

Taming the evolving black box:
A new architectural model for intelligent services

Alex Cummaudo

BSc *Swinburne*, BIT(Hons)

<ca@deakin.edu.au>

A thesis submitted in partial fulfilment of the requirements for the

Doctor of Philosophy



Applied Artificial Intelligence Institute

Deakin University

Melbourne, Australia

December 20, 2019

Abstract

Application developers are eager to integrate Machine Learning (ML) into their software, with a plethora of vendors providing pre-packaged components—typically under the ‘AI’ banner—to entice them. Such components are marketed as developer ‘friendly’ ML and easy for them to integrate (being ‘just another’ component added to their toolchain). These components are, however, non-trivial: in particular, developers unknowingly add the risk of mixing nondeterministic ML behaviour into their applications that, in turn, impact the quality of their software. Prior research advocates that a developer’s conceptual understanding is critical to effective interpretation of reusable components. However, these ready-made AI components do not present sufficient detail to allow developers to acquire this conceptual understanding. In this study, by use of a mixed-methods approach of survey and action research, we investigate if the application developers’ deterministic approach to software development clashes with the mindset needed to incorporate probabilistic components. Our goal

is to develop a framework to better document such AI components that improves both the quality of the software produced and the developer productivity behind it.

Declarations

I certify that the thesis entitled "*Taming the evolving black box: A new architectural model for intelligent services*" submitted for the degree of Doctor of Philosophy complies with all statements below.

- (i) I am the creator of all or part of the whole work(s) (including content and layout) and that where reference is made to the work of others, due acknowledgement is given.
- (ii) The work(s) are not in any way a violation or infringement of any copyright, trademark, patent, or other rights whatsoever of any person.
- (iii) That if the work(s) have been commissioned, sponsored or supported by any organisation, I have fulfilled all of the obligations required by such contract or agreement.
- (iv) That any material in the thesis which has been accepted for a degree or diploma by any university or institution is identified in the text.

(v) All research integrity requirements have been complied with.

Alex Cummaudo

BSc *Swinburne*, BIT(Hons)

December 20, 2019

Dedicated to all my teachers.

Acknowledgements

Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Contents

Abstract	iii
Declaration	v
Acknowledgements	ix
Contents	ix
List of Publications	xvii
List of Abbreviations	xvii
List of Figures	xxi
List of Tables	xxiv
List of Listings	xxvi
	xi

I Preface	1
1 Introduction	3
1.1 Motivation: Current Developer Mindsets'	8
1.1.1 The Impact on Software Quality	10
1.1.2 Motivating Scenarios	11
1.2 Research Outcomes	17
1.3 Concluding Remarks	27
1.4 Contributions	27
1.5 Structure	27
2 Background	33
2.1 Software Quality	35
2.1.1 Validation and Verification	37
2.1.2 Quality Attributes and Models	41
2.1.3 Reliability in Computer Vision	44
2.2 Probabilistic and Nondeterministic Systems	46
2.2.1 Interpreting the Uninterpretable	46
2.2.2 Explanation and Communication	49
2.2.3 Mechanics of Model Interpretation	51
2.3 Application Programming Interfaces	52
2.3.1 API Usability	53
3 Research Methodology	57
3.1 Research Questions Revisited	58

3.1.1	Knowledge Questions	59
3.1.2	Design Questions	60
3.2	Philosophical Stances	60
3.3	Research Design	62
3.3.1	Review of Relevant Research Methods	63
3.3.2	Review of Data Collection Techniques for Field Studies . . .	66
3.4	Proposed Experiments	67
3.4.1	Experiment I: Develop Initial Framework	67
3.4.2	Developing the Initial Framework	71
3.4.3	Experiment II: Validate Initial Framework	72
3.5	Empirical Validity	74
3.5.1	Threats to Internal Validity	75
3.5.2	Threats to External Validity	77
3.5.3	Threats to Construct Validity	79
II	Publications	81
4	Identifying Evolution in Computer Vision Services	83
4.1	Introduction	84
4.2	Motivating Example	87
4.3	Related Work	89
4.3.1	External Quality	89
4.3.2	Internal Quality	92

4.4	Method	94
4.5	Findings	98
4.5.1	Consistency of top labels	98
4.5.2	Consistency of confidence	102
4.5.3	Evolution risk	105
4.6	Recommendations	106
4.6.1	Recommendations for intelligent service users	106
4.6.2	Recommendations for intelligent service providers	108
4.7	Threats to Validity	111
4.7.1	Internal Validity	111
4.7.2	External Validity	112
4.7.3	Construct Validity	112
4.8	Conclusions & Future Work	113
5	Systematic Mapping Study of API Documentation Knowledge	115
5.1	Introduction	116
5.2	Related Work	118
5.3	Method	119
5.3.1	Systematic Mapping Study	120
5.3.2	Development of the Taxonomy	126
5.4	Taxonomy	128
5.5	Threats to Validity	130
5.6	Conclusions & Future Work	131

III Postface	139
6 Conclusion	141
References	165
IV Appendices	171
A Additional Materials	173
A.1 The Development, Documentation and Usage of Web APIs	175
B Authorship Statements	185
C Ethics Clearance	199
D Primary Sources from Systematic Mapping Study	203

List of Publications

1. A. Cummaudo, R. Vasa, J. Grundy, M. Abdelrazek, and A. Cain, “Losing Confidence in Quality: Unspoken Evolution of Computer Vision Services,” in *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. Cleveland, OH, USA: IEEE, Sep. 2019, pp. 333–342
2. A. Cummaudo, R. Vasa, and J. Grundy, “What should I document? A preliminary systematic mapping study into API documentation knowledge,” in *13th International Symposium on Empirical Software Engineering and Measurement (ESEM)*. Porto de Galinhas, Recife, Brazil: IEEE, Sep. 2019, pp. 1–6
3. T. Ohtake, A. Cummaudo, M. Abdelrazek, R. Vasa, and J. Grundy, “Merging Intelligent API Responses Using a Proportional Representation Approach,” in *Empirical Software Engineering and Verification*. Cham: Springer International Publishing, Apr. 2019, pp. 391–406

List of Abbreviations

A²I² Applied Artificial Intelligence Institute. 66, 71, 72, 77, 78

AI Artificial Intelligence. iii, iv, 3–6, 46–50, 181

API Application Programming Interface. 5–21, 25, 27, 29–31, 33, 38–40, 52–55, 58, 59, 62, 65, 67, 69, 71–73, 79, 175, 179

AWS Amazon Web Services. 4

BYOML Build Your Own Machine Learning. 6, 29

CDSS Clinical Decision Support System. 11, 15

CIS Cloud Intelligence Service. 7–11, 13, 16–21, 26, 30, 33–40, 43, 44, 46, 52, 54, 55, 58–60, 62, 64, 65, 67, 69, 72, 73, 79

CNN Convolutional Neural Network. 15, 16, 46

CRUD Create, read, update, and delete. 179

cvCIS Computer Vision Cloud Intelligence Service. 8, 9, 11–13, 18, 20, 27, 31, 33, 36, 39, 44, 53, 66, 68, 69, 73

DCE Distributed Computing Environment. 175

HITL Human-in-the-loop. 15, 16

HTTP Hypertext Transfer Protocol. 29, 175, 177–180

IDE Integrated Development Environment. 4

IDL Interface Definition Language. 175, 176, 179

JSON JavaScript Object Notation. 7

ML Machine Learning. iii, 4, 6, 7, 14, 29, 40, 48

NN Neural Network. 45, 47, 48, 51

QoS Quality of Service. 176

RAML RESTful API Modeling Language. 179

REST REpresentational State Transfer. 7, 176, 178–180

ROI Region of Interest. 15, 16

RPC Remote Procedure Call. 175

SE Software Engineer. 50

SOA Service-Level Agreement. 176

SOA Service-Oriented Architecture. 175, 176

SOAP Simple Object Access Protocol. 7, 175–179

SQuaRE Systems and software Quality Requirements and Evaluation. 43

SVM Support Vector Machine. 47, 51

URI Uniform Resource Identifier. 178

V&V Verification & Validation. 34, 35, 37–41

WADL Web Application Description Language. 179

WS Web Service. 175

WSDL Web Services Description Language. 176

XML eXtendable Markup Language. 7, 175, 177

List of Figures

1.1	Differences between data- and rule-driven cloud services	5
1.2	The spectrum of machine learning	6
1.3	Overview of cloud intelligence services	8
1.4	CancerAssist Context Diagram	16
2.1	Mindset clashes within the development, use and nature of a CIS . .	35
2.2	Leakage of internal and external quality in CISs	40
2.3	Overview of software quality models	42
2.4	Adversarial examples in computer vision	45
2.5	Deterministic versus nondeterministic systems	47
2.6	Theory of AI communication	51
3.1	High-level overview of the proposed experiments	68

4.1	Consistency of labels in CV services is rare	98
4.2	Top labels for images between CV services do not intersect	100
4.3	CV services can return multiple top labels	101
4.4	Cumulative distribution of top label confidences	104
4.5	Cumulative distribution of intersecting top label confidences	104
4.6	Agreement of labels between multiple CV services do not share similar confidences	105
5.1	A systematic map of API documentation knowledge studies	126
A.1	SOAP versus REST search interest over time	176
A.2	Categorisation of AI-based products and services	181
A.3	Increasing interest in the developer community of computer vision APIs	182
A.4	Review of field study techniques	183

List of Tables

1.1	Differing characteristics of cloud services	28
1.2	Comparison of the machine learning spectrum	29
1.3	Varying confidence changes over time between 3 CV APIs	30
1.4	Definitions of ‘confidence’ in CV documentation	31
4.1	Characteristics of data in CV evolution assessment	96
4.2	Ratio of consistent labels in CV services	99
4.3	Evolution of top labels and confidence values	103
5.1	Summary of search results in API documentation knowledge	122
5.2	Data extraction in API documentation knowledge study	125
5.3	Taxonomy proposed in API documentation knowledge study	133
5.3	An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (<i>Continued</i>).	134

5.3 An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (<i>Continued</i>).	135
5.3 An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (<i>Continued</i>).	136
5.3 An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (<i>Continued</i>).	137

List of Listings

A.1	An example SOAP request	177
A.2	An example SOAP response	177
A.3	An example RESTful request	179
A.4	An example RESTful response	180

Part I

Preface

CHAPTER 1

3

4

5

Introduction

6

7 Within the last half-decade, we have seen an explosion of cloud-based services
8 typically marketed under an Artificial Intelligence (AI) banner. Vendors are rapidly
9 pushing out AI-based solutions, technologies and products that encapsulate half
10 a century worth of machine-learning research: a 2016 report by market research
11 company Forrester captured such growth into four key areas [143] as replicated in
12 Figure A.2. Application developers are eager to develop the next generation of
13 ‘AI-first’ software, that will reason, sense, think, act, listen, speak and execute every
14 whim in our web browser or smartphone app.

15 A contrast is how these systems shift away from the traditional software engi-
16 neering paradigm that application developers are comfortable with. Typical systems
17 built by application developers are *rule-driven* (or algorithm-driven), deterministi-

¹⁸ cally human engineered using source code to drive each step behind the application.
¹⁹ These rule-driven systems typically consume, utilise, and integrate libraries and
²⁰ frameworks, IDEs and other tooling, and cloud-based services such as Amazon Web
²¹ Services (AWS) [268]. The AI-first software is, however, not rule-driven but *data-*
²² *driven*, fuelled with large datasets used to train Machine Learning (ML) prediction
²³ algorithms and classifiers that present a typically probabilistic and nondeterministic
²⁴ behaviour.

²⁵ So, how does the application developer approach her AI-first application? To
²⁶ answer this, let us consider three approaches that could be used:

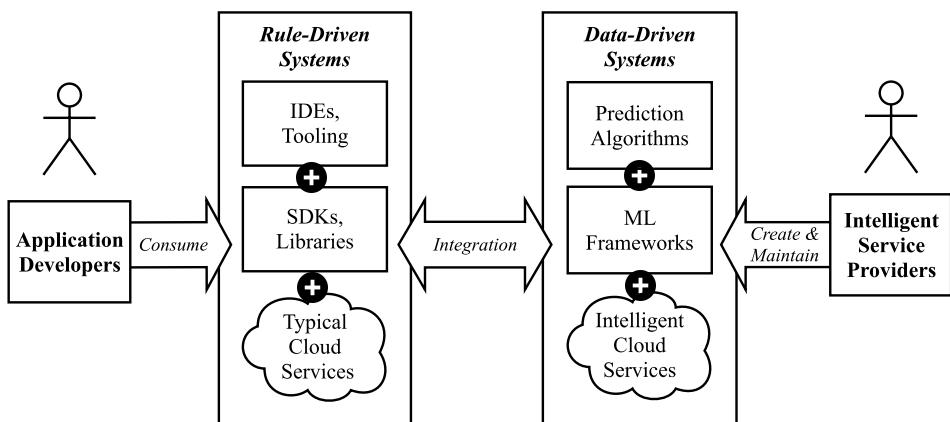
- ²⁷ 1. She writes her own ML classifier from scratch and trains it from her own
²⁸ curated dataset. This approach is laborious in time and demands formal
²⁹ training in ML and mathematical knowledge, but the tradeoff is that she has
³⁰ full autonomy over the models she creates.
- ³¹ 2. She downloads a pre-trained model and ‘plugs’ it into an existing ML frame-
³² work, such as Tensorflow [4]. While this approach is less demanding in time,
³³ it still requires her to revise and understand the ‘glue’ of the ML framework
³⁴ code¹ with her own application code.
- ³⁵ 3. She uploads her data to a pre-existing cloud-based service: a *Friendly Machine*
³⁶ *Learning* service. She doesn’t need to know anything behind the underlying
³⁷ ‘intelligence’ and how it works, is fast to integrate into her application and

¹Thus introducing a verbose list of ML terminology to her already-required developer vocabulary. See a list of 328 terms provided by Google here: <https://developers.google.com/machine-learning/glossary/>.

³⁸ all abstracted behind a web-based Application Programming Interface (API)
³⁹ call.

⁴⁰ At first sight, she perceives the data-driven cloud service as ‘just another’ cloud
⁴¹ service offered in her toolchain. Her perception is that just because this is another
⁴² cloud service, it should act and behave as any other typical service would. But
⁴³ internally, she isn’t aware that this particular cloud service isn’t a typical one;
⁴⁴ the data-driven service does not adapt to her rule-driven application in the way that
⁴⁵ she thinks (Figure 1.1). This is because the nature of typical cloud systems and
⁴⁶ data-driven ones are not exactly the same (Table 1.1).

Figure 1.1: The application developer’s rule-driven toolchain is distinct from data-driven toolchain. A developer must consume a typical, data-driven cloud service in a different way than an intelligent data-driven cloud service as they are not the same type of system.



⁴⁷ This ‘range’ of AI-first integration techniques partially reflects Google AI’s²

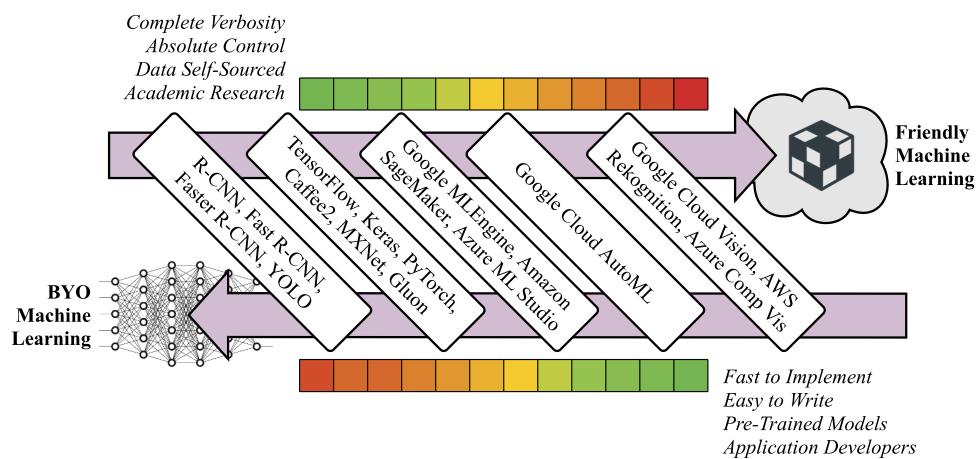
²Google AI was recently rebranded from Google Research, further highlighting how the ‘AI-first’ philosophy is increasingly becoming embedded in companies’ product lines and research and development teams. Spearheaded through work achieved at Google, Microsoft and Facebook, the emphasis can be seen through Google’s 2018 rebranding of *Google Research* to *Google AI* [263] or

⁴⁸ *Machine Learning Spectrum* [3, 130, 173], which encompasses the variety of skill, effort, users and types of outputs of integration techniques. On one extreme is the research of developing algorithms to achieve intelligence, produced chiefly in academia—coined as Build Your Own Machine Learning (BYOML) [2, 3, 173].

⁵² On the other, such intelligence becomes heavily abstracted as easy-to-use APIs, targeted mainly towards developers as ‘friendly’ ML. In the middle lies a broad mix of combining both cloud and locally-hosted solutions (with varying levels of automation to assist in development) that turn custom datasets into some form of predictive intelligence.

⁵⁷ All techniques in this spectrum are data-driven, and we illustrate their slightly varied characteristics further in Table 1.2 and examples of the computer vision spectrum in Figure 1.2.

Figure 1.2: Examples within the machine learning spectrum of computer vision. Benefits and drawbacks of each end of the spectrum are indicated with the colour scales.



how AI is leveraged *at scale* within Facebook’s infrastructure and platforms to serve its users with an AI-first attitude [177].

60 *In this study, we advocate that the (i) integration, (ii) documentation, and*
61 *(iii) quality attributes of data-driven cloud services juxtaposes the rule-driven nature*
62 *of end-applications as these intelligent cloud services are vastly different to their*
63 *traditional counterparts, and great care must therefore be considered by application*
64 *developers.*

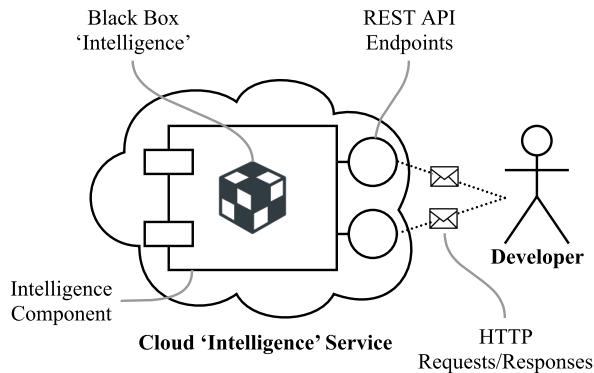
65 These cloud services have begun to gain traction within developer circles: Figure A.3 shows the increasing trend of posts since 2014 on Stack Overflow that
66 categorise popular computer vision cloud APIs.³ In academia, these ‘off-the-shelf’
67 and pre-packaged ML solutions present a varied nomenclature such as *Cognitive*
68 *Applications and Machine Learning Services* [101] or *Machine Learning as a Ser-*
69 *vice* [194]. Some services provide the infrastructure to rapidly begin training from
70 custom datasets (Google’s AutoML⁴ is one such example) while others provide pre-
71 trained datasets ‘ready-for-use’ in production without the need to train data. We refer
72 to these latter services under the broader term *Cloud Intelligence Service (CIS)*, and
73 diagrammatically express their usage within Figure 1.3.

75 The general workflow of a CIS is simple: a developer accesses a CIS component
76 via REST/SOAP API(s). For their given input, they receive an intelligent-like
77 response typically serialised as JSON/XML. We note the intelligence component
78 masks its ‘intelligence’ through a black-box: in recent years, there is a rise in
79 providing human-level intelligence via crowdsourcing Internet marketplaces such as

³Query run on 12 October 2018 using StackExchange Data Explorer. Refer to <https://data.stackexchange.com/stackoverflow/query/910188> for full query.

⁴<https://cloud.google.com/automl/> last accessed 7 December 2018.

Figure 1.3: Overview of Cloud Intelligence Services.



⁸⁰ Amazon Mechanical Turk [264] or ScaleAPI [265]. Thus, a CIS may be powered by
⁸¹ varying degrees of intelligence: human intelligence, machine learning, data mining
⁸² or even intelligence by brute-force.

⁸³ While there are different types of CISs evident (such as OCR transcription,
⁸⁴ text-to-speech and speech-to-text, object categorisation, object comparison, natural
⁸⁵ language processing etc.), we scope the work investigated in this study to Computer
⁸⁶ Vision Cloud Intelligence Services (cvCISs) [266–278]. The ubiquity of cvCISs
⁸⁷ is exemplified through evermore growing applications that use these APIs: aiding
⁸⁸ the vision-impaired [51, 192], accounting [146], data analytics [107], and student
⁸⁹ education [57]. Moreover, we refer to its growing adoption in developer circles
⁹⁰ within Figure A.3.

⁹¹ 1.1 Motivation: Current Developer Mindsets'

⁹² Figure A.3 shows an increasing trend to the adoption and discussion of CISs with de-
⁹³ velopers. These services are accessible through APIs and consist of an ‘intelligence’

⁹⁴ black box (Figure 1.3). When a term ‘black box’ is used, the input (or stimulus) is
⁹⁵ transformed to its outputs (or response) without any understanding of the internal
⁹⁶ architecture by which this transformation occurs, a theory arising from the electronic
⁹⁷ sciences and adapted to wider applications since the 1950s–60s [10, 33] to describe
⁹⁸ “systems whose internal mechanisms are not fully open to inspection” [10].

⁹⁹ In the world of machine learning and data mining, where we develop algorithms
¹⁰⁰ to make predictions in our datasets or discover patterns within them, these black
¹⁰¹ boxes are inherently probabilistic and stochastic; there is little room for certainty in
¹⁰² these results as such insight is purely statistical and associational [181] against its
¹⁰³ training dataset. As an example, a cvCIS returns the *probability* that a particular
¹⁰⁴ object (the response) exists in the raw pixels (the stimulus), and thus for a more
¹⁰⁵ certain (though not fully certain) distribution of overall confidence returned from
¹⁰⁶ the service, a developer must treat the problem stochastically by testing this case
¹⁰⁷ hundreds if not thousands of times to find a richer interpretation of the inference
¹⁰⁸ made. Developers (at present) do not need to treat their programs in any such
¹⁰⁹ stochastic way given their rule-driven mindset that computers will always make
¹¹⁰ certain outcomes.

¹¹¹ There are thus therefore three key factors to consider when implementing, testing
¹¹² and developing with a CIS: (i) the API usability, (ii) the nature of nondeterministic
¹¹³ and probabilistic systems, and (iii) how both impact on software quality.

1.1.1 The Impact on Software Quality

114 Do traditional techniques for documenting deterministic APIs also apply to non-deterministic systems? As APIs reflect a set of design choices made by their providers intended for use by the developer, does the mindset between the machine learning architect and the novice programmer match? Evaluations of API usability advocate for the accuracy, consistency and completeness of APIs and their documentation [186, 200] written by providers, while providers should consider mismatches between the developer's conceptual knowledge of the API its implementation [124]. However, consistency cannot be guaranteed in probabilistic systems, and the conceptual knowledge of such systems are still treated like black boxes. It is therefore imperative that CIS providers consider the impact of their API usability; if not, poor API usability hinders on the internal quality of development practices, slowing developers down to produce the software they need to create.

127 Moreover, CIS APIs are inherently non-deterministic in nature, but developers are still taught with the deterministic mindset that all API calls are the same. Simple arithmetic representations (e.g., $2 + 2 = 4$) will *always* result in 4; but a multi-layer perceptron neural network performing similar arithmetic representation [22] gives the probability where the target output (*exactly* 4) and the output inferred (*possibly* 4) matches as a percentage (or as an error where it does not match). That is, instead of an exact output, there is instead a *probabilistic* result: $2 + 2$ *may* equal 4 with a confidence of n . External quality must therefore be considered in the outcome of these systems, such as in the case of thresholding values, to consider whether or not the inference has a high enough confidence to justify its result to end-users.

¹³⁷ In order to fully understand this problem, there are multiple dimensions one must
¹³⁸ consider: the impact of software quality; the fact that these systems underneath are
¹³⁹ probabilistic and are stochastic; the cognitive biases of determinism in developers;
¹⁴⁰ the issue of consistency in API usage. While existing literature does extensively
¹⁴¹ explore software quality and API usability, these studies have only had emphasis on
¹⁴² deterministic systems and thus little work to date has investigated such factors on
¹⁴³ probabilistic systems that make up the core of cvCISs. We explore more of these
¹⁴⁴ facets in the motivating scenarios below.

¹⁴⁵ 1.1.2 Motivating Scenarios

¹⁴⁶ The market for intelligent services is increasing (Figure A.2) and as is developer
¹⁴⁷ uptake and enthusiasm in the software engineering community (Figure A.3). We
¹⁴⁸ investigate the impact of the mismatch between the developer's mindset and the
¹⁴⁹ service provider's mindset as little work has been presented in literature. How
¹⁵⁰ do developers work with a CIS, how usable are these cloud APIs, and how well
¹⁵¹ do developers understand the non-deterministic and stochastic nature of a services
¹⁵² backed by machine-learnt models?

¹⁵³ To illustrate the context of use, we present the two scenarios of varying risk: (i)
¹⁵⁴ a fictional software developer named Tom who wishes to develop an inherently low-
¹⁵⁵ risk photo detection application for his friends and family; and (ii) a high-risk cancer
¹⁵⁶ Clinical Decision Support System (CDSS) that uses patient scans to recommend to
¹⁵⁷ surgeons if the patient should be sent to surgery.

158 Motivating Scenario I: Tom's *PhotoSharer* App

159 Tom wants to develop a social media photo-sharing app on iOS and Android, *Photo-*
160 *toSharer*, that analyses photos taken on smartphones as they are taken. Tom wants
161 the app to categorise photos into scenes (e.g., day vs. night, landscape vs. indoors),
162 generate brief descriptions of each photo, and catalogue photos of his friends as well
163 as common objects (e.g., all photos with his Border Collie dog, all photos taken on
164 a beach on a sunny day). His app will then share all of this analysed intelligence of
165 his photos with his friends on a social-media-like platform, where his friends can
166 search and view the photos.

167 Rather than building a computer vision engine from scratch, which would take
168 far too much time and effort, Tom thinks he can achieve this using one of the
169 common cvCISs. Tom comes from a typical software engineering background and
170 has insufficient knowledge of key computer vision terminology and no understanding
171 of its underlying techniques. However, inspired by easily accessible cloud APIs that
172 offer computer vision analysis, he chooses to use these. Built upon his experience of
173 using other similar cloud services, he decides on one of the cvCIS APIs, and expects
174 a static result always and consistency between similar APIs. Analogously, when Tom
175 invokes the iOS Swift substring method "doggy".prefix(3), he rightfully expects
176 it to be consistent with the Android Java equivalent "doggy".substring(0, 2).
177 Consistent, here, means two things: (i) that 'dog' will *always* be returned every
178 time he invokes the method in either language (i.e., a static response); and (ii) that
179 'dog' will *always* be returned regardless of what programming language or string
180 library is used, given the deterministic nature of the 'substring' construct (i.e., results

181 for substring are API-agnostic).

182 More concretely, in Table 1.3, we illustrate how three (anonymised) cvCIS
183 providers fail to provide similar consistency to that of the substring example above.

184 If Tom uploads a photo of a border collie⁵ to three different providers in August
185 2018 and January 2019, he would find that each provider is different in both the

186 vocabulary used between. The confidence values and labels within the *same* provider
187 also vary within a matter of five months. The evolution of the confidence changes

188 is not explicitly documented by the providers (i.e., when the models change) nor
189 do they document what confidence even means. Their current tautological nature

190 of the definition of these changing confidence values (as presented in the API
191 documentation) provides no insight for Tom to understand why there was a change

192 in confidence, which we show in Table 1.4, unless he *knows* that the underlying
193 models change with them. Thus, the deterministic problem of a substring compared

194 to the nondeterministic nature of the CIS is not so simple to comprehend unless he
195 is explicitly told.

196 To make an assessment of these APIs, he tries his best to read through the
197 documentation of some cvCIS APIs, but he has no guiding framework to help him

198 choose the right one. Some of the questions that come to mind include:

199 • What does confidence mean?

200 • Are these APIs consistent in how they respond?

201 • Will he need a combination of multiple cvCIS APIs to solve this task?

5The image used for these results can be found at <https://www.akc.org/dog-breeds/border-collie/>.

²⁰² • How does he know when there is a defect in the response? How can he report
²⁰³ it?

²⁰⁴ • How does he know what labels the API can pick up, and what labels it can't?

²⁰⁵ • How does it describe his photos and detect the faces?

²⁰⁶ • How can he interpret the results if he disagrees with it to help improve his
²⁰⁷ app?

²⁰⁸ • Does he understand that the API uses a machine learnt model? Does he know
²⁰⁹ what a ML model even is?

²¹⁰ • If so, does he know when the models update? What is the release cycle?

²¹¹ Dazzled by this, he does some brief reading on Wikipedia but is confused by
²¹² the immense technical detail to take in. He would like some form of guiding
²¹³ communication framework to assist him and in software engineering terms aligned
²¹⁴ to his background.

²¹⁵ Although Tom generally anticipates some imperfections, he has no prior bench-
²¹⁶ mark to guide him on what to expect. He understands that the app is not always
²¹⁷ going to be perfect: perhaps some photos of his dog may be missed because the
²¹⁸ dog is in the background and not the foreground, or his friends can't find the photos
²¹⁹ of their recent trip to the beach because it wasn't sunny enough for the beach to be
²²⁰ recognised. These imperfections appear to be low-risk, but may become socially
²²¹ awkward when in use; for instance, if some of Tom's friends have low self-esteem
²²² and use the app, they may be sensitive to being misidentified or even mislabelled.
²²³ Privacy issues also come into play especially if certain friends have only access to

²²⁴ certain photos that they are (supposedly) in; e.g., photos from a holiday with Tom
²²⁵ and his partner, however if the API identifies Tom's partner as a work colleague,
²²⁶ then Tom's partner's privacy is at risk.

²²⁷ **Motivating Scenario II: Cancer Detection CDSS**

²²⁸ Recent works in the oncology domain have used deep-learning Convolutional Neural
²²⁹ Networks (CNNs) to detect Region of Interests (ROIs) in image scans of tissue (e.g.,
²³⁰ [66, 91, 142]), flagging these regions for doctors to review. Trials of such algorithms
²³¹ have been able to accurately detect cancer at higher rates than humans, and thus
²³² incorporating such capabilities into a CDSS is closer within reach. Some studies have
²³³ suggested that practitioner over-reliance may erode independent decision-making
²³⁴ [41, 109]; therefore the risks in developing CDSSs powered by intelligent services
²³⁵ become paramount.

²³⁶ In Figure 1.4 we present a context diagram for a fictional CDSS named *CancerAssist*. CancerAssist is used by a team of busy pathologists who review patient
²³⁷ lymph node scans and discuss and recommend, on consensus, if the patient should or
²³⁸ should not be sent to surgery. When consensus is made, the lead pathologist enters
²³⁹ the verdict into CancerAssist—running passively in the background—to ensure no
²⁴⁰ oversight has been made in the team's discussions. When a conflict exists between the
²⁴¹ team's verdict and CancerAssist's verdict, the system produces the scan with ROIs it
²⁴² thinks the team should review. Where the team override the output of CancerAssist,
²⁴³ this helps to reinforce CancerAssist's internal model as a Human-in-the-loop (HITL)
²⁴⁴ learning process.

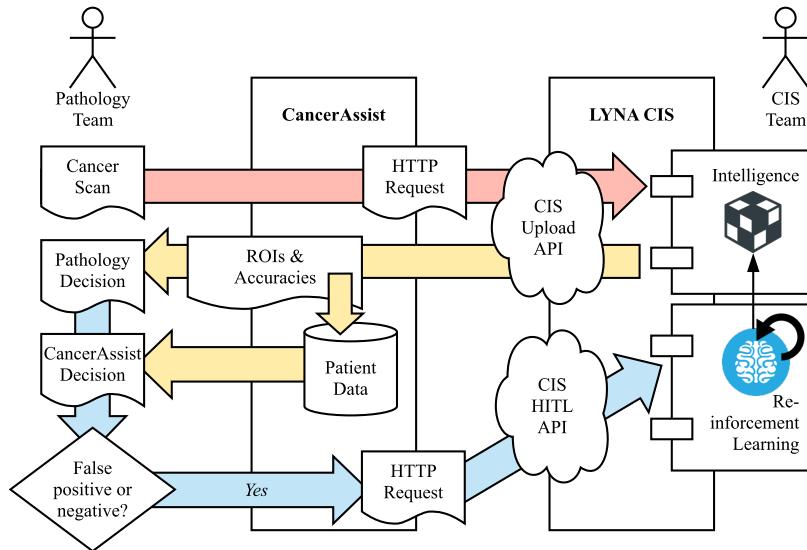


Figure 1.4: CancerAssist Context Diagram. **Key:** Red Arrows = Scan Input; Yellow Arrows = Decision Output; Blue Arrows = HITL Feedback Input.

²⁴⁶ Powering CancerAssist is Google AI's Lymph Node Assistant (LYNA) [142],
²⁴⁷ a CNN based on the Inception-v3 model [127, 232]. To provide intelligence to
²⁴⁸ CancerAssist, LYNA is hosted on a CIS, and thus the developers of CancerAssist call
²⁴⁹ the relevant CIS API endpoints, in conjunction with extra information such as patient
²⁵⁰ data and medical history, to produce the verdict. In the case of a positive verdict, the
²⁵¹ relevant ROIs CancerAssist has found are highlighted with their respective bounding
²⁵² boxes and their respective cancer detection accuracies.

²⁵³ The developer of CancerAssist has no interaction with the Data Science team
²⁵⁴ maintaining the LYNA CIS. As a result, they are unaware when updates to the
²⁵⁵ model occur, nor do they know what training data they could provide to test their
²⁵⁶ own system. The default assumptions are that the training data used to power the
²⁵⁷ intelligence is near-perfect for universal situations; i.e., the algorithm chosen is

²⁵⁸ the correct one for all the ontology tests that need to be assessed in the given use
²⁵⁹ case of CancerAssist. Thus, unlike deterministic systems—where the developer can
²⁶⁰ manually test and validate the outcomes of the APIs—this is impossible for non-
²⁶¹ deterministic systems such as CancerAssist and its underlying CIS. The ramifications
²⁶² of not being able to test such a system and putting it out into production may prove
²⁶³ fatal to patients.

²⁶⁴ Certain questions in the production of CancerAssist and its use of a CIS may
²⁶⁵ come into mind:

- ²⁶⁶ • When is the model updated and how do the CIS team communicate that the
²⁶⁷ model is updated?
- ²⁶⁸ • What benchmark test set of data do I use to ensure that the changed model
²⁶⁹ doesn't affect other results?
- ²⁷⁰ • How do we know that the assumptions made by the CIS team who train model
²⁷¹ are correct?

²⁷² Thus, improved documentation and additional metadata may be needed to better
²⁷³ improve communication between developers and CIS providers. Such claims are
²⁷⁴ further expanded upon in the following section.

²⁷⁵ 1.2 Research Outcomes

In this thesis, we present a framework to improve the documentation quality

of Computer Vision Cloud Intelligence Services (cvCISs) and their APIs. We demonstrate that developers currently lack the understanding of how these services function and our framework mitigates this by presenting a solution to improve the documentation quality of the APIs and improve the existing techniques used to integrate these services into developer's end-applications.

²⁷⁶ The goals of this study aim to provide a snapshot of current developer practices
²⁷⁷ towards the usage of CISs. This allows us to develop a guiding framework and
²⁷⁸ recommendations for application developers and CISs providers alike. Our anchor-
²⁷⁹ ing perspective is software quality—specifically, validation and verification—with
²⁸⁰ such systems and what best practices within the field of software engineering can
²⁸¹ be applied to assist in consumption of CISs. Based on the motivating case studies
²⁸² in Section 1.1, we articulate three Research Hypotheses (RH1–3) below and seven
²⁸³ Research Questions (RQs) based on both empirical and non-empirical software
²⁸⁴ engineering methodology [217, 218].

RH1: *Existing CISs present insufficient API documentation for general use.*

Research Hypothesis API documentation of intelligent services are inade-

RH1: Existing CISs present insufficient API documentation for general use. (cont)

quate and insufficient given the disparity of mindsets between the application

RH1: Existing CISs present insufficient API documentation for general use. (cont)

developers and CIS providers. Chiefly, application developers have limited general understanding of the ‘magic’ that occurs behind these probabilistic ‘intelligent’ APIs. We do not know what key aspects of the documentation matter to them, nor what they do or do not understand of the existing documentation.

