

Todo list

■ Important: Lay abstract.	iii
■ Important: Replace reading notes.	xv
■ OG ISO/IEC, 2005	2
■ OG IEEE, 1996	3
■ OG 2005	3
■ See #22	4
■ See #155	7
■ See #176	7
■ OG Myers 1976	8
■ OG ISO/IEC 2014	8
■ OG Black, 2009	8
■ See #21	9
■ See #44, #119, and #39	9
■ See #119	9
■ See #52	9
■ Q #1: Is this precise enough?	11
■ Present tense in footnote?	12
■ See #155	12
■ See #155	12
■ Mention detecting other cycles as flaws as Future Work	13
■ See #140	14
■ See #140 and #151	14
■ See #137 and #138	14
■ <i>Later:</i> Ensure this is up to date	15
■ See #176	15
■ See #89	19
■ Q #2: Better name for this?	21
■ See #55	21
■ Q #3: Should glossary headers be capitalized?	22
■ Q #4: Is my justification for inferring “Type” sufficient?	24
■ Q #5: Does this make sense? I tried rewording it for brevity/clarity	25
■ <i>Later:</i> Ensure this is up to date	25
■ See #137 and #138	27
■ Q #6: Is this clear enough?	27
■ See #157	27
■ See #184	28

	<i>Later</i> : Ensure this is up to date	28
	See #57 , #81 , #88 , and #125	28
	Q #7 : These issue refs might be useful in our final documents.	28
	See #83	34
	See #83 and #176	36
	<i>Later</i> : Ensure this is up to date	36
	Q #8 : Is this too strong of a synonym for “ground truth”?	43
	OG Hetzel88	47
	See #21	48
	<i>Later</i> : Ensure this is up to date	50
	find more academic sources	52
	Same label to different phantomsections; is that OK?	56
	find more academic sources	57
	See #35	59
	OG IEEE 2013	61
	Important : Fill in	63
	Important : Fill in	66
	investigate OG sources	73
	Should I include the definition of Constraints ?	73
	cite Dr. Smith	73
	add refs to ‘underlying Theory’ comment and ‘not all outputs be IMs’ comment	73
	add constraints	73
	A justification for why we decided to do this should be added	75
	add acronym?	75
	is this punctuation right?	75
	See #54	77
	OG ISO 26262	77
	see ISO 29119-11	78
	Investigate	78
	OG [11, 6]	79
	OG Halfond and Orso, 2007	79
	OG Artzi et al., 2008	79
	Investigate!	79
	Add paragraph/section number?	81
	Add example	81
	Add source(s)?	81
	Important : “Important” notes.	83
	Generic inlined notes.	83
	<i>Later</i> : TODO notes for later! For finishing touches, etc.	83
	<i>Easy</i> : Easier notes.	83
	<i>Needs time</i> : Tedious notes.	83
	Q #9 : Questions I might have?	83
	Investigate further	96
	See #41 and #44	98
	See #39 , #44 , and #28	98

■ See #63	98
■ OG Alalfi et al., 2010	98
■ OG Artzi et al., 2008	98
■ OG [19]	99
■ Q #10: Is this orthogonal? Investigate legacy testing	100
■ Q #11: Is this orthogonal? Investigate role-based testing	100
■ Q #12: Is this orthogonal? Investigate scenario-based testing	100
■ Q #13: Is this orthogonal? Investigate scenario-based testing	100
■ OG ISO/IEC, 2009	101
■ OG ISO/IEC, 2009	101
■ FIND SOURCES	105
■ OG 2012	105
■ OG Beizer	106
■ OG Hetzel88	106
■ more examples?	106
■ Are these separate approaches?	107
■ See #14	108
■ OG ISO1984	108
■ Q #14: What’s the short form of “annex”?	109
■ OG IREB Glossary	110
■ Q #15: Does this merit counting this as an Ambiguity as well as a Contradiction?	112
■ OG Beizer	112
■ Is this a def flaw?	112
■ OG [14]	112
■ Q #16: I ignore “syntax errors, runtime errors, and logical errors” (Washizaki, 2024, p. 16-13, cf. p. 18-13) since they seem to be in dif- ferent domains. Does that make sense? How should I document this?	112
■ OG ISO/IEC, 2005	112
■ Q #17: Should I be more specific?	113
■ Q #18: Should we classify this as a Synonym Flaw at all?	114
■ Q #19: Shouldn’t this be “is”, referring to “test item”?	114
■ Better way to handle/display this?	121
■ See #59	128

PUTTING SOFTWARE TESTING TERMINOLOGY TO THE TEST

PUTTING SOFTWARE TESTING TERMINOLOGY TO THE TEST

By SAMUEL CRAWFORD, B.Eng.

A THESIS
SUBMITTED TO THE DEPARTMENT OF COMPUTING AND SOFTWARE
AND THE SCHOOL OF GRADUATE STUDIES
OF MCMASTER UNIVERSITY
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
MASTER OF APPLIED SCIENCE

Master of Applied Science (2025)
(Department of Computing and Software)

McMaster University
Hamilton, Ontario

TITLE: Putting Software Testing Terminology to the Test
AUTHOR: Samuel Crawford, B.Eng.
SUPERVISOR: Dr. Carette and Dr. Smith
PAGES: xvi, 131

Lay Abstract

Important: Lay abstract.

Abstract

Despite the prevalence and importance of software testing, it lacks a standardized and consistent taxonomy, instead relying on a large body of literature with many flaws—even within individual documents! This hinders precise communication, contributing to misunderstandings when planning and performing testing. In this thesis, we explore the current state of software testing terminology by:

1. identifying established standards and prominent testing resources,
2. capturing relevant testing terms from these sources, along with their definitions and relationships (both explicit and implicit), and
3. constructing graphs to visualize and analyze these data.

This process uncovers 561 test approaches and four in-scope methods for deriving test approaches, such as those related to 77 software qualities. We also build a tool for generating graphs that illustrate relations between test approaches and track flaws captured by this tool and manually through the research process. This reveals 295 flaws, including nine terms used as synonyms to two (or more) disjoint test approaches and 18 pairs of test approaches that may either be synonyms or have a parent-child relationship. This also highlights notable confusion surrounding functional testing, operational acceptance testing, recovery testing, 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.

Acknowledgements

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. A special “thank you” to Christopher William Schankula for help with \LaTeX and various friends for discussing software testing with me and providing many of the approaches in [Section 3.5](#) (ChatGPT also provided pointers to the potential existence of some of these approaches). Finally, Dr. Spencer Smith and Dr. Jacques Carette have been great supervisors and valuable sources of guidance and feedback.

Contents

Todo list	i
Lay Abstract	iii
Abstract	iv
Acknowledgements	v
Contents	vi
List of Figures	x
List of Tables	xi
List of Source Codes	xii
List of Abbreviations and Symbols	xiii
Reading Notes	xv
Declaration of Academic Achievement	xvi
1 Introduction	1
2 Terminology	7
2.1 Test Approaches	8
2.1.1 Approach Categories	8
2.1.2 Synonym Relations	11
2.1.3 Parent-Child Relations	11
2.2 Flaws	12
2.2.1 Flaw Manifestations	12
2.2.2 Flaw Domains	12
2.2.3 Additional Flaw Terminology	14
2.3 Explicitness	15
2.4 Credibility	16
2.5 Source Tiers	17
2.5.1 Established Standards	18
2.5.2 Terminology Collections	18

2.5.3	Textbooks	18
2.5.4	Papers and Other Documents	19
3	Methodology	20
3.1	Identifying Sources	21
3.2	Identifying Relevant Terms	21
3.3	Recording Relevant Information	22
3.3.1	Recording Implicit Information	25
3.3.2	Recording Flaws	27
3.4	Undefined Terms	28
3.5	Stopping Criteria	28
4	Tools	30
4.1	Approach Relation Visualization	30
4.2	Flaw Analysis	33
4.2.1	Automated Flaw Analysis	34
4.2.2	Augmented Flaw Analysis	36
4.3	LaTeX Commands	37
5	Observed Flaws	43
5.1	Flaws by Manifestation	46
5.1.1	Mistakes	46
5.1.2	Omissions	47
5.1.3	Contradictions	47
5.1.4	Ambiguities	47
5.1.5	Overlaps	48
5.1.6	Redundancies	48
5.2	Flaws by Domain	48
5.2.1	Approach Category Flaws	48
5.2.2	Synonym Relation Flaws	50
5.2.3	Parent-Child Relation Flaws	50
5.3	Operational (Acceptance) Testing	52
5.4	Recovery Testing	53
5.5	Scalability Testing	54
5.6	Compatibility Testing	55
6	Recommendations	56
6.1	Operational (Acceptance) Testing	57
6.2	Recovery Testing	57
6.3	Scalability Testing	58
6.4	Performance(-related) Testing	59
6.5	Compatibility Testing	61
7	Threats to Validity	63
7.1	Reliability	63
7.2	Internal Validity	64

7.3	External Validity	66
7.4	Construct Validity	66
8	Future Work	67
8.1	Iterating Over Undefined Terms	67
8.2	Missing Terms	69
8.3	Detecting More Flaws	69
9	Development Process	71
9.1	Improvements to Manual Test Code	72
9.1.1	Testing with Mocks	72
9.2	The Use of Assertions in Code	73
9.3	Generating Requirements	73
10	Research	75
10.1	Existing Taxonomies, Ontologies, and the State of Practice	75
10.2	Definitions	76
10.2.1	Documentation	76
10.3	General Testing Notes	77
10.3.1	Test Oracles	77
10.3.2	Generating Test Cases	78
10.4	Examples of Metamorphic Relations	80
10.5	Roadblocks to Testing	80
10.5.1	Roadblocks to Testing Scientific Software	81
11	Extras	82
11.1	Writing Directives	82
11.2	HREFs	82
11.3	Pseudocode Code Snippets	83
11.4	TODOs	83
	Bibliography	84
A	Detailed Scope Analysis	95
A.1	Hardware Testing	95
A.2	V&V of Other Artifacts	96
A.3	Static Testing	97
A.4	Vague Terminology	98
A.5	Language-specific Approaches	98
A.6	Orthogonally Derived Approaches	99
B	Sources by Tier	102
B.1	Established Standards	102
B.2	Terminology Collections	102
B.3	Textbooks	102
B.4	Papers and Other Documents	102

C	Full Lists of Flaws	103
C.1	Full Lists of Flaws by Manifestation	103
C.1.1	Full List of Mistakes	103
C.1.2	Full List of Omissions	106
C.1.3	Full List of Contradictions	107
C.1.4	Full List of Ambiguities	112
C.1.5	Full List of Overlaps	114
C.1.6	Full List of Redundancies	115
C.2	Full Lists of Flaws by Domain	116
C.2.1	Multiple Categorizations	116
C.2.2	Intransitive Synonyms	121
C.2.3	Synonym and Parent-Child Overlaps	122
C.3	Inferred Flaws	126
C.3.1	Inferred Multiple Categorizations	126
C.3.2	Inferred Intransitive Synonyms	127
C.3.3	Inferred Parent Flaws	127
D	Code Snippets	129

List of Figures

2.1	Summary of how many sources comprise each source tier.	17
3.1	Procedure for recording test approaches in our glossary.	23
3.2	ISO/IEC and IEEE (2022, p. 1)’s glossary entry for “A/B testing”.	24
3.3	A/B testing’s inclusion in Firesmith (2015, p. 58)’s taxonomy.	24
4.1	Example generated visualizations of parent-child relations.	31
4.2	Example generated visualizations of synonym relations.	32
4.3	The sets of authors of established standards.	35
4.4	Example generated visualizations containing flaws.	36
5.1	Identified flaws by the source tier responsible. Some bars are omitted as they correspond to comparisons we do not make; see Section 4.2.	44
5.2	Visually meaningful synonym relations given explicitly by the literature.	51
5.3	Pairs of test approaches with a parent-child <i>and</i> synonym relation given explicitly by the literature.	52
6.1	Visualizations of relations between terms related to recovery testing.	58
6.2	Visualizations of relations between terms related to scalability testing.	60
6.3	Visualization of proposed relations between terms related to performance-related testing.	61
8.1	Breakdown of how many test approaches are undefined.	68

List of Tables

1.1	Selected entries from our test approach glossary with “Notes” column excluded for brevity.	6
2.1	Categories of test approaches given by ISO/IEC and IEEE.	10
2.2	Observed flaw manifestations.	13
2.3	Observed flaw domains.	14
3.1	Breakdown of keywords used for recording and analyzing implicitness.	26
4.1	Example glossary entries demonstrating how we track parent-child relations.	31
4.2	Example glossary entries demonstrating how we track synonym relations.	32
4.3	L ^A T _E X macros for calculated values.	37
4.4	L ^A T _E X macros for referencing well-defined sections.	40
4.5	L ^A T _E X macros for reused text.	41
5.1	Breakdown of identified flaws by manifestation and source tier.	45
5.2	Breakdown of identified flaws by domain and source tier.	45
5.3	Different kinds of mistakes found in the literature.	46
5.4	Different kinds of ambiguities found in the literature.	48
5.5	Summary of pairs of categories assigned to a test approach.	49
7.1	Categories of testing given by other sources.	65
C.1	Test approaches with more than one category.	116
C.2	Pairs of test approaches with a parent-child <i>and</i> synonym relation.	123
C.3	Test approaches inferred to have more than one category.	126

List of Source Codes

11.1 Pseudocode: exWD	82
11.2 Pseudocode: exPHref	83
D.1 Tests for main with an invalid input file	129
D.2 Projectile’s choice for constraint violation behaviour in code	130
D.3 Projectile’s manually created input verification requirement	130
D.4 “MultiDefinitions” (MultiDefn) Definition	130
D.5 Pseudocode: Broken QuantityDict Chunk Retriever	130

List of Abbreviations and Symbols

c-use/C-use	Computation data Use
CT	Continuous Testing
DAC	Differential Assertion Checking
DbIPend	Double Pendulum
DOM	Document Object Model
du-path/DU-path	Definition-Use path
EMSEC	EManations SECurity
FIST	Fault Injection Security Tool
GlassBR	Glass BReaking
HREF	Hypertext REference
IEC	International Electrotechnical Commission
ISTQB	International Software Testing Qualifications Board
LCSAJ	Linear Code Sequence and Jump
ML	Machine Learning
MR	Metamorphic Relation
OAT	Operational Acceptance/Orthogonal Array Testing
OT	Operational Testing
p-use/P-use	Predicate data Use
PAR	Product Anomaly Report
PDF	Portable Document Format
PIR	Product Incident Report
QAI	Quality Assurance Institute
RAC	Runtime Assertion Checking
RQ	Research Question
SglPend	Single Pendulum
SoS	System of Systems
SRS	Software Requirements Specification
SSP	Slope Stability analysis Program

SUT	System Under Test
SV	Software Verification
SWEBOK Guide	Guide to the SoftWare Engineering Body Of Knowledge
SWHS	Solar Water Heating System
TOAT	Taguchi's Orthogonal Array Testing
V&V	Verification and Validation

Reading Notes

Before reading this thesis, I encourage you to read through these notes, keeping them in mind while reading.

- The source code of this thesis is [publicly available](#).
- This thesis template is primarily intended for usage by the computer science community¹. However, anyone is free to use it.
- I’ve tried my best to make this template conform to the thesis requirements as per [those set forth in 2021 by McMaster University](#). However, you should double-check that your usage of this template is compliant with whatever the “current” rules are.

Important: Replace reading notes.

¹Hence why there are some \LaTeX macros for “code” snippets.

Declaration of Academic Achievement

This research and analysis was performed by Samuel Crawford under the guidance, supervision, and recommendations of Dr. Spencer Smith and Dr. Jacques Carette. The resulting contributions are three glossaries—one for each of test approaches, software qualities, and supplementary terms (see [Section 3.3](#))—as well as the tools for data visualization and automated analysis outlined in [Chapter 4](#). These are all available on [an open-source repo](#) for independent analysis and, ideally, extension as more test approaches are discovered and documented.

Chapter 1

Introduction

As with all fields of science and technology, software development should be approached systematically and rigorously. [Peters and Pedrycz](#) claim that “to be successful, development of software systems requires an engineering approach” that is “characterized by a practical, orderly, and measured development of software” (2000, 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” ([Naur and Randell, 1969](#), 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” ([McClure, 2001](#)). 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 ([Naur and Randell, 1969](#), 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). Furthermore, “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. [Kaner et al.](#) give 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 (2011, 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” ([Peters and Pedrycz, 2000](#), 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!

We noticed this lack of a standard language for software testing while working on our own software framework, Drasil (Carette et al., 2021), with the goal of “generating all of the software artifacts for (well understood) research software”. Currently, these include Software Requirements Specifications (SRSs), READMEs, Makefiles, and code in up to six languages, depending on the specific case study (Hunt et al., 2021). To improve the quality, functionality, and maintainability of Drasil, we want to add test cases to this list (Smith, 2024; Hunt et al., 2021). This process would be a part of our “continuous integration system, [sic] so that the generated code for the case studies is automatically tested with each build” and would be a big step forward, since we currently only “test that the generated code compiles” (Smith, 2024). However, before we can include test cases as a generated artifact, the underlying domain—software testing—needs to be “well understood”, which requires a “stable knowledge base” (Hunt et al., 2021).

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. (2020) 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. (2017) prioritize organizing test approaches over defining them,
- Firesmith (2015) similarly defines relations between test approaches but not the approaches themselves, and
- Unterkalmsteiner et al. (2014) 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 2.5.1) and terminology collections (see Section 2.5.2) 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” (ISO/IEC and IEEE, 2010, p. 128; van Vliet, 2000, p. 399).
- **Fault:** “an incorrect step, process, or data definition in a computer program” (ISO/IEC and IEEE, 2010, p. 140) inserted when a developer makes an error (pp. 128, 140; Washizaki, 2024, p. 12-3; van Vliet, 2000, pp. 399–400).
- **Failure:** the inability of a system “to perform a required function or ... within previously specified limits” (ISO/IEC and IEEE, 2019a, p. 7; 2010, p. 139; similar in van Vliet, 2000, p. 400) that is “externally visible” (ISO/IEC and IEEE, 2019a, p. 7; similar in van Vliet, 2000, p. 400) and caused by a fault (Washizaki, 2024, p. 12-3; van Vliet, 2000, p. 400).

- **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” (ISO/IEC and IEEE, 2010, p. 96).

This distinction is sometimes important, but not always (Bourque and Fairley, 2014, 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” (Washizaki, 2024, p. 12-3), and “defect” may be used as a “generic term that can refer to either a fault (cause) or a failure (effect)” (ISO/IEC and IEEE, 2017, p. 124; 2010, p. 96). Software testers may even choose to ignore these nuances completely! Patton (2006, 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”!

These decisions are not inherently wrong, since they may be useful in certain contexts or for certain teams (see Section 2.1.2 for more detailed discussion). Problems start to arise when teams need to make these decisions in the first place. Patton (2006, p. 14) notes that “a well-known computer company spent weeks in discussion with its engineers before deciding to rename Product Anomaly Reports (PARs) to Product Incident Reports (PIRs)”, a process that required “countless dollars” and updating “all the paperwork, software, forms, and so on”. While consistency and clear terminology may have been valuable to the company, “it’s unknown if [this decision] made any difference to the programmer’s or tester’s productivity” (p. 14). A potential way to avoid similar resource sinks would be to prescribe a standard terminology. Perhaps multiple sets of terms could be designed with varying levels of specificity so a company would only have to determine which one best suits their needs.

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

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

These kinds of inconsistencies can lead to miscommunications—such as that previously mentioned by Kaner et al. (2011, p. 7)—and are prominent in the literature. ISO/IEC and IEEE categorize experience-based testing as both a test design technique and a test practice in the same figure (2022, Fig. 2)! The structure of tours can be defined as either quite general (ISO/IEC and IEEE, 2022, p. 34) or “organized around a special focus” (Hamburg and Mogyorodi, 2024). Load testing is performed with loads “between anticipated conditions of low, typical, and peak usage” (ISO/IEC and IEEE, 2022, p. 5) or loads that are as large as possible (Patton, 2006, p. 86). Alpha testing is performed by “users within the organization developing the software” (ISO/IEC and IEEE, 2017, p. 17), “a small, selected group of potential users” (Washizaki, 2024, p. 5-8), or “roles outside the development organization” conducted “in the developer’s test environment” (Hamburg and Mogyorodi, 2024). It is clear that there is a notable gap in the literature, one which we attempt to describe. 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. The main way we accomplish this is by identifying “flaws” or “inconsistencies” in the literature, or areas where there is room for improvement. We track these flaws according to both *what* information is wrong and *how* (described in more detail in Section 2.2), which allows us to analyze them more thoroughly and reproducibly.

Based on this observed gap in software testing terminology and our original motivation for this research, we only consider the component of Verification and Validation (V&V) that tests code itself. However, some test approaches are only used to testing *other* artifacts, while others can be used for both! In these cases, we only consider the subsections that focus on code. For example, reliability testing and maintainability testing can start *without* code by “measur[ing] structural attributes of representations of the software” (Fenton and Pfleeger, 1997, p. 18), but only reliability and maintainability testing performed on code *itself* is in scope of this research. This is a high-level overview of what is in scope; see Appendix A for more detailed discussion on what we include and exclude.

This document describes our process, as well as our results, in more detail. We start by documenting the 561 test approaches mentioned by 71 sources (described in Section 2.5), recording their names, categories¹, definitions, synonyms², parents³, and flaws⁴ (see Section 3.3.2) as applicable. We also record any other relevant notes, such as prerequisites, uncertainties, other sources. We follow the procedure laid out in Chapter 3 and use these Research Questions (RQs) as a guide:

1. What testing approaches do the literature describe?

¹Defined in Section 2.1.1.

²Defined in Section 2.1.2.

³Defined in Section 2.1.3.

⁴Defined in Section 2.2.

2. Are these descriptions consistent?
3. Can we systematically resolve any of these inconsistencies?

An excerpt of this recorded information (excluding other notes for brevity), is given in [Table 1.1](#). We then create tools to support our analysis of our findings ([Chapter 4](#)). Despite the amount of well understood and organized knowledge, the literature is still quite flawed ([Chapter 5](#)). This reinforces the need for a proper taxonomy! We then provide some potential solutions covering some of these flaws ([Chapter 6](#)).

Table 1.1: Selected entries from [our test approach glossary](#) with “Notes” column excluded for brevity.

Name	Approach Category	Definition	Parent(s)	Synonym(s)
A/B Testing	Practice (ISO/IEC and IEEE, 2022 , p. 22), Type (inferred from usability testing)	Testing “that allows testers to determine which of two systems or components performs better” (ISO/IEC and IEEE, 2022 , pp. 1, 36)	Statistical Testing (ISO/IEC and IEEE, 2022 , pp. 1, 36), Usability Testing (Firesmith, 2015 , p. 58)	Split-Run Testing (ISO/IEC and IEEE, 2022 , pp. 1, 36)
Big-Bang Testing	Technique (Sharma et al., 2021 , pp. 601, 603, 605–606), Level (inferred from integration testing)	“Testing in which ... [components of a system] are combined all at once into an overall system, rather than in stages” (ISO/IEC and IEEE, 2017 , p. 45)	Integration Testing (ISO/IEC and IEEE, 2017 , p. 45; Washizaki, 2024 , p. 5-7; Sharma et al., 2021 , p. 603; Kam, 2008 , p. 42; Peters and Pedrycz, 2000 , p. 488, Tab. 12.8)	Sometimes spelled without a hyphen (Peters and Pedrycz, 2000 , p. 489)
Classification Tree Method	Technique (ISO/IEC and IEEE, 2022 , p. 22; 2021c , pp. 2, 12, Fig. 2; Hamburg and Mogyorodi, 2024 ; Firesmith, 2015 , p. 47)	Testing “based on exercising classes in a classification tree” (ISO/IEC and IEEE, 2021c , p. 2)	Specification-based Testing (ISO/IEC and IEEE, 2022 , p. 22; 2021c , pp. 2, 12, Fig. 2; Hamburg and Mogyorodi, 2024 ; Firesmith, 2015 , p. 47), Model-based Testing (ISO/IEC and IEEE, 2022 , p. 13; 2021c , pp. 6, 12)	Classification Tree Technique (Hamburg and Mogyorodi, 2024), Classification Tree Testing (Firesmith, 2015 , p. 47)

Chapter 2

Terminology

Our research aims to describe the current state of software testing literature, including its flaws. Since we critique the lack of clarity, consistency, and robustness in the literature, we need to hold ourselves to a high standard in these areas when defining and using terms. For example, since we focus on how the literature describes “test approaches”, we first define this term ([Section 2.1](#)). Likewise, before we can constructively describe the flaws in the literature, we need to define what we mean by “flaw” ([Section 2.2](#)). To further prevent bias, we only use classifications and relations already implicitly present in the literature instead of inventing our own; for example, test approaches can have categories ([Section 2.1.1](#)), synonyms ([Section 2.1.2](#)), and parent-child relations ([Section 2.1.3](#)). We also observe flaws having both manifestations ([Section 2.2.1](#)) and domains ([Section 2.2.2](#)) and use these terms to refer to these implicit concepts in the literature. All of these classifications and relations follow logically from the literature and as such are technically “results” of our research, but we define them here for clarity since we use them throughout this thesis.

See [#155](#)

Since the literature is flawed, we need to be careful with what information we take at face value. We do this by tracking the nuance, or “explicitness”, of information ([Section 2.3](#)) found in sources and the “credibility” of these sources themselves ([Section 2.4](#)). We then use this heuristic of credibility to group our identified sources into “tiers” ([Section 2.5](#)). Defining these terms helps reduce the effect of our preconceptions on our analysis (or at least makes it more obvious to future researchers), as there may be other equally valid ways to analyze the literature and its flaws. To be clear: we do *not* prescribe what terminology software testers *should* use, we simply observe the terminology that the literature uses and try to use it as consistently and logically as possible.

See [#176](#)

2.1 Test Approaches

OG Myers 1976

Software testing is defined as “the process of executing a program with the intent of finding errors” (Peters and Pedrycz, 2000, p. 438). For each test, the main steps are to:

1. identify the goal(s) of the test,
2. decide on an approach,
3. develop the tests,
4. determine the expected results,
5. run the tests, and
6. compare the expected results to the actual results (p. 443).

OG ISO/IEC 2014

The end goal is to evaluate “some aspect of the system or component” based on the results of step 6 (ISO/IEC and IEEE, 2022, p. 10; 2021c, p. 6; 2017, p. 465). When this evaluation reveals errors, “the faults causing them are what can and must be removed” (Washizaki, 2024, p. 5-3).

Of course, the approach chosen in step 2 influences what kinds of test cases should be developed and executed in later steps, so it is important that test approaches are defined correctly, consistently, and unambiguously. A “test approach” is a “high-level test implementation choice” (ISO/IEC and IEEE, 2022, p. 10) used to “pick the particular test case values” (2017, p. 465) used in step 5. The only approach that can “fully” test a system (exhaustive testing) is infeasible in most non-trivial situations (2022, p. 4; Washizaki, 2024, p. 5-5; Peters and Pedrycz, 2000, pp. 439, 461; van Vliet, 2000, p. 421), so multiple approaches are needed (ISO/IEC and IEEE, 2022, p. 18) to “suitably cover any system” (p. 33). This is why this process should be repeated and there are so many test approaches described in the literature.

2.1.1 Approach Categories

OG Black, 2009

Since there are so many test approaches, it is helpful to categorize them. The literature provides many ways to do so, but we use the one given by ISO/IEC and IEEE (2022) because of its wide usage. This schema divides test approaches into levels, types, techniques, and practices (2022, Fig. 2; see Table 2.1). These categories seem to be pervasive throughout the literature, particularly “level” and “type”. For example, six non-IEEE sources also give unit testing, integration testing, system testing, and acceptance testing as examples of test levels (Washizaki, 2024, pp. 5-6 to 5-7; Hamburg and Mogyorodi, 2024; Perry, 2006, pp. 807–808; Peters and Pedrycz, 2000, pp. 443–445; Kulešovs et al., 2013, p. 218; Gerrard, 2000a, pp. 9, 13). These categories seem to be orthogonal based on their definitions and usage. For example, “a test type can be performed at a single test level or across several test levels” (ISO/IEC and IEEE, 2022, p. 15; 2021a, p. 8; 2021c, p. 7),

and test practices can be “defined ... for a specific level or type of testing” (2021b, p. 9). We may assess this assumption more rigorously in the future, but for now, it implies that one can derive a specific test approach by combining multiple test approaches from different categories; see [Appendix A.6](#) for more detailed discussion.

We loosely describe these categories based on what they specify as follows:

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

While [ISO/IEC and IEEE \(2022\)](#)’s schema includes “static testing” as a test approach category, we omit it from our scope as it seems non-orthogonal to the others and thus less helpful for grouping test approaches. Other categorization schemas (see [Section 7.2](#)) may consider static testing orthogonal, and some may consider it out-of-scope entirely (see [Appendix A.3](#))! We also introduce a “supplemental” category of “artifact” since some terms can refer to both the application of a test approach *and* the resulting document(s). Additionally, [Souza et al. \(2017, p. 3\)](#) give “test procedures”, which may be considered equivalent to “test approaches” (see [Table 7.1](#)), as examples of testing artifacts. Therefore, we do *not* consider approaches categorized as an artifact *and* another category as flaws in [Section 5.2.1](#). Finally, we also record the test category of “process” for completeness, although this seems to be a higher level classification and these approaches will likely be excluded during later analysis.

See [#21](#)

See [#44](#), [#119](#),
and [#39](#)

See [#119](#)

See [#52](#)

Table 2.1: Categories of test approaches given by ISO/IEC and IEEE.

Term	Definition	Examples
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 (ISO/IEC and IEEE, 2022, p. 12; 2021a, p. 6; 2021c, p. 6) with their “own documentation and resources” (2017, p. 469)	unit/component testing, integration testing, system testing, acceptance testing (2022, p. 12; 2021a, p. 6; 2021c, p. 6; 2017, p. 467)
Test Practice	A “conceptual framework that can be applied to ... [a] test process to facilitate testing” (2022, p. 14; 2017, p. 471; OG IEEE 2013)	scripted testing, exploratory testing, automated testing (2022, p. 20)
Test Technique ^b	A “procedure used to create or select a test model ..., identify test coverage items ..., and derive corresponding test cases” (2022, p. 11; 2021a, p. 5; similar in 2017, p. 467; OG 2013) that “generate evidence that test item requirements have been met or that defects are present in a test item” (2021c, p. vii) “typically used to achieve a required level of coverage” (2021a, p. 5)	equivalence partitioning, boundary value analysis, branch testing (2022, p. 11; 2021a, p. 5)
Test Type	“Testing that is focused on specific quality characteristics” (2022, p. 15; 2021c, p. 7; 2017, p. 473; OG IEEE 2013)	security testing, usability testing, performance testing (2022, p. 15; 2021a, p. 8; 2017, p. 473)

^a Also called “test phase” (see [Overlap 2](#)) or “test stage” (see [Contradiction 22](#)).^b Also called “test design technique” (ISO/IEC and IEEE, 2022, p. 11; 2021a, p. 5; [Hamburg and Mogyorodi, 2024](#)).

2.1.2 Synonym Relations

The same approach often has many names. For example, “specification-based testing” is also called “black-box testing” (ISO/IEC and IEEE, 2022, p. 9; 2021c, p. 8; 2017, p. 431; Washizaki, 2024, p. 5-10; Hamburg and Mogyorodi, 2024; Firesmith, 2015, pp. 46–47¹; Sakamoto et al., 2013, p. 344; van Vliet, 2000, p. 399). 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.

We can formally define the synonym relation S on the set T of terms used by the literature to describe test approaches based on how synonyms are used in natural language. S is symmetric and transitive, and although pairs of synonyms in natural language are implied to be distinct, a relation that is symmetric and transitive is provably reflexive; this implies that all terms are trivially synonyms of themselves. Since S is symmetric, transitive, *and* reflexive, it is an equivalence relation, reflecting the role of synonyms in natural language where they can be used interchangeably. While synonyms may emphasize different aspects or express mild variations, their core meaning is nevertheless the same.

