

Taming the Evolving Black Box:
Towards Improved Integration and Documentation of
Intelligent Web 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

March 11, 2020

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: Towards Improved Integration and Documentation of Intelligent Web 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)
March 11, 2020

To my family, friends, and teachers.

Acknowledgements

A long journey of 20 years education has led me to this thesis, and there are so many people who've helped me get to this point that deserve thanking. To start, I must thank my family; you have always been there for me in times good and bad. I'm especially grateful to my mother, Rosa, my father, Paul, my siblings, Marc and Lisa, and my nonna, Michelina, for your love and support. I also thank my nieces Nina and Lucy (though too young to read this now!) for bringing us all so much joy in these last three years. I thank all my teachers, particularly Natalie Heath, whose hard efforts paid off in my tertiary education, and, of course, all those who assisted me along and help shape this PhD. Firstly, to Professor Rajesh Vasa, thank you for your many revisions, patience, ideas and efforts to shape this work: your years of phenomenal guidance—both as a supervisor and as a mentor—has reshaped my worldview to far wider perspectives and I now approach thought and problem-solving in a remarkably new light. Secondly, I thank Professor John Grundy for your efforts and for being such an approachable and hard-working supervisor, always willing to provide feedback and guidance, and help me get over the finish line. I also thank Dr Scott Barnett for the many fruitful discussions shared, for the interest you have shown in my work, and the joint collaborations we achieved in the last two years; Associate Professor Mohamed Abdelrazek for your help in shaping many of the works within this thesis; and, lastly, Associate Professor Andrew Cain, who not only taught me the realm of programming back in undergraduate days, but who first suggested a PhD was within my reach, of which I had never fathomed. I must thank everyone at the Applied Artificial Intelligence Institute for creating such an enjoyable environment to work in, especially Jake Renzella, Rodney Pilgrim, Mahdi Babaei, and Reuben Wilson for their friendship over these years and for all the coffee runs, conversations, and ideas shared. And, lastly, to Tom Fellowes: thank you for being by my side throughout this journey.

This chapter is now over, the next chapter awaits...

— Alex Cummaudo
July 2020

Contents

Abstract	iii
Declaration	v
Acknowledgements	ix
Contents	ix
List of Publications	xv
List of Abbreviations	xvii
List of Figures	xxi
List of Tables	xxiv
List of Listings	xxvi
I Preface	1
1 Introduction	3
1.1 Research Context	5
1.2 Motivating Scenarios	7
1.3 Research Motivation	12
1.4 Research Goals	14
1.5 Research Methodology	16
1.6 Thesis Organisation	16
1.6.1 Part I: Preface	17
1.6.2 Part II: Publications	17
1.6.3 Part III: Postface	19

1.6.4	Part IV: Appendices	19
1.7	Research Contributions	20
1.7.1	Contribution 1: Landscape Analysis & Preliminary Solutions	21
1.7.2	Contribution 2: Improving Documentation Attributes	22
1.7.3	Contribution 3: Service Integration Architecture	23
2	Background	25
2.1	Software Quality	26
2.1.1	Validation and Verification	27
2.1.2	Quality Attributes and Models	29
2.1.3	Reliability in Computer Vision	31
2.2	Probabilistic and Nondeterministic Systems	33
2.2.1	Interpreting the Uninterpretable	34
2.2.2	Explanation and Communication	35
2.2.3	Mechanics of Model Interpretation	35
2.3	Application Programming Interfaces	36
2.3.1	API Usability	37
3	Research Methodology	39
3.1	Research Questions Revisited	39
3.1.1	Empirical Research Questions	41
3.1.2	Non-Empirical Research Questions	41
3.2	Philosophical Stances	42
3.3	Research Methods	43
3.3.1	Review of Relevant Research Methods	43
3.3.2	Review of Data Collection Techniques for Field Studies	45
3.4	Research Design	45
3.4.1	Landscape Analysis of Computer Vision Services	47
3.4.2	Utility of API Documentation in Computer Vision Services	47
3.4.3	Developer Issues concerning Computer Vision Services	47
3.4.4	Designing Improved Integration Strategies	48
II	Publications	49
4	Identifying Evolution in Computer Vision Services	51
4.1	Introduction	51
4.2	Motivating Example	53
4.3	Related Work	54
4.3.1	External Quality	54
5	Interpreting Pain-Points in Computer Vision Services	55
5.1	Introduction	55
5.2	Motivation	57
5.3	Background	59
5.4	Method	60

5.4.1	Data Extraction	60
5.4.2	Data Filtering	62
5.4.3	Data Analysis	63
5.5	Findings	65
5.5.1	Post classification and reliability analysis	65
5.5.2	Developer Frustrations	66
5.5.3	Statistical Distribution Analysis	68
5.6	Discussion	68
5.6.1	Answers to Research Questions	68
5.6.2	The Developer’s Learning Approach	70
5.6.3	Implications	72
5.7	Threats to Validity	75
5.7.1	Internal Validity	75
5.7.2	External Validity	75
5.7.3	Construct Validity	76
5.8	Conclusions	76
6	Ranking Computer Vision Service Issues using Emotion	77
6.1	Introduction	77
6.2	Emotion Mining from Text	79
6.3	Methodology	79
6.3.1	Data Set Extraction from Stack Overflow (SO)	80
6.3.2	Question Type & Emotion Classification	82
6.4	Findings	84
6.5	Discussion	85
6.6	Implications	87
6.7	Threats to Validity	87
6.7.1	Internal Validity	87
6.7.2	External Validity	88
6.7.3	Construct Validity	88
6.8	Conclusions	88
7	Better Documenting Computer Vision Services	91
7.1	Introduction	91
7.2	Related Work	94
7.2.1	API Usability and Documentation Knowledge	94
7.2.2	Adapting the System Usability Scale	95
7.2.3	Computer Vision Services	95
7.3	Taxonomy Development	95
7.3.1	Systematic Mapping Study	96
7.3.2	Development of the Taxonomy	100
7.4	API Documentation Knowledge Taxonomy	102
7.5	Validating our Taxonomy	105
7.5.1	Survey Study	105

7.5.2	Empirical application of the taxonomy against Computer Vision Services	106
7.6	Taxonomy Analysis	107
7.6.1	In-Literature Scores for Taxonomy Categories	107
7.6.2	In-Practice Scores for Taxonomy Categories	108
7.6.3	Contrasting In-Literature to In-Practice Scores	109
7.6.4	Triangulating ILS and IPS with Computer Vision	110
7.6.5	Areas of Improvement for CVS Documentation	111
7.7	Threats to Validity	115
7.7.1	Internal Validity	115
7.7.2	External Validity	115
7.7.3	Construct Validity	116
7.8	Conclusions & Future Work	117
8	Using a Facade Pattern to combine Computer Vision Services	119
8.1	Introduction	119
8.1.1	Motivating Scenario: Intelligent vs Traditional Web Services	120
8.1.2	Research Motivation	121
8.2	Merging API Responses	121
8.2.1	API Facade Pattern	122
8.2.2	Merge Operations	122
8.2.3	Merging Operators for Labels	123
8.3	Graph of Labels	123
8.3.1	Labels and synsets	124
8.3.2	Connected Components	124
8.4	API Results Merging Algorithm	127
8.4.1	Mapping Labels to Synsets	127
8.4.2	Deciding Total Number of Labels	127
8.4.3	Allocating Number of Labels to Connected Components .	127
8.4.4	Selecting Labels from Connected Components	129
8.4.5	Conformance to properties	129
8.5	Evaluation	129
8.5.1	Evaluation Method	129
8.5.2	Naive Operators	130
8.5.3	Traditional Proportional Representation Operators	132
8.5.4	New Proposed Label Merge Technique	132
8.5.5	Performance	132
8.6	Conclusions and Future Work	133
9	Supporting Safe Usage of Intelligent Web Services	135
9.1	Introduction	135
9.2	Motivating Example	137
9.3	Threshy	139
9.4	Related work	140
9.4.1	Decision Boundary Estimation	140

9.4.2	Tooling for ML Frameworks	141
9.5	Conclusions & Future Work	142
10	FSE Paper	143
III	Postface	145
11	Conclusions & Future Work	147
	References	168
	List of Online Artefacts	169
IV	Appendices	173
A	Additional Materials	175
A.1	On the Development, Documentation and Usage of Web APIs	177
A.2	Additional Figures	180
A.3	Reference Architecture Source Code	183
B	Supplementary Materials to Chapter 7	185
B.1	Detailed Overview of Our Proposed Taxonomy	187
B.2	Sources of Documentation	191
B.3	List of Primary Sources	194
B.4	Survey Questions	196
C	Authorship Statements	201
D	Ethics Clearance	215

List of Publications

Below lists publications arising from work completed in this PhD.

1. A. Cummaudo, R. Vasa, J. Grundy, M. Abdelrazek, and A. Cain, “Losing Confidence in Quality: Unspoken Evolution of Computer Vision Services,” in *Proceedings of the 35th IEEE International Conference on Software Maintenance and Evolution*. Cleveland, OH, USA: IEEE, December 2019. DOI 10.1109/ICSME.2019.00051. ISBN 978-1-72-813094-1 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 *Proceedings of the 13th International Symposium on Empirical Software Engineering and Measurement*. Porto de Galinhas, Recife, Brazil: IEEE, October 2019. DOI 10.1109/ESEM.2019.8870148. ISBN 978-1-72-812968-6. ISSN 1949-3789 pp. 1–6
3. A. Cummaudo, R. Vasa, S. Barnett, J. Grundy, and M. Abdelrazek, “Interpreting Cloud Computer Vision Pain-Points: A Mining Study of Stack Overflow,” in *Proceedings of the 42nd International Conference on Software Engineering*. Seoul, Republic of Korea: IEEE, October 2020, In Press
4. A. Cummaudo, S. Barnett, R. Vasa, and J. Grundy, “Threshy: Supporting Safe Usage of Intelligent Web Services,” 2020, Unpublished
5. A. Cummaudo, R. Vasa, and J. Grundy, “Assessing API documentation knowledge for computer vision services,” 2020, Unpublished
6. A. Cummaudo, S. Barnett, R. Vasa, J. Grundy, and M. Abdelrazek, “Beware the evolving ‘intelligent’ web service! An integration architecture tactic to guard AI-first components,” 2020, Unpublished
7. T. Ohtake, A. Cummaudo, M. Abdelrazek, R. Vasa, and J. Grundy, “Merging intelligent API responses using a proportional representation approach,” in *Proceedings of the 19th International Conference on Web Engineering*. Daejeon, Republic of Korea: Springer, June 2019. DOI 10.1007/978-3-03-019274-7_28. ISBN 978-3-03-019273-0. ISSN 1611-3349 pp. 391–406
8. M. K. Curumsing, A. Cummaudo, U. M. Graestch, S. Barnett, and R. Vasa, “Ranking Computer Vision Service Issues using Emotion,” 2020, Unpublished

List of Abbreviations

A²I² Applied Artificial Intelligence Institute. 45, 47

AI artificial intelligence. iii, 3–5, 8, 12–14, 34, 35, 51, 52, 55–57, 59, 60, 70, 73, 74, 76, 77, 82, 121, 122, 181

API application programming interface. xiii, xxiv, 4–16, 18, 22, 23, 25, 26, 28, 29, 36–38, 40–42, 44, 47, 52, 53, 55–62, 65–78, 81, 82, 84–87, 91–96, 98–105, 108–117, 119–122, 124, 126, 127, 129, 130, 133, 135, 177, 179

BYOML Build Your Own Machine Learning. 5, 6

CC connected component. 124, 127–129, 132

CDSS clinical decision support system. 7, 10

CNN convolutional neural network. 10, 11, 33

CRUD create, read, update, and delete. 179

CV computer vision. 5, 7, 23, 32, 37, 51, 53, 55, 56, 60, 61, 67, 70, 71, 74, 76, 92, 95, 96, 105, 113, 114

CVS computer vision service. xxiv, 7–10, 12, 14–23, 25, 27, 28, 31, 37, 40, 41, 43–45, 47, 48, 51–55, 57, 62, 67, 71, 75–77, 79–81, 84, 88, 91–93, 95, 102, 104, 106, 107, 110–117, 119, 121–123, 127, 129, 133, 182

DCE distributed computing environment. 177

HITL human-in-the-loop. 11

HTTP Hypertext Transfer Protocol. 6, 177–179

IDL interface definition language. 177, 179

IRR inter-rater reliability. 75

IWS intelligent web service. 5–7, 9–12, 14, 15, 17, 18, 25–28, 31, 33, 36–38, 51–53, 55–61, 63, 65, 66, 68–78, 81, 82, 88, 91, 92, 117, 119–121, 133, 135

JSON JavaScript Object Notation. 7

ML machine learning. iii, 3–6, 8, 9, 12, 13, 18, 21, 29, 34, 37, 51, 52, 56–58, 60, 61, 73, 74, 87, 119, 120, 133

NN neural network. 12, 32, 34, 36

PaaS Platform as a Service. 7, 11

QoS quality of service. 177

RAML RESTful API Modeling Language. 179

REST REpresentational State Transfer. 7, 52, 55, 76, 120, 178, 179

ROI region of interest. 10, 11

RPC remote procedure call. 177

SDK software development kit. 53, 96, 111

SE software engineering. 14, 15, 17, 18, 35, 52, 57, 60, 63, 73, 76, 95, 96, 98–101, 105, 106, 115

SLA service-level agreement. 177

SMS systematic mapping study. 18, 19, 22, 93–96, 101, 107, 116, 117

SO Stack Overflow. xiii, 5, 15, 17, 18, 22, 23, 40, 41, 44, 47, 48, 55, 57–61, 63, 65, 66, 68–71, 73–80, 82–85, 87, 88, 182

SOA service-oriented architecture. 177

SOAP Simple Object Access Protocol. 7, 177–179

SOLO Structure of the Observed Learning Outcome. 70–73, 75

SQA service quality assurance. 53, 54

SQuaRE Systems and software Quality Requirements and Evaluation. 30

SUS System Usability Scale. 18, 92, 93, 95, 105, 106, 108, 116

SVM support vector machine. 34, 36

URI uniform resource identifier. 179

V&V verification & validation. 25–29

WADL Web Application Description Language. 179

WS web service. 6, 28, 120, 121, 177, 179

WSDL Web Services Description Language. 177

XML eXtendable markup language. 7, 177

List of Figures

1.1	Differences between data- and rule-driven cloud services	4
1.2	The spectrum of machine learning	5
1.3	Overview of intelligent web services	7
1.4	CancerAssist Context Diagram	11
1.5	Overview publication coherency	21
2.1	Mindset clashes within the development, use and nature of a IWS . .	26
2.2	Leakage of internal and external quality in	29
2.3	Overview of software quality models	30
2.4	Adversarial examples in computer vision	32
2.5	Deterministic versus nondeterministic systems	33
2.6	Theory of AI communication	36
3.1	Review of field study techniques	46
5.1	Traits of intelligent web services compared to DIY ML	58
5.2	Trend of Stack Overflow posts discussing computer vision services .	59
5.3	Comparing documentation-specific and generalised classifications of Stack Overflow posts	66
5.4	Alignment of Bloom and SOLO taxonomies against computer vision issues	72
6.1	Distribution of Stack Overflow question types	84
6.2	Proportion of emotions per question type	85
7.1	Systematic mapping study search results, by years	97
7.2	Filtering steps used in the systematic mapping study	97
7.3	A systematic map of API documentation knowledge studies	101
7.4	Our proposed API documentation knowledge taxonomy	103
7.5	Comparison of ILS and IPS values for each category (grouped by dimensions) presented as a relative percentage.	109

7.6	Comparison of the weighted ILS and IPS values for the three computer vision services (CVSs) assessed.	110
8.1	Overview of the proposed facade	122
8.2	Graph of associated synsets against two different endpoints	125
8.3	Label counts per API assessed	126
8.4	Connected components vs. images	126
8.5	Allocation to connected components	128
8.6	F-measure comparison	130
9.1	Example case study of evaluating model performance in two different models	136
9.2	Example pipeline of a computer vision system	138
9.3	UI workflow of Threshy	139
9.4	Architecture of Threshy	141
A.1	SOAP versus REST search interest over time	178
A.2	Categorisation of AI-based products and services	181
A.3	Increasing interest in the developer community of computer vision services	182

List of Tables

1.1	Differing characteristics of cloud services	6
1.2	Comparison of the machine learning spectrum	6
1.3	Varying confidence changes over time between three computer vision services	9
1.4	Definitions of ‘confidence’ in CV documentation	10
1.5	List of publications resulting from this thesis	20
3.1	Classification of research questions in this thesis	40
5.1	Taxonomies used in our Stack Overflow mining study	64
5.2	Example Stack Overflow posts aligning to Bloom’s and SOLO taxonomies	71
6.1	Our interpretations from a Stack Overflow question type taxonomy .	81
6.2	Frequency of emotions per question type.	84
6.3	Assigning emotion to Stack Overflow questions	86
7.1	Summary of search results in API documentation knowledge	98
7.2	Data extraction in API documentation knowledge study	100
7.3	Weighted ILS Scoring.	108
7.4	Weighted IPS Scoring.	108
7.5	Labels assigned to ILS and IPS values	109
8.1	Statistics for the number of labels	124
8.2	First allocation iteration	129
8.3	Second allocation iteration	129
8.4	Third allocation iteration	129
8.5	Fourth allocation iteration	129
8.6	Matching to human-verified labels	131
8.7	Evaluation results of the facade	131
8.8	Average of evaluation result of the facade	131

List of Listings

A.1	An example SOAP request	177
A.2	An example SOAP response	178
A.3	An example RESTful request	179
A.4	An example RESTful response	179

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 encapsulating half a
10 century worth of machine-learning research.¹ Application developers are eager to
11 develop the next generation of ‘AI-first’ software, that will reason, sense, think, act,
12 listen, speak and execute every whim in our web browser or smartphone app.

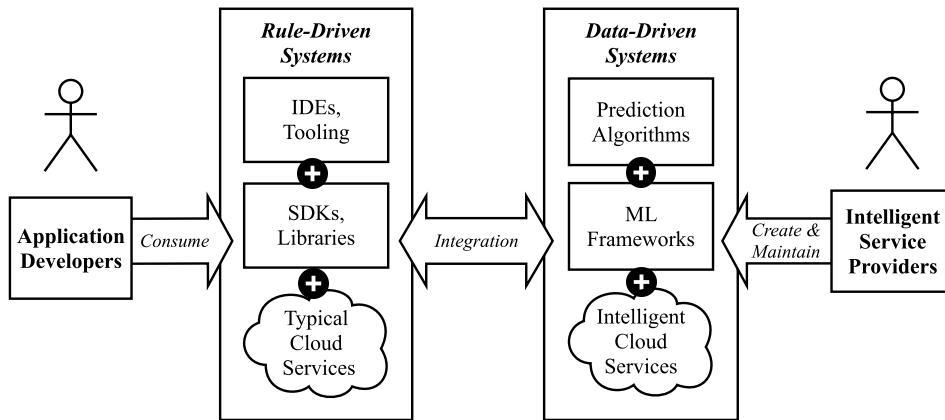
13 However, application developers, accustomed to traditional software engineering
14 paradigms, may not be aware of AI-first’s consequences. Application developers
15 build *rule-driven* applications, where every line of source code evaluates to produce
16 deterministic outcomes. AI-first software is, however, not rule-driven but *data-
driven*. Large datasets train machine learning (ML) prediction classifiers that result
18 in probabilistic confidences of results and nondeterministic behaviour if it continually
19 learns more data with time. Furthermore, developing AI-first applications requires
20 both code *and data*, and an application developer can approach developing from
21 three (non-traditional) perspectives, further expanded in Section 1.1:

- 22 1. The application developer writes an ML classifier from scratch and trains it
23 from a handcrafted and curated dataset. This approach is laborious in time and
24 demands formal training in ML and mathematical knowledge, but the tradeoff
25 is that they have full autonomy in the models they creates.
- 26 2. The application developer downloads a pre-trained model and ‘plugs’ it into
27 an existing ML framework, such as Tensorflow [1]. While this approach is
28 less demanding in time, it requires them to revise and understand how to ‘glue’
29 components of the ML framework together² into their application’s code.

¹A 2016 report by market research company Forrester captured such growth into four key areas [193], as reproduced in Figure A.2.

²Thus introducing a verbose list of ML terminology to her developer vocabulary. See a list of 328 terms provided by Google here: <https://developers.google.com/machine-learning/glossary/>. Last accessed 7 December 2018.

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.



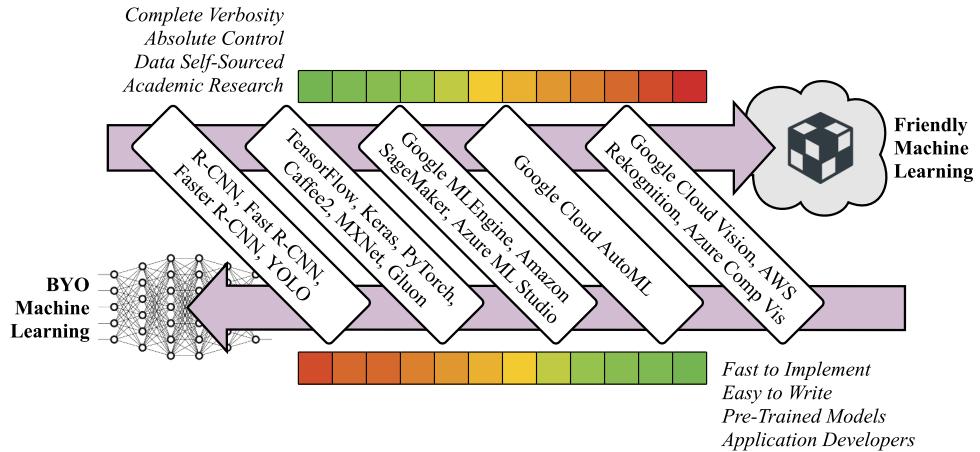
30 3. The application developer uses a data-driven and cloud-based service. They
 31 don’t need to know anything behind the underlying ‘intelligence’ and how it
 32 functions. It is fast to integrate into their applications, and the application pro-
 33 gramming interfaces (APIs) offered abstracts the technical know-how behind
 34 a web call.

35 The documentation of the service alludes that the data-driven service is as similar to
 36 other cloud services offered by the provider. Because this is ‘another’ cloud service,
 37 the application developer *assumes* it would act and behave as any other typical
 38 service would. But does this assumption—and a lack of appreciation of ML—lead
 39 to developer pain-points and miscomprehension? If so, how can the service providers
 40 improve their documentation to alleviate this? Do these data-driven services share
 41 similarities to the runtime behaviour of traditional cloud services? And if not,
 42 how best can the application developer integrate the data-driven service into their a
 43 rule-driven application to produce AI-first software?

44 Figure 1.1 provides an illustrative overview between the context clashing of rule-
 45 driven applications and data-driven cloud services, and we contrast characteristics
 46 of typical cloud systems and data-driven ones in Table 1.1.

In this thesis, we advocate that the integration and developer comprehen-
 sion of data-driven cloud services differ from the rule-driven nature of
 end-applications. As ‘intelligent’ components these contrast to traditional
 counterparts, and application developers need to take into account a greater
 appreciation of these factors.

Figure 1.2: Examples within the machine learning spectrum of computer vision. Colour scales indicates the benefits (green) and drawbacks (red) of each end of the spectrum.



47 1.1 Research Context

48 As described, the application developer has three key approaches in producing AI-
 49 first software. This ‘range’ of AI-first integration techniques partially reflects Google
 50 AI’s³ *machine learning spectrum* [180, 205, 232], which encompasses the variety
 51 of skill, effort, users and types of outputs of integration techniques. One extreme
 52 involves the academic research of developing algorithms and self-sourcing data
 53 to achieve intelligence—coined as Build Your Own Machine Learning (BYOML)
 54 [156, 205, 232]. The other extreme involves off-the-shelf, ‘friendlier’ (abstracted)
 55 intelligence with easy-to-use APIs targeted towards applications developers. The
 56 middle-ground involves a mix of the two, with varying levels of automation to assist
 57 in development, that turns custom datasets into predictive intelligence. We illustrate
 58 the slightly varied characteristics within this spectrum in Table 1.2 and Figure 1.2.

59 These data-driven ‘friendly’ services are gaining traction within developer circles:
 60 we show an increasing trend of Stack Overflow posts mentioning a mix of
 61 intelligent computer vision (CV) services in Figure A.3.⁴ Academia provides var-
 62 ied nomenclature for these services, such as *Cognitive Applications* and *Machine*
 63 *Learning Services* [317] or *Machine Learning as a Service* [256]. For the context
 64 of this thesis, we will refer to such services under broader term of **intelligent web**
 65 **services (IWSs)**, and diagrammatically express their usage within Figure 1.3.

66 While there are many types of IWSs available to software developers,⁵ the

³Google 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 on an AI-first attitude we see through Google’s 2018 rebranding of *Google Research* to *Google AI* [142] is evident. A further example includes how Facebook leverage AI *at scale* within their infrastructure and platforms [236].

⁴Query run on 12 October 2018 using StackExchange Data Explorer. Refer to <https://data.stackexchange.com/stackoverflow/query/910188> for full query.

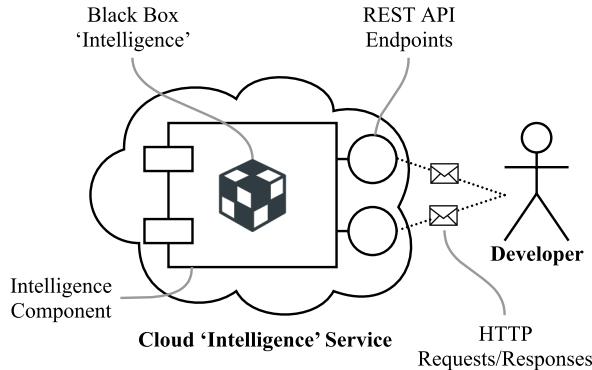
⁵Such as optical character recognition, text-to-speech and speech-to-text transcription, object

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

Intelligent web service	Typical web service
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	✓	✓			
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				✓	✓

Figure 1.3: Overview of IWSs.

⁶⁷ general workflow of using an IWS is more-or-less the same: a developer accesses
⁶⁸ an IWS component via REST/SOAP API(s), which is (typically) available as a
⁶⁹ cloud-based Platform as a Service (PaaS).⁶⁻⁷ For a given input, developers receive
⁷⁰ an ‘intelligent’ response and an associated confidence value that represents the
⁷¹ likelihood of that result. This is typically serialised as a JSON/XML response
⁷² object.

☞ Within this thesis, we scope our investigation to a mature subset of IWSs that provide computer vision intelligence [338, 341, 354, 355, 356, 358, 362, 371, 372, 374, 376, 422, 423]. For the context of this thesis, we will refer to such services as **CVSs**.

1.2 Motivating Scenarios

⁷³ The market for computer vision services (CVSs) is increasing (Figure A.2) and as
⁷⁴ is developer uptake and enthusiasm in the software engineering community (Figure
⁷⁵ A.3). However, the impact to software quality (internal and external) due to
⁷⁶ a mismatch of the application developer’s deterministic mindset and the service
⁷⁷ provider’s nondeterministic mindset is of concern.
⁷⁸

⁷⁹ To illustrate the context of use, we present the two scenarios of varying risk: (i) a
⁸⁰ fictional software developer, named Tom, who wishes to develop an inherently low-
⁸¹ risk photo detection application for his friends and family; and (ii) a high-risk cancer
⁸² clinical decision support system (CDSS) that uses patient scans to recommend if

categorisation, facial analysis and recognition, natural language processing etc.

⁸³We note, however, that a development team may use a similar approach *internally* within a product line or service that may not necessarily reflect a PaaS model.

⁸⁴A number of services provide the platform infrastructure to rapidly begin training from custom datasets, such as Google’s AutoML (<https://cloud.google.com/automl/>, last accessed 7 December 2018). Others provide pre-trained datasets ‘ready-for-use’ in production without the need to train data.

⁸³ surgeons should send their patients to surgery. Both describe scenarios where AI-
⁸⁴ first components has substantiative impact to end-users when the software engineers
⁸⁵ developing with them misunderstand the nuances of ML, ultimately adversely affecting
⁸⁶ external quality. Moreover, due to lack of comprehension, this hinders developer
⁸⁷ experience, productivity, and understanding/appreciation of AI-based components.

⁸⁸ 1.2.0.1 *Motivating Scenario I: Tom's PhotoSharer App*

⁸⁹ Tom wants to develop a social media photo-sharing app on iOS and Android, *Photo-
⁹⁰ Sharer*, that analyses photos taken on smartphones. Tom wants the app to categorise
⁹¹ photos into scenes (e.g., day vs. night, landscape vs. indoors), generate brief de-
⁹² scriptions of each photo, and catalogue photos of his friends and common objects
⁹³ (e.g., photos with his Border Collie dog, photos taken on a beach on a sunny day with
⁹⁴ his partner). His app will shares this analysed photo intelligence with his friends on
⁹⁵ a social-media platform, where his friends can search and view the photos.

⁹⁶ Instead of building a computer vision engine from scratch, which takes too much
⁹⁷ time and effort, Tom thinks he can achieve this using one of the common CVSs. Tom
⁹⁸ comes from a typical software engineering background and has insufficient knowl-
⁹⁹ edge of key computer vision terminology and no understanding of its underlying
¹⁰⁰ techniques. However, inspired by easily accessible cloud APIs that offer computer
¹⁰¹ vision analysis, he chooses to use these. Built upon his experience of using other
¹⁰² similar cloud services, he decides on one of the CVS APIs, and expects a static result
¹⁰³ always and consistency between similar APIs. Analogously, when Tom invokes the
¹⁰⁴ iOS Swift substring method "doggy".prefix(3), he expects it to be consistent
¹⁰⁵ with the Android Java equivalent "doggy".substring(0, 2). Consistent, here,
¹⁰⁶ means two things: (i) that calling substring or prefix on 'dog' will *always*
¹⁰⁷ return in the same way every time he invokes the method; and (ii) that the result is
¹⁰⁸ *always* 'dog' regardless of the programming language or string library used, given
¹⁰⁹ the deterministic nature of the 'substring' construct (i.e., results for substring are
¹¹⁰ API-agnostic).

¹¹¹ More concretely, in Table 1.3, we illustrate how three (anonymised) CVS
¹¹² providers fail to provide similar consistency to that of the substring example above.
¹¹³ If Tom uploads a photo of a border collie⁸ to three different providers in August
¹¹⁴ 2018 and January 2019, he would find that each provider is different in both the vo-
¹¹⁵ cabulary used between. The confidence values and labels within the *same* provider
¹¹⁶ varies within a matter of five months. The evolution of the confidence changes is
¹¹⁷ not explicitly documented by the providers (i.e., when the models change) nor do
¹¹⁸ they document what confidence means. Service providers use a tautological nature
¹¹⁹ when defining what the confidence confidence values are (as presented in the API
¹²⁰ documentation) provides no insight for Tom to understand why there was a change
¹²¹ in confidence, which we show in Table 1.4, unless he *knows* that the underlying
¹²² models change with them. Furthermore, they do not provide detailed understanding
¹²³ on how to select a threshold cut-off for a confidence value. Therefore, he's left with
¹²⁴ no understanding on how best to tune for image classification in this instance. The

⁸The image used for these results is <https://www.akc.org/dog-breeds/border-collie/>.

Table 1.3: First six responses of image analysis for a Border Collie sent to three CVS 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

¹²⁵ deterministic problem of a substring compared to the nondeterministic nature of the
¹²⁶ IWS is, therefore, non-trivial.

¹²⁷ To make an assessment of these APIs, he tries his best to read through the
¹²⁸ documentation of different CVS APIs, but he has no guiding framework to help him
¹²⁹ choose the right one. A number of questions come to mind:

- ¹³⁰ • What does ‘confidence’ mean?
- ¹³¹ • Which confidence is acceptable in this scenario?
- ¹³² • Are these APIs consistent in how they respond?
- ¹³³ • Are the responses in APIs static and deterministic?
- ¹³⁴ • Would a combination of multiple CVS APIs improve the response?
- ¹³⁵ • How does he know when there is a defect in the response? How can he report
¹³⁶ it?
- ¹³⁷ • How does he know what labels the API knows, and what labels it doesn’t?
- ¹³⁸ • How does it describe his photos and detect the faces?
- ¹³⁹ • Does he understand that the API uses a machine learnt model? Does he know
¹⁴⁰ what a ML model is?
- ¹⁴¹ • Does he know when models update? What is the release cycle?

¹⁴² Although Tom generally anticipates these CVSs to not be perfect, he has no
¹⁴³ prior benchmark to guide him on what to expect. The imperfections appear to be
¹⁴⁴ low-risk, but may become socially awkward when in use; for instance, if Tom’s
¹⁴⁵ friends have low self-esteem and use the app, they may be sensitive to the app not
¹⁴⁶ identifying them or mislabelling them. Privacy issues come into play especially
¹⁴⁷ if certain friends have access to certain photos that they are (supposedly) in; e.g.,

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

API Provider	Definition(s) of Confidence
Provider A	“Score is the confidence score, which ranges from 0 (no confidence) to 1 (very high confidence).” [360]
	“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].” [361]
	“The overall score of the result. Range [0, 1]” [361]
Provider B	“Confidence score, between 0 and 1... if there insufficient confidence in the ability to produce a caption, the tags maybe [sic] the only information available to the caller.” [377]
	“The level of confidence the service has in the caption.” [375]
Provider C	“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.” [339]
	“[Provider C] also provide[s] a percentage score for how much confidence [Provider C] has in the accuracy of each detected label.” [340]

¹⁴⁸ photos from a holiday with Tom and his partner, however if the API identifies Tom’s
¹⁴⁹ partner as a work colleague, Tom’s partner’s privacy is at risk.

¹⁵⁰ Therefore, the level of risk and the determination of what constitutes an ‘error’ is
¹⁵¹ dependent on the situation. In the following example, an error caused by the service
¹⁵² may be more dangerous.

¹⁵³ 1.2.0.2 Motivating Scenario II: Cancer Detection CDSS

¹⁵⁴ Recent studies in the oncology domain have used deep-learning convolutional neural
¹⁵⁵ networks (CNNs) to detect region of interests (ROIs) in image scans of tissue (e.g.,
¹⁵⁶ [27, 127, 192]), flagging these regions for doctors to review. Trials of such algorithms
¹⁵⁷ have been able to accurately detect cancer at higher rates than humans, and thus
¹⁵⁸ incorporating such capabilities into a CDSS is closer within reach. Studies have
¹⁵⁹ suggested these systems may erode a practitioner’s independent decision-making
¹⁶⁰ [63, 153] due to over-reliance; therefore the risks in developing CDSSs powered by
¹⁶¹ IWSs become paramount.

¹⁶² In Figure 1.4 we present a context diagram for a fictional CDSS named *CancerAssist*. A team of busy pathologists utilise CancerAssist to review patient lymph
¹⁶³ node scans and discuss and recommend, on consensus, if the patient requires an
¹⁶⁴ operation. When the team makes a consensus, the lead pathologist enters the ver-

166 dict into CancerAssist—running passively in the background—to ensure there is
167 no oversight in the team’s discussions. When a conflict exists between the team’s
168 verdict and CancerAssist’s verdict, the system produces the scan with ROIs it thinks
169 the team should review. Where the team overrides the output of CancerAssist, this
170 reinforces CancerAssist’s internal model as a human-in-the-loop (HITL) learning
171 process.

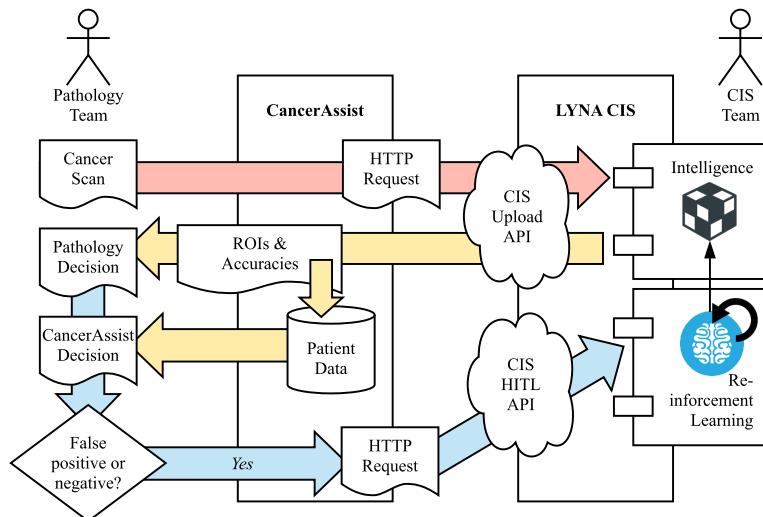


Figure 1.4: CancerAssist Context Diagram. *Key:* Red Arrows = Scan Input; Yellow Arrows = Decision Output; Blue Arrows = HITL Feedback Input.

172 Powering CancerAssist is Google AI’s Lymph Node Assistant (LYNA) [192],
173 a CNN based on the Inception-v3 model [177, 300]. To provide intelligence to
174 CancerAssist, the development team decide to host LYNA as an IWS using a cloud-
175 based PaaS solution. Thus, CancerAssist provides API endpoints integrated with
176 patient data and medical history, which produces the verdict. In the case of a positive
177 verdict, CancerAssist highlights the relevant ROIs found are with their respective
178 bounding boxes and their respective cancer detection accuracies.

179 The developer of CancerAssist has no interaction with the Data Science team
180 maintaining the LYNA IWS. As a result, they are unaware when updates to the
181 model occur, nor do they know what training data they provide to test their system.
182 The default assumptions are that the training data used to power the intelligence
183 is near-perfect for universal situations; i.e., the algorithm chosen is the correct one
184 for every assessable ontology tests in the given use case of CancerAssist. Thus,
185 unlike deterministic systems—where the developer can manually test and validate
186 the outcomes of the APIs—this is impossible for non-deterministic systems such
187 as CancerAssist and its underlying IWS. The ramifications of not being able to test
188 such a system and putting it out into production may prove fatal to patients.

189 Certain questions in the production of CancerAssist and its use of an IWS may
190 come into mind:

- 191 • When is the model updated and how do the IWS team communicate these

¹⁹² updates?

- ¹⁹³ • What benchmark test set of data ensures that the changed model doesn't affect other results?
- ¹⁹⁴
- ¹⁹⁵ • Are assumptions made by the IWS team who train the model correct?

¹⁹⁶ Thus, to improve communication between developers and IWS providers, developers require enhanced documentation, additional metadata, and guidance tooling.

¹⁹⁸ 1.3 Research Motivation

¹⁹⁹ Evermore applications are using IWSs as demonstrated by ubiquitous examples: ²⁰⁰ aiding the vision-impaired [83, 254], accounting [199], data analytics [151], and ²⁰¹ student education [89]. As our motivating examples have illustrated, these AI-based ²⁰² components—specifically CVSs—are accessible through APIs consisting of ‘black ²⁰³ box’ intelligence (Figure 1.3).⁹ Data science teams produce ML algorithms to make ²⁰⁴ predictions in our datasets and discover patterns within them. As these algorithms ²⁰⁵ are data-dependent, they are therefore inherently probabilistic and stochastic, which ²⁰⁶ results in four critical issues that motivate our thesis: (i) certainty in results, (ii) ²⁰⁷ evolution of datasets, (iii) selecting appropriate decision boundaries, and (iv) the ²⁰⁸ clarity of ML documentation that address items i–iii.

²⁰⁹ There is little room for certainty in these results as the insight is purely statistical ²¹⁰ and associational [242] against its training dataset. Developers who build these ²¹¹ applications **do not treat their programs with a stochastic or probabilistic** ²¹² **mindset, given that they are trained with a rule-driven mindset that computers** ²¹³ **make certain outcomes.** However, CVSs are data-driven, and therefore return the ²¹⁴ *probability* that a particular object exists in an input images’ pixels via confidence ²¹⁵ values. As an example, consider simple arithmetic representations (e.g., $2 + 2 =$ ²¹⁶ 4). The deterministic (rule-driven) mindset suggests that the result will *always* be ²¹⁷ 4. However, the non-deterministic (data-driven) mindset suggests that results are ²¹⁸ probable: target output (*exactly* 4) and the output inferred (*a likelihood of* 4) matches ²¹⁹ as a probable percentage (or as an error where it does not match).¹⁰ Instead of an ²²⁰ exact output, there is a *probabilistic* result: $2 + 2$ *may* equal 4 to a confidence of n . ²²¹ Thus, for a more certain (though not fully certain) distribution of overall confidence ²²² returned from the service, a developer must treat the problem stochastically by ²²³ testing this case hundreds if not thousands of times to find a richer interpretation of ²²⁴ the inference made and ensure reliability in its outcome.

²²⁵ Traditional software engineering principles advocate for software systems to be ²²⁶ versioned upon substantial change. Unfortunately **we find that the most prominent** ²²⁷ **cloud vendors providing these intelligent services (e.g., Microsoft Azure, Google** ²²⁸ **Cloud and Amazon Web Services) do not release new versioned endpoints of the**

⁹The ‘black box’ refers to a system that transforms input (or stimulus) to outputs (or response) without any understanding of the internal architecture by which this transformation occurs. This arises from a theory in the electronic sciences and adapted to wider applications since the 1950s–60s [12, 53] to describe “systems whose internal mechanisms are not fully open to inspection” [12].

¹⁰Blake et al. [36] produces a multi-layer perceptron neural network performing arithmetic representation.

229 APIs when the *internal model* changes [76]. In the context of computer vision, new
230 labels may be introduced or dropped, confidence values may differ, entire ontologies
231 or specific training parameters may change, but we hypothesise that is not effectively
232 communicated to developers. Broadly speaking, this can be attributed to a dichotomy
233 of release cycles from the data science and software engineering communities: the
234 data science iterations and work by which new models are trained and released runs
235 at a faster cycle than the maintenance cycle of traditional software engineering. Thus
236 we see cloud vendors integrating model changes without the *need* to update the API
237 version unless substantial code or schema changes are also introduced—the nuance
238 changes in the internal model does not warrant a shift in the API itself, and therefore
239 the version shift in a new model does not always propagate to a version shift in the
240 API endpoint. As demonstrated in Table 1.3, whatever input is uploaded at one time
241 may not necessarily be the same when uploaded at a later time. This again contrasts
242 the rule-driven mindset, where $2 + 2$ *always* equals 4. Therefore, in addition to the
243 certainty of a result in a single instance, the certainty of a result in *multiple instances*
244 may differ with time, which again impacts on the developers notion of reliable
245 software. Currently, it is impossible to invoke requests specific to a particular model
246 that was trained at a particular date in time, and therefore developers need to consider
247 how evolutionary changes of the services may impact their solutions *in production*.
248 Again, whether there is any noticeable behavioural changes from these changes is
249 dependent on the context of the problem domain—unless developers benchmark
250 these changes against their own domain-specific dataset and frequently check their
251 selected service against such a dataset, there is no way of knowing if substantive
252 errors have been introduced.

253 As the only response from these computer vision classifiers are a label and
254 confidence value; **the decision boundaries needs to always be appropriately con-**
255 sidered by client code for each use case and each model selected. The external
256 quality of such software needs to consider reliability in the case of thresholding con-
257 fidence values—that is whether the inference has an appropriate level of confidence
258 to justify a predicted (and reliable) result to end-users. Selecting this confidence
259 threshold is non-trivial; a ML course from Google suggests that “it is tempting
260 to assume that [a] classification threshold should always be 0.5, but thresholds
261 are problem-dependent, and are therefore values that you must tune.” [123]. Ap-
262 proaches to turning these values are considered for data scientists, but are not yet
263 well-understood for application developers with little appreciation of the nuances of
264 ML.

265 Similarly, developers should consider the internal quality of building AI-first
266 software. Reliable API usability and documentation advocate for the accuracy,
267 consistency and completeness of APIs and their documentation [247, 261] and
268 providers should consider mismatches between a developer’s conceptual knowledge
269 of the API its implementation [172]. **Unreliable APIs ultimately hinder developer**
270 **performance and thus reduces productivity**, in addition to producing potentially
271 unreliable software where documentation is not well-understood (or clear to the
272 developer).

273 Ultimately, these four issues present major threats to software reliability if left

²⁷⁴ unresolved. Given that such substantiative software engineering principles on re-
²⁷⁵ liability, versioning and quality are under-investigated within the context of IWSs,
²⁷⁶ we aim to explore guidance from the software engineering literature to investigate
²⁷⁷ what aspects in the development lifecycle could aide in mitigating these issues when
²⁷⁸ developing using AI-based components. This serves as our core motivation for this
²⁷⁹ work.

²⁸⁰ 1.4 Research Goals

²⁸¹ This thesis aims to investigate and better understand the nature of cloud-based
²⁸² computer vision services (CVSs)¹¹ as a concrete exemplar of intelligent web services
²⁸³ (IWSs). We identify the maturity, viability and risks of CVSs through the anchoring
²⁸⁴ perspective of *reliability* that affects the internal and external quality of software.
²⁸⁵ We adopt the McCall [203] and Boehm [38] interpretations of reliability via the sub-
²⁸⁶ characteristics of a service's *consistency* and *robustness* (or fault/error tolerance), and
²⁸⁷ the *completeness*¹² of its documentation. (A detailed discussion is further provided
²⁸⁸ in Section 2.1.) This thesis explores and contributes towards *four* key facets regarding
²⁸⁹ reliability in CVS usage and the completeness of its associated documentation. We
²⁹⁰ formulate four primary research questions (RQs) with seven sub-RQs, based on
²⁹¹ both empirical and non-empirical software engineering methodology [213], further
²⁹² discussed in Chapter 3.

²⁹³ Firstly, we investigate adverse implications that arise when using CVSs that
²⁹⁴ affects consistency and robustness (**Chapter 4**). We show how CVSs have a non-
²⁹⁵ deterministic runtime behaviour and evolve with unintended and non-trivial con-
²⁹⁶ sequences to developers. We demonstrate that these services have inconsistent
²⁹⁷ behaviour despite offering the same functionality and pose evolution risk that ef-
²⁹⁸ fects robustness of consuming applications when responses change given the same
²⁹⁹ (consistent) inputs. Thus, we conclude how the nature of these services (at present)
³⁰⁰ are not fully robust, consistent, and thus not reliable. Formally, we structure the
³⁰¹ following RQs:

❶ RQ1. What is the nature of cloud-based CVSs?