Research Goal To improve the effectiveness of the documentation in existing CIS providers, specifically of cvCIS APIs.

Research Questions

RQ1.1. What practices are in use for intelligent services’ API documentation?

RQ1.2. How do developers currently understand and interpret the documentation given a lack of formal training in artificial intelligence? That is, what do they understand and not understand, and what key aspects of the API documentation matter do developers as they see it?

RQ1.3. What additional information or attributes would developers prefer to be included in the API documentation?

Research Contribution An intelligent service API documentation quality assessment framework to evaluate how well the service has been documented for software engineers to use.

RH2: *Existing CISs present insufficient metadata for context-specificity.*

Research Hypothesis Intelligent service APIs respond with insufficient information for developers to operationalise the service into a business-driven application and, thus, additional metadata is needed to assist developers. Such metadata is likely to be needed as part of the response objects of the API.

Research Goal To improve the quality of *context-specific response data* from the API endpoints of intelligent services.

Research Questions

RQ2.1. What are current problems due to lack of return metadata?

RQ2.2. What additional metadata do developers desire to achieve implementing context-specific applications?

Research Contributions A list of metadata key-value-pairs that assist developers in using these APIs during the development of software that consume these services. In essence, improvements to the framework of Research Outcome 1: “*An intelligent service API documentation and metadata quality assessment framework*”.

RH3: *RH1 and RH2 improve quality, productivity or developer informativeness.*

Research Hypothesis The implication of hypotheses 1 and 2 suggest that

RH3: RH1 and RH2 improve quality, productivity or developer informativeness. (cont)

improving both the documentation and providing further metadata will improve



RH3: RH1 and RH2 improve quality, productivity or developer informativeness. (cont)

product quality (internal or external), and/or developer productivity and/or de-



RH3: RH1 and RH2 improve quality, productivity or developer informativeness. (cont)

veloper education in developing software with intelligent components.



RH3: RH1 and RH2 improve quality, productivity or developer informativeness. (cont)

Research Goal To confirm if improvements to API documentation and re-



RH3: RH1 and RH2 improve quality, productivity or developer informativeness. (cont)

sponse metadata are reflected as improvements to product quality, developer productivity and/or developer education.

Research Questions

RQ3.1. Does an improvement of documentation or metadata correlate to an improvement in software quality, developer productivity and/or developer informativeness?

RQ3.2. With respect to RQ3.1, the three aspects are explored:

- (a) Does the improvement cause increased product quality, as measured through improved external quality metrics?
- (b) Does the improvement cause increased developer productivity, as measured through improved internal quality metrics?
- (c) Does the improvement cause increased developer informativeness or increased confidence in developing CIS-powered applications?

Research Contribution A concrete sample solution or framework that improves such metrics, thereby confirming that our documentation and metadata quality assessment framework improves these facets.

1.3 Concluding Remarks

285 Ultimately, we seek to understand the conceptual understanding of software en-
286 gineers who operationalise stochastic and probabilistic systems, and furthermore
287 understand knowledge representation with these systems' API documentation. Our
288 motivation is to provide insight into current practices and compare the best practices
289 with actual practise. We strive for this to provide developers with a guiding frame-
290 work on how to best operationalise these systems via the form of some checklist or
291 tool they can use to ensure optimal software quality.

292 It is anticipated that the findings from this study in the cvCISs space will be gener-
293 alisable to other areas, such as time-series information, natural language processing
294 and others.

1.4 Contributions**1.5 Structure**

Table 1.1: Differing characteristics of intelligent and typical cloud services.

Intelligent Cloud Services	Typical Cloud Services
Probabilistic	Deterministic
Machine Learnt	Human Engineered
Data-Driven	Rule-Driven
Black-Box	Mostly Transparent

Table 1.2: Comparison of the machine learning spectrum.

Comparator	BYOML	ML F'work	Cloud ML	Auto-Cloud ML	Cloud API
Hosting					
Locally	✓	✓			
Cloud			✓	✓	✓
Output					
Custom Model	✓	✓	✓	✓	
HTTP Response					✓
Autonomy					
Low					✓
Medium					✓
High		✓	✓		
Highest	✓				
Time To Market					
Medium	✓	✓			
High			✓	✓	
Highest					✓
Data					
Self-Sourced	✓	✓	✓	✓	
Pre-Trained		✓			✓
Intended User					
Academics	✓	✓			
Data Scientist	✓	✓	✓	✓	
Developers				✓	✓

Table 1.3: First six responses of image analysis for a Border Collie sent to three computer vision CIS APIs providers five months apart. The specificity (to 3 s.f.) and vocabulary of each label in the response varies between all services, and—except for Provider B—changes over time. Any confidence changes greater than 1 per cent are highlighted in red.

Label	Provider A		Provider B		Provider C	
	Aug 2018	Jan 2019	Aug 2018	Jan 2019	Aug 2018	Jan 2019
Dog	0.990	0.986	0.999	0.999	0.992	0.970
Dog Like Mammal	0.960	0.962	-	-	-	-
Dog Breed	0.940	0.943	-	-	-	-
Border Collie	0.850	0.852	-	-	-	-
Dog Breed Group	0.810	0.811	-	-	-	-
Carnivoran	0.810	0.680	-	-	-	-
Black	-	-	0.992	0.992	-	-
Indoor	-	-	0.965	0.965	-	-
Standing	-	-	0.792	0.792	-	-
Mammal	-	-	0.929	0.929	0.992	0.970
Animal	-	-	0.932	0.932	0.992	0.970
Canine	-	-	-	-	0.992	0.970
Collie	-	-	-	-	0.992	0.970
Pet	-	-	-	-	0.992	0.970

Table 1.4: Tautological definitions of ‘confidence’ found in the API documentation of three common cvCIS providers.

API Provider	Definition(s) of Confidence
Provider A	<p><i>“Score is the confidence score, which ranges from 0 (no confidence) to 1 (very high confidence).” [279]</i></p>
	<p><i>“Deprecated. Use score instead. The accuracy of the entity detection in an image. For example, for an image in which the ‘Eiffel Tower’ entity is detected, this field represents the confidence that there is a tower in the query image. Range [0, 1].” [280]</i></p>
	<p><i>“The overall score of the result. Range [0, 1]” [280]</i></p>
Provider B	<p><i>“Confidence score, between 0 and 1... if there is insufficient confidence in the ability to produce a caption, the tags maybe [sic] the only information available to the caller.” [281]</i></p>
	<p><i>“The level of confidence the service has in the caption.” [282]</i></p>
Provider C	<p><i>“The response shows that the operation detected five labels (that is, beacon, building, lighthouse, rock, and sea). Each label has an associated level of confidence. For example, the detection algorithm is 98.4629% confident that the image contains a building.” [283]</i></p>
	<p><i>“[Provider C] also provides[s] a percentage score for how much confidence [Provider C] has in the accuracy of each detected label.” [284]</i></p>

CHAPTER 2

298

299

300

Background

301

302 In Chapter 1, we defined a common set of (artificial) intelligence-based cloud ser-
303 vices that we label Cloud Intelligence Services (CISs). Specifically, we scope the
304 primary body of this study’s work on Computer Vision Cloud Intelligence Services
305 (cvCISs) (e.g., Google Cloud Vision [266], AWS Rekognition [268], Azure Com-
306 puter Vision [267], Watson Visual Recognition [270] etc.). We claim developers
307 have a distinctly deterministic mindset ($2 + 2$ *always* equals 4) whereas a CIS’s
308 ‘intelligence’ component (a black box) may return probabilistic results ($2 + 2$ *might*
309 equal 4 *with a confidence of 95%*). Thus, there is a mindset mismatch between
310 probabilistic results (from the API provider) and results interpreted with certainty
311 (from the API consumer).

312 What affect does this mindset mismatch have on the developer’s approach to-

³¹³ wards building probabilistic software? What can we learn from common software
³¹⁴ engineering practices (e.g., [188, 223]) that apply to resolve this mismatch and
³¹⁵ thereby improve quality, such as Verification & Validation (V&V)? Chiefly, we an-
³¹⁶ chor this question around three lenses of software engineering: creating a CIS, using
³¹⁷ a CIS, and the nature of CISs themselves.

³¹⁸ Our chief concern lies with interaction and integration between CIS providers
³¹⁹ and consumers, the nature of applications built using a CIS, and the impact this has on
³²⁰ software quality. We triangulate this around three pillars, which we diagrammatically
³²¹ represent in Figure 2.1.

³²² **(1) The development of the CIS.** We investigate the internal quality attributes
³²³ of creating a CIS from the CIS *provider's* perspective. That is, we ask if
³²⁴ existing verification techniques are sufficient enough to ensure that the CIS
³²⁵ being developed actually satisfies the CIS consumer's needs and if the internal
³²⁶ perspective of creating the system with a non-deterministic mindset clashes
³²⁷ with the outside perspective (i.e., pillar 2).

³²⁸ **(2) The usage of the CIS.** We investigate the external quality attributes of using a
³²⁹ CIS from the CIS *consumer's* perspective. That is, we ask if existing validation
³³⁰ techniques are sufficient enough to ensure that the end-users can actually use
³³¹ a CIS to build their software in the ways they expect the CIS to work.

³³² **(3) The nature of a CIS.** We investigate what standard software engineering
³³³ practices apply when developing non-deterministic systems. That is, we
³³⁴ tackle what best practices exist when developing systems that are inherently
³³⁵ stochastic and probabilistic, i.e., the ‘black box’ intelligence itself.

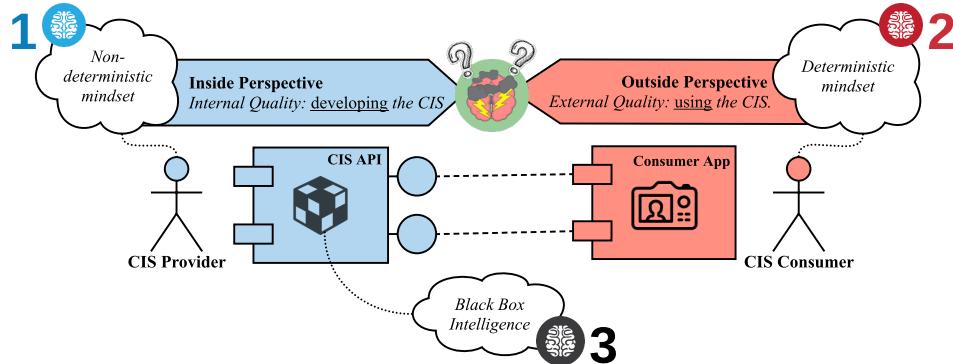


Figure 2.1: The three pillars by which we anchor the background: (1) developing a CIS with a non-deterministic mindset by the CIS provider; (2) the use of a CIS with a deterministic mindset by the CIS consumer; (3) the nature of a CIS itself.

³³⁶ Does a clash of deterministic consumer mindsets who use a CIS and the non-
³³⁷ deterministic provider mindsets who develop them exist? And what impact does this
³³⁸ have on the inside and outside perspective? Throughout this chapter, we will review
³³⁹ these core principles due to such mindset mismatch from the anchoring perspective
³⁴⁰ of software quality, particularly around V&V.

³⁴¹ 2.1 Software Quality

Quality... you know what it is, yet you don't know what it is.

³⁴² ROBERT PIRSIG, 1974 [187]

³⁴³ The philosophical viewpoint of ‘quality’ remains highly debated and there are multi-
³⁴⁴ ple facets to perceive this complex concept [83]. Transcendentally, a viewpoint like
³⁴⁵ that of Pirsig’s above shows that quality is not tangible but still recognisable; it’s hard
³⁴⁶ to explicitly define but you know when it’s missing. The International Organization

³⁴⁷ for Standardization provides a breakdown of seven universally-applicable principles
³⁴⁸ that defines quality for organisations, developers, customers and training providers
³⁴⁹ [105]. More pertinently, the 1986 ISO standard for quality was simply “the totality
³⁵⁰ of characteristics of an entity that bear on its ability to satisfy stated or implied
³⁵¹ needs” [104].

³⁵² Using this sentence, what characteristics exist for non-deterministic CISs like
³⁵³ that of a cvCIS? How do we know when the system has satisfied its ‘stated or implied
³⁵⁴ needs’ when the system can only give us uncertain probabilities in its outputs? Such
³⁵⁵ answers can be derived from related definitions—such as ‘conformance to specifica-
³⁵⁶ tion or requirements’ [48, 85], ‘meeting or exceeding customer expectation’ [176],
³⁵⁷ or ‘fitness for use’ [115]—but these then still depend on the solution description or
³⁵⁸ requirements specification, and thus the same questions still apply.

³⁵⁹ *Software* quality is somewhat more concrete. Pressman [188] adapted the
³⁶⁰ manufacturing-oriented view of quality from [21] and phrased software quality
³⁶¹ under three core pillars:

- ³⁶² • **effective software processes**, where the infrastructure that supports the cre-
³⁶³ ation of quality software needs is effective, i.e., poor checks and balances,
³⁶⁴ poor change management and a lack of technical reviews (all that lie in the
³⁶⁵ *process* of building software, rather than the software itself) will inevitably
³⁶⁶ lead to a poor quality product and vice-versa;
- ³⁶⁷ • **building useful software**, where quality software has fully satisfied the end-
³⁶⁸ goals and requirements of all stakeholders in the software (be it explicit or
³⁶⁹ implicit requirements) *in addition to* delivering these requirements in reliable

370 and error-free ways; and lastly

371 • **adding value to both the producer and user**, where quality software provides

372 a tangible value to the community or organisation using it to expedite a

373 business process (increasing profitability or availability of information) *and*

374 provides value to the software producers creating it whereby customer support,

375 maintenance effort, and bug fixes are all reduced in production.

376 In the context of a non-deterministic CIS, however, are any of the above actually

377 guaranteed? Given that the core of a system built using a CIS is fully dependent

378 on the *probability* that an outcome is true, what assurances must be put in place to

379 provide developers with the checks and balances needed to ensure that their software

380 is built with quality? For this answer, we re-explore the concept of Verification &

381 Validation (V&V).

382 2.1.1 Validation and Verification

383 To explain V&V, we analogously recount a tale given by Pham [185] on his works

384 on reliability. A high-school student sat a standardised test that was sent to 350,0000

385 students [233]. A multiple-choice algebraic equation problem used a variable, *a*,

386 and intended that students *assume* that the variable was non-negative. Without

387 making this assumption explicit, there were two correct answers to the multiple

388 choice answer. Up to 45,000 students had their scores retrospectively boosted by up

389 to 30 points for those who ‘incorrectly’ answered, however, outcomes of a student’s

390 higher education were, thereby, affected by this one oversight in quality assessment.

391 The examiners wrote a poor question due to poor process standards to check if

³⁹² their ‘correct’ answers were actually correct. The examiners “didn’t build the right
³⁹³ product” nor did they “build the product right” by writing an poor question and
³⁹⁴ failing to ensure quality standards, in the phrases Boehm [24] coined.

³⁹⁵ This story describes the issues with the cost of quality [23] and the importance
³⁹⁶ of V&V: just as the poorly written exam question had such a high toll the 45,000
³⁹⁷ unlucky students, so does poorly written software in production. As summarised
³⁹⁸ by Pressman [188], data sourced from Cigital Inc. [44] in a large-scale application
³⁹⁹ showed that the difference in cost to fix a bug in development versus system testing
⁴⁰⁰ is \$6,159 per error. In safety-critical systems, such as self-driving cars or clinical
⁴⁰¹ decision support systems, this cost skyrockets due to the extreme discipline needed
⁴⁰² to minimise error [235].

⁴⁰³ Formally, we refer to the IEEE Standard Glossary of Software Engineering
⁴⁰⁴ Terminology [102] for to define V&V:

⁴⁰⁵ **verification** The process of evaluating a system or component to determine
⁴⁰⁶ whether the products of a given development phase satisfy the
⁴⁰⁷ conditions imposed at the start of that phase.

⁴⁰⁸ **validation** The process of evaluating a system or component during or at the
⁴⁰⁹ end of the development process to determine whether it satisfies
⁴¹⁰ specified requirements.

⁴¹¹ Thus, in the context of a CIS, we have two perspectives on V&V: that of the API
⁴¹² provider and consumer (Figure 2.2).

⁴¹³ The verification process of API providers ‘leak’ out to the context of the devel-

⁴¹⁴ oper's project dependent on the CIS. Poor verification in the *internal quality* of the
⁴¹⁵ CIS will entail poor process standards, such as poor definitions and terminology used,
⁴¹⁶ support tooling and description of documentations [223]. Though it is common-
⁴¹⁷ place for providers to have a 'ship-first-fix-later' mentality of 'good-enough' software
⁴¹⁸ [242], the consequence of doing so leads to consumers absorbing the cost. Thus
⁴¹⁹ API providers must ensure that their verification strategies are rigorous enough for
⁴²⁰ the consumers in the myriad contexts they wish to use it in. Studies have considered
⁴²¹ V&V in the context of web services on the cloud [36, 37, 73, 94, 162, 164, 210, 257],
⁴²² though little have recently considered how adding 'intelligence' to these services af-
⁴²³ fects existing proposed frameworks and solutions. For a cvCIS, what might this
⁴²⁴ entail? Which assurances are given to the consumers, and how is that information
⁴²⁵ communicated? To verify if the service is working correctly, does that mean that
⁴²⁶ we need to deploy the system first to get a wider range of data, given the stochastic
⁴²⁷ nature of the black box?

⁴²⁸ Likewise, the validation perspective comes from that of the consumer. While the
⁴²⁹ former perspective is of creation, this perspective comes from end-user (developer)
⁴³⁰ expectation. As described in Chapter 1, a developer calls the CIS component using
⁴³¹ an API endpoint. Again, the mindset problem arises; does the developer know what
⁴³² to expect in the output? What are their expectations for their specific context? In
⁴³³ the area of non-deterministic systems of probabilistic output, can the developer be
⁴³⁴ assured that what they enter in a testing phase outcome the same result when in
⁴³⁵ production?

⁴³⁶ Therefore, just as the test answers with were both correct and incorrect at the

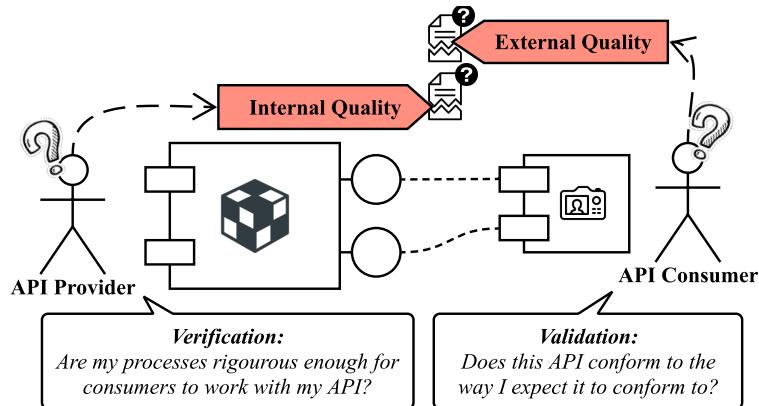


Figure 2.2: The ‘leakage’ of internal quality into the API consumer’s product and external quality imposing on the API provider.

⁴³⁷ same time, so is the same with CISs returning a probabilistic result: no result is
⁴³⁸ certain. While V&V has been investigated in the area of mathematical and earth
⁴³⁹ sciences for numerical probabilistic models and natural systems [172, 209], from
⁴⁴⁰ the software engineering literature, little work has been achieved to look at the
⁴⁴¹ surrounding area of probabilistic systems hidden behind API calls.

⁴⁴² Now that a developer is using a probabilistic system behind a deterministic API
⁴⁴³ call, what does it mean in the context of V&V? Do current verification approaches
⁴⁴⁴ and tools suffice, and if not, how do we fix it? From a validation perspective of
⁴⁴⁵ ML and end-users, after a model is trained and an inference is given and if the
⁴⁴⁶ output data point is incorrect, how will end users report a defect in the system?
⁴⁴⁷ Compared to deterministic systems where such tooling as defect reporting forms are
⁴⁴⁸ filled out (i.e., given input data in a given situation and the output data was X), how
⁴⁴⁹ can we achieve similar outputs when the system is non-deterministic? A key
⁴⁵⁰ problem with the probabilistic mindset is that once a model is ‘fixed’ by retraining

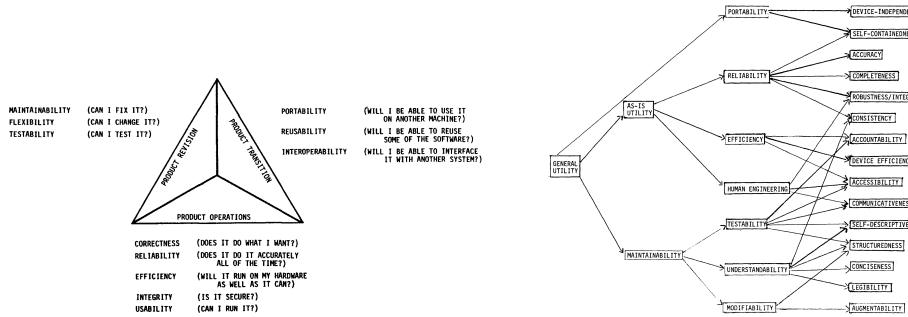
451 it, while one data-point may be fixed, others may now have been effected, thereby
452 not ensuring 100% validation. Thus, due to the unpredictable and blurry nature of
453 probabilistic systems, V&V must be re-thought out extensively.

454 2.1.2 Quality Attributes and Models

455 Similarly, quality models are used to capture internal and external quality attributes
456 via measurable metrics. Is a similar issue reflected from that of V&V due to
457 nondeterministic systems? As there is no ‘one’ definition of quality, there have been
458 differing perspectives with literature placing varying value on disparate attributes.

459 Quality attribute assessment models (like those shown in Figure 2.3) are an early
460 concept in software engineering, and systematically evaluating software quality
461 appears as early as 1968 [208]. Rubey and Hartwick’s 1968 study introduced the
462 phrase ‘attributes’ as a “prose expression of the particular quality of desired software”
463 (as worded by Boehm et al. [25]) and ‘metrics’ as mathematical parameters on a
464 scale of 0 to 100. Early attempts to categorise wider factors under a framework was
465 proposed by McCall, Richards, and Walters in the late 1970s [40, 148]. This model
466 described quality from the three perspectives of product revision (*how can we keep*
467 *the system operational?*), transition (*how can we migrate the system as needed?*)
468 and operation (*how effective is the system at achieving its tasks?*) (Figure 2.3a).

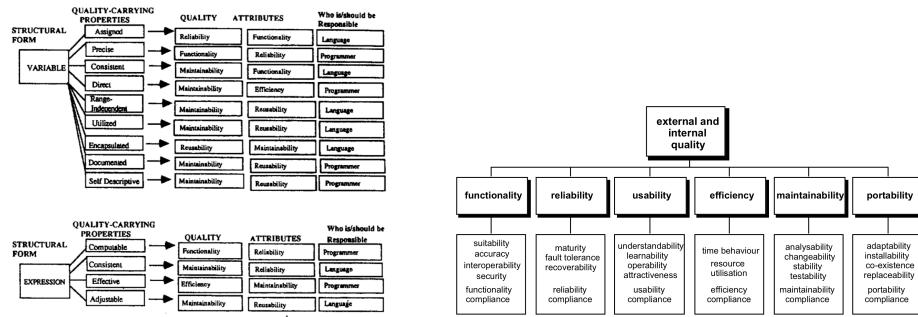
469 The model also introduced 11 attributes alongside numerous direct and indirect
470 measures to help quantify quality. This model was further developed by Boehm
471 et al. [25] who independently developed a similar model, starting with an initial set
472 of 11 software characteristics. It further defined candidate measurements of Fortran



(a) McCall's quality software factors (1977) (b) Bohem's software quality characteristics

[148].

tree (1978) [25].



(c) Dromey's quality-carrying properties and (d) ISO/IEC software product evaluation characteristics

programming languages (1995) [64].

characteristics (1999) [106].

Figure 2.3: A brief overview of the development of software quality models since 1977.

473 code to such characteristics, taking shape in a tree-like structure as in Figure 2.3b.

474 In the mid-1990s, Dromey's interpretation [64] defined a set of quality-carrying

475 properties with structural forms associated to specific programming languages and

476 conventions (Figure 2.3c). The model also supported quality defect identification

477 and proposed an improved auditing method to automate defect detection for code

478 editors in IDEs. As the need for quality models became prevalent, the International

479 Organization for Standardization standardised software quality under ISO/IEC-9126

480 [106] (the Software Product Evaluation Characteristics, Figure 2.3d), which has since

481 recently been revised to ISO/IEC-25010 with the introduction of the Systems and
482 software Quality Requirements and Evaluation (SQuaRE) model [103], separating
483 quality into *Product Quality* (consisting of eight quality characteristics and 31 sub-
484 characteristics) and *Quality In Use* (consisting of five quality characteristics and 9
485 sub-characteristics). An extensive review on the development of quality models in
486 software engineering is given in [5].

487 Of all the models described, there is one quality attribute that relates most with
488 our narrative of CIS quality: reliability. The definition of reliability is similar among
489 all quality models:

490 **McCall et al.** Extent to which a program can be expected to perform its in-
491 tended function with required precision [149].

492 **Boehm et al.** Code possesses the characteristic *reliability* to the extent that it
493 can be expected to perform its intended functions satisfactorily
494 [25].

495 **Dromey** Functionality implies reliability. The reliability of software is
496 therefore dependent on the same properties as functionality, that
497 is, the correctness properties of a program [64].

498 **ISO/IEC-9126** The capability of the software product to maintain a specified
499 level of performance when used under specified conditions [106].

500 These definitions strongly relate to the system's solution description in that
501 reliability is the ability to maintain its *functionality* under given conditions. But

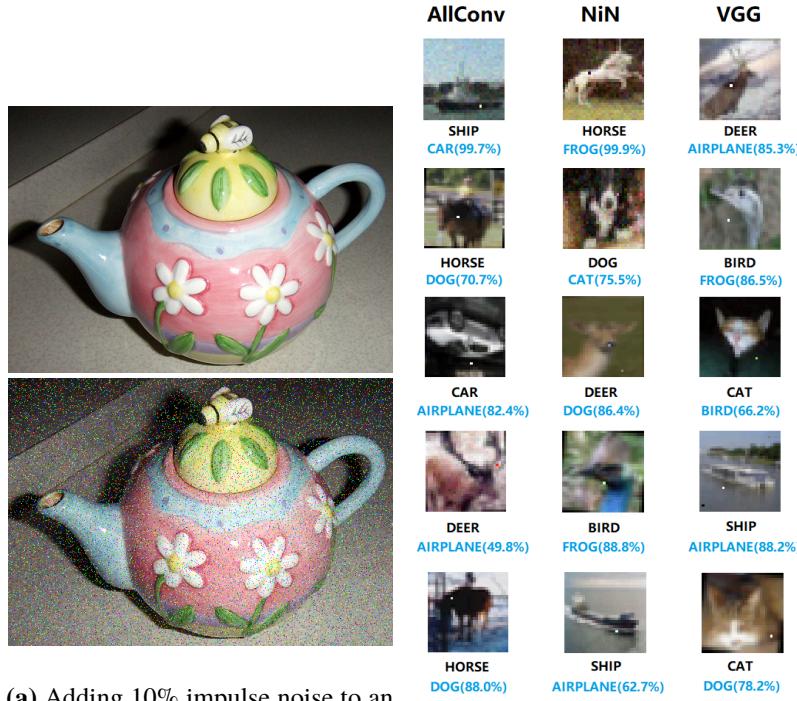
502 what defines reliability when the nature of a CIS in itself is inherently unpredictable
503 due to its probabilistic implementation? Can a non-deterministic system ever be
504 considered reliable when the output of the system is uncertain? How do developers
505 perceive these quality aspects of reliability in the context of such systems? A system
506 cannot be perceived as ‘reliable’ if the system cannot reproduce the same results due
507 to a probabilistic nature. Therefore, we believe the literature of quality models does
508 not suffice in the context of CIS reliability; a cvCIS can interpret an image of a dog
509 as a ‘Dog’ one day, but what if the next it interprets such image more specifically to
510 the breed, such as ‘Border Collie’? Does this now mean the system is unreliable?

511 Moreover, defining these systems in themselves is challenging when require-
512 ments specifications and solution descriptions are dependent on nondeterministic
513 and probabilistic algorithms. We discuss this further in Section 2.2.

514 **2.1.3 Reliability in Computer Vision**

515 Testing computer vision deep-learning reliability is an area explored typically
516 through the use of adversarial examples [231]. These input examples are where
517 images are slightly perturbed to maximise prediction error but are still interpretable
518 to humans. Refer to Figure 2.4.

519 Google Cloud Vision, for instance, fails to correctly classify adversarial examples
520 when noise is added to the original images [98]. Rosenfeld et al. [206] illustrated
521 that inserting synthetic foreign objects to input images (e.g., a cartoon elephant)
522 can alter classification output. Wang et al. [245] performed similar attacks on a
523 transfer-learning approach of facial recognition by modifying pixels of a celebrity’s



(a) Adding 10% impulse noise to an

image of a teapot changes Google

Cloud Vision's label from *teapot*(above) to *biology* (below) [98].

(b) One-pixel attacks applied to three

Neural Network (NN): AllConv, NiN

and VGG [225].



(c) Adversarial examples to trick face recognition from the source to target

images [245].

Figure 2.4: Sample adversarial examples in state-of-the-art computer vision studies.

524 face to be recognised as a different celebrity, all while still retaining the same human-
525 interpretable original celebrity. Su et al. [225] used the ImageNet database to show
526 that 41.22% of images drop in confidence when just a *single pixel* is changed in the
527 input image; and similarly, Eykholt et al. [69] recently showed similar results that
528 made a CNN interpret a stop road-sign (with mimiccid graffiti) as a 45mph speed
529 limit sign.

530 Thus, the state-of-the-art computer vision techniques may not be reliable enough
531 for safety critical applications (such as self-driving cars) as they do not handle inten-
532 tional or unintentional adversarial attacks. Moreover, as such adversarial examples
533 exist in the physical world [68, 129], “the real world may be adversarial enough”
534 [184] to fool such software.

535 2.2 Probabilistic and Nondeterministic Systems

536 Probabilistic and nondeterministic systems are those by which, for the same given
537 input, different outcomes may result. The underlying models that power a CIS
538 are treated as though they are nondeterministic; Chapter 2 introduces CISs as es-
539 sentially black-box behaviour that can change over time. As such, we adopt the
540 nondeterministic behaviour that they present.

541 2.2.1 Interpreting the Uninterpretable

542 As the rise of applied AI increases, the need for engineering interpretability around
543 models becomes paramount, chiefly from an external quality perspective that the
544 *reliability* of the system can be inspected by end-users. Model interpretability has

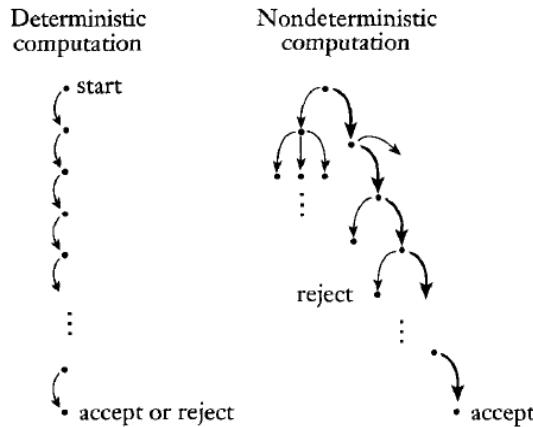


Figure 2.5: A deterministic system (left) always returns the same result in the same amount of steps. A nondeterministic system does not guarantee the same outcome, even with the same input data. Source: [72].

545 been stressed since early machine learning research in the late 1980s and 1990s (such
 546 as Quinlan [190] and Michie [155]), and although there has since been a significant
 547 body of work in the area [11, 12, 19, 26, 34, 53, 70, 79, 113, 137, 140, 147, 179,
 548 193, 207, 220, 241, 243], it is evident that ‘accuracy’ or model ‘confidence’ is still
 549 used as a primary criterion for AI evaluation [99, 108, 222]. Much research into
 550 Neural Network (NN) or Support Vector Machine (SVM) development stresses that
 551 ‘good’ models are those with high accuracy. However, is accuracy enough to justify
 552 a model’s quality?
 553 To answer this, we revisit what it means for a model to be accurate. Accuracy
 554 is an indicator for estimating how well a model’s algorithm will work with future
 555 or unforeseen data. It is quantified in the AI testing stage, whereby the algorithm
 556 is tested against cases known by humans to have ground truth but such cases are
 557 unknown by the algorithm. In production, however, all cases are unknown by both

⁵⁵⁸ the algorithm *and* the humans behind it, and therefore a single value of quality is
⁵⁵⁹ “not reliable if the future dataset has a probability distribution significantly different
⁵⁶⁰ from past data” [75], a problem commonly referred to as the *datashift* problem [228].
⁵⁶¹ Analogously, Freitas [75] provides the following description of the problem:

⁵⁶² *The military trained [a NN] to classify images of tanks into enemy
⁵⁶³ and friendly tanks. However, when the [NN] was deployed in the field
⁵⁶⁴ (corresponding to “future data”), it had a poor accuracy rate. Later,
⁵⁶⁵ users noted that all photos of friendly (enemy) tanks were taken on a
⁵⁶⁶ sunny (overcast) day. I.e., the [NN] learned to discriminate between
⁵⁶⁷ the colors of the sky in sunny vs. overcast days! If the [NN] had
⁵⁶⁸ output a comprehensible model (explaining that it was discriminating
⁵⁶⁹ between colors at the top of the images), such a trivial mistake would
⁵⁷⁰ immediately be noted.* [75]

⁵⁷¹ So, why must we interpret models? While the formal definition of what it means
⁵⁷² to be *interpretable* is still somewhat disparate (though some suggestions have been
⁵⁷³ proposed [140]), what is known is (i) there exists a critical trade-off between accuracy
⁵⁷⁴ and interpretability [59, 74, 89, 112, 118, 259], and (ii) a single quantifiable value
⁵⁷⁵ cannot satisfy the subjective needs of end-users [75]. As ever-growing domains
⁵⁷⁶ ML become widespread¹, these applications engage end-users for real-world goals,
⁵⁷⁷ unlike the aims in early ML research where the aim was to get AI working in the
⁵⁷⁸ first place. In safety-critical systems where AI provide informativeness to humans

¹In areas such as medicine [18, 34, 67, 109, 113, 134, 180, 195, 241, 256, 261], bioinformatics [58, 76, 110, 117, 230], finance [12, 56, 100] and customer analytics [137, 243].

⁵⁷⁹ to make the final call (see [38, 100, 119]), there is often a mismatch between the
⁵⁸⁰ formal objectives of the model (e.g., to minimise error) and complex real-world
⁵⁸¹ goals, where other considerations (such as the human factors and cognitive science
⁵⁸² behind explanations²) are not realised: model optimisation is only worthwhile if they
⁵⁸³ “actually solve the original [human-centred] task of providing explanation” [163]
⁵⁸⁴ to end-users. **Therefore, when human-decision makers must be interpretable**
⁵⁸⁵ **themselves [196], any AI they depend on must also be interpretable.**

⁵⁸⁶ Recently, discussion behind such a notion to provide legal implications of in-
⁵⁸⁷ terpretability is topical. Doshi-Velez et al. [62] discuss when explanations are not
⁵⁸⁸ provided from a legal stance—for instance, those affected by algorithmic-based de-
⁵⁸⁹ cisions have a ‘right to explanation’ [88, 244] under the European Union’s GDPR³.
⁵⁹⁰ But, explanations are not the only way to ensure AI accountability: theoretical
⁵⁹¹ guarantees (mathematical proofs) or statistical evidence can also serve as guarantees
⁵⁹² [62], however, in terms of explanations, what form they take and how they are proven
⁵⁹³ correct are still open questions [140].

⁵⁹⁴ 2.2.2 Explanation and Communication

⁵⁹⁵ From a software engineering perspective, explanations and interpretability are, by
⁵⁹⁶ definition, inherently communication issues: what lacks here is a consistent interface
⁵⁹⁷ between the AI system and the person using it. The ability to encode ‘common
⁵⁹⁸ sense reasoning’ [150] into programs today has been achieved, but *decoding* that
⁵⁹⁹ information is what still remains problematic. At a high level, Shannon and Weaver’s

²*Interpretations* and *explanations* are often used interchangeably.

³<https://www.eugdpr.org> last accessed 13 August 2018.

