing”, listed below. The relations between these test approaches (and other relevant ones) are shown in Figure 3a.

- 1) ISO/IEC and IEEE give “scalability testing” as a synonym of “capacity testing” [29, p. 39] while other sources differentiate between the two [7, p. 53], [45, pp. 22-23]
- 2) ISO/IEC and IEEE give the external modification of the system as part of “scalability” [29, p. 39], while [33] implies that it is limited to the system itself
- 3) The SWEBOK Guide V4’s definition of “scalability testing” [11, p. 5-9] is really a definition of usability testing!

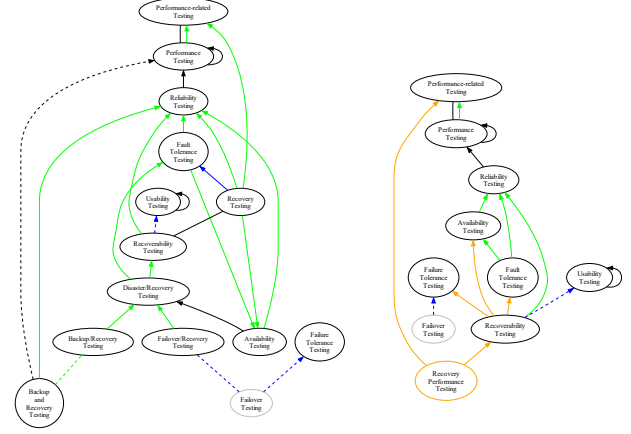
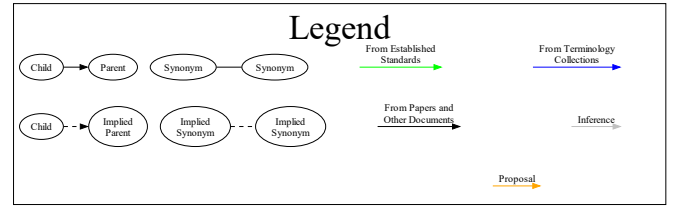
F. Compatibility Testing

“Compatibility testing” is defined as “testing that measures the degree to which a test item can function satisfactorily alongside other independent products in a shared environment (co-existence), and where necessary, exchanges information with other systems or components (interoperability)” [16, p. 3]. This definition is nonatomic as it combines the ideas of “co-existence” and “interoperability”. The term “interoperability testing” is not defined, but is used three times [16, pp. 22, 43] (although the third usage seems like it should be “portability testing”). This implies that “co-existence testing” and “interoperability testing” should be defined as their own terms, which is supported by definitions of “co-existence” and “interoperability” often being separate [14, pp. 73, 237], [17], the definition of “interoperability testing” from [14, p. 238], and the decomposition of “compatibility” into “co-existence” and “interoperability” by [33]! The “interoperability” element of “compatibility testing” is explicitly excluded by [29, p. 37], (incorrectly) implying that “compatibility testing” and “co-existence testing” are synonyms. Furthermore, the definition of “compatibility testing” in [43, p. 43] unhelpfully says “See interoperability testing”, adding another layer of confusion to the direction of their relationship.

V. Recommendations

We provide different recommendations for resolving various discrepancies (see Section IV). This was done with the goal of organizing them more logically and making them:

- 1) Atomic (e.g., disaster/recovery testing seems to have two disjoint definitions)
- 2) Straightforward (e.g., backup and recovery testing’s definition implies the idea of performance, but its name does not)
- 3) Consistent (e.g., backup/recovery testing and failover/recovery testing explicitly exclude an aspect included in its parent disaster/recovery testing)



(a) Graph of current relations. (b) Graph of proposed relations.

Fig. 2: Graphs of relations between terms related to recovery testing.

We give recommendations for the areas of recovery testing (Section V-A), scalability testing (Section V-B), and performance-related testing (Section V-C). Graphical representations (described in Section III-E) of these subsets are given in Figures 2 to 4, in which arrows representing relations between approaches are coloured based on the source tier (see Section III-A) that defines them. Any added approaches or relations are colored orange.

A. Recovery Testing

The following terms should be used in place of the current terminology to more clearly distinguish between different recovery-related test approaches. The result of the proposed terminology, along with their relations, is demonstrated in Figure 2b.

- Recoverability Testing: “Testing ... aimed at verifying software restart capabilities after a system crash or other disaster” [11, p. 5-9] including “recover[ing] the data directly affected and re-establish[ing] the desired state of the system” [33] (similar in [11, p. 7-10]) so that the system “can perform required functions” [14, p. 370]. “Recovery testing” will be a synonym, as in [43, p. 47], since it is the more prevalent term throughout various sources, although “recoverability testing” is preferred to indicate that this explicitly focuses on the ability to recover, not the performance of recovering.
- Failover Testing: Testing that “validates the SUT’s ability to manage heavy loads or unexpected failure

to continue typical operations” [11, p. 5-9] by entering a “backup operational mode in which [these responsibilities] ... are assumed by a secondary system” [17]. This will replace “failover/recovery testing”, since it is more clear, and since this is one way that a system can recover from failure, it will be a subset of “recovery testing”.