2.1.3 Parent-Child Relations

Many test approaches are multi-faceted and can be “specialized” into others; for example, load testing and stress testing are some subtypes of performance-related testing (described in more detail in Section 6.4). We refer to these “specializations” as “children” or “subapproaches” of their multi-faceted “parent(s)”. This nomenclature also extends to approach categories (such as “subtype”; see Section 2.1.1 and Table 2.1) and software qualities (“subquality”; see Section 3.3).

We can formally define the parent-child relation P on the set T of terms used by the literature to describe test approaches based on directed relations between approach pairs. This relation should be irreflexive, asymmetric, and transitive, making it a strict partial order. A consequence of this is that there should be no directed cycles, although since a given child approach c may have more than one parent approach p , undirected cycles may exist.

Parent-child relations often manifest when a “well-understood” test approach p is decomposed into smaller, independently performable approaches c_1, \dots, c_n , each with its own focus or nuance. This is frequently the case for hierarchies of approaches given in the literature (ISO/IEC and IEEE, 2022, Fig. 2; 2021c, Fig. 2; Firesmith, 2015). Another way for these relations to occur is when the completion of p indicates that “sufficient testing has been done” in regards to c (van Vliet, 2000, p. 402). While this only “compares the thoroughness of test techniques, not their ability to detect faults” (p. 434), it is sufficient to justify a parent-child relation between the two approaches. These relations may also be represented as hierarchies (ISO/IEC and IEEE, 2021c, Fig. F.1; van Vliet, 2000, Fig. 13.17).

¹Firesmith (2015) excludes the hyphen, calling it “black box testing”.

Q #1: Is this precise enough?

2.2 Flaws

Ideally, software testing literature would describe test approaches correctly, completely, consistently, and modularly, but this is not the case in reality. We use the term “flaw” to refer to any instance of the literature violating these ideals, *not* to instances of *software doing the same*.² We classify flaws by both their “manifestations” (*how* information is wrong; see [Section 2.2.1](#)) and their “domains” (*what* information is wrong; see [Section 2.2.2](#)). These are orthogonal classifications, since each flaw *manifests* in a particular *domain*, which we track by assigning each flaw one “key” for each classification. We list these keys in [Tables 2.2](#) and [2.3](#), respectively, and use them to count and analyze the flaws we uncover in [Chapter 5](#). We also introduce some terms to be used when discussing flaws based on how many sources contribute to them ([Section 2.2.3](#)).

2.2.1 Flaw Manifestations

Perhaps the most obvious example of something being “wrong” with the literature is that a piece of information it presents is incorrect—“wrong” in the literal sense. However, if our standards for correctness require clarity, consistency, and robustness, then there are many ways for a flaw to manifest. This is one view we take when observing, recording, and analyzing flaws: *how* information is “wrong”. We observe the “manifestations” described in [Table 2.2](#) throughout the literature, and give each a unique key for later analysis and discussion. We list them in descending order of severity, although this is partially subjective. While some may disagree with our ranking, it is clear that information being incorrect is worse than it being repeated. Our ordering has the benefit of serving as a “flowchart” for classifying flaws. For example, if a piece of information is not intrinsically incorrect, then there are five remaining manifestation types for the flaw. Note that some flaws involve information from multiple sources (contradictions and overlaps in particular). We do not categorize these flaws as “mistakes” if finding the ground truth requires analysis that has not been performed yet.

2.2.2 Flaw Domains

Another way to categorize flaws is by *what* information is wrong, which we call the flaw’s “domain”. We describe those we observe in [Table 2.3](#), and tracking these uncovers which knowledge domains are less standardized (and should therefore be approached with more rigour) than others. We explicitly define some of these

²When picking a word to describe these violations, we wanted to avoid words that are “overloaded with too many meanings” like “error” and “fault” ([Washizaki, 2024](#), p. 12-3; see [Chapter 1](#) for more detailed discussion). A small literature review revealed that established standards (see [Section 2.5.1](#)) use the term “flaw” to refer to requirements ([ISO/IEC and IEEE, 2022](#), p. 38), design (p. 43), “system security procedures ... and internal controls” ([IEEE, 2012](#), p. 194), but only once to refer to code itself (p. 92). [Washizaki \(2024, p. 7-9\)](#) also uses the word “flaw” in a way that aligns with our goal of analyzing what is wrong with software engineering’s testing literature.

Table 2.2: Observed flaw manifestations.

Manifestation	Description	Key
Mistake	Information is incorrect	WRONG ^a
Omission	Information that should be included is not	MISS ^b
Contradiction	Information from multiple places conflicts	CONTRA
Ambiguity	Information is unclear	AMBI
Overlap	Information is nonatomic or used in multiple contexts	OVER
Redundancy	Information is redundant	REDUN

^a We use **WRONG** here to avoid clashing with **MISS**.

^b We use **MISS** here to be more meaningful in isolation, as it implies the synonym of “missing”; **OMISS** is less intuitive and **OMIT** would be inconsistent with the keys being adjective-based.

domains in previous sections and thus present them in that same order. These are the domains in which we automatically detect flaws as described in [Section 4.2.1](#) and present in [Section 5.2](#), so these are the only ones that are hyperlinked. We automatically detect the following classes of flaws:

- Test approaches with more than one category that violate our assumption of orthogonality (see [Section 2.1.1](#))
- Synonyms that violate transitivity (see [Section 2.1.2](#)); if two distinct approaches share a synonym but are not synonyms themselves, at least one relation is incorrect or missing.
- Synonyms between independently defined approaches; if two separate approaches have their own definitions, nuances, etc. but are also labelled as synonyms, this indicates that:
 1. the terms are interchangeable and this relation is trivially reflexive (see [Section 2.1.2](#)),
 2. at least one of these terms is defined incorrectly, and/or
 3. this synonym relation is incorrect.
- Parent-child relations that violate irreflexivity, or approaches that are parents of themselves.
- Pairs of synonyms where one is a subapproach of the other; these relations cannot coexist as synonym relations are symmetric while parent-child relations are asymmetric as outlined in [Sections 2.1.2](#) and [2.1.3](#), respectively.

Mention detecting other cycles as flaws as Future Work

Despite their nuance, the remaining domains are relatively straightforward, so we define them briefly as follows instead of defining them more rigorously in their own sections. Terms can be thought of as definition-label pairs, but there is a

Table 2.3: Observed flaw domains.

Domain	Description	Key
Categories	Approach categories, defined in Section 2.1.1	CATS
Synonyms	Synonym relations, defined in Section 2.1.2	SYNS
Parents	Parent-child relations, defined in Section 2.1.3	PARS
Definitions	Definitions given to terms	DEFS
Labels	Labels or names given to terms	LABELS
Scope	Scope of the information	SCOPE
Traceability	Records of the source(s) of information	TRACE

meaningful distinction between definition flaws and label flaws. Definition flaws are quite self-explanatory, but label flaws are harder to detect, despite occurring independently. Examples of label flaws include terms that share the same acronym or contain typos or redundant information. Sometimes, an author may use one term when they mean another. One could argue that their “internal” definition of the term is the cause of this mistake, but we consider this a label flaw where the wrong label is used as we would change the *label* to fix it. Additionally, some information is presented with an incorrect scope and sometimes should not have been included at all! Finally, some traceability information is flawed, such as how one document cites another or even what information is included *within* a document.

2.2.3 Additional Flaw Terminology

Some flaws involve information from more than one source, but referring to this as a “flaw between two sources” is awkward. We instead refer to this kind of flaw as an “inconsistency” between the sources. This clearly indicates that there is disagreement between the sources, but also does not imply that either one is correct—the inconsistency could be with some ground truth if *neither* source is correct!

Other flaws only involve one source, but we make a distinction between “self-contained” flaws and “internal” flaws. Self-contained flaws are those that manifest by comparing a document to an assertion of ground truth. These may appear once in a document or consistently throughout it. Sometimes, these do not require an explicit comparison to ground truth; these often include omissions as the lack of information is contained within a single source and does not need to be cross-checked against an assertion of ground truth. On the other hand, internal flaws arise when a document disagrees with itself by containing two conflicting pieces of information; this includes many contradictions and overlaps. Internal flaws can even occur on the same page, such as when a source gives the same acronym to two distinct terms (see [Overlap 7](#) and [Overlap 8](#))!

2.3 Explicitness

When information is written in natural language, a considerable degree of nuance can get lost when interpreting or using it. We call this nuance “explicitness”, or how explicit a piece of information is (or is *not*). For example, a source may provide data from which the reader can logically draw a conclusion, but may not state or prove this conclusion explicitly. In the cases where information is *not* explicit, we record it (see [Section 3.3.1](#) for more detailed discussion) and present it using (at least) one of the following keywords: “implied”, “can be”, “sometimes”, “should be”, “ideally”, “usually”, “most”, “likely”, “often”, “if”, and “although”.

Most information provided by sources we investigate is given explicitly; all sources cited throughout this thesis support their respective claims explicitly unless specified otherwise, usually via one of these keywords. It is important to note that we use the term “implicit” (as well as “implied by” when describing sources of information) to refer to any instance of “not explicit” information for brevity. Any kind of information can be implicit, including the names, definitions, categories (see [Section 2.1.1](#)), synonyms (see [Section 2.1.2](#)), and parents (see [Section 2.1.3](#)) of identified test approaches.

As our research focuses on the flaws present in the literature, the explicitness of information affects how seriously we take it. We call flaws based on explicit information “objective”, since they are self-evident in the literature. On the other hand, we call flaws based on implicit information “subjective”, since some level of judgement is required to assess whether these flaws are *actually* problematic. By looking for the indicators of uncertainty mentioned above, we can automatically detect subjective flaws when generating graphs and performing analysis (see [Sections 4.1](#) and [4.2](#), respectively).

Throughout our research, we also infer some information through “surface-level” analysis that follows straightforwardly but is not stated, explicitly or otherwise, by a source. Although these data originate from our judgement, we document them for completeness, using the phrase “inferred from” when relevant. All data in our glossaries without a citation are inferred, such as algebraic testing ([Peters and Pedrycz, 2000](#), Fig. 12.2) being a child of mathematical-based testing. Additionally, some inferences are based on information given by a source, which we cite alongside these inferences. For example, [Gerrard](#) describes large scale integration testing and legacy system integration testing in ([2000b](#), p. 30) and ([2000a](#), Tab. 2; [2000b](#), Tab. 1), respectively. While he never explicitly says so, we infer that these approaches are children of integration testing and system integration testing, respectively. Similarly, some test approaches appear to be combinations of other (seemingly orthogonal) approaches (described in more detail in [Appendix A.6](#)), from which we can extrapolate other test approaches. For example, [Moghadam \(2019\)](#) uses the phrase “machine learning-assisted performance testing”; since performance testing is a known test approach, we infer the existence of the test approach “machine learning-assisted testing” and include it in [our test approach glossary](#) as such. We also infer that child approaches inherit their parents’ categories (see [Section 2.1.1](#)).

Later: Ensure this is up to date

See #176

2.4 Credibility

In the same way we distinguish between the explicitness of information from different sources, we also wish to distinguish between the “explicitness” of the sources themselves! Of course, we do not want to overload terms, so we define a source as more “credible” if it:

- has gone through a peer-review process,
- is written by numerous, well-respected authors,
- cites a (comparatively) large number of sources, and/or
- is accepted and used in the field of software.

Sources may meet only some of these criteria, so we use our judgement (along with the format of the sources themselves) when comparing them.

2.5 Source Tiers

For ease of discussion and analysis, we group the complete set of sources into “tiers” based on their format, method of publication, and our heuristic of credibility. In order of descending credibility, we define the following tiers:

1. established standards (Section 2.5.1),
2. terminology collections (Section 2.5.2),
3. textbooks (Section 2.5.3), and
4. papers and other documents (Section 2.5.4).

We provide a summary of how many sources comprise each tier in Figure 2.1 and list all sources in each tier in Appendix B. The “papers” tier is quite large since we often “snowball” on terminology itself when a term requires more investigation (e.g., its definition is missing or unclear). This includes performing a miniature literature review on this subset to “fill in” missing information (see Section 3.4) and potentially fully investigating these additional sources, as opposed to just the original subset of interest, based on their credibility and how much extra information they provide. We use standards the second most frequently due to their high credibility and broad scope; for example, the glossary portion of ISO/IEC and IEEE (2017) has 514 pages! Using these standards allows us to record many test approaches in a similar context from a source that is widely used and well-respected.



Figure 2.1: Summary of how many sources comprise each source tier.

2.5.1 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, so we consider them to be the most credible sources. 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” (ISO/IEC and IEEE, 2017, p. viii) “that can be used by any organization when performing any form of software testing” (2022, p. vii; similar in 2016, p. ix). Only standards for software development and testing are in scope for this research (see [Chapter 1](#) for a high-level overview of what is in scope; see [Appendix A](#) for more detailed discussion on what we include and exclude).

2.5.2 Terminology Collections

These are collections of software testing terminology built up from multiple sources, such as the established standards outlined in [Section 2.5.1](#). For example, the SWE-BOK 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” (Kaner et al., 2011, p. xix). Even though it is “published by the IEEE Computer Society”, it “reflects the current state of generally accepted, consensus-driven knowledge derived from the interaction between software engineering theory and practice” (Washizaki, 2025). Due to this combination of IEEE standards and state-of-the-practice observations, we designate it as a collection of terminology as opposed to an established standard. Collections such as this are often written by a large organization, such as the International Software Testing Qualifications Board (ISTQB), but not always. Firesmith (2015)’s taxonomy presents relations between many test approaches and Doğan et al. (2014)’s literature review cites many sources from which we can “snowball” if desired (see [Section 3.1](#)), so we include them in this tier as well.

2.5.3 Textbooks

We consider textbooks to be more credible than papers (see [Section 2.5.4](#)) because they are widely used as resources for teaching software engineering and industry frequently uses them as guides. Although textbooks have smaller sets of authors, they follow a formal review process before publication. Textbooks used at McMaster University (Patton, 2006; Peters and Pedrycz, 2000; van Vliet, 2000) served as the original (albeit ad hoc and arbitrary) starting point of this research, and we investigate other books as they arise. For example, Hamburg and Mogyorodi (2024) cite Gerrard and Thompson (2002) as the original source for their definition of “scalability” (see [Section 6.3](#)) which we verify by looking at this original source.

2.5.4 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 of these are journal articles and conference papers, we also include the following document types. Some of these are not peer-reviewed works but are still useful for observing how terms

See [#89](#) are used in practice:

- Report ([Kam, 2008](#); [Gerrard, 2000a;b](#))
- Thesis ([Bas, 2024](#))
- Website ([LambdaTest, 2024](#); [Pandey, 2023](#))
- Booklet ([Knüvener Mackert GmbH, 2022](#))
- ChatGPT (GPT-4o) (2024) (with its claims supported by [Rus et al. \(2008\)](#))³

³For more information on how and why we use ChatGPT, see ??.

Chapter 3

Methodology

Since we wish to analyze the software testing literature, we collect data (both correct and incorrect) from a wide variety of relevant documents, focusing on test approaches and supporting information. This results in a big glossary of software approaches, some glossaries of supplementary terms, and a list of flaws. To ensure this data can be analyzed and expanded thoroughly and consistently, we need a process that can be repeated for future developments in the field of software testing or by independent researchers seeking to verify our work. Our methodology for documenting the current (messy) state of software testing terminology consists of:

1. Identifying authoritative sources on software testing and “snowballing” from them ([Section 3.1](#))
2. Identifying all test approaches¹ and related testing terms that are used repeatedly and/or have complex definitions ([Section 3.2](#))
3. Recording all relevant data ([Section 3.3](#)), including implicit data ([Section 3.3.1](#)), for each term identified in step 2; test approach data are comprised of:
 - (a) Names
 - (b) Categories²
 - (c) Definitions
 - (d) Synonyms³
 - (e) Parents⁴
 - (f) Flaws⁵ ([Section 3.3.2](#))
 - (g) Other relevant notes (prerequisites, uncertainties, other sources, etc.)
4. Repeating steps 1 to 3 for any missing or unclear terms ([Section 3.4](#)) until some stopping criteria ([Section 3.5](#))

¹Defined in [Section 2.1](#).

²Defined in [Section 2.1.1](#).

³Defined in [Section 2.1.2](#).

⁴Defined in [Section 2.1.3](#).

⁵Defined in [Section 2.2](#).

3.1 Identifying 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 (ISO/IEC and IEEE, 2017) and two versions of the Guide to the SoftWare Engineering Body Of Knowledge (SWEBOK Guide)—the newest one (Bourque and Fairley, 2014) and one submitted for public review⁶ (Washizaki, 2024). 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” (Johnson, 2014). We apply this concept to “referrals” between sources. For example, Hamburg and Mogyorodi (2024) cite Gerrard and Thompson (2002) as the original source for their definition of “scalability” (see Section 6.3) which we verify by looking at this original source. We group all sources into the source tiers we define in Section 2.5 and list all sources in each tier in Appendix B.

3.2 Identifying Relevant Terms

Before we can consistently track software testing terminology used in the literature, we must first determine what to record. We use heuristics to guide this process to increase confidence that we identify all relevant terms, paying special attention to the following when investigating a new source:

- glossaries, taxonomies, hierarchies, and lists of terms,
- testing-related terms (e.g., terms containing “test(ing)”, “review(s)”, “audit(s)”, “attack(s)”, “validation”, or “verification”),
- terms that had emerged as part of already-discovered test approaches, *especially* those that were ambiguous or prompted further discussion (e.g., terms containing “performance”, “recovery”, “component”, “bottom-up”, “boundary”, or “configuration”), and
- terms that imply test approaches⁷, including:
 - coverage metrics that may imply related test techniques,
 - software qualities that may imply related test types, and
 - software requirements that may imply related test approaches.

⁶Washizaki (2024) has been published since we investigated these sources; if time permits, we will revisit this published version.

⁷Since these methods for deriving test approaches only arose as research progressed, some examples would have been missed during the first pass(es) of resources investigated earlier in the process. While reiterating over them would be ideal, this may not be possible due to time constraints.

We apply these heuristics to the entirety of most investigated sources, especially established standards (see [Section 2.5.1](#)). However, we only partially investigate some sources, 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 3.4](#): [ISO \(2022; 2015\)](#); [Dominguez-Pumar et al. \(2020\)](#); [Pierre et al. \(2017\)](#); [Trudnowski et al. \(2017\)](#); [Yu et al. \(2011\)](#); [Tsui \(2007\)](#); [Goralski \(1999\)](#).

3.3 Recording Relevant Information

Once we have identified which terms from the literature are relevant, we can then track them consistently. We do this by building glossaries, mainly focusing on the one most central to our research—[our test approach glossary](#)—where we record all test approaches we identify in the literature. We give each test approach its own row where we record its name and any given definitions, categories, synonyms, and parents (along with any other notes, such as questions, prerequisites, and other resources to investigate) following the procedure in [Figure 3.1](#). Note that only the name and category fields are required; all other fields may be left blank, although a lack of definition indicates that the approach should be investigated further to see if its inclusion is meaningful (see [Section 3.4](#)). When recording flawed data, we record it in our glossary as dubious information (see [Section 3.3.1](#)) and/or in a separate document if it is clearly flawed (see [Section 3.3.2](#)).

For example, when we first encounter “A/B Testing” in [ISO/IEC and IEEE \(2022, p. 1, 36\)](#) as shown in [Figure 3.2](#), we apply [our procedure as follows](#):

1. Create a new row with the name “A/B testing” and the category “Approach”.
2. Record the synonym “Split-Run Testing”.
3. Record the parent “Statistical Testing”.
4. Record the definition “Testing ‘that allows testers to determine which of two systems or components performs better’”; note that we abstract away information that we have previously captured (i.e., its synonym and parent).

Q #3: Should glossary headers be capitalized?



Figure 3.1: Procedure for recording test approaches in our glossary.

3.1

A/B testing

split-run testing

statistical *testing* (3.131) approach that allows testers to determine which of two systems or components performs better

Figure 3.2: ISO/IEC and IEEE (2022, p. 1)’s glossary entry for “A/B testing”.

We also include the relevant citation “(ISO/IEC and IEEE, 2022, pp. 1, 36)”⁸ with these entries (excluding its name) to track where they came from. This source also provides the following information, which we capture as follows:

1. Record the note “It ‘can be time-consuming, although tools can be used to support it’, ‘is a means of solving the test oracle problem by using the existing system as a partial oracle’, and is ‘not a test case generation technique as test inputs are not generated’ ” (p. 36).
2. Replace the category of “Approach” with the more specific “Practice” (Fig. 2); note that this is consistent with the exclusion of “Technique” as a possible category for this approach (p. 36).

As we investigate other sources, we learn more about this approach. Firesmith (2015, p. 58) includes it in his taxonomy as shown in Figure 3.3. We add to our entry for “A/B Testing” as follows:

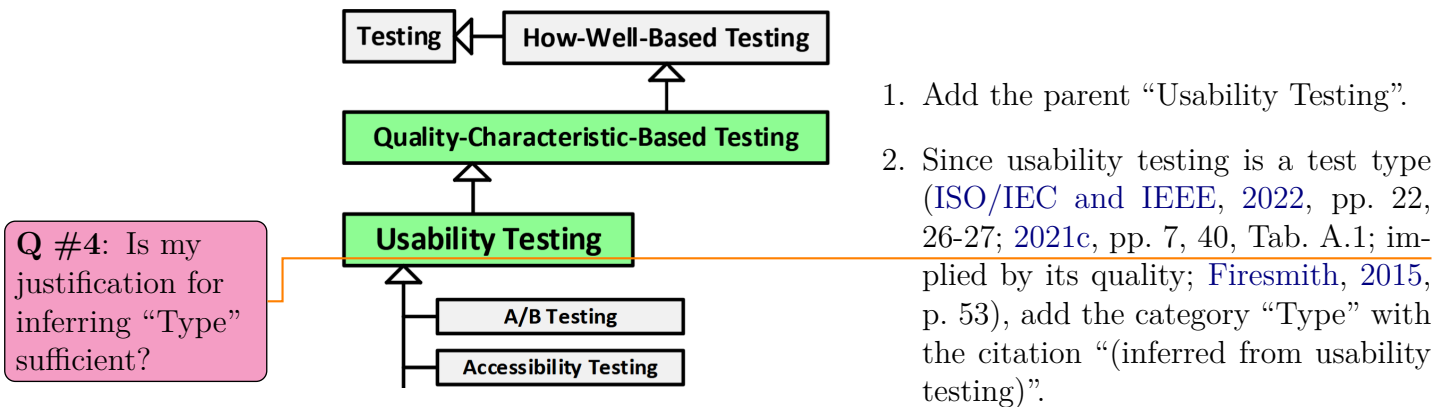


Figure 3.3: A/B testing’s inclusion in Firesmith (2015, p. 58)’s taxonomy.

This second change introduces an inference (defined in Section 2.3) that violates our assumption from Section 2.1.1 that categories are orthogonal, so we consider this to be an inferred flaw that we automatically detect and document (see Section 4.2.1 and Table C.3, respectively). This results in the corresponding row in Table 1.1, although we exclude the “Notes” column for brevity.

⁸In our test approach glossary, we *actually* cite this as “(IEEE, 2022, pp. 1, 36)” for brevity, as most standards are written by ISO/IEC and IEEE (see Figure 4.3); these glossary citations are handled by our choice of BibTeX keys, such as IEEE2022 in this example.

As implied throughout the previous example, we capture relevant citation information in our glossary in the author-year citation format. When citing the same authors or sources multiple times, we “reuse” information from previous citations when applicable. For example, the first row of [Table 4.1](#) contains the citation “(Author, 2022; 2021)”, which means that this information was present in two documents by Author: one written in 2022 and one in 2021. The following citation, “(2022)”, contains no author, which means it was written by the same one as the previous citation (Author). This reduces duplication and improves maintainability when recording this information. When processing these data when visualizing relations (see [Section 4.1](#)), we do so according to this logic (see the [relevant source code](#)) so we can consistently track the source(s) of data throughout our analysis.

We use this same procedure to track software qualities and supplementary terminology that is either shared by multiple approaches or too complicated to explain inline. We create a separate glossary for both qualities and supplementary terms, each with a similar format to [our test approach glossary](#). Since these terms do not have categories, the process of recording them is much simpler, only requiring us to record the name, definition, and synonym(s) of these terms, along with any additional notes. The only new information we capture is the “precedence” for a software quality to have an associated test type, since each test type measures a particular software quality (see [Table 2.1](#)). These precedences are instances where a given software quality is related to, is covered by, or is a child, parent, or prerequisite of another quality with an associated test type, as given by the literature.

Tracking information about software qualities helps us investigate the literature more thoroughly, since these data may become relevant based on information from other yet-uninvestigated sources. When the literature mentions (or implies) a test approach that corresponds to a software quality we have recorded, we first follow the procedure given in [Figure 3.1](#) with the information provided in the source that mentions it. We then remove the relevant data from our quality glossary and repeat our procedure with it to upgrade the quality to a test type in [our test approach glossary](#).

3.3.1 Recording Implicit Information

As described in [Section 2.3](#), the use of natural language introduces significant nuance that we need to document. Keywords such as “implied”, “can be”, “sometimes”, “should be”, “ideally”, “usually”, “most”, “likely”, “often”, “if”, and “although” indicate that information from the literature is *not* explicit. These keywords often appear directly within the literature, but even when they do not, we use them to track explicitness in [our test approach glossary](#) to provide a more complete summary of the state of software testing literature without getting distracted by less relevant details. We find the following non-mutually exclusive cases of implicit information from the literature:

1. **The information follows logically** from the source and information from

Q #5: Does this make sense? I tried rewording it for brevity/clarity

Later: Ensure this is up to date

others, but is not explicitly stated.

2. **The information is not universal** but still applies in certain cases.
3. **The information is conditional**, requiring certain prerequisites to be satisfied (a more specific case of information not being universal).
4. **The information is dubious**; while it is present in the literature, there is reason to doubt its accuracy.

When we encounter information that meets one of these criteria, we use an appropriate keyword to capture this nuance in [our test approach glossary](#) (see [Table 3.1](#)). This also helps us identify implicit information when performing later analysis. Despite “implicit” only describing the first of these cases, we use it (as well as “implied by” when describing sources of information) as a shorthand for all “not explicit” manifestations throughout this thesis for clarity.

Regarding the last entry in [Table 3.1](#), if a test approach in [our test approach glossary](#) has a name ending in “ (Testing)” (space included), then the word “Testing” might not be part of its name *or* it might not be a test approach at all! For example, the term “legacy system integration” is used in [Gerrard \(2000a, pp. 12–13, Tab. 2; 2000b, Tab. 1\)](#), but the more accurate “legacy system integration testing” is used in [\(2000b, pp. 30–31\)](#). In other cases where a term is *not* explicitly labelled as “testing”, we add the suffix “ (Testing)” (when it makes sense to do so) and consider the test approach to be implied.

Table 3.1: Breakdown of keywords used for recording and analyzing implicitness.

Keyword	Follows Logically	Not Universal	Conditional	Dubious
“implied”	X			
“can be”		X	X	
“sometimes”		X	X	
“should be”		X		X
“ideally”		X		X
“usually”		X		
“most”		X		
“likely”	X	X		X
“often”		X		
“if”		X	X	X
“although”				X
“ (Testing)”	X			

3.3.2 Recording Flaws

In [Section 4.2.1](#), we describe some subsets of flaws that we detect automatically by analyzing [our test approach glossary](#). Unfortunately, most are too complex and need to be tracked manually. We record these more detailed flaws in a separate document as \LaTeX enumeration items along with comments that we can parse automatically, allowing us to analyze them more broadly. These comments have the following format:

```
% Flaw count (MNFST, DMN): {A1} {A2} ... | {B1} ... | {C1} ...
```

MNFST and DMN are placeholders for the “keys” given in [Tables 2.2](#) and [2.3](#), respectively, that we use to track a flaw’s manifestation(s) and domain(s) (defined in [Section 2.2](#)). For example, the comment line for an incorrect synonym relation would start with `% Flaw count (WRONG, SYNS)` and one for a redundant label would start with `% Flaw count (REDUN, LABELS)`. We omit these keys from constructed examples of these comments without associated flaws throughout this chapter for brevity. Finally, A1, A2, B1, and C1 are each placeholders for a source involved in this example flaw; in general, there can be arbitrarily many. We represent each source by its \LaTeX key, and wrap each one in curly braces (with the exception of the ISTQB glossary due to its use of custom commands via `\citealias{}`) to mimic \LaTeX ’s citation commands for ease of parsing. We then separate each “group” of sources with a pipe symbol (`|`) so we can compare each pair of groups; in general, a flaw can have any number of groups of sources.

As mentioned in [Section 2.2.3](#), we make a distinction between “self-contained” flaws and “internal” flaws. We track self-contained flaws by recording the single source that the flaw is present in such as in the first line below. In contrast, we track internal flaws by recording the single source in multiple groups (as defined above). The second line is a standard example of this, while the third is more complex; in this case, source Y agrees with only one of the conflicting sources of information in X.

```
% Flaw count: {X}
% Flaw count: {X} | {X}
% Flaw count: {X} | {X} {Y}
```

We can also specify the “explicitness” (see [Section 2.3](#)) of a flaw by inserting the phrase “implied by” after the sources of explicit information and before those of implicit information, such as in the following example:

```
% Flaw count: {X} {Y} | {Z} implied by {X}
```

It is important to note that some flaws are multi-faceted, and thus can be viewed in multiple ways. For example, [Kam \(2008, p. 42\)](#) says “See *boundary value analysis*,” for the glossary entry of “boundary value testing” but does not include “boundary value analysis” in the glossary. This is trivially an example of a missing definition, but it is also an example of incorrect traceability information. Therefore, we include two comment lines, one for `(WRONG, TRACE)` and one for `(MISS, DEFS)`. The pair of keys that appears first determines where this flaw is

See [#137](#) and [#138](#)

Q #6: Is this clear enough?

See [#157](#)

displayed, although all pairs are counted as their own flaws; this flaw is displayed as an example of incorrect traceability information as [Mistake 25](#) since we determine this to be more meaningful.

Occasionally, we assert that a source from a less credible tier is more correct than a source from a more credible tier. For example, the terms “acceleration tolerance testing” and “acoustic tolerance testing” do not seem refer to software testing, but [Firesmith \(2015, p. 56\)](#) includes them regardless. Elsewhere, they seem to refer to testing the acoustic tolerance of rats ([Holley et al., 1996](#)) or the acceleration tolerance of astronauts ([Morgun et al., 1999, p. 11](#)), aviators ([Howe and Johnson, 1995, pp. 27, 42](#)), or catalysts ([Liu et al., 2023, p. 1463](#)), which which are not relevant to software testing. This flaw is supported by these additional papers found via a miniature literature review (described in [Section 3.4](#)) from a lower source tier than [Firesmith \(2015\)](#) (which is a terminology collection; see [Section 2.5](#)). However, this flaw is really based in [Firesmith \(2015\)](#) and not in these additional papers, but this would be counted as a flaw in these papers if they were included as detailed above. Therefore, we document these “assertion” sources separately to track them for traceability without counting them incorrectly in [Section 4.2.2](#). As an example, the following is the comment line for this specific example ([Mistake 10](#)):

```
% Flaw count (WRONG, LABELS): {Firesmith2015}
% Assertion: {LiuEtAl2023} {MorgunEtAl1999} {HolleyEtAl1996}
↪ {HoweAndJohnson1995}
```

3.4 Undefined Terms

The literature mentions many software testing terms without defining them. While this includes test approaches, software qualities, and more general software terms, we focus on the former as the main focus of our research. In particular, [ISO/IEC and IEEE \(2022\)](#) and [Firesmith \(2015\)](#) name many undefined test approaches. Once we exhaust the standards in [Section 2.5.1](#), we perform miniature literature reviews on these subsets to “fill in” the missing definitions (along with any relations), essentially “snowballing” on these terms as described in [Section 3.1](#). This process uncovers even more approaches, on which we can then repeat this process. While we are only able to do this partially due to time constraints, we outline what we accomplished and potential next steps in [Section 8.1](#).