RQ1.1. What is their runtime behaviour?

RQ1.2. What is their evolution profile?

³⁰² Secondly, we investigate the reliability of the documentation these services of-
³⁰³ fer through the lenses of its completeness. We collate prior knowledge of good
³⁰⁴ API documentation and assess the efficacy of such knowledge against practition-
³⁰⁵ ers (**Chapter 7**). We show that these service's behaviour and evolution is not
³⁰⁶ reliably documented adequately against this knowledge. Formally, we develop the
³⁰⁷ following RQs:

¹¹As these services are proprietary, we are unable to conduct source code or model analysis, and hence are not used in the investigation of this thesis.

¹²We treat the API documentation of a CVS as a first-class citizen.

② RQ2. Are CVS APIs sufficiently documented?

- RQ2.1.* What are the dimensions of a ‘*complete*’ API document, according to both literature and practitioners?
- RQ2.2.* What additional information or attributes do application developers need in CVS API documentation to make it more complete?

308 Thirdly, we investigate how software developers approach using these services
309 and directly assess developer pain-points resulting from the nature of CVSs and
310 their documentation (**Chapter 5**). We show that there is a statistically significant
311 difference in these complaints when contrasted against more established software
312 engineering domains (such as web or mobile development) as expressed as ques-
313 tions asked on Stack Overflow. We provide a number of exploratory avenues for
314 researchers, educators, software engineers and IWS providers to alleviate these com-
315 plaints based on this analysis. Further, using a data set consisting of 1,245 Stack
316 Overflow questions, we explore the emotional state of developers to understand
317 which aspects (i.e., pain-points) developers are most frustrated with (**Chapter 6**).
318 We formulate the following RQs:

**③ RQ3. Are CVSs more misunderstood than conventional software en-
gineering domains?**

- RQ3.1.* What types of issues do application developers face most when using CVSs, as expressed as questions on Stack Overflow?
- RQ3.2.* Which of these issues are application developers most frustrated with?
- RQ3.3.* Is the distribution CVS pain-points different to established software engineering domains, such as mobile or web development?

319 Lastly, we explore several strategies to help improve CVSs reliability. Firstly,
320 we investigate whether merging the responses of *multiple* CVSs can improve their
321 reliability and propose a novel algorithm—based on the proportional representation
322 method used in electoral systems—to merge labels and associated confidence values
323 from three providers (**Chapter 8**). Secondly, we develop an integration architecture
324 style (or facade) to guard against CVS evolution, and synthesise an integration
325 workflow that addresses the concerns raised by developers in addition to embedding
326 ‘complete’ documentation artefacts into the workflow’s design (**Chapters 9 and 10**).
327 Our final RQ is:

**④ RQ4. What strategies can developers employ to integrate their appli-
cations with CVSs while preserving robustness and reliability?**

328 1.5 Research Methodology

329 This thesis employs a mixed-methods approach using the concurrent triangulation
330 strategy [47, 202]. The research presented consists of both empirical and non-
331 empirical research design. This section provides a high-level overview of the re-
332 search methodology within this thesis. Further details are provided in Section 1.7
333 and Chapter 3.

334 Firstly, RQ1–RQ3 are all empirical, knowledge-based questions [97, 209] that
335 aim to provide the software engineering community with a greater understanding
336 of the phenomena surrounding CVSs from three perspectives: the nature of the ser-
337 vices themselves, how developers perceive these services and how service providers
338 can improve these services. We answer RQ1 using a longitudinal experiment that
339 assesses both the services’ responses and associated documentation (complement-
340 ing RQ2.2). We adopt qualitative and quantitative data collection; specifically (i)
341 structured observations to quantitatively analyse the results over time, and (ii) docu-
342 mentary research methods to inspect service documentation. Secondly, we perform
343 systematic mapping study following the guidelines of Kitchenham and Charters
344 [168] and Petersen et al. [244] to better understand how API documentation of these
345 services can be improved (i.e., more complete), which targets Item RQ2. Based on
346 the findings from this study, we use a systematic taxonomy development methodol-
347 ogy specifically targeted toward software engineering [311] that structures scattered
348 API documentation knowledge into a taxonomy. We then validate this taxonomy
349 against practitioners using survey research, adopting Brooke well-established Sys-
350 tematic Usability Score [50] surveying instrument and contextualising it within API
351 documentation utility, which answers RQ3.3. To answer RQ2.2, we perform an
352 empirical application of the taxonomy to three CVSs, and therefore assess where
353 improvements can be made. Thirdly, we adopt field survey research using repository
354 mining of developer discussion forums (i.e., Stack Overflow) to answer RQ3, and
355 classify these using both manual and automated techniques.

356 The second aspect of our research design involves non-empirical research, which
357 explores a design-based question [213] to answer RQ4. As the answers to our
358 first three RQs establish a greater understanding of the nature behind CVSs from
359 various perspectives, the strategies we design in RQ4 aims at designing more reliable
360 integration methods so that developers can better use these cloud-based services in
361 their applications.

362 1.6 Thesis Organisation

363 We organise the thesis into four parts. **Part I (The Preface)** includes introduc-
364 tory, background and methodology chapters. This is a *PhD by Publication*, and
365 **Part II (Publications)** comprises of seven publications resulting from this work
366 over Chapters 4 to 10; publications are included verbatim except for terminology
367 and formatting changes to better fit the suitability of a coherent thesis. **Part III (The**
368 **Postface)** includes the conclusion and future works chapter, as well as a list of aca-
369 demic studies and online artefacts referenced within the thesis. **Part IV (Appendices)**

³⁷⁰ includes all supplementary material, including mandatory authorship statements and
³⁷¹ ethics approval. Details of each chapter following this introductory chapter are pro-
³⁷² vided in the following section.

³⁷³ **1.6.1 Part I: Preface**

³⁷⁴ *1.6.1.1 Chapter 2: Background*

³⁷⁵ This chapter provides an overview of prior studies broadly around three key pillars:
³⁷⁶ the development of an IWS, the usage of an IWS, and the nature of an IWS. We use
³⁷⁷ the three perspectives of software quality (particularly, reliability), probabilistic and
³⁷⁸ non-deterministic systems, and explanation and communication theory to describe
³⁷⁹ prior work.

³⁸⁰ *1.6.1.2 Chapter 3: Research Methodology*

³⁸¹ This chapter provides a summative review of research methods and philosophical
³⁸² stances relevant to software engineering. We illustrate that the methods used within
³⁸³ our publications are sound via an analysis of the methodologies used in seminal
³⁸⁴ works referenced in this thesis.

³⁸⁵ **1.6.2 Part II: Publications**

³⁸⁶ *1.6.2.1 Chapter 4: Exploring the nature of CVSS*

³⁸⁷ This chapter was presented at the 2019 International Conference on Software
³⁸⁸ Maintenance and Evolution (ICSME) [76]. We describe an 11-month longitudinal
³⁸⁹ experiment assessing the behavioural (run-time) issues of three popular CVSSs:
³⁹⁰ Google Cloud Vision [362], Amazon Rekognition [341] and Azure Computer Vi-
³⁹¹ sion [376]. By using three different data sets—two of which we curate as additional
³⁹² contributions—we demonstrate how the services are inconsistent amongst each other
³⁹³ and within themselves. This study provides a detailed answer to RQ1: Despite
³⁹⁴ presenting conceptually-similar functionality, each service behaves and produces
³⁹⁵ slightly varied (inconsistent) results and demonstrates non-deterministic runtime
³⁹⁶ behaviour. We discuss potential evolution risks to consumers of such services as the
³⁹⁷ services provide non-static outputs for the same inputs, thereby having significant
³⁹⁸ impact to the robustness of consuming applications. Further details in the study
³⁹⁹ include a brief assessment into the lack of sufficient detail of these concerns in their
⁴⁰⁰ documentation.

⁴⁰¹ *1.6.2.2 Chapter 5: Understanding developer struggles when using CVSS*

⁴⁰² This chapter has been accepted for presentation at the 2020 International Conference
⁴⁰³ on Software Engineering (ICSE) [79]. We conduct a mining study of 1,425 Stack
⁴⁰⁴ Overflow questions that provide indications of the types frustrations that developers
⁴⁰⁵ face when integrating CVSSs into their applications. To gather what their pain-points
⁴⁰⁶ are, we use two classification taxonomies that also use Stack Overflow to understand

⁴⁰⁷ generalised and documentation-specific pain-points in mature software engineering
⁴⁰⁸ (SE) domains. This study answers RQ3 in detail and provides a validation to
⁴⁰⁹ our motivation of RQ2: we validate that the *completeness* of current CVS API
⁴¹⁰ documentation is a main concern for developers and there is insufficient explanation
⁴¹¹ into the errors and limitations of the service. We find that the documentation does
⁴¹² not adequately cover all aspects of the technical domain. In terms of integrating with
⁴¹³ the service, developers struggle most with simple errors and ways in which to use the
⁴¹⁴ APIs; this is in stark contrast to mature software domains. Our interpretation is that
⁴¹⁵ developers fail to understand the IWS lifecycle and the ‘whole’ system that wraps
⁴¹⁶ such services. We also interpret that developers have a shallower understanding
⁴¹⁷ of the core issues within CVSs (likely due to the nuances of ML as suggested in
⁴¹⁸ a discussion in the paper), which warrants an avenue for future work in software
⁴¹⁹ engineering education.

⁴²⁰ **1.6.2.3 Chapter 6: Ranking CVS pain-points by frustration**

⁴²¹ This chapter has been submitted to the the 2020 International Workshop on Emotion
⁴²² Awareness in Software Engineering (SEmotion) [81]. In this work, we use our
⁴²³ dataset consisting of the 1,425 SO questions from [79] to interpret the breakdown of
⁴²⁴ emotions developers express per classification of pain-points conducted in Chapter 5.
⁴²⁵ We find that the distribution of various emotions differ per question type, and
⁴²⁶ developers are most frustrated when the expectations of a CVS does not match the
⁴²⁷ reality of what these services actually provide, which shapes our answer for RQ3.2
⁴²⁸ and thus RQ3.

⁴²⁹ **1.6.2.4 Chapter 7: Investigating improvements to CVS API documentation**

⁴³⁰ This chapter was originally a short paper presented at the 2019 International Sym-
⁴³¹ posium on Empirical Software Engineering and Measurement (ESEM) [79]. To
⁴³² understand where to improve CVS documentation, we first need to investigate *what*
⁴³³ makes a good API document. This short paper initially answered one aspect of
⁴³⁴ RQ2.1: what *academic literature* suggests a good (complete) API document should
⁴³⁵ comprise of. By conducting an systematic mapping study resulting in 21 primary
⁴³⁶ studies, we systematically develop a taxonomy that combines the recommendations
⁴³⁷ of scattered work into a structured framework of 5 dimensions and 34 weighted cat-
⁴³⁸ egorisations. We then extend this work by triangulating the taxonomy with opinions
⁴³⁹ from developers using the System Usability Scale to assess the efficacy of these
⁴⁴⁰ recommendations (thereby answering the second aspect of RQ2.1). From this, we
⁴⁴¹ assess the how well CVS providers document their APIs via a heuristic validation
⁴⁴² of the taxonomy, using the three services from the ICSME publication to make rec-
⁴⁴³ ommendations where documentation should be more complete, thereby answering
⁴⁴⁴ RQ2.2 (and thus RQ2). The extended version of this chapter has been submitted to
⁴⁴⁵ the IEEE Transactions on Software Engineering (TSE) in [80] and is currently in
⁴⁴⁶ review.

447 1.6.2.5 Chapter 8: Merging responses of multiple CVSs

448 This chapter was presented at the 2019 International Conference on Web Engineering
449 (ICWE) [229]. Early exploration of CVSs showed that multiple services use
450 vastly different ontologies for the same input. As an initial strategy to improve
451 the reliability of these services, we explored if merging multiple responses using
452 WordNet [215] and a novel label merging algorithm based on the proportional rep-
453 resentation approach used in political voting could make any improvements. While
454 this approach resulted in a modest improvement to reliability, it did not consider to
455 the evolution issues or developer pain-points we later identified.

456 1.6.2.6 Chapter 9: Developing a confidence thresholding tool

457 This chapter has been submitted to the demonstrations track at FSE 2020 [77]. When
458 integrating with a CVS, developers need to select an appropriate confidence threshold
459 suited to their use case and determine whether a decision should be made. An issue,
460 however, is that these CVSs are not calibrated to the specific problem-domain datasets
461 and it is difficult for software developers to determine an appropriate confidence
462 threshold on their problem domain. This tool presents a workflow and supporting
463 tool for application developers to select decision thresholds suited to their domain
464 that—unlike existing tooling—is designed to be used in pre-development, pre-release
465 and production. This tooling forms part of a solution to RQ4 for developers to
466 maintain robustness and reliability in their systems.

467 1.6.2.7 Chapter 10: Developing a CVS integration architecture

468 This chapter has been submitted to the 2020 Joint European Software Engineer-
469 ing Conference and Symposium on the Foundations of Software Engineering [78].
470 *(TODO: Discuss findings.)*

471 1.6.3 Part III: Postface

472 In Chapter 11, we review the contributions made in this thesis and the relevance
473 and significance to identifying and resolving key issues when application developers
474 integrate with CVS. We evaluate these outcomes with reference to the research goals,
475 and discuss threats to validity of the work. Lastly, we discuss the various avenues
476 of research arising from this work. References from literature and a list of online
477 artefacts are provided after this concluding chapter.

478 1.6.4 Part IV: Appendices

479 Appendix A provides additional material referenced within this thesis but not pro-
480 vided in the body. The supplementary materials published with Chapter 7 are
481 reproduced in Appendix B, which also describes the list of primary sources arising
482 in the systematic mapping study we conducted. We provide mandatory coauthor
483 declaration forms describing the contribution breakdown for each publication within

Table 1.5: List of publications resulting from this thesis, separated by phenomena exploration (above) and solution design (below).

Ref.	Venue	Acronym	Rank ¹³	Published ¹⁴	Chapter	RQs
[76]	35 th International Conference on Software Maintenance and Evolution	ICSME	A	05 Dec 2019	Chapter 4	RQ1
[75]	13 th International Symposium on Empirical Software Engineering and Measurement	ESEM	A	17 Oct 2019	Excluded ¹⁵	RQ2.1
[79]	42 nd International Conference on Software Engineering	ICSE	A*	<i>In Press</i>	Chapter 5	RQ3
[81]	5 th International Workshop on Emotion Awareness in Software Engineering ¹⁶	SEmotion	A*	<i>In Review</i>	Chapter 6	RQ3.2
[80]	IEEE Transactions on Software Engineering	TSE	Q1	<i>In Review</i>	Chapter 7	RQ2
[229]	13 th International Conference on Web Engineering	ICWE	B	26 Apr 2019	Chapter 8	RQ4
[77]	28 th Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering	FSE(d) ¹⁷	A*	<i>In Review</i>	Chapter 9	RQ4
[78]	28 th Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering	FSE	A*	<i>In Review</i>	Chapter 10	RQ4

⁴⁸⁴ Appendix C. Appendix D contains copies of the ethics clearance for various experiments within this thesis.

⁴⁸⁶ 1.7 Research Contributions

⁴⁸⁷ The outcomes of answering the four primary research questions elaborated in Section 1.4 shapes three primary contributions this thesis offers to software engineering knowledge:

- ⁴⁹⁰ • An **improved understanding in the landscape of CVSs**, with respect to their runtime behaviour and evolutionary profiles.
- ⁴⁹¹ • A novel **service integration architecture** that helps developers with integrating their applications with CVSs.
- ⁴⁹² • A key list of attributes that should be documented, to assist CVS providers to better document their services.

⁴⁹³ In this section, we detail how each publication forms a coherent body of work and how each publication relates to the primary contributions made.

¹⁴Conference publications ranking measured using the CORE Conference Ranks (<http://www.core.edu.au/conference-portal>) and Journal publications rankings using the Scimago Ranking (<https://www.scimagojr.com/>). Rankings retrieved January 2020.

¹⁵Date of publication, if applicable.

¹⁶The extended version of this conference proceeding is provided in Chapter 7.

¹⁷We abbreviate this with an added ‘d’ (for the demonstrations track) to distinguish this paper from our full FSE 2020 paper.

498 After our exploratory analysis on the nature of CVSSs (Chapter 4), we proposed
 499 two sets of recommendations targeted towards two stakeholders: (i) the service
 500 *consumers* (i.e., application developers) and (ii) the service *providers*. Our sub-
 501 sequent publications arose as a two-fold investigation to develop two strategies in
 502 which developers and providers can, respectively, (i) better integrate these intelli-
 503 gent components into their applications, and (ii) how these services can be better
 504 documented. Table 1.5 provides a tabulated form of the publications and research
 505 questions addressed within this thesis; for ease of reference, we refer to the publica-
 506 tions in within this section in their abbreviated form as listed in Table 1.5. We also
 507 provide abbreviations for easier reference in this section. A high-level overview of
 508 the cohesiveness of our publications is provided in Figure 1.5.

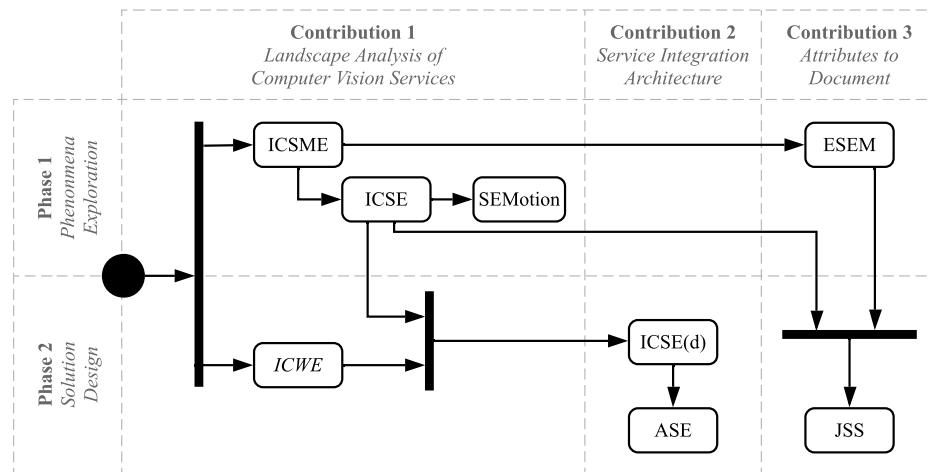


Figure 1.5: Activity diagram of the coherency of our publications, how our research was conducted, and relevant connections between publications. Our two-phase structure initial phenomena exploration and a proposed solutions to issues identified from the exploration. We map the contributions within each publication to the three primary contributions of the thesis.

509 1.7.1 Contribution 1: Landscape Analysis & Preliminary Solutions

510 The first two bodies of work in this paper are the ICSME and ICWE papers. These
 511 two works investigated a landscape analysis CVSSs from two perspectives: firstly, we
 512 conducted a longitudinal study to better understand the attributes associated with
 513 these services (ICSME)—particularly their evolution and behavioural profiles, and
 514 their potential impacts to software reliability—and tackled a preliminary solution
 515 facade to ‘merge’ responses of the services together (ICWE).

516 The ICSME paper confirmed our hypotheses that the services have a non-
 517 deterministic behavioural profile, and that the evolution occurring within the ML
 518 models powering these services are not sufficiently communicated to software en-
 519 gineers. This therefore led to follow up investigation into how developers perceive

520 these services, and thereby determine if they are frustrated due to this lack of communication.
521

522 Our ICWE paper explored one aspect identified from the ICSME paper that
523 we identified early on: that different services use different vocabularies to describe
524 semantically similar objects but in different ways (e.g., ‘border collie’ vs. ‘collie’),
525 despite offering functionally similar capabilities. We attempted to merge the re-
526 sponse labels from these services using a proportional representation approach, and
527 upon comparison with more naive merge approaches, we improved label-merge per-
528 formance by an F-measure of 0.015. However, while this was an interesting outcome
529 for a preliminary solution design, investigation from our following work suggested
530 that standardising ontologies between service providers becomes challenging and
531 normalising the entire ontological hierarchy of response labels would need to fall
532 under the responsibility of a certain body (that does not exist). Further, we did
533 not find sufficient evidence that developers would frequently switch between service
534 providers. Therefore, we opted for a shielded relay architecture in our later design
535 work.

536 **1.7.2 Contribution 2: Improving Documentation Attributes**

537 As mentioned, our ICSME paper found that evolutionary and non-deterministic
538 behavioural profile of are not adequately documented in the service’s APIs docu-
539 mentation. A recommendation concluding from this work was that service providers
540 should improve their documentation, however there lacked a strategy by which they
541 could do this, and our hypotheses that developers were actually frustrated by this
542 lack of communication was yet to be tested. This led to two follow-up further
543 investigations as presented in our ICSE and ESEM papers.

544 One aspect of our ICSE paper was to confirm whether developers are actually
545 frustrated with the service’s limited API documentation. By mining Stack Overflow
546 posts with reference to documentation issues, we adopted a 2019 documentation-
547 related taxonomy by Aghajani et al. [2] to classify posts, and found that 47.87%
548 of posts classified fell under the ‘completeness’ dimension of Aghajani et al.’s
549 taxonomy. This interpretation, therefore, warranted the recommendation proposed
550 in the ICSME paper to improve service documentation.

551 However, though improvements to more complete documentation was justified
552 from the ICSE paper, we needed to explore exactly *what* makes a ‘complete’ API
553 document. By conducting a systematic mapping study resulting in 4,501 results, we
554 curated 21 primary studies that outline the facets of API documentation knowledge.
555 From these studies, we distilled a documentation framework describing a priori-
556 tised order of the documentation assets API’s should document that is described
557 in our ESEM short paper. After receiving community feedback, we extended this
558 short paper with a follow-up experiment submitted to TSE. By conducting a sur-
559vey with developers, we assessed our API documentation taxonomy’s efficacy with
560 practitioner opinions, thereby producing a weighted taxonomy against *both* literature
561 and developer sources. Lastly, we triangulated both weightings against a heuristic
562 evaluation against common CVS providers’ documentation. This allowed us to de-

563 duce which specific areas in existing CVS providers' API documentation needed
564 improvement, which was a primary contribution from our TSE article.

565 1.7.3 Contribution 3: Service Integration Architecture

566 Two recommendations from our ICSME study encouraged developers to test their
567 applications with a representative ontology for their problem domain and to incorpo-
568 rate a specialised testing and monitoring techniques into their workflow. Strategies
569 on *how* to achieve this were explored in later studies. Following a similar approach
570 to our solution of improved API documentation, we validated the substantiveness of
571 our recommendations using our mining study of Stack Overflow (our ICSE paper)
572 to help inform us of generalised issues developers face whilst integrating CVSs into
573 their applications. To achieve this, we used a Stack Overflow post classification tax-
574 onomy proposed by Beyer et al. [34] into seven categories, where 28.9% and 20.37%
575 of posts asked issues regarding how to use the CVS API and conceptual issues be-
576 hind CVSs, respectively. Developers presented an insufficient understanding of the
577 non-deterministic runtime behaviour, functional capability, and limitations of these
578 services and are not aware of key computer vision terminology. When contrasted
579 to more conventional domains such as mobile-app development, the spread of these
580 issues vary substantially.

581 We proposed two technical solutions in our two FSE papers to help alleviate this
582 issue. 583 *(todo: Revise this... needs to be fleshed out)* Firstly, our FSE demonstrations
584 paper—FSE(d) for short—provides a workflow for developers to better select an
585 appropriate confidence threshold, and thus decision boundary, calibrated for their
586 particular use case. In our ESEC/FSE paper, we provide a reference architecture for
587 developers to guard against the non-deterministic issues that may ‘leak’ into their
588 applications. This architecture is a facade style, similar to the style proposed in our
589 ICWE paper, however, unlike the ICWE paper that uses proportional representation
590 approach to modify multiple sources, our FSE paper proposes a guarded relay,
591 whereby a single service is used, and the facade should maintain a lifecycle to
592 monitor evolution issues identified in ICSME and should be benchmarked against
593 the developer’s dataset (i.e., against the particular application domain) as suggested
in ICSE(d). These two primary contributions further serve as an answer to RQ4.

CHAPTER 2

594

595

596

Background

597

598 In Chapter 1, we defined a common set of (artificial) intelligence-based cloud ser-
599 vices that we label intelligent web services (IWSs). Specifically, we scope the
600 primary body of this study’s work on computer vision services (CVSs) (e.g., Google
601 Cloud Vision [362], AWS Rekognition [341], Azure Computer Vision [376], Watson
602 Visual Recognition [372] etc.). We claim developers have a distinctly determinis-
603 tic mindset ($2 + 2$ always equals 4) whereas an IWS’s ‘intelligence’ component (a
604 black box) may return probabilistic results ($2 + 2$ might equal 4 with a confidence
605 of 95%). Thus, there is a mindset mismatch between probabilistic results (from the
606 API provider) and results interpreted with certainty (from the API consumer).

607 What affect does this mindset mismatch have on the developer’s approach to-
608 wards building probabilistic software? What can we learn from common software
609 engineering practices (e.g., [249, 291]) that apply to resolve this mismatch and
610 thereby improve quality, such as verification & validation (V&V)? Chiefly, we an-
611 chor this question around three lenses of software engineering: creating an IWS,
612 using an IWS, and the nature of IWSs themselves.

613 Our chief concern lies with interaction and integration between IWS providers
614 and consumers, the nature of applications built using an IWS, and the impact this
615 has on software quality. We triangulate this around three pillars, which we diagram-
616 matically represent in Figure 2.1.

- 617 (1) **The development of the IWS.** We investigate the internal quality attributes
618 of creating an IWS from the IWS *provider’s* perspective. That is, we ask if
619 existing verification techniques are sufficient enough to ensure that the IWS
620 being developed actually satisfies the IWS consumer’s needs and if the internal
621 perspective of creating the system with a non-deterministic mindset clashes
622 with the outside perspective (i.e., pillar 2).
- 623 (2) **The usage of the IWS.** We investigate the external quality attributes of using
624 an IWS from the IWS *consumer’s* perspective. That is, we ask if existing
625 validation techniques are sufficient enough to ensure that the end-users can

626 actually use an IWS to build their software in the ways they expect the IWS to
 627 work.

628 **(3) The nature of an IWS.** We investigate what standard software engineering
 629 practices apply when developing non-deterministic systems. That is, we
 630 tackle what best practices exist when developing systems that are inherently
 631 stochastic and probabilistic, i.e., the ‘black box’ intelligence itself.

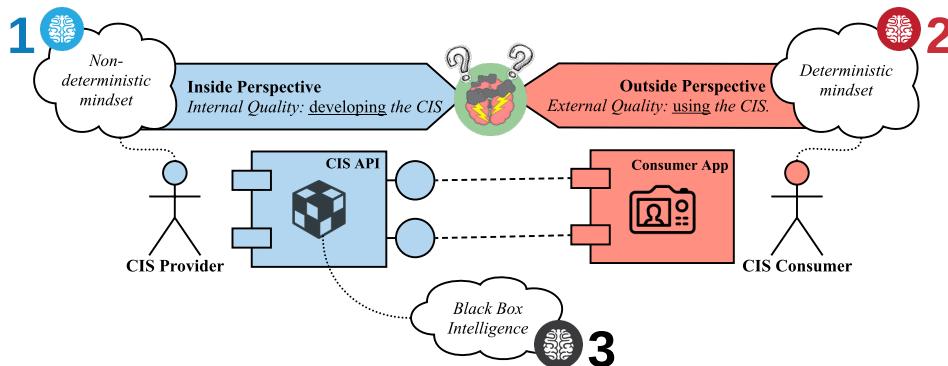


Figure 2.1: The three pillars by which we anchor the background: (1) developing an IWS with a non-deterministic mindset by the IWS provider; (2) the use of a IWS with a deterministic mindset by the IWS consumer; (3) the nature of a IWS itself.

632 Does a clash of deterministic consumer mindsets who use a IWS and the non-
 633 deterministic provider mindsets who develop them exist? And what impact does
 634 this have on the inside and outside perspective? Throughout this chapter, we will
 635 review these three core pillars due to such mindset mismatch from the anchoring per-
 636 spective of software quality, particularly around V&V and related quality attributes,
 637 probabilistic and nondeterministic software and the nature of APIs.

638 2.1 Software Quality

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

ROBERT PIRSIG, 1974 [248]

640 The philosophical viewpoint of ‘quality’ remains highly debated and there are mul-
 641 tiple facets to perceive this complex concept [115]. Transcendentally, a viewpoint
 642 like that of Pirsig’s above shows that quality is not tangible but still recognisable; it’s
 643 hard to explicitly define but you know when it’s missing. The International Orga-
 644 nization for Standardization provides a breakdown of seven universally-applicable
 645 principles that defines quality for organisations, developers, customers and training
 646 providers [148]. More pertinently, the 1986 ISO standard for quality was simply
 647 “the totality of characteristics of an entity that bear on its ability to satisfy stated or
 648 implied needs” [147].

Using this sentence, what characteristics exist for non-deterministic IWSs like that of a CVS? How do we know when the system has satisfied its ‘stated or implied needs’ when the system can only give us uncertain probabilities in its outputs? Such answers can be derived from related definitions—such as ‘conformance to specification or requirements’ [74, 119], ‘meeting or exceeding customer expectation’ [31], or ‘fitness for use’ [160]—but these then still depend on the solution description or requirements specification, and thus the same questions still apply.

Software quality is somewhat more concrete. Pressman [249] adapted the manufacturing-oriented view of quality from [32] and phrased software quality under three core pillars:

- **effective software processes**, where the infrastructure that supports the creation of quality software needs is effective, i.e., poor checks and balances, poor change management and a lack of technical reviews (all that lie in the *process* of building software, rather than the software itself) will inevitably lead to a poor quality product and vice-versa;
- **building useful software**, where quality software has fully satisfied the end-goals and requirements of all stakeholders in the software (be it explicit or implicit requirements) *in addition to* delivering these requirements in reliable and error-free ways; and lastly
- **adding value to both the producer and user**, where quality software provides a tangible value to the community or organisation using it to expedite a business process (increasing profitability or availability of information) *and* provides value to the software producers creating it whereby customer support, maintenance effort, and bug fixes are all reduced in production.

In the context of a non-deterministic IWS, however, are any of the above actually guaranteed? Given that the core of a system built using an IWS is fully dependent on the *probability* that an outcome is true, what assurances must be put in place to provide developers with the checks and balances needed to ensure that their software is built with quality? For this answer, we re-explore the concept of verification & validation (V&V).

2.1.1 Validation and Verification

To explain V&V, we analogously recount a tale given by Pham [246] on his works on reliability. A high-school student sat a standardised test that was sent to 350,0000 students [301]. A multiple-choice algebraic equation problem used a variable, a , and intended that students *assume* that the variable was non-negative. Without making this assumption explicit, there were two correct answers to the multiple choice answer. Up to 45,000 students had their scores retrospectively boosted by up to 30 points for those who ‘incorrectly’ answered, however, outcomes of a student’s higher education were, thereby, affected by this one oversight in quality assessment. The examiners wrote a poor question due to poor process standards to check if their ‘correct’ answers were actually correct. The examiners “didn’t build the right product” nor did they “build the product right” by writing an poor question and failing to ensure quality standards, in the phrases Boehm [40] coined.

692 This story describes the issues with the cost of quality [39] and the importance
 693 of V&V: just as the poorly written exam question had such a high toll the 45,000
 694 unlucky students, so does poorly written software in production. As summarised by
 695 Pressman [249], data sourced from Digital [68] in a large-scale application showed
 696 that the difference in cost to fix a bug in development versus system testing is
 697 \$6,159 per error. In safety-critical systems, such as self-driving cars or clinical
 698 decision support systems, this cost skyrockets due to the extreme discipline needed
 699 to minimise error [303].

700 Formally, we refer to the IEEE Standard Glossary of Software Engineering
 701 Terminology [145] for to define V&V:

- | | |
|--------------------------------|---|
| 702 <i>verification</i> | The process of evaluating a system or component to determine
703 whether the products of a given development phase satisfy the
704 conditions imposed at the start of that phase. |
| 705 <i>validation</i> | The process of evaluating a system or component during or at the
706 end of the development process to determine whether it satisfies
707 specified requirements. |

708 Thus, in the context of an IWS, we have two perspectives on V&V: that of the API
 709 provider and consumer (Figure 2.2).

710 The verification process of API providers ‘leak’ out to the context of the de-
 711 veloper’s project dependent on the IWS. Poor verification in the *internal quality*
 712 of the IWS will entail poor process standards, such as poor definitions and termi-
 713 nology used, support tooling and description of documentations [291]. Though
 714 it is commonplace for providers to have a ‘ship-first-fix-later’ mentality of ‘good-
 715 enough’ software [313], the consequence of doing so leads to consumers absorb-
 716 ing the cost. Thus API providers must ensure that their verification strategies
 717 are rigorous enough for the consumers in the myriad contexts they wish to use
 718 it in. Studies have considered V&V in the context of web services on the cloud
 719 [16, 58, 59, 104, 134, 220, 222, 332], though little have recently considered how
 720 adding ‘intelligence’ to these services affects existing proposed frameworks and
 721 solutions. For a CVS, what might this entail? Which assurances are given to the
 722 consumers, and how is that information communicated? To verify if the service is
 723 working correctly, does that mean that we need to deploy the system first to get a
 724 wider range of data, given the stochastic nature of the black box?

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

733 Therefore, just as the test answers with were both correct and incorrect at the
 734 same time, so is the same with IWSs returning a probabilistic result: no result is

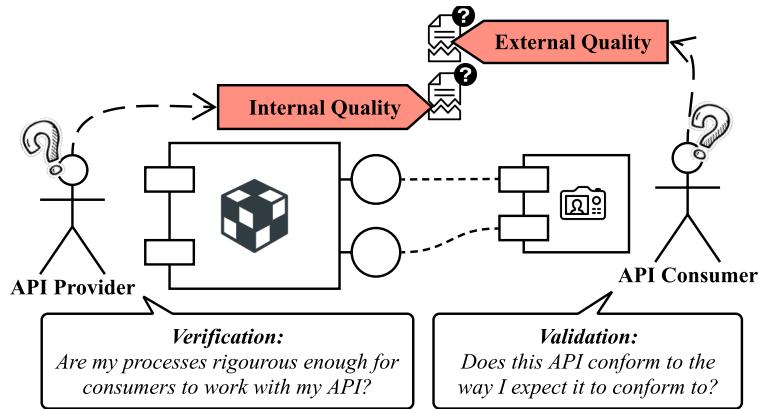


Figure 2.2: The ‘leakage’ of internal quality into the API consumer’s product and external quality imposing on the API provider.

735 certain. While V&V has been investigated in the area of mathematical and earth
 736 sciences for numerical probabilistic models and natural systems [231, 270], from
 737 the software engineering literature, little work has been achieved to look at the
 738 surrounding area of probabilistic systems hidden behind API calls.

739 Now that a developer is using a probabilistic system behind a deterministic API
 740 call, what does it mean in the context of V&V? Do current verification approaches
 741 and tools suffice, and if not, how do we fix it? From a validation perspective of
 742 ML and end-users, after a model is trained and an inference is given and if the
 743 output data point is incorrect, how will end users report a defect in the system?
 744 Compared to deterministic systems where such tooling as defect reporting forms are
 745 filled out (i.e., given input data in a given situation and the output data was X), how
 746 can we achieve similar outputs when the system is not non-deterministic? A key
 747 problem with the probabilistic mindset is that once a model is ‘fixed’ by retraining
 748 it, while one data-point may be fixed, others may now have been effected, thereby
 749 not ensuring 100% validation. Thus, due to the unpredictable and blurry nature of
 750 probabilistic systems, V&V must be re-thought out extensively.

751 2.1.2 Quality Attributes and Models

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

756 Quality attribute assessment models (like those shown in Figure 2.3) are an early
 757 concept in software engineering, and systematically evaluating software quality
 758 appears as early as 1968 [269]. Rubey and Hartwick’s 1968 study introduced the
 759 phrase ‘attributes’ as a “prose expression of the particular quality of desired software”
 760 (as worded by Boehm et al. [38]) and ‘metrics’ as mathematical parameters on a
 761 scale of 0 to 100. Early attempts to categorise wider factors under a framework was
 762 proposed by McCall, Richards, and Walters in the late 1970s [62, 203]. This model

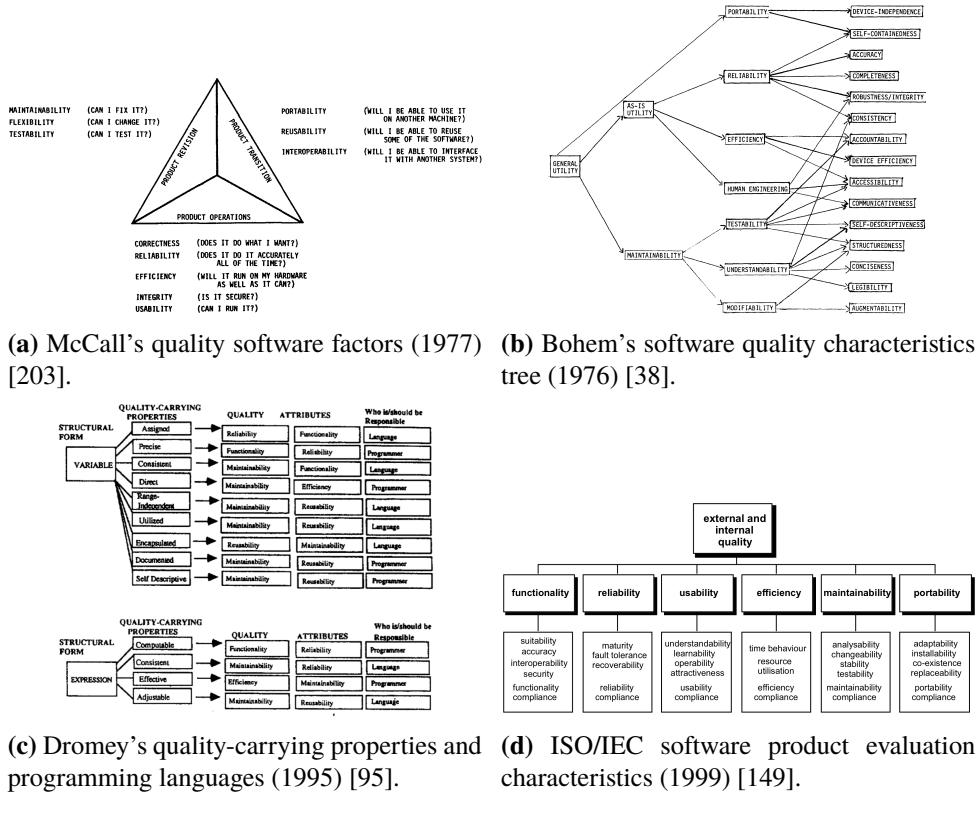


Figure 2.3: A brief overview of the development of software quality models since 1977.

described quality from the three perspectives of product revision (*how can we keep the system operational?*), transition (*how can we migrate the system as needed?*) and operation (*how effective is the system at achieving its tasks?*) (Figure 2.3a). The model also introduced 11 attributes alongside numerous direct and indirect measures to help quantify quality. This model was further developed by Boehm et al. [38] who independently developed a similar model, starting with an initial set of 11 software characteristics. It further defined candidate measurements of Fortran code to such characteristics, taking shape in a tree-like structure as in Figure 2.3b. In the mid-1990s, Dromey's interpretation [95] defined a set of quality-carrying properties with structural forms associated to specific programming languages and conventions (Figure 2.3c). The model also supported quality defect identification and proposed an improved auditing method to automate defect detection for code editors in IDEs. As the need for quality models became prevalent, the International Organization for Standardization standardised software quality under ISO/IEC-9126 [149] (the Software Product Evaluation Characteristics, Figure 2.3d), which has since recently been revised to ISO/IEC-25010 with the introduction of the Systems and software Quality Requirements and Evaluation (SQuaRE) model [146], separating quality into *Product Quality* (consisting of eight quality characteristics and 31 sub-characteristics) and *Quality In Use* (consisting of five quality characteristics and 9 sub-characteristics). An extensive review on the development of quality models in

⁷⁸³ software engineering is given in [5].

⁷⁸⁴ Of all the models described, there is one quality attribute that relates most
⁷⁸⁵ with our narrative of IWS quality: reliability. Reliability is the primary quality
⁷⁸⁶ factor investigated within this thesis (see Section 1.4). Both McCall and Boehm's
⁷⁸⁷ quality models have sub-characteristics of reliability relating to the primary research
⁷⁸⁸ questions that investigate the *robustness*, *consistency* and *completeness*¹ of CVSs
⁷⁸⁹ and its associated documentation. Moreover, the definition of reliability is similar
⁷⁹⁰ among all quality models:

⁷⁹¹ **McCall et al.** Extent to which a program can be expected to perform its in-
⁷⁹² tended function with required precision [203].

⁷⁹³ **Boehm et al.** Code possesses the characteristic *reliability* to the extent that
⁷⁹⁴ it can be expected to perform its intended functions satisfac-
⁷⁹⁵ torily [38].

⁷⁹⁶ **Dromey** Functionality implies reliability. The reliability of software is
⁷⁹⁷ therefore dependent on the same properties as functionality, that
⁷⁹⁸ is, the correctness properties of a program [95].

⁷⁹⁹ **ISO/IEC-9126** The capability of the software product to maintain a specified
⁸⁰⁰ level of performance when used under specified conditions [149].

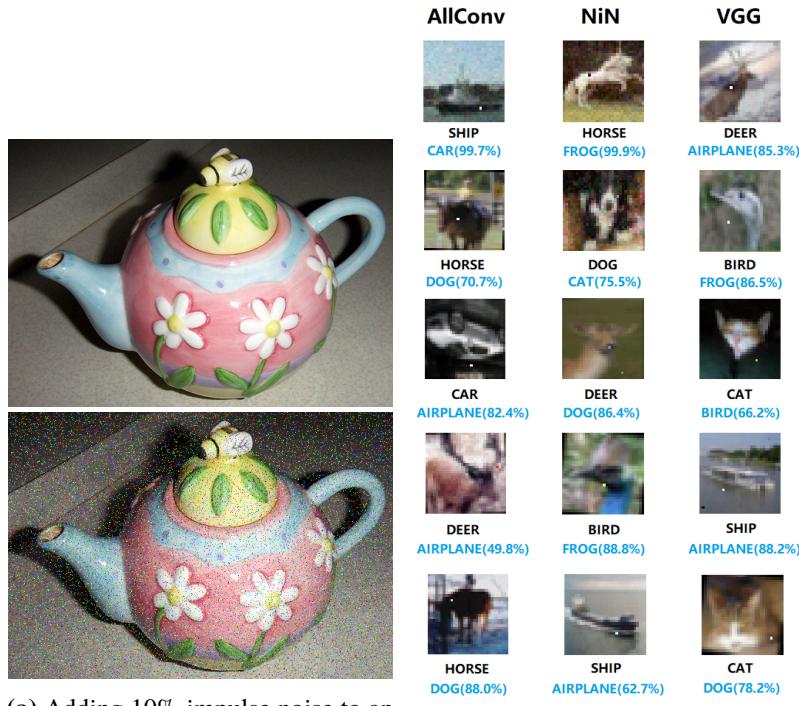
⁸⁰¹ These definitions strongly relate to the system's solution description in that
⁸⁰² reliability is the ability to maintain its *functionality* under given conditions. But what
⁸⁰³ defines reliability when the nature of an IWS in itself is inherently unpredictable
⁸⁰⁴ due to its probabilistic implementation? Can a non-deterministic system ever be
⁸⁰⁵ considered reliable when the output of the system is uncertain? How do developers
⁸⁰⁶ perceive these quality aspects of reliability in the context of such systems? A system
⁸⁰⁷ cannot be perceived as 'reliable' if the system cannot reproduce the same results due
⁸⁰⁸ to a probabilistic nature. Therefore, we believe the literature of quality models does
⁸⁰⁹ not suffice in the context of IWS reliability; a CVS can interpret an image of a dog
⁸¹⁰ as a 'Dog' one day, but what if the next it interprets such image more specifically to
⁸¹¹ the breed, such as 'Border Collie'? Does this now mean the system is unreliable?

⁸¹² Moreover, defining these systems in themselves is challenging when require-
⁸¹³ ments specifications and solution descriptions are dependent on nondeterministic
⁸¹⁴ and probabilistic algorithms. We discuss this further in Section 2.2.

⁸¹⁵ 2.1.3 Reliability in Computer Vision

⁸¹⁶ Testing computer vision deep-learning reliability is an area explored typically
⁸¹⁷ through the use of adversarial examples [299]. These input examples are where
⁸¹⁸ images are slightly perturbed to maximise prediction error but are still interpretable
⁸¹⁹ to humans. Refer to Figure 2.4.

¹In McCall's model, completeness is a sub-characteristic of the 'correctness' quality factor; however in Boehm's model it is a sub-characteristic of reliability. For consistency in this thesis, *completeness* is referred in the Boehm interpretation.



(c) Adversarial examples to trick face recognition from the source to target images [316].

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

820 Google Cloud Vision, for instance, fails to correctly classify adversarial examples
 821 when noise is added to the original images [140]. Rosenfeld et al. [267] illustrated
 822 that inserting synthetic foreign objects to input images (e.g., a cartoon elephant)
 823 can alter classification output. Wang et al. [316] performed similar attacks on a
 824 transfer-learning approach of facial recognition by modifying pixels of a celebrity's
 825 face to be recognised as a different celebrity, all while still retaining the same human-
 826 interpretable original celebrity. Su et al. [294] used the ImageNet database to show
 827 that 41.22% of images drop in confidence when just a *single pixel* is changed in the
 828 input image; and similarly, Eykholt et al. [99] recently showed similar results that
 829 made a CNN interpret a stop road-sign (with mimicked graffiti) as a 45mph speed
 830 limit sign.

831 Thus, the state-of-the-art computer vision techniques may not be reliable enough
 832 for safety critical applications (such as self-driving cars) as they do not handle intention-
 833 al or unintentional adversarial attacks. Moreover, as such adversarial examples
 834 exist in the physical world [99, 179], “the real world may be adversarial enough”
 835 [245] to fool such software.