- **Transfer Recovery Testing:** Testing to evaluate if, in the case of a failure, “operation of the test item can be transferred to a different operating site and ... be transferred back again once the failure has been resolved” [29, p. 37]. This replaces the second definition of “disaster/recovery testing”, since the first is just a description of “recovery testing”, and could potentially be considered as a kind of failover testing. This may not be intrinsic to the hardware/software (e.g., may be the responsibility of humans/processes).
- **Backup Recovery Testing:** Testing that determines the ability “to restor[e] from back-up memory in the event of failure” [29, p. 37]. The qualification that this occurs “without transfer[ing] to a different operating site or back-up system” [p. 37] could be made explicit, but this is implied since it is separate from transfer recovery testing and failover testing, respectively.
- **Recovery Performance Testing:** Testing “how well a system or software can recover ... [from] an interruption or failure” [11, p. 7-10] (similar in [33]) “within specified parameters of time, cost, completeness, and accuracy” [31, p. 2]. The distinction between the performance-related elements of recovery testing seemed to be meaningful, but was not captured consistently by the literature. This will be a subset of “performance-related testing” as “recovery testing” is in [16, p. 22]. This could also be extended into testing the performance of specific elements of recovery (e.g., failover performance testing), but this be too fine-grained and may better be captured as an orthogonally derived test approach.

B. Scalability Testing

The ambiguity around scalability testing found in the literature is resolved and/or explained by other sources! [29, p. 39] gives “scalability testing” as a synonym of “capacity testing”, defined as the testing of a system’s ability to “perform under conditions that may need to be supported in the future”, which “may include assessing what level of additional resources (e.g. memory, disk capacity, network bandwidth) will be required to support anticipated future loads”. This focus on “the future” is supported by [17], which defines “scalability” as “the degree to which a component or system can be adjusted for changing capacity”. In contrast, capacity testing focuses on the system’s present state, evaluating the “capability of a product to meet requirements for the maximum limits of a product parameter”, such as the number of concurrent users, transaction throughput, or database size

[33]. Because of this nuance, it makes more sense to consider these terms separate and not synonyms, as done by [7, p. 53] and [45, pp. 22-23].

Unfortunately, only focusing on future capacity requirements still leaves room for ambiguity. While the previous definition of “scalability testing” includes the external modification of the system, [33] describes it as testing the “capability of a product to handle growing or shrinking workloads or to adapt its capacity to handle variability”, implying that this is done by the system itself. The potential reason for this is implied by [11, p. 5-9]’s claim that one objective of elasticity testing is “to evaluate scalability”: [33]’s notion of “scalability” likely refers more accurately to “elasticity”! This also makes sense in the context of other definitions provided by [11]:

- **Scalability:** “the software’s ability to increase and scale up on its nonfunctional requirements, such as load, number of transactions, and volume of data” [p. 5-5]. Based on this definition, scalability testing is then a subtype of load testing and volume testing, as well as potentially transaction flow testing.
- **Elasticity Testing¹⁴:** testing that “assesses the ability of the SUT ... to rapidly expand or shrink compute, memory, and storage resources without compromising the capacity to meet peak utilization” [p. 5-9]. Based on this definition, elasticity testing is then a subtype of memory management testing (with both being a subtype of resource utilization testing) and stress testing.

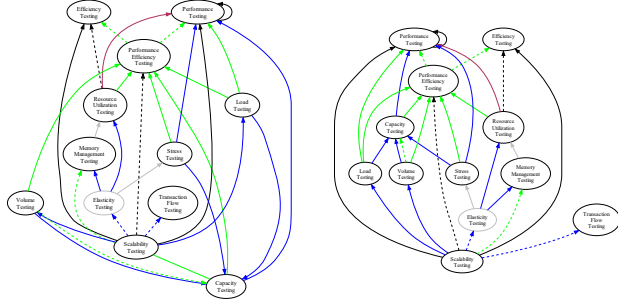
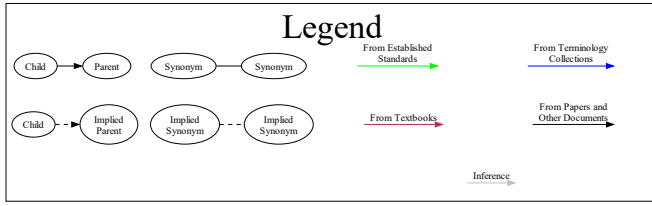
This distinction is also consistent with how the terms are used in industry: [47] says that scalability is the ability to “increase ... performance or efficiency as demand increases over time”, while elasticity allows a system to “tackle changes in the workload [that] occur for a short period”.

To make things even more confusing, the SWEBOK Guide V4 says “scalability testing evaluates the capability to use and learn the system and the user documentation” and “focuses on the system’s effectiveness in supporting user tasks and the ability to recover from user errors” [11, p. 5-9]. This seems to define “usability testing” with elements of functional and recovery testing, which is completely separate from the definitions of “scalability”, “capacity”, and “elasticity testing”! This definition should simply be disregarded, since it is inconsistent with the rest of the literature. The removal of the previous two synonym relations is demonstrated in Figure 3b.