3.5 Stopping Criteria

Unfortunately, we cannot continue looking for new approaches indefinitely. We therefore need a “stopping criteria” to let us know when we are “finished” looking for test approaches in the literature. Our original plan was to repeat step 4 until we got diminishing returns, implying that we achieved something close to a complete taxonomy! However, this was infeasible due to time constraints. We also considered property-based testing as a potential stopping point for our research as we

were surprised that it is not mentioned by any sources we investigated. With more time, we would find definitions for all uncovered terms as described in [Section 3.4](#), as well as those that did not arise as described in [Section 8.2](#), but unfortunately, we impose our stopping point artificially.

Chapter 4

Tools

To better understand our findings, we build tools to visualize relations between test approaches (Section 4.1) and automatically track their flaws (Section 4.2). Doing this manually would be error-prone due to the amount of data involved (for example, we identify 561 test approaches) and the number of situations where the underlying data would change, such as adding to it, further analyzing it, or correcting it. These all require tedious updates to the corresponding visualizations that may be overlooked or done incorrectly, so automating these processes allows for our results to be reproduced (by us or others) and account for new data. Automation also allows us to observe the impacts of smaller changes and helps us verify the tools themselves; for example, tracking a flaw manually should affect relevant flaw counts, which we can double-check. Additionally, we define \LaTeX macros (Section 4.3) to help achieve our goals of maintainability, traceability, and reproducibility.

4.1 Approach Relation Visualization

To better understand the relations between test approaches, we develop a tool to generate visualizations of these relations automatically. Since we use a consistent format to track synonym and parent-child relations (defined in Sections 2.1.2 and 2.1.3, respectively) between approaches in our test approach glossary, we can parse them systematically. For example, if the entries in Table 4.1 appear, then their parent-child relations are visualized as shown in Figure 4.1a. Overall, the parent-child relations between test approaches *should* result in something resembling a hierarchy (or multiple discrete hierarchies), although this is not the case due to flaws in the literature (see Section 2.2.2). Therefore, we visualize all parent-child relations as they are guaranteed to be visually meaningful.

Table 4.1: Example glossary entries demonstrating how we track parent-child relations.

Name^a	Parent(s)
A	B (Author, 2022; 2021), C (2022)
B	C (implied by Author, 2022)
C	D (implied by Author, 2017)
D (implied by Author, 2017)	

^a “Name” can refer to the name of a test approach, software quality, or other testing-related term, but we only visualize relations between test approaches.



Figure 4.1: Example generated visualizations of parent-child relations.

However, since each term is trivially a synonym of itself and there are many non-problematic synonyms that do not imply flaws (see [Section 2.1.2](#)), we only visualize the synonym relations that may indicate flaws given in [Section 2.2.2](#); i.e., intransitive synonyms and synonyms between independently defined approaches. We deduce these conditions from the information we parse from our glossary. For example, if the entries in [Table 4.2](#) appear, then they are visualized as shown in [Figure 4.2](#) (note that X does not appear since it is not defined independently and does not violate transitivity).

Since we also track the “explicitness” of information (defined in [Section 2.3](#)), we can represent explicit *and* implicit relations without double counting them during the analysis in [Section 4.2](#). If a relation is both explicit *and* implicit, we only display the implicit relation if its source is from a more credible source tier (see [Sections 2.4](#) and [2.5](#)). For example, if “StdAuthor” from [Table 4.2](#) is an author of an established standard, then we display the implicit relation from their document alongside the explicit one from “Author” as shown in [Figure 4.2](#). Implicit approaches and relations are denoted by dashed lines, as shown in [Figures 4.1a](#) and [4.2](#), while explicit approaches are *always* denoted by solid lines, even if they

Table 4.2: Example glossary entries demonstrating how we track synonym relations.

Name ^a	Synonym(s)
E	F (Author, 2022; implied by StdAuthor, 2021)
G	F (Author, 2017), H (implied by 2022)
H	X (StdAuthor, 2021)

^a “Name” can refer to the name of a test approach, software quality, or other testing-related term, but we only visualize relations between test approaches.

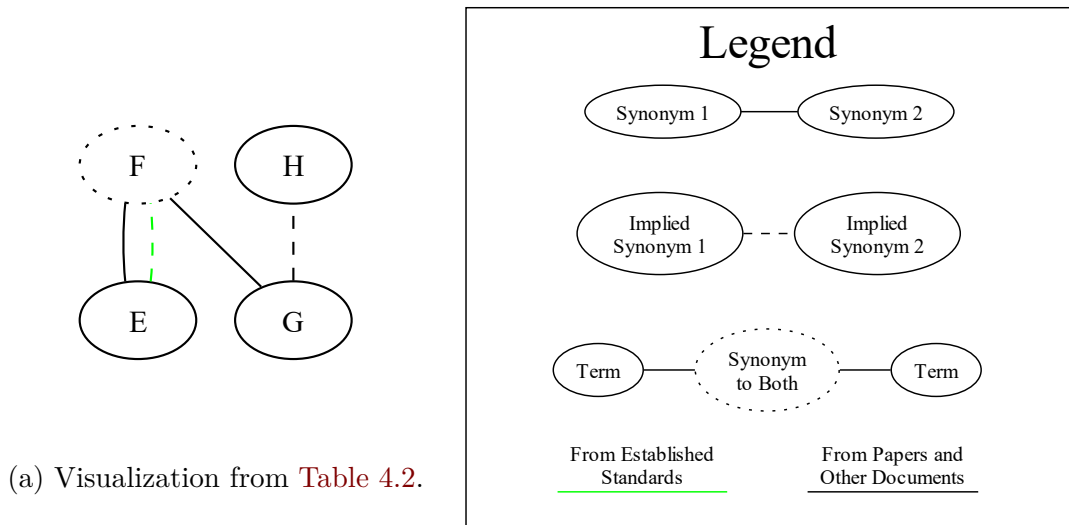


Figure 4.2: Example generated visualizations of synonym relations.

are also implicit. We can also generate “explicit” versions of these visualizations that exclude implicit approaches and relations; for example, [Figure 4.1b](#) is the explicit version of [Figure 4.1a](#).

We also colour each relation according to its source tier, which is possible because we cite all recorded relations as described in [Section 3.3](#). Each source tier gets its own colour, which we label for each relevant source tier in a given visualization’s legend (such as [Figure 4.2](#)), although we omit this colouring from [Figures 4.1](#) and [4.4](#) for clarity. We also only display the relation with the most credible source tier (except if there is a more credible implicit relation as we previously describe). Finally, we also colour inferences (see [Section 2.3](#)) grey and proposals (see [Chapter 6](#)) orange, such as in [Figures 6.1](#) to [6.3](#).

These visualizations tend to be large, so it is often useful to focus on specific subsets of them. For each approach category (defined in [Section 2.1.1](#)), we generate a visualization restricted to its approaches and the relations between them. We also generate a visualization of all static approaches along with the relations between them *and* between a static approach and a dynamic approach. This static-focused visualization is notable because static testing is sometimes considered to be a separate approach category (see [Contradiction 5](#)). Since dynamic approaches are our primary focus (see [Appendix A.3](#)), we include them in this static visualization, colouring their nodes grey to distinguish them. We can also generate more focused visualizations from a given subset of approaches, such as those pertaining to recovery or scalability. We use these visualizations to better understand the relations within these subsets of approaches, but we can also update them based on our recommendations in [Chapter 6](#) by specifying sets of approaches and relations to add or remove.

4.2 Flaw Analysis

In addition to analyzing specific flaws, it is also useful to examine them at a higher level. We automate subsets of this task where applicable ([Section 4.2.1](#)) and augment the remaining manual portion with automated tools ([Section 4.2.2](#)). This gives us an overview of:

- how many flaws (defined in [Section 2.2](#)) there are,
- how these flaws present themselves (see [Section 2.2.1](#)),
- in which knowledge domains these flaws occur (see [Section 2.2.2](#)),
- how explicit (see [Section 2.3](#)) these flaws are, and
- how responsible each source tier (defined in [Section 2.5](#)) is for these flaws.

To understand where flaws exist in the literature, we group them based on the source tier(s) responsible for them. We then count each flaw *once* per source tier if

it appears within it *and/or* between it and a more credible tier¹ (see [Sections 2.4](#) and [2.5](#)). This avoids counting the same flaw more than once for a given source tier, which would give the number of *occurrences* of all flaws instead of the more useful number of flaws *themselves*. When taking a more detailed look at the *sources* of flaws (as opposed to just the responsible source *tiers*) as we do in [Figure 5.1](#), we also count the following sources of flaws separately:

1. self-contained flaws (defined in [Section 2.2.3](#)),
2. internal flaws (defined in [Section 2.2.3](#)),
3. those between documents with the same set of authors, which includes
 - (a) the various combinations of authors of established standards (defined in [Section 2.5.1](#))—ISO, the International Electrotechnical Commission (IEC), and IEEE—as shown in [Figure 4.3](#) and
 - (b) the different versions of the Guides to the SoftWare Engineering Body Of Knowledge (SWEBOK Guides), which have different editors ([Washizaki, 2024](#); [Bourque and Fairley, 2014](#)) but are written by the same organization: the IEEE Computer Society ([Washizaki, 2025](#); see [Section 2.5.2](#)), and
4. those within a single source tier.

As before, we do not double count these sources of flaws, meaning that the maximum number of counted flaws possible within a *single* source tier in this more detailed view is four (one for each type). This only occurs if there is an example of each flaw source that is *not* ignored to avoid double counting; for example, while a single flaw within a single document would technically and trivially fulfill all four criteria, we would only count it once.

4.2.1 Automated Flaw Analysis

As outlined in [Section 2.2.2](#), we automatically detect synonym relations from [our test approach glossary](#) that violate transitivity to generate our visualizations. These relations are significant because they indicate potential flaws. We automatically detect and format these flaws to present them when discussing synonym relation flaws in [Section 5.2.2](#). For these and other kinds of flaws, we also generate the corresponding comments as described in [Section 3.3.2](#) we detect automatically before generating the corresponding L^AT_EX file to ensure these flaws also get analyzed. Since we already automatically detect one kind of flaw, the next logical step is then to detect more.

Parent-child relations that violate irreflexivity as outlined in [Section 2.1.3](#) (i.e., cases where a child is given as a parent of itself) are also trivial to automate by

¹If an inconsistency occurs between two source tiers and the more credible one is *incorrect*, we instead count it as an inconsistency between it and the asserted truth from the less credible source, as described in [Section 3.3.2](#).



Figure 4.3: The sets of authors of established standards.

looking for lines in the generated \LaTeX files that start with $I \rightarrow I$, where I is the label used for a test approach node in these visualizations. This process results in output similar to Figure 4.4a. We use a similar process to detect pairs of approaches with a synonym relation *and* a parent-child relation as described in Section 2.2.2. To find these pairs, we build a dictionary of each term’s synonyms to evaluate which synonym relations are notable enough to include in the visualization, and then check these mappings to see if one appears as a parent of the other. For example, if J and K are synonyms, a generated \LaTeX file with a parent line starting with $J \rightarrow K$ *or* $K \rightarrow J$ would be visualized as shown in Figure 4.4b. We present these two classes of flaws when discussing parent-child relation flaws in Section 5.2.3.

While just counting the total number of flaws (found automatically *or* manually) is trivial, tracking the source(s) of these flaws is more useful, albeit more involved. Since we consistently track the appropriate citations for each piece of information we record (see Tables 4.1 and 4.2 for examples of our citation format), we can use them to identify the offending source tier(s). This comes with the added benefit that we can format these citations to use with \LaTeX ’s citation commands in this thesis, including generating the comments used to analyze flaws as described in Section 4.2.2.

Alongside this citation information, we include keywords so we can assess how “explicit” a piece of information is (see Section 2.3). This is useful when counting



Figure 4.4: Example generated visualizations containing flaws.

flaws, since they can be both objective and subjective but should not be double counted as both! When presenting the numbers of flaws sorted by various criteria in [Chapter 5](#), we only count each flaw for its most “explicit” manifestation, similarly to how we visualize the relations between approaches as described in [Section 4.1](#).

4.2.2 Augmented Flaw Analysis

To perform more detailed analysis on the flaws we uncover, we use \LaTeX comments to capture information about the flaws themselves as outlined in [Section 3.3.2](#). We include these flaw comments when manually recording flaws from the literature ([Section 3.3](#)) and generate them when automatically detecting flaws from [our test approach glossary](#) ([Section 4.2.1](#)).

The main way we use these comments is to determine where each flaw originates. We compare the authors and years of each source involved with a given flaw to determine if it manifests within a single document and/or between documents with the same set of authors. Then, we group these sources into their tiers (see the [relevant source code](#)). We then distill these lists of sources down to sets of tiers and compare them against each other to determine how many times a given flaw manifests between source tiers, which we use when counting flaws in [Chapter 5](#).

We parse implicit information following the same rules given in [Section 4.2.1](#) for automatically detecting flaws. Note that we only count subjective flaws if there is not an equivalent objective flaw, as we do when visualizing relations (as described in [Section 4.1](#)). The following comment line from [Contradiction 9](#) is an example of a flaw that is both objective and subjective:

```
% Flaw count (CONTRA, DEFS): {IEEE2021c} {IEEE2017} |
↪ {vanVliet2000} implied by {IEEE2021c}
```

This indicates that the following flaws are present:

- an objective inconsistency between a textbook and a standard,

- a subjective flaw within a single document, and
- a subjective inconsistency between documents with the same set of authors (ISO/IEC and IEEE; see [Figure 4.3](#)).

This third flaw only affects our more nuanced breakdown of the sources of flaws in [Figure 5.1](#). Note that we do not double count the first flaw. We likewise do not double count flaws that reappear when comparing between pairs of groups; for example, we would only count the inconsistency between X and Z *once* in the following flaw comment:

```
% Flaw count: {X} | {X} {Y} | {Z}
```

In cases where a flaw can be viewed in multiple ways (see [Section 3.3.2](#)), we only represent the flaw once in the subsections of [Chapter 5](#) and the full lists in [Appendix C](#) based on the pair of keys that appears first. This allows us to decide which view is most central to understanding the flaw without affecting the results of our research. The choice of which flaw is more “meaningful” only affects its presentation, since we also count its other views as flaws (for example, in [Tables 5.1](#) and [5.2](#)), with the benefit of not introducing clutter by displaying it in full multiple times.

4.3 L^AT_EX Commands

To improve maintainability, traceability, and reproducibility, we define helper commands (also called “macros”) for content that is prone to change or used in multiple places. For example, we use Python scripts to calculate values based on our glossaries and save them to files to be assigned to corresponding L^AT_EX macros. We use these throughout our documents instead of manually updating these constantly changing values, which is prone to error. [Table 4.3](#) lists these macros, along with their current values and descriptions of what they represent. Our Python scripts convert numbers to their textual equivalents when necessary to follow IEEE guidelines.

Table 4.3: L^AT_EX macros for calculated values.

Macro	What it Counts	Value
<code>\approachCount{}</code>	Identified test approaches	561
<code>\qualityCount{}</code>	Identified software qualities	77
<code>\srcCount{}</code> ^a	Sources used in glossaries	71
<code>\flawCount{}</code> ^b	Identified flaws	295

Continued on next page

Table 4.3: L^AT_EX macros for calculated values. (Continued)

Macro	What it Counts	Value
<code>\TotalBefore{}</code> ^c	Test approaches identified before step 4 ^d	474
<code>\UndefBefore{}</code> ^c	Undefined test approaches identified before step 4 ^d	175
<code>\TotalAfter{}</code> ^c	Test approaches identified after step 4 ^d	561
<code>\UndefAfter{}</code> ^c	Undefined test approaches identified after step 4 ^d	191
<code>\multiCatCount{}</code>	Total number of approaches with multiple categories	28
<code>\multiCatMax{}</code>	Category with the most overlaps	technique
<code>\multiCatMaxCount{}</code>	Number of overlaps involving the previous category	23
<code>\multiSynCount{}</code>	Terms given as synonyms for multiple discrete terms	nine
<code>\parSynCount{}</code>	Pairs of test approaches with a child-parent <i>and</i> synonym relation	18
<code>\selfParCount{}</code>	Test approaches that are a parent of themselves	three

^a Calculated in L^AT_EX from source tier lists; see Section 4.3.

^b Alias for `\totalFlawDmnBrkdwn{15}`; see Section 4.3.

^c These macros are defined as counters to allow them to be used in calculations within L^AT_EX (such as in Section 3.4 and Fig. 8.1).

^d Step 4 of our methodology involves iterating over undefined terms and is described in more detail in Section 3.4.

Additionally, we count flaws based on their manifestation and domain, explicitness, and source tier (defined in Sections 2.2, 2.3, and 2.5, respectively). For each source tier, we create two files that each include both levels of explicitness: one for manifestations and one for domains. For example, flaws in standards are saved to `build/stdFlawMnfstBrkdwn.tex` by manifestation. We then assign these data to macros (such as `\stdFlawMnfstBrkdwn{}`) to populate Tables 5.1 and 5.2. For example, we access the number of objective and subjective mistakes in standards by using `\stdFlawMnfstBrkdwn{1}` and `\stdFlawMnfstBrkdwn{2}`, respectively. We follow a similar process for tracking the total numbers of flaws; this includes `\totalFlawMnfstBrkdwn{13}` and `\totalFlawDmnBrkdwn{15}` which are identical and track the total number of identified flaws.

Just as with calculated values, it is important that repeated text is updated consistently, which we accomplish by defining more macros. Some of these are generated by Python scripts in a similar fashion to calculated values, such as those for flaw manifestations, flaw domains, and source tiers as shown in Table 4.4. The latter are built by extracting all sources cited in our three glossaries, catego-

rizing, sorting, and formatting them (including handling edge cases), and saving them to a file. These are then assigned to `\stdSources{}`, `\metaSources{}`, `\textSources{}`, and `\paperSources{}` and include:

1. the source tier’s name,
2. the list of sources in the tier, and
3. the number of sources in the tier.

These are accessed by passing in the corresponding number in the above enumeration (e.g., `\paperSources{2}`). We use the first value for the subheadings in [Section 2.5](#), the first two for [Appendix B](#) and the third to build [Figure 2.1](#) and calculate `\srcCount{}` (see [Table 4.3](#)).

However, we create most of these macros for reused text manually when we first notice the reuse, including the source tier macros in [Table 4.4](#). Some of these macros account for context-specific formatting depending on how they are used, such as capitalization. These tend to be less well-defined, since they arise naturally from the writing process, so we omit these details from the manually created text macros in [Table 4.5](#), which are grouped based on the type of information they contain.

Table 4.4: L^AT_EX macros for referencing well-defined sections.

	Flaw Manifestations ^a	Flaw Domains ^b	Source Tiers ^c
Macros (Values)	<code>\wrong{}</code> (Mistakes)	<code>\cats{}</code> (Categories)	<code>\stds{}</code> (Established Standards)
	<code>\miss{}</code> (Omissions)	<code>\syms{}</code> (Synonyms)	<code>\metas{}</code> (Terminology Collections)
	<code>\contra{}</code> (Contradictions)	<code>\pars{}</code> (Parents)	<code>\texts{}</code> (Textbooks)
	<code>\ambi{}</code> (Ambiguities)		<code>\papers{}</code> (Papers and Other Documents)
	<code>\over{}</code> ^d (Overlaps)		<code>\papers*{}</code> (Papers and Others)
	<code>\redun{}</code> (Redundancies)		
Used In	Tables 2.2 and 5.1 (in both thesis and paper)	Tables 2.3 and 5.2 (in both thesis and paper)	Figures 2.1 and 5.1 and Tables 5.1 and 5.2

^a Defined in Section 2.2.1; we also define starred versions, such as `\wrong*{}` (Mistake), that use the singular noun for use in Table 2.2.

^b Defined in Section 2.2.2; we only include domains with their own section.

^c Defined in Section 2.5.

^d We overwrite the primitive T_EX command `\over{}` since we do not otherwise use it.

Table 4.5: L^AT_EX macros for reused text.

Type	Macro	Used in
Flaws ^a	<code>\bugPattonFlaw{}</code>	Chapter 1 and Mistake 19
	<code>\alphaFlaw{}</code>	Chapter 1 and Contradiction 20
	<code>\loadFlaw{}</code>	Chapter 1 and Contradiction 26
	<code>\tourFlaw{}</code>	Chapters 1 and 5 and Contradiction 19
	<code>\multiCatIntro{}</code>	Sections 2.2.2 and 5.2.1 and Appendices C.2.1 and C.3.1
	<code>\redBoxFlaw{}</code>	Sections 2.3 and 5.1.1 and Mistake 29
	<code>\seeRefMissing{}</code>	Section 3.3.2 and Mistake 25
	<code>\tolTestFlaw{}</code>	Sections 3.3.2 and 5.1.1 and Mistake 10
	<code>\accelTolTest{}</code>	<code>\tolTestFlaw{}</code> and Appendix A.1
	<code>\errorGuessFlaw{}</code>	Section 5.1.1 and Mistake 2
	<code>\parSheetTestFlaw{}</code>	Section 5.1.1 and Mistake 13
	<code>\parSynIntro{}</code>	Section 5.2.3 and Appendix C.2.3
	<code>\parSynFlaw{}</code>	<code>\parSynIntro{}</code> and Section 2.2.2
Footnotes	<code>\ftrnote{}</code>	Thesis (automated) and paper (manual) versions of Table C.2
	<code>\specfn{}</code>	
	<code>\ucstn{}</code>	
	<code>\notDefDistinctIEEE{}</code>	Contradiction 5 and Section 10.1
	<code>\distinctIEEE{}</code>	<code>\notDefDistinctIEEE{}</code> and Omission 7
Links	<code>\ourApproachGlossary{}</code>	Table 1.1, Sections 2.3, 3.3, 3.3.1, 3.3.2, 4.1, 4.2.1, 4.2.2, 5.2.2, and 6.1, and Appendix A.6
	<code>\seeSrcCode{}</code>	Sections 4.1 and 4.2.1
	<code>\recFigs{}</code>	Chapters 5 and 6
RQs	<code>\rqatext{}</code>	Chapter 1 and seminar slides
	<code>\rqbtext{}</code>	
	<code>\rqctext{}</code>	

Continued on next page

Table 4.5: L^AT_EX macros for reused text. (Continued)

Type	Macro	Used in
Misc.	<code>\supersAck{}</code>	Acknowledgements and seminar slides
	<code>\supers{}</code>	Declaration of Academic Achievement and <code>\supersAck{}</code>
	<code>\highLvlScope{}</code>	Chapter 1 and Section 2.5.1
	<code>\defRel{}</code>	Sections 2.1.2 and 2.1.3
	<code>\oneSrcDistinct{}</code>	Sections 2.2.3 and 4.2.2
	<code>\approachFields{}</code>	Sections 2.3 and 3.3
	<code>\impKeywords{}</code>	Sections 2.3 and 3.3.1
	<code>\orthTestIntro{}</code>	Section 2.3 and Appendix A.6
	<code>\listAllSrcs{}</code>	Sections 2.5 and 3.1
	<code>\addTextEx{}</code>	Sections 2.5.3 and 3.1

^a Defined in Section 2.2.

Chapter 5

Observed Flaws

After gathering all these data¹, we find many flaws. Figure 5.1 shows the source tiers (see Section 2.5) responsible for these flaws, which reveals a lot about software testing literature:

1. Established standards (Section 2.5.1) aren't actually standardized, since:
 - (a) other documents disagree with them *very* frequently and
 - (b) they are the most internally inconsistent source tier!
2. Less standardized (or “credible”; see Section 2.4) documents, such as terminology collections and textbooks (Sections 2.5.2 and 2.5.3, respectively) are also not followed to the extent they should be.
3. Documents across the board have flaws within the same document, between documents with the same author(s), or even with reality!

Q #8: Is this too strong of a synonym for “ground truth”?

To better understand and analyze these flaws, we group them by their manifestations and their domains as defined in Section 2.2. We present the total number of flaws by manifestation and by domain in Tables 5.1 and 5.2, respectively, where a given row corresponds to the number of flaws either within that source tier and/or with a “more credible” one (i.e., a previous row in the table; see Sections 2.4 and 2.5). We also group these flaws by their explicitness (defined in Section 2.3) by counting (Obj)ective and (Sub)jective flaws separately, since additional context may rectify them. Since we give each flaw a manifestation *and* a domain, the totals per source and grand totals in these tables are equal. From these tables, we can draw some conclusions about *how* the literature is flawed:

1. Contradictions are by *far* the most common way for a flaw to manifest, which makes sense: if two (groups of) authors do not communicate or work with different resources, there is a higher chance that they will disagree.
2. Approach categorizations are the most subjective and one of the most common flaw domains, likely due to the lack of standardization about what categories to use (see Section 7.2 for more detailed discussion).

¹Available in ApproachGlossary.csv, QualityGlossary.csv, and SuppGlossary.csv at <https://github.com/samm82/TestingTesting>.

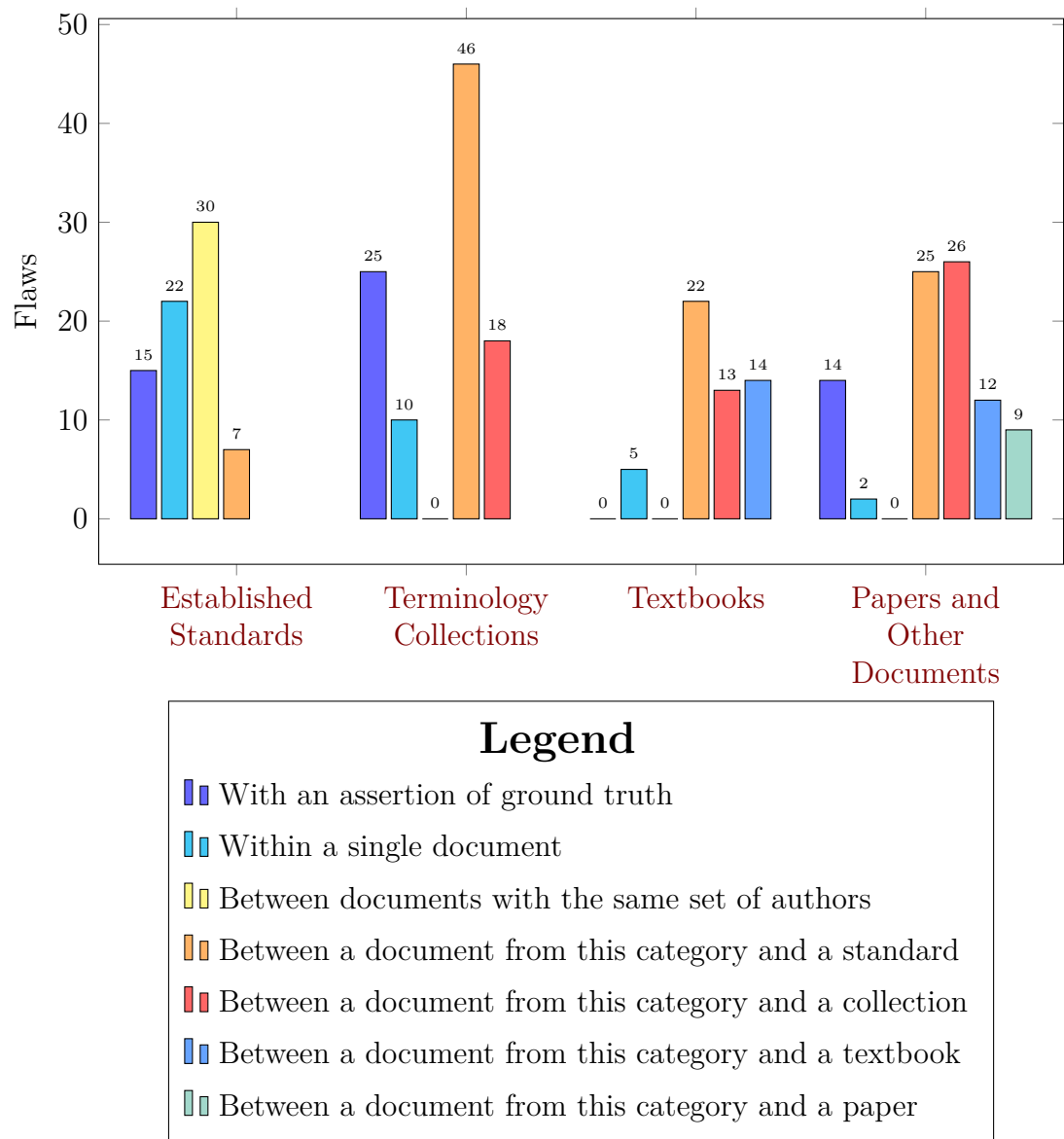


Figure 5.1: Identified flaws by the source tier responsible. Some bars are omitted as they correspond to comparisons we do not make; see [Section 4.2](#).

Table 5.1: Breakdown of identified flaws by manifestation and source tier.

Source Tier	Mistakes		Omissions		Contradictions		Ambiguities		Overlaps		Redundancies		Total
	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	
Established Standards	7	1	2	0	23	6	6	1	8	0	3	0	57
Terminology Collections	15	0	3	0	41	14	15	1	6	0	2	0	97
Textbooks	6	2	2	0	34	6	2	0	1	0	0	0	53
Papers and Others	16	1	6	0	31	21	9	0	0	1	3	0	88
Total	44	4	13	0	129	47	32	2	15	1	8	0	295

Table 5.2: Breakdown of identified flaws by domain and source tier.

Source Tier	Categories		Synonyms		Parents		Definitions		Labels		Scope		Trace.		Total
	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	Obj	Sub	
Established Standards	12	2	3	3	11	1	16	2	7	0	0	0	0	0	57
Terminology Collections	14	10	11	2	12	3	22	0	12	0	5	0	6	0	97
Textbooks	0	0	14	1	10	6	17	1	4	0	0	0	0	0	53
Papers and Others	14	13	19	6	12	3	11	0	8	1	0	0	1	0	88
Total	40	25	47	12	45	13	66	3	31	1	5	0	7	0	295

We summarize the flaws that we discover manually (see [Section 4.2.2](#)) in [Section 5.1](#) based on their manifestation (defined in [Section 2.2.1](#)). This lets us keep flaws that we automatically uncover (see [Section 4.2.1](#)) separate; we summarize these based on their domain (defined in [Section 2.2.2](#)) in [Section 5.2](#). We list *all* these flaws in [Appendix C](#) to balance completeness and brevity and denote implicit relations with the phrase “implied by” in [Tables C.1](#) to [C.3](#) as described in [Section 2.3](#). Moreover, certain “subsets” of testing contain many interconnected flaws, which we present in subsections as a “third view” to keep related information together. Therefore, the counts of flaws given in [Tables 5.1](#) and [5.2](#) can be thought of as the sum the flaws we describe by manifestation, domain, and “subset”. These subsets include operational (acceptance) testing ([Section 5.3](#)), recovery testing ([Section 5.4](#)), scalability testing ([Section 5.5](#)), and compatibility testing ([Section 5.6](#)). Finally, we also infer some flaws as described in [Section 2.3](#), which do not contribute to any counts due to their subjectivity; we list them in [Appendix C.3](#) for completeness.

5.1 Flaws by Manifestation

The following sections list observed flaws grouped by *how* they manifest as presented in [Section 2.2.1](#). These include mistakes ([Section 5.1.1](#)), omissions ([Section 5.1.2](#)), contradictions ([Section 5.1.3](#)), ambiguities ([Section 5.1.4](#)), overlaps ([Section 5.1.5](#)), and redundancies ([Section 5.1.6](#)).

5.1.1 Mistakes

There are many ways that information can be incorrect which we identify in [Table 5.3](#). We provide an example below for those that are less straightforward; see [Appendix C.1.1](#) for the full list of mistakes.

Information is incorrect based on an assertion from another source

Since errors are distinct from defects/faults ([ISO/IEC and IEEE, 2010](#), pp. 128, 140; [Washizaki, 2024](#), p. 12-3; [van Vliet, 2000](#), pp. 399–400), error guessing should instead be called “defect guessing” if it is based on a “checklist of potential defects”

Table 5.3: Different kinds of mistakes found in the literature.

Description	Count
Information is incorrect based on an assertion from another source	10
Information is provided with an incorrect scope	7
Information is not present where it is claimed to be	6
Information contains a minor mistake	3 ^a
Incorrect information makes other information incorrect	2

^a Comprises two typos and one duplication.