600 theory of communication [215] applies, just as others have done with similar issues
601 in the Software Engineer (SE) realm [158, 251] (albeit to the domain of visual
602 notations). Humans map the world in higher-level concepts easily when compared
603 to AI systems: while we think of a tree first (not the photons of light or atoms that
604 make up the tree), an algorithm simply sees pixels, and not the concrete object [62]
605 and the AI interprets the tree inversely to humans. Therefore, the interpretation or
606 explanation is done inversely: humans do not explain the individual neurons fired to
607 explain their predictions, and therefore the algorithmic transparent explanations of
608 AI algorithms (“*which neurons were fired to make this AI think this tree is a tree?*”)
609 do not work here.

610 Therefore, to the user (as mapped using Shannon and Weaver’s theory), an AI
611 pipeline (the communication *channel*) begins with a real-world concept, y , that acts
612 as an *information source*. This information source is fed in as a *message*, x , (as pixels)
613 to an AI system (the *transmitter*). The transmitter encodes the pixels to a prediction,
614 \hat{y} , the *signal* of the message. This signal is decoded by the *receiver*, an explanation
615 system, $e_x(x, \hat{y})$, that tailors the prediction with the given input data to the intended
616 end user (the *destination*) as an explanation, \tilde{y} , another type of *message*. Therefore,
617 the user only sees the channel as an input/output pipeline of real-world objects, y ,
618 and explanations, \tilde{y} , tailored to *them*, without needing to see the inner-mechanics of
619 a prediction \hat{y} . We present this diagrammatically in Figure 2.6.

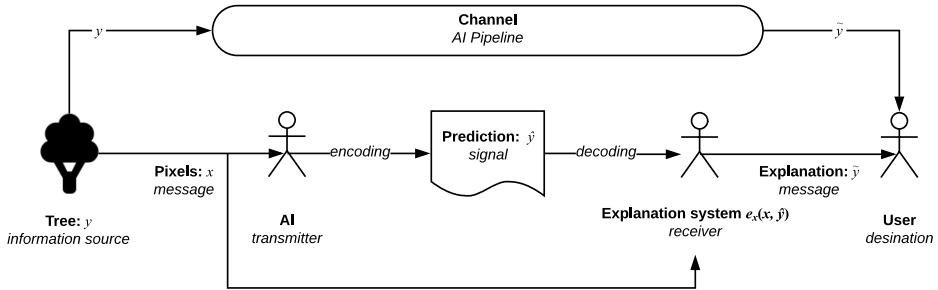


Figure 2.6: Theory of AI communication from information source, y , to intended user as explanations \tilde{y} .

620 2.2.3 Mechanics of Model Interpretation

621 How do we interpret models? Methods for developing interpretation models include:
 622 decision trees [31, 46, 92, 189, 203], decision tables [13, 137] and decision sets [131,
 623 163]; input gradients, gradient vectors or sensitivity analysis [12, 135, 193, 207, 213];
 624 exemplars [77, 120]; generalised additive models [38]; classification (*if-then*) rules
 625 [28, 45, 174, 238, 253] and falling rule lists [220]; nearest neighbours [147, 214,
 626 229, 250, 260] and Naïve Bayes analysis [18, 42, 78, 95, 126, 134, 154, 261].

627 Cross-domain studies have assessed the interpretability of these techniques
 628 against end-users, measuring response time, accuracy in model response and user
 629 confidence [6, 76, 93, 100, 147, 211, 226, 243], although it is generally agreed
 630 that decision rules and decision tables provide the most interpretation in non-linear
 631 models such as SVMs or NNs [76, 147, 243]. For an extensive survey of the benefits
 632 and fallbacks of these techniques, we refer to Freitas [75], Doshi-Velez et al. [62]
 633 and Doshi-Velez and Kim [61].

634 2.3 Application Programming Interfaces

635 Application Programming Interfaces (APIs) are the interface between a developer
636 needs and the software components at their disposal [8] by abstracting the underlying
637 component behind a subroutine, protocol or specific tool. Therefore, it is natural to
638 assess internal quality (and external quality if the software is in itself a service to
639 be used by other developers—in this case a CIS) is therefore directly related to the
640 quality the API offers [125].

641 Good APIs are known to be intuitive and require less documentation browsing
642 [186], thereby increasing developer productivity. Conversely, poor APIs are those
643 that are hard to interpret, thereby reducing developer productivity and product qual-
644 ity. The consequences of this have shown a higher demand of technical support (as
645 measured in [96]) that, ultimately, causes the maintenance to be far more expensive,
646 a phenomenon widely known in software engineering economics (see Section 2.1.1).

647 While there are different types of APIs, such as software library/framework
648 APIs for building desktop software, operating system APIs for interacting with the
649 operating system, remote APIs for communication of varying technologies through
650 common protocols, we focus on web APIs for communication of resources over
651 the web (being the common architecture of cloud-based services). Further infor-
652 mation on the development, usage and documentation of web APIs is provided in
653 Appendix A.1.

2.3.1 API Usability

654 If a developer doesn't understand the overarching concepts of the context behind
655 the API they wish to use, then they cannot formulate what gaps in their knowledge
656 is missing. For example, a developer that knows nothing about machine learning
657 techniques in computer vision cannot effectively formulate queries to help bridge
658 those gaps in their understanding to figure out more about the cvCIS they wish to
659 use.

660 Balancing the understanding of the information need (both conscious and un-
661 conscious), how to phrase that need and how to query it in an information retrieval
662 system is concept long studied in the information sciences [236]. In API design,
663 the most common form to convey knowledge to developers is through annotated
664 code examples and overviews to a platform's architectural and design decisions
665 [29, 60, 161, 201] though these studies have not effectively communicated *why* these
666 artefacts are important. What makes the developer *conceptually understand* these
667 artefacts?

668 Robillard and Deline [201] conducted a multi-phase, mixed-method approach to
669 create knowledge grounded in the professional experience of 440 software engineers
670 at Microsoft of varying experience to determine what makes APIs hard to learn,
671 the results of which previously published in an earlier report [200]. Their results
672 demonstrate that 'documentation-related obstacles' are the biggest hurdle in learning
673 new APIs. One of these implications are the *intent documentation* of an API (i.e.,
674 *what is the intent for using a particular API?*) and such documentation is required
675 only where correct API usage is not self-evident, where advanced uses of the API are

677 documented (but not the intent), and where performance aspects of the API impact
678 the application developed using it. They conclude that professional developers do
679 not struggle with learning the *mechanics* of the API, but in the *understanding* of how
680 the API fits in upwards to its problem domain and downward to its implementation:

681 *In the upwards direction, the study found that developers need help*
682 *mapping desired scenarios in the problem domain to the content of the*
683 *API, and in understanding what scenarios or usage patterns the API*
684 *provider intends and does not intend to support. In the downwards*
685 *direction, developers want to understand how the API's implementation*
686 *consumes resources, reports errors and has side effects.* [201]

687 These results particularly corroborate to that of previous studies where devel-
688 opers quote that they feel that existing learning content currently focuses on “*how*
689 to do things, not necessarily *why*” [169]. This thereby reiterates the conceptual
690 understanding of an API as paramount.

691 A later study by Ko and Riche [124] assessed the importance of a programmer’s
692 conceptual understanding of the background behind the task before implementing the
693 task itself, a notion that we find most relevant for users of CIS APIs. While the study
694 did not focus on developing web APIs (rather implementing a Bluetooth application
695 using platform-agnostic terminology), the study demonstrated how developers show
696 little confidence in their own metacognitive judgements to understand and assess the
697 feasibility of the intent of the API and understand the vocabulary and concepts within
698 the domain (i.e., wireless connectivity). This indecision over what search results

699 were relevant in their searches ultimately hindered their progress implementing the
700 functionality, again decreasing productivity. Ko and Riche suggest to improve API
701 usability by introducing the background of the API and its relevant concepts using
702 glossaries linked to tutorials to each of the major concepts, and then relate it back to
703 how to implement the particular functionality.

704 Thus, an analysis of the conceptual understanding of CIS APIs by a range of
705 developers (from beginner to professional) is critical to best understand any differ-
706 ences between existing studies and those that are nondeterministic. Our proposal is
707 to perform similar survey research (see Chapter 3) in the search for further insight
708 into the developer's approach toward existing CIS APIs.

CHAPTER 3

709

710

711

Research Methodology

712

713 Investigating software engineering practices is often a complex task as it is imperative
714 to understand the social and cognitive processes around software engineers and not
715 just the tools and processes used [65]. This chapter explores our research method-
716 ology by exploring five key elements of empirical software engineering research:
717 firstly, (i) we provide an extended focus to the study by reviewing our research
718 questions (see Section 1.2) anchored under the context of an existing classifica-
719 tion taxonomy, (ii) characterise our research goals through an explicit philosophical
720 stance, (iii) explain how the stance selected impacts our selection of research meth-
721 ods and data collection techniques (by dissecting our choice of methods used to
722 reach these research goals), (iv) discuss a set of criteria for assessing the validity of
723 our study design and the findings of our research, and lastly (v) discuss the practical

⁷²⁴ considerations of our chosen methods.

⁷²⁵ The foundations for developing this research methodology has been expanded
⁷²⁶ from that proposed by Easterbrook et al. [65], Wohlin and Aurum [254], Wohlin
⁷²⁷ et al. [255] and Shaw [216].

⁷²⁸ 3.1 Research Questions Revisited

⁷²⁹ In Section 1.2, we introduce three hypothesis of this study (RH1–RH3), namely:

⁷³⁰ (i) existing CIS APIs are poorly documented for general use (RH1); (ii) existing
⁷³¹ CIS APIs do not provide sufficient metadata when used in context-specific use cases
⁷³² (RH2); and (iii) the combination of improving documentation and metadata will
⁷³³ ultimately improve one of software quality, developer productivity and/or developer
⁷³⁴ understanding (RH3).

⁷³⁵ To discuss our research strategy, we revisit our research questions through the
⁷³⁶ classification technique discussed by Easterbrook et al. [65], a technique originally
⁷³⁷ proposed in the field of psychology by Meltzoff [152] but adapted to software engi-
⁷³⁸ neering. Our research study involves a mix of five *knowledge questions*, that focus
⁷³⁹ on existing practices and the ways in which they work, and two *design questions*, that
⁷⁴⁰ focuses on designing better ways to approach software engineering tasks [218]. Both
⁷⁴¹ classes of questions are respectively concerned with empirical and non-empirical
⁷⁴² software engineering that, in practice, are best combined in long-term software engi-
⁷⁴³ neering research studies (such as this one) as they assist in tackling the investigation
⁷⁴⁴ of a specific problem, approaches to solve that problem and finding what solutions

⁷⁴⁵ work best [252].

⁷⁴⁶ 3.1.1 Knowledge Questions

⁷⁴⁷ In total, five knowledge questions are posed in this study to help us understand the
⁷⁴⁸ way developers currently interact and work with a CIS API; two exploratory, one
⁷⁴⁹ base-rate, and two relationship and causality questions.

⁷⁵⁰ We begin by formulating two *exploratory questions* to attempt to better under-
⁷⁵¹ stand the phenomena of poor API documentation and metadata; both RQ1.1 and
⁷⁵² RQ2.1 respectively describe and classify what practices are in use for existing CIS
⁷⁵³ API documentation and what problems currently exist when no metadata is returned.
⁷⁵⁴ Answering these two questions assists in refining preciser terms of the phenomena,
⁷⁵⁵ ways in which we find evidence for them and ensuring the data found is valid.

⁷⁵⁶ By answering these questions, we have a clearer understanding of the phenom-
⁷⁵⁷ ena; we then follow up by posing an additional *base-rate question* that helps provide
⁷⁵⁸ a basis to confirm that the phenomena occurring is normal (or unusual) behaviour
⁷⁵⁹ by investigating the patterns of phenomena's occurrence. RQ1.2 is a descriptive-
⁷⁶⁰ process question to help us understand how the developer currently understands
⁷⁶¹ existing CIS API documentation, given their lack of formal extended training in
⁷⁶² artificial intelligence. This achieves us an insight into the developer's mindset and
⁷⁶³ regular thought patterns toward these APIs.

⁷⁶⁴ Lastly, we investigate the relationship between the improved documentation
⁷⁶⁵ and improvements to other aspects of the software development process. Chiefly,
⁷⁶⁶ RQ3.1 is concerned with whether any improvements to metadata or documentation

⁷⁶⁷ correlate to improvements in software quality, developer productivity, or developer
⁷⁶⁸ education (and is a *relationship establishment question*). If we establish such a
⁷⁶⁹ relationship, we refine the question and investigate the specific causes using three
⁷⁷⁰ *causality questions* defined under RQ3.2, namely by associating three classes of
⁷⁷¹ measurable metrics (internal quality metrics, external quality metrics, developer
⁷⁷² education insight metrics) to the improved documentation.

⁷⁷³ 3.1.2 Design Questions

⁷⁷⁴ RQ1.3 and RQ2.2 are both *design questions*; they are concerned with ways in which
⁷⁷⁵ we can improve a CIS by investigating what additional attributes are needed in both
⁷⁷⁶ the documentation and metadata that assist developers to achieve their goals. They
⁷⁷⁷ are not classified as knowledge questions as we investigate what *will be* and not *what*
⁷⁷⁸ *is*. By understanding the process by which developers desire additional attributes
⁷⁷⁹ of metadata and documentation, we can help shape improvements to the existing
⁷⁸⁰ design of a CIS.

⁷⁸¹ 3.2 Philosophical Stances

⁷⁸² Philosophical stances guide the researcher's action by fortifying what constitutes
⁷⁸³ 'valid truth' against a fundamental set of core beliefs [198]. In software engineer-
⁷⁸⁴ ing, four dominant philosophical stances are commonly characterised [47, 182]:
⁷⁸⁵ positivism (or post-positivism), constructivism (or interpretivism), pragmatism, and
⁷⁸⁶ critical theory (or advocacy/participatory). To construct such a 'validity of truth',
⁷⁸⁷ we will review these four philosophical stances in this section, and state the stance

788 that we explicitly adopt and our reasoning for this.

789 **Positivism** Positivists claim truth to be all observable facts, reduced piece-by-
790 piece to smaller components which is incrementally verifiable to form truth. We
791 do not base our work on the positivistic stance as the theories governing verifiable
792 hypothesis must be precise from the start of the research. Moreover, due to its
793 reductionist approach, it is difficult to isolate these hypotheses and study them in
794 isolation from context. As our hypotheses are not context-agnostic, we steer clear
795 from this stance.

796 **Constructivism** Constructivists see knowledge embedded within the human con-
797 text; truth is the *interpretive* observation by understanding the differences in human
798 thought between meaning and action [123]. That is, the interpretation of the theory
799 is just as important to the empirical observation itself. We partially adopt a con-
800 structivist stance as we attempt to model the developer's mindset, being an approach
801 that is rich in qualitative data on human activity.

802 **Pragmatism** Pragmatism is a less dogmatic approach that encourages the incom-
803 plete and approximate nature of knowledge and is dependent on the methods in which
804 the knowledge was extracted. The utility of consensually agreed knowledge is the
805 key outcome, and is therefore relative to those who seek utility in the knowledge—
806 what is the useful for one person is not so for the other. While we value the utility
807 of knowledge, it is difficult to obtain consensus especially on an ill-researched topic
808 such as ours, and therefore we do not adopt this stance.

⁸⁰⁹ **Critical Theory** This study chiefly adopts the philosophy of critical theory [35]. A
⁸¹⁰ key outcome of the study is to shift the developer's restrictive deterministic mindset
⁸¹¹ and shed light on developing a new framework actively with the developer community
⁸¹² that seeks to improve the process of using such APIs. In software engineering,
⁸¹³ critical theory is used to "actively [seek] to challenge existing perceptions about
⁸¹⁴ software practice" [65], and this study utilises such an approach to shift the mindset
⁸¹⁵ of CIS consumers and providers alike on how the documentation and metadata
⁸¹⁶ should not be written with the 'traditional' deterministic mindset at heart. Thus, our
⁸¹⁷ key philosophical approach is critical theory to seek out *what-can-be* using partial
⁸¹⁸ constructivism to model the current *what-is*.

⁸¹⁹

3.3 Research Design

⁸²⁰ Research methods are "a set of organising principles around which empirical data
⁸²¹ is collection and analysed" [65]. Creswell and Creswell [47] suggest that strong
⁸²² research design is reflected when the weaknesses of multiple methods complement
⁸²³ each other. Using a mixed-methods approach is therefore commonplace in software
⁸²⁴ engineering research, typically due to the human-oriented nature investigating how
⁸²⁵ software engineers work both individually (where methods from psychology may be
⁸²⁶ employed) and together (where methods from sociology may be employed).

⁸²⁷ Therefore, studies in software engineering are typically performed as field studies
⁸²⁸ where researchers and developers (or the artefacts they produce) are analysed either
⁸²⁹ directly or indirectly [219]. The mixed-methods approach combines five classes

830 of field study methods (or empirical strategies/studies) most relevant in empirical
831 software engineering research [65, 116, 255]: controlled experiments, case studies,
832 survey research, ethnographies, and action research. We chiefly adopt a mixed-
833 methods approach to our work using the *concurrent triangulation* mixed-methods
834 strategy [111] as it best compensates for weaknesses that exist in all research methods,
835 and employs the best strengths of others.

836 3.3.1 Review of Relevant Research Methods

837 Below we review some of the research methods most relevant to our research ques-
838 tions as refined in Section 3.1 as presented by Easterbrook et al. [65].

839 **Controlled Experiments** A controlled experiment is an investigation of a clear,
840 testable hypothesis that guides the researcher to decide and precisely measure how
841 at least one independent variable can be manipulated and effect at least one other
842 dependent variable. They determine if the two variables are related and if a cause-
843 effect relationship exists between them. The combination of independent variable
844 values is a *treatment*. It is common to recruit human subjects to perform a task and
845 measure the effect of a randomly assigned treatment on the subjects, though it is
846 not always possible to achieve full randomisation in real-life software engineering
847 contexts, in which case a *quasi-experiment* may be employed where subjects are not
848 randomly assigned to treatments.

849 While we have defined hypotheses (RH1–RH3), refining them into precise,
850 measurable variables is challenging due to the qualitative nature they present. A

851 well-defined population is also critical and must be easily accessible; the varied
852 range of beginner to expert software engineers with varied understanding of artifi-
853 cial intelligence concepts is required to perform controlled experiments, and thus
854 recruitment may prove challenging. Lastly, the controlled experiment is essentially
855 reductionist by affecting a small amount of variables of interest and controlling all
856 others. This approach is too clinical for the practical outcomes by which our research
857 goals aim for, and is therefore closely tied to the positivist stance.

858 **Case Studies** Case studies investigate phenomena in their real-life context and are
859 well-suited when the boundary between context and phenomena is unknown [258].
860 They offer understanding of how and why certain phenomena occur, thereby inves-
861 tigating ways cause-effect relationships can occur. They can be used to test existing
862 theories (*confirmatory case studies*) by refuting theories in real-world contexts in-
863 stead of under laboratory conditions or to generate new hypotheses and build theories
864 during the initial investigation of some phenomena (*exploratory case studies*).

865 Case studies are well-suited where the context of a situation plays a role in the
866 phenomenon being studied, which we specifically relate back to RH2 (RQ2.1 and
867 RQ2.2) in exploring whether the context of an application using a CIS requires the
868 CIS context-specific or of context-agnostic. They also lend themselves to purposive
869 sampling rather than random sampling, and thus we can selectively choose cases that
870 benefit the research goal of RH2 and (using our critical theorist stance) select cases
871 that will actively benefit our participant software engineering audience most to draw
872 attention to situations regarded as most problematic.

⁸⁷³ **Survey Research** Survey research identifies characteristics of a broad population
⁸⁷⁴ of individuals through direct data collection techniques such as interviews and ques-
⁸⁷⁵ tionnaires or independent techniques such as data logging. Defining that well-defined
⁸⁷⁶ population is critical, and selecting a representative sample from it to generalise the
⁸⁷⁷ data gathered usually assists in answering base-rate questions.

⁸⁷⁸ By identifying representative sample of the population, from beginner to ex-
⁸⁷⁹ perienced developers with varying understanding of CIS APIs, we can use survey
⁸⁸⁰ research to assist in answering our exploratory and base-rate research questions un-
⁸⁸¹ der RH1 and RH2 (see Section 3.1.1) in determining the qualitative aspects of how
⁸⁸² individual developers perceive and work with the existing APIs, either by directly
⁸⁸³ asking them or by mining third-party discussion websites such as Stack Overflow.
⁸⁸⁴ However, with direct survey research techniques, low response rates may prove
⁸⁸⁵ challenging, especially if no inducements can be offered for participation.

⁸⁸⁶ **Ethnographies** Ethnographies investigates the understanding of social interac-
⁸⁸⁷ tion within community through field observation [202]. Resulting ethnographies
⁸⁸⁸ help understand how software engineering technical communities build practices,
⁸⁸⁹ communication strategies and perform technical work collaboratively.

⁸⁹⁰ Ethnographies require the researcher to be highly trained in observational and
⁸⁹¹ qualitative data analysis, especially if the form of ethnography is participant observa-
⁸⁹² tion, whereby the researcher is embedded of the technical community for observation.
⁸⁹³ This may require the longevity of the study to be far greater than a couple of weeks,
⁸⁹⁴ and the researcher must remain part of the project for its duration to develop enough

⁸⁹⁵ local theories about how the community functions. While it assists in revealing
⁸⁹⁶ subtle but important aspects of work practices within software teams, this study
⁸⁹⁷ does not focus on the study of teams, and is therefore not a research method relevant
⁸⁹⁸ to this project.

⁸⁹⁹ **Action Research** Action researchers simultaneously solve real-world problems
⁹⁰⁰ while studying the experience of solving the problem [52] by actively seeking to
⁹⁰¹ intervene in the situation for the purpose of improving it. A precondition is to engage
⁹⁰² with a *problem owner* who is willing to collaborate in identifying and solving the
⁹⁰³ problem faced. The problem must be authentic (a problem worth solving) and must
⁹⁰⁴ have new knowledge outcomes for those involved. It is also characterised as an
⁹⁰⁵ iterative approach to problem solving, where the knowledge gained from solving the
⁹⁰⁶ problem has a desirable solution that empowers the problem owner and researcher.

⁹⁰⁷ This research is most associated to our adopted philosophical stance of critical
⁹⁰⁸ theory. As this project is being conducted under the Applied Artificial Intelligence
⁹⁰⁹ Institute (A²I²) collaboratively with engaged industry clients, we have identified a
⁹¹⁰ need for solving an authentic problem that industry faces. The desired outcome of
⁹¹¹ this project is to facilitate wider change in the usage and development of cvCISs;
⁹¹² thus, engaging action research as a primary method throughout the mixed-methods
⁹¹³ approach used in this research.

⁹¹⁴ 3.3.2 Review of Data Collection Techniques for Field Studies

⁹¹⁵ Singer et al. developed a taxonomy [136, 219] showcasing data collection techniques

916 in field studies that are used in conjunction with a variety of methods based on the
917 level of interaction between researcher and software engineer, if any. This taxonomy
918 is reproduced in Figure A.4.

919 3.4 Proposed Experiments

920 This section discusses two proposed experiments that we conduct in this study.
921 For each experiment, we describe an overview of the experiment grounded known
922 methods and techniques (Sections 3.3.1 and 3.3.2), our approach to analysing the
923 data, as well as linking the experiment back to a research hypothesis and question
924 (Section 1.2). A high-level overview of the proposed experiments and the major
925 bodies of work they encompass is presented in Figure 3.1.

926 3.4.1 Experiment I: Develop Initial Framework

927 Experiment I shapes a context-agnostic approach to understand current usage pat-
928 terns of CIS APIs and the ways by which developers interpret them. Briefly, this
929 experiment is comprised under two phases of field survey research: (i) repository
930 mining developer discussion forums (i.e., analysis of databases and documentation
931 analysis) to understand what developers currently complain about on these forums
932 and where their mismatch in understanding lies; (ii) conducting unstructured inter-
933 views and distributing a questionnaire to gather personal opinion based on individual
934 developer's anecdotal remarks.

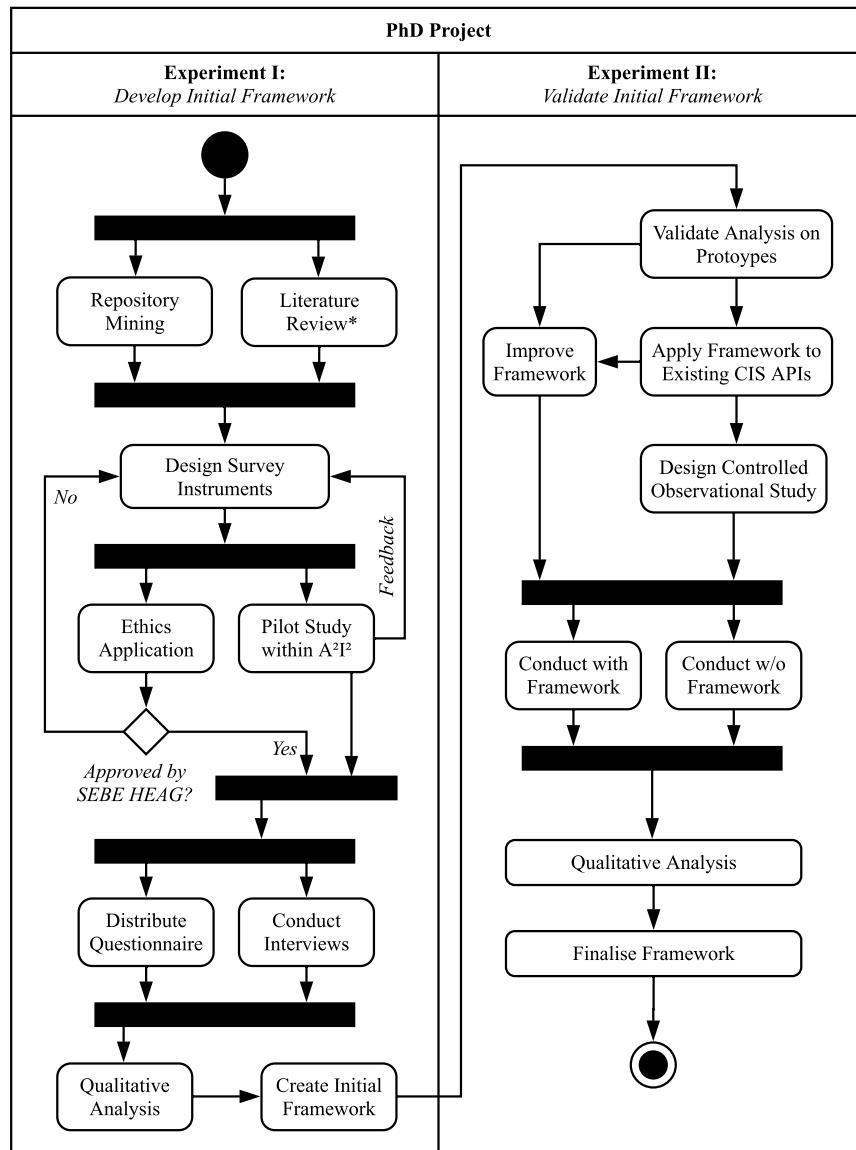


Figure 3.1: High-level activity diagram of the proposed experiments in this study. Literature review is ongoing.

935 Relevance and Motivation

936 Experiment I aims to better understand the existing mindsets that developers have

937 when approaching to use Computer Vision Cloud Intelligence Services (cvCISs).

938 This work therefore ties in to RH1; by understanding the developer mindset in how

⁹³⁹ they interpret cvCIS APIs, we are better informed to produce a framework that
⁹⁴⁰ increases the effectiveness of the documentation of those existing cvCIS providers.

⁹⁴¹ RH1 postulates that the software engineering community do not fully understand
⁹⁴² the ‘magic’ behind CIS APIs. As described in Section 1.2, they face a gap in their
⁹⁴³ understanding around the underlying architecture of pre-built, machine learnt APIs
⁹⁴⁴ (RQ1.2). Software developers are not well-supported by the CIS providers, and
⁹⁴⁵ therefore do not have a consistent set of common best practices when approaching
⁹⁴⁶ to use these APIs (RQ1.1). It is therefore necessary that CIS providers provide
⁹⁴⁷ additional information to gap this mismatched understanding (RQ1.3).

⁹⁴⁸ **Data Collection & Analysis**

⁹⁴⁹ **Phase 1: Repository Mining** Developers typically congregate in search of dis-
⁹⁵⁰ courses on issues they face in online forums, such as Stack Overflow and Quora, as
⁹⁵¹ well as writing their experiences in personal blogs such as Medium. The simplest of
⁹⁵² these platforms is Stack Overflow (a sub-community of the Stack Exchange family
⁹⁵³ of targeted communities) that specifically targets developer issues on using a simple
⁹⁵⁴ Q&A interface, where developers can discuss technical aspects and general software
⁹⁵⁵ development topics. Moreover, Stack Overflow is often acknowledged as *the ‘go-to’*
⁹⁵⁶ place for developers to find high-quality code snippets that assist in their problems
⁹⁵⁷ [227].

⁹⁵⁸ Thus, to begin validating CIS API usage and misunderstanding in a generalised
⁹⁵⁹ context (i.e., context-agnostic to the project at hand), we propose using repository
⁹⁶⁰ mining on Stack Overflow to help answer our research questions. Specifically,

⁹⁶¹ we select Stack Overflow due to its targeted community of developers¹ and the
⁹⁶² availability of its publicly available dataset released as ‘data dumps’ on the Stack
⁹⁶³ Exchange Data Explorer² and Google BigQuery³. Studies conducted have also used
⁹⁶⁴ Stack Overflow to mine developer discourse [7, 14, 15, 43, 139, 167, 175, 191, 204,
⁹⁶⁵ 221, 234, 246].

⁹⁶⁶ Due to the enormity of the data produced, we will use qualitative analysis on
⁹⁶⁷ the questions mined using assistive tools such as NVivo. For this, we will conduct a
⁹⁶⁸ thematic analysis on the themes of each question mined, the relevance of the question
⁹⁶⁹ to our research topic, and ensuring strict coding schemes (that reflect our research
⁹⁷⁰ goals) are adhered to. We refer to Singer et al. [219] and Miles et al. [156] on coding
⁹⁷¹ and analysing this qualitative data gathered.

⁹⁷² **Phase 2: Personal Opinion Surveys** We follow the triangulation approach pro-
⁹⁷³ posed by Jick [111] to corroborate the qualitative data of developers’ discussion
⁹⁷⁴ of Stack Overflow with secondary survey research, thereby validating what people
⁹⁷⁵ say on Stack Overflow with what is said and done in real life. Kitchenham and
⁹⁷⁶ Pfleeger [122] provide an introduction on methods used to conduct personal opinion
⁹⁷⁷ surveys which we adopt as an initial reference in (i) shaping our survey objectives

¹We also acknowledge that there are other targeted software engineering Stack Exchange communities such as Stack Exchange Software Engineering (<https://softwareengineering.stackexchange.com>), though (as of January 2019) this much smaller community consists of only 52,000 questions versus Stack Overflow’s 17 million.

²<https://data.stackexchange.com/stackoverflow> last accessed 17 January 2017.

³<https://console.cloud.google.com/marketplace/details/stack-exchange/stack-overflow> last accessed 17 January 2017.

⁹⁷⁸ around our research goals, (ii) designing a cross-sectional survey, (iii) developing
⁹⁷⁹ and evaluating our two survey instruments (consisting of a structured questionnaire
⁹⁸⁰ and semi-structured interview), (iv) evaluating our instruments, (v) obtaining the
⁹⁸¹ data and (vi) analysing the data.

⁹⁸² As is good practice in developing questionnaire instruments to evaluate their
⁹⁸³ reliability and validity [141], we evaluate our instrument design by asking colleagues
⁹⁸⁴ to critique it via pilot studies within A²I². This assists in identifying any problems
⁹⁸⁵ with the questionnaire itself and with any issues that may occur with the response rate
⁹⁸⁶ and follow-up procedures. We follow a similar approach by practicing the interview
⁹⁸⁷ instrument on colleagues within A²I².

⁹⁸⁸ Findings from the pilot study helps inform us for a widely distributed question-
⁹⁸⁹ naire and conducting interviews out in the field, where we recruit external software
⁹⁹⁰ engineers in industry through the industry contacts of A²I². Ethics approval from
⁹⁹¹ the Faculty of Science, Engineering and Built Environment Human Ethics Advisory
⁹⁹² Group (SEBE HEAG) will be required prior to externally conducting this survey
⁹⁹³ research (see Appendix C). The quantitative (survey) and qualitative (interview)
⁹⁹⁴ analysis allows us to shape the research outcome of RH1—an API documentation
⁹⁹⁵ quality assessment framework—and assists in stabilising our general understanding
⁹⁹⁶ of how developers use these existing APIs.

⁹⁹⁷ 3.4.2 Developing the Initial Framework

⁹⁹⁸ Our initial framework is developed using the qualitative and quantitative analyses
⁹⁹⁹ from the findings of Experiment I. As this is a creative phase in which we are

1000 developing a new framework, the exact process by which we develop the initial
1001 framework will come to light once more insight is determined. However it is
1002 anticipated discussion with other researchers and engineers at A²I² about the analyses
1003 of the findings (i.e., white-boarding sessions of potential ideas from the findings)
1004 will help develop our initial documentation framework. This framework will take
1005 the shape of a checklist or table, typical of information systems studies (e.g., [133]),
1006 that indicate what attributes should be best suited for what needs.

1007 **3.4.3 Experiment II: Validate Initial Framework**

1008 Experiment II extends the *generalised* context of Experiment I by evaluating how
1009 the findings of Experiment I translates to context-specific applications. We confirm
1010 that the generalised findings are (indeed) genuine by conducting action research in
1011 combination with an observational study on software engineers. This experiment
1012 is also compromised of: (i) development of prototypes using CIS APIs of differ-
1013 ing contexts; (ii) presenting a solution framework to developers to interpret the
1014 improvement of their understanding when using a CIS.

1015 **Relevance and Motivation**

1016 Experiment II aims to contextualise the findings from Experiment I; that is, if we
1017 add *varying contexts* to the applications we write using CIS APIs, what is needed to
1018 extend the *context-agnostic* framework developed in Experiment I? This work relates
1019 back to RH2; adding context-specific metadata to the endpoints of these APIs, we
1020 can highlight what issues exist when such metadata is not present (RQ2.1) and what

1021 types of metadata developers seek (RQ2.2).

1022 Moreover, the implication of the first two hypotheses suggest that applying an API
1023 documentation and metadata quality assessment framework may have an effect on
1024 other aspects within the software engineering process (RH3). Thus, this experiment
1025 also confirms if our framework makes an improvement to software quality, developer
1026 productivity and/or developer informativeness (RQ3.1 and RQ3.2).

1027 **Data Collection & Analysis**

1028 To confirm findings of the method within RH1 is genuine, we shift from reviewing
1029 the documentation from a general stance to a specialised (context-specific) stance in
1030 the use of these APIs.

1031 This is firstly achieved by using existing CIS APIs to develop basic ‘prototypes’,
1032 each having differing contexts. The number of prototypes to develop and the use
1033 cases they have will be informed by the results of Experiment I, and therefore cannot
1034 yet be described at this stage. Our action research in developing the prototypes will
1035 help inform any potential gaps that exist in the findings of RH1, especially with
1036 regards to context-specificity, and therefore improves the metadata component of our
1037 framework (as per the outcome of RH2).

1038 This outcome will also help us design the next stage of the experiment, con-
1039 sisting of a comparative controlled study [212] to capture firsthand behaviours and
1040 interactions toward how software engineers approach using a CIS with and without
1041 our framework applied. We will provide improved documentation and metadata re-
1042 spondes of a set of popular cvCISs that is documented with the additional metadata

¹⁰⁴³ and whose information is organised using our framework.

¹⁰⁴⁴ We then recruit 20 developers of varying experience (from beginner programmer
¹⁰⁴⁵ to principal engineer) to complete five tasks under an observational, comparative
¹⁰⁴⁶ controlled study, 10 of which will (a) develop with the *new* framework, and the
¹⁰⁴⁷ other 10 will (b) develop with the *as-is/existing* documentation. From this, we
¹⁰⁴⁸ compare if the framework makes improvements by capturing metrics and recording
¹⁰⁴⁹ the observational sessions for qualitative analysis. We use visual modelling to
¹⁰⁵⁰ analyse the qualitative data using matrices [55], maps and networks [156] as these
¹⁰⁵¹ help illustrate any causal, temporal or contextual relationships that may exist to map
¹⁰⁵² out the developer's mindset and difference in approaching the two sets of designs of
¹⁰⁵³ the same tasks.