836 2.2 Probabilistic and Nondeterministic Systems

837 Probabilistic and nondeterministic systems are those by which, for the same given
 838 input, different outcomes may result. The underlying models that power an IWS
 839 are treated as though they are nondeterministic; Chapter 2 introduces IWSs as
 840 essentially black-box behaviour that can change over time. As such, we adopt the
 841 nondeterministic behaviour that they present.

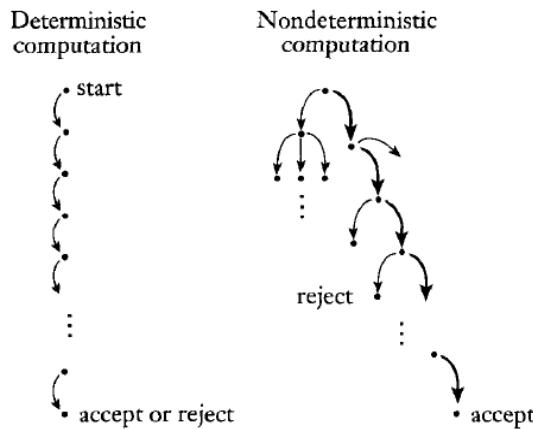


Figure 2.5: A deterministic system (left) always returns the same result in the same amount of steps. A nondeterministic system does not guarantee the same outcome, even with the same input data. Source: [103].

2.2.1 Interpreting the Uninterpretable

As the rise of applied AI increases, the need for engineering interpretability around models becomes paramount, chiefly from an external quality perspective that the *reliability* of the system can be inspected by end-users. Model interpretability has been stressed since early machine learning research in the late 1980s and 1990s (such as Quinlan [250] and Michie [214]), and although there has since been a significant body of work in the area [14, 29, 44, 55, 86, 101, 110, 118, 158, 187, 190, 200, 240, 255, 268, 288, 312, 314], it is evident that ‘accuracy’ or model ‘confidence’ is still used as a primary criterion for AI evaluation [143, 152, 290]. Much research into neural network (NN) or support vector machine (SVM) development stresses that ‘good’ models are those with high accuracy. However, is accuracy enough to justify a model’s quality?

To answer this, we revisit what it means for a model to be accurate. Accuracy is an indicator for estimating how well a model’s algorithm will work with future or unforeseen data. It is quantified in the AI testing stage, whereby the algorithm is tested against cases known by humans to have ground truth but such cases are unknown by the algorithm. In production, however, all cases are unknown by both the algorithm *and* the humans behind it, and therefore a single value of quality is “not reliable if the future dataset has a probability distribution significantly different from past data” [106], a problem commonly referred to as the *datashift* problem [274]. Analogously, Freitas [106] provides the following description of the problem:

The military trained [a NN] to classify images of tanks into enemy and friendly tanks. However, when the [NN] was deployed in the field (corresponding to “future data”), it had a poor accuracy rate. Later, users noted that all photos of friendly (enemy) tanks were taken on a sunny (overcast) day. I.e., the [NN] learned to discriminate between the colors of the sky in sunny vs. overcast days! If the [NN] had output a comprehensible model (explaining that it was discriminating between colors at the top of the images), such a trivial mistake would immediately be noted. [106]

So, why must we interpret models? While the formal definition of what it means to be *interpretable* is still somewhat disparate (though some suggestions have been proposed [190]), what is known is (i) there exists a critical trade-off between accuracy and interpretability [91, 105, 125, 157, 164, 334], and (ii) a single quantifiable value cannot satisfy the subjective needs of end-users [106]. As ever-growing domains ML become widespread², these applications engage end-users for real-world goals, unlike the aims in early ML research where the aim was to get AI working in the first place. In safety-critical systems where AI provide informativeness to humans to make the final call (see [60, 144, 166]), there is often a mismatch between the formal objectives of the model (e.g., to minimise error) and complex real-world goals, where other considerations (such as the human factors and cognitive science

²In areas such as medicine [28, 55, 98, 153, 158, 183, 241, 257, 312, 330, 336], bioinformatics [90, 107, 155, 163, 298], finance [14, 88, 144] and customer analytics [187, 314].

behind explanations³) are not realised: model optimisation is only worthwhile if they “actually solve the original [human-centred] task of providing explanation” [221] to end-users. **Therefore, when human-decision makers must be interpretable themselves [258], any AI they depend on must also be interpretable.**

Recently, discussion behind such a notion to provide legal implications of interpretability is topical. Doshi-Velez et al. [94] discuss when explanations are not provided from a legal stance—for instance, those affected by algorithmic-based decisions have a ‘right to explanation’ [197, 315] under the European Union’s GDPR⁴. But, explanations are not the only way to ensure AI accountability: theoretical guarantees (mathematical proofs) or statistical evidence can also serve as guarantees [94], however, in terms of explanations, what form they take and how they are proven correct are still open questions [190].

2.2.2 Explanation and Communication

From a software engineering perspective, explanations and interpretability are, by definition, inherently communication issues: what lacks here is a consistent interface between the AI system and the person using it. The ability to encode ‘common sense reasoning’ [204] into programs today has been achieved, but *decoding* that information is what still remains problematic. At a high level, Shannon and Weaver’s theory of communication [281] applies, just as others have done with similar issues in the SE realm [216, 325] (albeit to the domain of visual notations). Humans map the world in higher-level concepts easily when compared to AI systems: while we think of a tree first (not the photons of light or atoms that make up the tree), an algorithm simply sees pixels, and not the concrete object [94] and the AI interprets the tree inversely to humans. Therefore, the interpretation or explanation is done inversely: humans do not explain the individual neurons fired to explain their predictions, and therefore the algorithmic transparent explanations of AI algorithms (“*which neurons were fired to make this AI think this tree is a tree?*”) do not work here.

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

2.2.3 Mechanics of Model Interpretation

How do we interpret models? Methods for developing interpretation models include: decision trees [48, 72, 131, 194, 251], decision tables [15, 187] and decision sets

³Interpretations and explanations are often used interchangeably.

⁴<https://www.eugdpr.org> last accessed 13 August 2018.

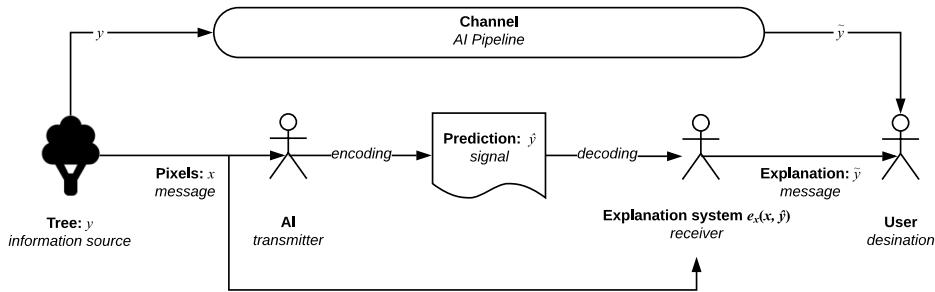


Figure 2.6: Theory of AI communication from information source, y , to intended user as explanations \tilde{y} .

[181, 221]; input gradients, gradient vectors or sensitivity analysis [14, 184, 255, 268, 279]; exemplars [108, 167]; generalised additive models [60]; classification (*if-then*) rules [45, 69, 234, 307, 327] and falling rule lists [288]; nearest neighbours [200, 252, 280, 323, 335] and Naïve Bayes analysis [28, 64, 100, 109, 135, 174, 183, 336].

Cross-domain studies have assessed the interpretability of these techniques against end-users, measuring response time, accuracy in model response and user confidence [6, 107, 132, 144, 200, 273, 295, 314], although it is generally agreed that decision rules and decision tables provide the most interpretation in non-linear models such as SVMs or NNs [107, 200, 314]. For an extensive survey of the benefits and fallbacks of these techniques, we refer to Freitas [106], Doshi-Velez et al. [94] and Doshi-Velez and Kim [93].

2.3 Application Programming Interfaces

Application programming interfaces (APIs) are the interface between a developer needs and the software components at their disposal [10] by abstracting the underlying component behind a subroutine, protocol or specific tool. Therefore, it is natural to assess internal quality (and external quality if the software is in itself a service to be used by other developers—in this case an IWS) is therefore directly related to the quality the API offers [173].

Good APIs are known to be intuitive and require less documentation browsing [247], thereby increasing developer productivity. Conversely, poor APIs are those that are hard to interpret, thereby reducing developer productivity and product quality. The consequences of this have shown a higher demand of technical support (as measured in [136]) that, ultimately, causes the maintenance to be far more expensive, a phenomenon widely known in software engineering economics (see Section 2.1.1).

While there are different types of APIs, such as software library/framework APIs for building desktop software, operating system APIs for interacting with the operating system, remote APIs for communication of varying technologies through common protocols, we focus on web APIs for communication of resources over

952 the web (being the common architecture of cloud-based services). Further information
953 on the development, usage and documentation of web APIs is provided in
954 Appendix A.1.

955 2.3.1 API Usability

956 If a developer doesn't understand the overarching concepts of the context behind the
957 API they wish to use, then they cannot formulate what gaps in their knowledge is
958 missing. For example, a developer that knows nothing about ML techniques in CV
959 cannot effectively formulate queries to help bridge those gaps in their understanding
960 to figure out more about the CVS they wish to use.

961 Balancing the understanding of the information need (both conscious and unconscious), how to phrase that need and how to query it in an information retrieval
962 system is concept long studied in the information sciences [305]. In API design,
963 the most common form to convey knowledge to developers is through annotated
964 code examples and overviews to a platform's architectural and design decisions
965 [46, 92, 218, 262] though these studies have not effectively communicated *why* these
966 artefacts are important. What makes the developer *conceptually understand* these
967 artefacts?

968 Robillard and Deline [262] conducted a multi-phase, mixed-method approach to
969 create knowledge grounded in the professional experience of 440 software engineers
970 at Microsoft of varying experience to determine what makes APIs hard to learn,
971 the results of which previously published in an earlier report [261]. Their results
972 demonstrate that 'documentation-related obstacles' are the biggest hurdle in learning
973 new APIs. One of these implications are the *intent documentation* of an API (i.e.,
974 *what is the intent for using a particular API?*) and such documentation is required
975 only where correct API usage is not self-evident, where advanced uses of the API are
976 documented (but not the intent), and where performance aspects of the API impact
977 the application developed using it. They conclude that professional developers do
978 not struggle with learning the *mechanics* of the API, but in the *understanding* of how
979 the API fits in upwards to its problem domain and downward to its implementation:

980 *In the upwards direction, the study found that developers need help
981 mapping desired scenarios in the problem domain to the content of the
982 API, and in understanding what scenarios or usage patterns the API
983 provider intends and does not intend to support. In the downwards
984 direction, developers want to understand how the API's implementation
985 consumes resources, reports errors and has side effects. [262]*

986 These results particularly corroborate to that of previous studies where developers
987 quote that they feel that existing learning content currently focuses on "how
988 to do things, not necessarily why" [228]. This thereby reiterates the conceptual
989 understanding of an API as paramount.

990 A later study by Ko and Riche [172] assessed the importance of a programmer's
991 conceptual understanding of the background behind the task before implementing the
992 task itself, a notion that we find most relevant for users of IWS APIs. While the study

994 did not focus on developing web APIs (rather implementing a Bluetooth application
995 using platform-agnostic terminology), the study demonstrated how developers show
996 little confidence in their own metacognitive judgements to understand and assess the
997 feasibility of the intent of the API and understand the vocabulary and concepts within
998 the domain (i.e., wireless connectivity). This indecision over what search results
999 were relevant in their searches ultimately hindered their progress implementing the
1000 functionality, again decreasing productivity. Ko and Riche suggest to improve API
1001 usability by introducing the background of the API and its relevant concepts using
1002 glossaries linked to tutorials to each of the major concepts, and then relate it back to
1003 how to implement the particular functionality.

1004 Thus, an analysis of the conceptual understanding of IWS APIs by a range of
1005 developers (from beginner to professional) is critical to best understand any differ-
1006 ences between existing studies and those that are nondeterministic. Our proposal is
1007 to perform similar survey research (see Chapter 3) in the search for further insight
1008 into the developer's approach toward existing IWS APIs.

CHAPTER 3

1009

1010

1011

Research Methodology

1012

1013 < *TODO: Revise this entire chapter for tense issues: JG - I did wonder about*
1014 *TENSE in Ch 3 - I didn't change but to think about - all this work*
1015 *is DONE so use either past (my pref) or present. Not "we propose*
1016 *to use..." etc etc all throughout. Especially for a by-papers thesis.*
1017 *Could revised to "we proposed to use ..." but I would suggest "We*
1018 *used ..." (my pref - past) or "We use..." (present). >*

1019 Investigating software engineering practices is often a complex task as it is im-
1020 perative to understand the social and cognitive processes around software engineers
1021 and not just the tools and processes used [97]. This chapter explores our research
1022 methodology by exploring five key elements of empirical software engineering re-
1023 search: firstly, (i) we provide an extended focus to the study by reviewing our research
1024 questions (see Section 1.4) anchored under the context of an existing research ques-
1025 tion classification taxonomy, (ii) characterise our research goals through an explicit
1026 philosophical stance, (iii) explain how the stance selected impacts our selection of
1027 research methods and data collection techniques (by dissecting our choice of meth-
1028 ods used to reach these research goals), (iv) discuss a set of criteria for assessing the
1029 validity of our study design and the findings of our research, and lastly (v) discuss
1030 the practical considerations of our chosen methods.

1031 The foundations for developing this research methodology has been expanded
1032 from that proposed by Easterbrook et al. [97], Wohlin and Aurum [328], Wohlin
1033 et al. [329] and Shaw [283].

1034 3.1 Research Questions Revisited

1035 To discuss our research strategy, we revisit our four primary and seven secondary
1036 research questions (RQs) through the classification technique discussed by East-
1037 erbrook et al. [97], a technique originally proposed in the field of psychology by

¹⁰³⁸ Meltzoff and Cooper [209] but adapted to software engineering. A summary of the
¹⁰³⁹ classifications made to our research questions are presented in Table 3.1.

¹⁰⁴⁰ Our research study involves a mix of nine *empirical*¹ RQs, that focus on observing
¹⁰⁴¹ and analysing existing phenomena, and two *non-empirical* RQs, that focuses
¹⁰⁴² on designing better approaches to solve software engineering tasks [213]. The use
¹⁰⁴³ of empirical *and* non-empirical RQs are best combined in long-term software en-
¹⁰⁴⁴ gineering research studies where the phenomena are under-explored, as is the case
¹⁰⁴⁵ with CVSs. Further, these approaches help propose solutions to issues found in the
¹⁰⁴⁶ phenomena studied [326]. We discuss both our empirical and non-empirical RQs in
¹⁰⁴⁷ Sections 3.1.1 and 3.1.2 below.

Table 3.1: A summary of our research questions classified using the strategies presented by Easterbrook et al. [97] and Meltzoff and Cooper [209].

#	RQ	Primary/ Secondary	RQ Classification
RQ1	What is the nature of cloud-based CVSs?	Primary	EMPIRICAL ↔ Exploratory ↔ Description/Classification
RQ1.1	What is their runtime behaviour?		EMPIRICAL ↔ Exploratory ↔ Description/Classification
RQ1.2	What is their evolution profile?		EMPIRICAL ↔ Exploratory ↔ Description/Classification
RQ2	Are CVS APIs sufficiently documented?	Primary	EMPIRICAL ↔ Exploratory ↔ Existence
RQ2.1	What are the dimensions of a ‘complete’ API doc- ument, according to both literature and practitioners?	Secondary	EMPIRICAL ↔ Exploratory ↔ Composition
RQ2.2	What additional information or attributes do appli- cation developers need in CVS API documentation to make it more complete?	Secondary	NON-EMPIRICAL ↔ Design
RQ3	Are CVSs more misunderstood than conventional software engineering domains?	Primary	EMPIRICAL ↔ Exploratory ↔ Descriptive-Comparative
RQ3.1	What types of issues do application developers face most when using CVSs, as expressed as questions on Stack Overflow?	Secondary	EMPIRICAL ↔ Base-Rate ↔ Frequency/Distribution
RQ3.2	Which of these issues are application developers most frustrated with?	Secondary	EMPIRICAL ↔ Exploratory ↔ Description/Classification
RQ3.3	Is the distribution CVS pain-points different to es- tablished software engineering domains, such as mobile or web development?	Secondary	EMPIRICAL ↔ Base-Rate ↔ Frequency/Distribution
RQ4	What strategies can developers employ to integrate their applications with CVSs while preserving ro- bustness and reliability?	Primary	NON-EMPIRICAL ↔ Design

¹Or ‘knowledge’ questions, that extend our *knowledge* on certain phenomena.

1048 3.1.1 Empirical Research Questions

1049 In total, nine empirically-based RQs are posed in this study to help us understand the
1050 way developers currently interact and work with web services that provide computer
1051 vision. The majority of these questions are *exploratory* questions that contribute to
1052 a landscape analysis of these services (RQ1, RQ1.1 and RQ1.2), how well they are
1053 documented (RQ2), and the issues developers currently face when using them (RQ3).
1054 Our other exploratory questions complement the answers to these questions. For
1055 instance, to understand if CVSs are sufficiently documented (an *existence* exploratory
1056 question posed in RQ2), we need to understand the components of a ‘sufficient’ or
1057 ‘complete’ API document via RQ2.1 as proposed in both the literature and by
1058 software developers. While RQ2.1 does not directly relate to CVSs, answering it
1059 gives us an understanding the components of complete API documentation, and
1060 therefore, we can assess what aspects they are missing and where improvements
1061 can be made (RQ2.2). These questions are *descriptive and classification* questions
1062 that help describe and classify what practices are in use for existing CVS API
1063 documentation and the nature behind these services. Answering these exploratory
1064 questions assists in refining preciser terms of the phenomena, ways in which we find
1065 evidence for them and ensuring the data found is valid.

1066 By answering these questions, we have a clearer understanding of the phenomena;
1067 we then follow up by posing two additional *base-rate questions* that helps
1068 provide a basis to confirm that the phenomena occurring is normal (or unusual)
1069 behaviour by investigating the patterns of phenomena’s occurrence against other
1070 phenomena. RQ3.1 is a *frequency and distribution* question to help us understand
1071 what types of issues developers often encounter most, given a lack of formal extended
1072 training in artificial intelligence. This achieves us an insight into the developer’s
1073 mindset and regular thought patterns toward these APIs. We can then contrast
1074 this distribution using our second base-rate question (RQ3.3), that assesses the
1075 distributional differences between these intelligent components and non-intelligent
1076 (conventional) software components. Combined, these two questions can help us
1077 answer how the issues raised against CVSs are different to normal Stack Overflow
1078 issues—our *descriptive-comparative* question posed in RQ3—and, similarly, we can
1079 classify and rank which issues developers find most frustrating (RQ3.2).

1080 3.1.2 Non-Empirical Research Questions

1081 RQ2.2 and RQ4 are both non-empirically-based *design questions*; they are con-
1082 cerned with ways in which we can improve a CVS by investigating what additional
1083 attributes are needed in both the documentation of CVSs and in the integration
1084 architectures developers can employ to improve reliability and robustness in their
1085 applications. They are not classified as empirical questions as we investigate what
1086 *will be* and not *what is*. By understanding the process by which developers desire
1087 additional attributes of documentation and integration strategies, we can help shape
1088 improvements to the existing designs of using CVSs.

1089 **3.2 Philosophical Stances**

1090 ⟨ *todo: JG: do you really need this section? :-)* ⟩ ⟨ *todo: AC: I am not sure – I*
1091 *thought it would be good to anchor the research per advice from Raj* ⟩

1092 Philosophical stances guide the researcher’s action by fortifying what constitutes
1093 ‘valid truth’ against a fundamental set of core beliefs [260]. In software engineer-
1094 ing, four dominant philosophical stances are commonly characterised [73, 243]:
1095 positivism (or post-positivism), constructivism (or interpretivism), pragmatism, and
1096 critical theory (or advocacy/participatory). To construct such a ‘validity of truth’,
1097 we will review these four philosophical stances in this section, and state the stance
1098 that we explicitly adopt and our reasoning for this.

1099 **3.2.0.1 Positivism**

1100 Positivists claim truth to be all observable facts, reduced piece-by-piece to smaller
1101 components which is incrementally verifiable to form truth. We do not base our
1102 work on the positivistic stance as the theories governing verifiable hypothesis must
1103 be precise from the start of the research. Moreover, due to its reductionist approach,
1104 it is difficult to isolate these hypotheses and study them in isolation from context.
1105 As our hypotheses are not context-agnostic, we steer clear from this stance.

1106 **3.2.0.2 Constructivism**

1107 Constructivists see knowledge embedded within the human context; truth is the
1108 *interpretive* observation by understanding the differences in human thought between
1109 meaning and action [171]. That is, the interpretation of the theory is just as important
1110 to the empirical observation itself. We partially adopt a constructivist stance as we
1111 attempt to model the developer’s mindset, being an approach that is rich in qualitative
1112 data on human activity.

1113 **3.2.0.3 Pragmism**

1114 Pragmatism is a less dogmatic approach that encourages the incomplete and approx-
1115 imate nature of knowledge and is dependent on the methods in which the knowledge
1116 was extracted. The utility of consensually agreed knowledge is the key outcome, and
1117 is therefore relative to those who seek utility in the knowledge—what is the useful
1118 for one person is not so for the other. While we value the utility of knowledge, it is
1119 difficult to obtain consensus especially on an ill-researched topic such as ours, and
1120 therefore we do not adopt this stance.

1121 **3.2.0.4 Critical Theory**

1122 This study chiefly adopts the philosophy of critical theory [8]. A key outcome of
1123 the study is to shift the developer’s restrictive deterministic mindset and shed light
1124 on developing a new framework actively with the developer community that seeks
1125 to improve the process of using such APIs. In software engineering, critical theory
1126 is used to “actively [seek] to challenge existing perceptions about software practice”

[97], and this study utilises such an approach to shift the mindset of CVS consumers and providers alike on how the documentation and metadata should not be written with the ‘traditional’ deterministic mindset at heart. Thus, our key philosophical approach is critical theory to seek out *what-can-be* using partial constructivism to model the current *what-is*.

3.3 Research Methods

Research methods are “a set of organising principles around which empirical data is collection and analysed” [97]. Creswell [73] suggests that strong research design is reflected when the weaknesses of multiple methods complement each other. Using a mixed-methods approach is therefore commonplace in software engineering research, typically due to the human-oriented nature investigating how software engineers work both individually (where methods from psychology may be employed) and together (where methods from sociology may be employed).

Therefore, studies in software engineering are typically performed as field studies where researchers and developers (or the artefacts they produce) are analysed either directly or indirectly [287]. The mixed-methods approach combines five classes of field study methods (or empirical strategies/studies) most relevant in empirical software engineering research [97, 162, 329]: controlled experiments, case studies, survey research, ethnographies, and action research. We chiefly adopt a mixed-methods approach to our work using the *concurrent triangulation* mixed-methods strategy [202] as it best compensates for weaknesses that exist in all research methods, and employs the best strengths of others [73].

3.3.1 Review of Relevant Research Methods

Below we review some of the research methods most relevant to our research questions as refined in Section 3.1 as presented by Easterbrook et al. [97].

3.3.1.1 Controlled Experiments

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

While we have well-defined RQs, refining them into precise, *measurable* variables is challenging due to the qualitative nature they present. A well-defined population is also critical and must be easily accessible; the varied range of beginner to expert software engineers with varied understanding of artificial intelligence concepts is required to perform controlled experiments, and thus recruitment may

¹¹⁶⁷ prove challenging. Lastly, the controlled experiment is essentially reductionist by
¹¹⁶⁸ affecting a small amount of variables of interest and controlling all others. This
¹¹⁶⁹ approach is too clinical for the practical outcomes by which our research goals aim
¹¹⁷⁰ for, and is therefore closely tied to the positivist stance.

¹¹⁷¹ **3.3.1.2 Case Studies**

¹¹⁷² Case studies investigate phenomena in their real-life context and are well-suited
¹¹⁷³ when the boundary between context and phenomena is unknown [333]. They offer
¹¹⁷⁴ understanding of how and why certain phenomena occur, thereby investigating ways
¹¹⁷⁵ cause-effect relationships can occur. They can be used to test existing theories
¹¹⁷⁶ (*confirmatory case studies*) by refuting theories in real-world contexts instead of
¹¹⁷⁷ under laboratory conditions or to generate new hypotheses and build theories during
¹¹⁷⁸ the initial investigation of some phenomena (*exploratory case studies*).

¹¹⁷⁹ Case studies are well-suited where the context of a situation plays a role in
¹¹⁸⁰ the phenomenon being studied. They also lend themselves to purposive sampling
¹¹⁸¹ rather than random sampling, and thus it is possible to selectively choose cases that
¹¹⁸² benefit our research goals and (using our critical theorist stance) select cases that
¹¹⁸³ will actively benefit our participant software engineering audience most to draw
¹¹⁸⁴ attention to situations regarded as problematic in CVS.

¹¹⁸⁵ **3.3.1.3 Survey Research**

¹¹⁸⁶ Survey research identifies characteristics of a broad population of individuals through
¹¹⁸⁷ direct data collection techniques such as interviews and questionnaires or indepen-
¹¹⁸⁸ dent techniques such as data logging. Defining that well-defined population is
¹¹⁸⁹ critical, and selecting a representative sample from it to generalise the data gathered
¹¹⁹⁰ usually assists in answering base-rate questions.

¹¹⁹¹ By identifying representative sample of the population, from beginner to ex-
¹¹⁹² perienced developers with varying understanding of CVS APIs, we can use survey
¹¹⁹³ research to assist in answering our exploratory and base-rate RQs (see Section 3.1.1)
¹¹⁹⁴ in determining the qualitative aspects of how individual developers perceive and
¹¹⁹⁵ work with the existing APIs, either by directly asking them, or by mining third-party
¹¹⁹⁶ discussion websites such as Stack Overflow (SO). Similarly, we can use this strategy
¹¹⁹⁷ to assess the developer’s understanding on what makes API documentation sufficient
¹¹⁹⁸ by assessing whether specific factors suggested from literature are useful according
¹¹⁹⁹ to developers. However, with direct survey research techniques, low response rates
¹²⁰⁰ may prove challenging, especially if no inducements can be offered for participation.

¹²⁰¹ **3.3.1.4 Ethnographies**

¹²⁰² Ethnographies investigates the understanding of social interaction within community
¹²⁰³ through field observation [264]. Resulting ethnographies help understand how soft-
¹²⁰⁴ ware engineering technical communities build practices, communication strategies
¹²⁰⁵ and perform technical work collaboratively.

1206 Ethnographies require the researcher to be highly trained in observational and
1207 qualitative data analysis, especially if the form of ethnography is participant observation.
1208 whereby the researcher is embedded of the technical community for observation.
1209 This may require the longevity of the study to be far greater than a couple of weeks,
1210 and the researcher must remain part of the project for its duration to develop enough
1211 local theories about how the community functions. While it assists in revealing
1212 subtle but important aspects of work practices within software teams, this study
1213 does not focus on the study of teams, and is therefore not a research method relevant
1214 to this project.

1215 **3.3.1.5 Action Research**

1216 Action researchers simultaneously solve real-world problems while studying the
1217 experience of solving the problem [84] by actively seeking to intervene in the
1218 situation for the purpose of improving it. A precondition is to engage with a
1219 *problem owner* who is willing to collaborate in identifying and solving the problem
1220 faced. The problem must be authentic (a problem worth solving) and must have
1221 new knowledge outcomes for those involved. It is also characterised as an iterative
1222 approach to problem solving, where the knowledge gained from solving the problem
1223 has a desirable solution that empowers the problem owner and researcher.

1224 This research is most associated to our adopted philosophical stance of critical
1225 theory. As this project is being conducted under the Applied Artificial Intelligence
1226 Institute (A^2I^2) collaboratively with engaged industry clients, we have identified a
1227 need for solving an authentic problem that industry faces. The desired outcome
1228 of this project is to facilitate wider change in the usage and development of CVSSs;
1229 thus, engaging action research as a potential method throughout the mixed-methods
1230 approach used in this research.

1231 **3.3.2 Review of Data Collection Techniques for Field Studies**

1232 Singer et al. developed a taxonomy [185, 287] showcasing data collection techniques
1233 in field studies that are used in conjunction with a variety of methods based on the
1234 level of interaction between researcher and software engineer, if any. This taxonomy
1235 is reproduced in Figure 3.1.

1236 **3.4 Research Design**

1237 This section discusses an overview of the design of methods used within the experi-
1238 ments conducted under this thesis. For each experiment, we describe an overview of
1239 the experiment grounded known methods and techniques (Sections 3.3.1 and 3.3.2)
1240 and our approach to analysing the data, as well as relating the selecting method back
1241 to a specific RQ. Details of each experiment presented in this thesis, the coherency
1242 between them, and where they can be found are given in Sections 1.6 and 1.7.

Figure 3.1: Questions asked by software engineering researchers (column 2) that can be answered by field study techniques. (From [287].)

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.

1243 **3.4.1 Landscape Analysis of Computer Vision Services**

1244 To understand the behavioural and evolutionary profiles of CVSs (i.e., RQ1), we
1245 employ a longitudinal study based around a dynamic system analysis [287]. Specifically,
1246 we employ structured observations of three services using the same dataset
1247 to understand how the responses from these services change with time. Lastly, we
1248 employ documentation analysis to assess the overall ‘picture’ of how these services
1249 are documented. Further details on this experiment is given in **Chapter 4, ??**.

1250 **3.4.2 Utility of API Documentation in Computer Vision Services**

1251 To assess whether these services are sufficiently documented (i.e., RQ2), we conduct
1252 a systematic mapping study [168, 244] of the various academic sources detailing API
1253 documentation knowledge. We then consolidate this information into a structured
1254 taxonomy following a systematic taxonomy development method specific to software
1255 engineering studies [311].

1256 We then follow the triangulation approach proposed by Mayring [202] to validate
1257 the taxonomy by use of a personal opinion survey. Kitchenham and Pfleeger [169]
1258 provide an introduction on methods used to conduct personal opinion surveys which
1259 we adopt as an initial reference in (i) shaping our survey objectives around our
1260 research goals, (ii) designing a cross-sectional survey, (iii) developing and evaluating
1261 our survey instrument, (iv) evaluating our instruments, (v) obtaining the data and
1262 (vi) analysing the data. We adapt Brooke’s systematic usability scale [50] technique
1263 by basing our research questions against a known surveying instrument.

1264 As is good practice in developing questionnaire instruments to evaluate their
1265 reliability and validity [191], we evaluate our instrument design by asking colleagues
1266 to critique it via pilot studies within A²I². This assists in identifying any problems
1267 with the questionnaire itself and with any issues that may occur with the response
1268 rate and follow-up procedures.

1269 Findings from the pilot study helps inform us for a widely distributed question-
1270 naire using snow-balling sampling. Ethics approval from the Faculty of Science,
1271 Engineering and Built Environment Human Ethics Advisory Group (SEBE HEAG)
1272 has been approved to externally conducting this survey research (see Appendix D).
1273 Further details on API these methods are detailed within **Chapter 7, Section 7.3**.

1274 **3.4.3 Developer Issues concerning Computer Vision Services**

1275 Developers typically congregate in search of discourses on issues they face in online
1276 forums, such as Stack Overflow (SO) and Quora, as well as writing their experiences
1277 in personal blogs such as Medium. The simplest of these platforms is SO (a sub-
1278 community of the Stack Exchange family of targeted communities) that specifically
1279 targets developer issues on using a simple Q&A interface, where developers can
1280 discuss technical aspects and general software development topics. Moreover, SO
1281 is often acknowledged as *the ‘go-to’ place* for developers to find high-quality code
1282 snippets that assist in their problems [296].

1283 Thus, to begin understanding the issues developers face when using CVSSs and
1284 whether there is a substantial difference to conventional domains (i.e., RQ3), we
1285 propose using repository mining on SO to help answer our research questions.
1286 Specifically, we select SO due to its targeted community of developers² and the
1287 availability of its publicly available dataset released as ‘data dumps’ on the Stack
1288 Exchange Data Explorer³ and Google BigQuery⁴. Studies conducted have also used
1289 SO to mine developer discourse [7, 17, 23, 66, 189, 225, 235, 253, 265, 289, 302,
1290 318]. Further details on how we approached the design for this study can be found
1291 in **Chapters 5 and 6, Section 5.4 and ??.**

1292 **3.4.4 Designing Improved Integration Strategies**

1293 Our improved integration strategies (i.e., RQ4) evolved organically over the dura-
1294 tion of this research project through the use of industry case studies and action
1295 research. We develop several iterative prototypes and use a mix of statistical and
1296 technical-expert assessment to analyse whether our improved integration strатегis
1297 can prove useful to developers. Further details about these approaches are detailed
1298 in **Chapters 8 to 10, Section 8.5.1 and ????. < *todo: Add more detail later* >**

²We 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 SO’s 17 million.

³<https://data.stackexchange.com/stackoverflow> last accessed 17 January 2017.

⁴<https://console.cloud.google.com/marketplace/details/stack-exchange/stack-overflow> last accessed 17 January 2017.

1299

Part II

1300

Publications

CHAPTER 4

1301

1302

1303

Identifying Evolution in Computer Vision Services[†]

1304

1305 **Abstract** Recent advances in artificial intelligence (AI) and machine learning (ML), such
1306 as computer vision (CV), are now available as intelligent web services (IWSs) and their
1307 accessibility and simplicity is compelling. Multiple vendors now offer this technology as
1308 cloud services and developers want to leverage these advances to provide value to end-users.
1309 However, there is no firm investigation into the maintenance and evolution risks arising from
1310 use of these IWSs; in particular, their behavioural consistency and transparency of their
1311 functionality. We evaluated the responses of three different IWSs (specifically CV) over 11
1312 months using 3 different data sets, verifying responses against the respective documentation
1313 and assessing evolution risk. We found that there are: (1) inconsistencies in how these
1314 services behave; (2) evolution risk in the responses; and (3) a lack of clear communication
1315 that documents these risks and inconsistencies. We propose a set of recommendations to
1316 both developers and IWS providers to inform risk and assist maintainability.

4.1 Introduction

1317 The availability of intelligent web services (IWSs) has made artificial intelligence
1318 (AI) tooling accessible to software developers and promises a lower entry barrier for
1319 their utilisation. Consider state-of-the-art computer vision (CV) analysers, which
1320 require either manually training a deep-learning classifier, or selecting a pre-trained
1321 model and deploying these into an appropriate infrastructure. Either are laborious
1322 in time, and require non-trivial expertise along with a large data set when training
1323 or customisation is needed. In contrast, IWSs providing CV (i.e., computer vision
1324 services or CVSs such as [341, 353, 354, 355, 358, 362, 370, 371, 372, 376, 388,

[†]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 *Proceedings of the 35th IEEE International Conference on Software Maintenance and Evolution*. Cleveland, OH, USA: IEEE, December 2019. DOI 10.1109/ICSME.2019.00051. ISBN 978-1-72-813094-1 pp. 333–342. Terminology has been updated to fit this thesis.

1326 389, 422, 423]) abstract these complexities behind a web application programming
1327 interface (API) call. This removes the need to understand the complexities required
1328 of machine learning (ML), and requires little more than the knowledge on how to
1329 use RESTful endpoints. The ubiquity of these services is exemplified through their
1330 rapid uptake in applications such as aiding the vision-impaired [83, 254].

1331 While IWSs have seen quick adoption in industry, there has been little work
1332 that has considered the software quality perspective of the risks and impacts posed
1333 by using such services. In relation to this, there are three main challenges: (1)
1334 incorporating stochastic algorithms into software that has traditionally been deter-
1335 ministic; (2) the general lack of transparency associated with the ML models; and
1336 (3) communicating to application developers.

1337 ML typically involves use of statistical techniques that yield components with
1338 a non-deterministic external behaviour; that is, for the same given input, different
1339 outcomes may result. However, developers, in general, are used to libraries and small
1340 components behaving predictably, while systems that rely on ML techniques work
1341 on confidence intervals¹ and probabilities. For example, the developer’s mindset
1342 suggests that an image of a border collie—if sent to three intelligent computer vision
1343 services (CVSs)—would return the label ‘dog’ consistently with time regardless
1344 of which service is used. However, one service may yield the specific dog breed,
1345 ‘border collie’, another service may yield a permutation of that breed, ‘collie’, and
1346 another may yield broader results, such as ‘animal’; each with results of varying
1347 confidence values.² Furthermore, the third service may evolve with time, and
1348 thus learn that the ‘animal’ is actually a ‘dog’ or even a ‘collie’. The outcomes
1349 are thus behaviourally inconsistent between services providing conceptually similar
1350 functionality. As a thought exercise, consider if the sub-string function were created
1351 using ML techniques—it would perform its operation with a confidence where the
1352 expected outcome and the AI inferred output match as a *probability*, rather than a
1353 deterministic (constant) outcome. How would this affect the developers’ approach
1354 to using such a function? Would they actively take into consideration the non-
1355 deterministic nature of the result?

1356 Myriad software quality models and software engineering (SE) practices advo-
1357 cate maintainability and reliability as primary characteristics; stability, testability,
1358 fault tolerance, changeability and maturity are all concerns for quality in software
1359 components [139, 249, 291] and one must factor these in with consideration to
1360 software evolution challenges [121, 122, 211, 212, 306]. However, the effect this
1361 non-deterministic behaviour has on quality when masked behind an IWS is still
1362 under-explored to date in SE literature, to our knowledge. Where software depends
1363 on IWSs to achieve functionality, these quality characteristics may not be achieved,
1364 and developers need to be wary of the unintended side effects and inconsistency that
1365 exists when using non-deterministic components. A CVS may encapsulate deep-
1366 learning strategies or stochastic methods to perform image analysis, but developers

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 CVSs.

1367 are more likely to approach IWSs with a mindset that anticipates consistency. Al-
1368 though the documentation does hint at this non-deterministic behaviour (i.e., the
1369 descriptions of ‘confidence’ in various CVSSs suggest they are not always confi-
1370 dent, and thus not deterministic [339, 360, 377]), the integration mechanisms offered
1371 by popular vendors do not seem to fully expose the nuances, and developers are not
1372 yet familiar with the trade-offs.

1373 Do popular CVSSs, as they currently stand, offer consistent behaviour, and if
1374 not, how is this conveyed to developers (if it is at all)? If CVSSs are to be used in
1375 production services, do they ensure quality under rigorous service quality assurance
1376 (SQA) frameworks [139]? What evolution risk [121, 122, 211, 212] do they pose
1377 if these services change? To our knowledge, few studies have been conducted to
1378 investigate these claims. This paper assesses the consistency, evolution risk and
1379 consequent maintenance issues that may arise when developers use IWSs. We
1380 introduce a motivating example in Section 4.2, discussing related work and our
1381 methodology in Section 4.3 and ???. We present and interpret our findings in ???.
1382 We argue with quantified evidence that these IWSs can only be considered with a
1383 mature appreciation of risks, and we make a set of recommendations in ???.

1384 4.2 Motivating Example

1385 Consider Rosa, a software developer, who wants to develop a social media photo-
1386 sharing mobile app that analyses her and her friends photos on Android and iOS.
1387 Rosa wants the app to categorise photos into scenes (e.g., day vs. night, outdoors
1388 vs. indoors), generate brief descriptions of each photo, and catalogue photos of her
1389 friends as well as common objects (e.g., all photos with a dog, all photos on the
1390 beach).

1391 Rather than building a CV engine from scratch, Rosa thinks she can achieve this
1392 using one of the popular CVSSs (e.g., [341, 353, 354, 355, 358, 362, 370, 371, 372,
1393 376, 388, 389, 422, 423]). However, Rosa comes from a typical software engineering
1394 background with limited knowledge of the underlying deep-learning techniques
1395 and implementations as currently used in CV. Not unexpectedly, she internalises a
1396 mindset of how such services work and behave based on her experience of using
1397 software libraries offered by various SDKs. This mindset assumes that different
1398 cloud vendor image processing APIs more-or-less provide similar functionality,
1399 with only minor variations. For example, cloud object storage for Amazon S3 is
1400 both conceptually and behaviourally very similar to that of Google Cloud Storage
1401 or Azure Storage. Rosa assumes the CVSSs of these platforms will, therefore, likely
1402 be very similar. Similarly, consider the string libraries Rosa will use for the app.
1403 The conceptual and behavioural similarities are consistent; a string library in Java
1404 (Android) is conceptually very similar to the string library she will use in Swift
1405 (iOS), and likewise both behave similarly by providing the same results for their
1406 respective sub-string functionality. However, **unlike the cloud storage and string**
1407 **libraries, different CVSSs often present conceptually similar functionality but**
1408 **are behaviourally very different.** IWS vendors also hide the depth of knowledge
1409 needed to use these effectively—for instance, the training data set and ontologies

¹⁴¹⁰ used to create these services are hidden in the documentation. Thus, Rosa isn't even
¹⁴¹¹ exposed to this knowledge as she reads through the documentation of the providers
¹⁴¹² and, thus, Rosa makes the following assumptions:

- ¹⁴¹³ • **"I think the responses will be consistent amongst these CVSSs."** When Rosa
¹⁴¹⁴ uploads a photo of a dog, she would expect them all to respond with 'dog'. If
¹⁴¹⁵ Rosa decides to switch which service she is using, she expects the ontologies
¹⁴¹⁶ to be compatible (all CVSSs *surely* return dog for the same image) and therefore
¹⁴¹⁷ she can expect to plug-in a different service should she feel like it making only
¹⁴¹⁸ minor code modifications such as which endpoints she is relying on.
- ¹⁴¹⁹ • **"I think the responses will be constant with time."** When Rosa uploads the
¹⁴²⁰ photo of a dog for testing, she expects the response to be the same in 10 weeks
¹⁴²¹ time once her app is in production. Hence, in 10 weeks, the same photo of the
¹⁴²² dog should return the same label.

¹⁴²³ 4.3 Related Work

¹⁴²⁴ If we were to view CVSSs through the lenses of an SQA framework, robustness,
¹⁴²⁵ consistency, and maintainability often feature as quality attributes in myriad soft-
¹⁴²⁶ ware quality models (e.g., [149]). Software quality is determined from two key
¹⁴²⁷ dimensions: (1) in the evaluation of the end-product (external quality) and (2) the
¹⁴²⁸ assurances in the development processes (internal quality) [249]. We discuss both
¹⁴²⁹ perspectives of quality within the context of our work in this section.

¹⁴³⁰ 4.3.1 External Quality

CHAPTER 5

1431

1432

1433

Interpreting Pain-Points in Computer Vision Services[†]

1434

1435 **Abstract** Intelligent web services (IWSs) are becoming increasingly more pervasive; ap-
1436 plication developers want to leverage the latest advances in areas such as computer vision
1437 (CV) to provide new services and products to users, and large technology firms enable
1438 this via RESTful APIs. While such APIs promise an easy-to-integrate on-demand machine
1439 intelligence, their current design, documentation and developer interface hides much of the
1440 underlying machine learning techniques that power them. Such APIs look and feel like
1441 conventional APIs but abstract away data-driven probabilistic behaviour—the implications
1442 of a developer treating these APIs in the same way as other, traditional cloud services, such
1443 as cloud storage, is of concern. The objective of this study is to determine the various
1444 pain-points developers face when implementing systems that rely on the most mature of
1445 these intelligent web services, specifically those that provide CV. We use Stack Overflow
1446 to mine indications of the frustrations that developers appear to face when using com-
1447 puter vision services, classifying their questions against two recent classification taxonomies
1448 (documentation-related and general questions). We find that, unlike mature fields like mo-
1449 bile development, there is a contrast in the types of questions asked by developers. These
1450 indicate a shallow understanding of the underlying technology that empower such systems.
1451 We discuss several implications of these findings via the lens of learning taxonomies to
1452 suggest how the software engineering community can improve these services and comment
1453 on the nature by which developers use them.

5.1 Introduction

1454 The availability of recent advances in artificial intelligence (AI) over simple RESTful
1455 end-points offers application developers new opportunities. These new intelligent

[†]This chapter is originally based on A. Cummaudo, R. Vasa, S. Barnett, J. Grundy, and M. Abd-
elrazek, “Interpreting Cloud Computer Vision Pain-Points: A Mining Study of Stack Overflow,” in
Proceedings of the 42nd International Conference on Software Engineering. Seoul, Republic of
Korea: IEEE, October 2020, In Press. Terminology has been updated to fit this thesis.

¹⁴⁵⁷ web services (IWSs) are AI components that abstract complex machine learning
¹⁴⁵⁸ (ML) and AI techniques behind simpler API calls. In particular, they hide (either
¹⁴⁵⁹ explicitly or implicitly) any data-driven and non-deterministic properties inherent
¹⁴⁶⁰ to the process of their construction. The promise is that software engineers can
¹⁴⁶¹ incorporate complex machine learnt capabilities, such as computer vision (CV), by
¹⁴⁶² simply calling an API end-point.

¹⁴⁶³ The expectation is that application developers can use these AI-powered services
¹⁴⁶⁴ like they use other conventional software components and cloud services (e.g., object
¹⁴⁶⁵ storage like AWS S3). Furthermore, the documentation of these AI components is
¹⁴⁶⁶ still anchored to the traditional approach of briefly explaining the end-points with
¹⁴⁶⁷ some information about the expected inputs and responses. The presupposition
¹⁴⁶⁸ is that developers can reason and work with this high level information. These
¹⁴⁶⁹ services are also marketed to suggest that application developers do not need to fully
¹⁴⁷⁰ understand how these components were created (i.e., assumptions in training data
¹⁴⁷¹ and training algorithms), the ways in which the components can fail, and when such
¹⁴⁷² components should and should not be used.