(ISO/IEC and IEEE, 2021c, p. 29) or “fault guessing” if it is a “fault-based technique” (Bourque and Fairley, 2014, p. 4-9) that “anticipate[s] the most plausible faults in each SUT” (Washizaki, 2024, 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 van Vliet, 2000, p. 399).

Information is provided with an incorrect scope

“Par sheet testing” from (Hamburg and Mogyorodi, 2024) seems to refer to the specific example mentioned in [Mistake 11](#) and does not seem more widely applicable, since a “PAR sheet” is “a list of all the symbols on each reel of a slot machine” (Bluejay, 2024). Conversely, the terms “acceleration tolerance testing” and “acoustic tolerance testing” do not seem refer to software testing, but Fire-smith (2015, p. 56) includes them regardless. Elsewhere, they seem to refer to testing the acoustic tolerance of rats (Holley et al., 1996) or the acceleration tolerance of astronauts (Morgun et al., 1999, p. 11), aviators (Howe and Johnson, 1995, pp. 27, 42), or catalysts (Liu et al., 2023, p. 1463), which which are not relevant to software testing.

Incorrect information makes other information incorrect

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, (see [Mistake 28](#)) casts doubt on the claim that “red-box testing” is a synonym for “acceptance testing” ([Sneed and Göschl, 2000, p. 18](#)) (see [Mistake 29](#)).

OG Hetzel88

5.1.2 Omissions

We find four cases where a definition is omitted, one where a category is omitted, and one where a term (along with its relations) is omitted; we list these in [Appendix C.1.2](#).

5.1.3 Contradictions

There are many cases where multiple sources of information (sometimes within the same document!) disagree. We find this happen with six categories, six synonym relations, seven parent-child relations, 17 definitions, and two labels. These can be found in [Appendix C.1.3](#).

5.1.4 Ambiguities

Some information given in the literature is unclear; there is definitely something “wrong”, but we cannot deduce the intent of the original author(s). We identify the kinds of ambiguous information given in [Table 5.4](#); see [Appendix C.1.4](#) for the full list of ambiguities.

Table 5.4: Different kinds of ambiguities found in the literature.

Description	Count
A term is defined ambiguously	8
A term is used inconsistently	3
The distinction between two terms is unclear	2

5.1.5 Overlaps

While information given in the literature should be atomic, this is not always the case. We find three definitions that overlap, two terms with multiple definitions, and three terms that share acronyms. We list these in [Appendix C.1.5](#); note that we track two of the terms with multiple definitions as the same flaw since they are related.

5.1.6 Redundancies

We find redundancies in two parent-child relations, two definitions, and two labels as listed in [Appendix C.1.6](#).

5.2 Flaws by Domain

The following sections present flaws that we detect automatically (see [Section 4.2.1](#)) grouped by *what* information is flawed as presented in [Section 2.2.2](#). We also provide more detailed information on specific areas of these domains that may require further investigation. The domains we focus on here are test approach categories ([Section 5.2.1](#)), synonym relations ([Section 5.2.2](#)), and parent-child relations ([Section 5.2.3](#)).

5.2.1 Approach Category Flaws

While the IEEE categorization of testing approaches described in [Table 2.1](#) 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 (see [Table 7.1](#)). The most blatant example of this is that [ISO/IEC and IEEE \(2017, p. 286\)](#) 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” ([ISO/IEC and IEEE, 2021c, 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; we track these differences so we can analyze them more systematically.

In particular, we automatically detect test approaches with more than one category that violate our assumption of orthogonality (see [Section 2.1.1](#)). We identify 28 such cases that we summarize in [Table 5.5](#) and list along with their sources in [Table C.1](#) for completeness. These flaws mostly involve the category of “test technique”, which shows up in 23 of them. This may simply be because authors use the term “test technique” in a more general sense where we would use the term “test approach”. A more specific reason for this is how the line between “test technique” and “test practice” can be blurred when these categories are *not* transitive. For example, “some test practices, such as exploratory testing or model-based testing are sometimes [incorrectly] referred to as ‘test techniques’ ... as they are not themselves providing a way to create test cases, but instead use test design techniques to achieve that” (ISO/IEC and IEEE, 2022, p. 11; 2021a, p. 5). This seems to be the case for the following approaches, which may explain these apparent “flaws”:

- Exploratory testing ([2022](#), p. 33; [2021c](#), p. viii; [2013](#), p. 13)
- Experience-based testing ([2022](#), p. 4; [2021c](#), pp. viii, 4; [2013](#), p. 33)
- Scripted testing ([2022](#), p. 33)
- Ad hoc testing ([Washizaki, 2024](#), p. 5-14; [Hamburg and Mogyorodi, 2024](#); [Kam, 2008](#), p. 42)
- Model-based testing ([Engström and Petersen, 2015](#), pp. 1–2; [Kam, 2008](#), p. 4)

As described in [Section 2.3](#), we infer that child approaches inherit their parents’ categories. However, there seem to be exceptions to this, which may indicate that categories are *not* transitive or that they are *except* in certain cases. For example, the practice of experience-based testing has many subtechniques, as described above, but it *also* has subpractices, including tours ([ISO/IEC and IEEE, 2022](#), p. 34) and exploratory testing (although the latter is categorized inconsistently; see [Table C.1](#) and above discussion). Similarly, the conflicting categorizations of beta testing in [Table C.1](#) may propagate to its children closed beta testing and open beta testing. When we infer these flaws, we exclude them from [Tables 5.5](#) and [C.1](#) and instead include them in [Table C.3](#) for completeness.

Table 5.5: Summary of pairs of categories assigned to a test approach.

Categories	Count
Technique/Level	2
Technique/Practice	8
Technique/Type	13
Level/Type	5
Total	28

5.2.2 Synonym Relation Flaws

While synonyms do not inherently signify a flaw (as we discuss in [Section 2.1.2](#)), the software testing literature is full of incorrect and ambiguous synonyms that do. As described in [Section 2.2.2](#), we pay special attention to synonyms between independently defined approaches (which may be flaws) and to intransitive synonyms (which definitely *are* flaws). We present explicit (see [Section 2.3](#)) synonym relations that fit either of these criteria in [Figure 5.2](#), which we automatically generate from [our test approach glossary](#) and manually modify for legibility. These relations are given as described by the literature and are therefore flawed. We provide the full list of synonyms that violate transitivity (along with their sources) in [Appendix C.2.2](#) and discuss and other kinds of flawed synonym relations in [Sections 5.2.3](#), [5.5](#), and [5.6](#).

Later: Ensure this is up to date

5.2.3 Parent-Child Relation Flaws

Parent-child relations are also not immune to flaws. For example, some approaches are given as parents of themselves which violates the irreflexivity of this relation as defined in [Section 2.1.3](#). We identify the following three examples through automatic analysis of our generated graphs (see [Section 4.2.1](#)):

1. Performance Testing ([Gerrard, 2000a](#), Tab. 2; [2000b](#), Tab. 1)
2. System Testing ([Firesmith, 2015](#), p. 23)
3. Usability Testing ([Gerrard, 2000a](#), Tab. 2; [2000b](#), Tab. 1)

Interestingly, performance testing is *not* described as a subapproach of usability testing by ([Gerrard, 2000a;b](#)), which would have been more meaningful information to capture.



Figure 5.2: Visually meaningful synonym relations given explicitly by the literature.

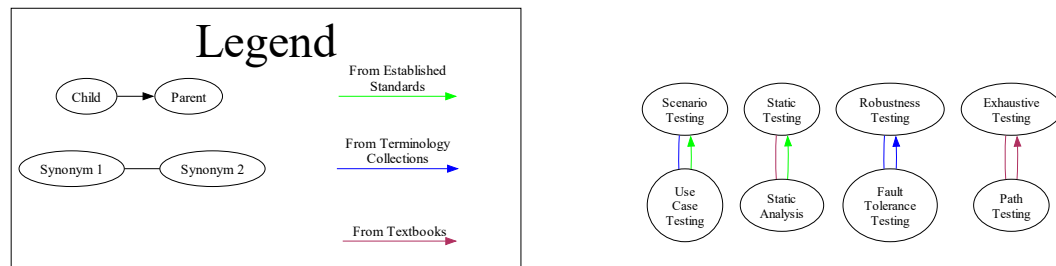


Figure 5.3: Pairs of test approaches with a parent-child *and* synonym relation given explicitly by the literature.

There are also pairs of synonyms where one is a subapproach of the other; these relations cannot coexist as synonym relations are symmetric while parent-child relations are asymmetric as outlined in [Sections 2.1.2](#) and [2.1.3](#), respectively. We identify 18 of these pairs through automatic analysis of our generated visualizations as described in [Section 4.2.1](#). We visualize the pairs where both relations are explicit in [Figure 5.3](#) and list all identified pairs in [Table C.2](#), as well as pairs where we infer a flaw in [Appendix C.3.3](#) for completeness. Of particular note is the relation between path testing and exhaustive testing. While [van Vliet \(2000, p. 421\)](#) claims that path testing done completely “is equivalent to exhaustively testing the program”², this overlooks the effects of input data ([ISO/IEC and IEEE, 2021c](#), pp. 129–130; [Patton, 2006](#), p. 121; [Peters and Pedrycz, 2000](#), 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” ([ISO/IEC and IEEE, 2022](#), p. 4, emphasis added; similar in [Hamburg and Mogyorodi, 2024](#); [Patton, 2006](#), p. 121).

5.3 Operational (Acceptance) Testing

(CONTRA, LABELS)

There are two names that the literature gives to this test approach:

- *Operational Acceptance Testing (OAT)* ([ISO/IEC and IEEE, 2022](#), p. 22; [Hamburg and Mogyorodi, 2024](#)) and
- *Operational Testing (OT)* ([ISO/IEC, 2018](#); [ISO/IEC and IEEE, 2017](#), p. 303; [Washizaki, 2024](#), p. 6-9, in the context of software engineering operations; [Bourque and Fairley, 2014](#), pp. 4-6, 4-9).

(CONTRA, SYNS)

[Firesmith \(2015, p. 30\)](#) lists the above terms separately, but they are considered synonyms elsewhere ([LambdaTest, 2024](#); [Bocchino and Hamilton, 1996](#)); since

²The contradictory definitions of path testing given in [Contradiction 18](#) add another layer of complexity to this claim.

[Firesmith](#) does not define these terms, it is hard to evaluate his distinction.

5.4 Recovery Testing

“Recovery testing” is “testing ... aimed at verifying software restart capabilities after a system crash or other disaster” ([Washizaki, 2024](#), p. 5-9) including “re-cover[ing] the data directly affected and re-establish[ing] the desired state of the system” ([ISO/IEC, 2023a](#); similar in [Washizaki, 2024](#), p. 7-10) so that the system “can perform required functions” ([ISO/IEC and IEEE, 2017](#), p. 370). However, the literature also describes similar test approaches with vague or non-existent distinctions between them. We describe these approaches and their flaws here and present the relations between them in [Figure 6.1a](#).

- *Recoverability testing* evaluates “how well a system or software can recover data during an interruption or failure” ([Washizaki, 2024](#), p. 7-10; similar in [ISO/IEC, 2023a](#)) and “re-establish the desired state of the system” ([2023a](#)). [Kam \(2008, p. 47\)](#) gives this as a synonym for “recovery testing”.
- *Disaster/recovery testing* evaluates if a system can “return to normal operation after a hardware or software failure” ([ISO/IEC and IEEE, 2017](#), 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” ([2021c](#), p. 37).
 - (OVER, DEFS) 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, making this term nonatomic.
- *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” ([ISO/IEC and IEEE, 2013](#), p. 2). This may be what is meant by “recovery testing” in the context of performance-related testing ([2022](#), Fig. 2).
- *Backup/recovery testing* determines the ability of a system “to restor[e] from back-up memory in the event of failure, without transfer[ing] to a different operating site or back-up system” ([ISO/IEC and IEEE, 2021c](#), p. 37).
 - (CONTRA, PARS) This given as a subtype of “disaster/recovery testing” which tests if “operation of the test item can be transferred to a different operating site” ([2021c](#), p. 37), even though this is *explicitly* excluded from its definition on the same page!
 - (OVER, LABELS) Its name is also quite similar to “backup and recovery testing”, adding further confusion.
- *Failover/recovery testing* determines the ability “to mov[e] to a back-up system in the event of failure, without transfer[ing] to a different operating site” ([ISO/IEC and IEEE, 2021c](#), p. 37).

- (CONTRA, PARS) This is also given as a subtype of “disaster/recovery testing” which tests if “operation of the test item can be transferred to a different operating site” (p. 37), even though this is *explicitly* excluded from its definition on the same page!
- (AMBI, PARS) While not explicitly related to recovery, *failover testing* “validates the SUT’s ability to manage heavy loads or unexpected failure to continue typical operations” (Washizaki, 2024, p. 5-9) by entering a “backup operational mode in which [these responsibilities] ... are assumed by a secondary system” (Hamburg and Mogyorodi, 2024). Its name implies that it is a child of “failover/recovery testing” but its definition makes it more broad (as it includes handling “heavy loads” where failover/recovery testing does not) which may reverse the direction of this relation.
- (AMBI, SYNS) Firesmith (2015, p. 56) uses the term “failover and recovery testing” which may be a synonym of “failover/recovery testing”.
- *Restart & recovery (testing)* is listed as a test approach by Gerrard (2000a, Fig. 5) but is not defined (MISS, DEFS) and may simply be a synonym to “recovery testing” (AMBI, SYNS).

5.5 Scalability Testing

(CONTRA, SYNS)

ISO/IEC and IEEE (2021c, p. 39) give “scalability testing” as a synonym of “capacity testing” while other sources differentiate between the two (Firesmith, 2015, p. 53; Bas, 2024, pp. 22–23).

(CONTRA, DEFS)

ISO/IEC and IEEE (2021c, p. 39) also include the external modification of the system as part of “scalability” but ISO/IEC (2023a) describe 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.

(WRONG, LABELS)

Washizaki (2024, p. 5-9) 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”. This seems to define “usability testing” with elements of functional and recovery testing.

5.6 Compatibility Testing

(OVER, DEFS)

“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)” (ISO/IEC and IEEE, 2022, p. 3). This definition is nonatomic as it combines the ideas of “co-existence” and “interoperability”.

(WRONG, SYNS)

The “interoperability” element of “compatibility testing” is explicitly excluded by ISO/IEC and IEEE (2021c, p. 37), (incorrectly) implying that “compatibility testing” and “co-existence testing” are synonyms.

(AMBI, SYNS)

Furthermore, the definition of “compatibility testing” in Kam (2008, p. 43) unhelpfully says “see *interoperability testing*”, adding another layer of confusion to the direction of their relationship.

(WRONG, LABELS)

ISO/IEC and IEEE (2022, pp. 22, 43) say “interoperability testing helps confirm that applications can work on multiple operating systems and devices”, but this seems to instead describe “portability testing”, which evaluates the “capability of a product to be adapted to changes in its requirements, contexts of use, or system environment” (ISO/IEC, 2023a; similar in ISO/IEC and IEEE, 2022, p. 7; 2017, pp. 184, 329; Hamburg and Mogyorodi, 2024), such as being “transferred from one hardware ... environment to another” (ISO/IEC and IEEE, 2021c, p. 39).

Chapter 6

Recommendations

As we have shown in [Chapter 5](#), “testing is a mess” (Mosser, 2023, priv. comm.)! It will take a lot of time, effort, expertise, and training to organize these terms (and their relations) logically. However, the hardest step is often the first one, so we attempt to give some examples of how this “rationalization” can occur. These changes often arise when we notice an issue with the current state of the terminology and think about what *we* would do to make it better. We do not claim that these are correct, unbiased, or exclusive, just that they can be used as an inspiration for those wanting to pick up where we leave off.

When redefining terms, we seek to make 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 ; failover/recovery testing, failover and recovery testing, and failover testing are all given separately)
3. Consistent (e.g., backup/recovery testing and failover/recovery testing explicitly exclude an aspect included in its parent disaster/recovery testing)

Likewise, we seek to eliminate classes of flaws that can be detected automatically, such as test approaches that are given as synonyms to multiple distinct approaches ([Appendix C.2.2](#)) or as parents of themselves ([Section 5.2.3](#)), or pairs of approaches with both a parent-child *and* synonym relation ([Section 5.2.3](#)).

We give recommendations for the areas of operational (acceptance) testing ([Section 6.1](#)), recovery testing ([Section 6.2](#)), scalability testing ([Section 6.3](#)), performance-related testing ([Section 6.4](#)), and compatibility testing ([Section 6.5](#)). We provide graphical representations (see [Section 4.1](#)) of these subsets when helpful in [Figures 6.1 to 6.3](#), in which arrows representing relations between approaches are coloured based on the source tier (see [Section 2.5](#)) that defines them. We colour all proposed approaches and relations orange. We also include inferred approaches and relations in grey for completeness, although they are not explicitly given in the literature (see [Section 2.3](#)).

Same label to different phantomsections; is that OK?

6.1 Operational (Acceptance) Testing

Since this terminology is not standardized (see [Section 5.3](#)), we propose that “Operational Acceptance Testing (OAT)” and “Operational Testing (OT)” are treated as synonyms for a type of acceptance testing ([ISO/IEC and IEEE, 2022](#), p. 22; [Hamburg and Mogyorodi, 2024](#)) that focuses on “non-functional” attributes of the system ([LambdaTest, 2024](#)). Indeed, this is how we track this approach in [our test approach glossary](#)! We define it as “test[ing] to determine the correct installation, configuration and operation of a module and that it operates securely in the operational environment” ([ISO/IEC, 2018](#)) or to “evaluate a system or component in its operational environment” ([ISO/IEC and IEEE, 2017](#), p. 303), particularly “to determine if operations and/or systems administration staff can accept [it]” ([Hamburg and Mogyorodi, 2024](#)).

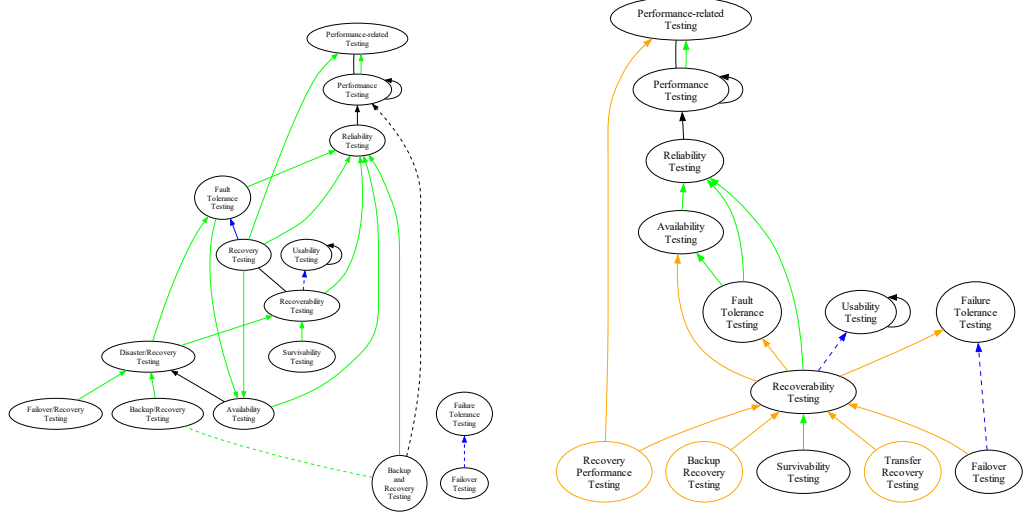
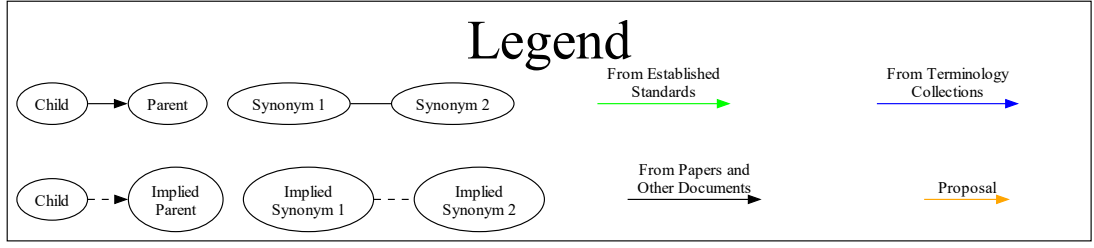
find more academic sources

6.2 Recovery Testing

To remedy the flaws we describe in [Section 5.4](#), we recommend that the literature uses these terms more consistently, resulting in the improved graph in [Figure 6.1b](#). The following proposals “recapture” information from the literature more consistently:

1. Prefer the term “recoverability testing” over “recovery testing” to indicate its focus on a system’s *ability* to recover, not its *performance* of recovering ([Kam, 2008](#), p. 47). “Recovery testing” may be an acceptable synonym, since it seems to be more prevalent in the literature.
2. Introduce the term *recovery performance testing* when evaluating performance metrics of a system’s recovery as a subapproach of recoverability testing and performance-related testing¹ ([ISO/IEC and IEEE, 2022](#), Fig. 2; [2013](#), p. 2).
3. Introduce separate terms for the different methods of recovery which are all subapproaches of recoverability testing:
 - (a) from backup memory (*backup recovery testing*) ([ISO/IEC and IEEE, 2021c](#), p. 37; [2013](#), p. 2),
 - (b) from a back-up system (*failover testing*) ([Washizaki, 2024](#), p. 5-9; [Hamburg and Mogyorodi, 2024](#)), or
 - (c) by transferring operations elsewhere (*transfer recovery testing*) ([ISO/IEC and IEEE, 2021c](#), p. 37).

¹See [Section 6.4](#).



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

Figure 6.1: Visualizations of relations between terms related to recovery testing.

6.3 Scalability Testing

We describe the issues with scalability testing terminology in [Section 5.5](#) and show the relations between related approaches (as described by the literature) in [Figure 6.2a](#). These flaws are resolved and/or explained by other sources! Taking this extra information into account results in [Figure 6.2b](#) and provides a more accurate description of scalability testing.

(CONTRA, SYNS)

ISO/IEC and IEEE (2021c, p. 39) define “scalability testing” as the testing of a system’s ability to “perform under conditions that may need to be supported in the future”. This focus on “the future” is supported by [Hamburg and Mogyorodi \(2024\)](#), who define “scalability” as “the degree to which a component or system can be adjusted for changing capacity”; the original source they reference agrees, defining it as “the measure of a system’s ability to be upgraded to accommodate increased loads” ([Gerrard and Thompson, 2002](#), p. 381). 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” ([ISO/IEC, 2023a](#)). Therefore, it these terms should *not* be synonyms, as done by [Firesmith](#)

(2015, p. 53) and Bas (2024, pp. 22–23).

(CONTRA, DEFS)

There seems to be an underlying reason that sources disagree on whether external modification of the system is part of scalability testing. Washizaki (2024, p. 5-9) claim that one objective of elasticity testing is “to evaluate scalability”, implying that (ISO/IEC, 2023a)’s notion of “scalability” actually refers to “elasticity”! This makes sense in the context of other definitions provided by Washizaki (2024):

- **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 Testing²:** 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 consistent with how the terms are used in industry: Pandey (2023) 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”.

(WRONG, LABELS)

Washizaki (2024, p. 5-9)’s definition of “scalability testing” is completely separate from the definitions of “scalability”, “capacity”, and “elasticity” discussed above! Therefore, this definition and its related synonym relations should simply be disregarded since they are inconsistent with the rest of the literature.

6.4 Performance(-related) Testing

“Performance testing” is defined as testing “conducted to evaluate the degree to which a test item ... accomplishes its designated functions” (ISO/IEC and IEEE, 2022, p. 7; 2021a, p. 2; 2017, p. 320; similar in 2021c, pp. 38-39; Moghadam, 2019, p. 1187). It does this by “measuring the performance metrics” (Moghadam, 2019, p. 1187; similar in Hamburg and Mogyorodi, 2024) (such as the “system’s capacity for growth” (Gerrard, 2000b, p. 23)), “detecting the functional problems appearing under certain execution conditions” (Moghadam, 2019, p. 1187), and “detecting violations of non-functional requirements under expected and stress conditions” (Moghadam, 2019, p. 1187; similar in Washizaki, 2024, p. 5-9). It is performed either ...

²This definition seems correct but its source is unclear; see [Mistake 9](#).

See [#35](#)

1. “within given constraints of time and other resources” (ISO/IEC and IEEE, 2022, p. 7; 2017, p. 320; similar in Moghadam, 2019, p. 1187), or
2. “under a ‘typical’ load” (ISO/IEC and IEEE, 2021c, 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’ ” (2021c, p. 38), along with other approaches like load and capacity testing (ISO/IEC and IEEE, 2022, p. 22). Note that “performance, load and stress testing might considerably overlap in many areas” (Moghadam, 2019, p. 1187). In contrast, Washizaki (2024, 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 (Washizaki, 2024, p. 5-9; Moghadam, 2019, p. 1187), as noted in Section 5.2.2. 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.



Figure 6.2: Visualizations of relations between terms related to scalability testing.

defined independently of its children: “co-existence testing” and “interoperability testing”.

Chapter 7

Threats to Validity

A case study is valid if its “results are true and not biased by the researchers’ subjective point of view” (Runeson and Höst, 2009, p. 153). Since our research is a case study on software testing terminology, we do our best to accomplish this. However, to present and explain our findings, we make some decisions with which others may disagree. The benefit of our work being open-source means that others can make their own decisions, modify our data/code appropriately, and observe how this changes the results. Despite this and our best efforts, some threats to validity remain; we use Runeson and Höst (2009, pp. 153–154)’s definitions of threats to validity and list them in their respective subsections as follows:

1. **Reliability:** If another researcher were to conduct the case study later on and get the same result, then the case study is reliable. (Section 7.1)
2. **Internal Validity:** “When ... investigating whether one factor affects an investigated factor[,] there is a risk that the investigated factor is also affected by a third factor”. This threatens internal validity “if the researcher is not aware of the third factor and/or does not know to what extent it affects the investigated factor”. (Section 7.2)
3. **External Validity:** For results of a study to be externally valid, they should be generalizable and “of interest to other people outside the investigated case”. (Section 7.3)
4. **Construct Validity:** This refers to the extent to which “the operational measures that are studied really represent what the researcher[s] have in mind and what is investigated”, or how well the studied data aligns with the data that the researchers *intended* to study. (Section 7.4)

7.1 Reliability

Important: Fill in

7.2 Internal Validity

- **Overloaded Terms:** The ambiguity of natural language means that terms are often overloaded. In [Table 2.1](#), we describe the categories we use, but these terms may not be used in the same way universally; “test type”, “test technique”, “test practice”, and “test approach” could all reasonably be used as synonyms. For example, [Kam \(2008, p. 45, emphasis added\)](#) defines interface testing as “an integration *test type* that is concerned with testing ... interfaces”, but since he does not define “test type”, this may not have special significance.
- **Similar Approach Categories:** While our categorization in [Table 2.1](#) is used fairly consistently, other sources (such as [Washizaki, 2024](#); [Barbosa et al., 2006](#)) propose similar yet distinct categories that clash or overlap with them. We give these alternate categorizations in [Table 7.1](#), which could simply map to their “IEEE Equivalents” or could provide new perspectives and be useful in some contexts, either in place of or in tandem with ours.
- **Alternate Approach Categories:** Similarly, we find many other criteria for categorizing test approaches in the literature. In general, these are defined less systematically but are more fine-grained, seeming to “specialize” our categories from [Table 2.1](#). The existence of these categorizations is not inherently wrong, as they may be useful for specific teams or in certain contexts. For example, functional testing and structural testing “use different sources of information and have been shown to highlight different problems”, and deterministic testing and random testing have “conditions that make one approach more effective than the other” ([Washizaki, 2024, p. 5-16](#)). Unfortunately, even these alternate categories are not used consistently (see [Contradiction 13](#))!
- **Unclear Category Boundaries:** Some of these categories can also be grouped into “classes” or “families” such as the classes of combinatorial ([ISO/IEC and IEEE, 2021c, p. 15](#)) and data flow testing (p. 3) and the family of performance-related testing ([Moghadam, 2019, p. 1187](#))¹. Note that “there is a lot of overlap between different classes of testing” ([Fire-smith, 2015, p. 8](#)), meaning that “one category [of test techniques] might deal with combining two or more techniques” ([Washizaki, 2024, p. 5-10](#)). A side effect of this is that it is difficult to “untangle” these classes; for example, take the following sentence: “whitebox fuzzing extends dynamic test generation based on symbolic execution and constraint solving from unit testing to whole-application security testing” ([Godefroid and Luchaup, 2011, p. 23](#))!

¹The original source describes “performance testing ... as a family of performance-related testing techniques”, but it makes more sense to consider “performance-related testing” as the “family” with “performance testing” being one of the variabilities (see [Section 6.4](#)).

Table 7.1: Categories of testing given by other sources.

Term	Definition	Examples	IEEE Equivalent
Level (objective-based) ^a	Test levels based on the purpose of testing (Washizaki, 2024, p. 5-6) that “determine how the test suite is identified ... regarding its consistency ... and its composition” (p. 5-2)	conformance testing, installation testing, regression testing, performance testing, security testing (Washizaki, 2024, pp. 5-7 to 5-9)	Type
Phase	none given by Perry (2006) or Barbosa et al. (2006)	unit testing, integration testing, system testing, regression testing (Perry, 2006, p. 221; Barbosa et al., 2006, p. 3)	Level
Procedure	The basis for how testing is performed that guides the process; “categorized in [to] testing methods, testing guidances ^b and testing techniques” (Barbosa et al., 2006, p. 3)	none given generally by Barbosa et al. (2006); see “Technique”	Approach
Process	“A sequence of testing steps” (Barbosa et al., 2006, p. 2) “based on a development technology and ... paradigm, as well as on a testing procedure” (p. 3)	none given by Barbosa et al. (2006)	Practice
Stage	An alternative to the “traditional ... test stages” based on “clear technical groupings” (Gerrard, 2000a, p. 13)	desktop development testing, infrastructure testing, post-deployment monitoring (Gerrard, 2000a, p. 13)	Level
Technique	“Systematic procedures and approaches for generating or selecting the most suitable test suites” (Washizaki, 2024, p. 5-10)	specification-based testing, structure-based testing, fault-based testing ^c (Washizaki, 2024, pp. 5-10, 5-13 to 5-15)	Technique

^a See Contradiction 22.^b Testing methods and guidances are omitted from this table since Barbosa et al. (2006) do not define or give examples of them.^c Synonyms for these examples are used by Souza et al. (2017, p. 3; OG Mathur, 2012) and Barbosa et al. (2006, p. 3).

7.3 External Validity

While we document issues with “standardized” software testing terminology so that they can be addressed in the future, some dismiss the importance of standardized terminology for various reasons, including the following:

- **Limitations of Standardized Terminology:** [Schoots \(2014\)](#) holds that “common terminology is dangerous” and “to be able to truly understand each other, we need to ask questions and discuss in depth”. However, these in-depth discussions are *not* mutually exclusive with common terminology! Having a shared understanding for how terms are defined allows for common ground during these sorts of discussions with the chance to adapt them to serve the context of a given team or project (which we give examples of in [Sections 2.1.2](#) and [7.2](#)).
- **Standards Being Mandated:** [Schoots \(2014\)](#) also states that while he wishes “standards would be guidelines, ... reality shows standards become mandatory often”. He supports this with examples from [Soundararajan \(2015\)](#) where “contracts and bids from large companies” often “reference[] ISO linking to industry best practices”. This claim, however, overlooks:
 1. the possibility of renegotiating contracts and
 2. the notion of “tailored conformance” to standards, which is mentioned throughout the family of standards these authors critique ([ISO/IEC and IEEE, 2021a](#), pp. 9-10; [2021b](#), pp. 5, 17, 37; [2021c](#), p. 7) and was perhaps introduced later to address concerns such as these.

7.4 Construct Validity

Important: Fill in

Chapter 8

Future Work

Ideally, our work would have resulted in a complete taxonomy of software testing terminology based on the literature. Unfortunately, we were not able to fully accomplish this due to time constraints. With more time, we would continue iterating over undefined terms (Section 8.1) and investigate terms we *expected* to find but never did (Section 8.2). Additionally, there is much more work we could do to analyze our data on the software testing literature, such as detecting more classes of flaws (and identifying them!) (Section 8.3).