¹⁰⁵⁴ 3.5 Empirical Validity

¹⁰⁵⁵ In Section 3.2, we state that this study primarily adopts a critical theorist stance.
¹⁰⁵⁶ Critical theorists assess research quality by the utility of the knowledge gained [65].
¹⁰⁵⁷ Lau [133] established criteria on validating information systems research specifically
¹⁰⁵⁸ for action research unifying four dimensions of the research (conceptual foundation,
¹⁰⁵⁹ study design, research process, and role expectations) against 22 varying criteria.
¹⁰⁶⁰ We also partially adopt constructivism as we attempt to model the developer mindset
¹⁰⁶¹ rich in qualitative data, to which eight strategies identified by Creswell and Creswell
¹⁰⁶² [47] cover.
¹⁰⁶³ We identify possible threats to internal- (study design), external- (generality of

1064 results), and construct-validity (theoretical understanding) in the following sections,
1065 and describe how we mitigate these threats.

1066 **3.5.1 Threats to Internal Validity**

1067 **Hawthorne Effect**

1068 Observational field techniques involving participants often run a risk producing mis-
1069 aligned results from laboratory versus environmental (practical) conditions. This is
1070 commonly known as the Hawthorne effect [63, 199] and careful consideration of this
1071 effect must be reflected when designing our controlled observation (Section 3.4.3).
1072 We aim to carefully explain the purpose and protocol to research participants, en-
1073 couraging them to act as much as possible as in their practical conditions. We also
1074 encourage the ‘think-aloud’ to participants protocol to reinforce this. By highlight-
1075 ing the Hawthorne effect to them, we anticipate that participants will be aware of the
1076 condition, and therefore avoid doing things that do not reflect real-world action.

1077 **Misleading Statements in Interviews**

1078 Similarly, threats to the interview survey instrument exist where participants do not
1079 often report differences in behaviour from what they actually do in practice [219].
1080 We anticipate that conducting interviews in a semi-structured manner may assist in
1081 following up with unexpected statements (as opposed to structured interviews) and
1082 additionally corroborate findings using Jick’s concurrent triangulation method [111]
1083 to verify potentially misleading statements from participants with questionnaire
1084 results and observation findings.

1085 Participant Observation Accuracy

1086 Conducting participant observations is a skill that requires training. While every
1087 effort will be instilled to ensure all relevant observations are noted, it is impossible
1088 for a single observer to note every possible interaction that occurs in all observations
1089 made. Therefore, to validate the consistency of data collected, we may require
1090 rater agreement exercises [114] and we will likely use a form of recording device
1091 (with participant consent) to ensure all information is transcribed correctly in the
1092 interview.

1093 Unintended Interviewee Bias

1094 Interviewers should introduce the research by which participants are involved in by
1095 describing an expiation of the research being conducted. However, the amount of
1096 information described may impact the bias instilled on the interviewee. For example,
1097 if the participant does not understand the goals of the study or feel that they are of
1098 the ‘right target’, then it is likely that they may choose not to be involved in the
1099 study or give misleading answers. On the other hand, if interviewees are told too
1100 much information, then they may filter responses and leave out vital data that the
1101 interviewer may be interested in. To mitigate this, varying levels of information will
1102 be ‘tested’ against colleagues to determine the right level of how much information
1103 is divulged at the beginning of the interview.

1104 Poor Questionnaire Responses

1105 Unless significant inducements are offered, Singer et al. [219] report that a con-
1106 sistent response rate of 5% has been found in software engineering questionnaires
1107 distributed and in information systems the median response rates for surveys are 60%
1108 [16]. We observe that low response rates may adversely effect the findings of our
1109 survey, typically as software engineers find little time to do them. Tackling this is-
1110 sue can be resolved by carefully designing succinct, unambiguous and well-worded
1111 questions that we will verify against our colleagues and within the pilot study in
1112 A^2I^2 , wherein any adjustments made from the pilot study due to unexpected poor
1113 quality of the questionnaire will be reported and explained. We also adopt research
1114 conducted in the field of questionnaire design, such as ensuring all scales are worded
1115 with labels [128] or using a summating rating scale [224] to address a specific topic
1116 of interest if people are to make mistakes in their response or answer in different
1117 ways at different times. Similar effects exist to that above where what occurs in
1118 reality is not what is reflected in our results; we refer to our concurrent triangulation
1119 approach to gap this risk.

1120 3.5.2 Threats to External Validity**1121 Representative Sample Size**

1122 Our results must generalise by ensuring a representative subset of the target popu-
1123 lation is found. If results do not generalise, then all that is found is potentially of
1124 little more value than personal anecdote [122]. Therefore, designing a well-defined

¹¹²⁵ sample frame to determine which developers we wish to target is empirical. For this,
¹¹²⁶ we refer to Kitchenham and Pfleeger [122].

¹¹²⁷ **Student Cohorts**

¹¹²⁸ External validity is typically undermined when students are recruited in software en-
¹¹²⁹ gineering research, which is common practice [65]. Analytical argument is required
¹¹³⁰ to describe why results on students are reflective of results found on developers in
¹¹³¹ industry. Therefore, we anticipate that—through industry contacts at A²I²—we will
¹¹³² be able to contact developers in industry, thereby minimising our reliance to use
¹¹³³ students as participants.

¹¹³⁴ **Concurrent Triangulation Strategy**

¹¹³⁵ A drawback with the concurrent triangulation strategy is that multiple sources of
¹¹³⁶ data are concurrently collected within the same time. Collecting and analysing data
¹¹³⁷ *sequentially*, instead of concurrently, allows for time to analyse data between studies,
¹¹³⁸ thus allowing each analysis to adapt as more emerging results are explored. Easter-
¹¹³⁹ brook et al. [65] states that the challenge in this approach is that it may be difficult
¹¹⁴⁰ for researchers to compare results of multiple analyses or resolve contradictions that
¹¹⁴¹ begin to arise when this is performed concurrently. A mitigation strategy, should
¹¹⁴² this occur, would be to seek out further sources of evidence, or even re-conduct a
¹¹⁴³ follow-up study if time permits.

1144 3.5.3 Threats to Construct Validity**1145 Developer Informativeness**

1146 RH3 describes that if we improve the documentation of CIS APIs, then developers
1147 are more informed/educated in what they do. This therefore increases their pro-
1148 ductivity and the quality of the applications they build. However, the construct of
1149 ‘informativeness’ is difficult to capture with standalone metrics, and using simple
1150 quantitative metrics such as time taken to complete a task or lines of code to imple-
1151 ment it may not reflect that a developer is more ‘informed’. Therefore, we propose
1152 further investigation into understanding how to measure informativeness of software
1153 engineers to ensure that this construct validity does not impact our results too greatly.

1154

Part II

1155

Publications

CHAPTER 4

1156

1157

1158

Identifying Evolution in Computer Vision Services[†]

1159

1160 **Abstract** Recent advances in artificial intelligence (AI) and machine learning (ML), such
1161 as computer vision, are now available as intelligent services and their accessibility and
1162 simplicity is compelling. Multiple vendors now offer this technology as cloud services
1163 and developers want to leverage these advances to provide value to end-users. However,
1164 there is no firm investigation into the maintenance and evolution risks arising from use
1165 of these intelligent services; in particular, their behavioural consistency and transparency
1166 of their functionality. We evaluated the responses of three different intelligent services

[†]This chapter is originally based on A. Cummaudo, R. Vasa, J. Grundy, M. Abdelrazek, and A. Cain, “Losing Confidence in Quality: Unspoken Evolution of Computer Vision Services,” in *2019 IEEE International Conference on Software Maintenance and Evolution (ICSME)*. Cleveland, OH, USA: IEEE, Sep. 2019, pp. 333–342. Terminology has been updated to fit this thesis.

¹¹⁶⁷ (specifically computer vision) over 11 months using 3 different data sets, verifying responses
¹¹⁶⁸ against the respective documentation and assessing evolution risk. We found that there are:
¹¹⁶⁹ (1) inconsistencies in how these services behave; (2) evolution risk in the responses; and (3) a
¹¹⁷⁰ lack of clear communication that documents these risks and inconsistencies. We propose a
¹¹⁷¹ set of recommendations to both developers and intelligent service providers to inform risk
¹¹⁷² and assist maintainability.

¹¹⁷³ 4.1 Introduction

¹¹⁷⁴ The availability of intelligent services has made AI tooling accessible to software
¹¹⁷⁵ developers and promises a lower entry barrier for their utilisation. Consider state-
¹¹⁷⁶ of-the-art computer vision analysers, which require either manually training a deep-
¹¹⁷⁷ learning classifier, or selecting a pre-trained model and deploying these into an
¹¹⁷⁸ appropriate infrastructure. Either are laborious in time, and require non-trivial
¹¹⁷⁹ expertise along with a large data set when training or customisation is needed.
¹¹⁸⁰ In contrast, intelligent services providing computer vision (e.g., [266–275, 285–
¹¹⁸¹ 288]) abstract these complexities behind a web API call. This removes the need
¹¹⁸² to understand the complexities required of ML, and requires little more than the
¹¹⁸³ knowledge on how to use RESTful endpoints. The ubiquity of these services is
¹¹⁸⁴ exemplified through their rapid uptake in applications such as aiding the vision-
¹¹⁸⁵ impaired [51, 192].

¹¹⁸⁶ While intelligent services have seen quick adoption in industry, there has been
¹¹⁸⁷ little work that has considered the software quality perspective of the risks and

1188 impacts posed by using such services. In relation to this, there are three main chal-
1189 lenges: (1) incorporating stochastic algorithms into software that has traditionally
1190 been deterministic; (2) the general lack of transparency associated with the ML
1191 models; and (3) communicating to application developers.

1192 ML typically involves use of statistical techniques that yield components with
1193 a non-deterministic external behaviour; that is, for the same given input, different
1194 outcomes may result. However, developers, in general, are used to libraries and small
1195 components behaving predictably, while systems that rely on ML techniques work
1196 on confidence intervals¹ and probabilities. For example, the developer's mindset
1197 suggests that an image of a border collie—if sent to three intelligent computer
1198 vision services—would return the label ‘dog’ consistently with time regardless of
1199 which service is used. However, one service may yield the specific dog breed,
1200 ‘border collie’, another service may yield a permutation of that breed, ‘collie’, and
1201 another may yield broader results, such as ‘animal’; each with results of varying
1202 confidence values.² Furthermore, the third service may evolve with time, and
1203 thus learn that the ‘animal’ is actually a ‘dog’ or even a ‘collie’. The outcomes
1204 are thus behaviourally inconsistent between services providing conceptually similar
1205 functionality. As a thought exercise, consider if the sub-string function were created
1206 using ML techniques—it would perform its operation with a confidence where the
1207 expected outcome and the AI inferred output match as a *probability*, rather than a

¹Varied terminology used here. Probability, confidence, accuracy and score may all be used interchangeably.

²Indeed, we have observed this phenomenon using a picture of a border collie sent to various computer vision services.

¹²⁰⁸ deterministic (constant) outcome. How would this affect the developers' approach
¹²⁰⁹ to using such a function? Would they actively take into consideration the non-
¹²¹⁰ deterministic nature of the result?

¹²¹¹ Myriad software quality models and SE practices advocate maintainability and
¹²¹² reliability as primary characteristics; stability, testability, fault tolerance, changeabil-
¹²¹³ ity and maturity are all concerns for quality in software components [97, 188, 223]
¹²¹⁴ and one must factor these in with consideration to software evolution challenges [54,
¹²¹⁵ 87, 153, 237, 239]. However, the effect this non-deterministic behaviour has on
¹²¹⁶ quality when masked behind an intelligent service is still under-explored to date in
¹²¹⁷ SE literature, to our knowledge. Where software depends on intelligent services to
¹²¹⁸ achieve functionality, these quality characteristics may not be achieved, and devel-
¹²¹⁹ opers need to be wary of the unintended side effects and inconsistency that exists
¹²²⁰ when using non-deterministic components. A computer vision service may encap-
¹²²¹ sulate deep-learning strategies or stochastic methods to perform image analysis, but
¹²²² developers are more likely to approach intelligent services with a mindset that antic-
¹²²³ ipates consistency. Although the documentation does hint at this non-deterministic
¹²²⁴ behaviour (i.e., the descriptions of ‘confidence’ in various computer vision services
¹²²⁵ suggest they are not always confident, and thus not deterministic [289–291]), the
¹²²⁶ integration mechanisms offered by popular vendors do not seem to fully expose the
¹²²⁷ nuances, and developers are not yet familiar with the trade-offs.

¹²²⁸ Do popular computer vision services, as they currently stand, offer consistent
¹²²⁹ behaviour, and if not, how is this conveyed to developers (if it is at all)? If computer
¹²³⁰ vision services are to be used in production services, do they ensure quality un-

1231 der rigorous Software Quality Assurance (SQA) frameworks [97]? What evolution
1232 risk [54, 87, 153, 239] do they pose if these services change? To our knowledge,
1233 few studies have been conducted to investigate these claims. This paper assesses
1234 the consistency, evolution risk and consequent maintenance issues that may arise
1235 when developers use intelligent services. We introduce a motivating example in Sec-
1236 tion 4.2, discussing related work and our methodology in Sections 4.3 and 4.4. We
1237 present and interpret our findings in Section 4.5. We argue with quantified evidence
1238 that these intelligent services can only be considered with a mature appreciation of
1239 risks, and we make a set of recommendations in Section 4.6.

1240 **4.2 Motivating Example**

1241 Consider Rosa, a software developer, who wants to develop a social media photo-
1242 sharing mobile app that analyses her and her friends photos on Android and iOS.
1243 Rosa wants the app to categorise photos into scenes (e.g., day vs. night, outdoors
1244 vs. indoors), generate brief descriptions of each photo, and catalogue photos of her
1245 friends as well as common objects (e.g., all photos with a dog, all photos on the
1246 beach).

1247 Rather than building a computer vision engine from scratch, Rosa thinks she can
1248 achieve this using one of the popular computer vision services (e.g., [266–275, 285–
1249 288]). However, Rosa comes from a typical software engineering background with
1250 limited knowledge of the underlying deep-learning techniques and implementations
1251 as currently used in computer vision. Not unexpectedly, she internalises a mindset

1252 of how such services work and behave based on her experience of using software
1253 libraries offered by various SDKs. This mindset assumes that different cloud vendor
1254 image processing APIs more-or-less provide similar functionality, with only minor
1255 variations. For example, cloud object storage for Amazon S3 is both conceptually
1256 and behaviourally very similar to that of Google Cloud Storage or Azure Storage.

1257 Rosa assumes the computer vision services of these platforms will, therefore, likely
1258 be very similar. Similarly, consider the string libraries Rosa will use for the app.

1259 The conceptual and behavioural similarities are consistent; a string library in Java
1260 (Android) is conceptually very similar to the string library she will use in Swift
1261 (iOS), and likewise both behave similarly by providing the same results for their
1262 respective sub-string functionality. However, **unlike the cloud storage and string**
1263 **libraries, different computer vision services often present conceptually similar**
1264 **functionality but are behaviourally very different.** Intelligent service vendors
1265 also hide the depth of knowledge needed to use these effectively—for instance,
1266 the training data set and ontologies used to create these services are hidden in
1267 the documentation. Thus, Rosa isn't even exposed to this knowledge as she reads
1268 through the documentation of the providers and, thus, Rosa makes the following
1269 assumptions:

- 1270 • **“I think the responses will be consistent amongst these computer vision**
1271 **services.”** When Rosa uploads a photo of a dog, she would expect them all to
1272 respond with ‘dog’. If Rosa decides to switch which service she is using, she
1273 expects the ontologies to be compatible (all computer vision services *surely*
1274 return dog for the same image) and therefore she can expect to plug-in a

1275 different service should she feel like it making only minor code modifications
1276 such as which endpoints she is relying on.

1277 • “**I think the responses will be constant with time.**” When Rosa uploads the
1278 photo of a dog for testing, she expects the response to be the same in 10 weeks
1279 time once her app is in production. Hence, in 10 weeks, the same photo of the
1280 dog should return the same label.

1281 4.3 Related Work

1282 If we were to view computer vision services through the lenses of an SQA framework,
1283 robustness, consistency, and maintainability often feature as quality attributes in
1284 myriad software quality models (e.g., [106]). Software quality is determined from
1285 two key dimensions: (1) in the evaluation of the end-product (external quality) and
1286 (2) the assurances in the development processes (internal quality) [188]. We discuss
1287 both perspectives of quality within the context of our work in this section.

1288 4.3.1 External Quality

1289 **Robustness for safety-critical applications** A typical focus of recent work has
1290 been to investigate the robustness of deep-learning within computer vision tech-
1291 nique implementation, thereby informing the effectiveness in the context of the
1292 end-product. The common method for this has been via the use of adversarial exam-
1293 ples [231], where input images are slightly perturbed to maximise prediction error
1294 but are still interpretable to humans.

1295 Google Cloud Vision, for instance, fails to correctly classify adversarial examples

¹²⁹⁶ when noise is added to the original images [98]. Rosenfeld et al. [206] illustrated
¹²⁹⁷ that inserting synthetic foreign objects to input images (e.g., a cartoon elephant)
¹²⁹⁸ can completely alter classification output. Wang et al. [245] performed similar
¹²⁹⁹ attacks on a transfer-learning approach of facial recognition by modifying pixels of
¹³⁰⁰ a celebrity’s face to be recognised as a completely different celebrity, all while still
¹³⁰¹ retaining the same human-interpretable original celebrity. Su et al. [225] used the
¹³⁰² ImageNet database to show that 41.22% of images drop in confidence when just a
¹³⁰³ *single pixel* is changed in the input image; and similarly, Eykholt et al. [69] recently
¹³⁰⁴ showed similar results that made a CNN interpret a stop road-sign (with mimicked
¹³⁰⁵ graffiti) as a 45mph speed limit sign.

¹³⁰⁶ The results suggest that current state-of-the-art computer vision techniques may
¹³⁰⁷ not be robust enough for safety critical applications as they do not handle intentional
¹³⁰⁸ or unintentional adversarial attacks. Moreover, as such adversarial examples exist in
¹³⁰⁹ the physical world [68, 129], “the natural world may be adversarial enough” [184]
¹³¹⁰ to fool AI software. Though some limitations and guidelines have been explored in
¹³¹¹ this area, the perspective of *intelligent* services is yet to be considered and specific
¹³¹² guidelines do not yet exist when using computer vision services.

¹³¹³ **Testing strategies in ML applications** Although much work applies ML tech-
¹³¹⁴ niques to automate testing strategies, there is only a growing emphasis that considers
¹³¹⁵ this in the opposite sense; that is, testing to ensure the ML product works correctly.
¹³¹⁶ There are few reliable test oracles that ensure if an ML has been implemented to
¹³¹⁷ serve its algorithm and use case purposefully; indeed, “the non-deterministic nature

1318 of many training algorithms makes testing of models even more challenging” [9].

1319 Murphy et al. [160] proposed a SE-based testing approach on ML ranking algorithms

1320 to evaluate the ‘correctness’ of the implementation on a real-world data set and prob-

1321 lem domain, whereby discrepancies were found from the formal mathematical proofs

1322 of the ML algorithm and the implementation.

1323 Recently, Braiek and Khomh [27] conducted a comprehensive review of testing

1324 strategies in ML software, proposing several research directions and recommenda-

1325 tions in how best to apply SE testing practices in ML programs. However, much

1326 of the area of this work specifically targets ML engineers, and not application de-

1327 velopers. Little has been investigated on how application developers perceive and

1328 understand ML concepts, given a lack of formal training; we note that other testing

1329 strategies and frameworks proposed (e.g., [30, 159, 166]) are targeted chiefly to the

1330 ML engineer, and not the application developer.

1331 However, Arpteg et al. [9] recently demonstrated (using real-world ML projects)

1332 the developmental challenges posed to developers, particularly those that arise when

1333 there is a lack of transparency on the models used and how to troubleshoot ML

1334 frameworks using traditional SE debugging tools. This said, there is no further

1335 investigations into challenges when using the higher, ‘ML friendly’ layers (e.g.,

1336 intelligent services) of the ‘machine learning spectrum’ [173], rather than the ‘lower

1337 layers’ consisting of existing ML frameworks and algorithms targeted toward the

1338 ML community.

1339 4.3.2 Internal Quality

1340 **Quality metrics for cloud services** Computer vision services are based on cloud
1341 computing fundamentals under a subset of the Platform as a Service (PaaS) model.

1342 There has been work in the evaluation of PaaS in terms of quality attributes [80]:
1343 these attributes are exposed using Service Level Agreements (SLA) between vendors
1344 and customers, and customers denote their demanded Quality of Service (QoS) to
1345 ensure the cloud services adhere to measurable KPI attributes.

1346 Although, popular services, such as cloud object storage, come with strong QoS
1347 agreement, to date intelligent services do not come with deep assurances around
1348 their performance and responses, but do offer uptime guarantees. For example,
1349 how can Rosa demand a QoS that ensures all photos of dogs uploaded to her app
1350 guarantee the specific dog breeds are returned so that users can look up their other
1351 friend's 'border collie's? If dog breeds are returned, what ontologies exist for breeds?
1352 Are they consistent with each other, or shortened? ('Collie' versus 'border collie';
1353 'staffy' versus 'staffordshire bull terrier')? For some applications, these unstated
1354 QoS metrics specific to the ML service may have significant legal ramifications.

1355 **Web service documentation and documenting ML** From the *developer's* per-
1356 spective, little has been achieved to assess intelligent service quality or assure quality
1357 of these computer vision services. Web services and their interfaces (APIs) are the
1358 bridge between developers' needs and the software components [8]; therefore, assess-
1359 ing such computer vision services from the quality of their APIs is thereby directly
1360 related to the development quality [125]. Good APIs should be intuitive and require

1361 less documentation browsing [186], thereby increasing productivity. Conversely,
1362 poor APIs that are hard to understand and work with reduce developer productivity,
1363 thereby reducing product quality. This typically leads to developers congregating on
1364 forums such as Stack Overflow, leading to a repository of unstructured knowledge
1365 likely to concern API design [247]. The consequences of addressing these concerns
1366 in development leads to a higher demand in technical support (as measured in [96])
1367 that, ultimately, causes the maintenance to be far more expensive, a phenomenon
1368 widely known in software engineering economics [24]. Rosa, for instance, isn't
1369 aware of technical ML concepts; if she cannot reason about what search results are
1370 relevant when browsing the service and understanding functionality, her productivity
1371 is significantly decreased. Conceptual understanding is critical for using APIs, as
1372 demonstrated by Ko and Riche, and the effects of maintenance this may have in the
1373 future of her application is unknown.

1374 Recent attempts to document attributes and characteristics on ML models have
1375 been proposed. Model cards were introduced by Mitchell et al. [157] to describe how
1376 particular models were trained and benchmarked, thereby assisting users to reason
1377 if the model is right for their purposes and if it can achieve its stated outcomes.
1378 Gebru et al. [84] also proposed datasheets, a standardised documentation format to
1379 describe the need for a particular data set, the information contained within it and
1380 what scenarios it should be used for, including legal or ethical concerns.

1381 However, while target audiences for these documents may be of a more technical
1382 AI level (i.e., the ML engineer), there is still no standardised communication format
1383 for application developers to reason about using particular intelligent services, and

1384 the ramifications this may have on the applications they write is not fully conveyed.

1385 Hence, our work is focused on the application developer perspective.

1386 4.4 Method

1387 This study organically evolved by observing phenomena surrounding computer vi-
1388 sion services by assessing both their documentation and responses. We adopted a
1389 mixed methods approach, performing both qualitative and quantitative data collec-
1390 tion on these two key aspects by using documentary research methods for inspecting
1391 the documentation and structured observations to quantitatively analyse the results
1392 over time. This, ultimately, helped us shape the following research hypotheses which
1393 this paper addresses:

1394 [RH1] Computer vision services do not respond with consistent outputs
1395 between services, given the same input image.

1396 [RH2] The responses from computer vision services are non-deterministic
1397 and evolving, and the same service can change its top-most response
1398 over time given the same input image.

1399 [RH3] Computer vision services do not effectively communicate this evolu-
1400 tion and instability, introducing risk into engineering these systems.

1401 We conducted two experiments to address these hypotheses against three popular
1402 computer vision services: AWS Rekognition [268], Google Cloud Vision [266],
1403 Azure Computer Vision [267]. Specifically, we targeted the AWS DetectLabels
1404 endpoint [290], the Google Cloud Vision annotate:images endpoint [289] and

1405 Azure's analyze endpoint [291]. For the remainder of this paper, we de-identify
1406 our selected computer vision services by labelling them as services A, B and C
1407 but do not reveal mapping to prevent any implicit bias. Our selection criteria for
1408 using these particular three services are based on the weight behind each service
1409 provider given their prominence in the industry (Amazon, Google and Microsoft),
1410 the ubiquity of their hosting cloud platforms as industry leaders of cloud computing
1411 (i.e., AWS, Google Cloud and Azure), being in the top three most adopted cloud
1412 vendors in enterprise applications in 2018 [197] and the consistent popularity of
1413 discussion amongst developers in developer communities such as Stack Overflow.
1414 While we choose these particular cloud computer vision services, we acknowledge
1415 that similar services [269–275] also exist, including other popular services used in
1416 Asia [285–288] (some offering 3D image analysis [292]). We reflect on the impacts
1417 this has to our study design in Section 4.7.

1418 Our study involved an 11-month longitudinal study which consisted of two 13
1419 week and 17 week experiments from April to August 2018 and November 2018 to
1420 March 2019, respectively. Our investigation into documentation occurred on August
1421 28 2018. In total, we assessed the services with three data sets; we first ran a pilot
1422 study using a smaller pool of 30 images to confirm the end-points remain stable,
1423 re-running the study with a larger pool of images of 1,650 and 5,000 images. Our
1424 selection criteria for these three data sets were that the images had to have varying
1425 objects, taken in various scenes and various times. Images also needed to contain
1426 disparate objects. Our small data set was sourced by the first author by taking photos
1427 of random scenes in an afternoon, whilst our second data set was sourced from

Table 4.1: Characteristics of our datasets and responses.

Data set	Small	Large	COCOVal17
# Images/data set	30	1,650	5000
# Unique labels found	307	3506	4507
Number of snapshots	9	22	22
Avg. days b/n requests	12 Days	8 Days	8 Days

¹⁴²⁸ various members of our research group from their personal photo libraries. We also
¹⁴²⁹ wanted to include a data set that was publicly available prior to running our study,
¹⁴³⁰ so for this data set we chose the COCO 2017 validation data set [138]. We have
¹⁴³¹ made our other two data sets available online ([293]). We collected results and their
¹⁴³² responses from each service’s API endpoint using a python script [294] that sent
¹⁴³³ requests to each service periodically via cron jobs. Table 4.1 summarises various
¹⁴³⁴ characteristics about the data sets used in these experiments.

¹⁴³⁵ We then performed quantitative analyses on each response’s labels, ensuring all
¹⁴³⁶ labels were lowercased as case changed for services A and C over the evaluation
¹⁴³⁷ period. To derive at the consistency of responses for each image, we considered only
¹⁴³⁸ the ‘top’ labels per image for each service and data set. That is, for the same image i
¹⁴³⁹ over all images in data set D where $i \in D$ and over the three services, the top labels
¹⁴⁴⁰ per image (T_i) of all labels per image L_i (i.e., $T_i \subseteq L_i$) is that where the respective
¹⁴⁴¹ label’s confidences are consistently the highest of all labels returned. Typically, the
¹⁴⁴² top labels returned is a set containing only one element—that is, only one unique
¹⁴⁴³ label consistently returned with the highest label ($|T_i| = 1$)—however there are cases

1444 where the top labels contains multiple elements as their respective confidences are
1445 *equal* ($|T_i| > 1$).

1446 We measure response consistency under 6 aspects:

1447 **(1) Consistency of the top label between each service.** Where the same image of,
1448 for example, a dog is sent to the three services, the top label for service A may
1449 be ‘animal’, B ‘canine’ and C ‘animal’. Therefore, service B is inconsistent.

1450 **(2) Semantic consistency of the top labels.** Where a service has returned multi-
1451 ple top labels ($|T_i| > 1$), there may lie semantic differences in what the service
1452 thinks the image best represents. Therefore, there is conceptual inconsistency
1453 in the top labels for a service even when the confidences are equal.

1454 **(3) Consistency of the top label’s confidence per service.** The top label for
1455 an image does not guarantee a high confidence. Therefore, there may be
1456 inconsistencies in how confident the top labels for all images in a service is.

1457 **(4) Consistency of confidence in the intersecting top label between each ser-**
1458 **vice.** The spread of a top intersecting label, e.g. ‘cat’, may not have the same
1459 confidences per service even when all three services agree that ‘cat’ is the top
1460 label. Therefore, there is inconsistency in the confidences of a top label even
1461 where all three services agree.

1462 **(5) Consistency of the top label over time.** Given an image, the top label in one
1463 week may differ from the top label the following week. Therefore, there is
1464 inconsistency in the top label itself due to model evolution.

1465 **(6) Consistency of the top label’s confidence over time.** The top label of an



Figure 4.1: The only consistent label for the above image is ‘people’ for services C and B.

The top label for A is ‘conversation’ and this label is not registered amongst the other two services.

¹⁴⁶⁶ image may remain static from one week to the next for the same service, but
¹⁴⁶⁷ its confidence values may change with time. Therefore, there is inconsistency
¹⁴⁶⁸ in the top label’s confidence due to model evolution.

¹⁴⁶⁹ For the above aspects of consistency, we calculated the spread of variation for the
¹⁴⁷⁰ top label’s confidences of each service for every 1 percent point; that is, the frequency
¹⁴⁷¹ of top label confidences within 100–99%, 99–98% etc. The consistency of top label’s
¹⁴⁷² and their confidences between each service was determined by intersecting the labels
¹⁴⁷³ of each service per image and grouping the intersecting label’s confidences together.
¹⁴⁷⁴ This allowed us to determine relevant probability distributions. For reproducibility,
¹⁴⁷⁵ all quantitative analysis is available online [295].

¹⁴⁷⁶ 4.5 Findings

¹⁴⁷⁷ 4.5.1 Consistency of top labels

Table 4.2: Ratio of the top labels (to images) that intersect in each data set for each permutation of service.

Service	Small	Large	COCOVal17	μ	σ
$A \cap B \cap C$	3.33%	2.73%	4.68%	2.75%	0.0100
$A \cap B$	6.67%	11.27%	12.26%	10.07%	0.0299
$A \cap C$	20.00%	13.94%	17.28%	17.07%	0.0304
$B \cap C$	6.67%	12.97%	20.90%	13.51%	0.0713

1478 **Consistency across services** Table 4.2 presents the consistency of the top labels
1479 between data sets, as measured by the cardinality of the intersection of all three ser-
1480 vices’ set of top labels divided by the number of images per data set. A combination
1481 of services present varied overlaps in their top labels; services A and C provide the
1482 best overlap for all three data sets, however the intersection of all three irrespective
1483 of data sets is low.

1484 The implication here is that, without semantic comparison (see Section 4.7),
1485 service vendors are not ‘plug-and-play’. If Rosa uploaded the sample images in
1486 this paper to her application to all services, she would find that only Figure 4.1
1487 responds with ‘person’ for services B and C in their respective set of top labels.
1488 However, if she decides to then adopt service A, then Figure 4.1’s top label becomes
1489 ‘conversation’; the ‘person’ label does not appear within the top 15 labels for service
1490 A and, conversely, the ‘conversation’ label does not appear in the other services top
1491 15.

1492 Should she decide if the performance of a particular service isn’t to her needs,



Figure 4.2: *Left:* The top labels for each service do not intersect, with each having a varied ontology: $T_i = \{ A = \{ \text{'black'} \}, B = \{ \text{'indoor'} \}, C = \{ \text{'slide'}, \text{'toy'} \} \}$. (Service C returns *both* ‘slide’ and ‘toy’ with equal confidence.) *Right:* The top labels for each service focus on disparate subjects in the image: $T_i = \{ A = \{ \text{'carrot'} \}, B = \{ \text{'indoor'} \}, C = \{ \text{'spoon'} \} \}$.

¹⁴⁹³ then the vocabulary used for these labels becomes inconsistent for all other images;
¹⁴⁹⁴ that is, the top label sets per service for Figure 4.2a shows no intersection at all.
¹⁴⁹⁵ Furthermore, the part of the image each service focuses on may not be consistent
¹⁴⁹⁶ for their top labels; in Figure 4.2b, service A’s top label focuses on the vegetable
¹⁴⁹⁷ (‘carrot’), service C focuses on the ‘spoon’, while service B’s focus is that the image
¹⁴⁹⁸ is ‘indoor’s. It is interesting to note that service B focuses on the scene matter
¹⁴⁹⁹ (indoors) rather than the subject matter. (Furthermore, we do not actually know if
¹⁵⁰⁰ the image in Figure 4.2b was taken indoors.)

¹⁵⁰¹ Hence, developers should ensure that the vocabulary used by a particular service
¹⁵⁰² is right for them before implementation. As each service does not work to the
¹⁵⁰³ same standardised model, trained with disparate training data, and tuned differently,
¹⁵⁰⁴ results will differ despite the same input. This is unlike deterministic systems: for
¹⁵⁰⁵ example, switching from AWS Object Storage to Google Cloud Object storage will



Figure 4.3: *Left:* Service C is 98.49% confident of the following labels: { ‘beverage’, ‘chocolate’, ‘cup’, ‘dessert’, ‘drink’, ‘food’, ‘hot chocolate’ }. However, it is up to the developer to decide which label to persist with as all are returned. *Right:* Service B persistently returns a top label set of { ‘book’, ‘several’ }. Both are semantically correct for the image, but disparate in what the label is to describe.

1506 conceptually provide the same output (storing files) for the same input (uploading
1507 files). However, computer vision services do not agree on the top label for images,
1508 and therefore developers are likely to be vendor locked, making changes between
1509 services non-trivial.

Semantic consistency where $|T_i| > 1$ Service C returns two top labels for Figure 4.2a; ‘slide’ and ‘toy’. More than one top label is typically returned in service C (80.00%, 56.97%, and 81.66% of all images for all three data sets, respectively) though this also occurs in B in the large (4.97% of all images) and COCOVal17 data sets (2.38%). Semantic inconsistencies of what this label conceptually represents becomes a concern as these labels have confidences of *equal highest* consistency. Thus, some services are inconsistent in themselves and cannot give a guaranteed answer of what exists in an image; services C and B have multiple top labels, but the

1518 respective services cannot ‘agree’ on what the top label actually is. In Figure 4.3a,
1519 service C presents a reasonably high confidence for the set of 7 top labels it returns,
1520 however there is too much diversity ranging from a ‘hot chocolate’ to the hypernym
1521 ‘food’. Both are technically correct, but it is up to the developer to decide the level
1522 of hypernymy to label the image as. We also observe a similar effect in Figure 4.3b,
1523 where the image is labelled with both the subject matter and the number of subjects
1524 per image.

1525 Thus, a taxonomy of ontologies is unknown; if a ‘border collie’ is detected in
1526 an image, does this imply the hypernym ‘dog’ is detected, and then ‘mammal’, then
1527 ‘animal’, then ‘object’? Only service B documents a taxonomy for capturing what
1528 level of scope is desired, providing what it calls the ‘86-category’ concept as found
1529 in its how-to guide:

1530 “*Identify and categorize an entire image, using a category taxonomy
1531 with parent/child hereditary hierarchies. Categories can be used alone,
1532 or with our new tagging models.*” [296]

1533 Thus, even if Rosa implemented conceptual similarity analysis for the image, the
1534 top label set may not provide sufficient information to derive at a conclusive answer,
1535 and if simply relying on only one label in this set, information such as the duplicity
1536 of objects (e.g., ‘several’ in Figure 4.3b) may be missed.

1537 **4.5.2 Consistency of confidence**

1538 **Consistency of top label’s confidence** In Figure 4.4, we see that there is high
1539 probability that top labels have high confidences for all services. In summary, one in

Table 4.3: Ratio of the top labels (to images) that remained the top label but changed confidence values between intervals.

Service	Small	Large	COCOVal17	$\mu(\delta_c)$	$\sigma(\delta_c)$	Median(δ_c)	Range(δ_c)
A	53.33%	59.19%	44.92%	9.62e-8	6.84e-8	5.96e-8	[5.96e-8, 6.56e-7]
B	0.00%	0.00%	0.02%	-	-	-	-
C	33.33%	41.36%	15.60%	5.35e-7	8.76e-7	3.05e-7	[1.27e-7, 1.13e-5]

¹⁵⁴⁰ nine images uploaded to any service will return a top label confident to at least 97%.

¹⁵⁴¹ However, there is higher probability for service A returning a lower confidence,

¹⁵⁴² followed by B. The best performing service is C, with 90% of requests having a top

¹⁵⁴³ label confident to $\gtrapprox 95\%$, when compared to $\gtrapprox 87\%$ and $\gtrapprox 93\%$ for services A and

¹⁵⁴⁴ B, respectively.