C. Performance(-related) Testing

“Performance testing” is defined as testing “conducted to evaluate the degree to which a test item accomplishes its designated functions” [14, p. 320], [16, p. 7] (similar in [29, pp. 38-39], [68, p. 1187]). It does this by “measuring the performance metrics” [68, p. 1187] (similar in [17]) (such as

¹⁴While this definition seems correct, it only cites a single source that doesn’t contain the words “elasticity” or “elastic”!



(a) Graph of current relations. (b) Graph of proposed relations.

Fig. 3: Graphs of relations between terms related to scalability testing.

the “system’s capacity for growth” [44, p. 23]), “detecting the functional problems appearing under certain execution conditions” [68, p. 1187], and “detecting violations of non-functional requirements under expected and stress conditions” [68, p. 1187] (similar in [11, p. 5-9]). It is performed either ...

- 1) “within given constraints of time and other resources” [14, p. 320], [16, p. 7] (similar in [68, p. 1187]), or
- 2) “under a ‘typical’ load” [29, p. 39].

It is listed as a subset of performance-related testing, which is defined as testing “to determine whether a test item performs as required when it is placed under various types and sizes of ‘load’” [29, p. 38], along with other approaches like load and capacity testing [16, p. 22]. Note that “performance, load and stress testing might considerably overlap in many areas” [68, p. 1187]. In contrast, [11, p. 5-9] gives “capacity and response time” as examples of “performance characteristics” that performance testing would seek to “assess”, which seems to imply that these are subapproaches to performance testing instead. This is consistent with how some sources treat “performance testing” and “performance-related testing” as synonyms [11, p. 5-9], [68, p. 1187], as noted in Section IV-B2. This makes sense because of how general the concept of “performance” is; most definitions of “performance testing” seem to treat it as a category of tests.

However, it seems more consistent to infer that the definition of “performance-related testing” is the more general one often assigned to “performance testing” performed “within given constraints of time and other resources” [14, p. 320], [16, p. 7] (similar in [68, p. 1187]), and “performance testing” is a subapproach of this performed “under a ‘typical’ load” [29, p. 39]. This has other implications for

relations between these types of testing; for example, “load testing” usually occurs “between anticipated conditions of low, typical, and peak usage” [14, p. 253], [17], [16, p. 5], [29, p. 39], so it is a child of “performance-related testing” and a parent of “performance testing”.

After these changes, some finishing touches remain. The “self-loops” mentioned in Section IV-B3 provide no new information and can be removed. Similarly, the term “soak testing” can be removed. Since it is given as a synonym to both “endurance testing” and “reliability testing” (see Section IV-B2), it makes sense to just use these terms instead of one that is potentially ambiguous. These changes (along with those from Sections V-A and V-B made implicitly) result in the relations shown in Figure 4.

VI. Conclusion

While a good starting point, the current literature on software testing has much room to grow. The many discrepancies create unnecessary barriers to software testing. While there is merit to allowing the state-of-the-practice terminology to descriptively guide how terminology is used, there may be a need to prescriptively structure terminology to intentionally differentiate between and organize various test approaches. Future work in this area will continue to investigate the current use of terminology, in particular Undefined Terms, determine if IEEE’s current Approach Categories are sufficient, and rationalize the definitions of and relations between terms.

Acknowledgment

ChatGPT was used to help generate supplementary Python code for constructing graphs and generating \LaTeX code, including regex. ChatGPT and GitHub Copilot were both used for assistance with \LaTeX formatting. ChatGPT and ProWritingAid were both used for proofreading. Jason Balaci’s McMaster thesis template provided many helper \LaTeX functions. Finally, Drs. Spencer Smith and Jacques Carette have been great supervisors in the past and have, both then and now, provided me with valuable guidance and feedback.