8.1 Iterating Over Undefined Terms

As we explain in Section 3.4, our methodology includes performing miniature literature reviews on undefined terms to record their missing definitions (and any relations). We were able to do this for the following approaches, although some are out of scope, such as EManations SEcURITY (EMSEC) testing, aspects of Orthogonal Array Testing (OAT) and loop testing (see Appendix A.1), and HTML testing (see Appendix A.5). We investigate the following terms (and their respective related terms) in the sources given:

- **Assertion Checking:** Lahiri et al. (2013); Chalin et al. (2006); Berdine et al. (2006)
- **Loop Testing¹:** Dhok and Ramanathan (2016); Godefroid and Luchaup (2011); Preuße et al. (2012); Forsyth et al. (2004)
- **EMSEC Testing:** Zhou et al. (2012); ISO (2021)
- **Asynchronous Testing:** Jard et al. (1999)
- **Performance(-related) Testing:** Moghadam (2019)
- **Web Application Testing:** Doğan et al. (2014); Kam (2008)

¹ISO (2022) and ISO (2015) were used as reference for terms but not fully investigated, Pierre et al. (2017) and Trudnowski et al. (2017) were added as potentially in scope, and Dominguez-Pumar et al. (2020) and Goralski (1999) were added as out-of-scope examples.

- **HTML Testing:** Choudhary et al. (2010); Sneed and Göschl (2000); Gerrard (2000b)
- **Document Object Model (DOM) Testing:** Bajammal and Mesbah (2018)
- **Sandwich Testing:** Sharma et al. (2021); Sangwan and LaPlante (2006)
- **Orthogonal Array Testing**²: Mandl (1985); Valcheva (2013)
- **Backup Testing**³: Bas (2024)
- **SoS (Integration) Testing:** ISO/IEC and IEEE (2019b)

Applying our procedure shown in [Figure 3.1](#) to these sources uncovers 87 new approaches and 71 new definitions. These definitions are either for existing undefined approaches or new uncovered approaches; while not every new approach is presented alongside a definition, if we assume that each of these definitions is for a new approach, we can deduce that about 82% of added test approaches are defined. This indicates that this procedure leads to a higher proportion of defined terms (63% vs. 66%), as shown in [Figure 8.1](#), which helps verify that our procedure constructively uncovers *and* defines new terminology. With repeated iterations, this ratio would approach 100%, resulting in a (plausibly) complete taxonomy.

²Yu et al. (2011) and Tsui (2007) were added as out-of-scope examples.

³See [Section 5.4](#).



Figure 8.1: Breakdown of how many test approaches are undefined.

8.2 Missing Terms

In addition to these undefined terms, some terms do not appear in the literature at all! While most test approaches arise as a result of our snowballing approach, we each have preexisting knowledge of what test approaches exist (a form of experience-based testing, if you will). Test approaches that arise independently of snowballing may serve as starting points for continued research if we do not find them in the literature using our iterative approach. The following terms come from previous knowledge, conversations with colleagues, research for other projects, or ad hoc cursory research to see what other test approaches exist:

- | | |
|--|------------------------------------|
| 1. Chaos engineering | 10. Lunchtime attacks ^c |
| 2. Chosen-ciphertext attacks | 11. Parallel testing |
| 3. Concolic testing | 12. Property-based testing |
| 4. Concurrent testing ^a | 13. Pseudo-random bit testing |
| 5. Context-driven testing | 14. Rubber duck testing |
| 6. Destructive testing | 15. Sanity testing |
| 7. Dogfooding | 16. Scream testing |
| 8. Implementation-based testing ^b | 17. Shadow testing |
| 9. Interaction-based testing | 18. Situational testing |

^aThis seems to be distinct from “concurrency testing”.

^bThis may or may not be distinct from “implementation-oriented testing.”

^cIn previous meetings, Dr. Smith mentioned that with the number of test approaches that suggest that people just like to label everything as “testing”, he would not be surprised if something like “Monday morning testing” existed. While independently researching chosen-ciphertext attacks out of curiosity, this prediction of a time-based test approach came true with “lunchtime attacks”.

8.3 Detecting More Flaws

In addition to the classes of flaws we *do* detect automatically, we could detect many more if time permitted. We currently detect parent-child relations that violate ir-reflexivity (see [Section 4.2.1](#)) which can be thought of as cycles with length $n = 1$. Since parent-child relations should also be transitive (see [Section 2.1.3](#)), cycles of *any* size are flaws. Given the current way we generate visualizations of these relations (see [Section 4.1](#)), detecting cycles where $n = 2$ would be straightforward: if a parent-child relation *and* its inverse (i.e., $A \rightarrow B$ and $B \rightarrow A$ for test approaches with labels A and B) both exist in the generated \LaTeX file for a visualization, (at least) one of these parent-child flaws is incorrect since they contradict each other and we have found a cycle. The main reason this would be time-consuming would be deciding on how to format these findings, writing code to do so automatically

for use in this thesis, and resolving any issues that arise during this process. Detecting larger cycles would also be possible but would likely require the use of an additional tool to analyze the graph of parent-child relations.

Unfortunately, writing code to detect flaws is only half the battle. First, we need to identify classes of flaws we even want to detect, which may never be complete as new test approaches emerge with potentially new relations. The classes of flaws we detect emerged over time based on our observations of the literature. For example, our understanding of the “standard” test approach categories described in [Section 2.1.1](#) led to us being able to detect approaches that are categorized inconsistently in [Section 5.2.1](#) and [Appendix C.2.1](#). Performing more thorough analysis, both on the literature and on our collected data, will likely reveal more classes of flaws that can be tracked automatically but is infeasible for the purposes of this thesis.

Chapter 9

Development Process

The following is a rough outline of the steps I have gone through this far for this project:

- Start developing system tests (this was pushed for later to focus on unit tests)
- Test inputting default values as `floats` and `ints`
- Check constraints for valid input
- Check constraints for invalid input
- Test the calculations of:
 - `t_flight`
 - `p_land`
 - `d_offset`
 - `s`
- Test the writing of valid output
- Test for projectile going long
- Integrate system tests into existing unit tests
- Test for assumption violation of `g`
 - Code generation could be flawed, so we can't assume assumptions are respected
 - Test cases shouldn't necessarily match what is done by the code; for example, `g = 0` shouldn't really give a `ZeroDivisionError`; it should be a `ValueError`
 - This inspired the potential for [The Use of Assertions in Code](#)
- Test that calculations stop on a constraint violation; this is a requirement should be met by the software (see [Section 9.3](#))

- Test behaviour with empty input file
- Start creation of test summary (for `InputParameters` module)
 - It was difficult to judge test case coverage/quality from the code itself
 - This is not really a test plan, as it doesn't capture the testing philosophy
 - Rationale for each test explains why it supports coverage and how Drasil derived (would derive) it
- Start researching testing
- Implement generation of `__init__.py` files ([#3516](#))
- Start the **Generating Requirements** subproject

9.1 Improvements to Manual Test Code

Even though this code will eventually be generated by Drasil, it is important that it is still human-readable, for the benefit of those reading the code later. This is one of the goals of Drasil (see [#3417](#) for an example of a similar issue). As such, the following improvements were discovered and implemented in the manually created testing code:

- use `pytest`'s parameterization
- reuse functions/data for consistency
- improve import structure
- use `conftest` for running code before all tests of a module

9.1.1 Testing with Mocks

When testing code, it is common to first test lower-level modules, then assume that these modules work when testing higher-level modules. An example would be using an input module to set up test cases for a calculation module after testing the input module. This makes sense when writing test cases manually since it reduces the amount of code that needs to be written and still provides a reasonably high assurance in the software; if there is an issue with the input module that affects the calculation module tests, the issue would be revealed when testing the input module.

However, since these test cases will be generated by Drasil, they can be consistently generated with no additional effort. This means that the testing of each module can be done completely independently, increasing the confidence in the tests.

9.2 The Use of Assertions in Code

While assertions are often only used when testing, they can also be used in the code itself to enforce constraints or preconditions; they act like documentation that determines behaviour! For example, they could be used to ensure that assumptions about values (like the value for gravitational acceleration) are respected by the code, which gives a higher degree of confidence in the code. This process is known as “assertion checking” ([Lahiri et al., 2013](#)) .

investigate OG sources

9.3 Generating Requirements

I structured my manually created test cases around Projectile’s functional requirements, as these are the most objective aspects of the generated code to test automatically. One of these requirements was “Verify-Input-Values”, which said “Check the entered input values to ensure that they do not exceed the data constraints. If any of the input values are out of bounds, an error message is displayed and the calculations stop.” This led me to develop a test case to ensure that if an input constraint was violated, the calculations would stop ([Source Code D.1](#)).

However, this test case failed, since the actual implementation of the code did *not* stop upon an input constraint violation. This was because the code choice for what to do on a constraint violation ([Source Code D.2](#)) was “disconnected” from the manually written requirement ([Source Code D.3](#)), as described in [#3523](#).

Should I include the definition of Constraints?

This problem has been encountered before ([#3259](#)) and presented a good opportunity for generation to encourage reusability and consistency. However, since it makes sense to first verify outputs before actually outputting them and inserting generated requirements among manually created ones seemed challenging, it made sense to first generate an output requirement.

While working on Drasil in the summer of 2019, I implemented the generation of an input requirement across most examples ([#1844](#)). I had also attempted to generate an output requirement, but due to time constraints, this was not feasible. The main issue with this change was the desire to capture the source of each output for traceability; this source was attached to the `InstanceModel` (or rarely, `DataDefinition`) and not the underlying `Quantity` that was used for a program’s outputs. The way I had attempted to do this was to add the reference as a `Sentence` in a tuple.

Taking another look at this four years later allowed us to see that we should be storing the outputs of a program as their underlying models, allowing us to keep the source information with it . While there is some discussion about how this might change in the future , for now, all outputs of a program should be `InstanceModels`. Since this change required adding the `Referable` constraints to the output field of `SystemInformation`, the outputs of all examples needed to be updated to satisfy this constraint; this meant that generating the output requirement of each example was nearly trivial once the outputs were specified correctly. After modifying `DataDefinitions` in GlassBR that were outputs to be `InstanceModels` ([#3569](#); [#3583](#)), reorganizing the requirements of SWHS ([#3589](#); [#3607](#)), and clarifying

cite Dr. Smith

add refs to ‘underlying Theory’ comment and ‘not all outputs be IMs’ comment

add constraints

the outputs of SWHS ([#3589](#)), SglPend ([#3533](#)), DblPend ([#3533](#)), GamePhysics ([#3609](#)), and SSP ([#3630](#)), the output requirement was ready to be generated.

Chapter 10

Research

It was realized early on in the process that it would be beneficial to understand the different kinds of testing (including what they test, what artifacts are needed to perform them, etc.). This section provides some results of this research, as well as some information on why and how it was performed.

A justification for why we decided to do this should be added

10.1 Existing Taxonomies, Ontologies, and the State of Practice

One thing we may want to consider when building a taxonomy/ontology is the semantic difference between related terms. For example, one ontology found that the term “‘IntegrationTest’ is a kind of Context (with semantic of stage, but not a kind of Activity)” while “‘IntegrationTesting’ has semantic of Level-based Testing that is a kind of Testing Activity [or] ... of Test strategy” (Tebes et al., 2019, p. 157).

add acronym?

In (Souza et al., 2017), the ontology (ROoST) is made to answer a series of questions, including “What is the test level of a testing activity?” and “What are the artifacts used by a testing activity?” (Souza et al., 2017, pp. 8-9). The question “How do testing artifacts relate to each other?” (Souza et al., 2017, p. 8) is later broken down into multiple questions, such as “What are the test case inputs of a given test case?” and “What are the expected results of a given test case?” (Souza et al., 2017, p. 21). *These questions seem to overlap with the questions we were trying to ask about different testing techniques.* Tebes et al. (2019, pp. 152-153) may provide some sources for software testing terminology and definitions (this seems to include those suggested by Dr. Carette) in addition to a list of ontologies (some of which have been investigated).

is this punctuation right?

One software testing model developed by the Quality Assurance Institute (QAI) includes the test environment (“conditions ...that both enable and constrain how testing is performed”, including mission, goals, strategy, “management support, resources, work processes, tools, motivation”), test process (testing “standards and procedures”), and tester competency (“skill sets needed to test software in a test environment”) (Perry, 2006, pp. 5-6).

Another source introduced the notion of an “intervention”: “an act performed

(e.g. use of a technique¹ or a process change) to adapt testing to a specific context, to solve a test issue, to diagnose testing or to improve testing” (Engström and Petersen, 2015, p. 1) and noted that “academia tend[s] to focus on characteristics of the intervention [while] industrial standards categorize the area from a process perspective” (Engström and Petersen, 2015, p. 2). It provides a structure to “capture both a problem perspective and a solution perspective with respect to software testing” (Engström and Petersen, 2015, pp. 3-4), but this seems to focus more on test interventions and challenges rather than approaches (Engström and Petersen, 2015, Fig. 5).

10.2 Definitions

- Test case: “the specification of all the entities that are essential for the execution, such as input values, execution and timing conditions, testing procedure, and the expected outcomes” (Washizaki, 2024, pp. 5-1 to 5-2)
- Verification: “the process of evaluating a system or component to determine whether the products of a given development phase satisfy the conditions imposed at the start of that phase” (van Vliet, 2000, p. 400)
- Validation: “the process of evaluating a system or component during or at the end of the development process to determine whether it satisfies specified requirements” (van Vliet, 2000, p. 400)
- Test Suite Reduction: the process of reducing the size of a test suite while maintaining the same coverage (Barr et al., 2015, p. 519); can be accomplished through mutation testing
- Test Case Reduction: the process of “removing side-effect free functions” from an individual test case to “reduc[e] test oracle costs” (Barr et al., 2015, p. 519)

10.2.1 Documentation

- Verification and Validation (V&V) Plan: a document for the “planning of test activities” described by IEEE Standard 1012 (van Vliet, 2000, p. 411)
- Test Plan: “a document describing the scope, approach, resources, and schedule of intended test activities” in more detail than the V&V Plan (van Vliet, 2000, pp. 412-413); should also outline entry and exit conditions for the testing activities as well as any risk sources and levels (Peters and Pedrycz, 2000, p. 445)
- Test Design documentation: “specifies ... the details of the test approach and identifies the associated tests” (van Vliet, 2000, p. 413)

¹Not formally defined, but distinct from the notion of “test technique” described in Table 2.1.

- Test Case documentation: “specifies inputs, predicted outputs and execution conditions for each test item” ([van Vliet, 2000](#), p. 413)
- Test Procedure documentation: “specifies the sequence of actions for the execution of each test” ([van Vliet, 2000](#), p. 413)
- Test Report documentation: “provides information on the results of testing tasks”, addressing software verification and validation reporting ([van Vliet, 2000](#), p. 413)

10.3 General Testing Notes

- The scope of testing is very dependent on what type of software is being tested, as this informs what information/artifacts are available, which approaches are relevant, and which tacit knowledge is present. For example, a method table is a tool for tracking the “test approaches, testing techniques and test types that are required depending ... on the context of the test object” ([Hamburg and Mogyorodi, 2024](#)), although this is more specific to the automotive domain
- If faults exist in programs, they “must be considered faulty, even if we cannot devise test cases that reveal the faults” ([van Vliet, 2000](#), p. 401)
- “There is no established consensus on which techniques ... are the most effective. The only consensus is that the selection will vary as it should be dependent on a number of factors” ([ISO/IEC and IEEE, 2021c](#), p. 128; similar in [van Vliet, 2000](#), p. 440), and it is advised to use many techniques when testing (p. 440). This supports the principle of *independence of testing*: the “separation of responsibilities, which encourages the accomplishment of objective testing” ([Hamburg and Mogyorodi, 2024](#))

10.3.1 Test Oracles

A test oracle is a “source of information for determining whether a test has passed or failed” ([ISO/IEC and IEEE, 2022](#), p. 13) or that “the SUT behaved correctly ... and according to the expected outcomes” and can be “human or mechanical” ([Washizaki, 2024](#), p. 5-5). Oracles provide either “a ‘pass’ or ‘fail’ verdict”; otherwise, “the test output is classified as inconclusive” ([Washizaki, 2024](#), p. 5-5). This process can be “deterministic” (returning a Boolean value) or “probabilistic” (returning “a real number in the closed interval $[0, 1]$ ”) ([Barr et al., 2015](#), p. 509). Probabilistic test oracles can be used to reduce the computation cost (since test oracles are “typically computationally expensive”) ([Barr et al., 2015](#), p. 509) or in “situations where some degree of imprecision can be tolerated” since they “offer a probability that [a given] test case is acceptable” ([Barr et al., 2015](#), p. 510). The SWEBOK Guide V4 lists “unambiguous requirements specifications, behavioral models, and code annotations” as examples ([Washizaki, 2024](#), p. 5-5), and [Barr et al.](#) provides four categories ([2015](#), p. 510):

- Specified test oracle: “judge[s] all behavioural aspects of a system with respect to a given formal specification” (Barr et al., 2015, p. 510)
- Derived test oracle: any “artefact[] from which a test oracle may be derived—for instance, a previous version of the system” or “program documentation”; this includes regression testing, metamorphic testing (Barr et al., 2015, p. 510), and invariant detection (either known in advance or “learned from the program”) (Barr et al., 2015, p. 516)
 - This seems to prove “relative correctness” as opposed to “absolute correctness” (Lahiri et al., 2013, p. 345) since this derived oracle may be wrong!
 - “Two versions can be checked for semantic equivalence to ensure the correctness of [a] transformation” in a process that can be done “incrementally” (Lahiri et al., 2013, p. 345)
 - Note that the term “invariant” may be used in different ways (see (Chalin et al., 2006, p. 348))
- Pseudo-oracle: a type of derived test oracle that is “an alternative version of the program produced independently” (by a different team, in a different language, etc.) (Barr et al., 2015, p. 515). *We could potentially use the programs generated in other languages as pseudo-oracles!*
- Implicit test oracles: detect “‘obvious’ faults such as a program crash” (potentially due to a null pointer, deadlock, memory leak, etc.) (Barr et al., 2015, p. 510)
- “Lack of an automated test oracle”: for example; a human oracle generating sample data that is “realistic” and “valid”, (Barr et al., 2015, pp. 510-511), crowdsourcing (Barr et al., 2015, p. 520), or a “Wideband Delphi”: “an expert-based test estimation technique that ... uses the collective wisdom of the team members” (Hamburg and Mogyorodi, 2024)

see ISO 29119-11

10.3.2 Generating Test Cases

- “Impl[ies] a reduction in human effort and cost, with the potential to impact the test coverage positively”, and a given “policy could be reused in analogous situations which leads to even more efficiency in terms of required efforts” (Moghadam, 2019, p. 1187)
- A selected “test adequacy criterion can be used in the test selection process [as a ‘test case generator’]. If a 100% statement coverage has not been achieved yet, an additional test case is selected that covers one or more statements yet untested. This generative view is used in many test tools” (van Vliet, 2000, p. 402). This may be similar to the “test data generators” mentioned on (p. 410)

Investigate

- “Dynamic test generation consists of running a program while simultaneously executing the program symbolically in order to gather constraints on inputs from conditional statements encountered along the execution ([Godefroid and Luchaup, 2011](#), p. 23)

OG [11, 6]

- “Generating tests to detect [loop inefficiencies]” is difficult due to “virtual call resolution”, reachability conditions, and order-sensitivity ([Dhok and Ramathan, 2016](#), p. 896)
- Can be facilitated by “testing frameworks such as JUnit [that] automate the testing process by writing test code” ([Sakamoto et al., 2013](#), p. 344)
- Assertion checking requires “auxiliary invariants”, and while “many ... can be synthesized automatically by invariant generation methods, the undecidable nature (or the high practical complexity) of assertion checking precludes complete automation for a general class of user-supplied assertions” ([Lahiri et al., 2013](#), p. 345)
 - Differential Assertion Checking (DAC) can be supported by “automatic invariant generation” ([Lahiri et al., 2013](#), p. 345)

OG Halfond and Orso, 2007

- *Automated interface discovery* can be used “for test-case generation for web applications” ([Doğan et al., 2014](#), p. 184)

OG Artzi et al., 2008

- “Concrete and symbolic execution” can be used in “a dynamic test generation technique ... for PHP applications” ([Doğan et al., 2014](#), p. 192)

- COBRA is a tool that “generates test cases automatically and applies them to the simulated industrial control system in a SiL Test” ([Preuße et al., 2012](#), p. 2)
- Test case generation is useful for instances where one kind of testing is difficult, but can be generated from a different, simpler kind (e.g., asynchronous testing from synchronous testing ([Jard et al., 1999](#)))
- Since some values may not always be applicable to a given scenario (e.g., a test case for zero doesn’t make sense if there is a constraint that the value in question cannot be zero), the user should likely be able to select categories of tests to generate instead of Drasil just generating all possible test cases based on the inputs ([Smith and Carette, 2023](#))

Investigate!

- “Test suite augmentation techniques specialise in identifying and generating” new tests based on changes “that add new features”, but they could be extended to also augment “the expected output” and “the existing *oracles*” ([Barr et al., 2015](#), p. 516)
- The Fault Injection Security Tool (FIST) “automates fault injection analysis of software using program inputs, fault injection functions, and assertions in programs written in C and C++” ([Ghosh and Voas, 1999](#), pp. 40–41). “For

example, Booleans are corrupted to their opposite value during execution, integers are corrupted using a random function with a uniform distribution centered around their current value, character strings are corrupted using random values” (p. 41). Unfortunately, “we do not have automatic learning systems for protecting software states, though it is the subject of ongoing research” (p. 44).

10.4 Examples of Metamorphic Relations (MRs)

- The distance between two points should be the same regardless of which one is the “start” point (ISO/IEC and IEEE, 2021c, p. 22)
- “If a person smokes more cigarettes, then their expected age of death will probably decrease (and not increase)” (ISO/IEC and IEEE, 2021c, p. 22)
- “For a function that translates speech into text[,] ... the same speech at different input volume levels ... [should result in] the same text” (ISO/IEC and IEEE, 2021c, p. 22)
- The average of a list of numbers should be equal (within floating-point errors) regardless of the list’s order (Kanevala and Yueh Chen, 2019, p. 67)
- For matrices, if $B = B_1 + B_2$, then $A \times B = A \times B_1 + A \times B_2$ (Kanevala and Yueh Chen, 2019, pp. 68-69)
- Symmetry of trigonometric functions; for example, $\sin(x) = \sin(-x)$ and $\sin(x) = \sin(x + 360^\circ)$ (Kanevala and Yueh Chen, 2019, p. 70)
- Modifying input parameters to observe expected changes to a model’s output (e.g., testing epidemiological models calibrated with “data from the 1918 Influenza outbreak”); by “making changes to various model parameters ... authors identified an error in the output method of the agent based epidemiological model” (Kanevala and Yueh Chen, 2019, p. 70)
- Using machine learning to predict likely MRs to identify faults in mutated versions of a program (about 90% in this case) (Kanevala and Yueh Chen, 2019, p. 71)

10.5 Roadblocks to Testing

- Intractability: it is generally impossible to test a program exhaustively (ISO/IEC and IEEE, 2022, p. 4; Washizaki, 2024, p. 5-5; Peters and Pedrycz, 2000, pp. 439, 461; van Vliet, 2000, p. 421)
- Adequacy: to counter the issue of intractability, it is desirable “to reduce the cardinality of the test suites while keeping the same effectiveness in terms of coverage or fault detection rate” (Washizaki, 2024, p. 5-4) which is difficult to

do objectively; see also “minimization”, the process of “removing redundant test cases” (Washizaki, 2024, p. 5-4)

- Undecidability (Peters and Pedrycz, 2000, p. 439): it is impossible to know certain properties about a program, such as if it will halt (i.e., the Halting Problem (Gurfinkel, 2017, p. 4)), so “automatic testing can’t be guaranteed to always work” for all properties (Nelson, 1999)

Add paragraph/section number?

10.5.1 Roadblocks to Testing Scientific Software (Kanewala and Yueh Chen, 2019, p. 67)

- “Correct answers are often unknown”: if the results were already known, there would be no need to develop software to model them (Kanewala and Yueh Chen, 2019, p. 67); in other words, complete test oracles don’t exist “in all but the most trivial cases” (Barr et al., 2015, p. 510), and even if they are, the “automation of mechanized oracles can be difficult and expensive” (Washizaki, 2024, p. 5.5)
- “Practically difficult to validate the computed output”: complex calculations and outputs are difficult to verify (Kanewala and Yueh Chen, 2019, p. 67)
- “Inherent uncertainties”: since scientific software models scenarios that occur in a chaotic and imperfect world, not every factor can be accounted for (Kanewala and Yueh Chen, 2019, p. 67)
- “Choosing suitable tolerances”: difficult to decide what tolerance(s) to use when dealing with floating-point numbers (Kanewala and Yueh Chen, 2019, p. 67)
- “Incompatible testing tools”: while scientific software is often written in languages like FORTRAN, testing tools are often written in languages like Java or C++ (Kanewala and Yueh Chen, 2019, p. 67)

Out of this list, only the first two apply. The scenarios modelled by Drasil are idealized and ignore uncertainties like air resistance, wind direction, and gravitational fluctuations. There are not any instances where special consideration for floating-point arithmetic must be taken; the default tolerance used for relevant testing frameworks has been used and is likely sufficient for future testing. On a related note, the scientific software we are trying to test is already generated in languages with widely-used testing frameworks.

Add example

Add source(s)?

Chapter 11

Extras

Writing Directives

- What macros do I want the reader to know about?

11.1 Writing Directives

I enjoy writing directives (mostly questions) to navigate what I should be writing about in each chapter. You can do this using:

Source Code 11.1: Pseudocode: exWD

```
\begin{writingdirectives}
  \item What macros do I want the reader to know about?
\end{writingdirectives}
```

Personally, I put them at the top of chapter files, just after chapter declarations.

11.2 HREFs

For PDFs, we have (at least) 2 ways of viewing them: on our computers, and printed out on paper. If you choose to view through your computer, reading links (as they are linked in this example, inlined everywhere with “clickable” links) is fine. However, if you choose to read it on printed paper, you will find trouble clicking on those same links. To mitigate this issue, I built the “porthref” macro (see `macros.tex` for the definition) to build links that appear as clickable text when “compiling for computer-focused reading,” and adds links to footnotes when “compiling for printing-focused reading.” There is an option (`compilingforprinting`) in the `manifest.tex` file that controls whether PDF builds should be done for

computers or for printers. For example, by default, [McMaster](#) is made with clickable functionality, but if you change the `manifest.tex` option as mentioned, then you will see the link in a footnote (try it out!).

Source Code 11.2: Pseudocode: `exPHref`

```
\porthref{McMaster}{https://www.mcmaster.ca/}
```

11.3 Pseudocode Code Snippets

For pseudocode, you can also use the pseudocode environment, such as that used in [Source Code D.5](#).

11.4 TODOs

While writing, I plastered my thesis with notes for future work because, for whatever reason, I just didn't want to, or wasn't able to, do said work at that time. To help me sort out my notes, I used the `todonotes` [package](#) with a few extra macros (defined in `macros.tex`). For example,...

Important notes:

Important: “Important” notes.

Generic inlined notes:

Generic inlined notes.

Notes for later:

Some “easy” notes:

Easy: Easier notes.

Tedious work:

Needs time: Tedious notes.

Questions:

Later: TODO notes for later! For finishing touches, etc.

Q #9: Questions I might have?

Bibliography

- Mominul Ahsan, Stoyan Stoyanov, Chris Bailey, and Alhussein Albarbar. Developing Computational Intelligence for Smart Qualification Testing of Electronic Products. *IEEE Access*, 8:16922–16933, January 2020. ISSN 2169-3536. doi: 10.1109/ACCESS.2020.2967858. URL <https://www.webofscience.com/api/gateway?GWVersion=2&SrcAuth=DynamicDOIArticle&SrcApp=WOS&KeyAID=10.1109%2FACCESS.2020.2967858&DestApp=DOI&SrcAppSID=USW2EC0CB9ABcVz5BcZ70BCflmtJ&SrcJTitle=IEEE+ACCESS&DestDOIRegistrantName=Institute+of+Electrical+and+Electronics+Engineers>. Place: Piscataway.
- Paul Ammann and Jeff Offutt. *Introduction to Software Testing*. Cambridge University Press, Cambridge, United Kingdom, 2nd edition, 2017. ISBN 978-1-107-17201-2. URL <https://eopcw.com/find/downloadFiles/11>.
- Mohammad Bajammal and Ali Mesbah. Web Canvas Testing Through Visual Inference. In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*, pages 193–203, Västerås, Sweden, 2018. IEEE. ISBN 978-1-5386-5012-7. doi: 10.1109/ICST.2018.00028. URL <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=8367048>.
- Ellen Francine Barbosa, Elisa Yumi Nakagawa, and José Carlos Maldonado. Towards the Establishment of an Ontology of Software Testing. volume 6, pages 522–525, San Francisco, CA, USA, January 2006.
- Luciano Baresi and Mauro Pezzè. An Introduction to Software Testing. *Electronic Notes in Theoretical Computer Science*, 148(1):89–111, February 2006. ISSN 1571-0661. doi: 10.1016/j.entcs.2005.12.014. URL <https://www.sciencedirect.com/science/article/pii/S1571066106000442>.
- Earl T. Barr, Mark Harman, Phil McMinn, Muzammil Shahbaz, and Shin Yoo. The Oracle Problem in Software Testing: A Survey. *IEEE Transactions on Software Engineering*, 41(5):507–525, 2015. doi: 10.1109/TSE.2014.2372785.
- Mykola Bas. *Data Backup and Archiving*. Bachelor thesis, Czech University of Life Sciences Prague, Praha-Suchbát, Czechia, March 2024. URL https://theses.cz/id/60licg/zaverecna_prace_Archive.pdf.
- Josh Berdine, Cristiano Calcagno, and Peter W. O’Hearn. Smallfoot: Modular Automatic Assertion Checking with Separation Logic. In Frank S. de Boer,

- Marcello M. Bonsangue, Susanne Graf, and Willem-Paul de Roever, editors, *Formal Methods for Components and Objects*, pages 115–137, Berlin, Heidelberg, 2006. Springer. ISBN 978-3-540-36750-5. doi: 10.1007/11804192_6.
- Michael Bluejay. Slot Machine PAR Sheets, May 2024. URL <https://easy.vegas/games/slots/par-sheets>.
- Chris Bocchino and William Hamilton. Eastern Range Titan IV/Centaur-TDRSS Operational Compatibility Testing. In *International Telemetry Conference Proceedings*, San Diego, CA, USA, October 1996. International Foundation for Telemetry. ISBN 978-0-608-04247-3. URL https://repository.arizona.edu/bitstream/handle/10150/607608/ITC_1996_96-01-4.pdf?sequence=1&isAllowed=y.
- Pierre Bourque and Richard E. Fairley, editors. *Guide to the Software Engineering Body of Knowledge, Version 3.0*. IEEE Computer Society Press, Washington, DC, USA, 2014. ISBN 0-7695-5166-1. URL www.swebok.org.
- Jacques Carette, Spencer Smith, Jason Balaci, Ting-Yu Wu, Samuel Crawford, Dong Chen, Dan Szymczak, Brooks MacLachlan, Dan Scime, and Maryyam Niazi. Drasil, February 2021. URL <https://github.com/JacquesCarette/Drasil/tree/v0.1-alpha>.
- Patrice Chalin, Joseph R. Kiniry, Gary T. Leavens, and Erik Poll. Beyond Assertions: Advanced Specification and Verification with JML and ESC/Java2. In Frank S. de Boer, Marcello M. Bonsangue, Susanne Graf, and Willem-Paul de Roever, editors, *Formal Methods for Components and Objects*, pages 342–363, Berlin, Heidelberg, 2006. Springer. ISBN 978-3-540-36750-5. doi: 10.1007/11804192_16.
- ChatGPT (GPT-4o). Defect Clustering Testing, November 2024. URL <https://chatgpt.com/share/67463dd1-d0a8-8012-937b-4a3db0824dcf>.
- Shauvik Roy Choudhary, Husayn Versee, and Alessandro Orso. A Cross-browser Web Application Testing Tool. In *2010 IEEE International Conference on Software Maintenance*, pages 1–6, Timisoara, Romania, September 2010. IEEE. ISBN 978-1-4244-8629-8. doi: 10.1109/ICSM.2010.5609728. URL <http://ieeexplore.ieee.org/abstract/document/5609728>. ISSN: 1063-6773.
- Alan Dennis, Barbara Haley Wixom, and Roberta M. Roth. *System Analysis and Design*. John Wiley & Sons, 5th edition, 2012. ISBN 978-1-118-05762-9. URL https://www.uoitc.edu.iq/images/documents/informatics-institute/Competitive_exam/Systemanalysisanddesign.pdf.
- Monika Dhok and Murali Krishna Ramanathan. Directed Test Generation to Detect Loop Inefficiencies. In *Proceedings of the 2016 24th ACM SIGSOFT International Symposium on Foundations of Software Engineering, FSE 2016*, pages 895–907, New York, NY, USA, November 2016. Association for Computing Machinery. ISBN 978-1-4503-4218-6. doi: 10.1145/2950290.2950360. URL <https://dl.acm.org/doi/10.1145/2950290.2950360>.