¹⁵⁴⁵ Therefore, Rosa could generally expect that the top labels she receives in her

¹⁵⁴⁶ images do have high confidence. That is, each service will return a top label that

¹⁵⁴⁷ they are confident about. This result is expected, considering that the ‘top’ label

¹⁵⁴⁸ is measured by the highest confidence, though it is interesting to note that some

¹⁵⁴⁹ services are generally more confident than others in what they present back to users.

¹⁵⁵⁰ **Consistency of intersecting top label’s confidence** Even where all three services

¹⁵⁵¹ do agree on a set of top labels, the disparity of how much they agree by is still of

¹⁵⁵² importance. Just because three services agree that an image contains consistent top

¹⁵⁵³ labels, they do not always have a small spread of confidence. In Figure 4.6, the three

¹⁵⁵⁴ services agree with $\sigma = 0.277$, significantly larger than that of all images in general

¹⁵⁵⁵ $\sigma = 0.0831$. Figure 4.5 displays the cumulative distribution of all intersecting top

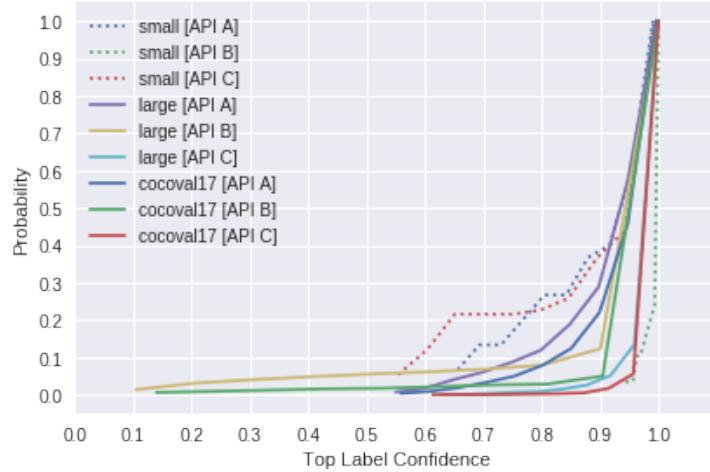


Figure 4.4: Cumulative distribution of the top labels' confidences. One in nine images return a top label(s) confident to $\gtrsim 97\%$, though there is a wider distribution for service A.

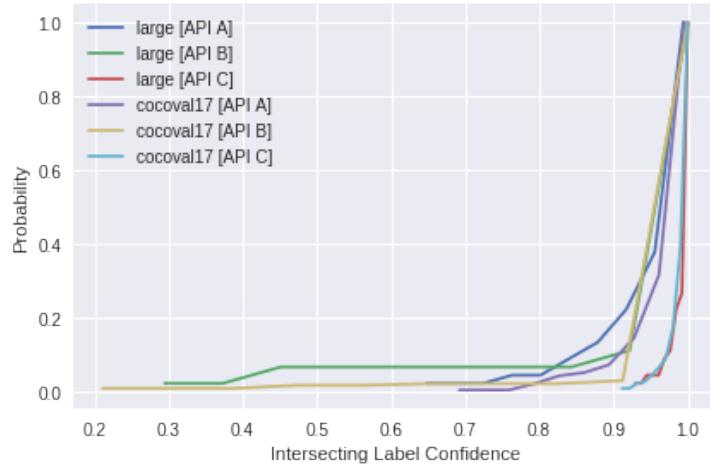


Figure 4.5: Cumulative distribution of intersecting top labels' confidences. The small data set is intentionally removed due to low intersections of labels (see Table 4.2).



Figure 4.6: All three services agree the top label for the above image is ‘food’, but the confidences to which they agree by vary significantly. Service C is most confident to 94.93% (in addition with the label ‘bread’); service A is the second most confident to 84.32%; service B is the least confident with 41.39%.

1556 labels’ confidence values, presenting slightly similar results to that of Figure 4.4.

1557 4.5.3 Evolution risk

1558 **Label Stability** Generally, the top label(s) did not evolve in the evaluation period.

1559 16.19% and 5.85% of images did change their top label(s) in the Large and COCO-

1560 Val17 data sets in service A. Thus, top labels are stable but not guaranteed to be

1561 constant.

1562 **Confidence Stability** Similarly, where the top label(s) remained the same from

1563 one interval to the next, the confidence values were stable. Table 4.3 displays the

1564 proportion of images that changed their top label’s confidence values with various

1565 statistics on the confidence deltas between snapshots (δ_c). However, this delta is so

1566 minuscule that we attribute such changes to statistical noise.

4.6 Recommendations

4.6.1 Recommendations for intelligent service users

1568 Test with a representative ontology for the particular use case Rosa should ensure that in her testing strategies for the app she develops, there is an ontology focus for the types of vocabulary that are returned. Additionally, we noted that there was a sudden change in case for services A and C; for all comparative purposes of labels, each label should be lower-cased.

1574 Incorporate a specialised intelligent service testing methodology into the development lifecycle Rosa can utilise the different aspects of consistency as outlined in this paper as part of her quality strategy. To ensure results are correct over time, we recommend developers create a representative data set of the intended application's data set and evaluate these changes against their chosen service frequently. This will help identify when changes, if any, have occurred if vendors do not provide a line of communication when this occurs.

1581 Intelligent services are not ‘plug-and-play’ Rosa will be locked into whichever vendor she chooses as there is inherent inconsistency between these services in both the vocabulary and ontologies that they use. We have demonstrated that very few services overlap in their vocabularies, chiefly because they are still in early development and there is yet to be an established, standardised vocabulary that can be shared amongst the different vendors. Issues such as those shown in Section 4.5.1 can therefore be avoided.

1588 Throughout this work, we observed that the terminologies used by the vari-
1589 ous vendors are different. Documentation was studied, and we note that there is
1590 inconsistency between the ways techniques are described to users. We note the
1591 disparity between the terms ‘detection’, ‘recognition’, ‘localisation’ and ‘analysis’.
1592 This applies chiefly to object- and facial-related techniques. Detection applies to
1593 facial detection, which gives bounding box coordinates around all faces in an image.
1594 Similarly, localisation applies the same methodology to disparate objects in an image
1595 and labels them. In the context of facial ‘recognition’, this term implies that a face
1596 is *recognised* against a known set of faces. Lastly, ‘analysis’ applies in the context
1597 of facial analysis (gender, eye colour, expression etc.); there does not exist a similar
1598 analysis technique on objects.

1599 We notice similar patterns with object ‘tagging’, ‘detection’ and ‘labelling’.
1600 Service A uses ‘Entity Detection’ for object categorisation, service B uses ‘Image
1601 Tagging’, and service C uses the term ‘Detect Labels’ : conceptually, these provide
1602 the same functionality but the lack of consistency used between all three providers is
1603 concerning and leaves room for confusion with developers during any comparative
1604 analyses. Rosa may find that she wants to label her images into day/night scenes, but
1605 this in turn means the ‘labelling’ of varying objects. There is therefore no consistent
1606 standards to use the same terminology for the same concepts, as there are in other
1607 developer areas (such as Web Development).

1608 **Avoid use in safety-critical systems** We have demonstrated in this paper that both
1609 labels and confidences are stable but not constant; there is still an evolution risk posed

1610 to developers that may cause unknown consequences in applications dependent on
1611 these computer vision services. Developers should avoid their use in safety critical
1612 systems due to the lack of visible changes.

1613 **4.6.2 Recommendations for intelligent service providers**

1614 **Improve the documentation** Rosa does not know that service A returns back
1615 ‘carrot’ for its top response, with service C returning ‘spoon’ (Figure 4.2b). She
1616 is unable to tell the service’s API where to focus on the image. Moreover, how
1617 can she toggle the level of specificity in her results? She is frustrated that service
1618 C can detect ‘chocolate’, ‘food’ and also ‘beverage’ all as the same top label in
1619 Figure 4.3a: what label is she to choose when the service is meant to do so for her,
1620 and how does she get around this? Thus, we recommend vendors to improve the
1621 documentation of services by making known the boundary set of the training data
1622 used for the algorithms. By making such information publicly available, developers
1623 would be able to review the service’s specificity for their intended use case (e.g.,
1624 maybe Rosa is satisfied her app can catalogue ‘food’ together, and in fact does not
1625 want specific types of foods (‘hot chocolate’) catalogued). We also recommend that
1626 vendors publish usage guidelines that include details of priors and how to
1627 evaluate the specific service results.

1628 Furthermore, we did not observe that the vendors documented how some images
1629 may respond with multiple labels of the exact same confidence value. It is not clear
1630 from the documentation that response objects can have duplicate top values, and
1631 tutorials and examples provided by the vendors do not consider this possibility. It

1632 is therefore left to the developer to decide which label from this top set of labels
1633 best suits for their particular use case; the documentation should describe that a rule
1634 engine may need to be added in the developer’s application to verify responses. The
1635 implications this would have on maintenance would be significant.

1636 **Improve versioning** We recommend introducing a versioning system so that a
1637 model can be used from a specific date in production systems: when Rosa tests her
1638 app today, she would like the service to remain *static* the same for when her app is
1639 deployed in production tomorrow. Thus, in a request made to the vendor, Rosa could
1640 specify what date she ran her app’s QA testing on so that she knows that henceforth
1641 these model changes will not affect her app.

1642 **Improve Metadata in Response** Much of the information in these services is
1643 reduced to a single confidence value within the response object, and the details about
1644 training data and the internal AI architecture remains unknown; little metadata is
1645 provided back to developers that encompass such detail. Early work into model
1646 cards and datasheets [84, 157] suggests more can be done to document attributes
1647 about ML systems, however at a minimum from our work, we recommend including
1648 a reference point via the form of an additional identifier. This identifier must also
1649 permit the developers to submit the identifier to another API endpoint should the
1650 developer wish to find further characteristics about the AI empowering the intelligent
1651 service, reinforcing the need for those presented in model cards and datasheets. For
1652 example, if Rosa sends this identifier she receives in the response object to the
1653 intelligent service descriptor API, she could find out additional information such as

¹⁶⁵⁴ the version number or date when the model was trained, thereby resolving potential
¹⁶⁵⁵ evolution risk, and/or the ontology of labels.

¹⁶⁵⁶ **Apply constraints for predictions on all inputs** In this study, we used some
¹⁶⁵⁷ images with intentionally disparate, and noisy objects. If services are not fully
¹⁶⁵⁸ confident in the responses they give back, a form of customised error message
¹⁶⁵⁹ should be returned. For example, if Rosa uploads an image of 10 various objects
¹⁶⁶⁰ on a table, rather than returning a list of top labels with varying confidences, it may
¹⁶⁶¹ be best to return a ‘too many objects’ exception. Similarly, if Rosa uploads a photo
¹⁶⁶² that the model has had no priors on, it might be useful to return an ‘unknown object’
¹⁶⁶³ exception than to return a label it has no confidence of. We do however acknowledge
¹⁶⁶⁴ that current state of the art computer vision techniques may have limits in what they
¹⁶⁶⁵ can and cannot detect, but this limitation can be exposed in the documentation to the
¹⁶⁶⁶ developers.

¹⁶⁶⁷ A further example is sending a one pixel image to the service, analogous to
¹⁶⁶⁸ sending an empty file. When we uploaded a single pixel white image to service A,
¹⁶⁶⁹ we received responses such as ‘microwave oven’, ‘text’, ‘sky’, ‘white’ and ‘black’
¹⁶⁷⁰ with confidences ranging from 51–95%. Prior checks should be performed on all
¹⁶⁷¹ input data, returning an ‘insufficient information’ error where any input data is below
¹⁶⁷² the information of its training data.

4.7 Threats to Validity**4.7.1 Internal Validity**

1673 Not all computer vision services were assessed. As suggested in Section 4.4, we note

1674 that there are other computer vision services such as IBM Watson. Many services

1675 from Asia were also not considered due to language barriers (of the authors) in

1676 assessing these services. We limited our study to the most popular three providers

1677 (outside of Asia) to maintain focus in this body of work.

1678 A custom confidence threshold was not set. All responses returned from each of

1679 the services were included for analysis; where confidences were low, they were still

1680 included for analysis. This is because we used the default thresholds of each API to

1681 hint at what real-world applications may be like when testing and evaluating these

1682 services.

1683 The label string returned from each service was only considered. It is common

1684 for some labels to respond back that are conceptually similar (e.g., ‘car’ vs. ‘automo-

1685 bile’) or grammatically different (e.g., ‘clothes’ vs. ‘clothing’). While we could have

1686 employed more conceptual comparison or grammatical fixes in this study, we chose

1687 only to compare lowercased labels and as returned. We leave semantic comparison

1688 open to future work.

1689 Only introductory analysis has been applied in assessing the documentation of

1690 these services. Further detailed analysis of documentation quality against a rigorous

1691 documentation quality framework would be needed to fortify our analysis of the

1692 evolution of these services’ documentation.

1695 4.7.2 External Validity

1696 The documentation and services do change over time and evolve, with many allowing
1697 for contributions from the developer community via GitHub. We note that our
1698 evaluation of the documentation was conducted on a single date (see Section 4.4)
1699 and acknowledge that the documentation may have changed from the evaluation date
1700 to the time of this publication. We also acknowledge that the responses and labelling
1701 may have evolved too since the evaluation period described and the date of this
1702 publication. Thus, this may have an impact on the results we have produced in this
1703 paper compared to current, real-world results. To mitigate this, we have supplied the
1704 raw responses available online [297].

1705 Moreover, in this paper we have investigated *computer vision* services. Thus,
1706 the significance of our results to other domains such as natural language processing
1707 or audio transcription is, therefore, unknown. Future studies may wish to repeat our
1708 methodology on other domains to validate if similar patterns occur; we remain this
1709 open for future work.

1710 4.7.3 Construct Validity

1711 It is not clear if all the recommendations proposed in Section 4.6 are feasible
1712 or implementable in practice. Construct validity defines how well an experiment
1713 measures up to its claims; the experiments proposed in this paper support our three
1714 hypotheses but these have been conducted in a clinical condition. Real-world case
1715 studies and feedback from developers and providers in industry would remove the
1716 controlled nature of our work.

4.8 Conclusions & Future Work

¹⁷¹⁷ This study explored three popular computer vision services over an 11 month longi-
¹⁷¹⁹ tudinal experiment to determine if these services pose any evolution risk or incon-
¹⁷²⁰ sistency. We find that these services are generally stable but behave inconsistently;
¹⁷²¹ responses from these services do change with time and this is not visible to the devel-
¹⁷²² opers who use them. Furthermore, the limitations of these systems are not properly
¹⁷²³ conveyed by vendors. From our analysis, we present a set of recommendations for
¹⁷²⁴ both intelligent service vendors and developers.

¹⁷²⁵ Standardised software quality models (e.g., [106]) target maintainability and
¹⁷²⁶ reliability as primary characteristics. Quality software is stable, testable, fault
¹⁷²⁷ tolerant, easy to change and mature. These computer vision services are, however,
¹⁷²⁸ in a nascent stage, difficult to evaluate, and currently are not easily interchangeable.
¹⁷²⁹ Effectively, the intelligent service response objects are shifting in material ways to
¹⁷³⁰ developers, albeit slowly, and vendors do not communicate this evolution or modify
¹⁷³¹ API endpoints; the endpoint remains static but the content returned does not despite
¹⁷³² the same input.

¹⁷³³ There are many potential directions stemming from this work. To start, we plan
¹⁷³⁴ to focus on preparing a more comprehensive datasheet specifically targeted at what
¹⁷³⁵ should be documented to application developers, and not data scientists. Reapplying
¹⁷³⁶ this work in real-world contexts, that is, to get real developer opinions and study
¹⁷³⁷ production grade systems, would also be beneficial to understand these phenomena
¹⁷³⁸ in-context. This will help us clarify if such changes are a real concern for developers

¹⁷³⁹ (i.e., if they really need to change between services, or the service evolution has real
¹⁷⁴⁰ impact on their applications). We also wish to refine and systematise the method
¹⁷⁴¹ used in this study and develop change detectors that can be used to identify evolution
¹⁷⁴² in these services that can be applied to specific ML domains (i.e., not just computer
¹⁷⁴³ vision), data sets, and API endpoints, thereby assisting application developers in their
¹⁷⁴⁴ testing strategies. Moreover, future studies may wish to expand the methodology
¹⁷⁴⁵ applied by refining how the responses are compared. As there does not yet exist a
¹⁷⁴⁶ standardised list of terms available between services, labels could be *semantically*
¹⁷⁴⁷ compared instead of using exact matches (e.g., by using stem words and synonyms
¹⁷⁴⁸ to compare similar meanings of these labels), similar to previous studies [170].

¹⁷⁴⁹ This paper has highlighted only some high-level issues that may be involved
¹⁷⁵⁰ in using these evolving services. The laws of software evolution suggest that for
¹⁷⁵¹ software to be useful, it must evolve [153, 237]. There is, therefore, a trade-off, as
¹⁷⁵² we have shown, between consistency and evolution in this space. For a component
¹⁷⁵³ to be stable, any changes to dependencies it relies on must be communicated. We
¹⁷⁵⁴ are yet to see this maturity of communication from intelligent service providers.
¹⁷⁵⁵ Thus, developers must be cautious between integrating intelligent components into
¹⁷⁵⁶ their applications at the expense of stability; as the field of AI is moving quickly,
¹⁷⁵⁷ we are more likely to see further instability and evolution in intelligent services as a
¹⁷⁵⁸ consequence.

CHAPTER 5

1759

1760

1761

Systematic Mapping Study of API Documentation Knowledge[†]

1762

1763

Abstract Good API documentation facilitates the development process, improving productivity and quality. While the topic of API documentation quality has been of interest for the last two decades, there have been few studies to map the specific constructs needed to create a good document. In effect, we still need a structured taxonomy against which to capture knowledge. This study reports emerging results of a systematic mapping study. We capture key conclusions from previous studies that assess API documentation quality, and synthesise the results into a single framework. By conducting a systematic review of 21 key works,

[†]This chapter is originally based on A. Cummaudo, R. Vasa, and J. Grundy, “What should I document? A preliminary systematic mapping study into API documentation knowledge,” in *13th International Symposium on Empirical Software Engineering and Measurement (ESEM)*. Porto de Galinhas, Recife, Brazil: IEEE, Sep. 2019, pp. 1–6. Terminology has been updated to fit this thesis.

¹⁷⁷⁰ we have developed a five dimensional taxonomy based on 34 categorised weighted recom-
¹⁷⁷¹ mendations. All studies utilise field study techniques to arrive at their recommendations,
¹⁷⁷² with seven studies employing some form of interview and questionnaire, and four conduct-
¹⁷⁷³ ing documentation analysis. The taxonomy we synthesise reinforces that usage description
¹⁷⁷⁴ details (code snippets, tutorials, and reference documents) are generally highly weighted as
¹⁷⁷⁵ helpful in API documentation, in addition to design rationale and presentation. We propose
¹⁷⁷⁶ extensions to this study aligned to developer's utility for each of the taxonomy's categories.

¹⁷⁷⁷ 5.1 Introduction

¹⁷⁷⁸ Improving the quality of API documentation is highly valuable to the software
¹⁷⁷⁹ development process; good documentation facilitates productivity and thus quality
¹⁷⁸⁰ is better engineered into the system [151]. Where an application developer integrates
¹⁷⁸¹ new pieces of functionality (via APIs) into a system, their productivity is affected
¹⁷⁸² either by inadequate skills ("I've never used an API like this, so must learn from
¹⁷⁸³ scratch") or, where their skills are adequate, an imbalanced cognitive load that
¹⁷⁸⁴ causes excessive context switching ("I have the skills for this, but am confused or
¹⁷⁸⁵ misunderstand"). In the latter case, what causes this confusion and how to mitigate
¹⁷⁸⁶ it via improved API documentation is an area that has been explored; prior studies
¹⁷⁸⁷ have provided recommendations based on both qualitative and quantitative analysis
¹⁷⁸⁸ of developer's opinions. These recommendations and guidelines propose ways by
¹⁷⁸⁹ which developers, managers and solution architects can construct systems better.

¹⁷⁹⁰ However, to date there has been little attempt to systematically capture this

1791 knowledge about API documentation from various studies into a readily accessible,
1792 consolidated format, that assists API designers to prepare better documentation.
1793 While previous works have covered certain aspects of API usage, many have lacked
1794 a systematic review of literature and do not offer a taxonomy to consolidate these
1795 guidelines together. For example, some studies have considered the technical imple-
1796 mentation improving API usability or tools to generate (or validate) API documen-
1797 tation from its source code (e.g., [144, 168, 248]); there still lacks a consolidated
1798 effort to capture the knowledge and artefacts best suited to *manually write* API doc-
1799 umentation. This paper presents outcomes from a preliminary work to address this
1800 gap and offers two key contributions:

- 1801 • a systematic mapping study (SMS) consisting of 21 studies that capture what
1802 knowledge or artefacts should be contained within API documentation; and,
- 1803 • a structured taxonomy based on the consolidated recommendations of these
1804 21 studies.

1805 After performing our SMS on what API knowledge should be captured in documentation—
1806 to assist API designers—we propose a five dimensional taxonomy consisting of:
1807 (1) Usage Description; (2) Design Rationale; (3) Domain Concepts; (4) Support
1808 Artefacts; and (5) Documentation Presentation.

1809 This paper is structured as thus: Section 5.2 presents related work in the area;
1810 Section 5.3 is divided into two subsections, the first describing how primary sources
1811 were selected in a SMS, with the second describing the development of our taxon-
1812 omy from these sources; Section 5.4 present our primary studies and our proposed
1813 taxonomy; Section 5.5 describes the threats to validity of this work and Section 5.6

¹⁸¹⁴ provides concluding remarks and the future directions of this study.

¹⁸¹⁵ 5.2 Related Work

¹⁸¹⁶ Systematic mapping studies have previously been explored in the area of API usabil-
¹⁸¹⁷ ity and developer experience. Nybom et al. [168] recently performed a systematic
¹⁸¹⁸ mapping study on 36 API documentation generation tools and approaches. Presented
¹⁸¹⁹ is an analysis of state-of-the-art of the tools developed, what kind of documentation is
¹⁸²⁰ generated by them, and the dependencies they require to generate this documentation.
¹⁸²¹ Their findings highlight a recent effort on the development of API documentation
¹⁸²² by producing example code snippets and/or templates on how to use the API or
¹⁸²³ bootstrap developers to begin using the APIs. A secondary focus is closely followed
¹⁸²⁴ by tools that produce natural language descriptions that can be produced within de-
¹⁸²⁵ veloper documentation. However, Nybom et al. produce a systematic mapping study
¹⁸²⁶ on the types of *tooling* that exists to assist in producing and validating API documen-
¹⁸²⁷ tation. While this is a systematic study with key insights into the types of tooling
¹⁸²⁸ produced, there is still a gap for a systematic mapping study in what *guidelines* have
¹⁸²⁹ been produced by the literature in developing natural-language documentation itself,
¹⁸³⁰ which our work has addressed.

¹⁸³¹ Watson [248] performed a heuristic assessment of 11 high-level universal design
¹⁸³² elements of API documentation against 35 popular APIs. He demonstrated that many
¹⁸³³ of these popular APIs fail to grasp even the basic of these elements; for example,
¹⁸³⁴ 25% of the documentation sets did not provide any basic overview documentation.

¹⁸³⁵ However, the heuristic used within this study consists of just 11 elements and is based
¹⁸³⁶ on only three seminal works. Our work extends these heuristics and structures them
¹⁸³⁷ into a consolidated, hierarchical taxonomy using a systematic taxonomy development
¹⁸³⁸ method for SE.

¹⁸³⁹ A taxonomy of knowledge patterns within API reference documentation by
¹⁸⁴⁰ Maalej and Robillard [144] classified 12 distinct knowledge types. Evaluation of the
¹⁸⁴¹ taxonomy against JDK 6 and .NET 4.0 showed that, while functionality and structure
¹⁸⁴² of the API is well-communicated, core concepts and rationale about the API are
¹⁸⁴³ quite rare to find. Moreover, they demonstrated that low-value ‘non-information’
¹⁸⁴⁴ (documentation that provides uninformative boilerplate text with no insight into the
¹⁸⁴⁵ API at all) is substantially present in the documentation of methods and fields in these
¹⁸⁴⁶ APIs. Their findings recommend that developers factor their 12 distinct knowledge
¹⁸⁴⁷ types into the process of code documentation and prevent documenting low-value
¹⁸⁴⁸ documentation. The development of their taxonomy consisted of questions to model
¹⁸⁴⁹ knowledge and information, thereby capturing the reason about disparate information
¹⁸⁵⁰ units independent to context; a key difference to this paper is the systematic taxonomy
¹⁸⁵¹ approach utilised.

¹⁸⁵² 5.3 Method

¹⁸⁵³ Our taxonomy development consisted of two phases. Firstly, we conducted a SMS
¹⁸⁵⁴ to identify and analyse API documentation studies, following the guidelines of
¹⁸⁵⁵ Kitchenham and Charters [121] and Petersen et al. [183]. Following this, we followed

¹⁸⁵⁶ the software engineering (SE) taxonomy development method devised by Usman
¹⁸⁵⁷ et al. [240] on our findings from the SMS.

¹⁸⁵⁸ 5.3.1 Systematic Mapping Study

¹⁸⁵⁹ **Research Questions (RQs)** To guide our SMS, we developed the following RQs:

¹⁸⁶⁰ **RQ1** What knowledge do API documentation studies contribute?

¹⁸⁶¹ **RQ2** How is API documentation studied?

¹⁸⁶² The intent behind RQ1 is to collate as much of the insight provided by the
¹⁸⁶³ literature on how API providers should best document their work. This helped us
¹⁸⁶⁴ shape and form the taxonomy provided in Section 5.4. RQ2 addresses methodologies
¹⁸⁶⁵ by which these studies come to these conclusions to identify gaps in literature where
¹⁸⁶⁶ future studies can potentially focus.

¹⁸⁶⁷ **Automatic Filtering** Informed by similar previous studies in SE [82, 86, 240], we
¹⁸⁶⁸ begin by defining the SWEBOK [102] knowledge areas (KAs) to assist in the search
¹⁸⁶⁹ and mapping process of an SMS. Our search query was built using related KAs,
¹⁸⁷⁰ relevant synonyms, and the term ‘software engineering’ (for comprehensiveness) all
¹⁸⁷¹ joined with the OR operator. Due to the lack of a standard definition of an API,
¹⁸⁷² we include the terms: ‘API’ and its expanded term; software library, component
¹⁸⁷³ and framework; and lastly SDK and its expanded term. These too were joined with
¹⁸⁷⁴ the OR operator, appended with an AND. Lastly, the term ‘documentation’ was
¹⁸⁷⁵ appended with an AND. Our final search string was:

("software design" OR "software architecture" OR "software construction" OR "software development" OR "software maintenance" OR "software engineering process" OR "software process" OR "software lifecycle" OR "software methods" OR "software quality" OR "software engineering professional practice" OR "software engineering") AND (api OR "application programming interface" OR "software library" OR "software component" OR "software framework" OR sdk OR "software development kit") AND (documentation)

¹⁸⁷⁶ The query was then executed on all available metadata (title, abstract and key-
¹⁸⁷⁷ words) on three primary sources to search for relevant studies in May 2019. Web of
¹⁸⁷⁸ Science¹ (WoS), Compendex/Inspec² (C/I) and Scopus³ were chosen due to their rel-
¹⁸⁷⁹ evance in SE literature (containing the IEEE, ACM, Springer and Elsevier databases)
¹⁸⁸⁰ and their ability to support advanced queries [32, 121]. A total 4,501 results⁴ were
¹⁸⁸¹ found, with 549 being duplicates. Table 5.1 displays our results in further detail
¹⁸⁸² (duplicates not omitted).

¹⁸⁸³ **Manual Filtering** A follow-up manual filtering to select primary studies was per-
¹⁸⁸⁴ formed on the 4,501 results using the following inclusion criteria (IC) and exclusion
¹⁸⁸⁵ criteria (EC):

¹⁸⁸⁶ **IC1** Studies must be relevant to API documentation: specifically, we
¹⁸⁸⁷ exclude studies that deal with improving the technical API usability
¹⁸⁸⁸ (e.g., improved usage patterns);

¹<http://apps.webofknowledge.com> last accessed 23 May 2019.

²<http://www.engineeringvillage.com> last accessed 23 May 2019.

³<http://www.scopus.com> last accessed 23 May 2019.

⁴Raw results can be located at <http://bit.ly/2KxBLs4>

Table 5.1: Summary of our search results and publication types

Publication type	WoS	C/I	Scopus	Total
Conference Paper	27	442	2353	2822
Journal Article	41	127	1236	1404
Book	23	17	224	264
Other	0	5	6	11
Total	91	591	3819	4501

- ¹⁸⁸⁹ **IC2** Studies must propose new knowledge or recommendations to document APIs;
- ¹⁸⁹⁰ **IC3** Studies must be relevant to SE as defined in SWEBOK;
- ¹⁸⁹² **EC1** Studies where full-text is not accessible through standard institutional databases;
- ¹⁸⁹⁴ **EC2** Studies that do not propose or extend how to improve the official, natural language documentation of an API;
- ¹⁸⁹⁶ **EC3** Studies proposing a third-party tool to enhance existing documentation or generate new documentation using data mining (i.e., not proposing strategies to improve official documentation);
- ¹⁸⁹⁹ **EC4** Studies not written in English;
- ¹⁹⁰⁰ **EC5** Studies not peer-reviewed.
- ¹⁹⁰¹ After exporting metadata of search results to a spreadsheet, a three-phase curation process was conducted. The first author read the publication source (to omit non-
- ¹⁹⁰²

1903 SE papers quickly), author keywords and title of all 4,501 studies (514 that were
1904 duplicates), and abstract. As we considered multiple databases, some studies were
1905 repeated. However, the DOIs and titles were sorted and reviewed, retaining only
1906 one copy of the paper from a single database. Moreover, as there was no limit
1907 to the year range in our query, some studies were republished in various venues.
1908 These, too, were handled with title similarity matching, wherein only the first paper
1909 was considered. Where the inclusion or exclusion criteria could not be determined
1910 from the abstract alone, the paper was automatically shortlisted. Any doubt in a
1911 study automatically included it into the second phase. This resulted in 133 studies
1912 being shortlisted to the second phase. We rejected 427 studies that were unrelated
1913 to SE, 3,235 were not directly related to documenting APIs (e.g., to enhance coding
1914 techniques that improve the overall developer usability of the API), 182 proposed new
1915 tools to enhance API documentation or used machine learning to mine developer's
1916 discussion of APIs, and 10 were not in English.

1917 The shortlisted studies were then re-evaluated by re-reading the abstract, the
1918 introduction and conclusion. Performing this second phase removed a further 64
1919 studies that were on API usability or non API-related documentation (i.e., code
1920 commenting); we further refined our exclusion criteria to better match the research
1921 outcomes of this goal (chiefly including the word ‘natural language’ documentation
1922 in EC2) which removed studies focused to improve technical documentation of
1923 APIs such as data types and communication schemas. Additionally, 26 studies were
1924 removed as they were related to introducing new tools (EC3), 3 were focused on tools
1925 to mine API documentation, 7 studies where no recommendations were provided,

1926 2 further duplicate studies, and a further 10 studies where the full text was not
1927 available, not peer reviewed or in English. Books are commonly not peer-reviewed
1928 (EC5), however no books were shortlisted within these results. **This resulted in 21**
1929 **primary studies for further analysis. The mapping of primary study identifiers**
1930 **to references S1–21 can be found in Appendix D.**

1931 Intra-rater reliability of our 133 shortlisted papers was tested using the test-retest
1932 approach [121] by re-evaluating a random sample of 10% (13 total) of the studies
1933 shortlisted above a week after initial studies were shortlisted. Using the Cohen's
1934 kappa coefficient as a metric for reliability, $\kappa = 0.7547$, indicating substantial
1935 agreement [132].

1936 **Data Extraction** Of the 21 primary studies, we conducted abstract key-wording
1937 adhering to Petersen et al.'s guidelines [183] to develop a classification scheme. An
1938 initial set of keywords were applied for each paper in terms of their methodologies and
1939 research approaches (RQ2), based on an existing classification schema by Wieringa
1940 and Heerkens [252]: evaluation, validation, personal experience and philosophical
1941 papers.

1942 After all primary studies had been assigned keywords, we noticed that **all papers**
1943 **used field study techniques**, and thus we consolidated these keywords using Singer
1944 et al.'s framework of SE field study techniques [219]. Singer et al. captures both
1945 study techniques *and* methods to collect data within the one framework, namely:
1946 *direct techniques*, including brainstorming and focus groups, interviews and ques-
1947 tionnaires, conceptual modelling, work diaries, think-aloud sessions, shadowing and

Table 5.2: Data extraction form

Data item(s)	Description
Citation metadata	Title, author(s), years, publication venue, publication type
Key recommendation(s)	As per IC2, the study must propose at least one recommendation on what should be captured in API documentation
Evaluation method	Did the authors evaluate their recommendations? If so, how?
Primary technique	The primary technique used to devise the recommendation(s)
Secondary technique	As above, if a second study was conducted
Tertiary technique	As above, if a third study was conducted
Research type	The research type employed in the study as defined by Wieringa and Heerkens's taxonomy

¹⁹⁴⁸ observation, participant observation; *indirect techniques*, including instrumenting systems, fly-on-the-wall; and *independent techniques*, including analysis of work databases, tool use logs, documentation analysis, and static and dynamic analysis.

¹⁹⁵¹ Table 5.2 describes our data extraction form, which was used to collect relevant data from each paper. Figure 5.1 maps each study to one (or more, if applicable) of ¹⁹⁵² methodologies plotted against Wieringa and Heerkens's research approaches.

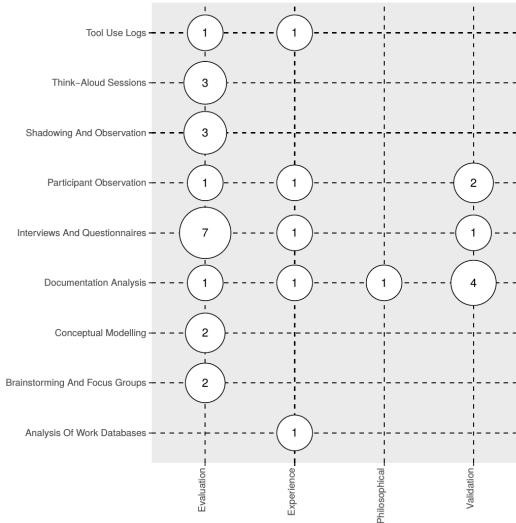


Figure 5.1: Systematic map: field study technique vs research type

1954 **5.3.2 Development of the Taxonomy**

1955 Usman et al. concludes that a majority of SE taxonomies are developed in an ad-hoc
 1956 way [240], and proposes a systematic approach to develop taxonomies in SE that
 1957 extends previous efforts by including lessons learned from more mature fields. In
 1958 this subsection, we outline the 4 phases and 13 steps taken to develop our taxonomy
 1959 based on Usman et al.'s technique.

1960 **Planning phase** The planning phase of the technique involves the following six
 1961 steps:

- 1962 **(1) defining the SE KA:** The software engineering KA, as defined by the SWEBOK,
 1963 is software construction;
- 1964 **(2) defining the objective:** The main objective of the proposed taxonomy is to
 1965 define a set of categories that enables to classify different facets of natural-

1966 language API *documentation* knowledge (not API *usability* knowledge) as
1967 reported in existing literature;

1968 (3) *defining the subject matter*: The subject matter of our proposed taxonomy is
1969 documentation artefacts of APIs;

1970 (4) *defining the classification structure*: The classification structure of our pro-
1971 posed taxonomy is *hierarchical*;

1972 (5) *defining the classification procedure*: The procedure used to classify the
1973 documentation artefacts is qualitative;

1974 (6) *defining the data sources*: The basis of the taxonomy is derived from field
1975 study techniques (see Section 5.3.1).

1976 **Identification and extraction phase** The second phase of the taxonomy devel-
1977 opment involves (7) *extracting all terms and concepts* from relevant literature, as
1978 selected from our 21 primary studies. These terms are then consolidated by (8) *per-*
1979 *forming terminology control*, as some terms may refer to different concepts and
1980 vice-versa.

1981 **Design phase** The design phase identified the core dimensions and categories
1982 within the extracted data items. The first step is to (9) *identify and define taxonomy*
1983 *dimensions*; for this study we utilised a bottom-up approach to identify the dimen-
1984 sions, i.e. extracting the categories first and then nominating which dimensions
1985 these categories fit into using an iterative approach. As a bottom-up approach was
1986 utilised, step (9) also encompassed the second stage of the design phase, which is to

1987 (10) *identify and describe the categories* of each dimension. Thirdly, we (11) *iden-*
1988 *tify and describe relationships* between dimensions and categories, which can be
1989 skipped if the relationships are too close together, as is the case of our grouping
1990 technique which allows for new dimensions and categories to be added. The last
1991 step in this phase is to (12) *define guidelines for using and updating the taxonomy*,
1992 however as this taxonomy still an emerging result, guidelines to update and use the
1993 taxonomy are anticipated future work.