¹⁴⁷³ The nuances of ML and AI powering IWSs have to be appreciated, as there are
¹⁴⁷⁴ real-world consequences to software quality for applications that depend on them if
¹⁴⁷⁵ they are ignored [76]. This is especially true when ML and AI are abstracted and
¹⁴⁷⁶ masked behind a conventional-looking API call, yet the mechanisms behind the API
¹⁴⁷⁷ are data-dependent, probabilistic and potentially non-deterministic [229]. We are
¹⁴⁷⁸ yet to discover what long-term impacts exist during development and production due
¹⁴⁷⁹ to poor documentation that do not capture these traits, nor do we know the depth of
¹⁴⁸⁰ understanding application developers have for these components. Given the way AI-
¹⁴⁸¹ powered services are currently presented, developers are also likely to reason about
¹⁴⁸² these new services much like a string library or a cloud data storage service. That
¹⁴⁸³ is, they may not fully consider the implications of the underlying statistical nature
¹⁴⁸⁴ of these new abstractions or the consequent impacts on productivity and quality.

¹⁴⁸⁵ Typically, when developers are unable to correctly align to the mindset of the
¹⁴⁸⁶ API designer, they attempt to resolve issues by (re-)reading the API documentation.
¹⁴⁸⁷ If they are still unable to resolve these issues on their own after some internet
¹⁴⁸⁸ searching, they consider online discussion platforms (e.g., Stack Overflow, GitHub
¹⁴⁸⁹ Issues, Mailing Lists) where they seek technological advice from their peers [3].
¹⁴⁹⁰ Capturing what developers discuss on these platforms offers an insight into the
¹⁴⁹¹ frustrations developers face when using different software components as shown
¹⁴⁹² by recent works [33, 165, 265, 293, 319]. However, to our knowledge, no studies
¹⁴⁹³ have yet analysed what developers struggle with when using the new generation of
¹⁴⁹⁴ *intelligent* services. Given the re-emergent interest in AI and the anticipated value
¹⁴⁹⁵ from this technology [193], a better understanding of issues faced by developers
¹⁴⁹⁶ will help us improve the quality of services. Our hypothesis is that application
¹⁴⁹⁷ developers do not fully appreciate the probabilistic nature of these services, nor do
¹⁴⁹⁸ they have sufficient appreciation of necessary background knowledge—however, we
¹⁴⁹⁹ do not know the specific areas of concern. The motivation for our study is to inform
¹⁵⁰⁰ API designers on which aspects to focus in their documentation, education, and
potentially refine the design of the end-points.

1502 This study involves an investigation of 1,825 Stack Overflow (SO) posts regarding
1503 one of the most mature types of IWSs—computer vision services (CVSs)—dating
1504 from November 2012 to June 2019. We adapt existing methodologies of prior SO
1505 analyses [33, 302] to extract posts related to CVSs. We then apply two existing SO
1506 question classification schemes presented at ICPC and ICSE in 2018 and 2019 [3, 34].
1507 These previous studies focused on mobile apps and web applications. Although not
1508 a direct motivation, our work also serves as a validation of the applicability of these
1509 two issue classification taxonomies [3, 34] in the context of IWSs (hence potential
1510 for generalisation). Additionally our work is the first—to our knowledge—to *test*
1511 the applicability of these taxonomies in a new study.

1512 The taxonomies in previous works focus on the specific aspects from the domain
1513 (e.g. API usage, specificity within the documentation etc.) and as such do not
1514 deeply consider the learning gap of an application developer. To explore the API
1515 learning implications raised by our SO analysis, we applied an additional lens of
1516 two taxonomies from the field of pedagogy. This was motivated by the need to offer
1517 an insight into the work needed to help developers learn how to use these relatively
1518 new services.

1519 The key findings of our study are:

- 1520 • The primary areas that developers raise as issues reflect a relatively primitive
1521 understanding of the underlying concepts of data-driven ML approaches used.
1522 We note this via the issues raised due to conceptual misunderstanding and
1523 confusion in interpreting errors,
- 1524 • Developers predominantly encounter a different distribution of issue types than
1525 were reported in previous studies, indicating the complexity of the technical
1526 domain has a non-trivial influence on intelligent API usage; and
- 1527 • Most of these issues can be resolved with better documentation, based on our
1528 analysis.

1529 The paper also offers a data-set as an additional contribution to the research
1530 community and to permit replication [357]. The paper structure is as follows:
1531 Section 5.2 provides motivational examples to highlight the core focus of our study;
1532 Section 5.3 provides a background on prior studies that have mined SO to gather
1533 insight into the software engineering (SE) community; Section 5.4 describes our
1534 study design in detail; Section 5.5 presents the findings from the SO extraction;
1535 Section 5.6 offers an interpretation of the results in addition to potential implications
1536 that arise from our work; Section 5.7 outlines the limitations of our study; concluding
1537 remarks are given in Section 5.8.

1538 5.2 Motivation

1539 “Intelligent” services are often available as a cloud end-point and provide devel-
1540 opers a friendly approach to access recent AI/ML advances without being experts
1541 in the underlying processes. Figure 5.1 highlights how these services abstract
1542 away much of the technical know-how needed to create and operationalise these
1543 IWSs [232]. In particular, they hide information about the training algorithm and

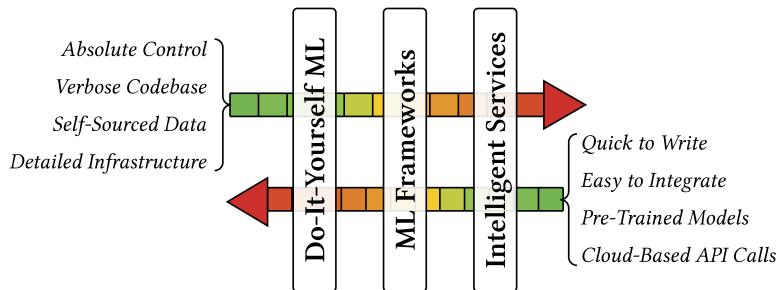


Figure 5.1: Some traits of Intelligent Services vs. ‘Do-It-Yourself’ ML. Green-to-red arrows indicate the presence of these traits. *Adapted from Ortiz [232].*

1544 data-sets used in training, the evaluation procedures, the optimisations undertaken,
 1545 and—surprisingly—they often do not offer a properly versioned end-point [76, 229].
 1546 That is, the cloud vendors may change the behaviour of the services without sufficient
 1547 transparency.

1548 The trade-off towards ease of use for application developers, coupled with the
 1549 current state of documentation (and assumed developer background) has a cost as
 1550 reflected in the increasing discussions on developer communities such as SO (see
 1551 Figure 5.2). To illustrate the key concerns, we list below a few up-voted questions:

- 1552 • **unsure of ML specific vocabulary:** “*Though it’s now not SO clear to me
 1553 what ‘score’ actually means.”* [399]; “*I’m trying out the [IWS], and there’s a
 1554 score field that returns that I’m not sure how to interpret [it].”* [413]
- 1555 • **frustrated about non-deterministic results:** “*Often the API has troubles
 1556 in recognizing single digits... At other times Vision confuses digits with
 1557 letters.”* [412]; “*Is there a way to help the program recognize numbers better,
 1558 for example limit the results to a specific format, or to numbers only?”* [409]
- 1559 • **unaware of the limitations behind the services:** “*Is there any API available
 1560 where we can recognize human other body parts (Chest, hand, legs and other
 1561 parts of the body), because as per the Google vision API it’s only able to detect
 1562 face of the human not other parts.”* [393]
- 1563 • **seeking further documentation:** “*Does anybody know if Google has pub-
 1564 lished their full list of labels ([‘produce’, ‘meal’, . . .]) and where I
 1565 could find that? Are those labels structured in any way? - e.g. is it known
 1566 that ‘food’ is a superset of ‘produce’, for example.”* [396]

1567 The objective of our study is to better understand the nature of the questions
 1568 that developers raise when using IWSs, in order to inform the service designers
 1569 and documenters. In particular, the knowledge we identify can be used to improve
 1570 the documentation, educational material and (potentially) the information contained
 1571 in the services’ response objects—these are the main avenues developers have to
 1572 learn and reason about when using these services. There is previous work that has
 1573 investigated issues raised by developers [3, 34, 302]. We build on top of this work
 1574 by adapting the study methodology and apply the taxonomies offered to identify the
 1575 nature of the issues and this results in the following research questions in this paper:

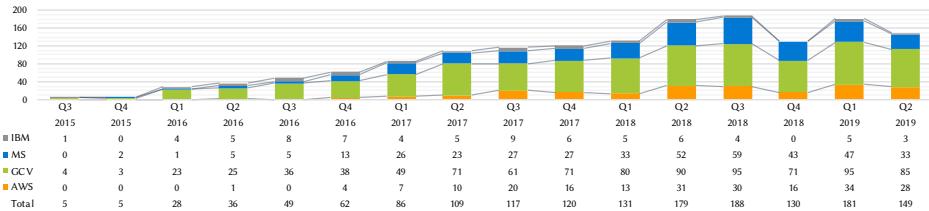


Figure 5.2: Trend of posts, where IBM = IBM Watson Visual Recognition, MS = Azure Computer Vision, AWS = AWS Rekognition and GCV = Google Cloud Vision. Three MS posts from Q4 2012, Q3 2013 and Q4 2013 have been removed for graph clarity.

- 1576 **RQ1. How do developers mis-comprehend IWSs as presented within SO**
 1577 **pain-points?** While the AI community is well aware in the the nuances that
 1578 empower IWSs, such services are being released for application developers
 1579 who may not be aware of their limitations or how they work. This is
 1580 especially the case when machine intelligence is accessed via web-based
 1581 APIs where such details are not fully exposed.
 1582 **RQ2. Are the distribution of issues similar to prior studies?** We compare
 1583 how the distributions of previous studies' of posts about conventional,
 1584 deterministic API services differ from those of IWSs. By assessing the
 1585 distribution of IWSs' issues against similar studies that focus on mobile
 1586 and web development, we identify whether a new taxonomy is needed
 1587 specific to AI-based services, and if gaps specific to AI knowledge exist
 1588 that need to be captured in these taxonomies.

1589 5.3 Background

1590 The primary goal of analysing issues is to better understand the root causes. Hence,
 1591 a good issue classification taxonomy should ideally capture the underlying causal
 1592 aspects (instead of pure functional groupings) [65]. Although this idea (of cause
 1593 related classification) is not new (Chillarege advocated for it in this TSE paper in
 1594 1992), this is not a universally followed approach when studying online discussions
 1595 and some recent works have largely classified issues into the “*what is*” and not
 1596 “*how to fix it*” [23, 33, 309]. They typically (manually) classify discussion into
 1597 either *functional areas* (e.g., Website Design/CSS, Mobile App Development, .NET
 1598 Framework, Java [23]) or *descriptive areas* (e.g., Coding Style/Practice, Problem/-
 1599 Solution, Design, QA [23, 309]). As a result, many of these studies do not give
 1600 us a prioritised means of targeted attack on how to *resolve* these issues with, for
 1601 example, improved documentation. Interestingly, recent taxonomies that studied SO
 1602 data (Aghajani et al. [3] and Beyer et al. [34]) were causal in nature and developed to
 1603 understand discussions related to mobile and web applications. However, issues that
 1604 arise when developers use IWSs have not been studied, nor do we know if existing
 1605 issue classification taxonomies are sufficient in this domain.

1606 Researchers studying APIs have also attempted to understand developer's opin-
 1607 ions towards APIs [309], categorise the questions they ask about these APIs [23,

1608 25, 34, 265], and understand API related documentation and usage issues [3, 4, 7,
1609 23, 141, 302]. These studies often employ automation to assist in the data analysis
1610 stages of their research. Latent Dirichlet Allocation [7, 23, 265, 309] is applied for
1611 topic modelling and other ML techniques such as Random Forests [34], Conditional
1612 Random Fields [4] or Support Vector Machines [34, 141] are also used.

1613 However, automatic techniques are tuned to classify into *descriptive* categories,
1614 that is, they help paint a landscape of *what is*, but generally do not address the
1615 causal factors to address the issues in great detail. For example, functional areas
1616 such as ‘Website Design’ [23], ‘User Interface’ [33] or ‘Design’ [310] result from
1617 such analyses. These automatic approaches are generally non-causal, making it hard
1618 to address reasons for *why* developers are asking such questions. However, not all
1619 studies in the space use automatic techniques; other studies employ manual thematic
1620 analysis [3, 25, 302] (e.g., card sorting) or a combination of both [33, 34, 265, 308].
1621 Our work uses a manual approach for classification, and we use taxonomies that
1622 are more causally aligned allowing our findings to be directly useful in terms of
1623 addressing the issues.

1624 Evidence-based SE [170] has helped shape the last 15 years worth of research,
1625 but the reliability of such evidence has been questioned [159, 161, 284]. Replication
1626 studies, especially in empirical works, can give us the confidence that existing results
1627 are adaptable to new domains; in this context, we extend (to IWSs) and work with
1628 study methods developed in previous works.

1629 5.4 Method

1630 5.4.1 Data Extraction

1631 This study initially attempted to capture SO posts on a broad range of many IWSs by
1632 identifying issues related to four popular IWS cloud providers: Google Cloud [362],
1633 AWS [341], Azure [376] and IBM Cloud [372]. We based our selection criteria on
1634 the prominence of the providers in industry (Google, Amazon, Microsoft, IBM) and
1635 their ubiquity in cloud platform services. Additionally, in 2018, these services were
1636 considered the most adopted cloud vendors for enterprise applications [259].

1637 However, during the filtering stage (see Section 5.4.2), we decided to focus
1638 on a subset of these services, CV, as these are one of the more mature and sta-
1639 ble ML/AI-based services with widespread and increasing adoption in the de-
1640 veloper community (see Figure 5.2). We acknowledge other services beyond the
1641 four analysed provide similar capabilities [354, 355, 358, 371, 422, 423] and only
1642 English-speaking services have been selected, excluding popular services from Asia
1643 (e.g., [352, 353, 370, 388, 389])—see Section 5.7. For comprehensiveness, we
1644 explain below our initial attempts to extract *all* IWSs.

1645 5.4.1.1 Defining a list of IWSs

1646 As there exists no global ‘list’ of IWSs to search on, we needed to derive a *corpus*
1647 of *initial terms* to allow us to know *what* to search for on the Stack Exchange Data

1648 Explorer¹ (SEDE). We began by looking at different brand names of cloud services
1649 and their permutations (e.g., Google Cloud Services and GCS) as well as various
1650 ML-related products (e.g., Google Cloud ML). To do this, we performed extensive
1651 Google searches² in addition to manually reviewing six ‘overview’ pages of the
1652 relevant cloud platforms. We identified 91 initial IWSs to incorporate into our
1653 search terms³.

1654 5.4.1.2 *Manual search for relevant, related terms*

1655 We then ran a manual search² on each term to determine if these terms were relevant.
1656 We did this by querying each term within SO’s search feature, reviewing the titles
1657 and body post previews of the first three pages of results (we did not review the
1658 answers, only the questions). We also noted down the user-defined *Tags* of each post
1659 (up to five per question); by clicking into each tag, we could review similar tags (e.g.,
1660 ‘project-oxford’ for ‘azure-cognitive-services’) and check if the tag had synonyms
1661 (e.g., ‘aws-lex’ and ‘amazon-lex’). We then compiled a *corpus of tags* consisting of
1662 31 terms.

1663 5.4.1.3 *Developing a search query*

1664 We recognise that searching SEDE via *Tags* exclusively can be ineffective (see [23,
1665 302]). To mitigate this, we produced a *corpus of title and body terms*. Such terms
1666 are those that exist within the title and body of the posts to reflect the ways in which
1667 individual developers commonly use to refer to different IWSs. To derive at such
1668 a list, we performed a search^{2,3} of the 31 tags above in SEDE, filtering out posts
1669 that were not answers (i.e., questions only) as we wanted to see how developers
1670 phrase their questions. For each search, we extracted a random sample of 100
1671 questions (400 total for each service) and reviewed each question. We noted many
1672 patterns in the permutations of how developers refer to these services, such as:
1673 common misspellings (‘bind’ vs. ‘bing’); brand misunderstanding (‘Microsoft CV’
1674 vs. ‘Azure CV’); hyphenation (‘Auto-ML’ vs. ‘Auto ML’); UK and US English
1675 (‘Watson Analyser’ vs. ‘Watson Analyzer’); and, the use of apostrophes, plurals,
1676 and abbreviations (‘Microsoft’s Computer Vision API’, ‘Microsoft Computer Vision
1677 Services’, ‘GCV’ vs. ‘Google Cloud Vision’). We arrived at a final list of 229 terms
1678 compromising all of the IWSs provided by Google, Amazon, Microsoft and IBM as
1679 of January 2019³.

1680 5.4.1.4 *Executing our search query*

1681 Our next step was to perform a case-insensitive search of all 229 terms within the
1682 body or title of posts. We used Google BigQuery’s public data-set of SO posts⁴ to
1683 overcome SEDE’s 50,000 row limit and to conduct a case-insensitive search. This

¹<http://data.stackexchange.com/stackoverflow>

²This search was conducted on 17 January 2019

³For reproducibility, this is available at <http://bit.ly/2ZcwNJO>.

⁴<http://bit.ly/2LrN7OA>

¹⁶⁸⁴ search was conducted on 10 May 2019, where we extracted 21,226 results. We then
¹⁶⁸⁵ performed several filtering steps to cleanse our extracted data, as explained below.

¹⁶⁸⁶ 5.4.2 Data Filtering

¹⁶⁸⁷ 5.4.2.1 Refining our inclusion/exclusion criteria

¹⁶⁸⁸ We performed an initial manual filtering of the 50 most recent posts (sorted by
¹⁶⁸⁹ descending *CreationDate* values) of the 21,226 posts above, assessing the suitability
¹⁶⁹⁰ of the results and to help further refine our inclusion and exclusion criteria. We
¹⁶⁹¹ did note that some abbreviations used in the search terms (e.g., ‘GCV’, ‘WCS’⁵),
¹⁶⁹² resulting in irrelevant questions in our result set. We therefore removed abbreviations
¹⁶⁹³ from our search query and consolidated all overlapping terms (e.g., ‘Google Vision
¹⁶⁹⁴ API’ was collapsed into ‘Google Vision’).

¹⁶⁹⁵ We also recognised that 21,226 results would be non-trivial to analyse without
¹⁶⁹⁶ automated techniques. As we wanted to do manual qualitative analysis, we reduced
¹⁶⁹⁷ our search space to 27 search terms of just the *CVSs* within the original corpus of
¹⁶⁹⁸ 229 terms. These were Google Cloud Vision [362], AWS Rekognition [341], Azure
¹⁶⁹⁹ Computer Vision [376], and IBM Watson Visual Recognition [372]. This resulted
¹⁷⁰⁰ in 1,425 results that were extracted on 21 June 2019. The query used and raw results
¹⁷⁰¹ are available online in our supplementary materials [357].

¹⁷⁰² 5.4.2.2 Duplicates

¹⁷⁰³ Within 1,425 results, no duplicate questions were noted, as determined by unique
¹⁷⁰⁴ post ID, title or timestamp.

¹⁷⁰⁵ 5.4.2.3 Automated and manual filtering

¹⁷⁰⁶ To assess the suitability and nature of the 1,425 questions extracted, the first author
¹⁷⁰⁷ began with a manual check on a randomised sample of 50 questions. As the questions
¹⁷⁰⁸ were exported in a raw CSV format (with HTML tags included in the post’s body), we
¹⁷⁰⁹ parsed the questions through an ERB templating engine script⁶ in which the ID, title,
¹⁷¹⁰ body, tags, created date, and view, answer and comment counts were rendered for
¹⁷¹¹ each post in an easily-readable format. Additionally, SQL matches in the extraction
¹⁷¹² process were also highlighted in yellow (i.e., in the body of the post) and listed at
¹⁷¹³ the top of each post. These visual cues helped to identify 3 false positive matches
¹⁷¹⁴ where library imports or stack traces included terms within our corpus of 26 CVS
¹⁷¹⁵ terms. For example, `aws-java-sdk-rekognition:jar` is falsely matched as a
¹⁷¹⁶ dependency within an unrelated question. As such exact matches would be hard to
¹⁷¹⁷ remove without the use of regular expressions, and due to the low likelihood (6%)
¹⁷¹⁸ of their appearance, we did not perform any followup automatic filtering.

⁵Watson Cognitive Services

⁶We make this available for future use at: <http://bit.ly/2NqBB70>

1719 5.4.2.4 Classification

1720 Our 1,425 posts were then split into 4 additional random samples (in addition to the
1721 random sample of 50 above). 475 posts were classified by the first author and three
1722 other research assistants, software engineers with at least 2 years industry experience,
1723 assisted to classify the remaining 900. This left a total of 1,375 classifications
1724 made by four people plus an additional 450 classifications made from reliability
1725 analysis, in which the remaining 50 posts were classified nine times (as detailed in
1726 Section 5.4.3.1). Thus, a total of 1,825 classifications were made from the original
1727 1,425 posts extracted.

1728 Whilst we could have chosen to employ topic modelling, these are too descriptive
1729 in nature (as discussed in Section 5.3). Moreover, we wanted to see if prior
1730 taxonomies can be applied to IWSs (as opposed to creating a new one) and compare
1731 if their distributions are similar. Therefore, we applied the two existing taxonomies
1732 described in Section 5.3 to each post; (i) a documentation-specific taxonomy that
1733 addresses issues directly resulting from documentation, and (ii) a generalised taxonomy
1734 that covers a broad range of SO issues in a well-defined SE area (specifically
1735 mobile app development). Aghajani et al.'s documentation-specific taxonomy (Tax-
1736 onomy A) is multi-layered consisting of four dimensions and 16 sub-categories [3].
1737 Similarly, Beyer's SO generalised post classification taxonomy (Taxonomy B) con-
1738 sists of seven dimensions [34]. We code each dimension with a number, X, and each
1739 sub-category with a letter y: (Xy). We describe both taxonomies in detail within
1740 Table 5.1. Where a post was included in our results but not applicable to IWSs (see
1741 Section 5.4.2.3) or not applicable to a taxonomy dimension/category, then the post
1742 was flagged for removal in further analysis. Table 5.1 presents *our understanding* of
1743 the respective taxonomies; our intent is not to methodologically replicate Aghajani
1744 et al. or Beyer et al.'s studies in the IWS domain, rather to acknowledge related
1745 work in the area of SO classification and reduce the need to synthesise a new taxonomy.
1746 We baseline all coding against *our interpretation only*. Our classifications are
1747 therefore independent of the previous authors' findings.

1748 5.4.3 Data Analysis**1749 5.4.3.1 Reliability of Classification**

1750 To measure consistency of the categories assigned by each rater to each post, we
1751 utilised both intra- and inter-rater reliability [206]. As verbatim descriptions from
1752 dimensions and sub-categories were considered quite lengthy from their original
1753 sources, all raters met to agree on a shared interpretation of the descriptions, which
1754 were then paraphrased as discussed in the previous subsection and tabulated in
1755 Table 5.1. To perform statistical calculations of reliability, each category was as-
1756 signed a nominal value and a random sample of 50 posts were extracted. Two-phase
1757 reliability analysis followed.

1758 Firstly, intra-rater agreement by the first author was conducted twice on 28 June
1759 2019 and 9 August 2019. Secondly, inter-rater agreement was conducted with the
1760 remaining four co-authors in addition to three research assistants within our research

Table 5.1: Descriptions of dimensions (■) and sub-categories (↔) from both taxonomies used.

A Documentation-specific classification (Aghajani et al. [3])	
A-1	■ Information Content (What)
A-1a	↔ <i>Correctness</i>
A-1b	↔ <i>Completeness</i>
A-1c	↔ <i>Up-to-dateness</i>
A-2	■ Information Content (How)
A-2a	↔ <i>Maintainability</i>
A-2b	↔ <i>Readability</i>
A-2c	↔ <i>Usability</i>
A-2d	↔ <i>Usefulness</i>
A-3	■ Process-Related
A-3a	↔ <i>Internationalisation</i>
A-3b	↔ <i>Contribution-Related</i>
A-3c	↔ <i>Configuration-Related</i>
A-3d	↔ <i>Implementation-Related</i>
A-3e	↔ <i>Traceability</i>
A-4	■ Tool-Related
A-4a	↔ <i>Tooling Bugs</i>
A-3b	↔ <i>Tooling Discrepancy</i>
A-3c	↔ <i>Tooling Help Required</i>
A-3d	↔ <i>Tooling Migration</i>
B Generalised classification (Beyer et al. [34])	
B-1	■ API usage
B-2	■ Discrepancy
B-3	■ Errors
B-4	■ Review
B-5	■ Conceptual
B-6	■ API change
B-7	■ Learning

1761 group in mid-August 2019. Thus, the 50 posts were classified an additional nine
1762 times, resulting in 450 classifications for reliability analysis. We include these
1763 classifications in our overall analysis.

1764 At first, we followed methods of reliability analysis similar to previous SO
1765 studies (e.g., [302]) using the percentage agreement metric that divides the number
1766 of agreed categories assigned per post by the total number of raters [206]. However,
1767 percentage agreement is generally rejected as an inadequate measure of reliability
1768 analysis [70, 128, 176] in statistical communities. As we used more than 2 coders
1769 and our reliability analysis was conducted under the same random sample of 50
1770 posts, we applied *Light's Kappa* [186] to our ratings, which indicates an overall
1771 index of agreement. This was done using the `irr` computational R package [112]
1772 as suggested in [128].

1773 **5.4.3.2 Distribution Analysis**

1774 In order to compare the distribution of categories from our study with previous studies
1775 we carried out a χ^2 test. We selected a χ^2 test as the following assumptions [285]
1776 are satisfied: (i) the data is categorical, (ii) all counts are greater than 5, and (iii)
1777 we can assume simple random sampling. The null hypothesis describes the case
1778 where each population has the same proportion of observations and the alternative
1779 hypothesis is where at least one of the null hypothesis statements is false. We chose
1780 a significance value, α , of 0.05 following a standard rule of thumb. As to the best
1781 of our knowledge this is the first statistical comparison using Taxonomy A and B on
1782 SO posts. To report the effect size we selected Cramer's Phi, ϕ_c which is well suited
1783 for use on nominal data [285].

1784 **5.5 Findings**

1785 We present our findings from classifying a total of 1,825 SO posts aimed at answering
1786 RQs 1 and 2. 450 posts were classified using Taxonomies A and B for reliability
1787 analysis as described in Section 5.4.3.1 and the remaining 1,375 posts were classified
1788 as per Section 5.4.2.4. A summary of our classification using Taxonomies A and B
1789 is shown in Figure 5.3.

1790 **5.5.1 Post classification and reliability analysis**

1791 When undertaking the classification, we found that 238 issues (13.04%) did not
1792 relate to IWSs directly. For example, library dependencies were still included in
1793 a number of results (see Section 5.4.2.3), and we found there to be many posts
1794 discussing Android's Mobile Vision API as Google (Cloud) Vision. These issues
1795 were flagged and ignored for further analysis (see Section 5.4.2.4).

1796 For our reliability analysis, we classified a total of 450 posts of which 70 posts
1797 were flagged as irrelevant. Landis and Koch [182] provide guidelines to interpret
1798 kappa reliability statistics, where $0.00 \leq \kappa \leq 0.20$ indicates *slight* agreement and
1799 $0.21 \leq \kappa \leq 0.40$ indicates *fair* agreement. Despite all raters meeting to agree
1800 on a shared interpretation of the taxonomies (see Section 5.4.3.1) our inter-rater

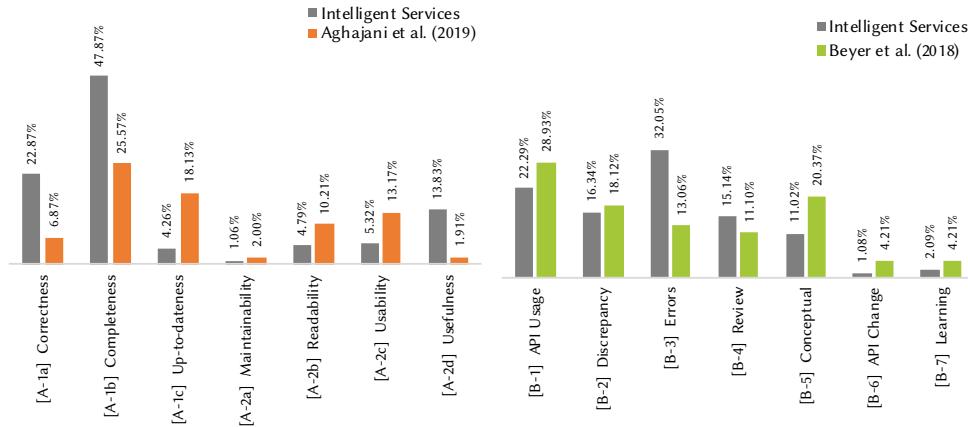


Figure 5.3: *Left:* Documentation-specific classification taxonomy results highlights a mostly similar distribution to that of Aghajani et al.’s findings [3]. *Right:* Generalised classification taxonomy results highlight differences from more mature fields (i.e., Android APIs in Beyer et al. [34]) to less mature fields (i.e., IWSs).

measures aligned *slightly* (0.148) for Taxonomy A and *fairly* (0.295) for Taxonomy B. We report further in Section 5.7.

5.5.2 Developer Frustrations

We found Beyer et al.’s high-level abstraction taxonomy (Taxonomy B) was able to classify 86.52% of posts. 10.30% posts were assigned exclusively under Aghajani et al.’s documentation-specific taxonomy (Taxonomy A). We found that developers do not generally ask questions exclusive to documentation, and typically either pair documentation-related issues to their own code or context. The following two subsections further explain results from both Taxonomy A and B’s perspective.

5.5.2.1 Results from Aghajani et al.’s taxonomy

Results for Aghajani et al.’s low-level documentation taxonomy (Taxonomy A), indicates that most discussion on SO does not directly relate to documentation about an IWS. We did not find any process-related (A-3) or tool-related (A-4) questions as, understandably, the developers who write the documentation of the IWSs would not be posting questions of such nature on SO. One can *infer* documentation-related issues from posts (i.e., parts of the documentation *lacking* that may cause the issue posted). However, there are few questions that *directly* relate to documentation of IWSs.

Few developers question or ask questions directly about the API documentation, but some (47.87%) posts ask for additional information to understand the API (**completeness (A-1b)**), for example: “*Is there a full list of potential labels that Google’s Vision API will return?*” [396]; “*There seems to be very little to no documentation for AWS iOS text recognition inside an image*” [394].

22.87% of posts question the **accuracy (A-1a)** of certain parts of the cloud docu-

1825 mentation, especially in relation to incorrect quotas and limitations: “*Are the Cloud*
1826 *Vision API limits in documentation correct?*” [407], “*According to the Google Vision*
1827 *documentation, the maximum number of image files per request is 16. Elsewhere,*
1828 *however, I’m finding that the maximum number of requests per minute is as high as*
1829 *1800.*” [392].

1830 There are also many references (23.94%) addressing the confusing nature of
1831 some documentation, indicating that the **readability, usability and usefulness of**
1832 **the documentation (A-2b, A-2c and A-2d)** could be improved. For example, “*Am*
1833 *I encoding it correctly? The docs are quite vague.*” [390], “*The aws docs for this*
1834 *are really confusing.*” [419].

1835

5.5.2.2 Results from Beyer et al.’s taxonomy

1836 We found that a majority (32.05%) of posts are primarily **error-related questions**
1837 **(B-3)**, including a dump of the stack trace or exception message from the service’s
1838 programming-language SDK (usually Java, Python or C#) that relates to a specific
1839 error. For example: “*I can’t fix an error that’s causing us to fall behind.*” [416]; “*I’m*
1840 *using the Java Google Vision API to run through a batch of images... I’m now getting*
1841 *a channel closed and ClosedChannelException error on the request.*” [410].

1842 **API usage questions (B-1)** were the second highest category at 22.29% of
1843 posts. Reading the questions revealed that many developers present an insufficient
1844 understanding of the behaviour, functional capability and limitation of these services
1845 and the need for further data processing. For example, while Azure provides an
1846 image captioning service, this is not universal to all CVSSs: “*In Amazon Rekognition*
1847 *for image processing how do I get the caption for an image?*” [401]. Similarly,
1848 OCR-related and label-related questions often indicate interest in cross-language
1849 translation, where a separate translation service would be required: “*Can Google*
1850 *Cloud Vision generate labels in Spanish via its API?*” [415]; “[*How can I] specify*
1851 *language for response in Google Cloud Vision API*” [402]; “[*When I request a text*
1852 *detection of an image, it gives only English Alphabet characters (characters without*
1853 *accents) which is not enough for me. How can I get the UTF-32 characters?*” [397].

1854 It was commonplace to see questions that demonstrate a lack of depth in under-
1855 standing and appreciating how these services work, instead posting simple debugging
1856 questions. For instance, in the 11.02% of **conceptual-related questions (B-5)** that
1857 we categorised, we noticed causal links to a misunderstanding (or lack of awareness)
1858 of the vocabulary used within CV. For example: “*The problem is that I need to know*
1859 *not only what is on the image but also the position of that object. Some of those*
1860 *APIs have such feature but only for face detection.*” [408]; “[*I want to know if the new*
1861 *image has a face similar to the original image.... [the service] can identify faces,*
1862 *but can I use it to get similar faces to the identified face in other images?*” [400]. It
1863 is evident that some application developers are not aware of conceptual differences
1864 in CV such as object/face *detection* versus *localisation* versus *recognition*.

1865 In the 16.34% of **discrepancy-related questions (B-2)**, we see further unaware-
1866 ness from developers in how the underlying systems work. In OCR-related questions,
1867 developers do not understand the pre-processing steps required before an OCR is

1868 performed. In instances where text is separated into multiple columns, for example,
 1869 text is read top-down rather than left-to-right and segmentation would be required
 1870 to achieve the expected results. For example, “*it appears that the API is using some*
 1871 *kind of logic that makes it scan top to bottom on the left side and moving to right*
 1872 *side and doing a top to bottom scan.*” [414]; “*this method returns scanned text in*
 1873 *wrong sequence... please tell me how to get text in proper sequence.*” [420].

1874 A number of **review-related questions (B-4)** (15.14%) seem to provide some
 1875 further depth in understanding the context to which these systems work, where train-
 1876 ing data (or training stages) are needed to understand how inferences are made: “*How*
 1877 *can we find an exhaustive list (or graph) of all logos which are effectively recognized*
 1878 *using Google Vision logo detection feature?*” [418]; “*when object banana is detected*
 1879 *with accuracy greater than certain value, then next action will be dispatched... how*
 1880 *can I confidently define and validate the threshold value for each item?*” [404].

1881 **API change (B-6)** was shown in 1.08% of posts, with evolution of the services
 1882 occurring (e.g., due to new training data) but not necessarily documented “*Recently*
 1883 *something about the Google Vision API changed... Suddenly, the API started to*
 1884 *respond differently to my requests. I sent the same picture to the API today, and I*
 1885 *got a different response (from the past).*” [417].

1886 5.5.3 Statistical Distribution Analysis

1887 We obtained the following results $\chi^2 = 131.86$, $\alpha = 0.05$, $p \text{ value} = 2.2 \times 10^{-16}$ and
 1888 $\phi_c = 0.362$ from our distribution analysis with Taxonomy A to compare our study
 1889 with that of Aghajani et al. [3]. Comparing our study to Beyer et al. [34] produced the
 1890 following results $\chi^2 = 145.58$, $\alpha = 0.05$, $p \text{ value} = 2.2 \times 10^{-16}$ and $\phi_c = 0.252$.
 1891 These results show that we are able to reject the null hypothesis that the distribution
 1892 of posts using each taxonomy was the same as the comparison study. While there are
 1893 limited guidelines for interpreting ϕ_c when there is no prior information for effect
 1894 size [297], Sun et al. suggests the following: $0.07 \leq \phi_c \leq 0.20$ indicates a *small*
 1895 effect, $0.21 \leq \phi_c \leq 0.35$ indicates a *medium* effect, and $0.35 > \phi_c$ indicates a *large*
 1896 effect. Based on this criteria we obtained a *large* effect size for the documentation-
 1897 specific classification (Taxonomy A) and a *medium* effect size for the generalised
 1898 classification (Taxonomy B).

1899 5.6 Discussion

1900 5.6.1 Answers to Research Questions

1901 5.6.1.1 How do developers mis-comprehend IWSs as presented within SO 1902 pain-points? (RQ1)

1903 Upon meeting to discuss the discrepancies between our categorisation of IWS usage
 1904 SO posts, we found that our interpretations of the *posts themselves* were largely sub-
 1905 jective. For example, many posts presented multi-faceted dimensions for Taxonomy
 1906 B; Beyer et al. [34] argue that a post can have more than one question category and

1907 therefore multi-label classification is appropriate at times. We highlight this further
1908 in the threats to validity (Section 5.7).

1909 We have to define the context of IWSs to address RQ1. We use the concept
1910 of a “technical domain” [20] to define this context. A technical domain captures
1911 the domain-specific concerns that influence the non-functional requirements of a
1912 system [20]. In the context of IWSs, the technical domain includes exploration, data
1913 engineering, distributed infrastructure, training data, and model characteristics as
1914 first class citizens [20]. We would then expect to see posts on SO related to these
1915 core concerns.

1916 In Figure 5.3, for the documentation-specific classification, the majority of posts
1917 were classified as **Completeness (A1-b)** related (47.87%). An interpretation for this
1918 is that the documentation does not adequately cover the technical domain concerns.
1919 Comments by developers such as “*I'm searching for a list of all the possible image*
1920 *labels that the Google Cloud Vision API can return?*” [395] indicates the documen-
1921 *tation does not adequately describe the training data for the API—developers do*
1922 *not know the required usage assumptions. Another quote from a developer, “Can*
1923 *Google Cloud Vision generate labels in Spanish via its API? ... [Does the API]*
1924 *allow to select which language to return the labels in?”* [415] points to a lack of
1925 *details relating to the characteristics of the models used by the API. It would seem*
1926 *that developers are unaware of aspects of the technical domain concerns.*

1927 The next most frequent category is **Correctness (A-1a)** with 22.87% of posts. In
1928 the context of the technical domain there are many limits that developers need to be
1929 aware of: range and increments of a model score [76]; required data pre-processing
1930 steps for optimal performance; and features provided by the models (as explained in
1931 Section 5.5.2.2). Considering the relation between technical concerns and software
1932 quality, developers are right to question providers on correctness; “*Are the Cloud*
1933 *Vision API limits in documentation correct?*” [407].

1934 5.6.1.2 *Are the distribution of issues similar to prior studies? (RQ2)*

1935 Visual inspection of Figure 5.3 shows that the distributions for the documentation-
1936 specific classification and the generalised classification are different (compared to
1937 prior studies). As a sanity check we conducted a χ^2 test and calculated the effect
1938 size ϕ_c . We were able to reject the null hypothesis for both classification schemes,
1939 that the distribution of issues were the same as the previous studies (see Section 5.5).
1940 We now discuss the most prominent differences between our study and the previous
1941 studies.

1942 In the context of IWS SO posts, Taxonomy B suggests that Errors (B-3) are
1943 discussed most amongst developers. These results are in contrast to similar studies
1944 made in more *mature* API domains, such as Mobile Development [21, 22, 33, 34, 265]
1945 and Web Development [308]. Here, API Usage (B-1) is much more frequently
1946 discussed, followed by Conceptual (B-5), Discrepancy (B-2) and Errors (B-3). We
1947 argue in the following section that an improved developer understanding can be
1948 achieved by educating them about the IWS lifecycle and the ‘whole’ system that
1949 wraps such services.

1950 In the Android study API usage questions (B-1) were the highest category
1951 (28.93% compared to 22.29% in our study). As stated in the analysis of the Error
1952 questions this discrepancy could be due to the maturity of the domain. However,
1953 another explanation could be the scope of the two individual studies. Beyer et al. [34]
1954 used a broad search strategy consisting of posts tagged Android. This search term
1955 fetches issues related to the entire Android platform which is significantly larger than
1956 searching for CV APIs using 229 search terms. As a consequence of more posts
1957 and more APIs there would be use cases resulting in additional posts related to API
1958 Usage (B-1).

1959 Applying existing SO taxonomies allowed us to better understand the distribution
1960 of the issues across different domains. In particular, the issues raised around IWSs
1961 appear to be primarily due to poor documentation, or insufficient explanation around
1962 errors and limitations. Hence, many of the concerns could be addressed by adding
1963 more details to the end-point descriptions, and by providing additional information
1964 around how these services are designed to work.

1965 5.6.2 The Developer’s Learning Approach

1966 In this subsection, we offer an explanation as to why developers are complaining
1967 about certain things when trying to use IWSs on SO (RQ1), as characterised through
1968 the use of prior SO classification frameworks (RQ2). This is described through
1969 the theoretical lenses of two learning taxonomies: Bloom’s context complexity and
1970 intellectual ability taxonomy, and the Structure of the Observed Learning Outcome
1971 (SOLO) taxonomy (i.e., the nature by which developer’s learn). We argue that the
1972 issues with using IWSs relating to the lower-levels of these learning taxonomies
1973 are easily solvable by slight fixes and improvements to the documentation of these
1974 services. However, the higher dimensions of these taxonomies demand far more
1975 rigorous mitigation strategies than documentation alone (potentially more structured
1976 education). Thus, many of the questions posted are from developers who are *learning*
1977 to *understand* the domain of IWSs and AI, and (hence) both SOLO and Bloom’s
1978 taxonomies are applicable for this discussion—as described below within the context
1979 of our domain—as pedagogical aides.

1980 5.6.2.1 Bloom’s Taxonomy

1981 The cognitive domain under Bloom’s taxonomy [37] consists of six objectives.
1982 Within the context of IWSs, developers are likely to ask questions due to causal
1983 links that exist in the following layers of Bloom’s taxonomy: (i) *knowledge*, where
1984 the developer does not remember or know of the basic concepts of CV and AI
1985 (in essence, they may think that AI is as smart as a human); (ii) *comprehension*,
1986 where the developer does not understand how to interpret basic concepts, or they
1987 are mis-understanding how they are used in context; (iii) *application*, where the
1988 developer is struggling to apply existing concepts within the context of their own
1989 situation; (iv) *analysis*, where the developer is unable to analyse the results from IWSs
1990 (i.e., understand response objects); (v) *evaluation*, where the developer is unable to
1991 evaluate issues and make use of best-practices when using IWSs; and (vi) *synthesise*,

¹⁹⁹² where the developer is posing creative questions to ask if new concepts are possible
¹⁹⁹³ with CVSSs.

¹⁹⁹⁴ 5.6.2.2 SOLO Taxonomy

¹⁹⁹⁵ The SOLO taxonomy [35] consists of five levels of understanding. The causal links
¹⁹⁹⁶ behind the SO questions we have found relate to the following layers of the SOLO
¹⁹⁹⁷ taxonomy: (i) *pre-structural*, where the developer has a question indicating incom-
¹⁹⁹⁸ petence or has little understanding of CV; (ii) *uni-structural*, where the developer
¹⁹⁹⁹ is struggling with one key aspect (i.e., a simple question about CV); (iii) *multi-
²⁰⁰⁰ structural*, where the developer is questioning multiple concepts (independently)
²⁰⁰¹ to understand how to build their system (e.g., system integration with the IWS);
²⁰⁰² (iv) *relational*, where the developer is comparing and contrasting the best ways to
²⁰⁰³ achieve something with IWSs; and (v) *extended abstract*, where the developer poses
²⁰⁰⁴ a question theorising, formulating or postulating a new concept within IWSs.

Table 5.2: Example Alignments of SO posts to Bloom's and the SOLO taxonomy.

Issue Quote	Bloom	SOLO
<i>"I'm using Microsoft Face API for a small project and I was trying to detect a face inside a .jpg file in the local system (say, stored in a directory D:\Image\abc.jpg)... but it does not work."</i> [411]	Knowledge	Pre-Structural
<i>"The problem is that the response JSON is rather big and confusing. It says a lot about the picture but doesn't say what the whole picture is of (food or something like that)." [391]</i>	Comprehension	Uni-Structural
<i>"The bounding box around individual characters is sometimes accurate and sometimes not, often within the same image. Is this a normal side-effect of a probabilistic nature of the vision algorithm, a bug in the Vision API, or of course an issue with how I'm interpreting the response?"</i> [398]	Comprehension	Multi-Structural
<i>"I'm working on image processing. SO far Google Cloud Vision and Clarifai are the best API's to detect objects from images and videos, but both API's doesn't support object detection from 360 degree images and videos. Is there any solution for this problem?"</i> [405]	Application	Uni-Structural
<i>"Before I train Watson, I can delete pictures that may throw things off. Should I delete pictures of: Multiple dogs, A dog with another animal, A dog with a person, A partially obscured dog, A dog wearing glasses, Also, would dogs on a white background make for better training samples? Watson also takes negative examples. Would cats and other small animals be good negative examples?"</i> [403]	Analysis	Relational

²⁰⁰⁵ 5.6.2.3 Aligning SO taxonomies to Bloom's and SOLO taxonomies

²⁰⁰⁶ To understand our findings with the lenses of pedagogical aids, we aligned Tax-
²⁰⁰⁷ onomies A and B to Bloom's and the SOLO taxonomies for a random sample of 50
²⁰⁰⁸ issues described in Section 5.4.3.1. To do this, we reviewed all 50 of these SO posted
²⁰⁰⁹ questions and applied both the Bloom and SOLO taxonomies. The primary author
²⁰¹⁰ assigned each of the 50 questions a level within the Bloom and SOLO taxonomies,

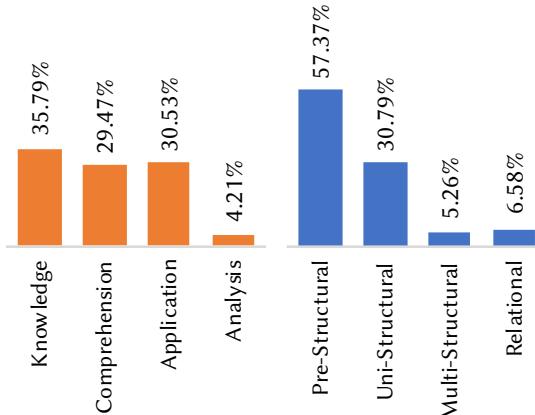


Figure 5.4: Alignment of Bloom (Orange) and SOLO (Blue) taxonomies against Taxonomy A and B dimensions against all 213 classifications made in the random sample of 50 posts.

removed out noise (i.e., false positive posts of no relevance to IWSs) and unassigned dimensions from reliability agreement, and then compared the relevant dimensions of Taxonomy A and B dimensions (not sub-categories). The comparison of alignments of posts to the five SOLO dimensions and six Bloom dimensions are shown in Figure 5.4. We acknowledge that this is only an approximation of the current state of the developer’s understanding of IWSs. This early model will require further studies to perform a more thorough analysis, but we offer this interpretation for early discussion.