- 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*, 135:106737, 2020. ISSN 0098-1354. doi: <https://doi.org/10.1016/j.compchemeng.2020.106737>. URL <https://www.sciencedirect.com/science/article/pii/S0098135419312736>.
- Serdar Doğan, Aysu Betin-Can, and Vahid Garousi. Web application testing: A systematic literature review. *Journal of Systems and Software*, 91:174–201, 2014. ISSN 0164-1212. doi: <https://doi.org/10.1016/j.jss.2014.01.010>. URL <https://www.sciencedirect.com/science/article/pii/S0164121214000223>.
- Emelie Engström and Kai 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)*, pages 1–4, 2015. doi: 10.1109/ICSTW.2015.7107470.
- Norman E. Fenton and Shari Lawrence Pfleeger. *Software Metrics: A Rigorous & Practical Approach*. PWS Publishing Company, Boston, MA, USA, 2nd edition, 1997. ISBN 0-534-95425-1.
- Donald G. Firesmith. A Taxonomy of Testing Types, 2015. URL <https://apps.dtic.mil/sti/pdfs/AD1147163.pdf>.
- 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)*, volume 1, pages 329–335, Aachen, Germany, 2004. IEEE. ISBN 0-7803-8399-0. doi: 10.1109/PESC.2004.1355765.
- Paul Gerrard. Risk-based E-business Testing - Part 1: Risks and Test Strategy. Technical report, Systeme Evolutif, London, UK, 2000a. URL https://www.agileconnection.com/sites/default/files/article/file/2013/XUS129342file1_0.pdf.
- Paul Gerrard. Risk-based E-business Testing - Part 2: Test Techniques and Tools. Technical report, Systeme Evolutif, London, UK, 2000b. URL wenku.uml.com.cn/document/test/EBTestingPart2.pdf.
- Paul Gerrard and Neil Thompson. *Risk-based E-business Testing*. Artech House computing library. Artech House, Norwood, MA, USA, 2002. ISBN 978-1-58053-570-0. URL <https://books.google.ca/books?id=54UKereAdJ4C>.
- Anup K Ghosh and Jeffrey M Voas. Inoculating software for survivability. *Communications of the ACM*, 42(7):38–44, July 1999. URL <https://dl.acm.org/doi/pdf/10.1145/306549.306563>. Publisher: ACM New York, NY, USA.
- Patrice Godefroid and Daniel Luchaup. Automatic Partial Loop Summarization in Dynamic Test Generation. In *Proceedings of the 2011 International Symposium on Software Testing and Analysis*, ISSTA '11, pages 23–33, New York, NY, USA, July 2011. Association for Computing Machinery. ISBN 978-1-4503-0562-4. doi:

10.1145/2001420.2001424. URL <https://dl.acm.org/doi/10.1145/2001420.2001424>.

W. Goralski. xDSL loop qualification and testing. *IEEE Communications Magazine*, 37(5):79–83, 1999. doi: 10.1109/35.762860.

Arie Gurfinkel. Testing: Coverage and Structural Coverage, 2017. URL <https://ece.uwaterloo.ca/~agurfink/ece653w17/assets/pdf/W03-Coverage.pdf>.

Matthias Hamburg and Gary Mogyorodi, editors. ISTQB Glossary, v4.3, 2024. URL https://glossary.istqb.org/en_US/search.

Daniel C Holley, Gary D Mele, and Sujata Naidu. NASA Rat Acoustic Tolerance Test 1994-1995: 8 kHz, 16 kHz, 32 kHz Experiments. Technical Report NASA-CR-202117, San Jose State University, San Jose, CA, USA, January 1996. URL <https://ntrs.nasa.gov/api/citations/19960047530/downloads/19960047530.pdf>.

R. Brian Howe and Robert Johnson. Research Protocol for the Evaluation of Medical Waiver Requirements for the Use of Lisinopril in USAF Aircrew. Interim Technical Report AL/AO-TR-1995-0116, Air Force Materiel Command, Brooks Air Force Base, TX, USA, November 1995. URL <https://apps.dtic.mil/sti/tr/pdf/ADA303379.pdf>.

Anthony Hunt, Peter Michalski, Dong Chen, Jason Balaci, and Spencer Smith. Drasil - Generate All the Things!, 2021. URL <https://jacquescarette.github.io/Drasil/>.

IEEE. IEEE Standard for System and Software Verification and Validation. *IEEE Std 1012-2012 (Revision of IEEE Std 1012-2004)*, 2012. doi: 10.1109/IEEESTD.2012.6204026.

Adisak Intana, Monchanok Thongthep, Phatcharee Thepnimit, Phaplak Saethapan, and Tanawat Monpipat. SYNTest: Prototype of Syntax Test Case Generation Tool. In *5th International Conference on Information Technology (InCIT)*, pages 259–264. IEEE, 2020. ISBN 978-1-72819-321-2. doi: 10.1109/InCIT50588.2020.9310968.

ISO. ISO 13849-1:2015 - Safety of machinery –Safety-related parts of control systems –Part 1: General principles for design. *ISO 13849-1:2015*, December 2015. URL <https://www.iso.org/obp/ui#iso:std:iso:13849:-1:ed-3:v1:en>.

ISO. ISO 21384-2:2021 - Unmanned aircraft systems –Part 2: UAS components. *ISO 21384-2:2021*, December 2021. URL <https://www.iso.org/obp/ui#iso:std:iso:21384:-2:ed-1:v1:en>.

ISO. ISO 28881:2022 - Machine tools –Safety –Electrical discharge machines. *ISO 28881:2022*, April 2022. URL <https://www.iso.org/obp/ui#iso:std:iso:28881:ed-2:v1:en>.

- ISO/IEC. 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*, March 2011. URL <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-1:v1:en>.
- ISO/IEC. ISO/IEC 2382:2015 - Information technology –Vocabulary. *ISO/IEC 2382:2015*, May 2015. URL <https://www.iso.org/obp/ui/en/#iso:std:iso-iec:2382:ed-1:v2:en>.
- ISO/IEC. ISO/IEC TS 20540:2018 - Information technology – Security techniques –Testing cryptographic modules in their operational environment. *ISO/IEC TS 20540:2018*, May 2018. URL <https://www.iso.org/obp/ui#iso:std:iso-iec:ts:20540:ed-1:v1:en>.
- 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*, November 2023a. URL <https://www.iso.org/obp/ui/#iso:std:iso-iec:25010:ed-2:v1:en>.
- ISO/IEC. ISO/IEC 25019:2023 - Systems and software engineering –Systems and software Quality Requirements and Evaluation (SQuaRE) –Quality-in-use model. *ISO/IEC 25019:2023*, November 2023b. URL <https://www.iso.org/obp/ui/en/#iso:std:iso-iec:25019:ed-1:v1:en>.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Systems and software engineering–Vocabulary. *ISO/IEC/IEEE 24765:2010(E)*, December 2010. doi: 10.1109/IEEESTD.2010.5733835.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Systems and software engineering –Software testing –Part 1: General concepts. *ISO/IEC/IEEE 29119-1:2013*, September 2013. doi: 10.1109/IEEESTD.2013.6588537.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Systems and software engineering –System life cycle processes. *ISO/IEC/IEEE 15288:2015*, May 2015. doi: 10.1109/IEEESTD.2015.7106435.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Software and systems engineering –Software testing –Part 5: Keyword-Driven Testing. *ISO/IEC/IEEE 29119-5:2016*, November 2016. doi: 10.1109/IEEESTD.2016.7750539.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Systems and software engineering–Vocabulary. *ISO/IEC/IEEE 24765:2017(E)*, September 2017. doi: 10.1109/IEEESTD.2017.8016712.
- 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*, March 2019a. doi: 10.1109/IEEESTD.2019.8657410.

- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Systems and software engineering –Guidelines for the utilization of ISO/IEC/IEEE 15288 in the context of system of systems (SoS). *ISO/IEC/IEEE 21840:2019*, December 2019b. doi: 10.1109/IEEESTD.2019.8929110.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Software and systems engineering –Software testing –Part 2: Test processes. *ISO/IEC/IEEE 29119-2:2021(E)*, October 2021a. doi: 10.1109/IEEESTD.2021.9591508.
- ISO/IEC and IEEE. ISO/IEC/IEEE International Standard - Software and systems engineering –Software testing –Part 3: Test documentation. *ISO/IEC/IEEE 29119-3:2021(E)*, October 2021b. doi: 10.1109/IEEESTD.2021.9591577.
- 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)*, October 2021c. doi: 10.1109/IEEESTD.2021.9591574.
- 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)*, January 2022. doi: 10.1109/IEEESTD.2022.9698145.
- Claude Jard, Thierry Jérón, Lénéaïck Tanguy, and César Viho. Remote testing can be as powerful as local testing. In Jianping Wu, Samuel T. Chanson, and Qiang Gao, editors, *Formal Methods for Protocol Engineering and Distributed Systems: Forte XII / PSTV XIX'99*, volume 28 of *IFIP Advances in Information and Communication Technology*, pages 25–40, Beijing, China, October 1999. Springer. ISBN 978-0-387-35578-8. doi: 10.1007/978-0-387-35578-8_2. URL https://doi.org/10.1007/978-0-387-35578-8_2.
- Timothy P. Johnson. Snowball Sampling: Introduction. In *Wiley StatsRef: Statistics Reference Online*. John Wiley & Sons, Ltd, 2014. ISBN 978-1-118-44511-2. doi: <https://doi.org/10.1002/9781118445112.stat05720>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat05720>. eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/9781118445112.stat05720>.
- Ben Kam. Web Applications Testing. Technical Report 2008-550, Queen’s University, Kingston, ON, Canada, October 2008. URL <https://research.cs.queensu.ca/TechReports/Reports/2008-550.pdf>.
- Cem Kaner, James Bach, and Bret Pettichord. *Lessons Learned in Software Testing: A Context-Driven Approach*. John Wiley & Sons, December 2011. ISBN 978-0-471-08112-8. URL <https://www.wiley.com/en-ca/Lessons+Learned+in+Software+Testing%3A+A+Context-Driven+Approach-p-9780471081128>.
- Upulee Kanewala and Tsong Yueh Chen. Metamorphic testing: A simple yet effective approach for testing scientific software. *Computing in Science & Engineering*, 21(1):66–72, 2019. doi: 10.1109/MCSE.2018.2875368.

- Knüvener Mackert GmbH. *Knüvener Mackert SPICE Guide*. Knüvener Mackert GmbH, Reutlingen, Germany, 7th edition, 2022. ISBN 978-3-00-061926-7. URL <https://knuevenermackert.com/wp-content/uploads/2021/06/SPICE-BOOKLET-2022-05.pdf>.
- Ivans Kuļšovs, Vineta Arnicanē, Guntis Arnicans, and Juris Borzovs. Inventory of Testing Ideas and Structuring of Testing Terms. 1:210–227, January 2013.
- Shuvendu K. Lahiri, Kenneth L. McMillan, Rahul Sharma, and Chris Hawblitzel. Differential Assertion Checking. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, ESEC/FSE 2013, pages 345–355, New York, NY, USA, August 2013. Association for Computing Machinery. ISBN 978-1-4503-2237-9. doi: 10.1145/2491411.2491452. URL <https://dl.acm.org/doi/10.1145/2491411.2491452>.
- LambdaTest. What is Operational Testing: Quick Guide With Examples, 2024. URL <https://www.lambdatest.com/learning-hub/operational-testing>.
- Danye Liu, Shaonan Tian, Yu Zhang, Chaoquan Hu, Hui Liu, Dong Chen, Lin Xu, and Jun Yang. Ultrafine SnPd nanoalloys promise high-efficiency electrocatalysis for ethanol oxidation and oxygen reduction. *ACS Applied Energy Materials*, 6(3):1459–1466, January 2023. doi: <https://doi.org/10.1021/acsaem.2c03355>. URL https://pubs.acs.org/doi/pdf/10.1021/acsaem.2c03355?casa_token=ItHfKXeQNbsAAAAA:8zEdU5hi2HfHsSony3ku-lbH902jkHpA-JZw8jleODzUvFtSdQRdbYhmVq47aX22igR52o2S22mnC88Mxw. Publisher: ACS Publications.
- Robert Mandl. Orthogonal Latin squares: an application of experiment design to compiler testing. *Communications of the ACM*, 28(10):1054–1058, October 1985. ISSN 0001-0782. doi: 10.1145/4372.4375. URL <https://doi.org/10.1145/4372.4375>.
- Robert M. McClure. Introduction, July 2001. URL <http://homepages.cs.ncl.ac.uk/brian.randell/NATO/Introduction.html>.
- Mahshid Helali 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*, ESEC/FSE 2019, pages 1187–1189, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 978-1-4503-5572-8. doi: 10.1145/3338906.3342484. URL <https://doi.org/10.1145/3338906.3342484>.
- 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. Technical report, NASA, Houston, TX, USA, 1999. URL <https://ntrs.nasa.gov/api/citations/2000085877/downloads/2000085877.pdf>.
- E. E. Mukhin, V. M. Nelyubov, V. A. Yukish, E. P. Smirnova, V. A. Solovei, N. K. Kalinina, V. G. Nagaitsev, M. F. Valishin, A. R. Belozeroval, S. A. Enin,

- A. A. Borisov, N. A. Deryabina, V. I. Khripunov, D. V. Portnov, N. A. Babinov, D. V. Dokhtarenko, I. A. Khodunov, V. N. Klimov, A. G. Razdobarin, S. E. Alexandrov, D. I. Elets, A. N. Bazhenov, I. M. Bukreev, An P. Chernakov, A. M. Dmitriev, Y. G. Ibragimova, A. N. Koval, G. S. Kurskiev, A. E. Litvinov, K. O. Nikolaenko, D. S. Samsonov, V. A. Senichenkov, R. S. Smirnov, S. Yu Tolstyakov, I. B. Tereschenko, L. A. Varshavchik, N. S. Zhiltsov, A. N. Mokeev, P. V. Chernakov, P. Andrew, and M. Kempenaars. Radiation tolerance testing of piezoelectric motors for ITER (first results). *Fusion Engineering and Design*, 176(article 113017), 2022. ISSN 0920-3796. doi: <https://doi.org/10.1016/j.fuseengdes.2022.113017>. URL <https://www.sciencedirect.com/science/article/pii/S0920379622000175>.
- Peter Naur and Brian, editors Randell. Software Engineering: Report on a conference sponsored by the NATO Science Committee, Garmisch, Germany, 7th to 11th October 1968. Brussels, Belgium, January 1969. Scientific Affairs Division, NATO. URL <http://homepages.cs.ncl.ac.uk/brian.randell/NATO/nato1968.PDF>.
- Randal C. Nelson. Formal Computational Models and Computability, January 1999. URL https://www.cs.rochester.edu/u/nelson/courses/csc_173/computability/undecidable.html.
- Pranav Pandey. Scalability vs Elasticity, February 2023. URL <https://www.linkedin.com/pulse/scalability-vs-elasticity-pranav-pandey/>.
- Bhupesh A. Parate, K.D. Deodhar, and V.K. Dixit. Qualification Testing, Evaluation and Test Methods of Gas Generator for IEDs Applications. *Defence Science Journal*, 71(4):462–469, July 2021. doi: 10.14429/dsj.71.16601. URL <https://publications.drdo.gov.in/ojs/index.php/dsj/article/view/16601>.
- Ron Patton. *Software Testing*. Sams Publishing, Indianapolis, IN, USA, 2nd edition, 2006. ISBN 0-672-32798-8.
- William E. Perry. *Effective Methods for Software Testing*. Wiley Publishing, Inc., Indianapolis, IN, USA, 3rd edition, 2006. ISBN 978-0-7645-9837-1.
- J.F. Peters and W. Pedrycz. *Software Engineering: An Engineering Approach*. Worldwide series in computer science. John Wiley & Sons, Ltd., 2000. ISBN 978-0-471-18964-0.
- Brian J. Pierre, Felipe Wilches-Bernal, David A. Schoenwald, Ryan T. Elliott, Jason C. Neely, Raymond H. Byrne, and Daniel J. Trudnowski. Open-loop testing results for the pacific DC intertie wide area damping controller. In *2017 IEEE Manchester PowerTech*, pages 1–6, 2017. doi: 10.1109/PTC.2017.7980834.
- Sebastian Preuß, Hans-Christian Lapp, and Hans-Michael Hanisch. Closed-loop System Modeling, Validation, and Verification. In *Proceedings of 2012 IEEE 17th International Conference on Emerging Technologies & Factory Automation*

(*ETFA 2012*), pages 1–8, Krakow, Poland, 2012. IEEE. ISBN 978-1-4673-4736-5. doi: 10.1109/ETFA.2012.6489679. URL <https://ieeexplore.ieee.org/abstract/document/6489679>.

Project Management Institute. *A Guide to the Project Management Body of Knowledge: PMBOK(R) Guide*. Project Management Institute, 5th edition, 2013. ISBN 1-935589-67-9.

Stuart C. Reid. Popular Misconceptions in Module Testing. In *Proceeding of the 13 International Conference on Testing Computer Software*, Washington, DC, USA, 1996.

Per Runeson and Martin Höst. Guidelines for conducting and reporting case study research in software engineering. *Empirical Software Engineering*, 14:131–164, April 2009. doi: <https://doi.org/10.1007/s10664-008-9102-8>. URL <https://link.springer.com/article/10.1007/s10664-008-9102-8>.

Vasile Rus, Sameer Mohammed, and Sajjan G Shiva. Automatic Clustering of Defect Reports. In *Proceedings of the Twentieth International Conference on Software Engineering & Knowledge Engineering (SEKE 2008)*, pages 291–296, San Francisco, CA, USA, July 2008. Knowledge Systems Institute Graduate School. ISBN 1-891706-22-5. URL <https://core.ac.uk/download/pdf/48606872.pdf>.

Kazunori Sakamoto, Kaizu Tomohiro, Daigo Hamura, Hironori Washizaki, and Yoshiaki Fukazawa. POGen: A Test Code Generator Based on Template Variable Coverage in Gray-Box Integration Testing for Web Applications. In Vittorio Cortellessa and Dániel Varró, editors, *Fundamental Approaches to Software Engineering*, pages 343–358, Berlin, Heidelberg, March 2013. Springer Berlin Heidelberg. ISBN 978-3-642-37057-1. URL https://link.springer.com/chapter/10.1007/978-3-642-37057-1_25.

Raghvinder S. Sangwan and Phillip A. LaPlante. Test-Driven Development in Large Projects. *IT Professional*, 8(5):25–29, October 2006. ISSN 1941-045X. doi: 10.1109/MITP.2006.122. URL <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=1717338>.

Huib Schoots. The ISO29119 debate, September 2014. URL <https://www.huibshoorts.nl/wordpress/?p=1800>.

Sheetal Sharma, Kartika Panwar, and Rakesh Garg. Decision Making Approach for Ranking of Software Testing Techniques Using Euclidean Distance Based Approach. *International Journal of Advanced Research in Engineering and Technology*, 12(2):599–608, February 2021. ISSN 0976-6499. doi: 10.34218/IJARET.12.2.2021.059. URL <https://iaeme.com/Home/issue/IJARET?Volume=12&Issue=2>.

Spencer Smith. Potential Projects, June 2024. URL <https://github.com/JacquesCarette/Drasil/wiki/Potential-Projects>.

Spencer Smith and Jacques Carette. Private Communication, July 2023.

Harry Sneed and Siegfried Göschl. A Case Study of Testing a Distributed Internet-System. *Software Focus*, 1:15–22, September 2000. doi: 10.1002/1529-7950(20009)1:13.3.CO;2-#. URL https://www.researchgate.net/publication/220116945_Testing_software_for_Internet_application.

Pradeep Soundararajan. An open letter to the President of the International Organization for Standardization about ISO 29119, December 2015. URL <https://moolya.com/blog/testing-stories/an-open-letter-to-the-president-of-the-international-organization-for-standardization-about-iso-29119/>.

Erica Souza, Ricardo Falbo, and Nandamudi Vijaykumar. ROoST: Reference Ontology on Software Testing. *Applied Ontology*, 12:1–32, March 2017. doi: 10.3233/AO-170177.

Ephraim Suhir, Laurent Bechou, Alain Bensoussan, and Johann Nicolics. Photovoltaic reliability engineering: quantification testing and probabilistic-design-reliability concept. In *Reliability of Photovoltaic Cells, Modules, Components, and Systems VI*, volume 8825, pages 125–138. SPIE, September 2013. doi: 10.1117/12.2030377. URL <https://www.spiedigitallibrary.org/conference-proceedings-of-spie/8825/88250K/Photovoltaic-reliability-engineering--quantification-testing-and-probabilistic-design-reliability/10.1117/12.2030377.full>.

Guido Tebes, Denis Peppino, Pablo Becker, Gerardo Matturro, Martín Solari, and Luis Olsina. A Systematic Review on Software Testing Ontologies. pages 144–160. August 2019. ISBN 978-3-030-29237-9. doi: 10.1007/978-3-030-29238-6_11.

Guido Tebes, Luis Olsina, Denis Peppino, and Pablo Becker. TestTDO: A Top-Domain Software Testing Ontology. pages 364–377, Curitiba, Brazil, May 2020. ISBN 978-1-71381-853-3.

Daniel Trudnowski, Brian Pierre, Felipe Wilches-Bernal, David Schoenwald, Ryan Elliott, Jason Neely, Raymond Byrne, and Dmitry Kosterev. Initial closed-loop testing results for the pacific DC intertie wide area damping controller. In *2017 IEEE Power & Energy Society General Meeting*, pages 1–5, 2017. doi: 10.1109/PESGM.2017.8274724.

Kwok-Leung Tsui. An Overview of Taguchi Method and Newly Developed Statistical Methods for Robust Design. *IIE Transactions*, 24(5):44–57, May 2007. doi: 10.1080/07408179208964244. URL <https://doi.org/10.1080/07408179208964244>. Publisher: Taylor & Francis.

Matheus A. Tunes, Sean M. Drewry, Jose D. Arregui-Mena, Sezer Picak, Graeme Greaves, Luigi B. Cattini, Stefan Pogatscher, James A. Valdez, Saryu Fensin, Osman El-Atwani, Stephen E. Donnelly, Tarik A. Saleh, and Philip D. Edmondson. Accelerated radiation tolerance testing of Ti-based MAX phases. *Materials Today Energy*, 30(article 101186), October 2022. ISSN 2468-6069. doi:

<https://doi.org/10.1016/j.mtener.2022.101186>. URL <https://www.sciencedirect.com/science/article/pii/S2468606922002441>.

Michael Unterkalmsteiner, Robert Feldt, and Tony Gorschek. A Taxonomy for Requirements Engineering and Software Test Alignment. *ACM Transactions on Software Engineering and Methodology*, 23(2):1–38, March 2014. ISSN 1049-331X, 1557-7392. doi: 10.1145/2523088. URL <http://arxiv.org/abs/2307.12477>. arXiv:2307.12477 [cs].

Petya Valcheva. Orthogonal Arrays and Software Testing. In Dimitar G. Velev, editor, *3rd International Conference on Application of Information and Communication Technology and Statistics in Economy and Education*, volume 200, pages 467–473, Sofia, Bulgaria, December 2013. University of National and World Economy. ISBN 978-954-644-586-5. URL <https://icaictsee-2013.unwe.bg/proceedings/ICAICTSEE-2013.pdf>.

Hans van Vliet. *Software Engineering: Principles and Practice*. John Wiley & Sons, Ltd., Chichester, England, 2nd edition, 2000. ISBN 0-471-97508-7.

Hironori Washizaki, editor. *Guide to the Software Engineering Body of Knowledge, Version 4.0*. January 2024. URL <https://waseda.app.box.com/v/SWEBOK4-book>.

Hironori Washizaki. Software Engineering Body of Knowledge (SWEBOK), February 2025. URL <https://www.computer.org/education/bodies-of-knowledge/software-engineering/>.

Wikibooks Contributors. *Haskell/Variables and functions*. Wikimedia Foundation, October 2023. URL https://en.wikibooks.org/wiki/Haskell/Variables_and_functions.

Han 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*, 26(3):1573–1580, August 2011. ISSN 0885-8950. doi: 10.1109/TPWRS.2010.2082576. URL <https://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=5620950>.

Kaiqiang Zhang, Chris Hutson, James Knighton, Guido Herrmann, and Tom Scott. Radiation Tolerance Testing Methodology of Robotic Manipulator Prior to Nuclear Waste Handling. *Frontiers in Robotics and AI*, 7(article 6), February 2020. ISSN 2296-9144. doi: 10.3389/frobt.2020.00006. URL <https://www.frontiersin.org/articles/10.3389/frobt.2020.00006>.

Changlin Zhou, Qun Yu, and Litao Wang. Investigation of the Risk of Electromagnetic Security on Computer Systems. *International Journal of Computer and Electrical Engineering*, 4(1):92, February 2012. URL <http://ijcee.org/papers/457-JE504.pdf>. Publisher: IACSIT Press.

Appendix A

Detailed Scope Analysis

As outlined in [Chapter 1](#), the scope of our research is limited to testing applied to code itself. Throughout our research, we identify many approaches that are out of scope based on this criteria.

A.1 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 we exclude from our research:

- Ergonomics testing and proximity-based testing ([Hamburg and Mogyorodi, 2024](#)) are out of scope, since they are used for testing hardware.
- EManations SECurity (EMSEC) testing ([ISO, 2021](#); [Zhou et al., 2012](#), p. 95), which deals with the “security risk” of “information leakage via electromagnetic emanation” ([Zhou et al., 2012](#), p. 95), is also out of scope.
- All the examples of domain-specific testing given by [Firesmith \(2015, p. 26\)](#) are focused on hardware, so these examples are out of scope. However, this might not be representative of *all* kinds of domain-specific testing (e.g., Machine Learning (ML) model testing seems domain-specific), so some subset of this approach may be in scope.
- Similarly, the examples of environmental tolerance testing given by [Firesmith \(2015, p. 56\)](#) do not seem to apply to software. For example, radiation tolerance testing seems to focus on hardware, such as motors ([Mukhin et al., 2022](#)), robots ([Zhang et al., 2020](#)), or “nanolayered carbide and nitride materials” ([Tunes et al., 2022](#), p. 1). Acceleration tolerance testing seems to focus on astronauts ([Morgun et al., 1999](#), p. 11), aviators ([Howe and Johnson, 1995](#), pp. 27, 42), or catalysts ([Liu et al., 2023](#), p. 1463) and acoustic tolerance testing on rats ([Holley et al., 1996](#)), which are even less related! Since these all focus on environment-specific factors that would not impact the code, these examples are out of scope. As with domain-specific testing,

a subset of environmental tolerance testing may be in scope, but since no candidates have been found, this approach is out of scope for now.

- [Knüvener Mackert GmbH \(2022\)](#) uses the terms “software qualification testing” and “system qualification testing” in the context of the automotive industry. While these may be in scope, the more general idea of “qualification testing” seems to refer to the process of making a hardware component, such as an electronic component ([Ahsan et al., 2020](#)), gas generator ([Parate et al., 2021](#)) or photovoltaic device, “into a reliable and marketable product” ([Suhir et al., 2013](#), p. 1). Therefore, it is currently unclear if this is in scope.
- Orthogonal Array Testing (OAT) can be used when testing software ([Mandl, 1985](#)) (in scope) but can also be used when testing hardware ([Valcheva, 2013](#), 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” ([Yu et al., 2011](#), p. 1573) or “product and manufacturing process design” ([Tsui, 2007](#), p. 44) and is thus also out of scope.
- Since control systems often have a software *and* hardware component ([ISO, 2015](#); [Preuße et al., 2012](#); [Forsyth et al., 2004](#)), only the software component is in scope. In some cases, it is unclear whether the “loops”¹ being tested are implemented by software or hardware, such as those in wide-area damping controllers ([Pierre et al., 2017](#); [Trudnowski et al., 2017](#)).
 - A related note: “path coverage” or “path testing” seems to be able to refer to either paths through code (as a subset of control-flow testing) ([Washizaki, 2024](#), p. 5-13) or through a model, such as a finite-state machine (as a subset of model-based testing) ([Doğan et al., 2014](#), p. 184).
- Physical testing (inferred from physical requirements ([ISO/IEC and IEEE, 2017](#), p. 322) and requirements-based testing; see [Section 3.2](#)) tests “a physical characteristic that a system or system component must possess” (p. 322).

A.2 V&V of Other Artifacts

While many artifacts produced by the software life cycle can be tested, only testing performed on the code *itself* is in scope. Therefore, we exclude the following test approaches either in full or in part:

- Design reviews and documentation reviews are out of scope, as they focus on the V&V of design ([ISO/IEC and IEEE, 2017](#), pp. 132) and documentation (p. 144), respectively.

¹Humorously, the testing of loops in chemical systems ([Dominguez-Pumar et al., 2020](#)) and copper loops ([Goralski, 1999](#)) are out of scope.

- Security audits can focus on “an organization’s ... processes and infrastructure” (Hamburg and Mogyorodi, 2024) (out of scope) or “aim to ensure that all of the products installed on a site are secure when checked against the known vulnerabilities for those products” (Gerrard, 2000b, p. 28) (in scope).
- Error seeding is the “process of intentionally adding known faults² 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” (ISO/IEC and IEEE, 2017, p. 165), which helps verify the actual code (in scope).
- Fault injection testing, where “faults are artificially introduced² into the SUT [System Under Test]”, can be used to evaluate the effectiveness of a test suite (Washizaki, 2024, 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” (ISO/IEC and IEEE, 2022, 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 a slightly modified version of the SUT” (Washizaki, 2024, 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 (Washizaki, 2024, p. 5-15) (in scope).
- Nontechnical testing (inferred from nontechnical requirements (ISO/IEC and IEEE, 2017, p. 293) and requirements-based testing; see Section 3.2) tests “product and service acquisition or development that is not a property of the product or service” (ISO/IEC and IEEE, 2017, p. 293).

A.3 Static Testing

Throughout the literature, static testing is more ambiguous than dynamic testing, with more ad hoc processes and inconsistent in/exclusion from the scope of software testing in general (see Contradiction 5). Furthermore, it seems less relevant to our original goal (the automatic generation of tests). In particular, many techniques require human intervention, either by design (such as code inspections) or to identify and resolve false positives (such as intentional exceptions to linting rules). Nevertheless, understanding the breadth of testing approaches requires a “complete” picture of how software can be tested and how the various approaches relate to one another. Parts of these static approaches may even be generated in

²While 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” (ISO/IEC and IEEE, 2017, p. 165) and “test the robustness of the system” (2022, p. 42), respectively. Therefore, these approaches are not considered synonyms, and the lack of this relation in the literature is not included in Section 5.2.2 as a synonym flaw.

the future! For these reasons, we keep static testing in scope for this stage of our work, even though static testing will likely be removed at a later stage of analysis based on our original motivation.