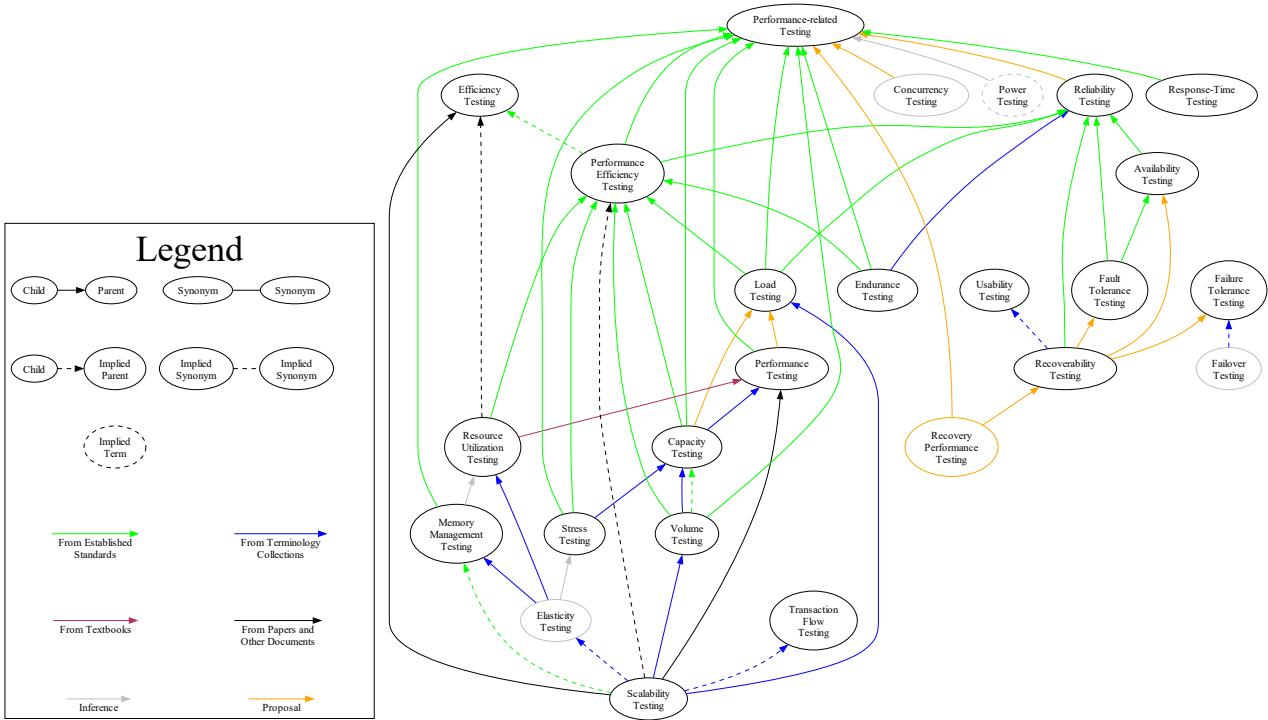


Fig. 4: Proposed relations between rationalized “performance-related testing” terms.

References

- [1] J. Peters and W. Pedrycz, *Software Engineering: An Engineering Approach*, ser. Worldwide series in computer science. John Wiley & Sons, Ltd., 2000.
- [2] P. Naur and B. Randell, “Software Engineering: Report on a conference sponsored by the NATO Science Committee, Garmisch, Germany, 7th to 11th October 1968,” Brussels, Belgium: Scientific Affairs Division, NATO, Jan. 1969. [Online]. Available: <http://homepages.cs.ncl.ac.uk/brian.randell/NATO/nato1968.PDF>
- [3] R. M. McClure, “Introduction,” Jul. 2001. [Online]. Available: <http://homepages.cs.ncl.ac.uk/brian.randell/NATO/Introduction.html>
- [4] C. Kaner, J. Bach, and B. Pettichord, *Lessons Learned in Software Testing: A Context-Driven Approach*. John Wiley & Sons, Dec. 2011. [Online]. Available: <https://www.wiley.com/en-ca/Lessons+Learned+in+Software+Testing%3A+A+Context-Driven+Approach-p-9780471081128>
- [5] G. Tebes, L. Olsina, D. Peppino, and P. Becker, “TestTDO: A Top-Domain Software Testing Ontology,” Curitiba, Brazil, May 2020, pp. 364–377.
- [6] E. Souza, R. Falbo, and N. Vijaykumar, “ROoST: Reference Ontology on Software Testing,” *Applied Ontology*, vol. 12, pp. 1–32, Mar. 2017.
- [7] D. G. Firesmith, “A Taxonomy of Testing Types,” Pittsburgh, PA, USA, 2015. [Online]. Available: <https://apps.dtic.mil/sti/pdfs/AD1147163.pdf>
- [8] M. Unterkalmsteiner, R. Feldt, and T. Gorschek, “A Taxonomy for Requirements Engineering and Software Test Alignment,” *ACM Transactions on Software Engineering and Methodology*, vol. 23, no. 2, pp. 1–38, Mar. 2014, arXiv:2307.12477 [cs]. [Online]. Available: <http://arxiv.org/abs/2307.12477>
- [9] ISO/IEC and IEEE, “ISO/IEC/IEEE International Standard - Systems and software engineering–Vocabulary,” ISO/IEC/IEEE 24765:2010(E), Dec. 2010.
- [10] H. van Vliet, *Software Engineering: Principles and Practice*, 2nd ed. Chichester, England: John Wiley & Sons, Ltd., 2000.
- [11] H. Washizaki, Ed., *Guide to the Software Engineering Body of Knowledge*, Version 4.0, Jan. 2024. [Online]. Available: <https://waseda.app.box.com/v/SWEBOK4-book>
- [12] ISO/IEC and IEEE, “ISO/IEC/IEEE International Standard - Systems and software engineering –Systems and software assurance –Part 1: Concepts and vocabulary,” ISO/IEC/IEEE 15026-1:2019, Mar. 2019.
- [13] P. Bourque and R. E. Fairley, Eds., *Guide to the Software Engineering Body of Knowledge*, Version 3.0. Washington, DC, USA: IEEE Computer Society Press, 2014. [Online]. Available: www.swebok.org
- [14] ISO/IEC and IEEE, “ISO/IEC/IEEE International Standard - Systems and software engineering–Vocabulary,” ISO/IEC/IEEE 24765:2017(E), Sep. 2017.
- [15] R. Patton, *Software Testing*, 2nd ed. Indianapolis, IN, USA: Sams Publishing, 2006.
- [16] ISO/IEC and IEEE, “ISO/IEC/IEEE International Standard - Systems and software engineering –Software testing –Part 1: General concepts,” ISO/IEC/IEEE 29119-1:2022(E), Jan. 2022.
- [17] M. Hamburg and G. Mogyorodi, editors, “ISTQB Glossary, v4.3,” 2024. [Online]. Available: https://glossary.istqb.org/en_US/search
- [18] N. E. Fenton and S. L. Pfleeger, *Software Metrics: A Rigorous & Practical Approach*, 2nd ed. Boston, MA, USA: PWS Publishing Company, 1997.
- [19] ISO, “ISO 21384-2:2021 - Unmanned aircraft systems –Part 2: UAS components,” ISO 21384-2:2021, Dec. 2021. [Online]. Available: <https://www.iso.org/obp/ui#iso:std:iso:21384-2:ed-1:v1:en>
- [20] C. Zhou, Q. Yu, and L. Wang, “Investigation of the Risk of Electromagnetic Security on Computer Systems,” *International Journal of Computer and Electrical Engineering*, vol. 4, no. 1, p. 92, Feb. 2012, publisher: IACSIT Press. [Online]. Available: <http://ijcee.org/papers/457-JE504.pdf>