As shown in Figure 5.4, the bulk of the posts fall in the lower constructs of Bloom’s and the SOLO taxonomy. This indicates that modification to certain documentation aspects can address many of these issues. For example, many issues can be ratified with better descriptions of response data and error messages: “*I was exploring google vision and in the specific function ‘detectCrops’, gives me the crop hints. what does this means exactly?*” [406]; “*I am making a very simple API call to the Google Vision API, but all the time it’s giving me error that ‘google.oauth2’ module not found.*” [421]

However, and more importantly, the higher-construct questions ranging from the middle of the third dimensions on are not as easily solvable through improved documentation (i.e., apply and multi-structural) which leaves 34.74% (Bloom’s) and 11.84% (SOLO) unaccounted for, resolvable only through improved education practices.

5.6.3 Implications

5.6.3.1 For Researchers

Investigate the evolution of post classification Analysing how the distribution of the reported issues changes over time would be an important study. This study could answer questions such as ‘*Does the evolution of IWSs follow the same pattern as previous software engineering trends such as mobile app or web development?*’ As

2038 with any new emerging field, it is key to analyse how developers perceive such issues
2039 over time. For instance, early issues with web or mobile app development matured
2040 as their respective domain matured, and we would expect similar results to occur
2041 in the IWSs space. Future researchers could plan for a longitudinal study, such as
2042 a long-term survey with developers to gather their insights in this evolving domain,
2043 reviewing case studies of projects that use intelligent web services from now into
2044 the future, or re-mining SO at a later date and comparing the results to this study.
2045 This will help assess evolving trends and characteristics, and determine how and if
2046 the nature of the developer's experience with IWSs (and AI in general) changes with
2047 time.

2048 **Investigate the impact of technical challenges on API usage** As discussed above,
2049 IWSs have characteristics that may influence API usage patterns and should be
2050 investigated as a further avenue of research. Further mining of open source software
2051 repositories that make use of IWSs could be assessed, thereby investigating if API
2052 patterns evolve with the rise of AI-based applications.

2053 *5.6.3.2 For Educators*

2054 **Education on high-level aspects of IWSs** As demonstrated in our analysis of their
2055 SO posts, many developers appear to be unaware of the higher-level concepts that
2056 exist within the AI and ML realm. This includes the need to pre- and post-process
2057 data, the data dependency and instability that exists in these services, and the specific
2058 algorithms that empower the underlying intelligence and hence their limitations and
2059 characteristics. However, most developers don't seem to complain about these factors
2060 due to the lack of documentation (i.e., via Taxonomy A). Rather, they are unaware
2061 that such information should be documentation and instead ask generalised and open
2062 questions (i.e., via Taxonomy B). Thus, documentation improvements alone may not
2063 be enough to solve these issues. This results in uncertainty during the preparation
2064 and operation (usage) of such services. Such high-level conceptual information is
2065 currently largely missing in developer documentation for IWSs. Furthermore, many
2066 of the background ML and AI algorithm information needed to understand and use
2067 intelligent systems in context are built within data science (not SE) communities.
2068 A possible road-map to mitigate this issue would be the development of a software
2069 engineer's 'crash-course' in ML and AI. The aim of such a course would encourage
2070 software engineers to develop an appreciation of the nuances and the inherent risks
2071 and implications that comes with using IWSs. This could be taught at an undergrad-
2072 uate level to prepare the next generation of developers of a 'programming 2.0' era.
2073 However, the key aspects and implications that are presented with AI would need
2074 to be well-understood before such a course is developed, and determining the best
2075 strategy to curate the content to developers would be best left to the SE education
2076 domain. Further investigation in applying educational taxonomies in the area (such
2077 as our attempts to interpret our findings using Bloom's and the SOLO taxonomies)
2078 would need to be thoroughly explored beforehand.

2079 *5.6.3.3 For Software Engineers*

2080 **Better understanding of intelligent API contextual usage** Our results show that
2081 developers are still learning to use these APIs. We applied two learning perspectives
2082 to interpret our results. In applying the two pedagogical taxonomies to our findings,
2083 we see that most issues seem to fall into the pre-structural and knowledge-based
2084 categories; little is asked of higher level concepts and a majority of issues do not
2085 offer complex analysis from developers. This suggests that developers are struggling
2086 as they are unaware of the vocabulary needed to actually use such APIs, further
2087 reinforcing the need for API providers to write overview documentation (as noted in
2088 prior work [75]) and not just simple endpoint documentation. This said, improved
2089 documentation isn't always enough—as suggested by our discussion in Section 5.6.2,
2090 software engineers should explore further education to attain a greater appreciation
2091 of the nuances of ML when attempting to use these services.

2092 *5.6.3.4 For Intelligent Service Providers*

2093 **Clarify use cases for IWSs** Inspecting SO posts revealed that there is a level of
2094 confusion around the capabilities of different IWSs. This needs to be clarified in
2095 associated API documentation. The complication with this comes with targeting
2096 the documentation such that software developers (who are untrained in the nuances
2097 of AI and ML as per Section 5.6.3.2) can digest it and apply it in-context to
2098 application development.

2099 **Technical domain matters** More needs to be provided than a simple endpoint
2100 description as conventional APIs offer by describing the whole framework by which
2101 the endpoint sits, giving further context. This said, compared to traditional APIs,
2102 we find that developers complain less about the documentation and more about
2103 shallower issues. All expected pre-processing and post-processing needs to be
2104 clearly explained. A possible mitigation to this could be an interactive tutorial that
2105 helps developers fully understand the technical domain using a hands-on approach.
2106 For example, websites offer interactive Git tutorials⁷ to help developers understand
2107 and explore the technical domain matters under version control in their own pace.

2108 **Clarify limitations** API developers need to add clear limitations of the existing
2109 APIs. Limitations include list of objects that can be returned from an endpoint. We
2110 found that the cognitive anchors of how existing, conventional API documentation
2111 is written has become ‘ported’ to the CV realm, however a lot more overview
2112 documentation than what is given at present (i.e., better descriptions of errors,
2113 improved context of how these systems work in etc.) needs to be given. Such
2114 documentation could be provided using interactive tutorials.

⁷For example, <https://learngitbranching.js.org>.

2115 5.7 Threats to Validity**2116 5.7.1 Internal Validity**

2117 As detailed in Section 5.4.3.1, Taxonomies A and B present slight and fair agreement,
2118 respectively, when inter-rater reliability was applied. The nature of our disagree-
2119 ments largely fell due to the subjectivity in applying either taxonomies to posts.
2120 Despite all coders agreeing to the shared interpretation of both taxonomies, both
2121 taxonomies are subjective in their application, which was not reported by either
2122 Aghajani et al. or Beyer et al.. In many cases, multi-label classification seemed ap-
2123 propriate, however both taxonomies use single-label mapping which we find results
2124 in too much subjectivity. This subjectivity, therefore, ultimately adversely affects
2125 inter-rater reliability (IRR) analysis. Thus, a future mitigation strategy for similar
2126 work should explore multi-label classification to avoid this issue; Beyer et al., for
2127 example, plan for multi-label classification as future work. However, these studies
2128 would need to consider the statistical challenges in calculating multi-rater, multi-
2129 label IRR for thorough reliability analysis in addressing subjectivity. The selection
2130 of SO posts used for our labelling, chiefly in the subjectivity of our classifications, is
2131 of concern. We mitigate this by an extensive review process assessing the reliability
2132 of our results as per Section 5.4.3.1. The classification of our posts into the SOLO
2133 and Bloom’s taxonomies was performed by the primary author only, and therefore
2134 no inter-rater reliability statistics were performed. However, we used these peda-
2135 gogy related taxonomies as a lens to gain an additional perspective to interpret our
2136 results. Future studies should attempt a more rigorous analysis of SO posts using
2137 Bloom’s and SOLO taxonomies. We only aligned posts to one category for each
2138 taxonomy and did not align these using multi-label classification. This brings more
2139 complexity to the analysis, and our attempts to repeat prior studies’ methodologies
2140 (see Section 5.3). Multi-label classification for IWSs SO posts is an avenue for future
2141 research.

2142 5.7.2 External Validity

2143 While every effort was made to select posts from SO relevant to CVSs, there are
2144 some cases where we may have missed some posts. This is especially due to the
2145 case where some developers mis-reference certain IWSs under different names (see
2146 Section 5.4.2.1).

2147 Our SOLO and Bloom’s taxonomy analysis has only been investigated through
2148 the lenses of IWSs, and not in terms of conventional APIs (e.g., Andriod APIs).
2149 Therefore, we are not fully certain how these results found would compare to other
2150 types of APIs. Two *existing* SO classification taxonomies were used rather than
2151 developing our own. We wanted to see if previous SO taxonomies could be applied
2152 to IWSs before developing a new, specific taxonomy, and these taxonomies were
2153 applied based on our interpretation (see Section 5.4.2.4) and may not necessarily
2154 reflect the interpretation of the original authors. Moreover, automated techniques
2155 such as topic modelling were not utilised as we found these produce descriptive
2156 classifications only (see Section 5.3). Hence, manual analysis was performed by

humans to ensure categories could be aligned back to causal factors. Only English-speaking IWSs were selected; the applicability of our analysis to other, non-English speaking services may affect results. Use of CV in this study is an illustrative example to focus on one area of the IWSs spectrum. While our narrow scope helps us obtain more concrete findings, we suggest that wider issues exist in other IWS domains may affect the generalisability of this study, and suggest future work be explored in this space.

5.7.3 Construct Validity

Some questions extracted from SO produced false positives, as mentioned in Sections 5.4.2.1 and 5.4.2.3 and Section 5.5. However, all non-relevant posts were marked as noise for our study, and thus did not affect our findings. Moreover, SO is known to have issues where developers simply ask basic questions without looking at the actual documentation where the answer exists. Such questions, although down-voted, were still included in our data-set analysis, but as these were SO few, it does not have a substantial impact on categorised posts.

5.8 Conclusions

CVSs offer powerful capabilities that can be added into the developer’s toolkit via simple RESTful APIs. However, certain technical nuances of CV become abstracted away. We note that this abstraction comes at the expense of a full appreciation of the technical domain, context and proper usage of these systems. We applied two recent existing SO classification taxonomies (from 2018 and 2019) to see if existing taxonomies are able to fully categorise the types of complaints developers have. IWSs have a diverging distribution of the types of issues developers ask when compared to more mature domains (i.e., mobile app development and web development). Developers are more likely to complain about shallower, simple debugging issues without a distinct understanding of the AI algorithms that actually empower the APIs they use. Moreover, developers are more likely to complain about the completeness and correctness of existing IWS documentation, thereby suggesting that the documentation approach for these services should be reconsidered. Greater attention to education in the use of AI-powered APIs and their limitations is needed, and our discussion offered in Section 5.6.2 motivates future work in resolving these issues in the SE education space.

CHAPTER 6

2189

2190

2191

Ranking Computer Vision Service Issues using Emotion[†]

2192

2193

Abstract Software developers are increasingly using intelligent web services to implement ‘intelligent’ features. Studies show that incorporating artificial intelligence (AI) into an application increases technical debt, creates data dependencies, and introduces uncertainty due to non-deterministic behaviour. However, we know very little about the emotional state of software developers who deal with such issues. In this paper, we do a landscape analysis of emotion found in 1,425 Stack Overflow (SO) posts about computer vision services. We investigate the application of an existing emotion classifier EmoTxt and manually verify our results. We found that the emotion profile varies for different question categories and that a new emotion schema is required to better represent the emotion present in SO questions. We propose an initial version of a new emotion classification scheme and confirm current findings that AI is insufficient for automatic classification of emotion.

2204

6.1 Introduction

2205

Recent advances in artificial intelligence have provided software engineers with new opportunities to incorporate complex machine learning capabilities, such as computer vision, through cloud-based intelligent web services (IWSs). These new set of services, typically offered as API calls are marketed as a way to reduce the complexity involved in integrating AI-components. However, recent work shows that software engineers struggle to use these IWSs [79].

2211

While seeking advice on the issues, software engineers tend to express their emotions (such as frustration or confusion) within the questions. Recognising the value of considering emotions, other researchers have investigated emotions expressed by software developers within communication channels [233] including Stack Overflow (SO) [57, 226]; the broad motivation of these works is to generally understand the

[†]This chapter is originally based on M. K. Curumsing, A. Cummaudo, U. M. Graestch, S. Barnett, and R. Vasa, “Ranking Computer Vision Service Issues using Emotion,” 2020, Unpublished. Terminology has been updated to fit this thesis.

2216 emotional landscape and improve developer productivity [111, 217, 233]. However,
2217 previous works have not directly focused on the nature of emotions expressed in
2218 questions related to IWSs. We also do not know if certain types of questions express
2219 stronger emotions.

2220 The machine-learnt behaviour of these IWSs is typically non-deterministic and,
2221 given the dimensions of data used, their internal inference process is hard to reason
2222 about [76]. Compounding the issue, documentation of these cloud systems does not
2223 explain the limits, nor how they were created (esp. data sets used to train them).
2224 This lack of transparency makes it difficult for even senior developers to properly
2225 reason about these systems, so their prior experience and anchors do not offer
2226 sufficient support [79]. In addition, adding machine learned behaviour to a system
2227 incurs ongoing maintenance concerns [277]. There is a need to better understand
2228 emotions expressed by developers to inform cloud vendors and help them improve
2229 their documentation and error messages provided by their services.

2230 This work builds on top of recent work that explored *what* pain-points developers
2231 face when using IWSs through a general analysis of 1,425 SO posts (questions) [79]
2232 using an existing SO issue classification taxonomy [34]. In this work, we consider
2233 the emotional state expressed within these pain-points, using the same data set of
2234 1,425 SO posts. We identify the emotions in each SO question, and investigate if
2235 the distribution of these emotions is similar across the various types of questions.

2236 In order to classify emotions from SO posts, we use EmoTxt, a recently proposed
2237 toolkit for emotion recognition from text [56, 57, 226]. EmoTxt has been trained
2238 and built on SO posts using the emotion classification model proposed by Shaver
2239 et al. [282]. The category of issue was manually determined in our prior work.

2240 The key findings of our study are:

- 2241 • The distribution of emotions is different across the taxonomy of issues.
- 2242 • A deeper analysis of the results, obtained from the EmoTxt classifier, suggests
2243 that the classification model needs further refinement. Love and joy, the
2244 least expected emotions when discussing API issues, are visible across all
2245 categories.
- 2246 • A different emotion classification scheme is required to better reflect the
2247 emotions within the questions.

2248 In order to promote future research and permit replication, we make our data
2249 set publicly available.¹ The paper structure is as follows: Section 6.2 provides
2250 an overview on prior work surrounding the classification of emotions from text;
2251 Section 6.3 describes our research methodology; Section 6.4 presents the results
2252 from the EmoTxt classifier; Section 6.5 provides a discussion of the results obtained;
2253 Section 6.6 highlights the implications of our study; Section 6.7 outlines the threats
2254 to validity; Section 6.8 presents the concluding remarks.

¹See <http://bit.ly/2RiULgW>.

6.2 Emotion Mining from Text

Several studies have investigated the role of emotions generally in software development [111, 233, 283, 331]. Work in the area of behavioural software engineering established the link between software developer's happiness and productivity [124]. Wrobel [331] investigated the impact that software developers' emotion has on the development process and found that frustration and anger were amongst the emotions that posed the highest risk to developer's productivity.

Recent studies focused on emotion mining from text within communication channels used by software engineers to communicate with their peers [111, 217, 226, 233]. Murgia et al. [217] and Ortu et al. [233] investigated the emotions expressed by developers within an issue tracking system, such as JIRA, by labelling issue comments and sentences written by developers using Parrott's framework. Gachechiladze et al. [111] applied the Shaver framework to detect anger expressed in comments written by developers in JIRA. The Collab team [56, 226] extended the work done by Ortu et al. [233] and developed a gold standard data set collected from SO posts consisting of questions, comments and feedback. This data set was manually annotated using the Shaver's emotion model. The Shaver's model consists of a tree-structured, three level, hierarchical classification of emotions. The top level consists of six basic emotions namely, love, joy, anger, sadness, fear and surprise [282]. The subsequent levels further refines the granularity of the previous level. One of their recent work [226] involved 12 raters to manually annotate 4,800 posts (where each post included the question, answer and comments) from SO. The same question was assigned to three raters to reduce bias and subjectivity. Each coder was requested to indicate the presence/absence of each of the six basic emotions from the Shaver framework. As part of their work they developed an emotion mining toolkit, EmoTxt [56]. The work conducted by the Collab team is most relevant to our study since their focus is on identifying emotion from SO posts and their toolkit is trained on a large data set of SO posts.

6.3 Methodology

As mentioned in our introduction, this paper uses the data set reported in Cummaudo et al.'s ICSE 2020 paper [79]. As this paper is in press, we reproduce a summary of the methodology used in constructing this data set methodology below. For full details, we refer to the original paper. Supplementary materials used for this work are provided for replication.¹

Our research methodology consisted of the following steps: (i) data extraction from SO resulting in 1,425 questions about intelligent computer vision services (CVSs); (ii) question classification using the taxonomy presented by Beyer et al. [34]; (iii) automatic emotion classification using EmoTxt based on Shaver et al.'s emotion taxonomy [282]; and (iv) manual classification of 25 posts to better understand developers emotion. We calculated the inter-rater reliability between EmoTxt and our manually classified questions in two ways: (i) to see the overall agreement between the three raters in applying the Shaver et al. emotions taxonomy, and (ii) to

²²⁹⁷ see the overall agreement with EmoTxt’s classifications. Further details are provided
²²⁹⁸ below.

²²⁹⁹ 6.3.1 Data Set Extraction from SO

²³⁰⁰ 6.3.1.1 Intelligent Service Selection

²³⁰¹ We contextualise this work within popular CVS providers: Google Cloud [362],
²³⁰² AWS [341], Azure [376] and IBM Cloud [372]. We chose these four providers given
²³⁰³ their prominence and ubiquity as cloud service vendors, especially in enterprise
²³⁰⁴ applications [259]. We acknowledge other services beyond the four analysed which
²³⁰⁵ provide similar capabilities [354, 355, 358, 371, 422, 423]. Additionally, only
²³⁰⁶ English-speaking services have been selected, excluding popular CVSs from Asia
²³⁰⁷ (e.g., [352, 353, 370, 388, 389]).

²³⁰⁸ 6.3.1.2 Developing a search query

²³⁰⁹ To understand the various ways developers refer to these services, we needed to find
²³¹⁰ search terms that are commonplace in question titles and bodies that discuss the
²³¹¹ service names. One approach is to use the *Tags* feature in SO. To discover which
²³¹² tags may be relevant, we ran a search² within SO against the various brand names of
²³¹³ these CVSs, reviewed the first three result pages, and recorded each tag assigned per
²³¹⁴ question.³ However, searching using tags alone on SO is ineffective (see [23, 302]).
²³¹⁵ To overcome this limitation, we ran a second query within the Stack Exchange Data
²³¹⁶ Explorer⁴ (SEDE) using these tags, we sampled 100 questions (per service), and
²³¹⁷ noted the permutations in how developers refer to each service⁵. We noted 229
²³¹⁸ permutations.

²³¹⁹ 6.3.1.3 Executing our search query

²³²⁰ Next, we needed to extract questions that make reference to any of these 229 per-
²³²¹ mutations. SEDE has a 50,000 row limit and does not support case-insensitivity,
²³²² however Google’s BigQuery does not. Therefore, we queried Google’s SO dataset
²³²³ on each of the 229 terms that may occur within the title or body of question posts,⁶
²³²⁴ which resulted in 21,226 questions.

²³²⁵ 6.3.1.4 Refining our inclusion/exclusion criteria

²³²⁶ To assess the suitability of these questions, we filtered the 50 most recent posts
²³²⁷ as sorted by their *CreationDate* values. This helped further refine the inclusion
²³²⁸ and exclusion criteria: for example, certain abbreviations in our search terms (e.g.,

²The query was run on January 2019.

³Up to five tags can be assigned per question.

⁴<http://data.stackexchange.com/stackoverflow>

⁵E.g., misspellings, misunderstanding of brand names, hyphenation, UK vs. US English, and varied uses of apostrophes, plurals, and abbreviations.

⁶See <http://bit.ly/2LrN70A>.

Table 6.1: Descriptions of dimensions from our interpretation of Beyer et al.’s SO question type taxonomy.

Dimension	Our Interpretation
API usage	Issue on how to implement something using a specific component provided by the API
Discrepancy	The questioner’s <i>expected behaviour</i> of the API does not reflect the API’s <i>actual behaviour</i>
Errors.....	Issue regarding an error when using the API, and provides an exception and/or stack trace to help understand why it is occurring
Review	The questioner is seeking insight from the developer community on what the best practices are using a specific API or decisions they should make given their specific situation
Conceptual.....	The questioner is trying to ascertain limitations of the API and its behaviour and rectify issues in their conceptual understanding on the background of the API’s functionality
API change.....	Issue regarding changes in the API from a previous version
Learning	The questioner is seeking for learning resources to self-learn further functionality in the API, and unlike discrepancy, there is no specific problem they are seeking a solution for

²³²⁹ ‘GCV’, ‘WCS’⁷) allowed for false positive questions to be included, which were removed. Furthermore, we consolidated all overlapping terms (e.g., ‘Google Vision ²³³⁰ **API**’ was collapsed into ‘Google Vision’) to enhance the query. Additionally, we ²³³¹ reduced our 221 search terms to just 27 search terms by focusing on CVSs *only*⁸ ²³³² which resulted in 1,425 questions. No duplicates were recorded as determined by ²³³³ the unique ID, title and timestamp of each question. ²³³⁴

²³³⁵ 6.3.1.5 Manual filtering

²³³⁶ The next step was to assess the suitability and nature of the 1,425 questions extracted. ²³³⁷ The second author ran a manual check on a random sample of 50 posts, which were ²³³⁸ parsed through a templating engine script⁹ in which the ID, title, body, tags, created ²³³⁹ date, and view, answer and comment counts were rendered for each post. Any match ²³⁴⁰ against the 27 search terms in the title or body of the post were highlighted, in which ²³⁴¹ three false positives were identified as either library imports or stack traces, such ²³⁴² as `aws-java-sdk-rekognition:jar`. In addition, we noted that there were false ²³⁴³ positive hits related to non-CVSs. We flagged posts of such nature as ‘noise’ and ²³⁴⁴ removed them from further classification.

⁷Watson Cognitive Services

⁸Our original data set aimed at extracting posts relevant to *all* IWSs, and not just CVSs. However, 21,226 questions were too many to assess without automated analysis, which was beyond the scope of our work.

⁹We make this available for future use at: <http://bit.ly/2NqBB70>.

2345 6.3.2 Question Type & Emotion Classification

2346 6.3.2.1 Manual classification of question category

2347 We classify our 1,425 posts using Beyer et al.'s taxonomy [34] as it was comprehensive and validated [79]. We split the posts into 4 additional random samples, in
 2348 addition to the random sample of 50 above. 475 posts were classified by the second
 2349 author and three other research assistants¹⁰ classified the remaining 900 (i.e., a total
 2350 of 1,375 classifications). An additional 450 classifications were assigned due to
 2351 reliability analysis, in which the remaining 50 posts were classified nine times by
 2352 various researchers in our group.¹¹

2353 Due to the nature of reliability analysis, multiple classifications (450) existed
 2354 for these 50 posts. Therefore, we applied a 'majority rule' technique to each post
 2355 allowing for a single classification assignment and therefore analysis within our re-
 2356 sults. When there was a majority then we used the majority classification; when
 2357 there was a tie, then we used the classification that was assigned the most out of the
 2358 entire 450 classifications. As an example, 3 raters classified a post as *API Usage*,
 2359 1 rater classified the same post as a *Review* question and 5 raters classified the post
 2360 as *Conceptual*, resulting in the post being classified as a *Conceptual* question. For
 2361 another post, three raters assigned *API Usage*, *Discrepancy* and *Learning* (respec-
 2362 tively), while 3 raters assigned *Review* and 3 raters assigned *Conceptual*. In this
 2363 case, *Review* and *Conceptual* were tied, but was resolved down to *Conceptual* as this
 2364 classification received 147 more votes than *Review* across all classifications made in
 2365 the sample of 50 posts.

2366 However, where a post was extracted from our original 1,425 posts but was either
 2367 a false positive, not applicable to IWSs (see Section 6.3.1.5), or not applicable to
 2368 a taxonomy dimension/category, then the post was flagged for removal in further
 2369 analysis. This was done 180 times, leaving a total of 1,245 posts.

2370 Our interpretation Beyer et al.'s taxonomy is provided in Table 6.1, which
 2371 presents a transcription of *our understanding* of the respective taxonomy. We
 2372 baselined all coding against *our interpretation only*, and thus our classifications
 2373 are therefore independent of Beyer et al.'s findings, since we baseline results via
 2374 Table 6.1's interpretation.

2376 6.3.2.2 Emotion classification using artificial intelligence (AI) techniques

2377 After extracting and classifying all posts, we then piped in the body of each question
 2378 into a script developed to remove all HTML tags, code snippets, blockquotes and
 2379 hyperlinks, as suggested by Novielli et al. [226]. We replicated and extended the
 2380 study conducted by Novielli et al. [226] on our data set derived from 1,425 SO posts,
 2381 consisting of questions only. Our study consisted of three main steps, namely, (1)
 2382 automatic emotion classification using EmoTxt, (2) manual annotation process and,
 2383 (3) comparison of the automatic classification result with the manually annotated
 2384 data set.

¹⁰Software engineers in our research group with at least 2 years industry experience

¹¹Due to space limitations, reliability analysis is omitted and is reported in [79].

2385 6.3.2.3 Emotion classification using EmoTxt

2386 We started with a file containing 1,245 non-noise SO questions, each with an as-
2387 sociated question type as classified using the strategy discussed in Section 6.3.2.1.
2388 We pre-processed this file by extracting the question ID and body text to meet the
2389 format requirements of the EmoTxt classifier [56]. This classifier was used as it
2390 was trained on SO posts as discussed in Section 6.2. We ran the classifier for each
2391 emotion as this was required by EmoTxt model. This resulted in 6 output prediction
2392 files (one file for each emotion: *Love, Joy, Surprise, Sadness, Fear, Anger*). Each
2393 question within these files referenced the question ID and a predicted classification
2394 (YES or NO) of the emotion. We then merged the emotion prediction files into an
2395 aggregate file with question text and Beyer et al.’s taxonomy classifications. This
2396 resulted in 796 emotion classifications. We further analysed the classifications and
2397 generated an additional classification of *No Emotion* for the 622 questions where
2398 EmoTxt predicted NO for all the emotion classification runs.

2399 Of the 796 questions with emotion detected, 143 questions had 2 or more
2400 emotions predicted: 1 question¹² had up to 4 emotions detected (*Surprise, Sadness,*
2401 *Joy and Fear*), 28 questions had up to 3 emotions detected, and the remaining 114
2402 had up to two emotions detected.

2403 6.3.2.4 Manual Annotation Process

2404 In order to evaluate and also better understand the process used by EmoTxt to
2405 classify emotions, we manually annotated a small sample of 25 SO posts, randomly
2406 selected from our data set. Each of these 25 posts were assigned to three raters who
2407 carried out the following three steps: (i) identify the presence of an emotion; (ii)
2408 if an emotion(s) exists, classify the emotion(s) under one of the six basic emotions
2409 proposed by the Shaver framework [282]; (iii) if no emotion is identified, annotate as
2410 neutral. We then collated all rater’s results and calculated Light’s Kappa (L_κ) [186]
2411 to measure the overall agreement *between* raters to measure the similarity in which
2412 independent raters classify emotions to SO posts. As L_κ does not support multi-class
2413 classification (i.e., multiple emotions) per subjects (i.e., per SO post), we binarised
2414 the results each emotion and rater as TRUE or FALSE to indicate presence, calculated
2415 the L_κ per emotion against the three raters, and averaged the result across all emotions
2416 to get an overall strength of agreement.

**2417 6.3.2.5 Comparing EmoTxt results with the results from Manual Classifi-
2418 cation**

2419 The next step involved comparing the ratings of the 25 SO posts that were manually
2420 annotated by the three raters with the results obtained for the same set of 25 SO
2421 posts from the EmoTxt classifier. Similar to Section 6.3.2.4, we used Cohen’s Kappa
2422 (C_κ) [70] to measure the consistency of classifications of EmoTxt’s classifications
2423 versus the manual classifications of each rater. We separated the classifications
2424 per emotion and calculated C_κ for each rater against EmoTxt and averaged these

¹²See <http://stackoverflow.com/q/55464541>.

²⁴²⁵ values for all emotions. After noticing poor results, the three raters involved in
²⁴²⁶ Section 6.3.2.4 were asked to compare and discuss the ratings from the EmoTxt
²⁴²⁷ classifier against the manual ratings.

²⁴²⁸ The findings from this process are presented and discussed in the next two
²⁴²⁹ sections.

²⁴³⁰ 6.4 Findings

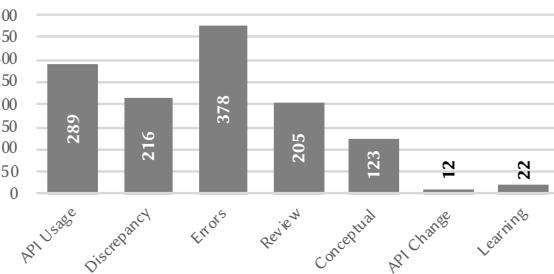


Figure 6.1: Distribution of SO question types.

²⁴³¹ Figure 6.1 displays the overall distribution of question types from the 1,245
²⁴³² posts classified in [79], when adjusted for majority ruling as per Section 6.3.2.1. It
²⁴³³ is evident that developers ask issues predominantly related to API errors when using
²⁴³⁴ CVSs and, additionally, how they can use the API to implement specific functionality.
²⁴³⁵ There are few questions related to version issues or self-learning.

Table 6.2: Frequency of emotions per question type.

Question Type	Fear	Joy	Love	Sadness	Surprise	Anger	No Emotion	Total
API Usage	50	22	34	18	59	13	135	331
Discrepancy	38	12	18	7	48	20	108	251
Errors	69	34	22	21	48	23	206	423
Review	34	16	15	16	42	14	98	235
Conceptual	26	10	10	7	21	5	59	138
API Change	4	2	2	1	1	1	5	16
Learning	3	4	2	0	4	0	11	24
Total	224	100	103	70	223	76	622	1418

²⁴³⁶ Table 6.2 displays the frequency of questions that were classified by EmoTxt
²⁴³⁷ when compared to our assignment of question types, while Figure 6.2 presents the
²⁴³⁸ emotion data proportionally across each type of question. *No Emotion* was the
²⁴³⁹ most prevalent across all question types, which is consistent with the findings of the
²⁴⁴⁰ Collab group during the training of the EmoTxt classifier. Interestingly, *API Change*
²⁴⁴¹ questions had a distinct distribution of emotions, where 31.25% of questions had *No*
²⁴⁴² *Emotion* compared to the average of 42.01%. This is likely due to the low sample
²⁴⁴³ size of *API Change* questions, with only 12 assignments, however the next highest
²⁴⁴⁴ set of emotive questions are found in the second largest sample (*API Usage*, at

59.21%) and so greater emotion detected is not necessarily proportional to sample size. Unsurprisingly, *Discrepancy* questions had the highest proportion of the *Anger* emotion, at 7.97%, compared to the mean of 4.74%, which is indicative of the frustrations developers face when the API does something unexpected. *Love*, an emotion which we expected least by software developers when encountering issues, was present across the different question types. The two highest emotions, by average, were *Fear* (16.67%) and *Surprise* (14.90%), while the two lowest emotions were *Sadness* (4.47%) and *Anger* (4.74%). *Joy* and *Love* were roughly the same and fell in between the two proportion ends, with means of 8.96% and 8.16%, respectively.

Results from our reliability analysis showed largely poor results. Guidelines of indicative strengths of agreement are provided by Landis and Koch [182], where $\kappa \leq 0.000$ is *poor agreement*, $0.000 < \kappa \leq 0.200$ is *slight agreement* and $0.200 < \kappa \leq 0.400$ is *fair agreement*. Our readings were indicative of poor agreement between raters ($C_\kappa = -0.003$) and slight agreement with EmoTxt ($L_\kappa = 0.155$). The strongest agreements found were for *No Emotion* both between each of our three raters ($L_\kappa = 0.292$) and each rater and EmoTxt ($C_\kappa = 0.086$), with fair and slight agreement respectively.

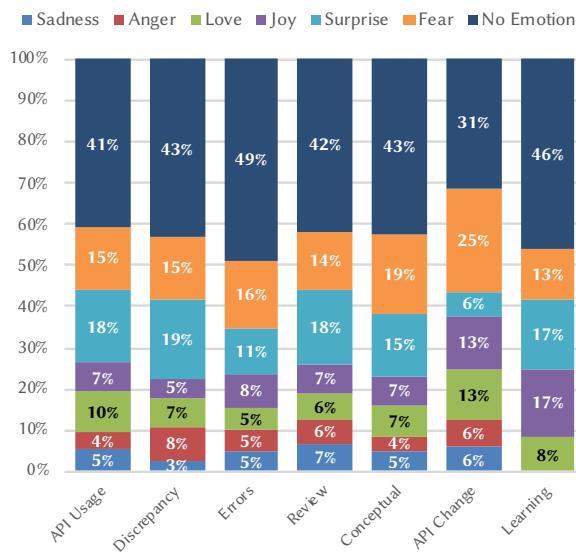


Figure 6.2: Proportion of emotions per question type.

6.5 Discussion

Our findings from the comparison between the manually annotated SO posts and the automatic classification revealed substantial discrepancies. Table 6.3 provide some sample questions from our data set and the emotion identified by EmoTxt within the text. A subset of questions analysed by our three raters do not indicate

Table 6.3: Sample questions comparing question type to emotion. Questions located at [https://stackoverflow.com/q/\[ID\]](https://stackoverflow.com/q/[ID]).

ID	Quote	Classification	Emotion
53249139	<i>"I'm trying to integrate my project with Google Vision API... I'm wondering if there is a way to set the credentials explicitly in code as that is more convenient than setting environment variables in each and every environment we are running our project on... I know for a former client version 1.22 that was possible... but for the new client API I was not able to find the way and documentation doesn't say anything in that regards."</i>	API Usage	Fear
40013910	<i>"I want to say something more about Google Vision API Text Detection, maybe any Google Expert here and can read this. As Google announced, their TEXT_DETECTION was fantastic... But for some of my pics, what happened was really funny... There must be something wrong with the text detection algorithm."</i>	Discrepancy	Anger
50500341	<i>"I just started using PYTHON and now i want to run a google vision cloud app on the server but I'm not sure how to start. Any help would be greatly appreciated."</i>	API Usage	Sadness
49466041	<i>"I am getting the following error when trying to access my s3 bucket... my hunch is it has something to do with the region...I have given almost all the permissions to the user I can think of.... Also the region for the s3 bucket appears to be in a place that can work with rekognition. What can I do?"</i>	Errors	Surprise
55113529	<i>"Following a tutorial, doing everything exactly as in the video... Hoping to figure this out as it is a very interesting concept...Thanks for the help... I'm getting this error:..."</i>	Errors	Joy
39797164	<i>"Seems that the Google Vision API has moved on and the open Sourced version has not....In my experiments this 'finds' barcodes much faster than using the processor that the examples show. Am I missing something somewhere?"</i>	API Change	Love

2468 the automatic (EmoTxt) emotion, and upon manual inspection of the text after poor
2469 results from our reliability analysis, an introspection of the data set sheds some light
2470 to the discrepancy. For example, question 55113529 shows no indication of *Joy*,
2471 rather the developer is expressing a state of confusion. The phrase “*Thanks for your*
2472 *help*” could be the reason why the miss-classification occurred if words like “thanks”
2473 were associated with joy. However, in this case, it seems unlikely that the developer
2474 is expressing joy as the developer has followed a tutorial but is still encountering
2475 an error. Similarly, question 39797164, classified as *Love* and question 50500341,
2476 classified as *Sadness* express a state of confusion and the urge to know more about the
2477 product; upon inspecting the entire question in context, it is difficult to consistently
2478 agree with the emotions as determined by EmoTxt, and further exploration into the
2479 behaviour and limitations of the model is necessary.

2480 Our results indicate further work is needed to refine the machine learning (ML)
2481 classifiers that mine emotions in the SO context. The question that arises is whether
2482 the classification model is truly reflective of real-world emotions expressed by soft-
2483 ware developers. As highlighted by Curumsing [82], the divergence of opinions with
2484 regards to the emotion classification model proposed by theorists raises doubts to
2485 the foundations of basic emotions. Most of the studies conducted in the area of emotion
2486 mining from text is based on an existing general purpose emotion framework
2487 from psychology [52, 226, 233]—none of which are tuned for software engineering
2488 domain. In our our study, we note the emotions expressed by software develop-
2489 ers within SO posts are quite narrow and specific. In particular, emotions such as
2490 frustration and confusion would be more appropriate over love and joy.

2491 6.6 Implications

2492 Based on our observations during the manual classification of SO posts and related
2493 work in the field [331], we propose a new taxonomy of emotions which is reflective
2494 of what software developers experience when encountering coding issues. We
2495 propose the following set of five emotions: (i) *Confusion*, an inability to understand
2496 something, e.g., “*why is the code not functioning?*” or “*where is the error?*”; (ii)
2497 *Frustration*, annoyance resulting from the inability to change or achieve something,
2498 e.g., “*I don’t understand why this code is not working.*”; (iii) *Curiosity*, an urge
2499 to learn more about the tool, e.g., “*I am looking for a way to do this...*”; (iv)
2500 *Contentedness*, where developers are satisfied with the current situation however
2501 there may be a small issue, e.g., “*It works pretty well, but...*”; and, (v) *Optimism*,
2502 hopeful that a solution can be found, e.g., “*I hope you can see what I’m doing
2503 wrong.*”.

2504 6.7 Threats to Validity

2505 6.7.1 Internal Validity

2506 The *API Change* and *Learning* question types were few in sample size (only 12 and
2507 22 questions, respectively). The emotion proportion distribution of these question

types are quite different to the others. Given the low number of questions, the sample is too small to make confident assessments. Furthermore, our assignment of Beyer et al.'s question type taxonomy was single-label; a multi-labelled approach may work better, however analysis of results would become more complex. A multi-labelled approach would be indicative for future work. Lastly, the study would be greatly improved with a reliability analysis of our proposed taxonomy; while we did resolve using majority voting (Section 6.3.2.4), no inter-rater reliability has been performed for this study. We plan to conduct reliability analysis, expand the number of raters, and increase the 25 question sample size in our future work for a more thorough analysis of our proposed taxonomy.

2518 **6.7.2 External Validity**

2519 EmoTxt was trained on questions, answers and comments, however our data set
2520 contained questions only. It is likely that our results may differ if we included other
2521 discussion items, however we wished to understand the emotion within developers'
2522 *questions* and classify the question based on the question classification framework
2523 by Beyer et al. [34]. Moreover, this study has only assessed frustrations within the
2524 context of a concrete domain of CVSs. The generalisability of this study to other
2525 IWSs, such as natural language processing services, or conventional web services,
2526 may be different. Furthermore, we only assessed four popular CVSs; expanding the
2527 data set to include more services, including non-English ones, would be insightful.
2528 We leave this to future work.

2529 **6.7.3 Construct Validity**

2530 Some posts extracted from SO were false positives. Whilst flagged for removal
2531 (Section 6.3.1.5), we cannot guarantee that all false positives were removed. Fur-
2532 thermore, SO is known to have questions that are either poorly worded or poorly
2533 detailed, and developers sometimes ask questions without doing any preliminary
2534 investigation. This often results in down-voted questions. We did not remove such
2535 questions from our data set, which may influence the measurement of our results.

2536 **6.8 Conclusions**

2537 In this paper we analysed SO posts for emotions using an automated tool and cross-
2538 checked it manually. We found that the distribution of emotion differs across the
2539 taxonomy of issues, and that the current emotion model typically used in recent
2540 works is not appropriate for emotions expressed within SO questions. Consistent
2541 with prior work [188], our results demonstrate that machine learning classifiers for
2542 emotion are insufficient; human assessment is required.

2543 Future work would include validating our proposed taxonomy of emotions
2544 through (1) a survey with software developers to identify the validity of the emotions
2545 present in the taxonomy; (2) manually classifying SO posts using the proposed emo-
2546 tion classification model to study the distribution of SO posts under each taxonomy

2547 of errors; and (3) extend the work to other communication channels used by software
2548 developers.

2549 CHAPTER 7

2550

2551

Better Documenting Computer Vision Services[†]

2552

2553 **Abstract** Using cloud-based computer vision services (CVSs) is gaining traction with
2554 developers for many applications for many reasons: developers can simply access these
2555 AI-components through familiar RESTful APIs, and need not orchestrate large training and
2556 inference infrastructures or curate and label large training datasets. However, while their
2557 APIs *seem* familiar to use, their non-deterministic run-time behaviour and evolution profile
2558 are not adequately communicated to developers, and this results in developers struggling
2559 to use such APIs in-practice. Therefore, improving these services' API documentation is
2560 paramount, as a more complete document facilities the development process of intelligent
2561 software. This study presents an analysis of what facets a 'complete' API document should
2562 have, as synthesised into a taxonomy from 21 academic studies via a systematic mapping
2563 study. We triangulate these findings from literature against 83 developers to assess the
2564 efficacy and utility in-practice of such knowledge. We produce two weighted 'scores'
2565 for each dimension in our taxonomy based on (i) the number of papers producing these
2566 outcomes and their citation count and (ii) the extent to which developers *agree* with the
2567 recommendations arising from these studies (based on our survey). Furthermore, we apply
2568 the taxonomy to three popular CVSs and assess their compliance, producing a third 'score'
2569 using the taxonomy to identify 12 suggested improvements to the API documentation of
2570 these intelligent web services.

2571

7.1 Introduction

2572 Improving API documentation quality is a valuable task for any API—an extensive
2573 API document facilitates productivity, and therefore improved quality is better en-
2574 gineered into a system [208]. Where application developers integrate new services
2575 (such as computer vision services (CVSs) [76]) into their systems via APIs, their

[†]This chapter is originally based on A. Cummaudo, R. Vasa, and J. Grundy, "Assessing API documentation knowledge for computer vision services," 2020, Unpublished. Terminology has been updated to fit this thesis.

productivity is affected either by inadequate skills (“*I’ve never used an API like this, so must learn from scratch*”) or, where their skills are adequate, an imbalanced cognitive load that causes excessive context switching (“*I have the skills for this, but am confused or misunderstand*”). This is commonly seen in the emerging computer vision (CV) web services space, where the documentation does not yet completely or correctly describe the APIs in full [79].

What causes a developer to be confused and how to mitigate it via an improved API document has been largely explored for conventional APIs. Various studies have provided a myriad of recommendations based on both qualitative and quantitative analysis of developer opinion. Such recommendations propose ways by which developers, managers and solution architects can construct systems better with improved documentation. However, while previous works have covered certain aspects of API usage, many have lacked a systematic review of literature and do not offer a taxonomy to consolidate these guidelines together. For example, some studies have considered the technical implementation improving API usability or tools to generate (or validate) API documentation from its source code (e.g., [196, 227, 320]); still lacks a consolidated effort to capture recommendations on how to *manually write* complete, correct, and effective API documentation. The works that *do* produce these recommendations from literature are largely scattered across multiple sources, and systematically capturing the information into a readily accessible, consolidated framework (designed to assist writing API documentation) must be validated in real-world circumstances to assess its efficacy with practitioners and existing documentation [75].

As a real-world use case, consider an intelligent web service (IWS)—such as CVSS—in which an AI-based component produces a non-deterministic result based on a machine-learnt data-driven algorithm, rather than a predictable, rule-driven one [76]. These services use machine intelligence to make predictions on images such as object labelling or facial recognition [341, 352, 353, 354, 355, 358, 362, 370, 371, 372, 376, 388, 389, 422, 423]. The impacts of poor and incomplete documentation results in developer complaints on online discussion forums such as Stack Overflow [79]. Many comments show that developers do not think in the non-deterministic mental model of the designers who created the CVSSs. They ask many varied questions from their peers to try and clarify their understanding.

This paper significantly extends our previous work [75] by evaluating our API documentation taxonomy in two additional contexts. In our previous work, we developed a weighted metric for each dimension and category based on how many literary sources agree that the aspects of our taxonomy should be implemented. We refer to this as an ‘in-literature’ agreement score. We build upon this facet but *in-practice* by assessing the efficacy of our taxonomy against developers using a survey built upon an interpretation of the System Usability Scale (SUS) [50]. We produce a second weighting for the dimensions and categories of the taxonomy, referred to as a ‘in-practice’ agreement score. We then compare both the in-literature and in-practice scores directly, thereby contrasting the statistical agreement the two have. Lastly, we assess the taxonomy against three popular CVSSs, namely Google Cloud Vision [362], Amazon Rekognition [341] and Azure Computer Vision [376].

2621 For each category in our taxonomy, we assess whether the respective service's
2622 documentation contains, partially-contains or does not contain the recommendation.
2623 From this, we triangulate each category's in-literature and in-practice score against
2624 the service's level of inclusion of the recommendation, thereby making a judgement
2625 as to where the services can improve their documentation to make them more
2626 complete.

2627 The primary contributions in this work are:

- 2628 • a systematic mapping study (SMS) consisting of 21 studies that capture what
2629 knowledge or artefacts should be contained within API documentation;
- 2630 • a five dimensional taxonomy consisting of 34 recommendations based on those
2631 consolidated from the 21 studies;
- 2632 • a score metric for each recommendation based on the number of papers that
2633 agree with the recommendation;
- 2634 • a score metric assessing the efficacy of the 34 recommendations that empiri-
2635 cally reflects what is important to document from a *practitioner* point of view;
2636 and,
- 2637 • a heuristic validation of each recommendation against CVSs, assessing where
2638 existing CVS API documentation needs improvement.