1994 **Validation phase** In the final phase of taxonomy development, taxonomy designers
1995 must (13) *validate the taxonomy* to assess its usefulness. Ideally, this is done by
1996 applying the taxonomy heuristically against developers or real-world case-studies.
1997 This remains a plan for future work (see Section 5.5).

1998 5.4 Taxonomy

1999 Our taxonomy consists of five dimensions (labelled A–E) that respectively cover:
2000 [A] **Usage Description** on *how* to use the API for the developer's intended use case;
2001 [B] **Design Rationale** on *when* the developer should choose this API for a particular
2002 use case; [C] **Domain Concepts** of the domain behind the API to understand *why* this
2003 API should be chosen for this domain; [D] **Support Artefacts** that describe *what ad-*
2004 *ditional documentation the API provides; and [E] Documentation Presentation* to
2005 help organise the *visualisation* of the above information. Further descriptions of the
2006 categories encompassing each dimension are given within Table 5.3, coded as [Xi],
2007 where i is the category identifier within a dimension, X , where $X \in \{A, B, C, D, E\}$.

2008 We expand these five dimensions into 34 categories (sub-dimensions) and
2009 Table 5.3 provides a weighting of these categories in the rightmost column as
2010 calculated as a percentage of the number of primary studies per category di-
2011 vided by the total of primary studies. The top five weighted categories (bolded in
2012 Table 5.3) highlight what most studies recommend documenting in API documen-
2013 tation, with the top three falling under the Usage Description dimension.

2014 The majority (71%) of studies advocate for **code snippets** as a necessary piece
2015 in the API documentation puzzle [A5]. While code snippets generally only reflect
2016 small portions of API functionality (limited to 15–30 LoC), this is complimented by
2017 **step-by-step tutorials** (57% of studies) that tie in multiple (disparate) components
2018 of API functionality, generally with some form of screenshots, demonstrating the
2019 development of a non-trivial application using the API step-by-step [A6]. The third
2020 highest category weighted was also under the Usage Description dimension, being
2021 **low-level reference documentation** at 52% [A2]. These three categories were the
2022 only categories to be weighted as majority categories (i.e., their weighting was above
2023 50%).

2024 The fourth and fifth highest weights are **an entry-level purpose/overview of**
2025 **the API** (48%) that gives a brief motivation as to why a developer should choose
2026 a particular API over another [B1] and **consistency in the look and feel** of the
2027 documentation throughout all of the API's official documentation (43%) [E6].

2028 5.5 Threats to Validity

2029 Threats to *internal validity* concern factors internal to our study that may affect
2030 results. Guidelines on producing systematic reviews [121] suggest that single re-
2031 searchers conducting their reviews should discuss the review protocol, inclusion
2032 decisions, data extraction with a third party. In this paper, we have presented the
2033 early outcomes of our systematic review, which has utilised the test-retest method-
2034 ology as a measure of reliability. MacDonell et al. [145] states that a defining
2035 characteristic of any SMS is to test the reliability of the review and extraction pro-
2036 cesses. We plan to mitigate this threat by conducting *inter*-relater reliability with
2037 the continuation of this work, using independent analysis and conflict resolution as
2038 per guidelines suggested by Garousi and Felderer [81]. Similarly, the development
2039 of our taxonomy would benefit from an inter-rater reliability categorisation of a
2040 sample of papers to both ensure that our weightings of categories are reliable and
2041 that the categories and dimensions fit the objectives of the taxonomy. Furthermore,
2042 a future user study (see Section 5.6) will be needed to assess whether the extracted
2043 information from API documentation actually impacts on developer productivity,
2044 and the usefulness of such a taxonomy should be evaluated.

2045 Threats to *external validity* represent the generalisation of the observations we
2046 have found in this study. While we have used a broad range of literature that
2047 encompasses API documentation guidelines, we acknowledge that not all papers
2048 contributing to API documentation may have been captured in the taxonomy. All
2049 efforts were made to include as many papers as possible given our filtering technique,

2050 though it is likely that some papers filtered out (e.g., papers not in English) may
2051 alter our conclusions, introducing conflicting recommendations. However, given the
2052 consistency of these trends within the studies that were sourced, we consider this a
2053 low likelihood.

2054 Threats to *construct validity* relates to the degree by which the data extrapolated
2055 in this study sufficiently measures its intended goals. Automatic searching was
2056 conducted in the SMS by choice of three popular databases (see Section 5.3.1).
2057 As a consequence of selecting multiple databases, duplicates were returned. This
2058 was mitigated by manually curating out all duplicate results from the set of studies
2059 returned. Additionally, we acknowledge that the lack manual searching of papers
2060 within particular venues may be an additional threat due to the misalignment of
2061 search query keywords to intended papers of inclusion. Thus, our conclusions are
2062 only applicable to the information we were able to extract and summarise, given the
2063 primary sources selected.

2064 5.6 Conclusions & Future Work

2065 API documentation is an aspect of quality of software, as it facilitates the developer's
2066 productivity and assists with evolution. Improving the quality of the documentation
2067 of third party APIs improves the quality of software using them.

2068 To date, we did not find a systematic literature review that offers a consolidated
2069 taxonomy of key recommendations. Moreover, there has been little work on map-
2070 ping the research produced in this space against the techniques used to arrive at

2071 the recommendations. Starting with 4,501 papers potentially relating to API doc-
2072 umentation, we identified 21 key relevant studies, and synthesise a taxonomy of
2073 the various documentation aspects that should improve API documentation quality.
2074 Furthermore, we also capture the most commonly used analysis techniques used in
2075 the academic literature. Figure 5.1 highlights that a majority of these studies employ
2076 interviews and questionnaires, and only some undertake structured documentation
2077 analysis.

2078 In future revisions of this work, we intend use our results as the input to a
2079 restricted systematic literature review in API documentation artefacts. In doing so,
2080 we will consider conducting the following:

- 2081 • improving reliability metrics of our study (see Section 5.5) with an inter-rater
2082 reliability method;
- 2083 • the development and applicability of our taxonomy will be further explored
2084 by triangulating the taxonomy against actual developers in industry to assess
2085 the efficacy of these recommendations—this will empirically reflect what is
2086 important from a *practitioner* point of view;
- 2087 • reviewing the techniques and evaluation of our selected studies to extract the
2088 effectiveness of the various approaches used in the conclusions;
- 2089 • conducting a heuristic validation of the taxonomy against intelligent APIs,
2090 given the current trend in SE that is exploring how machine learning and
2091 artificial intelligence-based applications may affect existing approaches;
- 2092 • arrive at a relevance ranking for each of the 34 categories, based on developer

2093 surveys and current weights.

2094 We believe the results of this preliminary empirical work may provide further
2095 insight for future follow-up user studies with developers. Whilst our aim is to
2096 eventually improve the quality of API documentation, the ultimate goal is improving
2097 the developer's experience when producing systems and, therefore, improving the
2098 efficacy and productivity at which software is produced within industry. We hope
2099 that API designers will utilise the taxonomy produced in this paper as a weighted
2100 checklist for what should be considered in their own APIs.

Table 5.3: An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy.

Key	Description	Primary	Total (%)
		Sources	
A1	Quick-start guides to rapidly get started using the API in a specific programming language.	S4, S9, S10	3/21 (14%)
A2	Low-level reference manual documenting all API components to review fine-grade detail.	S1, S3, S4, S8, S9, S10, S11, S12, S15, S16, S17	11/21 (52%)
A3	Explanations of the API's high-level architecture to better understand intent and context.	S1, S2, S4, S11, S14, S16, S19, S20	8/21 (38%)

Continued on next page...

Table 5.3: An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (*Continued*).

Key	Description	Primary	Total (%)
		Sources	
A4	Source code implementation and code comments (where applicable) to understand the API author's mindset.	S1, S4, S7, S12, S13, S17, S20	7/21 (33%)
A5	Code snippets (with comments) of no more than 30 LoC to understand a basic component functionality within the API.	S1, S2, S4, S5, S6, S7, S9, S10, S11, S14, S15, S16, S18, S20, S21	15/21 (71%)
A6	Step-by-step tutorials, with screenshots to understand how to build a non-trivial piece of functionality with multiple components of the API.	S1, S2, S4, S5, S7, S9, S10, S15, S16, S18, S20, S21	12/21 (57%)
A7	Downloadable source code of production-ready applications that use the API to understand implementation in a large-scale solution.	S1, S2, S5, S9, S15	5/21 (24%)
A8	Best-practices of implementation to assist with debugging and efficient use of the API.	S1, S2, S4, S5, S7, S8, S9, S14	8/21 (38%)
A9	An exhaustive list of all major components that exist within the API.	S4, S16, S19	3/21 (14%)

Continued on next page...

Table 5.3: An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (*Continued*).

Key	Description	Primary	Total (%)
		Sources	
A10	Minimum system requirements and dependencies to use the API.	S4, S7, S13, S17, S19	5/21 (24%)
A11	Instructions to install or begin using the API and details on its release cycle and updating it.	S4, S7, S8, S9, S11, S13, S16, S19	8/21 (38%)
A12	Error definitions that describe how to address a specific problem.	S1, S2, S4, S5, S9, S11, S13	7/21 (33%)
B1	A brief description of the purpose or overview of the API as a low barrier to entry.	S1, S2, S4, S5, S6, S8, S10, S11, S15, S16	10/21 (48%)
B2	Descriptions of the types of applications the API can develop.	S2, S4, S9, S11, S15, S18	6/21 (29%)
B3	Descriptions of the types of users who should use the API.	S4, S9	2/21 (10%)
B4	Descriptions of the types of users who will use the product the API creates.	S4	1/21 (5%)
B5	Success stories about the API used in production.	S4	1/21 (5%)
B6	Documentation to compare similar APIs within the context to this API.	S2, S6, S13, S18	4/21 (19%)

Continued on next page...

Table 5.3: An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (*Continued*).

Key Description	Primary Sources	Total (%)
B7 Limitations on what the API can and cannot provide.	S4, S5, S8, S9, S14, S16	6/21 (29%)
C1 Descriptions of the relationship between API components and domain concepts.	S3, S10	2/21 (10%)
C2 Definitions of domain-terminology and concepts, with synonyms if applicable.	S2, S3, S4, S6, S7, S10, S14, S16	8/21 (38%)
C3 Generalised documentation for non-technical audiences regarding the API and its domain.	S4, S8, S16	3/21 (14%)
D1 A list of FAQs.	S4, S7	2/21 (10%)
D2 Troubleshooting suggestions.	S4, S8	2/21 (10%)
D3 Diagrammatically representing API components using visual architectural representations.	S6, S13, S20	3/21 (14%)
D4 Contact information for technical support.	S4, S8, S19	3/21 (14%)
D5 A printed/printable resource for assistance.	S4, S6, S7, S9, S16	5/21 (24%)
D6 Licensing information.	S7	1/21 (5%)

Continued on next page...

Table 5.3: An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy (*Continued*).

Key	Description	Primary	Total (%)
		Sources	
E1	Searchable knowledge base.	S3, S4, S6, S10, S14, S17, S18	7/21 (33%)
E2	Context-specific discussion forum.	S4, S10, S11	3/21 (14%)
E3	Quick-links to other relevant documentation frequently viewed by developers.	S6, S16, S20	3/21 (14%)
E4	Structured navigational style (e.g., bread- crumbs).	S6, S10, S20	3/21 (14%)
E5	Visualised map of navigational paths to certain API components in the website.	S6, S14, S20	3/21 (14%)
E6	Consistent look and feel of documentation.	S1, S2, S3, S5, S6, S8, S10, S15, S20	9/21 (43%)

Part III

2101

Postface

2102

CHAPTER 6

2103

2104

2105

Conclusion

2106

2107

2108

References

2109

- 2110 [1] “SOAP, Representational state transfer - Explore - Google Trends.”
- 2111 [2] *Pivotal Cloud Foundry, Google ML, and Spring*, Dec. 2017.
- 2112 [3] *Machine learning with Google APIs*, Jan. 2019.
- 2113 [4] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving,
2114 and M. Isard, “Tensorflow: A system for large-scale machine learning,” in *12th USENIX
2115 Symposium on Operating Systems Design and Implementation (OSDI 16)*, 2016, pp. 265–283.
- 2116 [5] R. E. Al-Qutaish, “Quality models in software engineering literature: an analytical and com-
2117 parative study,” *Journal of American Science*, vol. 6, no. 3, pp. 166–175, 2010.
- 2118 [6] H. Allahyari and N. Lavesson, “User-oriented assessment of classification model under-
2119 ability,” in *11th scandinavian conference on Artificial intelligence*. IOS Press, 2011.
- 2120 [7] M. Allamanis and C. Sutton, “Why, when, and what: analyzing stack overflow questions by
2121 topic, type, and code,” in *Proceedings of the 10th Working Conference on Mining Software
2122 Repositories*. IEEE Press, 2013, pp. 53–56.
- 2123 [8] K. Arnold, “Programmers are people, too,” *ACM Queue*, vol. 3, no. 5, pp. 54–59, 2005.
- 2124 [9] A. Arpteg, B. Brinne, L. Crnkovic-Friis, and J. Bosch, “Software Engineering Challenges of

- 2125 Deep Learning,” in *2018 44th Euromicro Conference on Software Engineering and Advanced*
2126 *Applications (SEAA)*. IEEE, Oct. 2018, pp. 50–59.
- 2127 [10] W. R. Ashby and J. R. Pierce, “An Introduction to Cybernetics,” *Physics Today*, vol. 10, no. 7,
2128 pp. 34–36, Jul. 1957.
- 2129 [11] M. G. Augasta and T. Kathirvalavakumar, “Reverse engineering the neural networks for rule
2130 extraction in classification problems,” *Neural processing letters*, vol. 35, no. 2, pp. 131–150,
2131 2012.
- 2132 [12] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K.-R. MÃžller, “How
2133 to explain individual classification decisions,” *Journal of Machine Learning Research*, vol. 11,
2134 no. Jun, pp. 1803–1831, 2010.
- 2135 [13] B. Baesens, C. Mues, M. De Backer, J. Vanthienen, and R. Setiono, “Building Intelligent Credit
2136 Scoring Systems Using Decision Tables.” *ICEIS*, 2003.
- 2137 [14] K. Bajaj, K. Pattabiraman, and A. Mesbah, “Mining questions asked by web developers,” in
2138 *Proceedings of the 11th Working Conference on Mining Software Repositories*. ACM, 2014,
2139 pp. 112–121.
- 2140 [15] A. Barua, S. W. Thomas, and A. E. Hassan, “What are developers talking about? An analysis
2141 of topics and trends in Stack Overflow,” *Empirical Software Engineering*, vol. 19, no. 3, pp.
2142 619–654, 2014.
- 2143 [16] Y. Baruch, “Response rate in academic studies—A comparative analysis,” *Human relations*,
2144 vol. 52, no. 4, pp. 421–438, 1999.
- 2145 [17] L. Bass, P. Clements, and R. Kazman, *Software architecture in practice*. Addison-Wesley
2146 Professional, 2003.
- 2147 [18] R. Bellazzi and B. Zupan, “Predictive data mining in clinical medicine: current issues and
2148 guidelines,” *International journal of medical informatics*, vol. 77, no. 2, pp. 81–97, 2008.
- 2149 [19] A. Ben-David, “Monotonicity maintenance in information-theoretic machine learning algo-
2150 rithms,” *Machine Learning*, vol. 19, no. 1, pp. 29–43, 1995.
- 2151 [20] T. Berners-Lee, R. Fielding, and L. Masinter, “Uniform resource identifier (URI): Generic
2152 syntax,” Tech. Rep., 2004.

- 2153 [21] J. Bessin, "The Business Value of Quality," *IBM developerWorks, June*, vol. 15, 2004.
- 2154 [22] J. J. Blake, L. P. Maguire, B. Roche, T. M. McGinnity, and L. J. McDaid, "The Implementation
2155 of Fuzzy Systems, Neural Networks and Fuzzy Neural Networks using FPGAs." *Inf. Sci.*, 1998.
- 2156 [23] B. Boehm and V. R. Basili, "Software defect reduction top 10 list," *Foundations of empirical
2157 software engineering: the legacy of Victor R. Basili*, vol. 426, no. 37, 2005.
- 2158 [24] B. W. Boehm, *Software engineering economics*. Prentice-hall Englewood Cliffs (NJ), 1981,
2159 vol. 197.
- 2160 [25] B. W. Boehm, J. R. Brown, and H. Kaspar, "Characteristics of software quality," 1978.
- 2161 [26] O. Boz, "Extracting decision trees from trained neural networks," in *Proceedings of the eighth
2162 ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM,
2163 2002, pp. 456–461.
- 2164 [27] H. B. Braiek and F. Khomh, "On testing machine learning programs," *arXiv preprint
2165 arXiv:1812.02257*, 2018.
- 2166 [28] M. Bramer, *Principles of data mining*. Springer, 2007, vol. 180.
- 2167 [29] J. Brandt, P. J. Guo, J. Lewenstein, M. Dontcheva, and S. R. Klemmer, "Two studies of oppor-
2168 tunistic programming: interleaving web foraging, learning, and writing code," in *Proceedings
2169 of the SIGCHI Conference on Human Factors in Computing Systems*. ACM, 2009, pp. 1589–
2170 1598.
- 2171 [30] E. Breck, S. Cai, E. Nielsen, M. Salib, and D. Sculley, "What's your ML Test Score? A rubric
2172 for ML production systems," 2016.
- 2173 [31] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and Regression Trees*.
2174 CRC press, 1984.
- 2175 [32] P. Brereton, B. A. Kitchenham, D. Budgen, M. Turner, and M. Khalil, "Lessons from applying
2176 the systematic literature review process within the software engineering domain," *Journal of
2177 Systems and Software*, vol. 80, no. 4, pp. 571–583, Apr. 2007.
- 2178 [33] M. Bunge, "A General Black Box Theory," *Philosophy of Science*, vol. 30, no. 4, pp. 346–358,
2179 Oct. 1963.

- 2180 [34] A. Bussone, S. Stumpf, and D. O’Sullivan, “The Role of Explanations on Trust and Reliance in
2181 Clinical Decision Support Systems.” *ICHI*, 2015.
- 2182 [35] C. Calhoun, *Critical social theory: Culture, history, and the challenge of difference*. Wiley-
2183 Blackwell, 1995.
- 2184 [36] G. Canfora, “User-side testing of web services,” in *Software Maintenance and Reengineering*,
2185 2005. *CSMR 2005. Ninth European Conference on*. IEEE, 2005, p. 301.
- 2186 [37] G. Canfora and M. Di Penta, “Testing services and service-centric systems: Challenges and
2187 opportunities,” *It Professional*, vol. 8, no. 2, pp. 10–17, 2006.
- 2188 [38] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad, “Intelligible Models
2189 for HealthCare - Predicting Pneumonia Risk and Hospital 30-day Readmission.” *KDD*, pp.
2190 1721–1730, 2015.
- 2191 [39] F. Casati, H. Kuno, G. Alonso, and V. Machiraju, “Web Services-Concepts, Architectures and
2192 Applications,” 2003.
- 2193 [40] J. P. Cavano, J. A. McCall, J. P. Cavano, J. A. McCall, J. P. Cavano, and J. A. McCall, “A
2194 framework for the measurement of software quality,” *ACM SIGSOFT Software Engineering*
2195 *Notes*, vol. 3, no. 5, pp. 133–139, Nov. 1978.
- 2196 [41] C. V. Chambers, D. J. Balaban, B. L. Carlson, and D. M. Grasberger, “The effect of
2197 microcomputer-generated reminders on influenza vaccination rates in a university-based family
2198 practice center,” *The Journal of the American Board of Family Practice*, vol. 4, no. 1, pp. 19–26,
2199 1991.
- 2200 [42] J. Cheng and R. Greiner, “Learning bayesian belief network classifiers: Algorithms and system,”
2201 in *Conference of the Canadian Society for Computational Studies of Intelligence*. Springer,
2202 2001, pp. 141–151.
- 2203 [43] E. Choi, N. Yoshida, R. G. Kula, and K. Inoue, “What do practitioners ask about code clone?
2204 a preliminary investigation of stack overflow.” in *IWSC*, 2015, pp. 49–50.
- 2205 [44] Digital Inc., “Case study: Finding defects earlier yields enormous savings,” 2003.
- 2206 [45] P. Clark and R. Boswell, “Rule induction with CN2: Some recent improvements,” in *European*
2207 *Working Session on Learning*. Springer, 1991, pp. 151–163.

- 2208 [46] M. Craven and J. W. Shavlik, “Extracting Tree-Structured Representations of Trained Net-
2209 works.” *NIPS*, 1995.
- 2210 [47] J. W. Creswell and J. D. Creswell, *Research design: Qualitative, quantitative, and mixed
2211 methods approaches*. Sage publications, 2017.
- 2212 [48] P. B. Crosby, “Quality is free: The art of making quality free,” *New York*, 1979.
- 2213 [49] A. Cummaudo, R. Vasa, and J. Grundy, “What should I document? A preliminary systematic
2214 mapping study into API documentation knowledge,” in *13th International Symposium on Em-
2215 pirical Software Engineering and Measurement (ESEM)*. Porto de Galinhas, Recife, Brazil:
2216 IEEE, Sep. 2019, pp. 1–6.
- 2217 [50] A. Cummaudo, R. Vasa, J. Grundy, M. Abdelrazek, and A. Cain, “Losing Confidence in Quality:
2218 Unspoken Evolution of Computer Vision Services,” in *2019 IEEE International Conference on
2219 Software Maintenance and Evolution (ICSME)*. Cleveland, OH, USA: IEEE, Sep. 2019, pp.
2220 333–342.
- 2221 [51] H. da Mota Silveira and L. C. Martini, “How the New Approaches on Cloud Computer Vision
2222 can Contribute to Growth of Assistive Technologies to Visually Impaired in the Following
2223 Years,” *Journal of Information Systems Engineering and Management*, vol. 2, no. 2, pp. 1–3,
2224 2017.
- 2225 [52] R. Davison, M. G. Martinsons, and N. Kock, “Principles of canonical action research,” *Infor-
2226 mation systems journal*, vol. 14, no. 1, pp. 65–86, 2004.
- 2227 [53] K. Dejaeger, F. Goethals, A. Giangreco, L. Mola, and B. Baesens, “Gaining insight into student
2228 satisfaction using comprehensible data mining techniques,” *European Journal of Operational
2229 Research*, vol. 218, no. 2, pp. 548–562, 2012.
- 2230 [54] S. Demeyer and T. Mens, *Software Evolution*. Springer, 2008.
- 2231 [55] I. Dey, *Qualitative data analysis: A user friendly guide for social scientists*. Routledge, 2003.
- 2232 [56] V. Dhar, D. Chou, and F. Provost, “Discovering Interesting Patterns for Investment Decision
2233 Making with GLOWER—A Genetic Learner Overlaid with Entropy Reduction,” *Data Mining
2234 and Knowledge Discovery*, vol. 4, no. 4, pp. 251–280, 2000.

- 2235 [57] V. C. Dibia, M. Ashoori, A. Cox, and J. D. Weisz, “TJBot,” in *the 2017 CHI Conference*
2236 *Extended Abstracts*. New York, New York, USA: ACM Press, 2017, pp. 381–384.
- 2237 [58] M. Doderer, K. Yoon, J. Salinas, and S. Kwek, “Protein subcellular localization prediction
2238 using a hybrid of similarity search and error-correcting output code techniques that produces
2239 interpretable results,” *In silico biology*, vol. 6, no. 5, pp. 419–433, 2006.
- 2240 [59] P. Domingos, “Occam’s two razors: The sharp and the blunt,” in *KDD*, 1998, pp. 37–43.
- 2241 [60] B. Dorn and M. Guzdial, “Learning on the job: characterizing the programming knowledge
2242 and learning strategies of web designers,” in *Proceedings of the SIGCHI Conference on Human*
2243 *Factors in Computing Systems*. ACM, 2010, pp. 703–712.
- 2244 [61] F. Doshi-Velez and B. Kim, “Towards a rigorous science of interpretable machine learning,”
2245 2017.
- 2246 [62] F. Doshi-Velez, M. Kortz, R. Budish, C. Bavitz, S. Gershman, D. O’Brien, S. Schieber, J. Waldo,
2247 D. Weinberger, and A. Wood, “Accountability of AI Under the Law: The Role of Explanation,”
2248 *arXiv.org*, Nov. 2017.
- 2249 [63] S. W. Draper. The Hawthorne, Pygmalion, Placebo and other effects of expectation: some notes.
- 2250 [64] R. G. Dromey, “A model for software product quality,” *International Software Engineering*
2251 *Research Network*, vol. 21, no. 2, pp. 146–162, 1995.
- 2252 [65] S. Easterbrook, J. Singer, M.-A. Storey, and D. Damian, “Selecting empirical methods for soft-
2253 ware engineering research,” in *Guide to Advanced Empirical Software Engineering*. Springer
2254 Science & Business Media, Nov. 2007, pp. 285–311.
- 2255 [66] B. Ehteshami Bejnordi, M. Veta, P. Johannes van Diest, B. van Ginneken, N. Karssemeijer,
2256 G. Litjens, J. A. W. M. van der Laak, and the CAMELYON16 Consortium, M. Hermsen, Q. F.
2257 Manson, M. Balkenhol, O. Geessink, N. Stathonikos, M. C. van Dijk, P. Bult, F. Beca, A. H.
2258 Beck, D. Wang, A. Khosla, R. Gargoya, H. Irshad, A. Zhong, Q. Dou, Q. Li, H. Chen, H.-J. Lin,
2259 P.-A. Heng, C. Haß, E. Bruni, Q. Wong, U. Halici, M. Ü. Öner, R. Cetin-Atalay, M. Berseth,
2260 V. Khvatkov, A. Vylegzhanin, O. Kraus, M. Shaban, N. Rajpoot, R. Awan, K. Sirinukunwattana,
2261 T. Qaiser, Y.-W. Tsang, D. Tellez, J. Annuscheit, P. Hufnagl, M. Valkonen, K. Kartasalo,
2262 L. Latonen, P. Ruusuviuri, K. Liimatainen, S. Albarqouni, B. Mungal, A. George, S. Demirci,

- 2263 N. Navab, S. Watanabe, S. Seno, Y. Takenaka, H. Matsuda, H. Ahmady Phoulady, V. Kovalev,
2264 A. Kalinovsky, V. Liauchuk, G. Bueno, M. M. Fernandez-Carrobles, I. Serrano, O. Deniz,
2265 D. Racoceanu, and R. Venâncio, “Diagnostic Assessment of Deep Learning Algorithms for
2266 Detection of Lymph Node Metastases in Women With Breast Cancer,” *JAMA*, vol. 318, no. 22,
2267 pp. 2199–12, Dec. 2017.
- 2268 [67] W. Elazmeh, W. Matwin, D. O’Sullivan, W. Michalowski, and W. Farion, “Insights from pre-
2269 dicting pediatric asthma exacerbations from retrospective clinical data,” in *Evaluation Methods
2270 for Machine Learning II—Papers from 2007 AAAI Workshop*, 2007, pp. 10–15.
- 2271 [68] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno,
2272 and D. Song, “Robust Physical-World Attacks on Deep Learning Visual Classification,” in
2273 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp.
2274 1625–1634.
- 2275 [69] ———, “Robust Physical-World Attacks on Deep Learning Visual Classification,” in *Proceedings
2276 of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1625–1634.
- 2277 [70] A. J. Feelders, “Prior knowledge in economic applications of data mining,” in *European Con-
2278 ference on Principles of Data Mining and Knowledge Discovery*. Springer, 2000, pp. 395–400.
- 2279 [71] R. T. Fielding, “Architectural Styles and the Design of Network-based Software Architectures ,”
2280 Ph.D. dissertation, University of California, Irvine.
- 2281 [72] I. Finalyson, “Nondeterministic Finite Automata,” 2018.
- 2282 [73] H. Foster, S. Uchitel, J. Magee, and J. Kramer, “Model-based verification of web service com-
2283 positions,” in *Automated Software Engineering, 2003. Proceedings. 18th IEEE International
2284 Conference on*. IEEE, 2003, pp. 152–161.
- 2285 [74] A. A. Freitas, “A critical review of multi-objective optimization in data mining: a position
2286 paper,” *ACM SIGKDD Explorations Newsletter*, vol. 6, no. 2, pp. 77–86, 2004.
- 2287 [75] ———, “Comprehensible classification models,” *ACM SIGKDD Explorations Newsletter*, vol. 15,
2288 no. 1, pp. 1–10, Mar. 2014.
- 2289 [76] A. A. Freitas, D. C. Wieser, and R. Apweiler, “On the importance of comprehensible classi-

- 2290 fication models for protein function prediction,” *IEEE/ACM Transactions on Computational*
2291 *Biology and Bioinformatics (TCBB)*, vol. 7, no. 1, pp. 172–182, 2010.
- 2292 [77] B. J. Frey and D. Dueck, “Clustering by Passing Messages Between Data Points,” *Science*, vol.
2293 315, no. 5814, pp. 972–976, Feb. 2007.
- 2294 [78] N. Friedman, D. Geiger, and M. Goldszmidt, “Bayesian network classifiers,” *Machine Learning*,
2295 vol. 29, no. 2-3, pp. 131–163, 1997.
- 2296 [79] G. Fung, S. Sandilya, and R. B. Rao, “Rule extraction from linear support vector machines,” in
2297 *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery*
2298 *in data mining*. ACM, 2005, pp. 32–40.
- 2299 [80] S. K. Garg, S. Versteeg, and R. Buyya, “SMICloud: A Framework for Comparing and Ranking
2300 Cloud Services,” in *2011 IEEE 4th International Conference on Utility and Cloud Computing*
2301 (*UCC 2011*). IEEE, Nov. 2011, pp. 210–218.
- 2302 [81] V. Garousi and M. Felderer, “Experience-based guidelines for effective and efficient data ex-
2303 traction in systematic reviews in software engineering,” in *Proceedings of the 21st International*
2304 *Conference on Evaluation and Assessment in Software Engineering*, ser. EASE’17. New York,
2305 NY, USA: ACM, 2017, pp. 170–179.
- 2306 [82] V. Garousi, M. Felderer, and M. V. Mäntylä, “Guidelines for including grey literature and
2307 conducting multivocal literature reviews in software engineering,” *Information and Software*
2308 *Technology*, vol. 106, pp. 101 – 121, 2019.
- 2309 [83] D. A. Garvin and W. D. P. Quality, “Really Mean,” *Sloan management review*, vol. 25, 1984.
- 2310 [84] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. Daumeé III, and
2311 K. Crawford, “Datasheets for datasets,” *arXiv preprint arXiv:1803.09010*, pp. 1–17, 2018.
- 2312 [85] H. L. Gilmore, “Product conformance cost,” *Quality progress*, vol. 7, no. 5, pp. 16–19, 1974.
- 2313 [86] R. L. Glass, I. Vessey, and V. Ramesh, “Research in software engineering: an analysis of the
2314 literature,” *Information and Software Technology*, vol. 44, no. 8, pp. 491–506, 2002.
- 2315 [87] M. W. Godfrey and D. M. German, “The past, present, and future of software evolution,” in
2316 *2008 Frontiers of Software Maintenance*, Sep. 2008, pp. 129–138.

- 2317 [88] B. Goodman and S. R. Flaxman, "EU regulations on algorithmic decision-making and a "right
2318 to explanation"." *IEEE Transactions on Evolutionary Computation*, 2016.
- 2319 [89] P. D. Grünwald, *The minimum description length principle*. MIT press, 2007.
- 2320 [90] M. J. Hadley, "Web application description language (WADL)," 2006.
- 2321 [91] H. A. Haenssle, C. Fink, R. Schneiderbauer, F. Toberer, T. Buhl, A. Blum, A. Kalloo, A. B. H.
2322 Hassen, L. Thomas, A. Enk, L. Uhlmann, Reader study level-I and level-II Groups, C. Alt,
2323 M. Arenbergerova, R. Bakos, A. Baltzer, I. Bertlich, A. Blum, T. Bokor-Billmann, J. Bowling,
2324 N. Braghiroli, R. Braun, K. Buder-Bakhaya, T. Buhl, H. Cabo, L. Cabrijan, N. Covic, A. Classen,
2325 D. Deltgen, C. Fink, I. Georgieva, L.-E. Hakim-Meibodi, S. Hanner, F. Hartmann, J. Hartmann,
2326 G. Haus, E. Hoxha, R. Karls, H. Koga, J. Kreusch, A. Lallas, P. Majenka, A. Marghoob,
2327 C. Massone, L. Mekokishvili, D. Mestel, V. Meyer, A. Neuberger, K. Nielsen, M. Oliviero,
2328 R. Pampena, J. Paoli, E. Pawlik, B. Rao, A. Rendon, T. Russo, A. Sadek, K. Samhaber,
2329 R. Schneiderbauer, A. Schweizer, F. Toberer, L. Trennheuser, L. Vlahova, A. Wald, J. Winkler,
2330 P. Wölbing, and I. Zalaudek, "Man against machine: diagnostic performance of a deep learning
2331 convolutional neural network for dermoscopic melanoma recognition in comparison to 58
2332 dermatologists," *Annals of Oncology*, vol. 29, no. 8, pp. 1836–1842, May 2018.
- 2333 [92] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*, ser. Data
2334 Mining, Inference, and Prediction. Springer Science & Business Media, Jan. 2001.
- 2335 [93] B. Hayete and J. R. Bienkowska, "GOTrees - Predicting GO Associations from Protein Domain
2336 Composition Using Decision Trees." *Pacific Symposium on Biocomputing*, pp. 127–138, 2005.
- 2337 [94] R. Heckel and M. Lohmann, "Towards contract-based testing of web services," *Electronic Notes
2338 in Theoretical Computer Science*, vol. 116, pp. 145–156, 2005.
- 2339 [95] D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. Kadie, "Dependency
2340 networks for inference, collaborative filtering, and data visualization," *Journal of Machine
2341 Learning Research*, vol. 1, no. Oct, pp. 49–75, 2000.
- 2342 [96] M. Henning, "API design matters," *Commun. ACM*, vol. 52, no. 5, pp. 46–56, 2009.
- 2343 [97] J. W. Horch, *Practical guide to software quality management*. Artech House, 2003.
- 2344 [98] H. Hosseini, B. Xiao, and R. Poovendran, "Google's Cloud Vision API is Not Robust to Noise,"

- 2345 in *2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA)*.
2346 IEEE, Jan. 2018, pp. 101–105.
- 2347 [99] J. Huang and C. X. Ling, “Using AUC and accuracy in evaluating learning algorithms,” *IEEE
2348 Transactions on knowledge and Data Engineering*, vol. 17, no. 3, pp. 299–310, 2005.
- 2349 [100] J. Huysmans, K. Dejaeger, C. Mues, J. Vanthienen, and B. Baesens, “An empirical evaluation
2350 of the comprehensibility of decision table, tree and rule based predictive models,” *Decision
2351 Support Systems*, vol. 51, no. 1, pp. 141–154, Apr. 2011.
- 2352 [101] K. Hwang, *Cloud Computing for Machine Learning and Cognitive Applications: A Machine
2353 Learning Approach*. MIT Press, 2017.
- 2354 [102] IEEE, “IEEE Standard Glossary of Software Engineering Terminology,” 1990.
- 2355 [103] International Organization for Standardization, “Systems and software engineering – Systems
2356 and software Quality Requirements and Evaluation (SQuaRE) – System and software quality
2357 models,” 2011.
- 2358 [104] ——, “Quality – Vocabulary,” ISO/IEC, 1986.
- 2359 [105] ——, “Quality management systems – Fundamentals and vocabulary,” ISO/IEC, 2015.
- 2360 [106] ——, “**Information technology – Software product quality**,” Nov. 1999.
- 2361 [107] A. Iyengar, “Supporting Data Analytics Applications Which Utilize Cognitive Services,” in
2362 *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE,
2363 May 2017, pp. 1856–1864.
- 2364 [108] N. Japkowicz and M. Shah, *Evaluating learning algorithms: a classification perspective*.
2365 Cambridge University Press, 2011.
- 2366 [109] M. W. M. Jaspers, M. Smeulers, H. Vermeulen, and L. W. Peute, “Effects of clinical decision-
2367 support systems on practitioner performance and patient outcomes: a synthesis of high-quality
2368 systematic review findings,” *Journal of the American Medical Informatics Association*, vol. 18,
2369 no. 3, pp. 327–334, 2011.
- 2370 [110] T. Jiang and A. E. Keating, “AVID: an integrative framework for discovering functional rela-
2371 tionships among proteins,” *BMC bioinformatics*, vol. 6, no. 1, p. 136, 2005.
- 2372 [111] T. Jick, *Mixing qualitative and quantitative methods*, ser. triangulation in action, 1979.