- [21] R. Mandl, "Orthogonal Latin squares: an application of experiment design to compiler testing," *Communications of the ACM*, vol. 28, no. 10, pp. 1054–1058, Oct. 1985. [Online]. Available: <https://doi.org/10.1145/4372.4375>
- [22] P. Valcheva, "Orthogonal Arrays and Software Testing," in *3rd International Conference on Application of Information and Communication Technology and Statistics in Economy and Education*, D. G. Velez, Ed., vol. 200. Sofia, Bulgaria: University of National and World Economy, Dec. 2013, pp. 467–473. [Online]. Available: <https://icaictsee-2013.unwe.bg/proceedings/ICAICTSEE-2013.pdf>
- [23] H. Yu, C. Y. Chung, and K. P. Wong, "Robust Transmission Network Expansion Planning Method With Taguchi's Orthogonal Array Testing," *IEEE Transactions on Power Systems*, vol. 26, no. 3, pp. 1573–1580, Aug. 2011. [Online]. Available: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=5620950>
- [24] K.-L. Tsui, "An Overview of Taguchi Method and Newly Developed Statistical Methods for Robust Design," *IIE Transactions*, vol. 24, no. 5, pp. 44–57, May 2007, publisher: Taylor & Francis. [Online]. Available: <https://doi.org/10.1080/07408179208964244>
- [25] P. Ammann and J. Offutt, *Introduction to Software Testing*, 2nd ed. Cambridge, United Kingdom: Cambridge University Press, 2017. [Online]. Available: <https://eopcw.com/find/downloadFiles/11>
- [26] P. Gerrard, "Risk-based E-business Testing - Part 1: Risks and Test Strategy," *Systeme Evolutif*, London, UK, Tech. Rep., 2000. [Online]. Available: https://www.agileconnection.com/sites/default/files/article/file/2013/XUS129342file1_0.pdf
- [27] T. P. Johnson, "Snowball Sampling: Introduction," in *Wiley StatsRef: Statistics Reference Online*. John Wiley & Sons, Ltd, 2014, eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781118445112.stat05720>. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat05720>
- [28] P. Gerrard and N. Thompson, *Risk-based E-business Testing*, ser. Artech House computing library. Norwood, MA, USA: Artech House, 2002. [Online]. Available: <https://books.google.ca/books?id=54UKereAdJ4C>
- [29] ISO/IEC and IEEE, "ISO/IEC/IEEE International Standard - Software and systems engineering -Software testing -Part 4: Test techniques," ISO/IEC/IEEE 29119-4:2021(E), Oct. 2021.
- [30] —, "ISO/IEC/IEEE International Standard - Software and systems engineering -Software testing -Part 5: Keyword-Driven Testing," ISO/IEC/IEEE 29119-5:2016, Nov. 2016.
- [31] —, "ISO/IEC/IEEE International Standard - Systems and software engineering -Software testing -Part 1: General concepts," ISO/IEC/IEEE 29119-1:2013, Sep. 2013.
- [32] IEEE, "IEEE Standard for System and Software Verification and Validation," IEEE Std 1012-2012 (Revision of IEEE Std 1012-2004), 2012.
- [33] ISO/IEC, "ISO/IEC 25010:2023 - Systems and software engineering -Systems and software Quality Requirements and Evaluation (SQuARE) -Product quality model," ISO/IEC 25010:2023, Nov. 2023. [Online]. Available: <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-2:v1:en>
- [34] —, "ISO/IEC 25019:2023 - Systems and software engineering -Systems and software Quality Requirements and Evaluation (SQuARE) -Quality-in-use model," ISO/IEC 25019:2023, Nov. 2023. [Online]. Available: <https://www.iso.org/obp/ui/en/#iso:std:iso-iec:25019:ed-1:v1:en>
- [35] —, "ISO/IEC TS 20540:2018 - Information technology - Security techniques -Testing cryptographic modules in their operational environment," ISO/IEC TS 20540:2018, May 2018. [Online]. Available: <https://www.iso.org/obp/ui#iso:std:iso-iec:ts:20540:ed-1:v1:en>
- [36] —, "ISO/IEC 2382:2015 - Information technology -Vocabulary," ISO/IEC 2382:2015, May 2015. [Online]. Available: <https://www.iso.org/obp/ui/en/#iso:std:iso-iec:2382:ed-1:v2:en>
- [37] —, "ISO/IEC 25010:2011 - Systems and software engineering -Systems and software Quality Requirements and Evaluation (SQuARE) -System and software quality models," ISO/IEC 25010:2011, Mar. 2011. [Online]. Available: <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-1:v1:en>