2639 After performing our SMS on what API knowledge should be captured in doc-
2640 umentation to assist API designers, we propose our taxonomy consisting of the
2641 following dimensions: (1) Usage Description; (2) Design Rationale; (3) Domain
2642 Concepts; (4) Support Artefacts; and (5) Documentation Presentation. Following
2643 this, we adopted the SUS surveying technique to assess the overall utility of each
2644 of these recommendations, producing a survey consisting of 43 questions. This
2645 survey was then tested three times within our research group: firstly against three
2646 researchers for feedback on the survey's design, secondly against three software
2647 engineers in our research group with varying levels of experience for developers
2648 for test-retest reliability [169], thirdly against 22 software engineers in our research
2649 group for wider feedback on the survey. Given these feedback improvements, we
2650 surveyed 83 external developers between May 2019 to October 2019, and then anal-
2651 ysed the relevance of each recommendation from the practitioner's viewpoint. We
2652 also assessed the three CVSs for inclusion of each recommendation, and once our
2653 surveys were complete, determined a weighted 'score' of each service to see where
2654 improvements to their documentation was made.

2655 This paper is structured as thus: Section 7.2 presents related work in the areas
2656 of API usability, intelligent CVSs, and the SUS; Section 7.3 is divided into two
2657 subsections, the first describing how primary sources were selected in a SMS with the
2658 second describing the development of our taxonomy from these sources; Section 7.4
2659 presents the taxonomy; Section 7.5 describes how we developed a survey instrument
2660 of 43 questions to validate the taxonomy against developers, and assess its efficacy
2661 against the three popular CVSs selected to make 12 suggested improvements to
2662 the existing service API documentation; Section 7.6 presents the findings from our
2663 validation analysis and the weightings for the taxonomy; Section 7.7 describes the

2664 threats to validity of this work and Section 7.8 provides concluding remarks and the
2665 future directions of this study. Additional materials are provided in Appendix B.

2666 **7.2 Related Work**

2667 **7.2.1 API Usability and Documentation Knowledge**

2668 Use of the SMS approach has explored developer experience and API usability.
2669 A 2018 study reviewed 36 API documentation generation tools and approaches, and
2670 analysed the tools developed and their inputs and documentation outputs [227]. The
2671 findings from this study emphasise that the largest effort in API documentation tool-
2672 ling is to assist developers to generate either example code snippets and/or templates
2673 or natural language descriptions of the API directly from the program’s source code.
2674 These snippets or descriptions can then be placed in the API documentation, thereby
2675 increasing the efficiency at which API documentation can be written. Additionally,
2676 tools from 12 studies target the maintainability of existing APIs of existing APIs,
2677 with tools from 11 studies target the correctness and accuracy of the documentation
2678 by validating that what is written in the documentation is accurate to the technical
2679 structure of the API. From the end-developer’s perspective, some tools (17 studies)
2680 help target improvements to the developer’s understandability and learnability of
2681 new APIs by linking in examples directly with questions such as on Stack Overflow.

2682 However, the results from this study regards the *tooling* used to either assist in
2683 producing, validating or learning from API documentation. While this is a systematic
2684 study with key insights into the types of tooling produced, there is still a gap for a
2685 SMS in what *guidelines* have been produced by the literature in developing natural-
2686 language documentation itself and how well developers *agree* to those guidelines,
2687 which our work has addressed.

2688 Watson [320] performed a heuristic assessment from 35 popular APIs against 11
2689 high-level universal design elements of API documentation. This study highlighted
2690 how many APIs, even popular ones, fail to grasp these basic design elements.
2691 For example, 25% of the documentation sets did not provide any basic overview
2692 documentation to the API. The heuristics used within Watson’s study is based on
2693 only three seminal works and only contains 11 design elements—our study extends
2694 these heuristics and structures them into a consolidated, hierarchical taxonomy which
2695 we then validate against practitioners.

2696 A taxonomy of distinct knowledge patterns within reference documentation
2697 by Maalej and Robillard [196] classified 12 distinct knowledge types. The tax-
2698 onomy was then evaluated against the JDK 6 and .NET 4.0 frameworks, and showed
2699 that the functionality and structure of these APIs are well-communicated, although
2700 core concepts and rationale about the API are quite rarer to see. The authors also
2701 identified low-value ‘non-information’—described as documentation that provides
2702 uninformative boilerplate text with no insight into the API at all—which was sub-
2703 stantially present in the documentation of methods and fields in the two frameworks.
2704 They recommend that developers factor their 12 distinct knowledge types into the
2705 process of code documentation, thereby preventing low-value non-information. The

2706 development of their taxonomy consisted of questions to model knowledge and information,
2707 thereby capturing the reason about disparate information units independent
2708 to context; a key difference to this paper is the systematic taxonomy approach utilised.

2709 7.2.2 Adapting the System Usability Scale

2710 The SUS was first introduced by Brooke as early as 1986 as a “quick and dirty”
2711 survey scale to easily assess the overall usability of a product or service in a timely
2712 manner. Its popularity in the usability community demonstrated the need for a
2713 tool that can collect a quantifiable rating of usability from a participant’s subjective
2714 opinion, and was later published in [50]. Since, its adoption as an industry standard is
2715 widely demonstrated [19, 51] and studies have adopted its ease of use for generalised
2716 purposes.

2717 While translation of the SUS into other languages [41, 201, 272] is generally
2718 the most adapted form of Brooke’s original survey, some studies have proposed
2719 alternative measurement models to the SUS, such as separating the usability and
2720 learnability components of the survey into a two-dimensional structure [41]. Other
2721 adaptations of the SUS include a 2014 study that proposed a usability scale based
2722 on the SUS for Handheld Augmented Reality applications [271] conceptualised
2723 against comprehensibility and manipulability. However, few studies have designed
2724 questionnaires patterned from the SUS in other contexts, and to our knowledge, this
2725 study presents an initial attempt at doing so in the API documentation knowledge
2726 domain.

2727 7.2.3 Computer Vision Services

2728 Recent studies into cloud-based CVSSs have demonstrated that poor reliability and
2729 robustness in CV can ‘leak’ into end-applications if such aspects are not sufficiently
2730 appreciated by developers. A study by Hosseini et al. [140] showed that Google
2731 Cloud Vision’s labelling fails when as little as 10% noise is added to the image. Facial
2732 recognition classifiers are easily confused by modifying pixels of a face and using
2733 transfer learning to adapt one person’s face into another [316]. Our own prior work
2734 found that the non-deterministic evolution of these types of services is not adequately
2735 communicated to developers [76], resulting in lost developer productivity whereby
2736 developers ask fundamental questions about the concepts behind these services, how
2737 they work, and where better documentation can be found [79]. This paper continues
2738 this line of research by providing a means for service providers to better document
2739 their services using a taxonomy and suggested improvements.

2740 7.3 Taxonomy Development

2741 We developed our taxonomy under two primary phases. First, we conducted a SMS
2742 identifying API documentation studies, following guidelines by Kitchenham and
2743 Charters [168] and Petersen et al. [244] (Section 7.3.1). A high level overview of
2744 this first phase is given in Figure 7.2. Second, we followed a software engineering
2745 (SE) taxonomy development method by Usman et al. [311] (Section 7.3.2) based on

²⁷⁴⁶ the findings of our SMS, which involved an extensive validation involving real-world
²⁷⁴⁷ developers and contextualised with CV APIs (Section 7.5).

²⁷⁴⁸ 7.3.1 Systematic Mapping Study

²⁷⁴⁹ 7.3.1.1 Research Questions (RQs)

²⁷⁵⁰ The first step in producing our SMS was to pose two RQs:

- ²⁷⁵¹ • **RQ1:** What knowledge do API documentation studies contribute?
- ²⁷⁵² • **RQ2:** How is API documentation studied?

²⁷⁵³ Our intent behind RQ1 was to collect as many studies provided by literature on how
²⁷⁵⁴ API documentation should be written using natural language (i.e., not using assistive
²⁷⁵⁵ tooling). This helped us shape and form the taxonomy provided in Section 7.4.
²⁷⁵⁶ Secondly, RQ2's intent was to understand how the studies derive at their conclusions,
²⁷⁵⁷ thereby helping us identify gaps in literature where future studies can potentially
²⁷⁵⁸ focus.

²⁷⁵⁹ 7.3.1.2 Automatic Filtering

²⁷⁶⁰ As done in similar SE studies [114, 120, 311], we explored automatic filtering of
²⁷⁶¹ online databases. We defined which SWEBOK knowledge areas [145] were relevant
²⁷⁶² to devise a search query. Our search query was built using related knowledge areas,
²⁷⁶³ relevant synonyms, and the term 'software engineering' (for comprehensiveness) all
²⁷⁶⁴ joined with the OR operator. Due to the lack of a standard definition of an API,
²⁷⁶⁵ we include the terms: 'API' and its expanded term; software library, component
²⁷⁶⁶ and framework; and lastly software development kit (SDK). These too were joined
²⁷⁶⁷ with the OR operator, appended with an AND. Lastly, the term 'documentation' was
²⁷⁶⁸ appended with an AND. Our final search string was:

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

²⁷⁶⁹ We executed the query on all available metadata (title, abstract and keywords) in
²⁷⁷⁰ May 2019 against Web of Science¹ (WoS), Compendex/Inspec² (C/I) and Scopus³.
²⁷⁷¹ We selected three particular primary sources given their relevance in SE literature
²⁷⁷² (containing the IEEE, ACM, Springer and Elsevier databases) and their ability to
²⁷⁷³ support advanced queries [49, 168]. A total 4,501 results⁴ were found, with 549
²⁷⁷⁴ being duplicates. Table 7.1 displays our results in further detail (duplicates not
²⁷⁷⁵ omitted); Figure 7.1 shows an exponential trend of API documentation publications

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>.

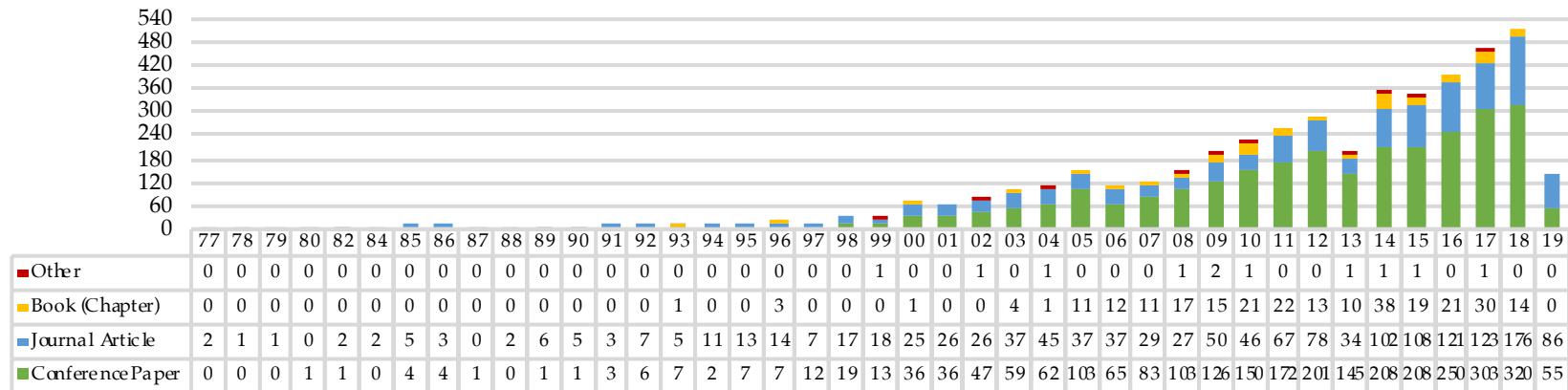


Figure 7.1: Search results by year and venue type.

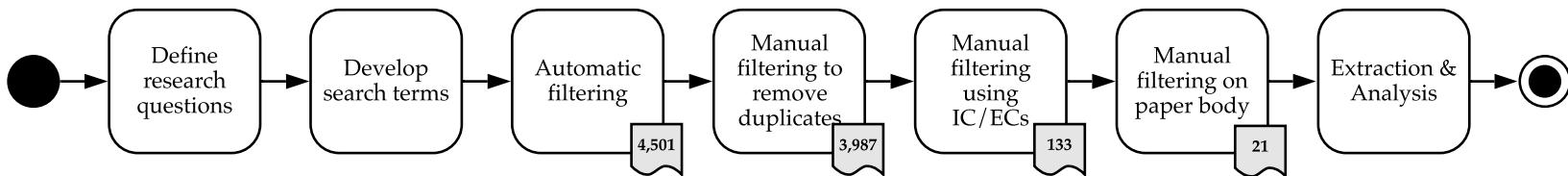


Figure 7.2: A high level overview of the filtering steps from defining and executing our search query to the data extraction of our primary studies. Number of accepted papers resulting from each filtering step is shown.

Table 7.1: 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

²⁷⁷⁶ produced within the last two decades. (As this search was conducted in May 2019,
²⁷⁷⁷ results taper in 2019.)

²⁷⁷⁸ 7.3.1.3 *Manual Filtering*

²⁷⁷⁹ A follow-up manual filtering stage followed the 4,501 results obtained by automatic
²⁷⁸⁰ filtering. As described below, we applied the following inclusion criteria (IC) and
²⁷⁸¹ exclusion criteria (EC) to each result:

- ²⁷⁸² **IC1** Studies must be relevant to API documentation: specifically, we exclude
²⁷⁸³ studies that deal with improving the technical API usability (e.g., improved
²⁷⁸⁴ usage patterns);
- ²⁷⁸⁵ **IC2** Studies must propose new knowledge or recommendations to document
²⁷⁸⁶ APIs;
- ²⁷⁸⁷ **IC3** Studies must be relevant to SE as defined in SWEBOK;
- ²⁷⁸⁸ **EC1** Studies where full-text is not accessible through standard institutional databases;
- ²⁷⁸⁹ **EC2** Studies that do not propose or extend how to improve the official, natural
²⁷⁹⁰ language documentation of an API;
- ²⁷⁹¹ **EC3** Studies proposing a third-party tool to enhance existing documentation or
²⁷⁹² generate new documentation using data mining (i.e., not proposing strategies
²⁷⁹³ to improve official documentation);
- ²⁷⁹⁴ **EC4** Studies not written in English;
- ²⁷⁹⁵ **EC5** Studies not peer-reviewed.

²⁷⁹⁶ Each of these ICs and ECs were applied to every paper after exporting all
²⁷⁹⁷ metadata of our results to a spreadsheet. The first author then curated the publications
²⁷⁹⁸ using the following revision process.

²⁷⁹⁹ Firstly, we read the publication source—to rapidly omit non-SE papers—as well
²⁸⁰⁰ as the author keywords, title, and abstract of all 4,501 studies. As some studies were
²⁸⁰¹ duplicated between our three primary sources, we needed to remove any repetitions.
²⁸⁰² We sorted and reviewed any duplicate DOIs and fuzzy-matched all very similar titles
²⁸⁰³ (i.e., changes due to punctuation between primary sources), thereby retaining only
²⁸⁰⁴ one copy of the paper from a single database. Similarly, as there was no limit do
²⁸⁰⁵ our date ranges, some studies were republished in various venues (i.e., same title but
²⁸⁰⁶ different DOIs). These were also removed using fuzzy-matching on the title, and the

2807 first instance of the paper's publication was retained. This second phase resulted in
2808 3,987 papers.

2809 Secondly, we applied our inclusion and exclusion criteria to each of the 3,987
2810 papers by reading the abstract. Where there was any doubt in applying the criteria
2811 to the abstract alone, we automatically shortlisted the study. We rejected 427 studies
2812 that were unrelated to SE, 3,235 were not directly related to documenting APIs
2813 (e.g., to enhance coding techniques that improve the overall developer usability of
2814 the API), 182 proposed new tools to enhance API documentation or used machine
2815 learning to mine developer's discussion of APIs, and 10 were not in English. This
2816 resulted in 133 studies being shortlisted to the final phase.

2817 Thirdly, we re-evaluated each shortlisted paper by re-reading the abstract, the
2818 introduction and conclusion. We removed a further 64 studies that were on API
2819 usability or non API-related documentation (i.e., code commenting). At this stage,
2820 we decided to refine our exclusion criteria to better match the research goals of this
2821 study by including the word 'natural language' documentation in EC2. This removed
2822 studies where the focus was to improve technical documentation of APIs such as
2823 data types and communication schemas. Additionally, we removed 26 studies as
2824 they were related to introducing new tools (EC3), 3 were focused on tools to mine
2825 API documentation, 7 studies where no recommendations were provided, 2 further
2826 duplicate studies, and a further 10 studies where the full text was not available,
2827 not peer reviewed or in English. Books are commonly not peer-reviewed (EC5),
2828 however no books were shortlisted within these results. This final stage resulted in
2829 21 primary studies for further analysis, and the mapping of primary study identifiers
2830 to references S1–21 can be found in Appendix B.3.

2831 As a final phase, we conducted reliability analysis of our shortlisting method.
2832 We conducted intra-rater reliability of our 133 shortlisted papers using the test-
2833 retest approach suggested by Kitchenham and Charters [168]. We re-evaluated a
2834 random sample of 10% of the 133 shortlisted papers a week after initial studies were
2835 shortlisted. This resulted in *substantial agreement* [182], measured using Cohen's
2836 kappa ($\kappa = 0.7547$).

2837 7.3.1.4 Data Extraction & Systematic Mapping

2838 Of the 21 primary studies, we conducted abstract key-wording adhering to Petersen
2839 et al.'s guidelines [244] to develop a classification scheme. An initial set of keywords
2840 were applied for each paper in terms of their methodologies and research approaches
2841 (RQ2), based on an existing classification schema used in the requirements engineering
2842 field by Wieringa et al. [326]. These are: *evaluation papers*, which evaluates
2843 existing techniques in-practice; *validation papers*, which investigates proposed tech-
2844 niques not yet implemented in-practice; *experience papers*, which do investigate or
2845 evaluate either proposed or existing techniques, but presents insightful experiences
2846 of authors that warrant communication to other practitioners; and *philosophical pa-
2847 pers*, which presents new conceptual frameworks that describes a language by which
2848 we can describes our observations of existing or new techniques, thereby implying
2849 a new viewpoint for understanding phenomena.

Table 7.2: 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 et al.'s taxonomy

2850 After all primary studies had been assigned keywords, we noticed that all papers
 2851 used field study techniques, and thus we consolidated these keywords using Singer
 2852 et al.'s framework of SE field study techniques [287]. Singer et al. captures both
 2853 study techniques *and* methods to collect data within the one framework, namely:
 2854 *direct techniques*, including brainstorming and focus groups, interviews and ques-
 2855 tionnaires, conceptual modelling, work diaries, think-aloud sessions, shadowing and
 2856 observation, participant observation; *indirect techniques*, including instrumenting
 2857 systems, fly-on-the-wall; and *independent techniques*, including analysis of work
 2858 databases, tool use logs, documentation analysis, and static and dynamic analysis.

2859 Table 7.2 describes our data extraction form, which was used to collect relevant
 2860 data from each paper. Figure 7.3 presents our systematic mapping, where each study
 2861 is mapped to one (or more, if applicable) of methodologies plotted against Wieringa
 2862 et al.'s research approaches. We find that a majority of these studies survey develop-
 2863 ers using direct techniques (i.e., interviews and questionnaires) and some performing
 2864 structured documentation analysis. Few studies report recent experiences, with the
 2865 majority of API documentation knowledge being evaluation research, and some val-
 2866 idation studies. There are few experience papers describing anecdotal evidence of
 2867 API documentation knowledge, and almost no philosophical papers that describe new
 2868 conceptual ways at approaching API documentation as a large majority of existing
 2869 work either evaluates existing (in-practice) strategies or validates the effectiveness
 2870 of new strategies.

2871 **7.3.2 Development of the Taxonomy**

2872 A majority of taxonomies produced in SE studies are often made extemporane-
 2873 ously [311]. For this reason, we decided to proceed with a systematic approach to
 2874 develop our taxonomy using the guidelines provided by Usman et al. [311], which
 2875 are extended from lessons learned in more mature domains. In this subsection, we
 2876 outline the 4 phases and 13 steps taken to develop our taxonomy based on Usman

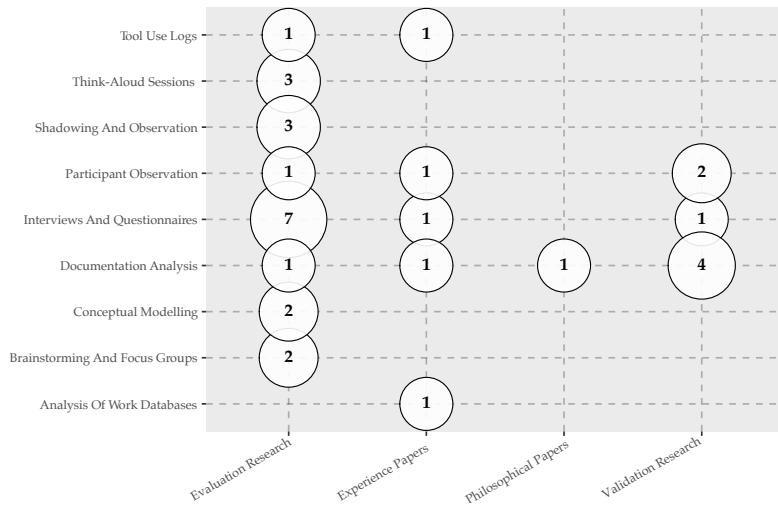


Figure 7.3: Systematic map: field study technique vs research type

²⁸⁷⁷ et al.'s technique. Usman et al.'s final *validation* phase is largely detailed within
²⁸⁷⁸ Section 7.5 after we present our taxonomy in Section 7.4.

²⁸⁷⁹ 7.3.2.1 Planning phase

²⁸⁸⁰ The preliminary phase involves answering the following:

- ²⁸⁸¹ **(1) define the SE knowledge area:** The SE knowledge area, as defined by the
²⁸⁸² SWEBOK, is software construction;
- ²⁸⁸³ **(2) define the objective:** The main objective of the proposed taxonomy is to define
²⁸⁸⁴ a set of categories that enables to classify different facets of natural-language
²⁸⁸⁵ API documentation knowledge (not API usability knowledge) as reported in
²⁸⁸⁶ existing literature;
- ²⁸⁸⁷ **(3) define the subject matter:** The subject matter of our proposed taxonomy is
²⁸⁸⁸ documentation artefacts of APIs;
- ²⁸⁸⁹ **(4) define the classification structure:** The classification structure of our proposed
²⁸⁹⁰ taxonomy is *hierarchical*;
- ²⁸⁹¹ **(5) define the classification procedure:** The procedure used to classify the docu-
²⁸⁹² mentation artefacts is qualitative;
- ²⁸⁹³ **(6) define the data sources:** The basis of the taxonomy is derived from field study
²⁸⁹⁴ techniques (see Section 7.3.1.4).

²⁸⁹⁵ 7.3.2.2 Identification and extraction phase

²⁸⁹⁶ The second phase of the taxonomy development involves **(7) extracting all terms**
²⁸⁹⁷ and **concepts** from relevant literature, which we have achieved from our SMS. These
²⁸⁹⁸ terms are then consolidated by **(8) performing terminology control**, as some terms
²⁸⁹⁹ may refer to different concepts and vice-versa.

2900 7.3.2.3 *Design phase*

2901 The design phase identified the core dimensions and categories within the extracted
2902 data items. The first step is to (9) *identify and define taxonomy dimensions*; for this
2903 study we utilised a bottom-up approach to identify each dimension, i.e., extracting the
2904 categories first and then nominating which dimensions these categories fit into using
2905 an iterative approach. As we used a bottom-up approach, step (9) also encompassed
2906 the second stage of the design phase, which is to (10) *identify and describe the*
2907 *categories* of each dimension. Thirdly, we (11) *identify and describe relationships*
2908 between dimensions and categories, which can be skipped if the relationships are
2909 too close together, as is the case of our grouping technique which allows for new
2910 dimensions and categories to be added. The last step in this phase is to (12) *define*
2911 *guidelines for using and updating the taxonomy*. The taxonomy is as simple as a
2912 checklist that can be heuristically applied to an API document, and each dimension
2913 is malleable and covers a broad spectrum of artefacts; while we do not anticipate
2914 any further dimensions to be added, new categories can easily be fitted into one of
2915 the dimensions (see Section 7.8). We provide guidelines for use in our application
2916 of the taxonomy against CVSs within Sections 7.4 and 7.6.

2917 7.3.2.4 *Validation phase*

2918 In the final phase of taxonomy development, taxonomy designers must (13) *validate*
2919 *the taxonomy* to assess its usefulness. Usman et al. [311] describe three approaches to
2920 validate taxonomies: (i) orthogonal demonstration, in which the taxonomy's orthog-
2921 onality is demonstrated against the dimensions and categories, (ii) benchmarking
2922 the taxonomy against similar classification schemes, or (iii) utility demonstration by
2923 applying the taxonomy heuristically against subject-matter examples. In our study,
2924 we adopt utility demonstration by use of a survey and heuristic application of the
2925 taxonomy against real-world case-studies (i.e., within the domain of CVSs). This is
2926 is discussed in greater detail within Section 7.5.

2927 **7.4 API Documentation Knowledge Taxonomy**

2928 Our taxonomy consists of five dimensions (labelled A–E). We expand these five di-
2929 mensions into 34 categories (sub-dimensions). Each dimension respectively covers:

- 2930 • **[A] Usage Description** on *how* to use the API for the developer's intended
2931 use case;
- 2932 • **[B] Design Rationale** on *when* the developer should choose this API for a
2933 particular use case;
- 2934 • **[C] Domain Concepts** of the domain behind the API to understand *why* this
2935 API should be chosen for this domain;
- 2936 • **[D] Support Artefacts** that describe *what* additional documentation the API
2937 provides; and
- 2938 • **[E] Documentation Presentation** to help organise the *visualisation* of the
2939 above information.

[A] Usage Description

- [A1] Quick-start guides  #3
- [A2] Low-level reference manual  #3  SH 
- [A3] Explanation of high level architecture 
- [A4] Introspection source code comments  SH 
- [A5] Code snippets of basic component function  #2  #1  VH 
- [A6] Step-by-step tutorials with multiple components  #2  SH
- [A7] Downloadable production-ready source code
- [A8] Best-practices of implementation
- [A9] An exhaustive list of all components
- [A10] Minimum system requirements to use the API
- [A11] Instructions to install/update the API and its release cycle  #4
- [A12] Error definitions describing how to address problems  #5

[B] Design Rationale

- [B1] Entry-point purpose of the API  #4 
- [B2] What the API can develop
- [B3] Who should use the API
- [B4] Who will use the applications built using the API 
- [B5] Success stories on the API
- [B6] Documentation comparing similar APIs to this API
- [B7] Limitations on what the API can/cannot provide  #1

[C] Domain Concepts

- [C1] Relationship between API components and domain concepts
- [C2] Definitions of domain terminology
- [C3] Documentation for nontechnical audiences 

[D] Support Artefacts

- [D1] FAQs 
- [D2] Troubleshooting hints 
- [D3] API diagrams
- [D4] Contact for technical support  NH 
- [D5] Printed guide
- [D6] Licensing information

[E] Documentation Presentation

- [E1] Searchable knowledge base 
- [E2] Context-specific discussion forums
- [E3] Quick-links to other relevant components 
- [E4] Structured navigation style 
- [E5] Visualised map of navigational paths 
- [E6] Consistent look and feel  #5 

Figure 7.4: Our proposed taxonomy on what artefacts should be documented in a complete API document.

2940 Further descriptions of the categories encompassing each dimension are given within
2941 Figure 7.4 and Appendix B.1, coded as $[Xi]$, where i is the category identifier within
2942 a dimension, $X \in \{A, B, C, D, E\}$.

2943 Appendix B.1 shows which of the primary sources (S1–21) provide the rec-
2944 ommendation described as well as an ‘in-literature score’ (ILS). This score is a
2945 weighting calculated as a percentage of the number of primary studies that make the
2946 recommendation divided by the total of primary studies, and indicates the overall
2947 level of agreement that academic sources suggest these documentation artefacts.
2948 This score is contrasted to the ‘in-practice score’ (IPS) which indicates the over-
2949 all level of agreement that *practitioners* think such documentation artefacts are
2950 needed. Further details about the ILS and IPS values, how they were calculated and
2951 analysed for each category, and a rigorous contrast between the two are provided
2952 Section 7.5.1.2 and Sections 7.6.1 to 7.6.3. For comparative purposes, we illustrate
2953 a colour scale (from red to green) to indicate the relevancy weight between ILS and
2954 IPS values in Appendix B.1: for example, while quick-start guides [A1] are few ref-
2955 erenced in academic sources at 14%, they are generally well-desired by practitioners
2956 88% agreement. We then provide three columns that assesses the presence of these
2957 documentation artefacts against three popular CVSSs: Google Cloud Vision, AWS’s
2958 Rekognition, and Azure Cloud Vision (abbreviated to GCV, AWS and ACV). A
2959 fully shaded circle (●) indicates that the documentation artefact was clearly found
2960 in the service, while a half-shaded circle (◐) indicates that the artefact was only
2961 partially present. An outlined circle (○) indicates that the service lacks the indicated
2962 documentation artefact within our taxonomy. This empirical assessment is further
2963 detailed in Section 7.6.5, which outlines concrete areas in the respective services’
2964 documentation where improvements could be made, as well as hyperlinks to the
2965 documentation where relevant.

2966 Figure 7.4 illustrates these findings, with underlines indicating key artefacts and
2967 various iconography to indicate specific results. The computer icon (💻) includes a
2968 ranking from 1–5 of the top five most recommended artefacts according to devel-
2969 opers, as calculated from their relevant IPS scores. Conversely, the book icon (📘)
2970 indicates the rankings of the top five most recommended artefacts according to liter-
2971 ature. For example, while literature suggests the most useful documentation artefact
2972 are API usage description code snippets [A5], in-practice, we find that developers
2973 prefer design rationale on what the limitations of API are [B7] with code snippets
2974 coming in second place. Where there is strong agreement between developers and
2975 literature (within a standard deviation of 0.15) we use the handshake icon (🤝) and
2976 list whether both agree if the category is Very Helpful (VH), Slightly Helpful (SH) or
2977 Not Helpful (NH). Further details on this explanation are provided in Section 7.6.3.
2978 Lastly, we provide iconography for the presence (✓) or non-presence (✗) of these
2979 artefacts in *all three* CVSSs assessed, per Section 7.6.2.

2980 7.5 Validating our Taxonomy**2981 7.5.1 Survey Study****2982 7.5.1.1 Designing the Survey**

2983 We followed the guidelines by Kitchenham and Pfleeger [169] on conducting personal opinion surveys in SE to validate our survey. In developing our survey instrument, we shaped questions around each of our 5 dimensions and 34 categories. To achieve this, we used Brooke's SUS [50] as inspiration and re-shaped the 34 categories around a question. Each dimension was marked a numeric question (3–7), and alphabetic sub-questions were marked for each sub-dimension or category.

2989 We used closed questioning where respondents could choose an answer on a 2990 5-point Likert-scale (1=*strongly disagree*, 2=*somewhat disagree*, 3=*neither agree* 2991 *nor disagree*, 4=*slightly agree* and 5=*strongly agree*). Like Brooke's study, each 2992 question alternated in positive and negative sentiment. Half of our questions were 2993 written where a likely common response would be in strong agreement and vice- 2994 versa for the other half, such that participants would have to “read each statement 2995 and make an effort to think whether they would agree or disagree with it” [50]. For 2996 example, the question regarding [B7] on API limitations was framed as: “*I believe it* 2997 *is important to know about what the limitations are on what the API can and cannot* 2998 *provide*” (Q4g), whereas the question regarding [C1] on domain concepts of the API 2999 was framed as: “*I wouldn't read through theory about the API's domain that relates* 3000 *theoretical concepts to API components and how both work together*” (Q5a).

3001 In addition, the remaining eight questions asked demographical information. 3002 An extra open question asked for further comments. The full survey is provided in 3003 Appendix B.4.

3004 7.5.1.2 Evaluating the Survey

3005 After the first pass at designing questions was completed, we evaluated our survey 3006 on three researchers within our research group for general feedback. This resulted 3007 in minor changes, such as slight re-wording of questions, clarifying the difference 3008 between web services and web APIs, and providing specific questions with examples 3009 (some with images). For example, the question regarding [A9] on an exhaustive list 3010 of all major components in the API was framed as “*I believe an exhaustive list of all* 3011 *major components in the API without excessive detail would be useful when learning* 3012 *an API*” (Q3i) with the example “e.g., a CV web API might list object detection, 3013 *object localisation, facial recognition, and facial comparison as its 4 components*”.

3014 After this, we conducted reliability analysis using a test-retest approach on three 3015 developers within our group seven weeks apart. This was calculated using the `irr` 3016 computational R package [112] (as suggested in [128]) and resulted in an average 3017 intra-class correlation of 0.63 which indicates a good overall index of agreement [67].

3018 **7.5.1.3 Recruiting Participants**

3019 Our target population for the study was application software developers with varying
3020 degrees of experience (including those who and who have not used CVSs or related
3021 tools before) and varying understanding of fundamental machine learning concepts.
3022 We began by recruiting software developers within our research group using a
3023 group-wide message sent on our internal messaging system. Of the 44 developers in
3024 our group's engineering cohort, 22 responses were returned, indicating an internal
3025 response rate of 50%.

3026 For external participant recruiting, we shared the survey on social media plat-
3027 forms and online-discussion forums relevant to software development. We adopted
3028 a non-probabilistic snowballing sampling where the participants, at the end of the
3029 survey, were encouraged to share the survey link to others using *AddThis*⁵. This
3030 resulted in 43 additional visits to the survey. Additionally, snowballing sampling was
3031 encouraged within members of our research group who shared the survey with an
3032 additional 21 participants. However, while there were a total of 86 respondents, only
3033 51 finished the survey, leaving 35 participants with partially completed responses.
3034 Our final response rate was therefore 59%, which is very close to median response
3035 rates of 60% [24] in information systems and 5% in SE [287].

3036 **7.5.1.4 Analysing Response Data**

3037 To analyse our response data, we used an adapted version of the SUS method to
3038 produce a score for each question's 5-point response. As per Brooke's methodology,
3039 we mapped the responses from their ordinal scale of 1–5 to 0–4, and subtracted that
3040 value by 1 for positive questions and subtracted the value from 5 for the negative
3041 questions [50]. Unlike Brooke's method, we averaged each response for every
3042 question and divided by four (i.e., now a 4-point scale) to obtain scores for each
3043 category. This is presented in Appendix B.1 under the 'in-practice score' (IPS) for
3044 each category.

3045 Demographics for our survey were consistent in terms of the experience levels of
3046 developers who responded. Most were professional programmers with 75% report-
3047 ing between 1–10 years of work experience. A majority of our respondents (33%)
3048 reported to be in mid-tier roles. Most worked in either consulting or information
3049 technology services, reported at 17% for both.

3050 **7.5.2 Empirical application of the taxonomy against Computer Vision
3051 Services**

3052 Once our taxonomy had been developed, we performed an empirical application
3053 against three CVSs: Google Cloud Vision [362], Amazon Rekognition [341] and
3054 Azure Computer Vision [376]. Our selection criteria in choosing these particular
3055 services to analyse is based on the prominence of the service providers in industry
3056 and the ubiquity of their cloud platforms (Google Cloud, Amazon Web Services,
3057 and Microsoft Azure) in addition to being the top three adopted vendors used for

⁵<https://www.addthis.com> last accessed 7 January 2020

3058 cloud-based enterprise applications [259]. In addition, we had conducted extensive
3059 investigation into the services’ non-deterministic runtime behaviour and evolution
3060 profile in prior work [76] and have also identified developers’ complaints about their
3061 incomplete documentation in a prior mining study on Stack Overflow [79].

3062 We began with an exploratory analysis of the presence of each dimension and
3063 its categories. Appendix B.2 displays all sources of documentation used; although
3064 we initially started on the respective services homepages [341, 362, 376], this search
3065 was expanded to other webpages hyperlinked. For each category, we listed the
3066 documentation’s presence as either fully present, partially present or not present
3067 at all. This is shown in Appendix B.1 with the indication of (half-)filled circles or
3068 circle outlines for Google Cloud Vision (abbreviated to GCV), Amazon Rekognition
3069 (abbreviated to AWS), and Azure Computer Vision (abbreviated to ACV). Notes were
3070 taken for each webpage justifying the presence, and exact sources of documentation
3071 were listed when (partially) present. PDFs of each webpage were downloaded
3072 between 14–18 March 2019 for analysis.

3073 Once our analysis was completed and results from the survey finalised, we then
3074 calculated *weighted* ILS and IPS values for each dimension’s category. This was done
3075 by multiplying the ILS and IPS values for each category (listed in Appendix B.1) by
3076 either 0, 0.5 or 1 for categories not present, partially present, or present (respectively)
3077 in each service. The ‘maximum’ ILS and IPS values indicate the highest possible
3078 score a service can be ranked as though *all* categories are present. Tables 7.3 and 7.4
3079 show the sum of weights for each category in its respective dimension, in addition to
3080 the maximum possible score. Again, we use the same abbreviations for each service
3081 as per Appendix B.1. The scores are normalised into percentages for comparative
3082 purposes as a ratio of the score over all dimensions for a particular service to
3083 the maximum possible score. For comparative purposes, these are illustrated in
3084 Figure 7.6.

3085 7.6 Taxonomy Analysis

3086 In this section, we analyse investigating the taxonomy from two perspectives. Firstly,
3087 we describe the ILS values, being an interpretation of the number of papers that con-
3088 clude the recommendations in each category and dimension, and the weighted ILS
3089 scores, being an application of the taxonomy specifically to CVSs. Secondly, we look
3090 at the results from our survey and their respective IPS values, being an interpretation
3091 of how well developers agree with these recommendations, and the weighted IPS
3092 scores, being the application of how application developers would agree with the
3093 documentation of the CVSs. We then contrast the difference between what literature
3094 recommends and how well developers agree with these recommendations.

3095 7.6.1 In-Literature Scores for Taxonomy Categories

3096 ILS values indicate the proportion of papers that recommend categories within our
3097 taxonomy of all 21 studies. The most highly recommended categories from our
3098 SMS fall under the Usage Description dimension. The majority (0.71) of studies

Table 7.3: Weighted ILS Scoring.

Dimension	GCV	AWS	ACV	Max
[A] Usage Description	2.64 (60%)	3.10 (71%)	3.02 (69%)	4.38
[B] Design Rationale	0.79 (55%)	0.95 (67%)	0.95 (67%)	1.43
[C] Domain Concepts	0.33 (54%)	0.14 (23%)	0.43 (69%)	0.62
[D] Support Artefacts	0.24 (31%)	0.52 (69%)	0.50 (66%)	0.76
[E] Documentation Presentation	1.05 (79%)	1.05 (79%)	0.98 (73%)	1.33
Total	5.05 (59%)	5.76 (68%)	5.88 (69%)	8.52

Table 7.4: Weighted IPS Scoring.

Dimension	GCV	AWS	ACV	Max
[A] Usage Description	4.84 (57%)	5.26 (62%)	5.62 (66%)	8.48
[B] Design Rationale	1.78 (43%)	2.51 (61%)	2.51 (61%)	4.13
[C] Domain Concepts	0.92 (51%)	0.55 (31%)	1.43 (80%)	1.80
[D] Support Artefacts	0.96 (28%)	1.80 (53%)	1.85 (55%)	3.36
[E] Documentation Presentation	2.66 (70%)	2.66 (70%)	2.38 (63%)	3.79
Total	11.17 (52%)	12.79 (59%)	13.79 (64%)	21.56

³⁰⁹⁹ advocate for code snippets as a necessary piece in the API documentation puzzle
³¹⁰⁰ [A5]. While code snippets generally only reflect small portions of API functionality
³¹⁰¹ (limited to 15–30 LoC), this is complimented by step-by-step tutorials (0.57) that tie
³¹⁰² in multiple (disparate) components of API functionality, generally with some form
³¹⁰³ of screenshots, demonstrating the development of a non-trivial application using the
³¹⁰⁴ API step-by-step [A6]. The third highest category scored was also under the Usage
³¹⁰⁵ Description dimension, being low-level reference documentation at 0.52 [A2]. These
³¹⁰⁶ three categories were the only categories to be scored as majority categories (i.e.,
³¹⁰⁷ their scores were above 0.50). The fourth and fifth highest scores are an entry-
³¹⁰⁸ level purpose/overview of the API (0.48) that gives a brief motivation as to why a
³¹⁰⁹ developer should choose a particular API over another [B1] and consistency in the
³¹¹⁰ look and feel of the documentation throughout all of the API’s official documentation
³¹¹¹ (0.43) [E6].

³¹¹² 7.6.2 In-Practice Scores for Taxonomy Categories

³¹¹³ IPS values indicate the extent to which developers ‘agree’ with the statements made
³¹¹⁴ in our survey, as calculated using the SUS technique [50]. These values are generally
³¹¹⁵ greater than the ILS values, since they are ranked by all survey participants and are not
³¹¹⁶ a ratio of the 21 primary studies. Unlike ILS scores, 28 categories scored above 0.50.
³¹¹⁷ The highest dimension corroborates that of the ILS scores; within the top five ranked
³¹¹⁸ ILS scores, Usage Description categories feature four times. However, developers
³¹¹⁹ generally find limitations on what the APIs can and cannot provide the most useful,
³¹²⁰ at 0.94, which falls under the Design Rationale dimension [B7]. Following this, the

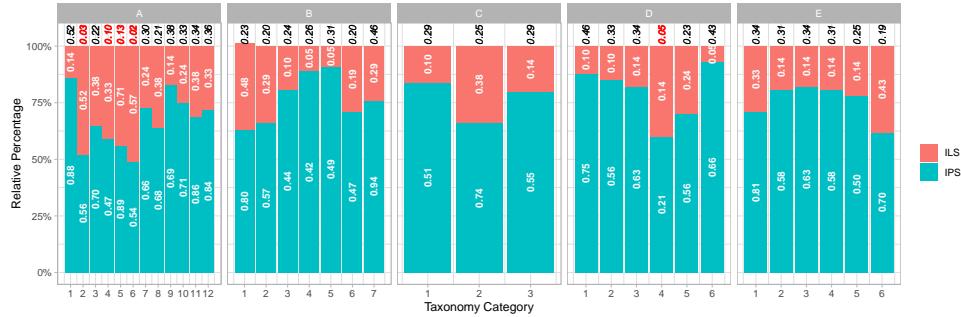


Figure 7.5: Comparison of ILS and IPS values for each category (grouped by dimensions) presented as a relative percentage.

Table 7.5: Labels assigned to ILS and IPS values

Description	Lower Score Bound	Upper Score Bound
<i>Not Helpful</i>	0.00	0.24
<i>Slightly Unhelpful</i>	0.25	0.49
<i>Slightly Helpful</i>	0.50	0.74
<i>Very Helpful</i>	0.75	1.00

3121 code-snippets [A5] is highly ranked (as per the ILS values) with developers agreeing
 3122 that code-snippets should be included in most API documentation. Quick-start
 3123 guides [A1] are the next most-useful category that developers advocate for, reported
 3124 at 0.88. Following this, the instructions on how to install the API or begin using the
 3125 API, its release cycle, and frequently it is updated [A11] is also important, ranking
 3126 fourth at 0.86. Lastly, error definitions describing how developers can address
 3127 problems [A12] were scored at 0.84.

3128 7.6.3 Contrasting In-Literature to In-Practice Scores

3129 Figure 7.5 highlights the relative percentage of each ILS and IPS value for all
 3130 subcategories, thereby indicating the relative agreement between the two. In this
 3131 graph, an ILS and/or IPS core approaching a relative percentage of 50% indicates
 3132 equal agreement whereby both developer's and literary references share a similar
 3133 distribution of recommendation agreement. Italicised labels above each column
 3134 indicates the standard deviation between the ILS and IPS values, where red labels
 3135 indicated a standard deviation less than 0.15 (i.e., developers and literature agree to
 3136 the values to a similar extent).

3137 Where the standard deviation between ILS and IPS values is less than 0.015
 3138 (as indicated by red labels above each column in Figure 7.5), then there is strong
 3139 alignment between both scores. However, of all 34 categories, only five cases of this
 3140 occur. Developers agree to the academic works that make the recommendations *to*
 3141 *the same relative proportion* as per the labels assigned in Table 7.5:

- 3142 • Having email addresses or phone numbers listed within an API is generally

3143 not helpful at all [D4],

3144 • Introspecting the source code comments of an API is only somewhat helpful
3145 [A4],

3146 • Low-level reference documentation with all objects and methods (etc.) docu-
3147 mented is slightly helpful [A2],

3148 • Following step-by-step tutorials are also slightly helpful [A6],

3149 • Code snippets are the most helpful [A5].

3150 The remaining categories in the dimension do not share strong association be-
3151 tween both developer opinions and the number of papers producing recomme-
3152 ndations. Due to the disparity between these ILS and IPS values, we do not report on
3153 their utility.

3154 7.6.4 Triangulating ILS and IPS with Computer Vision

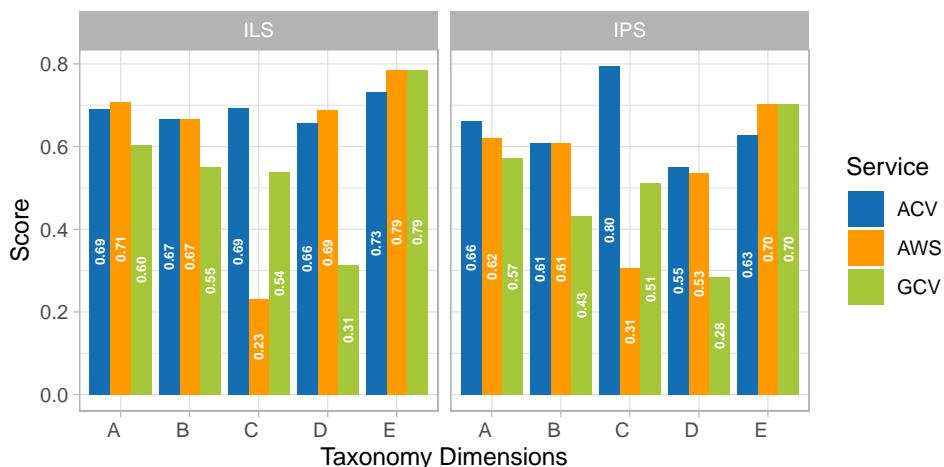


Figure 7.6: Comparison of the weighted ILS and IPS values for the three CVSs assessed.