- 2373 [112] Y. Jin, *Multi-objective machine learning*. Springer Science & Business Media, 2006, vol. 16.
- 2374 [113] U. Johansson and L. Niklasson, “Evolving decision trees using oracle guides,” in *IEEE Symposium on Computational Intelligence and Data Mining (CIDM)*. IEEE, 2009.
- 2375 [114] C. M. Judd, E. R. Smith, and L. H. Kidder, “Research Methods in Social Relations, Fort Worth: Holt, Rinehart and Winston,” 1991.
- 2376 [115] J. M. Juran, *Juran on planning for quality*. Collier Macmillan, 1988.
- 2377 [116] N. Juristo and A. M. Moreno, *Basics of Software Engineering Experimentation*. Springer Science & Business Media, Mar. 2013.
- 2378 [117] A. Karwath and R. D. King, “Homology induction: the use of machine learning to improve sequence similarity searches,” *BMC bioinformatics*, vol. 3, no. 1, p. 11, 2002.
- 2379 [118] K. A. Kaufman and R. S. Michalski, “Learning from inconsistent and noisy data: the AQ18 approach,” in *International Symposium on Methodologies for Intelligent Systems*. Springer, 1999, pp. 411–419.
- 2380 [119] B. Kim, *Interactive and Interpretable Machine Learning Models for Human Machine Collaboration*. Massachusetts Institute of Technology, 2015.
- 2381 [120] B. Kim, C. Rudin, and J. A. Shah, “The Bayesian Case Model - A Generative Approach for Case-Based Reasoning and Prototype Classification.” *NIPS*, 2014.
- 2382 [121] B. Kitchenham and S. Charters, “Guidelines for performing Systematic Literature Reviews in Software Engineering,” Tech. Rep., 2007.
- 2383 [122] B. A. Kitchenham and S. L. Pfleeger, “Personal opinion surveys,” in *Guide to Advanced Empirical Software Engineering*. Springer Science & Business Media, Nov. 2007, pp. 63–92.
- 2384 [123] H. K. Klein and M. D. Myers, “A set of principles for conducting and evaluating interpretive field studies in information systems,” *MIS quarterly*, pp. 67–93, 1999.
- 2385 [124] A. J. Ko and Y. Riche, “The role of conceptual knowledge in API usability,” in *2011 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. IEEE, 2011, pp. 173–176.
- 2386 [125] A. J. Ko, B. A. Myers, and H. H. Aung, “Six learning barriers in end-user programming

- systems,” in *2004 IEEE Symposium on Visual Languages-Human Centric Computing*. IEEE, 2004, pp. 199–206.
- [126] I. Kononenko, “Inductive and Bayesian learning in medical diagnosis,” *Applied Artificial Intelligence an International Journal*, vol. 7, no. 4, pp. 317–337, 1993.
- [127] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet Classification with Deep Convolutional Neural Networks.” 2012.
- [128] J. A. Krosnick, “Survey research,” *Annual Review of Psychology*, vol. 50, no. 1, pp. 537–567, 1999.
- [129] A. Kurakin, I. Goodfellow, and S. Bengio, “Adversarial examples in the physical world,” *arXiv preprint arXiv:1607.02533*, vol. cs.CV, 2016.
- [130] G. Laforge, “Machine Intelligence at Google Scale,” in *QCon*, Jun. 2018, pp. 1–58.
- [131] H. Lakkaraju, S. H. Bach, and J. Leskovec, “Interpretable Decision Sets - A Joint Framework for Description and Prediction.” *KDD*, pp. 1675–1684, 2016.
- [132] J. R. Landis and G. G. Koch, “The Measurement of Observer Agreement for Categorical Data,” *Biometrics*, vol. 33, no. 1, pp. 159–17, Mar. 1977.
- [133] F. Lau, “Toward a framework for action research in information systems studies,” *Information Technology & People*, vol. 12, no. 2, pp. 148–176, 1999.
- [134] N. Lavrač, “Selected techniques for data mining in medicine,” *Artificial intelligence in medicine*, vol. 16, no. 1, pp. 3–23, 1999.
- [135] T. Lei, R. Barzilay, and T. Jaakkola, “Rationalizing Neural Predictions,” *arXiv.org*, Jun. 2016.
- [136] T. C. Lethbridge, S. E. Sim, and J. Singer, “Studying Software Engineers: Data Collection Techniques for Software Field Studies,” *Empirical Software Engineering*, vol. 10, no. 3, pp. 311–341, Jul. 2005.
- [137] E. Lima, C. Mues, and B. Baesens, “Domain knowledge integration in data mining using decision tables: Case studies in churn prediction,” *Journal of the Operational Research Society*, vol. 60, no. 8, pp. 1096–1106, 2009.
- [138] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft COCO: Common Objects in Context,” in *Computer Vision – ECCV 2014*, D. Fleet,

- 2428 T. Pajdla, B. Schiele, and T. Tuytelaars, Eds. Cham: Springer International Publishing, 2014,

2429 pp. 740–755.

2430 [139] M. Linares-Vásquez, G. Bavota, M. Di Penta, R. Oliveto, and D. Poshyvanyk, “How do api
2431 changes trigger stack overflow discussions? a study on the android sdk,” in *proceedings of the*
2432 *22nd International Conference on Program Comprehension*. ACM, 2014, pp. 83–94.

2433 [140] Z. C. Lipton, “The Mythos of Model Interpretability.” *IEEE Transactions on Evolutionary*
2434 *Computation*, 2016.

2435 [141] M. S. Litwin and A. Fink, *How to measure survey reliability and validity*. Sage, 1995, vol. 7.

2436 [142] Y. Liu, T. Kohlberger, M. Norouzi, G. E. Dahl, J. L. Smith, A. Mohtashamian, N. Olson,
2437 L. H. Peng, J. D. Hipp, and M. C. Stumpe, “Artificial Intelligence-Based Breast Cancer Nodal
2438 Metastasis Detection,” *Archives of Pathology & Laboratory Medicine*, pp. arpa.2018–0147–
2439 OA–11, Oct. 2018.

2440 [143] D. Lo Giudice, C. Mines, A. LeClair, R. Curran, and A. Homan, “How AI Will Change Software
2441 Development And Applications,” Tech. Rep., Nov. 2016.

2442 [144] W. Maalej and M. P. Robillard, “Patterns of Knowledge in API Reference Documentation,”
2443 *International Software Engineering Research Network*, vol. 39, no. 9, pp. 1264–1282.

2444 [145] S. MacDonell, M. Shepperd, B. Kitchenham, and E. Mendes, “How reliable are systematic
2445 reviews in empirical software engineering?” *IEEE Transactions on Software Engineering*,
2446 vol. 36, no. 5, pp. 676–687, Sep. 2010.

2447 [146] T. E. Marshall and S. L. Lambert, “Cloud-based intelligent accounting applications: accounting
2448 task automation using IBM watson cognitive computing,” *Journal of Emerging Technologies*
2449 *in Accounting*, vol. 15, no. 1, pp. 199–215, 2018.

2450 [147] D. Martens, J. Vanthienen, W. Verbeke, and B. Baesens, “Performance of classification models
2451 from a user perspective,” *Decision Support Systems*, vol. 51, no. 4, pp. 782–793, 2011.

2452 [148] J. A. McCall, “Factors in software quality,” *US Rome Air development center reports*, 1977.

2453 [149] J. A. McCall, P. K. Richards, and G. F. Walters, “Factors in software quality. volume i. concepts
2454 and definitions of software quality,” Tech. Rep., 1977.

2455 [150] J. McCarthy, “Programs with Common Sense,” Cambridge, MA, USA, Tech. Rep., 1960.

- ²⁴⁵⁶ [151] L. McLeod and S. G. MacDonell, “Factors that affect software systems development project
²⁴⁵⁷ outcomes: A survey of research,” *ACM Computing Surveys (CSUR)*, vol. 43, no. 4, p. 24, 2011.
- ²⁴⁵⁸ [152] J. Meltzoff, *Critical thinking about research: Psychology and related fields*. American
²⁴⁵⁹ psychological association, 1998.
- ²⁴⁶⁰ [153] T. Mens, M. Wermelinger, S. Ducasse, S. Demeyer, R. Hirschfeld, and M. Jazayeri, “Challenges
²⁴⁶¹ in software evolution,” in *Eighth International Workshop on Principles of Software Evolution
(IWPSE’05)*, Sep. 2005, pp. 13–22.
- ²⁴⁶³ [154] D. Michie, D. J. Spiegelhalter, and C. C. Taylor, “Machine Learning and Statistical Classification
²⁴⁶⁴ of Artificial intelligence,” 1994.
- ²⁴⁶⁵ [155] D. Michie, “Machine Learning in the Next Five Years.” *EWSL*, 1988.
- ²⁴⁶⁶ [156] M. B. Miles, A. M. Huberman, M. A. Huberman, and M. Huberman, *Qualitative data analysis:
An expanded sourcebook*. sage, 1994.
- ²⁴⁶⁸ [157] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D.
²⁴⁶⁹ Raji, and T. Gebru, “Model cards for model reporting,” *arXiv preprint arXiv:1810.03993*, pp.
²⁴⁷⁰ 220–229, 2018.
- ²⁴⁷¹ [158] D. L. Moody, “The ‘Physics’ of Notations - Toward a Scientific Basis for Constructing Visual
²⁴⁷² Notations in Software Engineering.” *IEEE Trans. Software Eng.*, 2009.
- ²⁴⁷³ [159] C. Murphy and G. E. Kaiser, “Improving the dependability of machine learning applications,”
²⁴⁷⁴ 2008.
- ²⁴⁷⁵ [160] C. Murphy, G. E. Kaiser, and M. Arias, “An approach to software testing of machine learning
²⁴⁷⁶ applications,” 2007.
- ²⁴⁷⁷ [161] B. A. Myers, S. Y. Jeong, Y. Xie, J. Beaton, J. Stylos, R. Ehret, J. Karstens, A. Efeoglu, and D. K.
²⁴⁷⁸ Busse, “Studying the documentation of an API for enterprise Service-Oriented Architecture,”
²⁴⁷⁹ *Journal of Organizational and End User Computing (JOEUC)*, vol. 22, no. 1, pp. 23–51, 2010.
- ²⁴⁸⁰ [162] S. Nakajima, “Model-checking verification for reliable web service,” in *OOPSLA 2002 Work-
shop on Object-Oriented Web Services, Seattle, Washington*, 2002.
- ²⁴⁸² [163] M. Narayanan, E. Chen, J. He, B. Kim, S. Gershman, and F. Doshi-Velez, “How do Humans

- 2483 Understand Explanations from Machine Learning Systems? An Evaluation of the Human-
2484 Interpretability of Explanation.” *IEEE Transactions on Evolutionary Computation*, 2018.
- 2485 [164] S. Narayanan and S. A. McIlraith, “Simulation, verification and automated composition of web
2486 services,” in *Proceedings of the 11th international conference on World Wide Web*. ACM,
2487 2002, pp. 77–88.
- 2488 [165] B. J. Nelson, “Remote procedure call,” Ph.D. dissertation, Carnegie Mellon University, 1981.
- 2489 [166] Y. Nishi, S. Masuda, H. Ogawa, and K. Uetsuki, “A test architecture for machine learning prod-
2490 uct,” in *2018 IEEE International Conference on Software Testing, Verification and Validation*
2491 *Workshops (ICSTW)*. IEEE, 2018, pp. 273–278.
- 2492 [167] N. Novielli, F. Calefato, and F. Lanobile, “The challenges of sentiment detection in the social
2493 programmer ecosystem,” in *Proceedings of the 7th International Workshop on Social Software*
2494 *Engineering*. ACM, 2015, pp. 33–40.
- 2495 [168] K. Nybom, A. Ashraf, and I. Porres, “A Systematic Mapping Study on API Documentation
2496 Generation Approaches,” in *2018 44th Euromicro Conference on Software Engineering and*
2497 *Advanced Applications (SEAA)*, Prague, Czech Republic, pp. 462–469.
- 2498 [169] J. Nykaza, R. Messinger, F. Boehme, C. L. Norman, M. Mace, and M. Gordon, “What pro-
2499 grammers really want - results of a needs assessment for SDK documentation.” *SIGDOC*,
2500 2002.
- 2501 [170] T. Ohtake, A. Cummaudo, M. Abdelrazek, R. Vasa, and J. Grundy, “Merging Intelligent API
2502 Responses Using a Proportional Representation Approach,” in *Empirical Software Engineering*
2503 *and Verification*. Cham: Springer International Publishing, Apr. 2019, pp. 391–406.
- 2504 [171] Open Software Foundation, “Part 3: DCE Remote Procedure Call (RPC),” in *OSF DCE*
2505 *application development guide: revision 1.0*. Prentice Hall, Dec. 1991.
- 2506 [172] N. Oreskes, K. Shrader-Frechette, and K. Belitz, “Verification, Validation, and Confirmation of
2507 Numerical Models in the Earth Sciences,” *Science*, vol. 263, no. 5147, pp. 641–646, 1994.
- 2508 [173] A. L. M. Ortiz, “Curating Content with Google Machine Learning Application Programming
2509 Interfaces,” in *EIAPortugal*, Jul. 2017.
- 2510 [174] F. E. Otero and A. A. Freitas, “Improving the interpretability of classification rules discovered

- 2511 by an ant colony algorithm,” in *Proceedings of the 15th annual conference on Genetic and*
2512 *evolutionary computation*. ACM, 2013, pp. 73–80.
- 2513 [175] A. Pal, S. Chang, and J. A. Konstan, “Evolution of experts in question answering communities.”
2514 in *ICWSM*, 2012.
- 2515 [176] A. Parasuraman, V. A. Zeithaml, and L. L. Berry, “Servqual: A multiple-item scale for mea-
2516 suring consumer perceptions of service quality,” *Journal of retailing*, vol. 64, no. 1, pp. 12–29,
2517 1988.
- 2518 [177] R. Parekh, “Designing AI at Scale to Power Everyday Life,” in *the 23rd ACM SIGKDD*
2519 *International Conference*. New York, New York, USA: ACM Press, 2017, pp. 27–27.
- 2520 [178] C. Pautasso, O. Zimmermann, and F. Leymann, “**RESTful Web Services vs. “Big” Web Ser-**
2521 **vices: Making the Right Architectural Decision** ,” in *Proceedings of the 17th international*
2522 *conference on World Wide Web*. ACM, 2008, pp. 805–814.
- 2523 [179] M. Pazzani, “Comprehensible knowledge discovery: gaining insight from data,” in *First Federal*
2524 *Data Mining Conference and Exposition*, 1997, pp. 73–82.
- 2525 [180] M. J. Pazzani, S. Mani, and W. R. Shankle, “Acceptance of rules generated by machine learning
2526 among medical experts,” *Methods of information in medicine*, vol. 40, no. 05, pp. 380–385,
2527 2001.
- 2528 [181] J. Pearl, “The Seven Tools of Causal Inference with Reflections on Machine Learning,” 2018.
- 2529 [182] K. Petersen and C. Gencel, “Worldviews, Research Methods, and their Relationship to Validity in
2530 Empirical Software Engineering Research,” in *2013 Joint Conference of the 23nd International*
2531 *Workshop on Software Measurement and the 8th International Conference on Software Process*
2532 *and Product Measurement (IWSM-MENSURA)*. IEEE, Jan. 2019, pp. 81–89.
- 2533 [183] K. Petersen, R. Feldt, S. Mujtaba, and M. Mattsson, “Systematic Mapping Studies in Software
2534 Engineering.” in *EASE*, 2008, pp. 68–77.
- 2535 [184] Z. Pezzementi, T. Tabor, S. Yim, J. K. Chang, B. Drozd, D. Guttendorf, M. Wagner, and
2536 P. Koopman, “Putting image manipulations in context: robustness testing for safe perception,”
2537 in *2018 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR)*. IEEE,
2538 2018, pp. 1–8.

- 2539 [185] H. Pham, *Software reliability*. Springer Science & Business Media, 2000.
- 2540 [186] M. Piccioni, C. A. Furia, and B. Meyer, “An Empirical Study of API Usability,” in *2013 ACM /*
2541 *IEEE International Symposium on Empirical Software Engineering and Measurement*. IEEE,
2542 2013, pp. 5–14.
- 2543 [187] R. M. Pirsig, *Zen and the art of motorcycle maintenance: An inquiry into values*. William
2544 Morrow and Company, 1974.
- 2545 [188] R. S. Pressman, *Software engineering: a practitioner’s approach*. Palgrave Macmillan, 2005.
- 2546 [189] J. R. Quinlan, “C4. 5: Programming for machine learning,” *Morgan Kauffmann*, vol. 38, p. 48,
2547 1993.
- 2548 [190] ——, “Some elements of machine learning,” in *International Conference on Inductive Logic*
2549 *Programming*. Springer, 1999, pp. 15–18.
- 2550 [191] M. Rebouças, G. Pinto, F. Ebert, W. Torres, A. Serebrenik, and F. Castor, “An empirical
2551 study on the usage of the swift programming language,” in *Software Analysis, Evolution, and*
2552 *Reengineering (SANER), 2016 IEEE 23rd International Conference on*. IEEE, 2016, pp.
2553 634–638.
- 2554 [192] A. Reis, D. Paulino, V. Filipe, and J. Barroso, “Using Online Artificial Vision Services to
2555 Assist the Blind - an Assessment of Microsoft Cognitive Services and Google Cloud Vision.”
2556 *WorldCIST*, vol. 746, no. 12, pp. 174–184, 2018.
- 2557 [193] M. T. Ribeiro, S. Singh, and C. Guestrin, ““Why Should I Trust You?”,” in *the 22nd ACM*
2558 *SIGKDD International Conference*. New York, New York, USA: ACM Press, 2016, pp.
2559 1135–1144.
- 2560 [194] M. Ribeiro, K. Grolinger, and M. A. M. Capretz, “MLaaS: Machine Learning as a Service,”
2561 in *2015 IEEE 14th International Conference on Machine Learning and Applications (ICMLA)*.
2562 IEEE, Dec. 2015, pp. 896–902.
- 2563 [195] G. Richards, V. J. Rayward-Smith, P. H. Sönksen, S. Carey, and C. Weng, “Data mining for
2564 indicators of early mortality in a database of clinical records,” *Artificial intelligence in medicine*,
2565 vol. 22, no. 3, pp. 215–231, 2001.

- 2566 [196] G. Ridgeway, D. Madigan, T. Richardson, and J. O’Kane, “Interpretable Boosted Naïve Bayes
2567 Classification.” *KDD*, 1998.
- 2568 [197] RightScale Inc., “RightScale 2018 State of the Cloud Report,” Tech. Rep., 2018.
- 2569 [198] G. Ritzer and E. Guba, “The Paradigm Dialog,” *Canadian Journal of Sociology / Cahiers
2570 canadiens de sociologie*, vol. 16, no. 4, p. 446, 1991.
- 2571 [199] S. P. Robbins and T. Judge, *Essentials of organizational behavior*. Pearson,, 2014.
- 2572 [200] M. P. Robillard, “What makes APIs hard to learn? Answers from developers,” *IEEE Software*,
2573 vol. 26, no. 6, pp. 27–34, 2009.
- 2574 [201] M. P. Robillard and R. Deline, “A field study of API learning obstacles,” *Empirical Software
2575 Engineering*, vol. 16, no. 6, pp. 703–732, 2011.
- 2576 [202] H. Robinson, J. Segal, and H. Sharp, “Ethnographically-informed empirical studies of software
2577 practice,” *Information and Software Technology*, vol. 49, no. 6, pp. 540–551, 2007.
- 2578 [203] L. Rokach and O. Z. Maimon, *Data mining with decision trees: theory and applications*. World
2579 scientific, 2008, vol. 69.
- 2580 [204] C. Rosen and E. Shihab, “What are mobile developers asking about? A large scale study using
2581 stack overflow,” *Empirical Software Engineering*, vol. 21, no. 3, pp. 1192–1223, 2016.
- 2582 [205] W. Rosenberry, D. Kenney, and G. Fisher, *Understanding DCE*. O’Reilly & Associates, Inc.,
2583 1992.
- 2584 [206] A. Rosenfeld, R. Zemel, and J. K. Tsotsos, “The elephant in the room,” *arXiv preprint
2585 arXiv:1808.03305*, 2018.
- 2586 [207] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, “Right for the Right Reasons: Training Differ-
2587 entiable Models by Constraining their Explanations,” *arXiv.org*, Mar. 2017.
- 2588 [208] R. J. Rubey and R. D. Hartwick, “Quantitative measurement of program quality,” in *Proceedings
2589 of the 1968 23rd ACM national conference*. New York, New York, USA: ACM, 1968, pp.
2590 671–677.
- 2591 [209] J. Rutten, M. Kwiatkowska, G. Norman, and D. Parker, *Mathematical Techniques for Analyzing
2592 Concurrent and Probabilistic Systems*, P. Panangaden and F. van Breugel (eds.), ser. CRM
2593 Monograph Series. American Mathematical Society, 2004, vol. 23.

- 2594 [210] H. W. Schmidt, I. Crnkovic, G. T. Heineman, and J. A. Stafford, Eds., *A Framework for Contract-*
2595 *Based Collaborative Verification and Validation of Web Services*. Berlin, Heidelberg: Springer
2596 Berlin Heidelberg, 2007.
- 2597 [211] M. Schwabacher, P. Langley, and P. Norvig, “Discovering communicable scientific knowledge
2598 from spatio-temporal data,” *ICML*, 2001.
- 2599 [212] C. B. Seaman, “Qualitative methods,” in *Guide to Advanced Empirical Software Engineering*.
2600 Springer Science & Business Media, Nov. 2007, pp. 35–62.
- 2601 [213] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM:
2602 Visual Explanations from Deep Networks via Gradient-Based Localization,” in *2017 IEEE*
2603 *International Conference on Computer Vision (ICCV)*. IEEE, 2017, pp. 618–626.
- 2604 [214] S. Sen and L. Knight, “A genetic prototype learner,” in *IJCAI*. Citeseer, 1995, pp. 725–733.
- 2605 [215] C. E. Shannon and W. Weaver, *The mathematical theory of communication*. Urbana, IL: The
2606 University of Illinois Press, 1963.
- 2607 [216] M. Shaw, “Writing good software engineering research papers,” in *25th International Confer-*
2608 *ence on Software Engineering, 2003. Proceedings*. IEEE, 2003, pp. 726–736.
- 2609 [217] Shull, Forrest, Singer, Janice, and Sjøberg, Dag I K, *Guide to Advanced Empirical Software*
2610 *Engineering*, F. Shull, J. Singer, and D. I. K. Sjøberg, Eds. Springer Science & Business
2611 Media, Nov. 2007.
- 2612 [218] H. A. Simon, *The sciences of the artificial*. MIT press, 1996.
- 2613 [219] J. Singer, S. E. Sim, and T. C. Lethbridge, “Software engineering data collection for field
2614 studies,” in *Guide to Advanced Empirical Software Engineering*. Springer Science & Business
2615 Media, Nov. 2007, pp. 9–34.
- 2616 [220] S. Singh, M. T. Ribeiro, and C. Guestrin, “Programs as Black-Box Explanations,” *arXiv.org*,
2617 Nov. 2016.
- 2618 [221] V. S. Sinha, S. Mani, and M. Gupta, “Exploring activeness of users in QA forums,” in *Proceed-*
2619 *ings of the 10th Working Conference on Mining Software Repositories*. IEEE Press, 2013, pp.
2620 77–80.

- 2621 [222] M. Sokolova and G. Lapalme, “A systematic analysis of performance measures for classification
2622 tasks,” *Information Processing & Management*, vol. 45, no. 4, pp. 427–437, 2009.
- 2623 [223] I. Sommerville, *Software Engineering*, 9th ed. Addison-Wesley, 2011.
- 2624 [224] P. E. Spector, *Summated rating scale construction: An introduction*. Sage, 1992.
- 2625 [225] J. Su, D. V. Vargas, and K. Sakurai, “One pixel attack for fooling deep neural networks.” *IEEE*
2626 *Transactions on Evolutionary Computation*, 2019.
- 2627 [226] G. H. Subramanian, J. Nosek, S. P. Raghunathan, and S. S. Kanitkar, “A comparison of the
2628 decision table and tree,” *Communications of the ACM*, vol. 35, no. 1, pp. 89–94, 1992.
- 2629 [227] S. Subramanian, L. Inozemtseva, and R. Holmes, “Live API documentation,” in *the 36th*
2630 *International Conference*. New York, New York, USA: ACM Press, 2014, pp. 643–652.
- 2631 [228] M. Sugiyama, N. D. Lawrence, and A. Schwaighofer, *Dataset shift in machine learning*. The
2632 MIT Press, 2017.
- 2633 [229] N. R. Suri, V. S. Srinivas, and M. N. Murty, “A cooperative game theoretic approach to prototype
2634 selection,” in *European Conference on Principles of Data Mining and Knowledge Discovery*.
2635 Springer, 2007, pp. 556–564.
- 2636 [230] D. Szafron, P. Lu, R. Greiner, D. S. Wishart, B. Poulin, R. Eisner, Z. Lu, J. Anvik, C. Macdonell,
2637 and A. Fyshe, “Proteome Analyst: custom predictions with explanations in a web-based tool
2638 for high-throughput proteome annotations,” *Nucleic acids research*, vol. 32, no. 2, pp. W365–
2639 W371, 2004.
- 2640 [231] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus,
2641 “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013.
- 2642 [232] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the Inception Architec-
2643 ture for Computer Vision.” *25th IEEE Conference on Computer Vision and Pattern Recognition*,
2644 2016.
- 2645 [233] M. B. W. Tabor, “**Student Proves That S.A.T. Can Be: (D) Wrong**,” *New York Times*, Feb.
2646 1997.
- 2647 [234] A. Tahir, A. Yamashita, S. Licorish, J. Dietrich, and S. Counsell, “Can you tell me if it smells?:

- 2648 A study on how developers discuss code smells and anti-patterns in Stack Overflow,” *the 22nd*
2649 *International Conference*, pp. 68–78, 2018.
- 2650 [235] G. Tassey, *The Economic Impacts of Inadequate Infrastructure for Software Testing*. National
2651 Institute of Standards and Technology, Sep. 2002.
- 2652 [236] R. S. Taylor, “Question-negotiation and information-seeking in libraries (Vol. 29): College and
2653 Research Libraries,” 1968.
- 2654 [237] S. W. Thomas, B. Adams, A. E. Hassan, and D. Blostein, “Studying software evolution using
2655 topic models,” *Science of Computer Programming*, vol. 80, pp. 457 – 479, 2014.
- 2656 [238] S. Thrun, “Is Learning The n-th Thing Any Easier Than Learning The First?” p. 7, 1996.
- 2657 [239] Q. Tu *et al.*, “Evolution in open source software: A case study,” in *Proceedings 2000 Interna-*
2658 *tional Conference on Software Maintenance*. IEEE, 2000, pp. 131–142.
- 2659 [240] M. Usman, R. Britto, J. Börstler, and E. Mendes, “Taxonomies in software engineering: A
2660 Systematic mapping study and a revised taxonomy development method,” *Information and*
2661 *Software Technology*, vol. 85, pp. 43–59, May 2017.
- 2662 [241] A. Van Assche and H. Blockeel, “Seeing the forest through the trees: Learning a comprehensible
2663 model from an ensemble,” in *European conference on machine learning*. Springer, 2007, pp.
2664 418–429.
- 2665 [242] B. Venners, “Design by Contract: A Conversation with Bertrand Meyer,” *Artima Developer*,
2666 2003.
- 2667 [243] W. Verbeke, D. Martens, C. Mues, and B. Baesens, “Building comprehensible customer churn
2668 prediction models with advanced rule induction techniques,” *Expert Systems with Applications*,
2669 vol. 38, no. 3, pp. 2354–2364, 2011.
- 2670 [244] S. Wachter, B. Mittelstadt, and L. Floridi, “Why a Right to Explanation of Automated Decision-
2671 Making Does Not Exist in the General Data Protection Regulation,” *International Data Privacy*
2672 *Law*, vol. 7, no. 2, pp. 76–99, Jun. 2017.
- 2673 [245] B. Wang, Y. Yao, B. Viswanath, H. Zheng, and B. Y. Zhao, “With Great Training Comes Great
2674 Vulnerability - Practical Attacks against Transfer Learning.” *USENIX Security Symposium*,
2675 2018.

- 2676 [246] S. Wang, D. Lo, and L. Jiang, “An empirical study on developer interactions in StackOverflow,”
2677 in *Proceedings of the 28th Annual ACM Symposium on Applied Computing*. ACM, 2013, pp.
2678 1019–1024.
- 2679 [247] W. Wang, H. Malik, and M. W. Godfrey, “Recommending Posts concerning API Issues in
2680 Developer Q&A Sites,” in *2015 IEEE/ACM 12th Working Conference on Mining Software
2681 Repositories*, May 2015, pp. 224–234.
- 2682 [248] R. B. Watson, “Development and application of a heuristic to assess trends in API documenta-
2683 tion,” in *30th ACM international conference on Design of communication*. Seattle, Washington,
2684 USA: ACM, 2012, pp. 295–302.
- 2685 [249] S. Weerawarana, *Web Services Platform Architecture*, ser. SOAP, WSDL, WS-Policy, WS-
2686 Addressing, WS-BPEL, WS-Reliable Messaging, and More. Prentice-Hall PTR, 2005.
- 2687 [250] D. Wettschereck, D. W. Aha, and T. Mohri, “A review and empirical evaluation of feature
2688 weighting methods for a class of lazy learning algorithms,” *Artificial Intelligence Review*,
2689 vol. 11, no. 1-5, pp. 273–314, 1997.
- 2690 [251] H. Wickham, “A Layered Grammar of Graphics,” *Journal of Computational and Graphical
2691 Statistics*, vol. 19, no. 1, pp. 3–28, Jan. 2010.
- 2692 [252] R. J. Wieringa and J. M. Heerkens, “The methodological soundness of requirements engineering
2693 papers: a conceptual framework and two case studies,” *Requirements engineering*, vol. 11, no. 4,
2694 pp. 295–307, 2006.
- 2695 [253] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *Data Mining: Practical machine learning
2696 tools and techniques*. Morgan Kaufmann, 2016.
- 2697 [254] C. Wohlin and A. Aurum, “Towards a decision-making structure for selecting a research design
2698 in empirical software engineering,” *Empirical Software Engineering*, vol. 20, no. 6, pp. 1427–
2699 1455, May 2014.
- 2700 [255] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslén, *Experimentation
2701 in Software Engineering*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012.
- 2702 [256] M. L. Wong and K. S. Leung, *Data mining using grammar based genetic programming and
2703 applications*. Springer Science & Business Media, 2006, vol. 3.

- 2704 [257] X. Yi and K. J. Kochut, "A cp-nets-based design and verification framework for web services
2705 composition," in *Web Services, 2004. Proceedings. IEEE International Conference on.* IEEE,
2706 2004, pp. 756–760.
- 2707 [258] R. K. Yin, *Case study research and applications: Design and methods.* Sage publications,
2708 2017.
- 2709 [259] J. Zahálka and F. Železný, "An experimental test of Occam's razor in classification," *Machine
2710 Learning*, vol. 82, no. 3, pp. 475–481, 2011.
- 2711 [260] J. Zhang and R. Kasturi, "Extraction of Text Objects in Video Documents: Recent Progress," in
2712 *2008 The Eighth IAPR International Workshop on Document Analysis Systems (DAS).* IEEE,
2713 2008, pp. 5–17.
- 2714 [261] B. Zupan, J. Demšar, M. W. Kattan, J. R. Beck, and I. Bratko, "Machine learning for survival
2715 analysis: a case study on recurrence of prostate cancer," *Artificial intelligence in medicine*,
2716 vol. 20, no. 1, pp. 59–75, 2000.
- 2717 [262] M. zur Muehlen, J. V. Nickerson, and K. D. Swenson, "Developing web services choreography
2718 standards—the case of REST vs. SOAP," *Decision Support Systems*, vol. 40, no. 1, pp. 9–29,
2719 Jul. 2005.

2720

Online Artefacts

2721

- 2722
-
- 2723 [263] C. Howard, “Introducing Google AI,” <https://ai.googleblog.com/2018/05/introducing-google-ai.html>, May 2018.
- 2724 [264] “Amazon Mechanical Turk,” <https://www.mturk.com>, accessed: 15 October 2018.
- 2725 [265] “Scale: API for Training Data,” <https://www.scaleapi.com>, accessed: 15 October 2018.
- 2726 [266] “Vision API - Image Content Analysis | Cloud Vision API | Google Cloud,” <http://bit.ly/2TD9mBs>, accessed: 13 September 2018.
- 2727 [267] “Image Processing with the Computer Vision API | Microsoft Azure,” <http://bit.ly/2YqhkS6>, accessed: 13 September 2018.
- 2728 [268] “Amazon Rekognition,” <https://amzn.to/2TyT2BL>, accessed: 13 September 2018.
- 2729 [269] “Detecting labels in an image,” <http://bit.ly/2UlkW9K>, accessed: 13 September 2018.
- 2730 [270] “Watson visual recognition,” <https://ibm.co/2TBNIO4>, accessed: 13 September 2018.
- 2731 [271] “Image Recognition API & Visual Search Results,” <http://bit.ly/2UmNPCw>, accessed: 13 September 2018.
- 2732 [272] “Clarifai,” <http://bit.ly/2TB3kSa>, accessed: 13 September 2018.
- 2733 [273] “Image Recognition API | DeepAI,” <http://bit.ly/2TBNYgf>, accessed: 26 September 2018.

- ²⁷³⁸ [274] “Imagga - powerful image recognition apis for automated categorization & tagging in the cloud
²⁷³⁹ and on-premises,” <http://bit.ly/2TxsyRe>, accessed: 13 September 2018.
- ²⁷⁴⁰ [275] “Image recognition - talkwalker,” <http://bit.ly/2TyT7W5>, accessed: 13 September 2018.
- ²⁷⁴¹ [276] “Kairos: Serving businesses with face recognition,” <https://www.kairos.com>, accessed: 15
²⁷⁴² October 2018.
- ²⁷⁴³ [277] “The face recognition company - cognitec,” <http://www.cognitec.com>, accessed: 15 October
²⁷⁴⁴ 2018.
- ²⁷⁴⁵ [278] “Home - affectiva : Affectiva,” <https://www.affectiva.com>, accessed: 15 October 2018.
- ²⁷⁴⁶ [279] “Detect research-methodology | Google Cloud Vision API Documentation | Google Cloud,”
²⁷⁴⁷ <https://cloud.google.com/vision/docs/research-methodology>, accessed: 28 August 2018.
- ²⁷⁴⁸ [280] “Class EntityAnnotation | Google.Cloud.Vision.V1,” <https://googlecloudplatform.github.io/google-cloud-dotnet/docs/Google.Cloud.Vision.V1/api/Google.Cloud.Vision.V1.EntityAnnotation.html>, accessed: 28 August 2018.
- ²⁷⁵¹ [281] “How to call the Computer Vision API,” <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/vision-api-how-to-topics/howtocallvisionapi>, accessed:
²⁷⁵²
²⁷⁵³ 28 August 2018.
- ²⁷⁵⁴ [282] “azure-sdk-for-java/ImageTag.java,” <https://github.com/Azure/azure-sdk-for-java/blob/8d70f41fdb8b3a9297af5ba24551cf26f40ad4/cognitiveservices/data-plane/vision/computervision/src/main/java/com/microsoft/azure/cognitiveservices/vision/computervision/models/ImageTag.java#L24>, accessed: 28 August 2018.
- ²⁷⁵⁸ [283] “Detecting research-methodology in an image,” <https://docs.aws.amazon.com/rekognition/latest/dg/research-methodology-detect-research-methodology-image.html>, accessed: 28 Au-
²⁷⁵⁹
²⁷⁶⁰ gust 2018.
- ²⁷⁶¹ [284] “Detecting objects and scenes,” <https://docs.aws.amazon.com/rekognition/latest/dg/research-methodology.html>, accessed: 28 August 2018.
- ²⁷⁶³ [285] “Megvii,” <http://bit.ly/2WJYFzk>, accessed: 3 April 2019.
- ²⁷⁶⁴ [286] “Tuputech,” <http://bit.ly/2uF4IsN>, accessed: 3 April 2019.
- ²⁷⁶⁵ [287] “Yitu technology,” <http://bit.ly/2uGvxgf>, accessed: 3 April 2019.