- [38] ISO, "ISO 28881:2022 - Machine tools -Safety -Electrical discharge machines," ISO 28881:2022, Apr. 2022. [Online]. Available: <https://www.iso.org/obp/ui#iso:std:iso:28881:ed-2:v1:en>
- [39] —, "ISO 13849-1:2015 - Safety of machinery -Safety-related parts of control systems -Part 1: General principles for design," ISO 13849-1:2015, Dec. 2015. [Online]. Available: <https://www.iso.org/obp/ui#iso:std:iso:13849:-1:ed-3:v1:en>
- [40] S. Dogan, A. Betin-Can, and V. Garousi, "Web application testing: A systematic literature review," *Journal of Systems and Software*, vol. 91, pp. 174–201, 2014. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0164121214000223>
- [41] A. Dennis, B. H. Wixom, and R. M. Roth, *System Analysis and Design*, 5th ed. John Wiley & Sons, 2012. [Online]. Available: https://www.uoitc.edu.iq/images/documents/informatics-institute/Competitive_exam/Systemanalysisanddesign.pdf
- [42] W. E. Perry, *Effective Methods for Software Testing*, 3rd ed. Indianapolis, IN, USA: Wiley Publishing, Inc., 2006.
- [43] B. Kam, "Web Applications Testing," *Queen's University, Kingston, ON, Canada, Technical Report 2008-550*, Oct. 2008. [Online]. Available: <https://research.cs.queensu.ca/TechReports/Reports/2008-550.pdf>
- [44] P. Gerrard, "Risk-based E-business Testing - Part 2: Test Techniques and Tools," *Systeme Evolutif*, London, UK, Tech. Rep., 2000. [Online]. Available: wenku.uml.com.cn/document/test/EBTestingPart2.pdf
- [45] M. Bas, "Data Backup and Archiving," *Bachelor Thesis, Czech University of Life Sciences Prague, Praha-Suchdol, Czechia*, Mar. 2024. [Online]. Available: https://theses.cz/id/60licg/zaverecna_prace_Archive.pdf
- [46] LambdaTest, "What is Operational Testing: Quick Guide With Examples," 2024. [Online]. Available: <https://www.lambdatest.com/learning-hub/operational-testing>
- [47] P. Pandey, "Scalability vs Elasticity," Feb. 2023. [Online]. Available: <https://www.linkedin.com/pulse/scalability-vs-elasticity-pranav-pandey/>
- [48] Knüvener Mackert GmbH, Knüvener Mackert SPICE Guide, 7th ed. Reutlingen, Germany: Knüvener Mackert GmbH, 2022. [Online]. Available: <https://knuevenermackert.com/wp-content/uploads/2021/06/SPICE-BOOKLET-2022-05.pdf>
- [49] ChatGPT (GPT-4o), "Defect Clustering Testing," Nov. 2024. [Online]. Available: <https://chatgpt.com/share/67463dd1-d0a8-8012-937b-4a3db0824dcf>
- [50] V. Rus, S. Mohammed, and S. G. Shiva, "Automatic Clustering of Defect Reports," in *Proceedings of the Twentieth International Conference on Software Engineering & Knowledge Engineering (SEKE 2008)*. San Francisco, CA, USA: Knowledge Systems Institute Graduate School, Jul. 2008, pp. 291–296. [Online]. Available: <https://core.ac.uk/download/pdf/48606872.pdf>
- [51] M. Bajammal and A. Mesbah, "Web Canvas Testing Through Visual Inference," in *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*. Västerås, Sweden: IEEE, 2018, pp. 193–203. [Online]. Available: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8367048>
- [52] E. F. Barbosa, E. Y. Nakagawa, and J. C. Maldonado, "Towards the Establishment of an Ontology of Software Testing," vol. 6, San Francisco, CA, USA, Jan. 2006, pp. 522–525.
- [53] L. Baresi and M. Pezzè, "An Introduction to Software Testing," *Electronic Notes in Theoretical Computer Science*, vol. 148, no. 1, pp. 89–111, Feb. 2006. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S1571066106000442>
- [54] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, "The Oracle Problem in Software Testing: A Survey," *IEEE Transactions on Software Engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [55] J. Berdine, C. Calcagno, and P. W. O'Hearn, "Smallfoot: Modular Automatic Assertion Checking with Separation Logic," in *Formal Methods for Components and Objects*, F. S. de Boer,