3155 When applied in the context of CVSs, we see that Azure Computer Vision
3156 (ACV) and Amazon Rekognition (AWS) are better documented than Google Cloud
3157 Vision (GCV), particularly in Design Rationale and Usage Description. Figure 7.6
3158 highlights that Azure Computer Vision is especially well documented in Domain
3159 Concepts when measured using the weighted ILS and has the highest score of all
3160 services and dimensions. It is evident that Google Cloud Vision needs improved
3161 Design Rationale documentation and further Support Artefacts would be helpful.
3162 Generally speaking, Google Cloud Vision is less ‘complete’ than other services,
3163 except in Domain Concepts documentation and its Documentation Presentation.

3164 In the context of CVSs, IPS values share a similar distribution to ILS val-
3165 ues. Notably, in-practice, it seems developers prefer the documentation of Amazon
3166 Rekognition compared to the the in-literature weighted scoring of Azure Computer
3167 Vision (Figure 7.6). Except in the case of documenting Domain Concepts, Amazon

3168 Rekognition scores slightly higher than Azure Computer Vision except for the De-
3169 sign Rationale documentation where it is equal. Similar to the ILS scoring, Google
3170 Computer Vision has low compliance to the recommendations proposed, except in
3171 its Documentation Presentation.

3172 7.6.5 Areas of Improvement for CVS Documentation

3173 Triangulating the taxonomy developed from literary sources, the developer survey
3174 on this taxonomy to understand its efficacy in-practice, and applying the taxonomy to
3175 the CVS domain, we are able to assess the key areas of improvement in this domain.

3176 For this assessment, we select the ILS or IPS values for categories that are
3177 considered either somewhat or very helpful (i.e., a score greater than 0.50). We then
3178 match these against categories that are found to be partially or not present within
3179 each service. In total, we found 12 categories where improvements can be made
3180 across all dimensions except Documentation Presentation, detailed below .

3181 7.6.5.1 Issues regarding Usage Description

3182 **Quick-start guides [A1]:** Quick-start guides should provide a short tutorial that
3183 allows programmers to pick up the basics of an API in a programming language of
3184 their choice. For the services assessed, each offer various client SDKs (e.g., as Java
3185 or Python client libraries). Google Cloud Vision and Azure Computer Vision offer
3186 quick-start guides [365, 382] in which sets of articles target various SDKs or are
3187 client-agnostic with code snippets that can be changed to the client language/SDK
3188 of the developer's choice. Amazon Rekognition offers exercises in setting up the
3189 AWS SDK and using the command-line interface to interact with image analysis
3190 components [347], however this is client-agnostic nor does it provide details in how
3191 to get started with using the client SDKs.

 **Suggested improvement:** Ensure tutorials detail all client-libraries and how developers can produce a minimum working example using the service on their own computer using that client library. For each SDK offered, there should be details on how to install, authenticate and use a component using local data. For example, this may be as simple as using the service to determine if an image of a dog contains the label 'dog'.

3192 **Step-by-step tutorials [A6]:** Google Cloud Vision offers tutorials limited to one
3193 component. These do not sufficiently demonstrate how to combine *multiple compo-*
3194 *nents* of the API together and how developers should integrate it with a differ-
3195 ent platform, which a good step-by-step tutorial should detail. The official AWS
3196 Machine Learning blog [344] provides extensive tutorials (in some cases, with a
3197 suggested tutorial completion time of over an hour) that integrate multiple Amazon
3198 Rekognition components with other AWS components. Microsoft provide tutori-
3199 als [380, 385, 386] integrating multiple components within their service to mobile
3200 applications and the Azure platform.

 **Suggested improvement:** Ensure tutorials combine multiple components of the service together, are extensive, and require developers to spend a non-trivial amount of time to produce a basic application. For example, the tutorial may detail how to integrate the API into a smartphone application to achieve the following: (i) take a photo with the camera, (ii) detect if a person is within the image, (iii) analyse the visual features of the person.

3201 **Downloadable production-ready applications [A7]:** Microsoft provide a down-
3202 loadable application [384] that explores many components of the Azure Computer
3203 Vision API. The application is thoroughly documented with and also provides guid-
3204 ance on how to structure the architecture design of the program. While Rekognition
3205 and Google Cloud Vision also provide downloadable source code, they are largely
3206 under-documented, do not combine multiple components of the API together, and
3207 only use god-classes to handle all requests to the API [348, 367].

 **Suggested improvement:** Downloadable source code should be thoroughly docu-
mented, and should avoid the use of god-classes that demonstrate a single piece of the
service's functionality. Ideally, the architecture of a production-ready application should
be demonstrated to developers.

3208 **Understanding best-practices [A8]:** Google Cloud provides best-practices for its
3209 platform in both general and enterprise contexts [359, 368], but there is little advice
3210 provided to guide developers on how best to use Google Cloud Vision. Microsoft
3211 provides guidance on improving results of custom vision classifiers [381], but no
3212 further details on non-custom vision classifiers are found. We found the most detailed
3213 best-practices to be provided by Amazon Rekognition [346], which outlines more
3214 detailed strategies such as reducing data transfer by storing and referencing images
3215 on S3 Buckets or the attributes images should have in various scenarios (e.g., the
3216 angles of a person's face in facial recognition).

 **Suggested improvement:** Document best-practices for all major components of the
CVS. Guide developers on the types of input data that produce the best results, advisable
minimum image sizes and recommended file types, and suggest ways to overcome
limitations that improve usage and cost efficiency. Provide guidance in more than one
use case; give a range of scenarios that demonstrate different best practices for different
domains.

3217 **Exhaustive lists of all major API components [A9]:** Amazon provides a two-fold
3218 feature list that describes both the key features of Rekognition at a high-level [345]
3219 as well as a detailed, technical breakdown of each API operation provided within the
3220 service [343]. Microsoft also provide a list of high-level features that Azure Com-
3221 puter Vision can analyse [387] which provides hyperlinks to detailed descriptions of
3222 each feature. Google's Cloud Vision API provides a partial breakdown of the types
3223 of services provided, however this list is not fully complete, nor are there hyperlinks
3224 to more detailed descriptions of each of the features [369].

⌚ Suggested improvement: Document key features that the CV classifier can perform at a high level. This should be easy to find from the service's landing page. Each feature should be described with reference to more detailed descriptions of the feature's exact API endpoint and required inputs, outputs and possible errors.

3225 **Minimum system requirements and dependencies [A10]:** Although there is no
3226 dedicated webpage for this on any of the services investigated, there are listed
3227 dependencies for the client libraries in Google's and Azure's quick-start guides [365,
3228 378]. These may be embedded within the quick-start guide as developers are likely
3229 to encounter dependency issues when they first start using the API. We found it a
3230 challenge to discover similar documentation this in Amazon's documentation.

⌚ Suggested improvement: Any system requirements and dependency issues should be well-highlighted within the documentation's quick-start guide; developers are likely to encounter these issues within the early stages of using an API, and it is highly relevant to provide solutions to these issues within the quick-starts.

3231 **Installation and release cycle notes [A11]:** It is imperative that developers know
3232 what has changed between releases and how frequently the releases are exported.
3233 We found release notes for Amazon Computer Vision, although they are only major
3234 releases and have not been updated since 2017 [342] which does not account for
3235 evolution in the service's responses [76]. Google's and Microsoft's release notes are
3236 generally more frequently updated, therefore developers can get a sense of its release
3237 frequency [366, 383]. However, there are evolution issues that are not addressed.
3238 Installation instructions are detailed within Rekognition's developer guide, outlining
3239 how to sign up for an account, and install the AWS command-line interface [350].

⌚ Suggested improvement: Ensure release notes detail label evolution, including any new additional labels that may have been introduced within the service. Transparency around the changes made to the service should go beyond new features: document potential changes that may influence maintenance of a system using the CVS so that developers are aware of potential side-effects of upgrading to a newer release.

3240 7.6.5.2 Issues regarding Design Rationale

3241 **Limitations of the API [B7]:** The most detailed limitations documented were
3242 found on Rekognition's dedicated limitations page [349] that outlines functional
3243 limitations such as the maximum number of faces or words that can be detected
3244 in an image, the size requirements of images, and file type information. For the
3245 other services, functional limitations are generally found within each endpoint's API
3246 documentation, instead of within a dedicated page.

⌚ Suggested improvement: Document all functional limitations in a dedicated page that outline the maximum and minimum input requirements the classifier can handle. Documentation of the types of labels the service can provide is also desired.

3247 7.6.5.3 Issues regarding Domain Concepts

3248 Conceptual understanding of the API [C1]: Azure Computer Vision provides
‘concept’ pages describing the high-level concepts behind CV and where these
3249 functions are implemented within the APIs (e.g., [379]). We were unable to find
3250 similar conceptual documentation for the other services assessed.
3251

3252 *☞ Suggested improvement: Document the concepts behind CV; differentiate between
3253 foundational concepts such as object localisation, object recognition, facial localisation
3254 and facial analysis such that developers are able to make the distinction between them.
3255 Relate these concepts back to the API and provide references to where the APIs implement
3256 these concepts.*

3257 Definitions of domain-specific terminology [C2]: Terminologies relevant to ma-
3258 chine learning concepts powering these CVSs are well detailed within Google’s
3253 machine learning glossary [363], however few examples matching CV are imme-
3254 diately relevant. While this page is linked from the original Google Cloud Vision
3255 documentation, it may be too technical for application developers to grasp. A slightly
3256 better example of this is [387], where developers can understand CV terms in lay
3257 terms.
3258

3259 *☞ Suggested improvement: Current CVSs use a myriad of terminologies to refer to the
3260 same conceptual feature; for example, while Microsoft refers to object recognition as
3261 ‘image tagging’, Google refers to this as ‘label detection’. If a consolidation of terms
3262 is not possible, then CVSs should provide a glossary that provides synonyms for these
3263 terminologies so that developers can easily move between service providers without
3264 needing to relink terms back to concepts.*

3265 7.6.5.4 Issues regarding Support Artefacts

3266 Troubleshooting suggestions [D2]: The only troubleshooting tips found in our
3267 analysis were in Rekognition’s video service [351]. Further detailed instances of
3268 these troubleshooting tips could be expanded to non-video issues. For instance,
3269 if developers upload ‘noisy’ images, how can they inform the system of a specific
3270 ontology to use or to focus on parts of the foreground or background of the image?
3271 These are suggestions which we have proposed in prior work [76] that do not seem
3272 to be documented.

3273 *☞ Suggested improvement: Ensure troubleshooting tips provide advice for testing against
3274 different types of valid input images.*

3275 Diagrammatic overview of the API [D3]: None of the CVSs provide any overview
3276 of the API in terms of the features and processing steps on how they should be used.
3277 For instance, pre-processing and post-processing of input and response data should
3278 be considered and an understanding of how this fits into the ‘flow’ of an application
3279 highlighted. Moreover, no UML diagrams could be found.
3280

⇨ **Suggested improvement:** Provide diagrams illustrating the service within context of use, such as how it can be integrated with other service features or how a specific API endpoint may be used within a client application. Consider integrating interactive UML diagrams so that developers can easily explore various aspects of the documentation in a visual perspective.

3272 7.7 Threats to Validity

3273 7.7.1 Internal Validity

3274 Threats to *internal validity* represent internal factors of our study which affect
3275 concluded results. Kitchenham and Charters' guidelines on producing systematic
3276 reviews [168] suggest that researchers conducting reviews should discuss the review
3277 protocol, inclusion decisions, data extraction with a third party. Within this study,
3278 we discussed our protocols with other researchers within our research group and
3279 utilised test-retest reliability. Further assessments into reliability would involve an
3280 assessment of the review and extraction processes, which can be investigated using
3281 inter-rater reliability measures. Guidelines suggested by Garousi and Felderer [113]
3282 describe methods for independent analysis and conflict resolution could help resolve
3283 this.

3284 As stated in Section 7.3.2, we utilised a systematic SE taxonomy development
3285 method by Usman et al. [311]. Two additional taxonomy validation approaches
3286 proposed by Usman et al. were not considered in our work: benchmarking and
3287 orthogonality demonstration. To our knowledge, there are no other studies that
3288 classify existing API knowledge studies into a structured taxonomy, and therefore
3289 we are unable to benchmark our taxonomy against others. We would encourage the
3290 research community to conduct a replication of our work and investigate whether
3291 our taxonomy classification approaches are replicable to ensure that categories are
3292 reliable and the dimensions fit the objectives of the taxonomy. Moreover, we did
3293 not investigate orthogonality demonstration as our primary goals for this work were
3294 to investigate the efficacy of the taxonomy by practitioners and in-practice, with
3295 reference to our wider research area of intelligent CVSSs. Therefore, we solely
3296 adopted the utility demonstration approach in two detailed experiments (Sections 7.5
3297 and 7.6) to analyse the efficacy of our taxonomy and identify potential improvements
3298 for these services' API documentation.

3299 7.7.2 External Validity

3300 Threats to *external validity* concern the generalisation of our observations. Our
3301 systematic mapping study has used a broad range of sources however not all papers
3302 contributing to API documentation may have been found or captured within the
3303 taxonomy. While we attempted to include as many papers as we could find in our
3304 study, some papers may have been filtered out due to our exclusion criteria. For
3305 example, there are studies we found that were excluded as they were not written in
3306 English, and these excluding factors may alter our conclusions, introducing conflict-

3307 ing recommendations. However, given the consistency of these trends within the
3308 studies that were sourced, we consider this a low likelihood.

3309 Documentation of web APIs are non-static, and may evolve using contributions
3310 from both official sources and the developer community (e.g., via GitHub). We
3311 downloaded the three service’s API documentation in March of 2019—it is highly
3312 likely that new documentation may have been added since or modified since publi-
3313 cation. A recommendation to mitigate this would be to re-evaluate this study once
3314 intelligent CVSs have matured and become even more mainstream in developer
3315 communities.

3316 We also adopt research conducted in the field of questionnaire design, such as
3317 ensuring all scales are worded with labels [178] and have used a summatting rating
3318 scale [292] to address a specific topic of interest if people are to make mistakes in
3319 their response or answer in different ways at different times. This approach was
3320 also extended using the SUS methodology, in which positive and negative items
3321 were used—as multiple studies have shown [51, 272], this approach helps reduce
3322 poor-quality responses by minimising extreme responses and acquiescence biases.

3323 7.7.3 Construct Validity

3324 Threats to *construct validity* relates to the degree by which the data extrapolated
3325 in this study sufficiently measures its intended goals. Automatic searching was
3326 conducted in the SMS by choice of three popular databases (see Section 7.3.1).
3327 As a consequence of selecting multiple databases, duplicates were returned. This
3328 was mitigated by manually curating out all duplicate results from the set of studies
3329 returned. Additionally, we acknowledge that the lack manual searching of papers
3330 within particular venues may be an additional threat due to the misalignment of
3331 search query keywords to intended papers of inclusion. Thus, our conclusions are
3332 only applicable to the information we were able to extract and summarise, given the
3333 primary sources selected.

3334 While we have investigated the application of this taxonomy using a user study
3335 (Section 7.5.1), we would like to explore an observational study of developers
3336 to assess how improved and non-improved API documentation impacts developer
3337 productivity. The outcome of this work can help design a follow-up experiment,
3338 consisting of a comparative controlled study [278] that capture firsthand behaviours
3339 and interactions toward how software engineers approach using a CVS with and
3340 without our taxonomy applied. This can be achieved by providing ‘mock’ improved
3341 documentation with the suggested improvements included in this work. Such an ex-
3342 periment could recruit a sample of developers of varying experience (from beginner
3343 programmer to principal engineer) to complete a certain number of tasks under an
3344 observational, comparative controlled study, half of which will (a) develop using
3345 the improved ‘mock’ documentation, and the other half will (b) develop with the
3346 *as-is/existing* documentation. From this, we can compare if the framework makes
3347 improvements by capturing metrics and recording the observational sessions for
3348 qualitative analysis. Visual modelling can be adopted to analyse the qualitative data
3349 using matrices [87], maps and networks [275] as these help illustrate any causal, tem-

3350 poral or contextual relationships that may exist to map out the developer’s mindset
3351 and difference in approaching the two sets of designs of the same tasks.

3352 7.8 Conclusions & Future Work

3353 A good API document should facilitate a developer’s productivity, and is therefore
3354 associated to the quality of software produced; improving the quality of the docu-
3355 mentation of third-party APIs improves the quality of dependent software. However,
3356 there does not yet exist a consolidated taxonomy of key recommendations proposed
3357 by literature, and—more importantly—it is useful to know if what developers need
3358 *in-practice* differs to what documentation artefacts are anticipated by literature.
3359 Moreover, there has been little work on mapping the research produced in this space
3360 against the techniques used to arrive at the recommendations.

3361 This study prioritises which aspects of API documentation knowledge is both (i)
3362 suggested by literature, and (ii) is demanded *most* by developers. We conduct a
3363 SMS from a pool of 4,501 studies and identify 21 seminal studies. From this, we
3364 synthesise a taxonomy of the various documentation aspects that should improve
3365 API documentation quality. Furthermore, we also capture the most commonly used
3366 analysis techniques used in the academic literature. We then validate our taxonomy
3367 against developers to assess its efficacy with practitioners, and conduct a heuristic
3368 evaluation against three popular CVSs. We offer 12 detailed suggested improve-
3369 ments where these services currently have weaknesses, and where specifically they
3370 may be able to improve their documentation.

3371 Future extensions of our work may involve a restricted systematic literature
3372 review in API documentation artefacts, and many suggestions are further detailed
3373 in Section 7.7. Further, a review into the techniques of these primary studies may
3374 extend the mapping we conducted in this work, by evaluating the effectiveness of
3375 the various approaches used in each study and assessing these against the proposed
3376 conclusions of each study.

3377 The findings of our work provides a solid baseline for improving the documen-
3378 tation of non-deterministic software, such as CVSs. While our aim is to eventually
3379 improve the quality of API documentation, the ultimate goal is to improve the soft-
3380 ware engineer’s experience of non-deterministic IWSs. We hope the guidelines from
3381 this extensive study help both software developers and API providers alike by using
3382 our taxonomy as a go-to checklist for what should be considered in documenting any
3383 API.

CHAPTER 8

3384

3385

3386 Using a Facade Pattern to combine Computer Vision Services[†]

3387

3388 **Abstract** Intelligent computer vision services, such as Google Cloud Vision or Amazon
3389 Rekognition, are becoming evermore pervasive and easily accessible to developers to build
3390 applications. Because of the stochastic nature that ML entails and disparate datasets used in
3391 their training, the outputs from different computer vision services varies with time, resulting
3392 in low reliability—for some cases—when compared against each other. Merging multiple
3393 unreliable API responses from multiple vendors may increase the reliability of the overall
3394 response, and thus the reliability of the intelligent end-product. We introduce a novel
3395 methodology—inspired by the proportional representation used in electoral systems—to
3396 merge outputs of different intelligent computer vision API provided by multiple vendors.
3397 Experiments show that our method outperforms both naive merge methods and traditional
3398 proportional representation methods by 0.015 F-measure.

3399 8.1 Introduction

3400 With the introduction of intelligent web services (IWSs) that make machine learning
3401 (ML) more accessible to developers [256, 317], we have seen a large growth of
3402 intelligent applications dependent on such services [54, 117]. For example, consider
3403 the advances made in computer vision, where objects are localised within an image
3404 and labelled with associated categories. Cloud-based computer vision services
3405 (CVSs)—e.g., [341, 354, 358, 362, 371, 372, 376, 423]—are a subset of IWSs.
3406 They utilise ML techniques to achieve image recognition via a remote black-box
3407 approach, thereby reducing the overhead for application developers to understand
3408 how to implement intelligent systems from scratch. Furthermore, as the processing

[†]This chapter is originally based on T. Ohtake, A. Cummaudo, M. Abdelrazek, R. Vasa, and J. Grundy, “Merging intelligent API responses using a proportional representation approach,” in *Proceedings of the 19th International Conference on Web Engineering*. Daejeon, Republic of Korea: Springer, June 2019. DOI 10.1007/978-3-030-19274-7_28. ISBN 978-3-03-019273-0. ISSN 1611-3349 pp. 391–406. Terminology has been updated to fit this thesis.

3409 and training of the machine-learnt algorithms is offloaded to the cloud, developers
3410 simply send RESTful API requests to do the recognition. There are, however, inherit
3411 differences and drawbacks between traditional web services and IWSs, which we
3412 describe with the motivating scenario below.

3413 **8.1.1 Motivating Scenario: Intelligent vs Traditional Web Services**

3414 An application developer, Tom, wishes to develop a social media Android and iOS
3415 app that catalogues photos of him and his friends, common objects in the photo,
3416 and generates brief descriptions in the photo (e.g., all photos with his husky dog,
3417 all photos on a sunny day etc.). Tom comes from a typical software engineering
3418 background with little knowledge of computer vision and its underlying concepts.
3419 He knows that intelligent computer vision web APIs are far more accessible than
3420 building a computer vision engine from scratch, and opts for building his app using
3421 these cloud services instead.

3422 Based on his experiences using similar cloud services, Tom would expect consistency
3423 of the results from the same API and different APIs that provide the same (or
3424 similar) functionality. As an analogy, when Tom writes the Java substring method
3425 "doggy".substring(0, 2), he expects it to be the same result as the Swift equivalent
3426 "doggy".prefix(3). Each and every time he interacts with the substring
3427 method using either API, he gets "dog" as the response. This is because Tom is
3428 used to deterministic, rule-driven APIs that drive the implementation behind the
3429 substring method.

3430 Tom's deterministic mindset results in three key differentials between a traditional
3431 web service and an IWS:

3432 **(1) Given similar input, results differ between similar IWSs.** When Tom
3433 interacts with the API of an IWS, he is not aware that each API provider trains
3434 their own, unique ML model, both with disparate methods and datasets. These
3435 IWSs are, therefore, nondeterministic and data-driven; input images—even
3436 if they contain the same conceptual objects—often output different results.
3437 Contrast this to the substring example, where the rule-driven implementation provides
3438 certainty to the results, this is not guaranteed for IWSs. When Tom interacts with an IWS,
3439 he is likely to receive different results for the same input image.

3440 **(2) Intelligent responses are not certain.** When Tom interprets the response
3441 object of an IWS, he finds that there is a ‘confidence’ value or ‘score’. This
3442 is because the ML models that power IWSs are inherently probabilistic and
3443 stochastic; any insight they produce is purely statistical and associational [242].
3444 Unlike the substring example, where the rule-driven implementation provides
3445 certainty to the results, this is not guaranteed for IWSs. For example, a picture
3446 of a husky breed of dog is misclassified as a wolf. This could be due to
3447 adversarial examples [299] that ‘trick’ the model into misclassifying images
3448 when they are fully decipherable to humans. It is well-studied that such
3449 adversarial examples exist in the real world unintentionally [99, 179, 245].

3450 **(3) Intelligent APIs evolve over time.** Tom may find that responses to processing
3451 an image may change over time; the labels he processes in testing may evolve

3452 and therefore differ to when in production. In traditional web services, evo-
3453 lution in responses is slower, generally well-communicated, and usually rare
3454 (Tom would always expect "dog" to be returned in the substring example).
3455 This has many implications on software systems that depend on these APIs,
3456 such as confidence in the output and portability of the solution. Currently, if
3457 Tom switches from one API provider to another, or if he doesn't regularly test
3458 his app in production, he may begin to see a very different set of labels and
3459 confidence levels.

3460 **8.1.2 Research Motivation**

3461 These drawbacks bring difficulties to the intended API users like Tom. We identify a
3462 gap in the software engineering literature regarding such drawbacks, including: lack
3463 of best practices in using IWSs; assessing and improving the reliability of APIs for
3464 their use in end-products; evaluating which API is suitable for different developer
3465 and application needs; and how to mitigate risk associated with these APIs. We
3466 focus on improving reliability of CVSs for use in end-products. The key research
3467 questions in this paper are:

- 3468 **RQ1:** Is it possible to improve reliability by merging multiple CVS results?
3469 **RQ2:** Are there better algorithms for merging these results than currently in
3470 use?

3471 Previous attempts at overcoming low reliability include triple-modular redundan-
3472 cacy [195]. This method uses three modules and decides output using majority
3473 rule. However, in CVSs, it is difficult to apply majority rule: these APIs respond with
3474 a list of labels and corresponding scores. Moreover, disparate APIs ordinarily output
3475 different results. These differences make it hard to apply majority rule because the
3476 type of outputs are complex and disparate APIs output different results for the same
3477 input. Merging search results is another technique to improve reliability [286]. It
3478 normalises scores of different databases using a centralised sample database. Nor-
3479 malising scores makes it possible to merge search results into a single ranked list.
3480 However, search responses are disjoint, whereas they are not in the context of most
3481 CVSs.

3482 In this paper, we introduce a novel method to merge responses of CVSs, using
3483 image recognition APIs endpoints as our motivating example. Section 8.2 describes
3484 naive merging methods and requirements. Section 8.3 gives insights into the struc-
3485 ture of labels. Section 8.4 introduces our method of merging computer vision labels.
3486 Section 8.5 compares precision and recall for each method. Section 8.6 presents
3487 conclusions and future work.

3488 **8.2 Merging API Responses**

3489 Image recognition APIs have similar interfaces: they receive a single input (image)
3490 and respond with a list of labels and associated confidence scores. Similarly, other
3491 supervised-AI-based APIs do the same (e.g., detecting emotions from text and
3492 natural language processing [373, 424]). It is difficult to apply majority rule on such

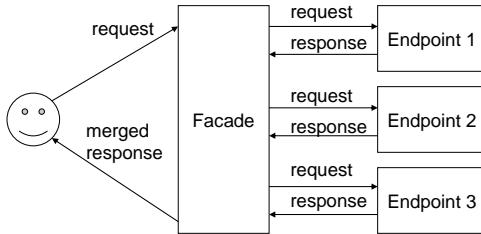


Figure 8.1: The user sends a request to the facade; this request is propagated to the relevant APIs. Responses are merged by the facade and returned back to the user.

disparate, complex outputs. While the outputs by *multiple* AI-based API endpoints is different and complex, the general format of the output is the same: a list of labels and associated scores.

8.2.1 API Facade Pattern

To merge responses from multiple APIs, we introduce the notion of an API facade. It is similar to a metasearch engine, but differs in their external endpoints. The facade accepts the input from one API endpoint (the facade endpoint), propagates that input to all user-registered concrete (external) API endpoints simultaneously, then ‘merges’ outputs from these concrete endpoints before sending this merged response to the API user. We demonstrate this process in Figure 8.1.

Although the model introduces more time and cost overhead, both can be mitigated by caching results. On the other hand, the facade pattern provides the following benefits:

- **Easy to modify:** It requires only small modifications to applications, e.g., changing each concrete endpoint URL.
- **Easy to customise:** It merges results from disparate and concrete APIs according to the user’s preference.
- **Improves reliability:** It enhances reliability of the overall returned result by merging results from different endpoints.

8.2.2 Merge Operations

The API facade is applicable to many use cases. However, this paper focuses on APIs that output a list of labels and scores, as is the case for CVSs. Merge operations involve the mapping of multiple lists and associated scores, produced by multiple APIs, to just one list. For instance, a CVS receives a bowl of fruit as the input image and outputs the following:

[[‘apple’, 0.9], [‘banana’, 0.8]]

where the first item is the label and the second item is the score. Similarly, another computer vision API outputs the following for the same image:

[[‘apple’, 0.7], [‘cherry’, 0.8]].

3522 Merge operations can, therefore, merge these two responses into just one response.
3523 Naive ways of merging results could make use of *max*, *min*, and *average* operations
3524 on the confidence scores. For example, *max* merges results to:

3525 $\text{[['apple', 0.9], ['banana', 0.8], ['cherry', 0.8]]};$

3526 *min* merges results to:

3527 $\text{[['apple', 0.7]]};$

3528 and *average* merges results to:

3529 $\text{[['apple', 0.8], ['banana', 0.4], ['cherry', 0.4]]}.$

3530 However, as the object's labels in each result are natural language, the operations
3531 do not exploit the label's semantics when conducting label merging. To improve
3532 the quality of the merged results, we consider the ontologies of these labels, as we
3533 describe below.

3534 8.2.3 Merging Operators for Labels

3535 Merge operations on labels are *n*-ary operations that map R^n to R , where $R_i =$
3536 $\{(l_{ij}, s_{ij})\}$ is a response from endpoint i and contains pairs of labels (l_{ij}) and scores
3537 (s_{ij}). Merge operations on labels have the following properties:

- 3538 • *identity* defines that merging a single response should output same response
3539 (i.e., $R = \text{merge}(R)$ is always true);
- 3540 • *commutativity* defines that the order of operands should not change the result
3541 (i.e., $\text{merge}(R_1, R_2) = \text{merge}(R_2, R_1)$ is always true);
- 3542 • *reflexivity* defines that merging multiple same responses should output same
3543 response (i.e., $R = \text{merge}(R, R)$ is always true); and,
- 3544 • *additivity* defines that, for a specific label, the merged response should have
3545 higher or equal score for the label if a concrete endpoint has a higher score.
3546 Let $R = \text{merge}(R_1, R_2)$ and $R' = \text{merge}(R'_1, R_2)$ be merged responses. R_1 and
3547 R'_1 are same, except R'_1 has a higher score for label l_x than R_1 . The additive
3548 score property requires that R' score for l_x should be greater than or equal to
3549 R score for l_x .

3550 The *max*, *min*, and *average* operations in Section 8.2.2 follow each of these rules
3551 as all operations calculate the score by applying these operations on each score.

3552 8.3 Graph of Labels

3553 CSVs typically return lists of labels and their associated scores. In most cases, the
3554 label can be a singular word (e.g., 'husky') or multiple words (e.g., 'dog breed').
3555 Lexical databases, such as WordNet [215], can therefore be used to describe the
3556 ontology behind these labels' meanings. Figure 8.2 is an example of a graph of

Table 8.1: Statistics for the number of labels, on average, per service identified.

Endpoint	Average number of labels	Has synset	No synset
Amazon Rekognition	11.42 ± 7.52	10.74 ± 7.10 (94.0%)	0.66 ± 0.87
Google Cloud Vision	8.77 ± 2.15	6.36 ± 2.22 (72.5%)	2.41 ± 1.93
Azure Computer Vision	5.39 ± 3.29	5.26 ± 3.32 (97.6%)	0.14 ± 0.37

3557 labels and synsets. A synset is a grouped set of synonyms for a word. In this image,
 3558 we consider two fictional endpoints, endpoints 1–2. We label red nodes as labels
 3559 from endpoint 1, yellow nodes as labels from endpoint 2, and blue nodes as synsets
 3560 for the associated labels from both endpoints. As actual graphs are usually more
 3561 complex, Figure 8.2 is a simplified graph to illustrate the usage of associating labels
 3562 from two concrete sources to synsets.

3563 8.3.1 Labels and synsets

3564 The number of labels depends on input images and concrete API endpoints used.
 3565 Table 8.1 and Figure 8.3 show how many labels are returned, on average per image,
 3566 from Google Cloud Vision [362], Amazon Rekognition [341] and Azure Computer
 3567 Vision [376] image recognition APIs. These statistics were calculated using 1,000
 3568 images from Open Images Dataset V4 [364] Image-Level Labels set.

3569 Labels from Amazon and Microsoft tend to have corresponding synsets, and
 3570 therefore these endpoints return common words that are found in WordNet. On the
 3571 other hand, Google’s labels have less corresponding synsets: for example, labels
 3572 without corresponding synsets are car models and dog breeds.¹

3573 8.3.2 Connected Components

3574 A connected component (CC) is a subgraph in which there are paths between any
 3575 two nodes. In graphs of labels and synsets, CCs are clusters of labels and synsets
 3576 with similar semantic meaning. For instance, there are two CCs in Figure 8.2. CC 1
 3577 in Figure 8.2 has ‘beverage’, ‘dessert’, ‘chocolate’, ‘hot chocolate’,
 3578 ‘drink’, and ‘food’ labels from the red first endpoint and ‘coffee’, ‘hot
 3579 chocolate’, ‘drink’, ‘caffeine’, and ‘tea’ labels from the yellow second
 3580 endpoint. Therefore, these labels are related to ‘drink’. On the other hand, CC 2
 3581 in Figure 8.2 has ‘cup’ and ‘coffee cup’ labels from the first red endpoint and
 3582 ‘cup’, ‘coffee cup’, and ‘tableware’ labels from the yellow second endpoint.
 3583 These labels are, therefore, related to ‘cup’.

3584 Figure 8.4 shows a distribution of number of CCs for the 1,000-image label
 3585 detections on Amazon Rekognition, Google Cloud Vision, and Azure Computer
 3586 Vision APIs. The average number of CCs is 9.36 ± 3.49 . The smaller number of
 3587 CCs means that most of labels have similar meanings, while a larger value means
 3588 that the labels are more disparate.

¹We noticed from our upload of 1,000 images that Google tries to identify objects in greater detail.

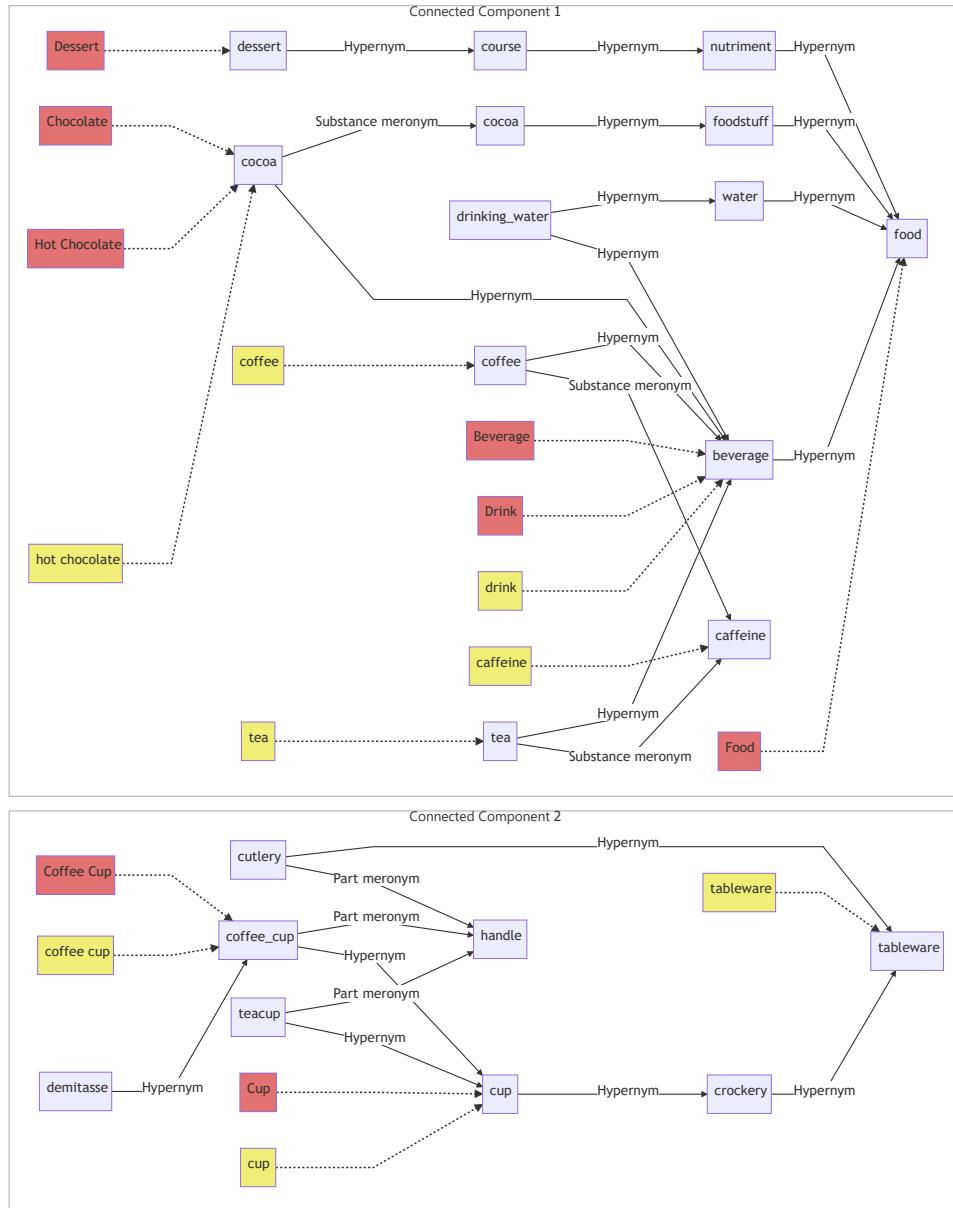


Figure 8.2: Graph of labels from two concrete endpoints (red and yellow) and their associated synsets related to both words (blue).

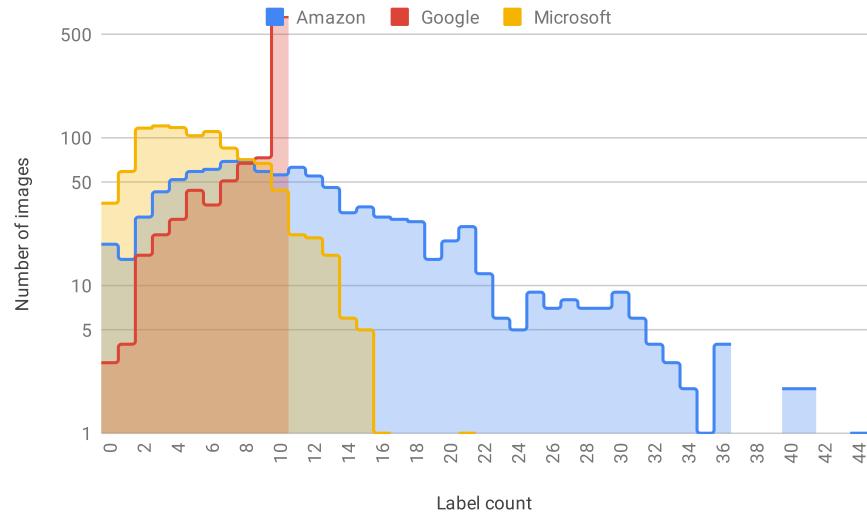


Figure 8.3: Number of labels responded from our input dataset to three concrete APIs assessed.

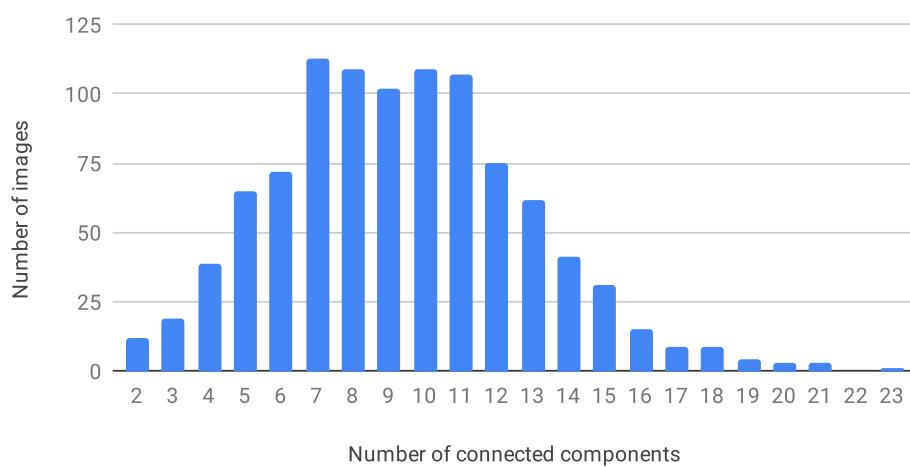


Figure 8.4: Number of connected components compared to the number of images.

3589 8.4 API Results Merging Algorithm

3590 Our proposed algorithm to merge labels consists of four parts: (1) mapping labels to
3591 synsets, (2) deciding the total number of labels, (3) allocating the number of labels
3592 to CCs, and (4) selecting labels from CCs.

3593 8.4.1 Mapping Labels to Synsets

3594 Labels returned in CVS responses are words (in natural language) that do not always
3595 identify their intended meanings. For instance, a label *orange* may represent the
3596 fruit, the colour, or the name of the longest river in South Africa. To identify the
3597 actual meanings behind a label, our facade enumerates all synsets corresponding to
3598 labels. It then finds the most likely synsets for labels by traversing WordNet links.
3599 For instance, if an API endpoint outputs the ‘*orange*’ and ‘*lemon*’ labels, the
3600 facade regards ‘*orange*’ as a related synset word of ‘*fruit*’. If an API endpoint
3601 outputs ‘*orange*’ and ‘*water*’ labels, the facade regards ‘*orange*’ as a ‘*river*’.

3602 8.4.2 Deciding Total Number of Labels

3603 The number of labels in responses from endpoints vary as described in Section 8.3.1.
3604 The facade decides the number of merged labels using the numbers of labels from
3605 each endpoint. We formulate the following equation to calculate the number of
3606 labels:

$$\min_i(|R_i|) \leq \frac{\sum_i |R_i|}{n} \leq \max_i(|R_i|) \leq \sum_i |R_i|$$

3607 where $|R|$ is number of labels and scores in response, and n is number of endpoints.
3608 In case of naive operations in Section 8.2.2, the following is true:

$$\begin{aligned} |\text{merge}_{\max}(R_1, \dots, R_n)| &\leq \min_i(|R_i|) \\ \max_i(|R_i|) &\leq |\text{merge}_{\min}(R_1, \dots, R_n)| \leq \sum_i |R_i| \\ \max_i(|R_i|) &\leq |\text{merge}_{\text{average}}(R_1, \dots, R_n)| \leq \sum_i |R_i|. \end{aligned}$$

3609 The proposal uses $\lfloor \sum_i |R_i| / n \rfloor$ to conform to the necessary condition described in
3610 Section 8.4.3.

3611 8.4.3 Allocating Number of Labels to Connected Components

3612 The graph of labels and synsets is then divided into several CCs. The facade decides
3613 how many labels are allocated for each CC. For example, in Figure 8.5, there are
3614 three CCs, where square-shaped nodes are labels in responses from endpoints. Text
3615 within these label nodes describe which endpoint outputs the label and score, for
3616 instance, “L-1a, 0.9” is label *a* from endpoint *L* with a score 0.9. Circle-shaped nodes
3617 represent synsets, where the edges between the label and synset nodes indicate the
3618 relationships between them. Edges between synsets are links in WordNet.

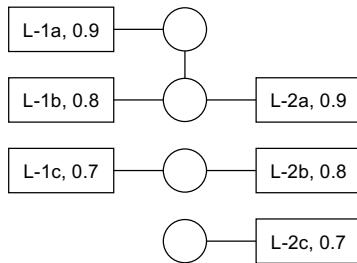


Figure 8.5: Allocation to connected components.

Allegorically, allocating the number of labels to CCs is similar to proportional representation in a political voting system, where CCs are the political parties and labels are the votes to a party. Several allocation algorithms are introduced in proportional representation, for instance, the D'Hondt and Hare-Niemeyer methods [224]. However, there are differences from proportional representation in the political context. For label merging, labels have scores and origin endpoints and such information may improve the allocation algorithm. For instance, CCs supported with more endpoints should have a higher allocation than CCs with fewer endpoints, and CCs with higher scores should have a higher allocation than CCs with lower scores. We introduce an algorithm to allocate the number of labels to CCs. This allocates more to a CC with more supporting endpoints and higher scores. The steps of the algorithm are:

- 3631 **Step I.** Sort scores separately for each endpoint.
- 3632 **Step II.** If all CCs have an empty score array or more, remove one, and go to Step II.
- 3633
- 3634 **Step III.** Select the highest score for each endpoint and calculate product of highest scores.
- 3635
- 3636 **Step IV.** A CC with the highest product score receives an allocation. This CC removes every first element from the score array.
- 3637
- 3638 **Step V.** If the requested number of allocations is complete, then stop allocation. Otherwise, go to Step II.
- 3639

Tables 8.2 to 8.5 are examples of allocation iterations. In Table 8.2, the facade sorts scores separately for each endpoint. For instance, the first CC in Figure 8.5 has scores of 0.9 and 0.8 from endpoint 1 and 0.9 from endpoint 2. All CCs have a non-empty score array or more, so the facade skips Step II. The facade then picks the highest scores for each endpoint and CC. CC 1 has the largest product of highest scores and receives an allocation. In Table 8.3, the first CC removes every first score in its array as it received an allocation in Table 8.2. In this iteration, the second CC has largest product of scores and receives an allocation. In Table 8.4, the second CC removes every first score in its array. At Step II, all the three CCs have an empty array. The facade removes one empty array from each CC. In Table 8.5, the first CC receives an allocation. The algorithm is applicable if total number of allocation is

Table 8.2: Allocation iteration 1.

Scores	Highest	Product	Allocated
[0.9, 0.8], [0.9]	[0.9, 0.9]	0.81	0+1
[0.7], [0.8]	[0.7, 0.8]	0.56	0
[], [0.7]	[N/A, 0.7]	N/A	0

Table 8.4: Allocation iteration 3.

Scores	Highest	Product	Allocated
[0.8], []	—	—	1
[], []	—	—	1
[], [0.7]	—	—	0

Table 8.3: Allocation iteration 2.

Scores	Highest	Product	Allocated
[0.8], []	[0.8, N/A]	N/A	1
[0.7], [0.8]	[0.7, 0.8]	0.56	0+1
[], [0.7]	[N/A, 0.7]	N/A	0

Table 8.5: Allocation iteration 4.

Scores	Highest	Product	Allocated
[0.8]	[0.8]	0.8	1+1
[]	[N/A]	N/A	1
[0.7]	[0.7]	0.7	0

3651 less than or equal to $\max_i(|R_i|)$ as scores are removed in Step II. The condition is a
3652 necessary condition.

3653 8.4.4 Selecting Labels from Connected Components

3654 For each CC, the facade applies the *average* operator from Section 8.2.2 and takes
3655 labels with n -highest scores up to allocation, as per Section 8.4.3.