See #41 and #44

A.4 Vague Terminology

Some terms are so vague that they do not provide any new, meaningful information. For example, the “systematic determination of the extent to which an entity meets its specified criteria” (ISO/IEC and IEEE, 2017, p. 167) is certainly relevant to testing software; while this definition of “evaluation” may be meaningful when defining software testing generally, it does not define a new approach or procedure and applies much more broadly than just to testing. We decided that the following terms are too vague to merit tracking in our glossaries or analyzing further:

See #39, #44, and #28

- **Evaluation:** the “systematic determination of the extent to which an entity meets its specified criteria” (ISO/IEC and IEEE, 2017, p. 167)
- **Product Analysis:** the “process of evaluating a product by manual or automated means to determine if the product has certain characteristics” (ISO/IEC and IEEE, 2017, p. 343)
- **Quality Audit:** “a structured, independent process to determine if project activities comply with organizational and project policies, processes, and procedures” (Project Management Institute, 2013, p. 247)
- **Software Product Evaluation:** a “technical operation that consists of producing an assessment of one or more characteristics of a software product according to a specified procedure” (ISO/IEC and IEEE, 2017, p. 424)

A.5 Language-specific Approaches

Specific programming languages are sometimes used to define test approaches. If the reliance on a specific programming language is intentional, then this really implies an underlying test approach that may be generalized to other languages.

These are therefore considered out-of-scope, including the following examples:

- “An approach ... for JavaScript testing (referred to as Randomized)” (Doğan et al., 2014, p. 192) is really just random testing used within JavaScript.
- “SQL statement coverage” (Doğan et al., 2014, Tab. 13) is really just statement coverage used specifically for SQL statements.
- Testing for “faults specific to PHP” is just a subcategory of fault-based testing, since “execution failures ... caused by missing an included file, wrong MySQL quer[ies] and uncaught exceptions” (Doğan et al., 2014, Tab. 27) are not exclusive to PHP.

OG Alalfi et al., 2010

OG Artzi et al., 2008

- While “HTML testing” is listed or implied by [Gerrard \(2000a, Tab. 2; 2000b, Tab. 1, p. 3\)](#) and [Patton \(2006, p. 220\)](#), it seems to be a combination of syntax testing, functionality testing, hyperlink testing/link checking, cross-browser compatibility testing, performance testing, content checking ([Gerrard, 2000b, p. 3](#)), and grey-box testing ([Patton, 2006, pp. 218–220](#)).

A.6 Orthogonally Derived Approaches

Some test approaches appear to be combinations of other (seemingly orthogonal) approaches. While the use of a combination term can sometimes make sense, such as when writing a paper or performing testing that focuses on the intersection between two test approaches, they are sometimes given the same “weight” as their atomic counterparts. For example, [Hamburg and Mogyorodi \(2024\)](#) include “formal reviews” and “informal reviews” in their glossary as separate terms, despite their definitions essentially boiling down to “reviews that follow (or do not follow) a formal process”, which do not provide any new information. These approaches are simply the combinations of “reviews” with “formal” and “informal testing”, respectively. If a source describes an orthogonally derived approach in more detail, such as security audits, we record it as a distinct approach in [our test approach glossary](#) with its related information. Otherwise, we consider it out of scope since its details are captured by its in-scope subapproaches. The following are examples of these orthogonally derived approaches, most of which are out of scope:

1. Black box conformance testing ([Jard et al., 1999, p. 25](#))
2. Black-box integration testing ([Sakamoto et al., 2013, pp. 345–346](#))
3. Checklist-based reviews ([Hamburg and Mogyorodi, 2024](#))
4. Closed-loop HiL verification ([Preuße et al., 2012, p. 6](#))
5. Closed-loop protection system testing ([Forsyth et al., 2004, p. 331](#))
6. Endurance stability testing ([Firesmith, 2015, p. 55](#))
7. End-to-end functionality testing ([ISO/IEC and IEEE, 2021c, p. 20](#); [Gerrard, 2000a, Tab. 2](#))
8. Formal reviews ([Hamburg and Mogyorodi, 2024](#); [Washizaki, 2024, p. 12–14](#))
9. Grey-box integration testing ([Sakamoto et al., 2013, p. 344](#))
10. Incremental integration testing ([Sharma et al., 2021, pp. 601, 603, 605–606](#))
11. Informal reviews ([Hamburg and Mogyorodi, 2024](#); [Washizaki, 2024, p. 12–14](#))
12. Infrastructure compatibility testing ([Firesmith, 2015, p. 53](#))
13. Invariant-based automatic testing ([Doğan et al., 2014, pp. 184–185, Tab. 21](#)), including for “AJAX user interfaces” (p. 191)

Q #10: Is this orthogonal? Investigate legacy testing	14. Legacy system integration (testing) (Gerrard, 2000a, Tab. 2) 15. Manual procedure testing (Firesmith, 2015, p. 47) 16. Manual security audits (Gerrard, 2000b, p. 28) 17. Model-based GUI testing (Doğan et al., 2014, Tab. 1; implied by Sakamoto et al., 2013, p. 356) 18. Model-based web application testing (implied by Sakamoto et al., 2013, p. 356) 19. Non-functional search-based testing (Doğan et al., 2014, Tab. 1) 20. Offline MBT (Hamburg and Mogyorodi, 2024) 21. Online MBT (Hamburg and Mogyorodi, 2024)
Q #11: Is this orthogonal? Investigate role-based testing	22. Role-based reviews (Hamburg and Mogyorodi, 2024) 23. Scenario walkthroughs (Gerrard, 2000a, Fig. 4) 24. Scenario-based reviews (Hamburg and Mogyorodi, 2024)
Q #12: Is this orthogonal? Investigate scenario-based testing	25. Security attacks (Hamburg and Mogyorodi, 2024) 26. Security audits (ISO/IEC and IEEE, 2021c, p. 40; Gerrard, 2000b, p. 28) 27. Statistical web testing (Doğan et al., 2014, p. 185) 28. Usability test script(ing) (Hamburg and Mogyorodi, 2024)
Q #13: Is this orthogonal? Investigate scenario-based testing	29. Web application regression testing (Doğan et al., 2014, Tab. 21) 30. White-box unit testing (Sakamoto et al., 2013, pp. 345–346)

There are some cases where the subapproaches of the “compound” approaches listed previously are *not* from separate categories. However, these cases can be explained by insufficient data or by edge cases that require special care. While we assume that the categories given in Table 2.1 are orthogonal, further analysis may disprove this. For now, all of these special cases are affected by at least one of the following conditions:

1. **At least one subapproach is categorized inconsistently.** When a subapproach has more than one category (see Section 5.2.1), it is unclear which one should be used to assess orthogonality.
2. **At least one subapproach’s category is inferred.** When the category of a test approach is not given by the literature but is inferred from related context (see Section 2.3), it is unclear if it can be used to assess orthogonality.

3. **At least one subapproach is only categorized as an approach.** Since “approach” is a catch-all categorization, it does not need to be orthogonal to its subcategories.

4. **A subapproach is explicitly based on another in the same category.**

An example of this is stability testing, which tests a “property that an object has with respect to a given failure mode if it cannot exhibit that failure mode” (ISO/IEC and IEEE, 2017, p. 434). This notion of “property” is similar to that of “quality” that the test type category is built on, so it is acceptable that is implied to be a test type by its quality (ISO/IEC and IEEE, 2017, p. 434) and by Firesmith (2015, p. 55).

OG ISO/IEC,
2009

OG ISO/IEC,
2009

Appendix B

Sources by Tier

The following lists of sources comprise the corresponding source tier as defined in [Section 2.5](#).

B.1 Established Standards

ISO/IEC and IEEE, 2022; 2021a;b;c; 2019a;b; 2017; 2016; 2015; 2013; 2010; IEEE, 2012; ISO/IEC, 2023a;b; 2018; 2015; 2011; ISO, 2022; 2015

B.2 Terminology Collections

Hamburg and Mogyorodi, 2024; Doğan et al., 2014; Firesmith, 2015; Project Management Institute, 2013; Washizaki, 2024; Bourque and Fairley, 2014

B.3 Textbooks

Ammann and Offutt, 2017; Dennis et al., 2012; Gerrard and Thompson, 2002; Kaner et al., 2011; Patton, 2006; Perry, 2006; Peters and Pedrycz, 2000; van Vliet, 2000

B.4 Papers and Other Documents

Bajammal and Mesbah, 2018; Barbosa et al., 2006; Baresi and Pezzè, 2006; Barr et al., 2015; Bas, 2024; Berdine et al., 2006; Chalin et al., 2006; ChatGPT (GPT-4o), 2024; Choudhary et al., 2010; Dhok and Ramanathan, 2016; Engström and Petersen, 2015; Forsyth et al., 2004; Gerrard, 2000a;b; Ghosh and Voas, 1999; Godefroid and Luchaup, 2011; Intana et al., 2020; Jard et al., 1999; Kam, 2008; Kanewala and Yueh Chen, 2019; Kuľšovs et al., 2013; Lahiri et al., 2013; LambdaTest, 2024; Mandl, 1985; Moghadam, 2019; Pandey, 2023; Preuße et al., 2012; Reid, 1996; Rus et al., 2008; Knüvener Mackert GmbH, 2022; Sakamoto et al., 2013; Sangwan and LaPlante, 2006; Sharma et al., 2021; Sneed and Göschl, 2000; Souza et al., 2017; Tsui, 2007; Valcheva, 2013; Yu et al., 2011

Appendix C

Full Lists of Flaws

The following are the full lists of manually identified flaws, first grouped by their manifestation ([Appendix C.1](#)), then by their domain ([Appendix C.2](#)) as defined in [Section 2.2](#). We then present inferred flaws ([Appendix C.3](#)) as described in [Section 2.3](#) for completeness, although these flaws do not contribute to any counts.

C.1 Full Lists of Flaws by Manifestation

We sort the following groups of flaws by their source tier (defined in [Section 2.5](#)) in descending order of credibility (defined in [Section 2.4](#)).

C.1.1 Full List of Mistakes

1. [ISO/IEC and IEEE \(2022, p. 5\)](#) give fuzz testing the tag “artificial intelligence”; while fuzz testing could certainly be implemented in this way, it does not seem to be a requirement.
2. Since errors are distinct from defects/faults ([ISO/IEC and IEEE, 2010, pp. 128, 140](#); [Washizaki, 2024, p. 12-3](#); [van Vliet, 2000, pp. 399–400](#)), error guessing should instead be called “defect guessing” if it is based on a “checklist of potential defects” ([ISO/IEC and IEEE, 2021c, p. 29](#)) or “fault guessing” if it is a “fault-based technique” ([Bourque and Fairley, 2014, p. 4-9](#)) that “anticipate[s] the most plausible faults in each SUT” ([Washizaki, 2024, 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 [van Vliet, 2000, p. 399](#)).
3. Similarly, “fault seeding” is not a synonym of “error seeding” as claimed by [ISO/IEC and IEEE \(2017, p. 165\)](#) and [van Vliet \(2000, p. 427\)](#). The term “error seeding”, also used by [Firesmith \(2015, p. 34\)](#), 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” ([ISO/IEC and IEEE, 2017, p. 165](#)) based

on the ratio between the number of new faults and the number of introduced faults that were discovered (van Vliet, 2000, p. 427).

4. “Functionality” is defined as the “capabilities of the various ... features provided by a product” (ISO/IEC and IEEE, 2017, p. 196). Hamburg and Mogyorodi (2024) say that it is a synonym of “functional suitability”, which refers to the “capability of a product to provide [specified] functions” (ISO/IEC, 2023a; similar in ISO/IEC and IEEE, 2017, p. 196; Hamburg and Mogyorodi, 2024) as opposed to the capabilities of those functions themselves.
5. A typo in (ISO/IEC and IEEE, 2021c, Fig. 2) means that “specification-based techniques” is listed twice, when the latter should be “structure-based techniques”.
6. ISO/IEC and IEEE (2017) use the same definition for “partial correctness” (p. 314) and “total correctness” (p. 480).
7. Hamburg and Mogyorodi (2024) classify ML model testing as a test level, which they define as “a specific instantiation of a test process”: a vague definition that does not match the one in Table 2.1.
8. Washizaki (2024, p. 5-4) says that quality improvement, along with quality assurance, is an aspect of testing that involves “defining methods, tools, skills, and practices to achieve the specific quality level and objective”; while testing that a system possesses certain qualities is in scope, actively improving the system in response is *not* part of testing.
9. Washizaki (2024, p. 5-9)’s definition of “elasticity testing” only cites a single source **that doesn’t contain the words “elasticity” or “elastic”!**
10. The terms “acceleration tolerance testing” and “acoustic tolerance testing” do not seem refer to software testing, but Firesmith (2015, p. 56) includes them regardless. Elsewhere, they seem to refer to testing the acoustic tolerance of rats (Holley et al., 1996) or the acceleration tolerance of astronauts (Morgun et al., 1999, p. 11), aviators (Howe and Johnson, 1995, pp. 27, 42), or catalysts (Liu et al., 2023, p. 1463), which which are not relevant to software testing.
11. The definition of “math testing” given by Hamburg and Mogyorodi (2024) is too specific to be useful, likely taken from an example instead of a general definition: “testing to determine the correctness of the pay table implementation, the random number generator results, and the return to player computations”.
12. A similar issue exists with multiplayer testing, where its definition specifies “the casino game world” (Hamburg and Mogyorodi, 2024).
13. “Par sheet testing” from (Hamburg and Mogyorodi, 2024) seems to refer to the specific example mentioned in Mistake 11 and does not seem more widely

applicable, since a “PAR sheet” is “a list of all the symbols on each reel of a slot machine” (Bluejay, 2024).

14. The source that Hamburg and Mogyorodi (2024) cite for the definition of “test type” does not seem to actually provide a definition.
15. The same is true for “visual testing” (Hamburg and Mogyorodi, 2024).
16. The same is true for “security attack” (Hamburg and Mogyorodi, 2024).
17. Reid (1996, p. 4) says “the use of the term ‘condition’ in branch condition testing can mislead the reader into thinking all conditions are exercised”, which Hamburg and Mogyorodi (2024) seem to do by giving “condition coverage” as a synonym of “branch condition coverage”.
18. Doğan et al. (2014, p. 184) claim that Sakamoto et al. (2013) define “prime path coverage”, but they do not.
19. The differences between the terms “error”, “failure”, “fault”, “defect” are significant and meaningful (ISO/IEC and IEEE, 2010, pp. 128, 139–140; Washizaki, 2024, p. 12-3; van Vliet, 2000, pp. 399–400), but Patton (2006, 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”!
20. Peters and Pedrycz (2000, Fig. 12.31) imply that decision coverage is a child of both c-use coverage *and* p-use coverage; this seems incorrect, since decisions are the result of p-uses and *not* c-uses (ISO/IEC and IEEE, 2021c, pp. 5, 27; 2017, p. 332; van Vliet, 2000, p. 424), and only the p-use relation is implied by ISO/IEC and IEEE (2021c, Fig. F.1).
21. Sharma et al. (2021, p. 601) seem to use the terms “grey-box testing” and “(stepwise) code reading” interchangeably, which would incorrectly imply that they are synonyms.

FIND SOURCES

22. Kam (2008, p. 46) gives “program testing” as a synonym of “component testing” but it would make more sense as a synonym of “system testing”, which is conducted on the system, or program, as a whole (ISO/IEC and IEEE, 2017, p. 456; Hamburg and Mogyorodi, 2024; Peters and Pedrycz, 2000, Tab. 12.3; van Vliet, 2000, p. 439; Sakamoto et al., 2013, pp. 343–344).
23. Kam (2008, p. 46) gives “mutation testing” as a synonym of “back-to-back testing”; while the two are related (ISO/IEC and IEEE, 2010, p. 30), the variants used in mutation testing are generated or designed to be detected as incorrect by the test suite (Washizaki, 2024, p. 5-15; similar in van Vliet, 2000, pp. 428–429) which is not a requirement of back-to-back testing.

OG 2012

OG Beizer

24. [Kam \(2008, p. 46\)](#) 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” ([Hamburg and Mogyorodi, 2024](#)) (i.e., negative testing) is one way to help test this!
25. [Kam \(2008, p. 42\)](#) says “See *boundary value analysis*,” for the glossary entry of “boundary value testing” but does not include “boundary value analysis” in the glossary.
26. [Kam \(2008\)](#) misspells “state-based” as “state-base” (pp. 13, 15) and “stated-base” (Tab. 1).
27. [Gerrard \(2000b, p. 28\)](#)’s definition of “manual security audits” may be too specific, only applying to “the products installed on a site” and “the known vulnerabilities for those products”.

OG Hetzel88

28. [Sneed and Göschl \(2000, p. 18\)](#) give “white-box testing”, “grey-box testing”, and “black-box testing” as synonyms for “module testing”, “integration testing”, and “system testing”, respectively, but this mapping is incorrect; for example, [ISO/IEC and IEEE \(2017, p. 444\)](#) say “structure-based [(or white-box)] testing is not restricted to use at component level and can be used at all levels” and [Sakamoto et al. \(2013, pp. 345–346\)](#) describe “black-box integration testing”.

more examples?

29. The previous flaw makes the claim that “red-box testing” is a synonym for “acceptance testing” ([Sneed and Göschl, 2000, p. 18](#)) lose credibility.

C.1.2 Full List of Omissions

1. [ISO/IEC and IEEE \(2021c\)](#) cite [Reid \(1996\)](#) as the source for their Fig. F.1 but they omit Linear Code Sequence and Jump (LCSAJ) testing with no explanation, both from this figure and from the document as a whole. They label the figure as a “partial ordering”, which might explain its omission from Fig. F.1, but [Reid \(1996, p. 7\)](#) already identifies that his hierarchy is incomplete as “it relates only a subset of the available test completion criteria, so other criteria ... should still be considered”.
2. The acronym for System of Systems (SoS) ([ISO/IEC and IEEE, 2019b](#)) is used but not defined by [Firesmith \(2015, p. 23\)](#).
3. The [Project Management Institute \(2013, p. 476\)](#) uses the acronym “QA” which is implied to refer to “quality assurance” as “QC” refers to “quality control” (p. 535), but this is not made explicit.
4. [Van Vliet \(2000, p. 425\)](#) defines many types of data flow coverage, including all-p-uses, all-p-uses/some-c-uses, and all-c-uses/some-p-uses, but excludes all-c-uses, which is implied by these definitions and defined elsewhere

(ISO/IEC and IEEE, 2021c, p. 27; 2017, p. 83; Peters and Pedrycz, 2000, p. 479).

5. Bas (2024, p. 16) lists “three [backup] location categories: local, offsite and cloud based [sic]” but does not define or discuss “offsite backups” (pp. 16–17).
6. Gerrard (2000a, Tab. 2) makes a distinction between “transaction verification” and “transaction testing” and uses the phrase “transaction flows” (Fig. 5) but doesn’t explain them.
7. Availability testing is not assigned to a test priority (Gerrard, 2000a, Tab. 2), despite the claim that “the test types¹ have been allocated a slot against the four test priorities” (p. 13); usability testing and/or performance testing would have been good candidates.

C.1.3 Full List of Contradictions

1. Regression testing and retesting are sometimes given as two distinct approaches (ISO/IEC and IEEE, 2022, p. 8; 2021a, p. 3; Firesmith, 2015, p. 34), but sometimes regression testing is defined as a form of “selective retesting” (ISO/IEC and IEEE, 2017, p. 372; Washizaki, 2024, pp. 5-8, 6-5, 7-5 to 7-6; Barbosa et al., 2006, p. 3). Moreover, the two possible variations of regression testing given by van Vliet (2000, p. 411) are “retest-all” and “selective retest”, which is possibly the source of the above misconception. This creates a cyclic relation between regression testing and selective retesting.
2. Accessibility testing is a subtype of usability testing (ISO/IEC and IEEE, 2022, p. 1; 2021c, Tab. A.1; 2017, p. 6; ISO/IEC, 2011; Firesmith, 2015, p. 58) but these two test types are listed at the same level by ISO/IEC and IEEE (2022, Fig. 2).
3. Integration testing, system testing, and system integration testing are all listed as separate test levels (ISO/IEC and IEEE, 2022, p. 12, Fig. 2; 2021b, p. 41–44, 46, 51, 58, 74; 2021c, p. 6), but system integration testing is listed as a child of both integration testing (Hamburg and Mogyorodi, 2024) and system testing (Firesmith, 2015, p. 23).
4. Similarly, the relations between component testing, integration testing, and component integration testing are unclear; in particular, Hamburg and Mogyorodi (2024) seem to give integration testing as a parent of component integration testing in the latter’s definition, but as a child in their graph of static and dynamic testing.

¹“Each type of test addresses a different risk area” (Gerrard, 2000a, p. 12), which is distinct from the notion of “test type” described in Table 2.1.

Are these separate approaches?

5. “Software testing” is often defined to exclude static testing (Firesmith, 2015, p. 13; Ammann and Offutt, 2017, p. 222; Peters and Pedrycz, 2000, p. 439), restricting “testing” to mean dynamic validation (Washizaki, 2024, p. 5-1) or verification “in which a system or component is executed” (ISO/IEC and IEEE, 2017, p. 427). However, “terminology is not uniform among different communities, and some use the term ‘testing’ to refer to static techniques² as well” (Washizaki, 2024, p. 5-2). This is done by ISO/IEC and IEEE (2022, pp. 16–17) and Gerrard (2000a, pp. 8–9); the former even explicitly *exclude* static testing in another document (2017, p. 440)!
6. When static testing *is* included as part of software testing, it is not categorized consistently. ISO/IEC and IEEE categorize it as a test level in (2021b, p. 43) but give it its own test approach category in (2022, p. 10, 23, Fig. 2).
7. ISO/IEC (2023a) and ISO/IEC and IEEE (2017, p. 196) both say that functional suitability is “concerned with whether the functions meet stated and implied needs”, but the former includes “the functional specification” as part of its scope while the latter explicitly excludes it.
8. A component is an “entity with discrete structure ... within a system considered at a particular level of analysis” (ISO/IEC, 2023b) 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 (ISO/IEC and IEEE, 2017, p. 82). For example, Hamburg and Mogyorodi (2024) define them as synonyms while Baresi and Pezzè (2006, p. 107) say “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 (ISO/IEC and IEEE, 2017, p. 419) and “can be tested in isolation” (Hamburg and Mogyorodi, 2024), “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.
9. While a computation data use is defined as the “use of the value of a variable in *any* type of statement” (ISO/IEC and IEEE, 2021c, p. 2; 2017, p. 83, emphasis added), it is often qualified to *not* be a predicate data use (van Vliet, 2000, p. 424; implied by ISO/IEC and IEEE, 2021c, p. 27).
10. ISO/IEC and IEEE define an “extended entry (decision) table” both as a decision table where the “conditions consist of multiple values rather than simple Booleans” (2021c, p. 18) and one where “the conditions and actions are generally described but are incomplete” (2017, p. 175).
11. ISO/IEC and IEEE (2021c, Fig. F.1) is an adaptation of Reid (1996, Fig. 2) and one of the changes they make is replacing “branch [coverage]” with “branch/decision coverage”. Reid notes that “the term decision coverage is

²Not formally defined, but distinct from the notion of “test technique” described in Table 2.1.

used interchangeably with that of branch coverage” but that comparing one to the other is not a direct mapping (p. 4). [ISO/IEC and IEEE](#) agree, saying “branch and decision coverage are closely related..., although lower levels of coverage may not be the same” ([2021c](#), p. 104) and separating the terms in (Fig. G.1, Secs. 5.3.2, 5.3.3, Annex C.2.2). However, (Fig. F.1) implies that these terms are synonyms, contradicting this separation and making decision testing’s relations ambiguous.

Q #14: What’s the short form of “annex”?

12. [ISO/IEC and IEEE](#) (2017, p. 83)’s definition of “all-c-uses testing”—testing that aims to execute all data “use[s] of the value of a variable in any type of statement”—is *much* more vague than the definition they give in ([2021c](#), p. 27; similar in [van Vliet, 2000](#), p. 425; [Peters and Pedrycz, 2000](#), p. 479)—testing that exercises “control flow sub-paths from each variable definition to each c-use of that definition (with no intervening definitions)”.
13. Since keyword-driven testing can be used for automated *or* manual testing ([ISO/IEC and IEEE, 2016](#), pp. 4, 6), the claim that “test cases can be either manual test cases or keyword test cases” (p. 6) is incorrect.
14. [ISO/IEC and IEEE](#) (2022, p. 36) say “A/B testing is not a test case generation technique as test inputs are not generated”, where “test case generation technique” may be a synonym of “test design technique”. However, the inclusion of A/B testing under the heading “Test design and execution” in the same document implies that it may be considered a test technique.³
15. The claim that “test cases can be either manual test cases or keyword test cases” ([ISO/IEC and IEEE, 2016](#), p. 6) implies that “keyword-driven testing” could be a synonym of “automated testing” instead of its child, which seems more reasonable (p. 4; [2022](#), p. 35).
16. Performance testing and security testing are given as subtypes of reliability testing by [ISO/IEC](#) (2023a), but these are all listed separately by [Firesmith](#) (2015, p. 53).
17. Similarly, random testing is a subtechnique of specification-based testing ([ISO/IEC and IEEE, 2022](#), pp. 7, 22; [2021c](#), pp. 5, 20, Fig. 2; [Washizaki, 2024](#), p. 5-12; [Hamburg and Mogyorodi, 2024](#)) but is listed separately by [Firesmith](#) (2015, p. 46).
18. Path testing “aims to execute all entry-to-exit control flow paths in a SUT’s control flow graph” ([Washizaki, 2024](#), p. 5-13; similar in [Patton, 2006](#), p. 119), but [ISO/IEC and IEEE](#) (2017, p. 316) add that it can also be “designed to execute ... selected paths.”
19. The structure of tours can be defined as either quite general ([ISO/IEC and IEEE, 2022](#), p. 34) or “organized around a special focus” ([Hamburg and Mogyorodi, 2024](#)).

³For simplicity, this implied categorization as “technique” is omitted from [Table C.1](#).

20. Alpha testing is performed by “users within the organization developing the software” (ISO/IEC and IEEE, 2017, p. 17), “a small, selected group of potential users” (Washizaki, 2024, p. 5-8), or “roles outside the development organization” conducted “in the developer’s test environment” (Hamburg and Mogyorodi, 2024).
21. “Conformance testing” is defined by Washizaki (2024, p. 5-7) as testing that “aims to verify that the SUT conforms to standards, rules, specifications, requirements, design, processes, or practices”, but this disagrees with the definition given by the Project Management Institute (2013, p. 523): testing that evaluates the degree to which “results ... fall within the limits that define acceptable variation for a quality requirement”. Washizaki’s definition instead seems to correspond to the definition of compliance testing given by Hamburg and Mogyorodi (2024) and Firesmith (2015, p. 33), which may explain why Kam (2008, p. 43) gives them as synonyms (along with his unhelpful definition of compliance testing: “testing to determine the compliance of the component or system”).
22. The terms “test level” and “test stage” are given as synonyms (Hamburg and Mogyorodi, 2024; implied by ISO/IEC and IEEE, 2015, p. 9; Gerrard, 2000a, p. 9), but Washizaki (2024, 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 2.1.3) to “test level” instead. However, the examples of “test stages” listed—unit testing, integration testing, system testing, and acceptance testing (Washizaki, 2024, pp. 5-6 to 5-7)—are commonly categorized as “test levels” (see Section 2.1.1).
23. “Operational acceptance testing” and “production acceptance testing” are given as synonyms by Hamburg and Mogyorodi (2024) but listed separately by Firesmith (2015, p. 30).
24. Retesting and regression testing seem to be categorized separately from the rest of the testing approaches (ISO/IEC and IEEE, 2022, pp. 15, 23; 2021a, p. 8; 2021b, p. 4) but this is not justified. Hamburg and Mogyorodi (2024) consider regression testing to be a test type and Barbosa et al. (2006, p. 3) consider it a test level; since it is not included as an example of a test level by the sources that describe them (see Section 2.1.1), this latter categorization is likely not universal at best and incorrect at worst.
25. While Patton (2006, p. 120) implies that condition testing is a subtechnique of path testing, van Vliet (2000, 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.
26. Load testing is performed with loads “between anticipated conditions of low, typical, and peak usage” (ISO/IEC and IEEE, 2022, p. 5) or loads that are as large as possible (Patton, 2006, p. 86).

27. State testing requires that “all states in the state model ... [are] ‘visited’ ” (ISO/IEC and IEEE, 2021c, p. 19), but Patton (2006, pp. 82–83) lists this as only one of its possible criteria.
28. System testing is “conducted on a complete, integrated system” (ISO/IEC and IEEE, 2017, p. 456; similar in Peters and Pedrycz, 2000, Tab. 12.3; van Vliet, 2000, p. 439), but Patton (2006, p. 109) says it can also be done on “at least a major portion” of the product.
29. “Walkthroughs” and “structured walkthroughs” are given as synonyms by Hamburg and Mogyorodi (2024) but Peters and Pedrycz (2000, p. 484) imply that they are different, saying a more structured walkthrough may have specific roles.
30. While Patton (2006, p. 92, emphasis added) says that reviews are “*the* process[es] under which static white-box testing is performed”, van Vliet (2000, pp. 418–419) gives correctness proofs as another example.
31. Different capitalizations of the abbreviations of “computation data use” and “predicate data use” are used: the lowercase “c-use” and “p-use” (ISO/IEC and IEEE, 2021c, pp. 3, 27–29, 35–36, 114–155, 117–118, 129; 2017, p. 124; Peters and Pedrycz, 2000, p. 477, Tab. 12.6; Reid, 1996, Fig. 2) and the uppercase “C-use” and “P-use” (van Vliet, 2000, pp. 424–425).
32. Similarly for “definition-use” (such as in “definition-use path”), both the lowercase “du” (ISO/IEC and IEEE, 2021c, pp. 3, 27, 29, 35, 119–121, 129; Peters and Pedrycz, 2000, pp. 478–479; Reid, 1996, Fig. 2) and the uppercase “DU” (van Vliet, 2000, p. 425) are used.
33. Van Vliet (2000, pp. 424–425) specifies that every successor of a data definition use needs to be executed as part of all-uses testing, but this condition is not included elsewhere (ISO/IEC and IEEE, 2021c, pp. 28–29; 2017, p. 120; Peters and Pedrycz, 2000, pp. 478–479).
34. All-du-paths testing is usually defined as exercising all “loop-free control flow sub-paths from each variable definition to every use (both p-use and c-use) of that definition (with no intervening definitions)” (ISO/IEC and IEEE, 2021c, p. 29; similar in 2017, p. 125; Washizaki, 2024, p. 5-13; Peters and Pedrycz, 2000, p. 479); however, paths containing simple cycles may also be required (van Vliet, 2000, p. 425).
35. Similarly, van Vliet (2000, pp. 432–433) says that all-du-paths testing is only stronger than all-uses testing if there are infeasible paths, but Washizaki (2024, p. 5-13) does not specify this caveat.
36. Acceptance testing is “usually performed by the purchaser ... with the ... vendor” (ISO/IEC and IEEE, 2017, p. 5), “may or may not involve the developers of the system” (Bourque and Fairley, 2014, p. 4-6), and/or “is often performed under supervision of the user organization” (van Vliet, 2000,

p. 439); these descriptions of who the testers are contradict each other *and* all introduce some uncertainty (“usually”, “may or may not”, and “often”, respectively).

Q #15: Does this merit counting this as an Ambiguity as well as a Contradiction?

OG Beizer

37. Although ad hoc testing is classified as a “technique” (Washizaki, 2024, p. 5-14), it is one in which “no recognized test design technique is used” (Kam, 2008, p. 42).

38. Kam (2008, p. 46) says “negative testing is related to the testers’ attitude rather than a specific test approach or test design technique”; while ISO/IEC and IEEE (2021c) seem to support this idea of negative testing being at a “higher” level than other approaches, they also imply that it is a test technique (pp. 10, 14).

C.1.4 Full List of Ambiguities

1. “Data definition” is defined as a “statement where a variable is assigned a value” (ISO/IEC and IEEE, 2021c, p. 3; 2017, p. 115; similar in 2012, p. 27; van Vliet, 2000, p. 424), but for functional programming languages such as Haskell with immutable variables (Wikibooks Contributors, 2023), this could cause confusion and/or be imprecise.