- [2766](#) [288] “Sensetime,” <http://bit.ly/2WH6RjF>, accessed: 3 April 2019.
- [2767](#) [289] “Detect Labels | Google Cloud Vision API Documentation | Google Cloud,” <http://bit.ly/2TD5kcy>, accessed: 28 August 2018.
- [2768](#)
- [2769](#) [290] “Detecting labels in an image,” <https://amzn.to/2TBNtTa>, accessed: 28 August 2018.
- [2770](#) [291] “How to call the Computer Vision API,” <http://bit.ly/2TD5oJk>, accessed: 28 August 2018.
- [2771](#) [292] “Deepglint,” <http://bit.ly/2uHHdPS>, accessed: 3 April 2019.
- [2772](#) [293] <http://bit.ly/2KlyhcD>, accessed: 27 March 2019.
- [2773](#) [294] <http://bit.ly/2G6ZOeC>, accessed: 27 March 2019.
- [2774](#) [295] <http://bit.ly/2G7saFJ>, accessed: 27 March 2019.
- [2775](#) [296] “What is Computer Vision?” <http://bit.ly/2TDgUnU>, accessed: 28 August 2018.
- [2776](#) [297] <http://bit.ly/2G5ZEEe>, accessed: 27 March 2019.
- [2777](#) [298] K. Ballinger, “Simplicity and utility, or, why soap lost,” <http://keithba.net/simplicity-and-utility-or-why-soap-lost>, Dec. 2014.
- [2778](#)
- [2779](#) [299] L. Mandel, “Describe REST Web services with WSDL 2.0,” <https://www.ibm.com/developerworks/webservices/library/ws-restwsdl/>, May 2008.

Part IV

Appendices

APPENDIX A

Additional Materials

A.1 The Development, Documentation and Usage of Web APIs

The development of web APIs (commonly referred to as a *web service*) and web APIs traces its roots back to the early 1990s, where the Open Software Foundation’s Distributed Computing Environment (DCE) introduced a collection of services and tools for developing and maintaining distributed systems using a client/server architecture [205]. This framework used the synchronous communication paradigm Remote Procedure Calls (RPCs) first introduced by Nelson [165] that allows procedures to be called in a remote address space as if it were local. Its communication paradigm, DCE/RPC [171], enables developers to write distributed software with underlying network code abstracted away. To bridge remote DCE/RPCs over components of different operating systems and languages, an Interface Definition Language (IDL) document served as the common service contract or *service interface* for software components.

This important leap toward language-agnostic distributed programming paved way for XML-RPC, enabling RPCs over HTTP (and thus the Web) encoded using XML (instead of octet streams [171]). As new functionality was introduced, this lead to the natural development of the Simple Object Access Protocol (SOAP), the backbone messaging connector for Web Service (WS) applications, a realisation of the Service-Oriented Architecture (SOA) [39] pattern. The SOA pattern prescribes that services are offered by service providers and consumed by service consumers in a platform- and language-agnostic manner and are used in large-scale enterprise

systems (e.g., banking, health). Key to the SOA pattern is that a service's quality attributes (see Section 2.1) can be specified and guaranteed using a Service-Level Agreement (SOA) whereby the consumer and provider agree upon a set level of service, which in some cases are legally binding [17]. This agreement can be measured using Quality of Service (QoS) parameters met by the service provider during the transportation layer (e.g., response time, cost of leasing resources, reliability guarantees, system availability and trust/security assurance [101, 249]). These attributes are included within SOAP headers; thus, QoS aspects are independent from the transport layer and instead exist at the application layer [178]. The IDL of SOAP is Web Services Description Language (WSDL), providing a description of how the web service is invoked, what parameters to expect, and what data structures are returned.

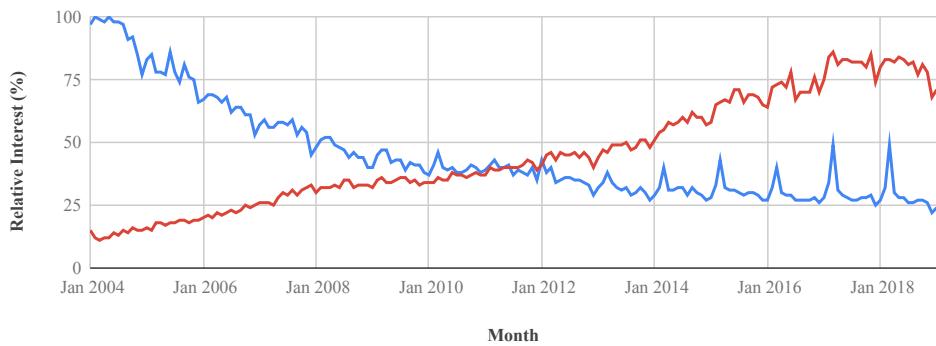


Figure A.1: Worldwide search interest for SOAP (blue) and REST (red) since 2004.

Source: [1]

While it is rich in metadata and verbosity, discussions on whether this was a benefit or drawback came about the mid-2000s [178, 262] whether the amount of data transfer paid off (especially for mobile clients where data usage was scarce).

Developer usability for debugging the SOAP ‘envelopes’ (messages POSTed over HTTP to the service provider component) was difficult, both due to the nature of XML’s wordiness and difficulty to test (by sending POST requests) in-browser. As a simple example, 25 lines (794 bytes) of HTTP communication is transferred to request a customer’s name from a record using SOAP (Listings A.1 and A.2).

Listing A.1: A SOAP HTTP POST consumer request to retrieve customer record #43456 from a web service provider. Source: [298].

```
1 POST /customers HTTP/1.1
2 Host: www.example.org
3 Content-Type: application/soap+xml; charset=utf-8
4
5 <?xml version="1.0"?>
6 <soap:Envelope
7   xmlns:soap="http://www.w3.org/2003/05/soap-envelope">
8   <soap:Body>
9     <m:GetCustomer
10       xmlns:m="http://www.example.org/customers">
11       <m:CustomerId>43456</m:CustomerId>
12     </m:GetCustomer>
13   </soap:Body>
14 </soap:Envelope>
```

Listing A.2: The SOAP HTTP service provider response for Listing A.1. Source: [298].

```

1 HTTP/1.1 200 OK
2 Content-Type: application/soap+xml; charset=utf-8
3
4 <?xml version='1.0' ?>
5 <env:Envelope
6   xmlns:env="http://www.w3.org/2003/05/soap-envelope" >
7   <env:Body>
8     <m:GetCustomerResponse
9       xmlns:m="http://www.example.org/customers">
10      <m:Customer>Foobar Quux, inc</m:Customer>
11    </m:GetCustomerResponse>
12  </env:Body>
13 </env:Envelope>
```

SOAP uses the architectural principle that web services (or the applications they provide) should remain *outside* the web, using HTTP only as a tunnelling protocol to enable remote communication [178]. That is, the HTTP is considered as a transport protocol solely. In 2000, Fielding [71] introduced REpresentational State Transfer (REST), which instead approaches the web as a medium to publish data (i.e., HTTP is part of the *application* layer instead). Hence, applications become amalgamated into the Web. Fielding bases REST on four key principles:

- **URIs identify resources.** Resources and services have a consistent global address space that aides in their discovery via URIs [20].

- **HTTP verbs manipulate those resources.** Resources are manipulated using the four consistent CRUD verbs provided by HTTP: POST, GET, PUT, DELETE.
- **Self-descriptive messages.** Each request provides enough description and context for the server to process that message.
- **Resources are stateless.** Every interaction with a resource is stateless.

Consider the equivalent example of Listings A.1 and A.2 but in a RESTful architecture (Listings A.3 and A.4) and it is clear why this style has grown more popular with developers (as we highlight in Figure A.1). Developers have since embraced RESTful API development, though the major drawback of RESTful services is its lack of a uniform IDL to facilitate development (though it is possible to achieve this using Web Application Description Language (WADL) [299]). Therefore, no RESTful service uses a standardised response document or invocation syntax. While there are proposals, such as WADL [90], RAML¹, API Blueprint², and the OpenAPI³ specification (initially based on Swagger⁴), there is still no consensus as there was for SOAP and convergence of these IDLs is still underway.

Listing A.3: An equivalent HTTP consumer request to that of Listing A.1, but using REST.

Source: [298].

¹<https://raml.org> last accessed 25 January 2019.

²<https://apiblueprint.org> last accessed 25 January 2019.

³<https://www.openapis.org> last accessed 25 January 2019.

⁴<https://swagger.io> last accessed 25 January 2019.

```
1 | GET /customers/43456 HTTP/1.1
2 | Host: www.example.org
```

Listing A.4: The REST HTTP service provider response for Listing A.3. Source: [298].

```
1 | HTTP/1.1 200 OK
2 | Content-Type: application/json; charset=utf-8
3 |
4 | {"Customer": "Foobar Quux, inc"}
```

Figure A.2: A Broad Range of AI-Based Products And Services Is Already Visible. (From [143].)

Category	Sample vendors and products	Typical use cases
Embedded AI Expert assistants leverage AI technology embedded in platforms and solutions.	<ul style="list-style-type: none"> • Amazon: Alexa • Apple: Siri • Facebook: Messenger • Google: Google Assistant (and more) • Microsoft: Cortana • Salesforce: MetaMind (acquisition) 	<ul style="list-style-type: none"> • Personal assistants for search, simple inquiry, and growing as expert assistance (composed problems, not just search) • Available on mobile platforms, devices, the internet of things • Voice, image recognition, various levels of NLP sophistication • Bots, agents
AI point solutions Point solutions provide specialized capabilities for NLP, vision, speech, and reasoning.	<ul style="list-style-type: none"> • 24[7]: 24[7] • Admantx: Admantx • Affectiva: Affdex • Assist: AssistDigital • Automated Insights: Wordsmith • Beyond Verbal: Beyond Verbal • Expert System: Cogito • HPE: Haven OnDemand • IBM: Watson Analytics, Explorer, Advisor • Narrative Science: Quill • Nuance: Dragon • Salesforce: MetaMind (acquisition) • Wise.io: Wise Support 	<ul style="list-style-type: none"> • Semantic text, facial/visual recognition, voice intonation, intelligent narratives • Various levels of NLP from brief text messaging, chat/conversational messaging, full complex text understanding • Machine learning, predictive analytics, text analytics/mining • Knowledge management and search • Expert advisors, reasoning tools • Customer service, support • APIs
AI platforms Platforms that offer various AI tech, including (deep) machine learning, as tools, APIs, or services to build solutions.	<ul style="list-style-type: none"> • CognitiveScale: Engage, Amplify • Digital Reasoning: Synthesys • Google: Google Cloud Machine Learning • IBM: Watson Developers, Watson Knowledge Studio • Intel: Saffron Natural Intelligence • IPsoft: Amelia, Apollo, IP Center • Microsoft: Cortana Intelligence Suite • Nuance: 360 platform • Salesforce: Einstein • Wipro: Holmes 	<ul style="list-style-type: none"> • APIs, cloud services, on-premises for developers to build AI solutions • Insights/advice building • Rule-based reasoning • Vertical domain advisors (e.g., fraud detection in banking, financial advisors, healthcare) • Cognitive services and bots
Deep learning Platforms, advanced projects, and algorithms for deep learning.	<ul style="list-style-type: none"> • Amazon: FireFly • Google: TensorFlow/DeepMind • LoopAI Labs: LoopAI • Numenta: Grok • Vicarious: Vicarious 	<ul style="list-style-type: none"> • Deep learning neural networks for categorization, clustering, search, image recognition, NLP, and more • Location pattern recognition • Brain neocortex simulation

Figure A.3: Increasing interest on Stack Overflow for these intelligent cloud services.

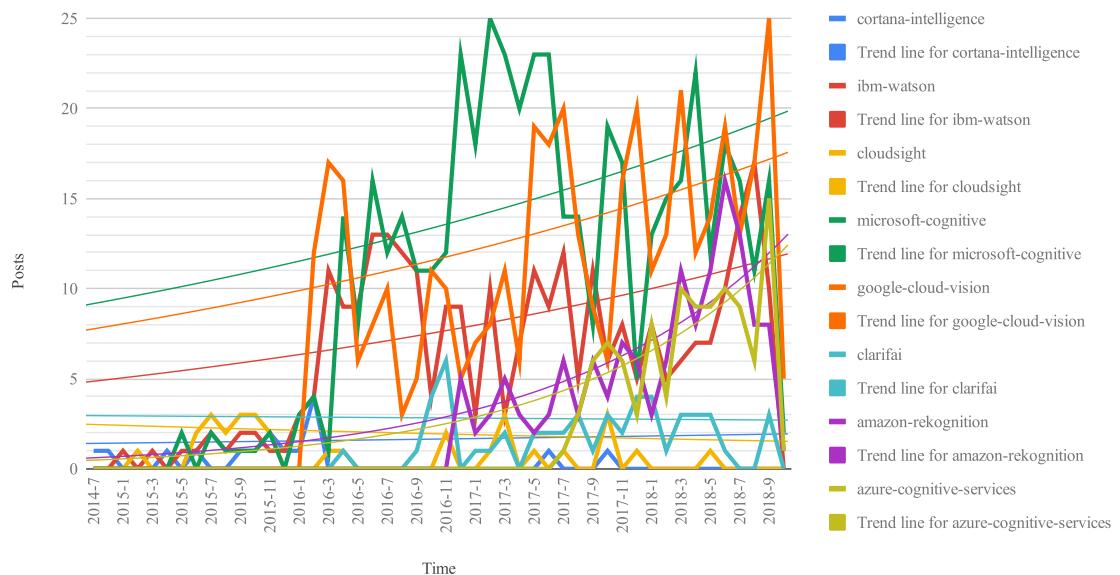


Figure A.4: Questions asked by software engineering researchers (column 2) that can be answered by field study techniques. (From [219].)

Technique	Used by researchers when their goal is to understand:	Volume of data	Also used by software engineers for
Direct techniques			
Brainstorming and focus groups	Ideas and general background about the process and product, general opinions (also useful to enhance participant rapport)	Small	Requirements gathering, project planning
Interviews and questionnaires	General information (including opinions) about process, product, personal knowledge etc.	Small to large	Requirements and evaluation
Conceptual modeling	Mental models of product or process	Small	Requirements
Work diaries	Time spent or frequency of certain tasks (rough approximation, over days or weeks)	Medium	Time sheets
Think-aloud sessions	Mental models, goals, rationale and patterns of activities	Medium to large	UI evaluation
Shadowing and observation	Time spent or frequency of tasks (intermittent over relatively short periods), patterns of activities, some goals and rationale	Small	Advanced approaches to use case or task analysis
Participant observation (joining the team)	Deep understanding, goals and rationale for actions, time spent or frequency over a long period	Medium to large	
Indirect techniques			
Instrumenting systems	Software usage over a long period, for many participants	Large	Software usage analysis
Fly on the wall	Time spent intermittently in one location, patterns of activities (particularly collaboration)	Medium	
Independent techniques			
Analysis of work databases	Long-term patterns relating to software evolution, faults etc.	Large	Metrics gathering
Analysis of tool use logs	Details of tool usage	Large	
Documentation analysis	Design and documentation practices, general understanding	Medium	Reverse engineering
Static and dynamic analysis	Design and programming practices, general understanding	Large	Program comprehension, metrics, testing, etc.

APPENDIX B

Authorship Statements

Deakin University Authorship Procedure

Schedule A: Authorship Statement

1. Details of the publication and executive author

Title of publication	Losing Confidence in Quality: Unspoken Evolution of Computer Vision Services
Publication details	Presented at the 35th IEEE International Conference on Software Maintenance and Evolution, Cleveland, USA, 2019
Name of executive author	Alex Cummaudo
School/Institute/Division if at Deakin	Applied Artificial Intelligence Institute
Organisation and address if non-Deakin	
Email or phone	ca@deakin.edu.au

2. Inclusion of publication in a thesis

Is it intended to include this publication in a higher degree by research (HDR) thesis? Yes
*If Yes, please complete Section 3
If No, go straight to Section 4.*

3. HDR thesis author's declaration

Name of HDR thesis author if different from above. (If the same, write "as above")	As above
School/Institute/Division if at Deakin	Applied Artificial Intelligence Institute
Thesis title	Provenance in Large-Scale Artificial Intelligence Solutions
If there are multiple authors, give a full description of HDR thesis author's contribution to the publication.	See page 2

I declare that the above is an accurate description of my contribution to this paper, and the contributions of other authors are as described below.

Signed:



Dated: 22 July 2019

4. Description of all author contributions

Name and affiliation of author 1

Alex Cummaudo
Applied Artificial Intelligence Institute
Deakin University

Contribution of author 1

Alex Cummaudo initiated the conception of the project. Additionally, he designed a detailed methodology, conducted all data collection via a data-collection instrument he designed and implemented and performed a majority of data analysis. He drafted the full manuscript and made further revisions, modifications and prepared the camera ready version for publication in the conference proceedings.

Name and affiliation of author 2

Rajesh Vasa
Applied Artificial Intelligence Institute
Deakin University

Contribution of author 2

Rajesh Vasa contributed to the initial conception of this project by providing high-level guidance over overview of what the project and its experiments should comprise of. Rajesh also contributed to detailed revisions of the initial manuscripts, and assisted in advising Alex Cummaudo on improved analytical insight into the collected results. Rajesh Vasa also assisted in shaping the paper to specifically target the conference audience. Rajesh Vasa is the primary supervisor of Alex Cummaudo.

Name and affiliation of author 3

John Grundy
Faculty of Information Technology
Monash University

Contribution of author 3

John Grundy provided high-level oversight of the project. He contributed to detailed reviews of the methodology and manuscript. John Grundy is the external supervisor of Alex Cummaudo.

Name and affiliation of author 4

Mohamed Abdelrazek
School of Information Technology
Deakin University

Contribution of author 4

Mohamed Abdelrazek made final edits and suggestions to the final draft of the manuscript before submitting for peer review. Mohamed Abdelrazek is an associate supervisor of Alex Cummaudo.

Name and affiliation of author 5

Andrew Cain
School of Information Technology
Deakin University

Contribution of author 5

Andrew Cain made edits and suggestions to the abstract and introduction paragraphs of the manuscript. Andrew Cain is an associate supervisor of Alex Cummaudo.

5. Author declarations

I agree to be named as one of the authors of this work, and confirm:

- i. that I have met the authorship criteria set out in the Deakin University Research Conduct Policy,
- ii. that there are no other authors according to these criteria,
- iii. that the description in Section 4 of my contribution(s) to this publication is accurate,
- iv. that the data on which these findings are based are stored as set out in Section 7 below.

If this work is to form part of an HDR thesis as described in Sections 2 and 3, I further

- v. consent to the incorporation of the publication into the candidate's HDR thesis submitted to Deakin University and, if the higher degree is awarded, the subsequent publication of the thesis by the university (subject to relevant Copyright provisions).

Author 1

Alex Cummaudo

Signed: 
Dated: 22 July 2019

Author 2

Rajesh Vasa

Signed: 
Dated: 22 July 2019

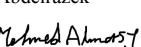
Author 3

John Grundy

Signed: 
Dated: 22 July 2019

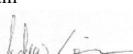
Author 4

Mohamed Abdelrazek

Signed: 
Dated: 22 July 2019

Author 5

Andrew Cain

Signed: 
Dated: 22 July 2019

6. Other contributor declarations

There are no other contributors for this publication to declare.

7. Data storage

The original data for this project are stored in the following locations. (The locations must be within an appropriate institutional setting. If the executive author is a Deakin staff member and data are stored outside Deakin University, permission for this must be given by the Head of Academic Unit within which the executive author is based.)

Data format	Comma separated values (CSV), iPython Notebook
Storage location	Deakin University Research Data Store (RDS) Location: RDS29448-Alex-Cummaudo-PhD/results/icsme19

8. Additional notices

This form must be retained by the executive author, within the school or institute in which they are based.

If the publication is to be included as part of an HDR thesis, a copy of this form must be included in the thesis with the publication.

Deakin University Authorship Procedure

Schedule A: Authorship Statement

1. Details of the publication and executive author

Title of publication	What should I document? A preliminary systematic mapping study into API documentation knowledge
Publication details	Presented at the 13th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), Porto de Galinhas, Brazil, 2019
Name of executive author	Alex Cummaudo
School/Institute/Division if at Deakin Organisation and address if non-Deakin	Applied Artificial Intelligence Institute
Email or phone	ca@deakin.edu.au

2. Inclusion of publication in a thesis

Is it intended to include this publication in a higher degree by research (HDR) thesis?
*If Yes, please complete Section 3
 If No, go straight to Section 4.*

3. HDR thesis author's declaration

Name of HDR thesis author if different from above. <i>(If the same, write "as above")</i>	As above
School/Institute/Division if at Deakin	Applied Artificial Intelligence Institute
Thesis title	Provenance in Large-Scale Artificial Intelligence Solutions
If there are multiple authors, give a full description of HDR thesis author's contribution to the publication.	See page 2

I declare that the above is an accurate description of my contribution to this paper, and the contributions of other authors are as described below.

Signed:



Dated: 22 July 2019

4. Description of all author contributions

Name and affiliation of author 1

Alex Cummaudo
Applied Artificial Intelligence Institute
Deakin University

Contribution of author 1

Alex Cummaudo devised the conception of this project and the intended objectives and hypotheses. Additionally, he designed a detailed methodology, conducted data collection with a custom tool he wrote himself and performed analysis. He drafted the manuscript and made further revisions, modifications and prepared the camera ready version for publication in the conference proceedings.

Name and affiliation of author 2

Rajesh Vasa
Applied Artificial Intelligence Institute
Deakin University

Contribution of author 2

Rajesh Vasa contributed to the initial conception of this project by providing high-level guidance over overview of what the project and its experiments should comprise of. Rajesh also contributed to detailed revisions of the initial manuscripts, and assisted in advising Alex Cummaudo on improved analytical insight into the collected results. Rajesh Vasa also assisted in shaping the paper to specifically target the conference audience. Rajesh Vasa is the primary supervisor of Alex Cummaudo.

Name and affiliation of author 3

John Grundy
Faculty of Information Technology
Monash University

Contribution of author 3

John Grundy provided high-level oversight of the project. He contributed to detailed reviews of the methodology and manuscript. John Grundy is the external supervisor of Alex Cummaudo.

5. Author declarations

I agree to be named as one of the authors of this work, and confirm:

- i. that I have met the authorship criteria set out in the Deakin University Research Conduct Policy,
- ii. that there are no other authors according to these criteria,
- iii. that the description in Section 4 of my contribution(s) to this publication is accurate,
- iv. that the data on which these findings are based are stored as set out in Section 7 below.

If this work is to form part of an HDR thesis as described in Sections 2 and 3, I further

- v. consent to the incorporation of the publication into the candidate's HDR thesis submitted to Deakin University and, if the higher degree is awarded, the subsequent publication of the thesis by the university (subject to relevant Copyright provisions).

Author 1

Alex Cummaudo

Signed: 
Dated: 22 July 2019

Author 2

Rajesh Vasa

Signed: 
Dated: 22 July 2019

Author 3

John Grundy

Signed: 
Dated: 22 July 2019

6. Other contributor declarations

There are no other contributors for this publication to declare.

7. Data storage

The original data for this project are stored in the following locations. (The locations must be within an appropriate institutional setting. If the executive author is a Deakin staff member and data are stored outside Deakin University, permission for this must be given by the Head of Academic Unit within which the executive author is based.)

Data format Comma separated values (CSV), Portable Document Format (PDF)

Storage location Deakin University Research Data Store (RDS)
Location: RDS29448-Alex-Cummaudo-PhD/results/esem19

8. Additional notices

This form must be retained by the executive author, within the school or institute in which they are based.

If the publication is to be included as part of an HDR thesis, a copy of this form must be included in the thesis with the publication.

Deakin University Authorship Procedure

Schedule A: Authorship Statement

1. Details of the publication and executive author

Title of publication	Merging Intelligent API Responses Using a Proportional Representation Approach
Publication details	Presented at the 19th International Conference on Web Engineering (ICWE), Daejeon, South Korea, 2019
Name of executive author	Tomohiro Otake
School/Institute/Division if at Deakin Organisation and address if non-Deakin	Faculty of Science, Engineering and Built Environment
Email or phone	tomohiro.otake@deakin.edu.au

2. Inclusion of publication in a thesis

Is it intended to include this publication in a higher degree by research (HDR) thesis?
*If Yes, please complete Section 3
 If No, go straight to Section 4.*

3. HDR thesis author's declaration

Name of HDR thesis author if different from above. (If the same, write "as above")	As above
School/Institute/Division if at Deakin	Applied Artificial Intelligence Institute
Thesis title	Provenance in Large-Scale Artificial Intelligence Solutions
If there are multiple authors, give a full description of HDR thesis author's contribution to the publication.	See page 2

I declare that the above is an accurate description of my contribution to this paper, and the contributions of other authors are as described below.

Signed:



Dated: 2 August 2019

4. Description of all author contributions

Name and affiliation of author 1

Tomohiro Otake
Faculty of Science, Engineering and Built Environment
Deakin University

Contribution of author 1

Tomohiro Otake designed a detailed methodology for data collection in the primary experiment of this work. He conducted all data collection via a data-collection instrument he designed and implemented and performed a majority of data analysis. He drafted the full manuscript and made further revisions, modifications and prepared the camera ready version for publication in the conference proceedings.

Name and affiliation of author 2

Alex Cummaudo
Applied Artificial Intelligence Institute
Deakin University

Contribution of author 2

Alex Cummaudo's primary contribution to this work was the conception and writing up of the motivating sections in the manuscript. He additionally contributed to detailed editing of the manuscripting to make further revisions and modifications and implemented reviewer feedback.

Name and affiliation of author 3

Mohamed Abdelrazek
Faculty of Science, Engineering and Built Environment
Deakin University

Contribution of author 3

Mohamed Abdelrazek contributed to the initial conception of this project by providing high-level guidance over overview of what the project and its experiments should comprise of. Mohamed also contributed to detailed revisions of the initial manuscripts, and assisted in advising Tomohiro Otake on improved analytical insight into the collected results, and implementing reviewer feedback.

Name and affiliation of author 4

Rajesh Vasa
Faculty of Science, Engineering and Built Environment
Deakin University

Contribution of author 4

Rajesh Vasa provided high-level oversight of the project. He contributed to detailed reviews of the methodology and manuscript.

Name and affiliation of author 5

John Grundy
Faculty of Information Technology
Monash University

Contribution of author 5

John Grundy provided high-level oversight of the project. He contributed to detailed reviews of the methodology and manuscript.

5. Author declarations

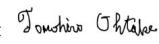
I agree to be named as one of the authors of this work, and confirm:

- i. that I have met the authorship criteria set out in the Deakin University Research Conduct Policy,
- ii. that there are no other authors according to these criteria,
- iii. that the description in Section 4 of my contribution(s) to this publication is accurate,
- iv. that the data on which these findings are based are stored as set out in Section 7 below.

If this work is to form part of an HDR thesis as described in Sections 2 and 3, I further

- v. consent to the incorporation of the publication into the candidate's HDR thesis submitted to Deakin University and, if the higher degree is awarded, the subsequent publication of the thesis by the university (subject to relevant Copyright provisions).

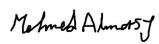
Author 1

Tomohiro Ohtake
 Signed: 
 Dated: 2 August 2019

Author 2

Alex Cummaudo
 Signed: 
 Dated: 2 August 2019

Author 3

Mohamed Abdelrazek
 Signed: 
 Dated: 2 August 2019

Author 4

Rajesh Vasa
 Signed: 
 Dated: 2 August 2019

Author 5

John Grundy
 Signed: 
 Dated: 2 August 2019

6. Other contributor declarations

There are no other contributors for this publication to declare.

7. Data storage

The original data for this project are stored in the following locations. (The locations must be within an appropriate institutional setting. If the executive author is a Deakin staff member and data are stored outside Deakin University, permission for this must be given by the Head of Academic Unit within which the executive author is based.)

Data format	Comma separated values (CSV)
Storage location	Deakin University Research Data Store (RDS) Location: RDS29448-Alex-Cummaudo-PhD/results/icwe19

8. Additional notices

This form must be retained by the executive author, within the school or institute in which they are based.

If the publication is to be included as part of an HDR thesis, a copy of this form must be included in the thesis with the publication.

APPENDIX C

Ethics Clearance



Rajesh Vasa and Alex Cummaudo
Applied Artificial Intelligence Institute (A²I²)
C.c Mohamed Abdelrazek, Andrew Cain

2 May 2019

Dear Rajesh and Alex

STEC-11-2019-CUMMAUDO titled "*Developer opinions towards the importance of web API documentation recommendations*"

Thank you for submitting the above project for consideration by the Faculty Human Ethics Advisory Group (HEAG). The HEAG recognised that the project complies with the National Statement on Ethical Conduct in Human Research (2007) and has approved it. You may commence the project upon receipt of this communication.

The approval period is for three years until **02/05/22**. It is your responsibility to contact the Faculty HEAG immediately should any of the following occur:

- Serious or unexpected adverse effects on the participants
- Any proposed changes in the protocol, including extensions of time
- Any changes to the research team or changes to contact details
- Any events which might affect the continuing ethical acceptability of the project
- The project is discontinued before the expected date of completion.

You will be required to submit an annual report giving details of the progress of your research. Please forward your first annual report on **02/05/20**. Failure to do so may result in the termination of the project. Once the project is completed, you will be required to submit a final report informing the HEAG of its completion.

Please ensure that the Deakin logo is on the Plain Language Statement and Consent Forms. You should also ensure that the project ID is inserted in the complaints clause on the Plain Language Statement, and be reminded that the project number must always be quoted in any communication with the HEAG to avoid delays. All communication should be directed to sciethic@deakin.edu.au

The Faculty HEAG and/or Deakin University Human Research Ethics Committee (HREC) may need to audit this project as part of the requirements for monitoring set out in the National Statement on Ethical Conduct in Human Research (2007).

If you have any queries in the future, please do not hesitate to contact me.

We wish you well with your research.

Kind regards

A handwritten signature in blue ink that reads "Teresa Treffry".

Teresa Treffry
Secretary, Human Ethics Advisory Group (HEAG)
Faculty of Science Engineering & Built Environment



Rajesh Vasa, Mohamed Abdelrazek, Andrew Cain, Scott Barnett, Alex Cummaudo
 Applied Artificial Intelligence Institute (A²I²) (G)

23rd July 2019

Dear Rajesh and research team

STEC-39-2019-CUMMAUDO titled "*Factors that impact the learnability, interpretability and adoption of intelligent services*".

Thank you for submitting the above project for consideration by the Faculty Human Ethics Advisory Group (HEAG). The HEAG recognised that the project complies with the National Statement on Ethical Conduct in Human Research (2007) and has approved it. You may commence the project upon receipt of this communication.

The approval period is for three years until 23/07/22. It is your responsibility to contact the Faculty HEAG immediately should any of the following occur:

- Serious or unexpected adverse effects on the participants
- Any proposed changes in the protocol, including extensions of time
- Any changes to the research team or changes to contact details
- Any events which might affect the continuing ethical acceptability of the project
- The project is discontinued before the expected date of completion.

You will be required to submit an annual report giving details of the progress of your research. Please forward your first annual report on 23/07/20. Failure to do so may result in the termination of the project. Once the project is completed, you will be required to submit a final report informing the HEAG of its completion.

Please ensure that the project number must always be quoted in any communication with the HEAG to avoid delays. All communication should be directed to sciethic@deakin.edu.au.

The Faculty HEAG and/or Deakin University Human Research Ethics Committee (HREC) may need to audit this project as part of the requirements for monitoring set out in the National Statement on Ethical Conduct in Human Research (2007).

If you have any queries in the future, please do not hesitate to contact me.

We wish you well with your research.

Kind regards

Rickie Morey

Rickie Morey
 Senior Research Administration Officer
 Representing the Human Ethics Advisory Group (HEAG)
 Faculty of Science Engineering & Built Environment

APPENDIX D

Primary Sources from Systematic Mapping Study

References

- [1] M. P. Robillard, "What makes APIs hard to learn? Answers from developers," *IEEE Software*, vol. 26, no. 6, pp. 27–34, 2009.
- [2] M. P. Robillard and R. Deline, "A field study of API learning obstacles," *Empirical Software Engineering*, vol. 16, no. 6, pp. 703–732, 2011.
- [3] A. J. Ko and Y. Riche, "The role of conceptual knowledge in API usability," in *2011 IEEE Symposium on Visual Languages and Human-Centric Computing (VL/HCC)*. IEEE, 2011, pp. 173–176.
- [4] J. Nykaza, R. Messinger, F. Boehme, C. L. Norman, M. Mace, and M. Gordon, "What programmers really want - results of a needs assessment for SDK documentation." *SIGDOC*, 2002.
- [5] R. Watson, M. Mark Stamnes, J. Jeannot-Schroeder, and J. H. Spyridakis, "API documentation and software community values," in *the 31st ACM international conference*. New York, New York, USA: ACM Press, 2013, pp. 165–10.
- [6] S. Y. Jeong, Y. Xie, J. Beaton, B. A. Myers, J. Stylos, R. Ehret, J. Karstens, A. Efeoglu, and D. K. Busse, "Improving documentation for eSOA APIs through user studies," in *International Symposium on End User Development*. Springer, 2009, pp. 86–105.
- [7] E. Aghajani, C. Nagy, O. L. Vega-Márquez, M. Linares-Vásquez, L. Moreno, G. Bavota, and M. Lanza, "Software Documentation Issues Unveiled," *ICSE*, vol. 2, no. 3, pp. 127–139, 2019.
- [8] S. Haselbock, R. Weinreich, G. Buchgeher, and T. Kriechbaum, "Microservice Design Space Analysis and Decision Documentation: A Case Study on API Management," *2018 IEEE 11th Conference on Service-Oriented Computing and Applications (SOCA)*, pp. 1–8, Nov. 2018.
- [9] S. Inzunza, R. Juárez-Ramírez, and S. Jiménez, "API Documentation," in *Trends and Advances in Information Systems and Technologies*. Cham: Springer, Cham, Mar. 2018, pp. 229–239.
- [10] M. Meng, S. Steinhardt, and A. Schubert, "Application Programming Interface Documentation: What Do Software Developers Want?" *Journal of Technical Writing and Communication*, vol. 48, no. 3, pp. 295–330, Aug. 2017.
- [11] R. S. Geiger, N. Varoquaux, C. Mazel-Cabasse, and C. Holdgraf, "The Types, Roles, and Practices of Documentation in Data Analytics Open Source Software Libraries," *Computer Supported Cooperative Work (CSCW)*, vol. 27, no. 3-6, pp. 767–802, May 2018.
- [12] A. Head, C. Sadowski, E. Murphy-Hill, and A. Knight, *When not to comment: questions and tradeoffs with API documentation for C++ projects*, ser. questions and tradeoffs with API documentation for C++ projects. New York, New York, USA: ACM, May 2018.

-
-
- [13] L. Aversano, D. Guardabascio, and M. Tortorella, "Analysis of the Documentation of ERP Software Projects," *Procedia Computer Science*, vol. 121, pp. 423–430, Jan. 2017.
 - [14] M. P. Robillard, A. Marcus, C. Treude, G. Bavota, O. Chaparro, N. Ernst, M. A. Gerosa, M. Godfrey, M. Lanza, M. Linares-Vásquez, G. C. Murphy, L. Moreno, D. Shepherd, and E. Wong, "On-demand Developer Documentation," in *2017 IEEE International Conference on Software Maintenance and Evolution (IC-SME)*. IEEE, 2017, pp. 479–483.
 - [15] R. B. Watson, "Development and application of a heuristic to assess trends in API documentation." *SIGDOC*, 2012.
 - [16] W. Maalej and M. P. Robillard, "Patterns of Knowledge in API Reference Documentation." *IEEE Trans. Software Eng.*, 2013.
 - [17] D. L. Parnas and S. A. Vilkomir, "Precise Documentation of Critical Software," in *10th IEEE High Assurance Systems Engineering Symposium (HASE'07)*. IEEE, Sep. 2007, pp. 237–244.
 - [18] C. Bottomley, "What part writer? What part programmer? A survey of practices and knowledge used in programmer writing," *IPCC 2005. Proceedings. International Professional Communication Conference, 2005.*, pp. 802–812, Jan. 2005.
 - [19] A. Taulavuori, E. Niemelä, and P. Kallio, "Component documentation—a key issue in software product lines," *Information and Software Technology*, vol. 46, no. 8, pp. 535–546, Jun. 2004.
 - [20] J. Kotula, "Using Patterns To Create Component Documentation." *IEEE Software*, 1998.
 - [21] S. G. McLellan, A. W. Roesler, J. T. Tempest, and C. Spinuzzi, "Building More Usable APIs." *IEEE Software*, 1998.