- M. M. Bonsangue, S. Graf, and W.-P. de Roever, Eds. Berlin, Heidelberg: Springer, 2006, pp. 115–137.
- [56] C. Bocchino and W. Hamilton, “Eastern Range Titan IV/Centaur-TDRSS Operational Compatibility Testing,” in *International Telemetering Conference Proceedings*. San Diego, CA, USA: International Foundation for Telemetering, Oct. 1996. [Online]. Available: https://repository.arizona.edu/bitstream/handle/10150/607608/ITC_1996_96-01-4.pdf?sequence=1&isAllowed=y
 - [57] P. Chalin, J. R. Kiniry, G. T. Leavens, and E. Poll, “Beyond Assertions: Advanced Specification and Verification with JML and ESC/Java2,” in *Formal Methods for Components and Objects*, F. S. de Boer, M. M. Bonsangue, S. Graf, and W.-P. de Roever, Eds. Berlin, Heidelberg: Springer, 2006, pp. 342–363.
 - [58] S. R. Choudhary, H. Versee, and A. Orso, “A Cross-browser Web Application Testing Tool,” in *2010 IEEE International Conference on Software Maintenance*. Timisoara, Romania: IEEE, Sep. 2010, pp. 1–6, iSSN: 1063-6773. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/5609728>
 - [59] M. Dhok and M. K. Ramanathan, “Directed Test Generation to Detect Loop Inefficiencies,” in *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering*, ser. FSE 2016. New York, NY, USA: Association for Computing Machinery, Nov. 2016, pp. 895–907. [Online]. Available: <https://dl.acm.org/doi/10.1145/2950290.2950360>
 - [60] E. Engström and K. Petersen, “Mapping software testing practice with software testing research — serp-test taxonomy,” in *2015 IEEE Eighth International Conference on Software Testing, Verification and Validation Workshops (ICSTW)*, 2015, pp. 1–4.
 - [61] P. Forsyth, T. Maguire, and R. Kuffel, “Real Time Digital Simulation for Control and Protection System Testing,” in *2004 IEEE 35th Annual Power Electronics Specialists Conference (IEEE Cat. No.04CH37551)*, vol. 1. Aachen, Germany: IEEE, 2004, pp. 329–335.
 - [62] P. Godefroid and D. Luchaup, “Automatic Partial Loop Summarization in Dynamic Test Generation,” in *Proceedings of the 2011 International Symposium on Software Testing and Analysis*, ser. ISSTA '11. New York, NY, USA: Association for Computing Machinery, Jul. 2011, pp. 23–33. [Online]. Available: <https://dl.acm.org/doi/10.1145/2001420.2001424>
 - [63] A. Intana, M. Thongthep, P. Thepnimit, P. Saethapan, and T. Monpipat, “SYNTTest: Prototype of Syntax Test Case Generation Tool,” in *5th International Conference on Information Technology (InCIT)*. IEEE, 2020, pp. 259–264.
 - [64] C. Jard, T. Jérón, L. Tanguy, and C. Viho, “Remote testing can be as powerful as local testing,” in *Formal Methods for Protocol Engineering and Distributed Systems: Forte XII / PSTV XIX'99*, ser. IFIP Advances in Information and Communication Technology, J. Wu, S. T. Chanson, and Q. Gao, Eds., vol. 28. Beijing, China: Springer, Oct. 1999, pp. 25–40. [Online]. Available: https://doi.org/10.1007/978-0-387-35578-8_2
 - [65] U. Kanewala and T. Yueh Chen, “Metamorphic testing: A simple yet effective approach for testing scientific software,” *Computing in Science & Engineering*, vol. 21, no. 1, pp. 66–72, 2019.
 - [66] I. Kuļšovs, V. Arnican, G. Arnicans, and J. Borzovs, “Inventory of Testing Ideas and Structuring of Testing Terms,” vol. 1, pp. 210–227, Jan. 2013.
 - [67] S. K. Lahiri, K. L. McMillan, R. Sharma, and C. Hawblitzel, “Differential Assertion Checking,” in *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ser. ESEC/FSE 2013. New York, NY, USA: Association for Computing Machinery, Aug. 2013, pp. 345–355. [Online]. Available: <https://dl.acm.org/doi/10.1145/2491411.2491452>
 - [68] M. H. Moghadam, “Machine Learning-Assisted Performance Testing,” in *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, ser. ESEC/FSE 2019. New York, NY, USA: Association for Computing Machinery, 2019, pp. 1187–1189. [Online]. Available: <https://doi.org/10.1145/3338906.3342484>