2. The distinctions between development testing (ISO/IEC and IEEE, 2017, p. 136), developmental testing (Firesmith, 2015, p. 30), and developer testing (p. 39; Gerrard, 2000a, p. 11) are unclear and seem miniscule.

Is this a def flaw?

3. Hamburg and Mogyorodi (2024) define “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.

4. While “error” is defined as “a human action that produces an incorrect result” (Washizaki, 2024, pp. 12-3; van Vliet, 2000, p. 399), Washizaki does not use this consistently, sometimes implying that errors can be intrinsic to software itself (2024, pp. 4-9, 6-5, 7-3, 12-4, 12-9, 12-13).

OG [14]

Q #16: I ignore “syntax errors, runtime errors, and logical errors” (Washizaki, 2024, p. 16-13, cf. p. 18-13) since they seem to be in different domains. Does that make sense? How should I document this?

5. Washizaki (2024, p. 1-1) defines “defect” as “an observable difference between what the software is intended to do and what it does”, but this seems to match the definition of “failure”: the inability of a system “to perform a required function or ... within previously specified limits” (ISO/IEC and IEEE, 2019a, p. 7; 2010, p. 139; similar in van Vliet, 2000, p. 400) that is “externally visible” (ISO/IEC and IEEE, 2019a, p. 7; similar in van Vliet, 2000, p. 400).

6. Similarly, Washizaki (2024, p. 4-11) says “fault tolerance is a collection of techniques that increase software reliability by detecting errors and then recovering from them or containing their effects if recovery is not possible.”

OG ISO/IEC, 2005

This should either be called “*error* tolerance” or be described as “detecting *faults* and then recovering from them”, since “error” and “fault” have distinct meanings (Washizaki, 2024, p. 5-3; van Vliet, 2000, pp. 399–400). The intent of this term-definition pair is unclear, as the strategies given—“backing up and retrying, using auxiliary code and voting algorithms, and replacing an erroneous value with a phony value that will have a benign effect” (Washizaki, 2024, p. 4-11)—could be used for errors or faults.

7. While ergonomics testing is out of scope (as it tests hardware, not software; see [Appendix A.1](#)), its definition of “testing to determine whether a component or system and its input devices are being used properly with correct posture” (Hamburg and Mogyorodi, 2024) seems to focus on how the system is *used* as opposed to the system *itself*.
8. Similarly, end-to-end testing is defined as testing “in which business processes are tested from start to finish under production-like circumstances” (Hamburg and Mogyorodi, 2024); it is unclear whether this tests the business processes themselves or the system’s role in performing them.
9. Hamburg and Mogyorodi (2024) describe the term “software in the loop” as a kind of testing, while the source they reference seems to describe “Software-in-the-Loop-Simulation” as a “simulation environment” that may support software integration testing (Knüvener Mackert GmbH, 2022, p. 153); is this a testing approach or a tool that supports testing?
10. While model testing is said to test the object under test, it seems to describe testing the models themselves (Firesmith, 2015, p. 20); using the models to test the object under test seems to be called “driver-based testing” (p. 33).
11. Similarly, it is ambiguous whether “tool/environment testing” refers to testing the tools/environment *themselves* or *using* them to test the object under test; the wording of its subtypes (Firesmith, 2015, p. 25) seems to imply the former.
12. The Project Management Institute (2013, p. 244; similar on p. 535) says “quality assurance work will fall under the conformance work category in the cost of quality framework”, but (Fig. 8-2) suggests that “conformance work” is a part of quality assurance. This introduces ambiguity at best and creates a cyclic parent-child relation (which violates our definition in [Section 2.1.3](#)) at worst.
13. Hamburg and Mogyorodi (2024) claim that code inspections are related to peer reviews but Patton (2006, pp. 94–95) makes them quite distinct.
14. Patton (2006, p. 119) says that branch testing is “the simplest form of path testing” which is also implied by ISO/IEC and IEEE (2021c, Fig. F.1) and van Vliet (2000, p. 433). This is true in the example Patton gives, but is not necessarily generalizable; one could test the behaviour at branches without

Q #17: Should I be more specific?

testing even a *subset* of complete paths, which ISO/IEC and IEEE (2017, p. 316) give as a definition of “path testing” (see [Contradiction 18](#))!

C.1.5 Full List of Overlaps

1. ISO/IEC and IEEE (2022, p. 34) give 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)?
2. ISO/IEC and IEEE (2017, pp. 469, 470; 2013, p. 9) say that “test level” and “test phase” are synonyms, both meaning a “specific instantiation of [a] test sub-process”, 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” (2022, p. 24) and “test phase” can also refer to the “period of time in the software life cycle” when testing occurs (2017, p. 470), usually after the implementation phase (pp. 420, 509; Perry, 2006, p. 56).
3. ISO/IEC and IEEE (2010, p. 128) define “error” as “a human action that produces an incorrect result”, but also as “an incorrect result” itself. Since faults are inserted when a developer makes an error (pp. 128, 140; Washizaki, 2024, p. 12-3; van Vliet, 2000, pp. 399–400), this means that they are “incorrect results”, making “error” and “fault” synonyms and the distinction between them less useful.
4. 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” (ISO/IEC and IEEE, 2010, p. 128; similar in Washizaki, 2024, 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” (ISO/IEC and IEEE, 2010, p. 30; similar in Hamburg and Mogyorodi, 2024), which seems to refer to this definition of “error”.
5. The SWEBOK Guide V4 defines “privacy testing” as testing that “assess[es] the security and privacy of users’ personal data to prevent local attacks” (Washizaki, 2024, p. 5-10). This seems to overlap (both in scope and name) with the definition of “security testing” in (ISO/IEC and IEEE, 2022, 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.

Q #18: Should we classify this as a Synonym Flaw at all?

Q #19: Shouldn’t this be “is”, referring to “test item”?

6. “Orthogonal array testing” (Washizaki, 2024, pp. 5-1, 5-11; implied by Valcheva, 2013, pp. 467, 473; Yu et al., 2011, pp. 1573–1577, 1580) and “operational acceptance testing” (Firesmith, 2015, p. 30) have the same acronym (“OAT”).
7. “Customer acceptance testing” and “contract(ual) acceptance testing” have the same acronym (“CAT”) (Firesmith, 2015, p. 30).
8. The same is true for “hardware-in-the-loop testing” and “human-in-the-loop testing” (“HIL”) (Firesmith, 2015, p. 23), although Preuße et al. (2012, p. 2) use “HiL” for the former.

C.1.6 Full List of Redundancies

1. ISO/IEC and IEEE (2021c, p. 4) define “exit point” as the “last executable statement within a test item”, then later note that it “is most commonly the last executable statement within the test item”.
2. ISO/IEC and IEEE (2017, p. 375) say that “dependability characteristics include availability and its inherent or external influencing factors, such as availability”.
3. ISO/IEC and IEEE (2017, p. 228) provide a definition for “inspections and audits”, despite also giving definitions for “inspection” (p. 227) and “audit” (p. 36); while the first term *could* be considered a superset of the latter two, this distinction doesn’t seem useful.
4. “Ethical hacking testing” is given as a synonym of penetration testing by Washizaki (2024, p. 13-4), which seems redundant; Gerrard (2000b, p. 28) uses the term “ethical hacking” which is clearer.
5. While correct, ISTQB’s definition of “specification-based testing” is not helpful: “testing based on an analysis of the specification of the component or system” (Hamburg and Mogyorodi, 2024).
6. The phrase “continuous automated testing” (Gerrard, 2000a, p. 11) is redundant since Continuous Testing (CT) is already a subapproach of automated testing (ISO/IEC and IEEE, 2022, p. 35; Hamburg and Mogyorodi, 2024).

C.2 Full Lists of Flaws by Domain

The following sections provide all of the data that we automatically detect (as described in [Section 4.2.1](#)) and summarize in [Section 5.2](#).

C.2.1 Multiple Categorizations

As mentioned in [Section 5.2.1](#), we automatically detect test approaches with more than one category that violate our assumption of orthogonality (see [Section 2.1.1](#)) and list them in [Table C.1](#).

Table C.1: Test approaches with more than one category.

Approach	Category 1	Category 2
Ad Hoc Testing	Practice (ISO/IEC and IEEE, 2013 , p. 33)	Technique (Washizaki, 2024 , p. 5-14)
Capacity Testing	Technique (ISO/IEC and IEEE, 2021c , p. 38–39)	Type (ISO/IEC and IEEE, 2022 , p. 22; 2013 , p. 2; implied by its quality (ISO/IEC, 2023a ; ISO/IEC and IEEE, 2021c , Tab. A.1); Firesmith, 2015 , p. 53)
Checklist-based Testing	Practice (ISO/IEC and IEEE, 2022 , p. 34)	Technique (Hamburg and Mogyorodi, 2024)
Data-driven Testing	Practice (ISO/IEC and IEEE, 2022 , p. 22)	Technique (Kam, 2008 , p. 43; OG Fewster and Graham)
End-to-end Testing	Type (Hamburg and Mogyorodi, 2024)	Technique (Firesmith, 2015 , p. 47; Sharma et al., 2021 , pp. 601, 603, 605–606)

Continued on next page

Table C.1: Test approaches with more than one category. (Continued)

Approach	Category 1	Category 2
Endurance Testing	Technique (ISO/IEC and IEEE, 2021c, p. 38–39)	Type (ISO/IEC and IEEE, 2013, p. 2; implied by Firesmith, 2015, p. 55)
Error Guessing	Practice (ISO/IEC and IEEE, 2013, p. 33)	Technique (ISO/IEC and IEEE, 2022, pp. 4, 34, Fig. 2; 2021c, pp. iii–iv, 4, 11, 29, 35, 122, 125, Fig. 2, Tab. A.2; 2013, pp. 3, 33; Washizaki, 2024, p. 5-13; Firesmith, 2015, p. 50)
Experience-based Testing	Technique (ISO/IEC and IEEE, 2022, Fig. 2; Firesmith, 2015, pp. 46, 50)	Practice (ISO/IEC and IEEE, 2022, Fig. 2; 2021c, p. viii; 2013, pp. iii, 31, 33)
Exploratory Testing	Technique (Washizaki, 2024, p. 5-14; Firesmith, 2015, p. 50)	Practice (ISO/IEC and IEEE, 2022, pp. 11, 20, 34, Fig. 2; 2021a, p. 5; 2021c, p. viii; 2013, pp. 13, 33; implied by 2022, p. 33)
Functional Testing	Technique (Barbosa et al., 2006, p. 3; inferred from specification-based testing)	Type (ISO/IEC and IEEE, 2022, pp. 15, 20, 22; 2021a, pp. 8, 11; 2021b, p. 41; 2021c, pp. 7, 38, Tab. A.1; 2017, p. 473; 2016, p. 4; Hamburg and Mogyorodi, 2024; implied by the quality of “correctness” (ISO/IEC and IEEE, 2017, p. 104; Washizaki, 2024, p. 3-13))
Load Testing	Technique (ISO/IEC and IEEE, 2021c, p. 38–39)	Type (ISO/IEC and IEEE, 2022, pp. 5, 20, 22; 2017, p. 253; OG IEEE 2013; Hamburg and Mogyorodi, 2024; implied by Firesmith, 2015, p. 54)

Continued on next page

Table C.1: Test approaches with more than one category. (Continued)

Approach	Category 1	Category 2
Model-based Testing	Technique (Souza et al., 2017, p. 3; implied by ISO/IEC and IEEE, 2017, p. 469)	Practice (ISO/IEC and IEEE, 2022, p. 11, Fig. 2; 2021a, p. 5; 2021c, p. viii; 2013, pp. iii, 31)
Penetration Testing	Technique (ISO/IEC and IEEE, 2021c, p. 40; Hamburg and Mogyorodi, 2024)	Type (ISO/IEC and IEEE, 2021b, pp. 41, 43; implied by Firesmith, 2015, p. 57; inferred from security testing)
Performance Testing	Technique (ISO/IEC and IEEE, 2021c, p. 38–39)	Type (ISO/IEC and IEEE, 2022, pp. 7, 22, 26–27; 2021a, pp. 2, 8, 11; 2021b, pp. 41, 43; 2021c, p. 7; implied by Firesmith, 2015, p. 53)
Regression Testing	Type (Hamburg and Mogyorodi, 2024)	Level (Barbosa et al., 2006, p. 3)
Specification-based Testing	Type (Hamburg and Mogyorodi, 2024)	Technique (ISO/IEC and IEEE, 2022, p. 22; 2021b, p. 45; 2021c, p. 8; Washizaki, 2024, p. 5-10; Hamburg and Mogyorodi, 2024; Firesmith, 2015, pp. 46–47; Souza et al., 2017, p. 3; Sakamoto et al., 2013, p. 344; implied by ISO/IEC and IEEE, 2022, pp. 2–4, 6–9)
Stress Testing	Technique (ISO/IEC and IEEE, 2021c, p. 38–39)	Type (ISO/IEC and IEEE, 2022, pp. 9, 22; 2017, p. 442; implied by Firesmith, 2015, p. 54)

Continued on next page

Table C.1: Test approaches with more than one category. (Continued)

Approach	Category 1	Category 2
Structure-based Testing	Type (Hamburg and Mogyorodi, 2024)	Technique (ISO/IEC and IEEE, 2022, p. 22; 2021b, p. 45; 2021c, p. 8; Washizaki, 2024, pp. 5-10, 5-13; Hamburg and Mogyorodi, 2024; Firesmith, 2015, pp. 46, 49; Souza et al., 2017, p. 3; Sakamoto et al., 2013, p. 344; implied by ISO/IEC and IEEE, 2022, pp. 2, 4, 6, 9; Barbosa et al., 2006, p. 3)
Alpha Testing	Type (implied by Firesmith, 2015, p. 58)	Level (ISO/IEC and IEEE, 2022, p. 22; inferred from acceptance testing)
Attacks	Practice (ISO/IEC and IEEE, 2022, p. 34; 2013, p. 33)	Technique (implied by Hamburg and Mogyorodi, 2024)
Beta Testing	Type (implied by Firesmith, 2015, p. 58)	Level (ISO/IEC and IEEE, 2022, p. 22; inferred from acceptance testing)
Integration Testing	Technique (implied by Sharma et al., 2021, pp. 601, 603, 605–606)	Level (ISO/IEC and IEEE, 2022, pp. 12, 20–22, 26–27; 2021a, pp. 6, 11; 2021b, Fig. 2, pp. 41, 43, 51; 2021c, p. 6; Washizaki, 2024, p. 5-7; Hamburg and Mogyorodi, 2024; Peters and Pedrycz, 2000, Tab. 12.3; van Vliet, 2000, p. 438; Souza et al., 2017, p. 3; Sakamoto et al., 2013, p. 343; Barbosa et al., 2006, p. 3)
Interface Testing	Type (Kam, 2008, p. 45)	Level (implied by ISO/IEC and IEEE, 2017, p. 235; Sakamoto et al., 2013, p. 343; inferred from integration testing)
Procedure Testing	Technique (implied by Firesmith, 2015, p. 47)	Type (ISO/IEC and IEEE, 2022, pp. 7, 22; 2021c, p. 39, Tab. A.1; 2017, p. 337; OG IEEE, 2013)
Survivability Testing	Technique (Ghosh and Voas, 1999, p. 39)	Type (implied by its quality (ISO/IEC, 2011); inferred from robustness testing and security testing)

Continued on next page

Table C.1: Test approaches with more than one category. (Continued)

Approach	Category 1	Category 2
Unit Testing	Technique (implied by Engström and Petersen, 2015, pp. 1–2)	Level (ISO/IEC and IEEE, 2022, pp. 12, 20–22, 26–27; 2021a, pp. 6, 11; 2021b, Fig. 2, pp. 41, 43, 51; 2021c, p. 6; 2017, p. 467; 2016, p. 4; Washizaki, 2024, p. 5-6; Hamburg and Mogyorodi, 2024; Peters and Pedrycz, 2000, Tab. 12.3; van Vliet, 2000, p. 438; Souza et al., 2017, p. 3; Sakamoto et al., 2013, p. 343; Barbosa et al., 2006, p. 3)
Volume Testing	Technique (ISO/IEC and IEEE, 2021c, p. 38–39)	Type (implied by Firesmith, 2015, p. 54; inferred from performance-related testing)
Infrastructure Testing	Type (implied by Firesmith, 2015, p. 57)	Level (implied by Gerrard, 2000a, p. 13; see Table 2.1)

C.2.2 Intransitive Synonyms

As described in [Section 2.2.2](#), intransitive synonyms are examples of flaws since they violate our definition of the synonym relation in [Section 2.1.2](#). We identify nine such cases through automatic analysis of generated visualizations, listed below (test approaches in *italics* are synonyms with each other, but not with other terms not in *italics*):

Better way to handle/display this?

1. Condition Testing:

- *Branch Condition Combination Testing* ([Patton, 2006](#), p. 120; [Sharma et al., 2021](#), Fig. 1)
- Branch Condition Testing ([Hamburg and Mogyorodi, 2024](#))
- *Decision Testing* ([Washizaki, 2024](#), p. 5-13)

2. Soak Testing:

- Endurance Testing ([ISO/IEC and IEEE, 2021c](#), p. 39)
- Reliability Testing⁴ ([Gerrard, 2000a](#), Tab. 2; [2000b](#), Tab. 1, p. 26)

3. Functional Testing:

- Behavioural Testing ([van Vliet, 2000](#), p. 399)
- Correctness Testing ([Washizaki, 2024](#), p. 5-7)
- Specification-based Testing ([ISO/IEC and IEEE, 2017](#), p. 196; [Kam, 2008](#), pp. 44–45; [van Vliet, 2000](#), p. 399; implied by [ISO/IEC and IEEE, 2021c](#), p. 129; [2017](#), p. 431)

4. Link Testing:

- Branch Testing (implied by [ISO/IEC and IEEE, 2021c](#), p. 24; [Reid, 1996](#), p. 4)
- Component Integration Testing ([Kam, 2008](#), p. 45)
- Integration Testing (implied by [Gerrard, 2000a](#), p. 13)

5. Decision Testing:

- Branch Condition Combination Testing ([Washizaki, 2024](#), p. 5-13)
- Branch Testing (often) ([Reid, 1996](#), p. 4)

6. Exhaustive Testing:

- Branch Condition Combination Testing (if “each subcondition is viewed as a single input”) ([Peters and Pedrycz, 2000](#), p. 464)

⁴Endurance testing is given as a child of reliability testing by [Firesmith \(2015, p. 55\)](#), although the terms are not synonyms.

- Path Testing⁵ (van Vliet, 2000, p. 421)

7. Static Verification:

- Static Assertion Checking⁶ (Chalin et al., 2006, p. 343)
- Static Testing (implied by Chalin et al., 2006, p. 343)

8. Testing-to-Fail:

- Forcing Exception Testing (implied by Patton, 2006, pp. 66–67, 78)
- Negative Testing (Patton, 2006, pp. 67, 78, 84–87)

9. Invalid Testing:

- Error Tolerance Testing (implied by Kam, 2008, p. 45)
- Negative Testing (Hamburg and Mogyorodi, 2024; implied by ISO/IEC and IEEE, 2021c, p. 10)

C.2.3 Synonym and Parent-Child Overlaps

As described in Section 5.2.3, there are also pairs of synonyms where one is a subapproach of the other; these relations cannot coexist as synonym relations are symmetric while parent-child relations are asymmetric as outlined in Sections 2.1.2 and 2.1.3, respectively. Below are all 18 of these pairs that we identify through automatic analysis of our generated visualizations as described in Section 4.2.1.

⁵This synonym relation is likely incorrect (see Section 5.2.3).

⁶Chalin et al. (2006, p. 343) list Runtime Assertion Checking (RAC) and Software Verification (SV) as “two complementary forms of assertion checking”; based on how the term “static assertion checking” is used by Lahiri et al. (2013, p. 345), it seems like this should be the complement to RAC instead.

Table C.2: Pairs of test approaches with a parent-child *and* synonym relation.

“Child” → “Parent”	Parent-Child Source(s)	Synonym Source(s)
All Transitions Testing → State Transition Testing	(ISO/IEC and IEEE, 2021c, p. 19)	(Kam, 2008, p. 15)
Co-existence Testing → Compatibility Testing	(ISO/IEC, 2023a; ISO/IEC and IEEE, 2022, p. 3; 2021c, Tab. A.1)	(ISO/IEC and IEEE, 2021c, p. 37)
Domain Testing → Specification-based Testing	(Peters and Pedrycz, 2000, Tab. 12.1)	(Washizaki, 2024, p. 5-10)
Fault Tolerance Testing → Robustness Testing ^a	(Firesmith, 2015, p. 56)	(Hamburg and Mogyorodi, 2024)
Functional Testing → Specification-based Testing ^b	(ISO/IEC and IEEE, 2021c, p. 38, Kam, 2008, p. 42)	(ISO/IEC and IEEE, 2017, p. 196; Kam, 2008, pp. 44–45, 48; van Vliet, 2000, p. 399; implied by ISO/IEC and IEEE, 2021c, p. 129; 2017, p. 431)
Orthogonal Array Testing → Pairwise Testing	(Mandl, 1985, p. 1055)	(Washizaki, 2024, p. 5-11; Valcheva, 2013, p. 473)
Path Testing → Exhaustive Testing	(Peters and Pedrycz, 2000, pp. 466–467, 476; implied by Patton, 2006, pp. 120–121)	(van Vliet, 2000, p. 421)
Performance Testing → Performance-related Testing	(ISO/IEC and IEEE, 2022, p. 22; 2021c, p. 38)	(Moghadam, 2019, p. 1187)

Continued on next page

Table C.2: Pairs of test approaches with a parent-child *and* synonym relation. (Continued)

“Child” → “Parent”	Parent-Child Source(s)	Synonym Source(s)
Static Analysis → Static Testing	(ISO/IEC and IEEE, 2022, pp. 9, 17, 25, 28; 2021a, pp. 3, 21; Hamburg and Mogyorodi, 2024; Gerrard, 2000a, Fig. 4, p. 12; 2000b, p. 3)	(Peters and Pedrycz, 2000, p. 438)
Structure-based Testing → Structural Testing	(Patton, 2006, pp. 105–121)	(ISO/IEC and IEEE, 2022, p. 9; 2017, pp. 443–444; Hamburg and Mogyorodi, 2024; implied by Barbosa et al., 2006, p. 3)
Use Case Testing → Scenario Testing ^c	(ISO/IEC and IEEE, 2021c, p. 20; OG Hass, 2008)	(Hamburg and Mogyorodi, 2024)
Branch Condition Combination Testing → Decision Testing	(implied by the caveat of “atomic conditions” in Hamburg and Mogyorodi, 2024)	(Washizaki, 2024, p. 5-13)
Branch Condition Combination Testing → Exhaustive Testing	(implied by Patton, 2006, p. 121)	(Peters and Pedrycz, 2000, p. 464)
Reviews → Structural Analysis	(Patton, 2006, p. 92)	(implied by Patton, 2006, p. 92)
Smoke Testing → Build Verification Testing	(implied by Hamburg and Mogyorodi, 2024)	(Washizaki, 2024, p. 5-14)

Continued on next page

Table C.2: Pairs of test approaches with a parent-child *and* synonym relation. (Continued)

“Child” → “Parent”	Parent-Child Source(s)	Synonym Source(s)
Structural Testing → Structure-based Testing ^d	(Peters and Pedrycz, 2000, p. 447)	(implied by Peters and Pedrycz, 2000, p. 447)
Exploratory Testing → Unscripted Testing	(ISO/IEC and IEEE, 2022, p. 33; 2021a, p. 2; 2017, p. 174; Firesmith, 2015, p. 45; implied by Washizaki, 2024, p. 5-14)	(implied by Kuřesovs et al., 2013, p. 214)
Beta Testing → User Testing	(implied by Firesmith, 2015, p. 39)	(implied by Firesmith, 2015, p. 39)

^a Fault tolerance testing may also be a subapproach of reliability testing (ISO/IEC and IEEE, 2017, p. 375; Washizaki, 2024, p. 7-10), which is distinct from robustness testing (Firesmith, 2015, p. 53).

^b Hamburg and Mogyorodi (2024) cite ISO/IEC and IEEE (2017, p. 431) for their definition of “functional testing” but exclude the transitive synonym relationship (see Section 2.1.2) they give to “specification-based testing”. These terms are also defined separately elsewhere (2022, Fig. 2; 2021c, pp. 8, 49, 125), further supporting that they are not synonyms.

^c ISO/IEC and IEEE (2022, Fig. 2) also list “use case testing” and “scenario testing” separately, further supporting that these terms are not synonyms.

^d Peters and Pedrycz (2000, p. 447) claim that “structural testing subsumes white box testing”, but also say “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!

C.3 Inferred Flaws

Throughout our research, we infer many potential flaws as described in [Section 2.3](#). Some of these have a conflicting source while others do not. Since these are more subjective and are based on our own judgement, we exclude them from any counts of the numbers of flaws but give them here for completeness.

C.3.1 Inferred Multiple Categorizations

As mentioned in [Section 5.2.1](#), we automatically detect test approaches with more than one category that violate our assumption of orthogonality (see [Section 2.1.1](#)). This includes those with categories that we infer based on our assumption that child approaches inherit their parents' categories (as described in [Section 2.3](#)). We list these approaches and their given and inferred categories (along with their relevant parent in parentheses) in [Table C.3](#).

Table C.3: Test approaches inferred to have more than one category.

Approach	Given Category	Inferred Category
A/B Testing	Practice (ISO/IEC and IEEE, 2022 , Fig. 2)	Type (usability testing)
Big-Bang Testing	Technique (Sharma et al., 2021 , pp. 601, 603, 605–606)	Level (integration testing)
Bottom-Up Testing	Technique (Sharma et al., 2021 , pp. 601, 603, 605–606)	Level (integration testing)
Fuzz Testing	Technique (Hamburg and Mogyorodi, 2024 ; Firesmith, 2015 , p. 51)	Practice (mathematical-based testing)
Memory Management Testing	Technique (ISO/IEC and IEEE, 2021c , p. 38–39)	Type (performance-related testing)
Privacy Testing	Technique (ISO/IEC and IEEE, 2021c , p. 40)	Type (security testing)
Sandwich Testing	Technique (Sharma et al., 2021 , pp. 601, 603, 605–606)	Level (integration testing)
Closed Beta Testing	Type (implied by Firesmith, 2015 , p. 58)	Level (beta testing)
Open Beta Testing	Type (implied by Firesmith, 2015 , p. 58)	Level (beta testing)

C.3.2 Inferred Intransitive Synonyms

In addition to the nine cases of intransitive synonyms we list in [Appendix C.2.1](#), we infer some synonym relations based on the terms’ definitions that create the following inferred flaws (some pairs of synonyms have sources; those that do not are ones we infer):

1. Production Acceptance Testing:

- Operational (Acceptance) Testing ([Hamburg and Mogyorodi, 2024](#))⁷
- Production Verification Testing⁸

2. Structure-based Testing:

- Implementation-oriented Testing
- Structural Testing (implied by [Peters and Pedrycz, 2000](#), p. 447)

3. Operational (Acceptance) Testing:

- Field Testing
- Qualification Testing

C.3.3 Inferred Parent Flaws

As discussed in [Section 5.2.3](#), some pairs of synonyms also have a parent-child relation, abusing the meaning of “synonym” and causing confusion. While [Table C.2](#) gives the cases where both relations are supported by the literature, some are less explicit. The following automatically generated lists contain examples where at least one of these conflicting relations is *not* explicitly supported by the literature but may, nonetheless, be correct. The relations in the first two lists are explicitly given in the literature but may be incorrect, while those in the third list are unsubstantiated by the literature and require more thought before a recommendation can be made.

Pairs given a parent-child relation

1. Programmer Testing → Developer Testing ([Firesmith, 2015](#), p. 39)

Pairs given a synonym relation

1. Dynamic Analysis → Dynamic Testing⁹ ([Peters and Pedrycz, 2000](#), p. 438; implied by [ISO/IEC and IEEE, 2017](#), p. 149; [Hamburg and Mogyorodi, 2024](#))
2. Structured Walkthroughs → Walkthroughs¹⁰ ([Hamburg and Mogyorodi, 2024](#))

⁷See [Contradiction 23](#).

⁸“Production acceptance testing” ([Firesmith, 2015](#), p. 30) seems to be the same as “production verification testing” ([ISO/IEC and IEEE, 2022](#), p. 22) but neither is defined.

⁹The proposed relation is inferred from its static counterpart ([ISO/IEC and IEEE, 2022](#), pp. 9, 17, 25, 28; [2021a](#), pp. 3, 21; [Hamburg and Mogyorodi, 2024](#); [Gerrard, 2000a](#), Fig. 4, p. 12; [2000b](#), p. 3).

¹⁰See [Contradiction 29](#).

3. Functionality Testing → Functional Suitability Testing (implied by [Hamburg and Mogyorodi, 2024](#))

Pairs that could have a parent/child *or* synonym relation

1. Computation Flow Testing → Computation Testing
2. Field Testing → Operational Testing
3. Monkey Testing → Fuzz Testing
4. Organization-based Testing → Role-based Testing¹¹
5. Scenario-based Evaluations → Scenario-based Testing
6. System Qualification Testing → System Testing

In addition to these flaws, ([Gerrard, 2000a](#), Tab. 2) does *not* give “functionality testing” as a parent of “end-to-end functionality testing” as we infer he should (see [Appendix A.6](#)).

¹¹The distinction between organization- and role-based testing in ([Firesmith, 2015](#), pp. 17, 37, 39) seems arbitrary, but further investigation may prove it to be meaningful.

Appendix D

Code Snippets

Source Code D.1: [Tests for main with an invalid input file](#)

```
# from
↳ https://stackoverflow.com/questions/54071312/how-to-pass-command-line-arguments
## \brief Tests main with invalid input file
# \par Types of Testing:
# Dynamic Black-Box (Behavioural) Testing
# Boundary Conditions
# Default, Empty, Blank, Null, Zero, and None
# Invalid, Wrong, Incorrect, and Garbage Data
# Logic Flow Testing
@mark.parametrize("filename", invalid_value_input_files)
@mark.xfail
def test_main_invalid(monkeypatch, filename):
    # from
    ↳ https://stackoverflow.com/questions/10840533/most-pythonic-way-to-delete-a-file
    try:
        remove(output_filename)
    except OSError as e: # this would be "except OSError, e:"
        ↳ before Python 2.6
        if e.errno != ENOENT: # no such file or directory
            raise # re-raise exception if a different error
            ↳ occurred

    assert not path.exists(output_filename)

    with monkeypatch.context() as m:
        m.setattr(sys, 'argv', ['Control.py',
            ↳ str(Path("test/test_input") / f"{filename}.txt")])
```

```
Control.main()
```

```
assert not path.exists(output_filename)
```

Source Code D.2: [Projectile’s choice for constraint violation behaviour in code](#)

```
srsConstraints = makeConstraints Warning Warning,
```

Source Code D.3: [Projectile’s manually created input verification requirement](#)

```
verifyParamsDesc = foldlSent [S "Check the entered", plural
  ↪ inValue,
  S "to ensure that they do not exceed the" +:+. namedRef (datCon
    ↪ [] []) (plural datumConstraint),
  S "If any of the", plural inValue, S "are out of bounds" `sC`
  S "an", phrase errMsg, S "is displayed" `S.andThe` plural
    ↪ calculation, S "stop"]
```

Source Code D.4: [“MultiDefinitions” \(MultiDefn\) Definition](#)

```
-- | 'MultiDefn's are QDefinition factories, used for showing one
  ↪ or more ways
-- we can define a QDefinition.
data MultiDefn e = MultiDefn{
  -- | UID
  _rUId :: UID,
  -- | Underlying quantity it defines.
  _qd :: QuantityDict,
  -- | Explanation of the different ways we can define a quantity.
  _rDesc :: Sentence,
  -- | All possible ways we can define the related quantity.
  _rvs :: NE.NonEmpty (DefiningExpr e)
}
```

Source Code D.5: Pseudocode: Broken QuantityDict Chunk Retriever

```
retrieveQD :: UID -> ChunkDB -> Maybe QuantityDict
retrieveQD u cdb = do
  (Chunk expectedQd) <- lookup u cdb
  pure expectedQd
```