 - [69] S. Preuß, H.-C. Lapp, and H.-M. Hanisch, “Closed-loop System Modeling, Validation, and Verification,” in *Proceedings of 2012 IEEE 17th International Conference on Emerging Technologies & Factory Automation (ETFA 2012)*. Krakow, Poland: IEEE, 2012, pp. 1–8. [Online]. Available: <https://ieeexplore.ieee.org/abstract/document/6489679>
 - [70] K. Sakamoto, K. Tomohiro, D. Hamura, H. Washizaki, and Y. Fukazawa, “POGen: A Test Code Generator Based on Template Variable Coverage in Gray-Box Integration Testing for Web Applications,” in *Fundamental Approaches to Software Engineering*, V. Cortellessa and D. Varró, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, Mar. 2013, pp. 343–358. [Online]. Available: https://link.springer.com/chapter/10.1007/978-3-642-37057-1_25
 - [71] R. S. Sangwan and P. A. LaPlante, “Test-Driven Development in Large Projects,” *IT Professional*, vol. 8, no. 5, pp. 25–29, Oct. 2006. [Online]. Available: <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=1717338>
 - [72] S. Sharma, K. Panwar, and R. Garg, “Decision Making Approach for Ranking of Software Testing Techniques Using Euclidean Distance Based Approach,” *International Journal of Advanced Research in Engineering and Technology*, vol. 12, no. 2, pp. 599–608, Feb. 2021. [Online]. Available: <https://iaeme.com/Home/issue/IJARET?Volume=12&Issue=2>
 - [73] H. Sneed and S. Göschl, “A Case Study of Testing a Distributed Internet-System,” *Software Focus*, vol. 1, pp. 15–22, Sep. 2000. [Online]. Available: https://www.researchgate.net/publication/220116945_Testing_software_for_Internet_application
 - [74] M. Dominguez-Pumar, J. M. Olm, L. Kowalski, and V. Jimenez, “Open loop testing for optimizing the closed loop operation of chemical systems,” *Computers & Chemical Engineering*, vol. 135, p. 106737, 2020. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/S0098135419312736>
 - [75] B. J. Pierre, F. Wilches-Bernal, D. A. Schoenwald, R. T. Elliott, J. C. Neely, R. H. Byrne, and D. J. Trudnowski, “Open-loop testing results for the pacific DC intertie wide area damping controller,” in *2017 IEEE Manchester PowerTech*, 2017, pp. 1–6.
 - [76] D. Trudnowski, B. Pierre, F. Wilches-Bernal, D. Schoenwald, R. Elliott, J. Neely, R. Byrne, and D. Kosterev, “Initial closed-loop testing results for the pacific DC intertie wide area damping controller,” in *2017 IEEE Power & Energy Society General Meeting*, 2017, pp. 1–5.
 - [77] W. Goralski, “xDSL loop qualification and testing,” *IEEE Communications Magazine*, vol. 37, no. 5, pp. 79–83, 1999.
 - [78] D. C. Holley, G. D. Mele, and S. Naidu, “NASA Rat Acoustic Tolerance Test 1994-1995: 8 kHz, 16 kHz, 32 kHz Experiments,” San Jose State University, San Jose, CA, USA, Tech. Rep. NASA-CR-202117, Jan. 1996. [Online]. Available: <https://ntrs.nasa.gov/api/citations/19960047530/downloads/19960047530.pdf>
 - [79] V. V. Morgun, L. I. Voronin, R. R. Kaspransky, S. L. Pool, M. R. Barratt, and O. L. Novinkov, “The Russian-US Experience with Development Joint Medical Support Procedures for Before and After Long-Duration Space Flights,” NASA, Houston, TX, USA, Tech. Rep., 1999. [Online]. Available: <https://ntrs.nasa.gov/api/citations/20000085877/downloads/20000085877.pdf>
 - [80] R. B. Howe and R. Johnson, “Research Protocol for the Evaluation of Medical Waiver Requirements for the Use of Lisinopril in USAF Aircrew,” Air Force Materiel Command, Brooks Air Force Base, TX, USA, Interim Technical Report AL/AO-TR-1995-0116, Nov. 1995. [Online]. Available: <https://apps.dtic.mil/sti/tr/pdf/ADA303379.pdf>
 - [81] D. Liu, S. Tian, Y. Zhang, C. Hu, H. Liu, D. Chen, L. Xu, and J. Yang, “Ultrafine SnPd nanoalloys promise high-efficiency electrocatalysis for ethanol oxidation and oxygen reduction,” *ACS Applied Energy Materials*, vol. 6, no. 3, pp. 1459–1466, Jan. 2023, publisher: ACS Publications. [Online]. Available: https://pubs.acs.org/doi/pdf/10.1021/acsaem.2c03355?casa_token=ItHfKxeQNbsAAAAA:8zEdU5hi2HfHsSony3ku-lbH902jkHpA-JZw8jleODzUvFtSdQRdbYhmVq47