3656 8.4.5 Conformance to properties

3657 Section 8.2.3 defines four properties: identity, commutativity, reflexivity, and additivity.
3658 Our proposed method conforms to these properties:

- 3659 • *identity*: the method outputs same result if there is one response;
- 3660 • *commutativity*: the method does not care about ordering of operands;
- 3661 • *reflexivity*: the allocations to CCs are same to number of labels in CCs; and
- 3662 • *additivity*: increases in score increases or does not change the allocation to
3663 the corresponding CC.

3664 8.5 Evaluation

3665 8.5.1 Evaluation Method

3666 To evaluate the merge methods, we merged CVS results from three representative
3667 image analysis API endpoints and compared these merged results against human-
3668 verified labels. Images and human-verified labels are sourced from 1,000 randomly-
3669 sampled images from the Open Images Dataset V4 [364] Image-Level Labels test
3670 set.

3671 The first three rows in Table 8.7 are the evaluation of original responses from
3672 each API endpoint. Precision, recall, and F-measure in Table 8.7 do not reflect
3673 actual values: for instance, it appears that Google performs best at first glance, but
3674 this is mainly because Google’s labels are similar to that of the Open Images label
3675 set.

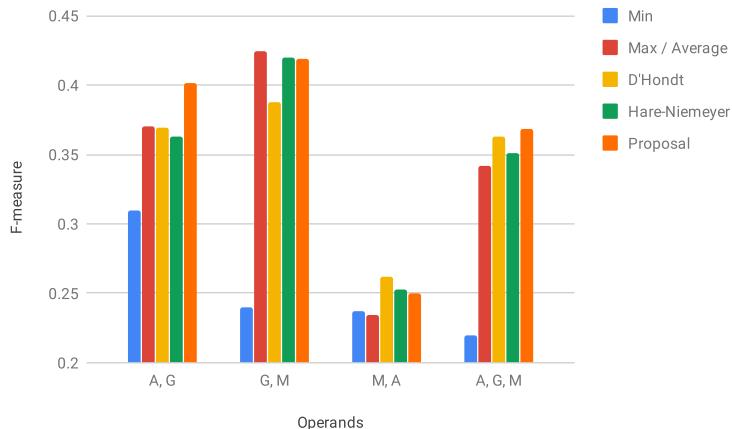


Figure 8.6: F-measure comparison.

3676 The Open Images Dataset uses 19,995 classes for labelling. The human-verified
 3677 labels for the 1,000 images contain 8,878 of these classes. Table 8.6 shows the
 3678 correspondence between each service's labels and the Open Images Dataset classes.
 3679 For instance, Amazon Rekognition outputs 11,416 labels in total for 1,000 images.
 3680 There are 1,409 unique labels in 11,416 labels. 1,111 labels out of 1,409 can be
 3681 found in Open Images Dataset classes. Rekognition's labels matches to Open Images
 3682 Dataset classes at 78.9% ratio, while Google has an outstanding matched percentage
 3683 of 94.1%. This high match is likely due to Google providing both Google Cloud
 3684 Vision and the Open Images Dataset—it is likely that they are trained on the same
 3685 data and labels. An endpoint with higher matched percentage has a more similar
 3686 label set to the Open Images Dataset classes. However, a higher matched percentage
 3687 does not mean imply *better quality* of an API endpoint; it will increase apparent
 3688 precision, recall, and F-measure only.

3689 The true and false positive (TP/FP) label averages and the TP/FP ratio is shown
 3690 in Table 8.7. Where the TP/FP ratio is larger, the scores are more reliable, however
 3691 it is possible to increase the TP/FP ratio by adding more false labels with low scores.
 3692 On the other hand, it is impossible to increase F-measure intentionally, because
 3693 increasing precision will decrease recall, and vice versa. Hence, the importance of
 3694 the F-measure statistic is critical for our analysis.

3695 Let R_A , R_G , and R_M be responses from Amazon Rekognition, Google Cloud
 3696 Vision, and Microsoft's Azure Computer Vision, respectively. There are four sets of
 3697 operands, i.e., (R_A, R_G) , (R_G, R_M) , (R_M, R_A) , and (R_A, R_G, R_M) . Table 8.7 shows
 3698 the evaluation of each operands set, Table 8.8 shows the averages of the four operands
 3699 sets, and Figure 8.6 shows the comparison of F-measure for each methods.

3700 8.5.2 Naive Operators

3701 Results of *min*, *max*, and *average* operators are shown in Tables 8.7 and 8.8 and Fig-
 3702 ure 8.6. The *min* operator is similar to *union* operator of set operation, and outputs
 3703 all labels of operands. The precision of the *min* operator is always greater than any

Table 8.6: Matching to human-verified labels.

Endpoint	Total	Unique	Matched	Matched %
Amazon Rekognition	11,416	1,409	1,111	78.9
Google Cloud Vision	8,766	2,644	2,487	94.1
Azure Computer Vision	5,392	746	470	63.0

Table 8.7: Evaluation results. A = Amazon Rekognition, G = Google Cloud Vision, M = Microsoft’s Azure Computer Vision.

Operands	Operator	Precision	Recall	F-measure	TP average	FP average	TP/FP ratio
A		0.217	0.282	0.246	0.848 ± 0.165	0.695 ± 0.185	1.220
G		0.474	0.465	0.469	0.834 ± 0.121	0.741 ± 0.132	1.126
M		0.263	0.164	0.202	0.858 ± 0.217	0.716 ± 0.306	1.198
A, G	Min	0.771	0.194	0.310	0.805 ± 0.142	0.673 ± 0.141	1.197
A, G	Max	0.280	0.572	0.376	0.850 ± 0.136	0.712 ± 0.171	1.193
A, G	Average	0.280	0.572	0.376	0.546 ± 0.225	0.368 ± 0.114	1.485
A, G	D’Hondt	0.350	0.389	0.369	0.713 ± 0.249	0.518 ± 0.202	1.377
A, G	Hare-Niemeyer	0.344	0.384	0.363	0.723 ± 0.242	0.527 ± 0.199	1.371
A, G	Proposal	0.380	0.423	0.401	0.706 ± 0.239	0.559 ± 0.190	1.262
G, M	Min	0.789	0.142	0.240	0.794 ± 0.209	0.726 ± 0.210	1.093
G, M	Max	0.357	0.521	0.424	0.749 ± 0.135	0.729 ± 0.231	1.165
G, M	Average	0.357	0.521	0.424	0.504 ± 0.201	0.375 ± 0.141	1.342
G, M	D’Hondt	0.444	0.344	0.388	0.696 ± 0.250	0.551 ± 0.254	1.262
G, M	Hare-Niemeyer	0.477	0.375	0.420	0.696 ± 0.242	0.591 ± 0.226	1.179
G, M	Proposal	0.414	0.424	0.419	0.682 ± 0.238	0.597 ± 0.209	1.143
M, A	Min	0.693	0.143	0.237	0.822 ± 0.201	0.664 ± 0.242	1.239
M, A	Max	0.185	0.318	0.234	0.863 ± 0.178	0.703 ± 0.229	1.228
M, A	Average	0.185	0.318	0.234	0.589 ± 0.262	0.364 ± 0.144	1.616
M, A	D’Hondt	0.271	0.254	0.262	0.737 ± 0.261	0.527 ± 0.223	1.397
M, A	Hare-Niemeyer	0.260	0.245	0.253	0.755 ± 0.251	0.538 ± 0.218	1.402
M, A	Proposal	0.257	0.242	0.250	0.769 ± 0.244	0.571 ± 0.205	1.337
A, G, M	Min	0.866	0.126	0.220	0.774 ± 0.196	0.644 ± 0.219	1.202
A, G, M	Max	0.241	0.587	0.342	0.857 ± 0.142	0.714 ± 0.210	1.201
A, G, M	Average	0.241	0.587	0.342	0.432 ± 0.233	0.253 ± 0.106	1.712
A, G, M	D’Hondt	0.375	0.352	0.363	0.678 ± 0.266	0.455 ± 0.208	1.492
A, G, M	Hare-Niemeyer	0.362	0.340	0.351	0.693 ± 0.260	0.444 ± 0.216	1.559
A, G, M	Proposal	0.380	0.357	0.368	0.684 ± 0.259	0.484 ± 0.200	1.414

Table 8.8: Average of the evaluation result.

Operator	Precision	Recall	F-measure	TP/FP ratio
Min	0.780	0.151	0.252	1.183
Max	0.266	0.500	0.344	1.197
Average	0.266	0.500	0.344	1.539
D’Hondt	0.361	0.335	0.346	1.382
Hare-Niemeyer	0.361	0.336	0.347	1.378
Proposal	0.358	0.362	0.360	1.289

3704 precision of operands, and the recall is always lesser than any precision of operands.
3705 *Max* and *average* operators are similar to *intersection* operator of set operations.
3706 Both operators output intersection of labels of operands and there is no clear relation
3707 to the precision and recall of operands. Since both operators have the same preci-
3708 sion, recall, and F-measure, Figure 8.6 groups them into one. The *average* operator
3709 performs well on the TP/FP ratio, where most of the same labels from multiple
3710 endpoints are TPs. In many cases of the four operand sets, all naive operators’
3711 F-measures are between F-measures of operands. None of naive operators therefore
3712 improve results by merging responses from multiple endpoints.

3713 8.5.3 Traditional Proportional Representation Operators

3714 There are many existing allocation algorithms in proportional representation, e.g.,
3715 the Niemeyer and Niemeyer method [224]. These methods may be replacements of
3716 those in Section 8.4.3. Other steps, i.e., Sections 8.4.1, 8.4.2 and 8.4.4, are the same
3717 as for our proposed technique. Tables 8.7 and 8.8 and Figure 8.6 show the result of
3718 these traditional proportional representation algorithms. Averages of F-measures by
3719 traditional proportional representation operators are almost equal to that of the *max*
3720 and *average* operators. It is worth noting that merging *M* and *A* responses results in
3721 a better F-measure than each F-measure of *M* and *A* individually. As these are not
3722 biased to human-verified labels, situations in the real-world usage should, therefore,
3723 be similar to the case of *M* and *A*. Hence, RQ1 is true.

3724 8.5.4 New Proposed Label Merge Technique

3725 As shown in Table 8.8, our proposed new method performs best in F-measure.
3726 Instead, the TP/FP ratio is less than *average*, the D’Hondt method, and Hare-
3727 Niemeyer method. As described in Section 8.5.1, we argue that F-measure is a
3728 more important measure than the TP/FP ratio (in this case). Therefore, RQ2 is
3729 true. Shown in Table 8.7, our proposed new method improves the results when
3730 merging *M* and *A* in non-biased endpoints. It is similar to traditional proportional
3731 representation operators, but does not perform as well. However, it performs better
3732 on other operand sets, and performs best overall as shown in Figure 8.6.

3733 8.5.5 Performance

3734 We used AWS EC2 m5.large instance (2 vCPUs, 2.5 GHz Intel Xeon, 8 GiB RAM);
3735 Amazon Linux 2 AMI (HVM), SSD Volume Type; Node.js 8.12.0. It takes 0.370
3736 seconds to merge responses from three endpoints. Computational complexity of the
3737 algorithm in Section 8.4.3 is $O(n^2)$, where n is total number of labels in responses.
3738 (The estimation assumes that the number of endpoints is a constant.) Complexity
3739 of Step I in Section 8.4.3 is $O(n \log n)$, as the worst case is that all n labels are from
3740 one single endpoint and all n labels are in one CC. Complexity of Step II to Step V
3741 is $O(n^2)$, as the number of CCs is less than or equal to n and number of iterations
3742 are less than or equal to n . As Table 8.1 shows, the averaged total number of three
3743 endpoints is 25.58. Most of time for merging is consumed by looking up WordNet

3744 synsets (Section 8.4.1). The API facade calls each APIs on actual endpoints in
3745 parallel. It takes about 5 seconds, which is much longer than 0.370 seconds taken
3746 for the merging of responses.

3747 8.6 Conclusions and Future Work

3748 In this paper, we propose a method to merge responses from CVSs. Our method
3749 merges API responses better than naive operators and other proportional represen-
3750 tation methods (i.e., D'Hondt and Hare-Niemeyer). The average of F-measure of
3751 our method marks 0.360; the next best method, Hare-Niemeyer, marks 0.347. Our
3752 method and other proportional representation methods are able to improve the F-
3753 measure from original responses in some cases. Merging non-biased responses
3754 results in an F-measure of 0.250, while original responses have an F-measure be-
3755 tween 0.246 and 0.242. Therefore, users can improve their applications' precision
3756 with small modification, i.e., by switching from a singular URL endpoint to a facade-
3757 based architecture. The performance impact by applying facades is small, because
3758 overhead in facades is much smaller than API invocation. Our proposal method
3759 conforms identity, commutativity, reflexivity, and additivity properties and these
3760 properties are advisable for integrating multiple responses.

3761 Our idea of a proportional representation approach can be applied to other IWSs.
3762 If the response of such a service is list consisting of an entity and score, and if there is a
3763 way to group entities, a proposal algorithm can be applied. The opposite approach is
3764 to improve results by inferring labels. Our current approach picks some of the labels
3765 returned by endpoints. IWSs are not only based on supervised ML—thus to cover a
3766 wide range of IWSs, it is necessary to classify and analyse each APIs and establish
3767 a method to improve results by merging. Currently graph structures of labels and
3768 synsets (Figure 8.2) are not considered when merging results. Propagating scores
3769 from labels could be used, losing the additivity property but improving results for
3770 users. There are many ways to propagate scores. For instance, setting propagation
3771 factors for each link type would improve merging and could be customised for users'
3772 preferences. It would be possible to generate an API facade automatically. APIs
3773 with the same functionality have same or similar signatures. Machine-readable API
3774 documentation, for instance, OpenAPI Specification, could help a generator to build
3775 an API facade.

CHAPTER 9

3776

3777

3778 Threshy: Supporting Safe Usage of Intelligent Web Services[†]

3779

3780 **Abstract** Increased popularity of ‘intelligent’ web services provides end-users with machine-
3781 learnt functionality at little effort to developers. However, these services require a decision
3782 threshold to be set which is dependent on problem-specific data. Developers lack a systematic
3783 approach for evaluating intelligent services and existing evaluation tools are predominantly
3784 targeted at data scientists for pre-development evaluation. This paper presents a workflow
3785 and supporting tool, Threshy, to help *software developers* select a decision threshold suited
3786 to their problem domain. Unlike existing tools, Threshy is designed to operate in multiple
3787 workflows including pre-development, pre-release, and support. Threshold configuration
3788 files exported by Threshy can be integrated into client applications and monitoring infrastruc-
3789 ture. Demo: <https://bit.ly/2YKeYhE>.

3790 9.1 Introduction

3791 Machine learning algorithm adoption is increasing in modern software. End users
3792 routinely benefit from machine-learnt functionality through personalised recom-
3793 mendations [71], voice-user interfaces [219], and intelligent digital assistants [43]. The
3794 easy accessibility and availability of intelligent web services (IWSs)¹ is contribut-
3795 ing to their adoption. These IWSs simplify the development of machine learning
3796 solutions as they (i) do not require specialised machine learning expertise to build
3797 and maintain, (ii) abstract away infrastructure related issues associated with machine
3798 learning [11, 277], and (iii) provide web APIs for ease of integration.

3799 However, unlike traditional web services, the functionality of these *intelligent*

[†]This chapter is originally based on A. Cummaudo, S. Barnett, R. Vasa, and J. Grundy, “Threshy: Supporting Safe Usage of Intelligent Web Services,” 2020, Unpublished. Terminology has been updated to fit this thesis.

¹Such as Azure Computer Vision (<https://azure.microsoft.com/en-au/services/cognitive-services/computer-vision/>), Google Cloud Vision (<https://cloud.google.com/vision/>), or Amazon Rekognition (<https://aws.amazon.com/rekognition/>).

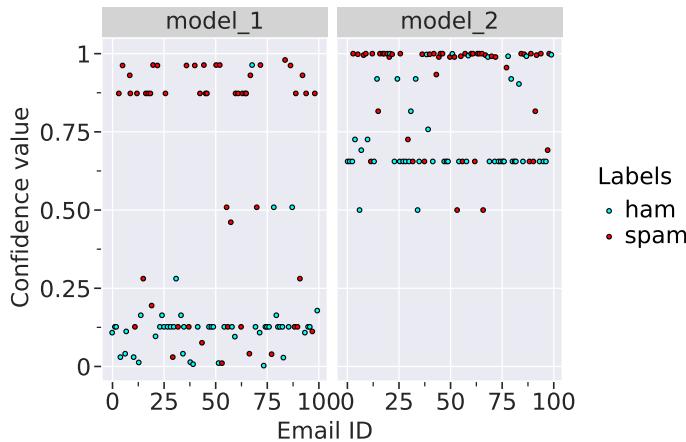


Figure 9.1: Predictions for 100 emails from two spam classifiers. Decision thresholds are classifier-dependent: a single threshold for both classifiers is *not* appropriate as ham emails are clustered at 0.12 (model_1) and at 0.65 (model_2). Developers must evaluate performance for *both* thresholds.

3800 services is dependent on a set of assumptions unique to machine learning [76].
 3801 These assumptions are based on the data used to train machine learning algorithms,
 3802 the choice of algorithm, and the choice of data processing steps—most of which
 3803 are not documented. For developers, these assumptions mean that the performance
 3804 characteristics of an intelligent service in any particular application problem domain
 3805 is not fully knowable. Intelligent services represent this uncertainty through a
 3806 confidence value associated with their predictions. Thus an evaluation procedure
 3807 must be followed as a part of using an intelligent service for an application.

3808 A typical evaluation process would involve a test data set (curated by the devel-
 3809 opers using the intelligent service) that is used to determine an appropriate threshold.
 3810 Choice of a decision threshold is a critical element of the evaluation procedure [129].
 3811 This is especially true for classification problems such as detecting if an image con-
 3812 tains cancer or identifying all of the topics in a document. Simple approaches
 3813 to selecting a threshold are often insufficient, as highlighted in Google’s machine
 3814 learning course: “*It is tempting to assume that [a] classification threshold should
 3815 always be 0.5, but thresholds are problem-dependent, and are therefore values that
 3816 you must tune.*”² As an example consider the predictions from two email spam
 3817 classifiers shown in Figure 9.1. The predicted safe emails, ‘ham’, are in two separate
 3818 clusters (a simple threshold set to approx. 0.2 for model 1 and 0.65 for model 2,
 3819 indicating that different decision thresholds may be required depending on the clas-
 3820 sifier. Also note that some emails have been misclassified; how many depends on
 3821 the choice of decision threshold. An appropriate threshold considers factors outside
 3822 algorithmic performance, such as financial cost and impact of wrong decisions. To
 3823 select an appropriate decision threshold, developers using intelligent services need
 3824 approaches to reason about and consider trade-offs between competing *cost fac-
 3825 tors*. These include impact, financial costs, and maintenance implications. Without

²See <https://bit.ly/36oMgWb>.

3826 considering these trade-offs, sub-optimal decision thresholds will be selected.

3827 The standard approach for tuning thresholds in classification problems involve
3828 making trade-offs between the number of false positives and false negatives using
3829 the receiver operating characteristic (ROC) curve. However, developers (i) need
3830 to realise that this trade-off between false positives and false negatives is a data
3831 dependent optimisation process [276], (ii) often need to develop custom scripts
3832 and follow a trial-and-error based approach to determine a threshold, (iii) must
3833 have appropriate statistical training and expertise, and (iv) be aware that multi-
3834 label classification require more complex optimisation methods when setting label
3835 specific costs. However, current intelligent services do not sufficiently guide or
3836 support software engineers through the evaluation process, nor do they make this
3837 need clear in the documentation.

3838 In this paper we present **Threshy**³, a tool to assist developers in selecting decision
3839 thresholds when using intelligent services. The motivation for developing Threshy
3840 arose from our consultancy work with industry. Unlike existing tooling (see Section
3841 9.4), **Threshy serves as a means to up-skill and educate software engineers**
3842 **in selecting machine-learnt decision thresholds**, for example, on aspects such as
3843 confusion matrices. Threshy provides a visually interactive interface for developers
3844 to fine-tune thresholds and explore trade-offs of prediction hits/misses. This exposes
3845 the need for optimisation of thresholds, which is dependent on particular use cases.

3846 Threshy improves developer productivity through automation of the threshold
3847 selection process by leveraging an optimisation algorithm to propose thresholds.
3848 The algorithm considers different cost factors providing developers with summary
3849 information so they can make more informed trade-offs. Developers also benefit
3850 from the workflow implemented in Threshy by providing a reproducible procedure
3851 for testing and tuning thresholds for any category of classification problem (binary,
3852 multi-class, and multi-label). Threshy has also been designed to work for different
3853 input data types including images, text and categorical values. The output, is a
3854 text file and can be integrated into client applications ensuring that the thresholds
3855 can be updated without code changes (if needed), and continuously monitored in a
3856 production setting.

3857 9.2 Motivating Example

3858 As a motivating example consider Nina, a fictitious developer, who has been em-
3859 ployed by Lucy’s Tomato Farm to automate the picking of tomatoes from their vines
3860 (when ripe) using computer vision and a harvesting robot. Lucy’s Farm grow five
3861 types of tomatoes (roma, cherry, plum, green, and yellow tomatoes). Nina’s robot—
3862 using an attached webcam—will crawl and take a photo of each vine to assess it
3863 for harvesting. Nina’s automated harvester needs to sort picked tomatoes into a
3864 respective container, and thus several business rules need to be encoded into the
3865 prediction logic to sort each tomato detected based on its *ripeness* (ripe or not ripe)
3866 and *type of tomato* (as above).

³Threshy is available for use at <http://bit.ly/a2i2threshy>.

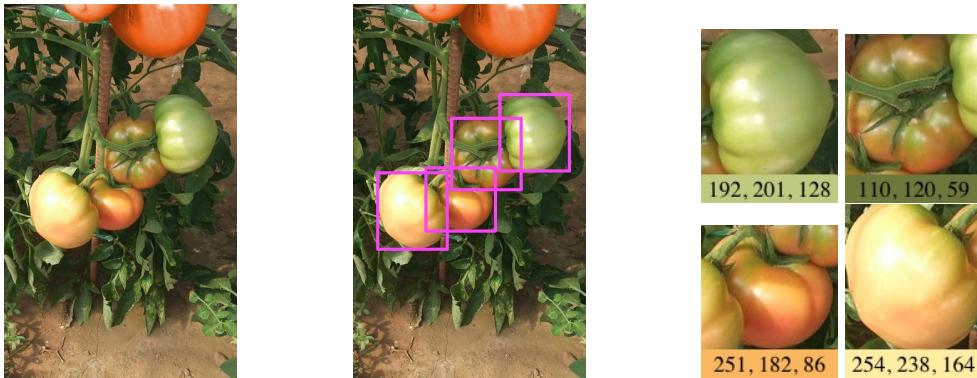


Figure 9.2: Pipeline of Nina’s harvesting robot. *Left:* Photo from harvesting robot’s webcam. *Centre:* Classification detecting different types of tomatoes. *Right:* Binary classification for ripeness (ripe/unripe) based on (R, G, B values).

3867 Nina uses a two-stage pipeline consisting of a multi-class and a binary classi-
 3868 fication model. She has decided to evaluate the viability of cloud based intelligent
 3869 services and use them if operationally effective. Figure 9.2 illustrates an example of
 3870 the the pipeline as listed below:

- 3871 1. **Classify tomato ‘type’.** This stage uses an object localisation service to detect
 3872 all tomato-like objects in the frame and classifies each tomato into one of the
 3873 following labels: [‘roma’, ‘cherry’, ‘plum’, ‘green’, ‘yellow’].
- 3874 2. **Assess tomato ‘ripeness’.** This stage uses a crop of the localised tomatoes
 3875 from the original frame to assess the crop’s colour properties (i.e., average
 3876 colour must have $R > 200$ and $G < 240$). This produces a binary classification
 3877 to deduce whether the tomato is ripe or not.

3878 Nina only has a minimal appreciation of the evaluation method to use for off-
 3879 the-shelf computer vision (classification) services. She also needs to consider the
 3880 financial costs of mis-classifying either the tomato type or the ripeness. Missing a
 3881 few ripe tomatoes isn’t a problem as the robot travels the field twice a week during
 3882 harvest season. However, picking an unripe tomato is expensive as Lucy cannot sell
 3883 them. Therefore, Nina needs a better (automated) way to assess the performance
 3884 of the service and set optimal thresholds for her picking robot, thereby maximising
 3885 profit.

3886 To assist in developing Nina’s pipeline, Lucy sampled a section of 1000 tomatoes
 3887 by taking a photo of each tomato, labelling its type, and assessing whether the vine
 3888 was ‘ripe’ or ‘not_ripe’. Nina ran the labelled images through an intelligent
 3889 service, with each image having a predicted type (multi-class) and ripeness (binary),
 3890 with respective confidence values.

3891 Nina combined the predictions, their respective confidence values, and Lucy’s
 3892 labelled ground truths into a CSV file which was then uploaded to Threshy. Nina
 3893 asked Lucy to assist in setting relevant costs for correct predictions and false predic-
 3894 tions. Threshy then recommended a choice of decision threshold which Nina then
 3895 fine tuned while considering the performance and cost implications.

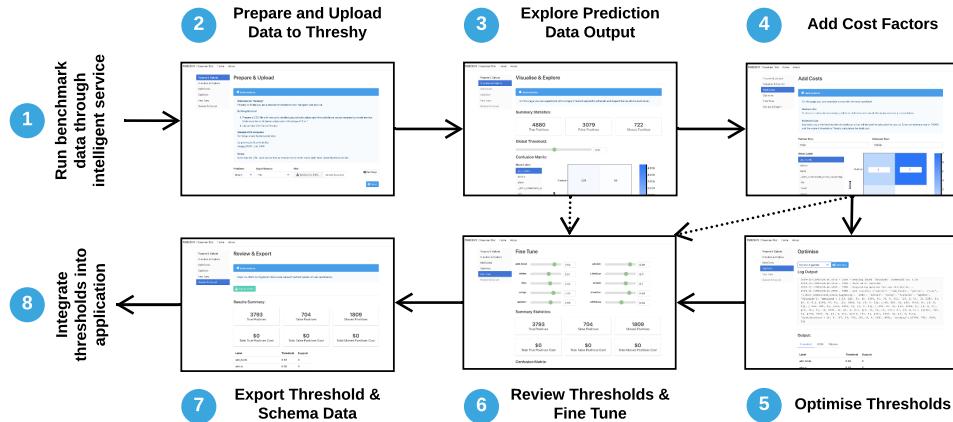


Figure 9.3: UI workflow for interacting with Threshy to optimise the thresholds for classification problem.

3896 9.3 Threshy

3897 Threshy is a tool to assist software engineers with setting decision thresholds when
 3898 integrating machine-learnt components in a system. Our tool also serves as a method
 3899 to inform and educate engineers about the nuances to consider. The novel features
 3900 of Threshy are:

- 3901 • Automating threshold selection using an optimisation algorithm (NSGA-II
 3902 [85]), optimising the results for each label.
- 3903 • Support for additional user defined weights when optimising thresholds such
 3904 as financial costs and impact to society (different type of cost). This allows
 3905 decision thresholds to be set within a business context as they differ from
 3906 application to application [96].
- 3907 • Handles nuances of classification problems such as dealing with multi-objective
 3908 optimisation, and metric selection—reducing errors of omission.
- 3909 • Support key classification problems including binary (e.g. email is either
 3910 spam or ham), multi-class (e.g. predicting the colour of a car), and multi-label
 3911 (e.g. assign multiple topics to a document). Existing tools ignore multi-label
 3912 classification.

3913 Setting thresholds in Threshy is an eight step process as shown in Figure 9.3.
 3914 Software engineers ① run a benchmark dataset through the machine-learnt com-
 3915 ponent to create a CSV file with true labels and predicted labels along with the
 3916 predicted confidence values. The CSV file is then ② uploaded for initial explo-
 3917 ration where engineers can ③ experiment with modifying a single global threshold
 3918 for the dataset. Developers may choose to exit at this point (as indicated by dotted
 3919 arrows in Figure 9.3). Optionally, the engineer ④ defines costs for missed pre-
 3920 dictions followed by selecting optimisation settings. The optional optimisation step of
 3921 Threshy ⑤ considers the performance and costs when deriving the thresholds. Fi-
 3922 nally, the engineer can ⑥ review and fine tune the calculated thresholds, associated

3923 costs, and ⑦ download generated threshold meta-data to be ⑧ integrated into their
3924 application.

3925 Threshy runs a client/server architecture with a thin-client (see Figure 9.4). The
3926 web-based application consists of an interactive front-end where developers upload
3927 benchmark results—consisting of both human annotated labels (ground truths) and
3928 machine predictions (from the intelligent service)—and use threshold tuners (via
3929 sliders) to present a data summary of the uploaded CSV. Predicted performances
3930 and costs are entered manually into the web interface by the developer. The back-end
3931 of Threshy asynchronously runs a data analyser, cost processor and metrics calculator
3932 when relevant changes are made to the front-end’s tuning sliders. Separating the
3933 two concerns allows for high intensity processing to be done on the server and not
3934 the front end.

3935 The data analyser provides a comprehensive overview of confusion matrices
3936 compatible for multi-label multi-class classification problems. When representing
3937 the confusion matrix, it is trivial to represent instances where multi-label multi-
3938 classification is not considered. For example, in the simplest case, a single row in
3939 the matrix represents a single label out of two classes, or each row has one label but
3940 it has multiple classes. However, a more challenging case to visualise the confusion
3941 arises when you have n labels and n classes; the true/false matches become too
3942 excessive to visualise as it is disproportionate to the true results. To deal with this
3943 issue, we condense the summary statistics down to three constructs: (i) number of
3944 true positives, (ii) false positives, (iii) missed positives. This therefore allows us to
3945 optimise against the true positives and minimise the other two constructs.

3946 Threshy is a fully self-contained repository containing implementation of the
3947 tool, scripting and exploratory notebooks, which we make available at <https://github.com/a2i2/threshy>.

3949 9.4 Related work

3950 9.4.1 Decision Boundary Estimation

3951 Optimal machine-learnt decision boundaries depend on identifying the operating
3952 conditions of the problem domain. A systematic study by Drummond and Holte
3953 [96] classifies four such operating conditions to determine a decision threshold: (i)
3954 the operating condition is known and thus the model trained matches perfectly; (ii)
3955 where the operating conditions are known but change with time, and thus the model
3956 must be adaptable to such changes; (iii) where there is uncertainty in the knowledge
3957 of the operating conditions certain changes in the operating condition are more likely
3958 than others; (iv) where there is no knowledge of the operating conditions and the
3959 conditions may change from the model in any possible way. Various approaches
3960 to determine appropriate thresholds exist for all four of these cases, such as cost-
3961 sensitive learning, ROC analysis, cost curves, and Brier scores.

3962 However, an *automated* attempt to calibrate decision threshold boundaries is
3963 not considered, and is largely pitched at a non-software engineering audience. A
3964 more recent study touches on this in model management for large-scale adversarial

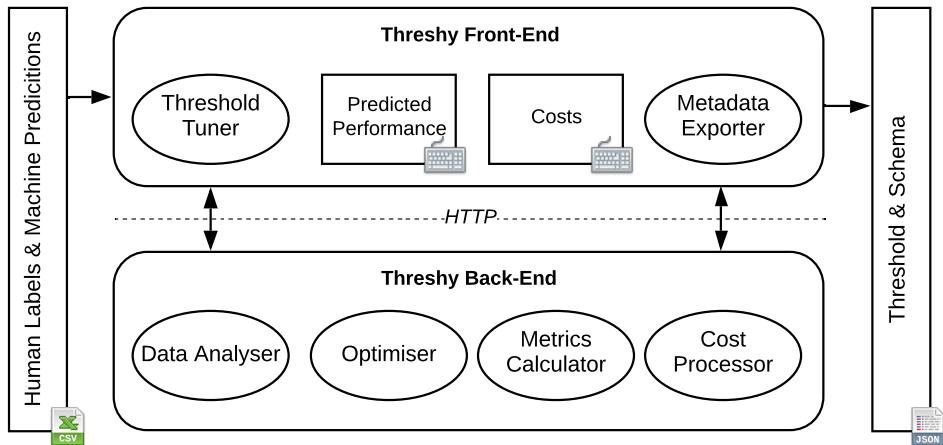


Figure 9.4: Architecture of Threshy.

3965 instances in Google’s advertising system [276], however this is only a single com-
 3966 ponent within the entire architecture, and is not a tool that is useful for developer’s
 3967 in varying contexts. Unlike this study, our work presents a ‘plug-and-play’ style
 3968 calibration method where any context/domain can have thresholds automatically
 3969 calibrated (in-context) *and* optimised for engineers; Threshy’s architecture and
 3970 design facilitates operating in a headless mode enabling use in monitoring and support
 3971 workflows.

3972 9.4.2 Tooling for ML Frameworks

3973 Support tools for ML frameworks generally fall into two categories; the first attempts
 3974 to illuminate the ‘black box’ by offering ways in which developers can better under-
 3975 stand the internals of the model to improve its performance. (For extensive analyses
 3976 and surveys into this area, see [138, 238].) However, a recent emphasis to probe only
 3977 inputs and outputs of a model has been explored, exploring off-the-shelf models
 3978 without knowledge of its unknowns (see Figure 9.1) to reflect the nature of real-
 3979 world development. Google’s *What-If Tool* [324] for Tensorflow provides a means
 3980 for data scientists to visualise, measure and assess model performance and fairness
 3981 with various hypothetical scenarios and data features; similarly, Microsoft’s *Gamut*
 3982 tool [137] provides an interface to test hypotheticals (although only on Generalized
 3983 Additive Models) and their *ModelTracker* tool [9] collates summary statistics on a
 3984 set of sample data to enable rich visualisation of model behaviour and access to key
 3985 performance metrics.

3986 However, these tools are largely focused toward pre-development model eval-
 3987 uation and are not designed for the software engineering workflow. They are also
 3988 targeted to data scientists and not engineers, and certain tools are tied to specific
 3989 machine learning frameworks (e.g., What-If and Tensorflow). Our work attempts to
 3990 bridge these gaps through a structured workflow with an automated tool targeted to
 3991 software developers. We also consider the need to have a consistent tool that works
 3992 across development, test, and production environments.

9.5 Conclusions & Future Work

3994 Primary contributions of this work include Threshy, a tool for automating threshold
3995 selection, and the overall meta-workflow proposed in Threshy that developers can
3996 use as a point of reference for calibrating thresholds. In future work, we plan to
3997 evaluate Threshy with software engineers to identify additional insights required to
3998 make decision thresholds in practice and add code synthesis for monitoring concept
3999 drift and for implementing decision thresholds.

CHAPTER 10

4000

4001

4002

4003

FSE Paper[†]

[†]This chapter is originally based on A. Cummaudo, S. Barnett, R. Vasa, J. Grundy, and M. Abd-elrazek, “Beware the evolving ‘intelligent’ web service! An integration architecture tactic to guard AI-first components,” 2020, Unpublished. Terminology has been updated to fit this thesis.

4004

Part III

4005

Postface

CHAPTER 11

4006

4007

4008

4009

Conclusions & Future Work

4010

References

4011

4012

- 4013 [1] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving,
4014 M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker,
4015 V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, “TensorFlow: A system for large-
4016 scale machine learning,” in *Proceedings of the 12th USENIX Symposium on Operating Systems
4017 Design and Implementation*. Savannah, GA, USA: ACM, 2016. ISBN 978-1-93-197133-1 pp.
4018 265–283.
- 4019 [2] E. Aghajani, C. Nagy, G. Bavota, and M. Lanza, “A Large-scale empirical study on linguistic
4020 antipatterns affecting apis,” in *Proceedings of the 34th International Conference on Software
4021 Maintenance and Evolution*. Madrid, Spain: IEEE, September 2018. DOI 10.1109/IC-
4022 SME.2018.00012. ISBN 978-1-53-867870-1 pp. 25–35.
- 4023 [3] E. Aghajani, C. Nagy, O. L. Vega-Marquez, M. Linares-Vasquez, L. Moreno, G. Bavota,
4024 and M. Lanza, “Software Documentation Issues Unveiled,” in *Proceedings of the 41st Inter-
4025 national Conference on Software Engineering*. Montreal, QC, Canada: IEEE, May 2019.
4026 DOI 10.1109/ICSE.2019.00122. ISBN 978-1-72-810869-8. ISSN 0270-5257 pp. 1199–1210.
- 4027 [4] M. Ahazanuzzaman, M. Asaduzzaman, C. K. Roy, and K. A. Schneider, “Classifying stack
4028 overflow posts on API issues,” in *Proceedings of the 25th International Conference on
4029 Software Analysis, Evolution and Reengineering*. Campobasso, Italy: IEEE, March 2018.
4030 DOI 10.1109/SANER.2018.8330213. ISBN 978-1-53-864969-5 pp. 244–254.
- 4031 [5] R. E. Al-Qutaish, “Quality Models in Software Engineering Literature: An Analytical and
4032 Comparative Study,” *Journal of American Science*, vol. 6, no. 3, pp. 166–175, 2010.
- 4033 [6] H. Allahyari and N. Lavesson, “User-oriented assessment of classification model under-
4034 standability,” in *Proceedings of the 11th Scandinavian Conference on Artificial Intelligence*, vol.
4035 227. Trondheim, Norway: IOS Press, May 2011. DOI 10.3233/978-1-60750-754-3-11.
4036 ISBN 978-1-60-750753-6. ISSN 0922-6389 pp. 11–19.
- 4037 [7] M. Allamanis and C. Sutton, “Why, when, and what: Analyzing stack overflow questions
4038 by topic, type, and code,” in *Proceedings of the 10th IEEE International Working Con-
4039 ference on Mining Software Repositories*. San Francisco, CA, USA: IEEE, May 2013.
4040 DOI 10.1109/MSR.2013.6624004. ISBN 978-1-46-732936-1. ISSN 2160-1852 pp. 53–56.
- 4041 [8] J. Alway and C. Calhoun, *Critical Social Theory: Culture, History, and the Challenge of
4042 Difference*. American Sociological Association, 1997, vol. 26, no. 1, DOI 10.2307/2076647.
- 4043 [9] S. Amershi, M. Chickering, S. M. Drucker, B. Lee, P. Simard, and J. Suh, “Modeltracker:
4044 Redesigning performance analysis tools for machine learning,” in *Proceedings of the 33rd
4045 Annual ACM Conference on Human Factors in Computing Systems*. Seoul, Republic of
4046 Korea: ACM, April 2015. DOI 10.1145/2702123.2702509. ISBN 978-1-45-033145-6 pp.
4047 337–346.
- 4048 [10] K. Arnold, “Programmers are People, Too,” *ACM Queue*, vol. 3, no. 5, pp. 54–59, 2005,
4049 DOI 10.1145/1071713.1071731. ISSN 1542-7749

- [11] A. Arpteg, B. Brinne, L. Crnkovic-Friis, and J. Bosch, "Software engineering challenges of deep learning," in *Proceedings of the 44th Euromicro Conference on Software Engineering and Advanced Applications*. Prague, Czech Republic: IEEE, August 2018. DOI 10.1109/SEAA.2018.00018. ISBN 978-1-53-867382-9 pp. 50–59.
- [12] W. R. Ashby and J. R. Pierce, "An Introduction to Cybernetics," *Physics Today*, vol. 10, no. 7, pp. 34–36, July 1957.
- [13] L. Aversano, D. Guardabascio, and M. Tortorella, "Analysis of the Documentation of ERP Software Projects," *Procedia Computer Science*, vol. 121, pp. 423–430, January 2017, DOI 10.1016/j.procs.2017.11.057. ISSN 1877-0509
- [14] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K. R. Müller, "How to explain individual classification decisions," *Journal of Machine Learning Research*, vol. 11, pp. 1803–1831, 2010. ISSN 1532-4435
- [15] B. Baesens, C. Mues, M. De Backer, J. Vanthienen, and R. Setiono, "Building intelligent credit scoring systems using decision tables," in *Proceedings of the 5th International Conference on Enterprise Information Systems*, vol. 2. Angers, France: IEEE, April 2003. DOI 10.1007/1-4020-2673-0_15. ISBN 9-72-988161-8 pp. 19–25.
- [16] X. Bai, Y. Wang, G. Dai, W. T. Tsai, and Y. Chen, "A framework for contract-based collaborative verification and validation of Web services," in *Proceedings of the 10th International Symposium of Component-Based Software Engineering*. Medford, MA, USA: Springer, July 2007. DOI 10.1007/978-3-540-73551-9_18. ISBN 978-3-54-073550-2. ISSN 0302-9743 pp. 258–273.
- [17] K. Bajaj, K. Pattabiraman, and A. Mesbah, "Mining questions asked by web developers," in *Proceedings of the 11th Working Conference on Mining Software Repositories*. Hyderabad, India: ACM, May 2014. DOI 10.1145/2597073.2597083. ISBN 978-1-45-032863-0 pp. 112–121.
- [18] K. Ballinger, "Simplicity and Utility, or, Why SOAP Lost," [Online] Available: <http://bit.ly/37vLms0>, December 2014, Accessed: 28 August 2018.
- [19] A. Bangor, P. T. Kortum, and J. T. Miller, "An empirical evaluation of the system usability scale," *International Journal of Human-Computer Interaction*, 2008, DOI 10.1080/10447310802205776. ISSN 10447318
- [20] S. Barnett, "Extracting technical domain knowledge to improve software architecture," Ph.D. dissertation, Swinburne University of Technology, Hawthorn, VIC, Australia, 2018.
- [21] S. Barnett, R. Vasa, and J. Grundy, "Bootstrapping Mobile App Development," in *Proceedings of the 37th International Conference on Software Engineering*. Florence, Italy: IEEE, May 2015. DOI 10.1109/ICSE.2015.216. ISBN 978-1-47-991934-5. ISSN 0270-5257 pp. 657–660.
- [22] S. Barnett, R. Vasa, and A. Tang, "A Conceptual Model for Architecting Mobile Applications," in *Proceedings of the 12th Working IEEE/IFIP Conference on Software Architecture*. Montreal, QC, Canada: IEEE, May 2015. DOI 10.1109/WICSA.2015.28. ISBN 978-1-47-991922-2 pp. 105–114.
- [23] A. Barua, S. W. Thomas, and A. E. Hassan, "What are developers talking about? An analysis of topics and trends in Stack Overflow," *Empirical Software Engineering*, vol. 19, no. 3, pp. 619–654, 2014, DOI 10.1007/s10664-012-9231-y. ISSN 1573-7616
- [24] Y. Baruch, "Response rate in academic studies - A comparative analysis," *Human Relations*, vol. 52, no. 4, pp. 421–438, 1999, DOI 10.1177/00182679905200401. ISSN 0018-7267
- [25] O. Barzilay, C. Treude, and A. Zagalsky, "Facilitating crowd sourced software engineering via stack overflow," in *Finding Source Code on the Web for Remix and Reuse*, 2014, no. 4, pp. 289–308. ISBN 978-1-46-146596-6
- [26] L. Bass, P. Clements, and R. Kazman, *Software Architecture in Practice*, 2nd ed. Addison-Wesley, 2003. ISBN 0-32-115495-9
- [27] B. E. Bejnordi, M. Veta, P. J. Van Diest, B. Van Ginneken, N. Karssemeijer, G. Litjens, J. A. W. M. Van Der Laak, M. Hermsen, Q. F. Manson, M. Balkenhol, O. Geessink, N. Stathonikos, M. C. R. F. Van Dijk, P. Bult, F. Beca, A. H. Beck, D. Wang, A. Khosla, R. Gargoya, H. Irshad, A. Zhong, Q. Dou, Q. Li, H. Chen, H. J. Lin, P. A. Heng, C. Haß, E. Bruni, Q. Wong, U. Halici, M. Ü. Öner, R. Cetin-Atalay, M. Berseth, V. Khvatkov, A. Vylegzhannin, O. Kraus, M. Shaban, N. Rajpoot, R. Awan, K. Sirinukunwattana, T. Qaiser, Y. W. Tsang, D. Tellez,

- 4105 J. Annuscheit, P. Hufnagl, M. Valkonen, K. Kartasalo, L. Latonen, P. Ruusuvuori, K. Li-
4106 imatainen, S. Albarqouni, B. Mungal, A. George, S. Demirci, N. Navab, S. Watanabe, S. Seno,
4107 Y. Takenaka, H. Matsuda, H. A. Phoulady, V. Kovalev, A. Kalinovsky, V. Liauchuk, G. Bueno,
4108 M. M. Fernandez-Carrobles, I. Serrano, O. Deniz, D. Racoceanu, and R. Venâncio, "Diagnostic
4109 assessment of deep learning algorithms for detection of lymph node metastases in women with
4110 breast cancer," *Journal of the American Medical Association*, vol. 318, no. 22, pp. 2199–2210,
4111 December 2017, DOI 10.1001/jama.2017.14585. ISSN 1538-3598
- 4112 [28] R. Bellazzi and B. Zupan, "Predictive data mining in clinical medicine: Current issues and
4113 guidelines," *International Journal of Medical Informatics*, vol. 77, no. 2, pp. 81–97, 2008,
4114 DOI 10.1016/j.ijmedinf.2006.11.006. ISSN 1386-5056
- 4115 [29] A. Ben-David, "Monotonicity Maintenance in Information-Theoretic Machine Learning Algo-
4116 rithms," *Machine Learning*, vol. 19, no. 1, pp. 29–43, 1995, DOI 10.1023/A:1022655006810.
4117 ISSN 1573-0565
- 4118 [30] T. Berners-Lee, R. Fielding, and L. Masinter, "Uniform resource identifier (URI): Generic
4119 syntax," Tech. Rep., 2004.
- 4120 [31] L. L. Berry, A. Parasuraman, and V. A. Zeithaml, "SERVQUAL: A multiple-item scale for
4121 measuring consumer perceptions of service quality," *Journal of Retailing*, vol. 64, no. 1, pp.
4122 12–40, 1988, DOI 10.1016/S0148-2963(99)00084-3. ISBN 00224359. ISSN 0022-4359
- 4123 [32] J. Bessin, "The Business Value of Quality," [Online] Available: <https://ibm.co/2u0UDK0>, June
4124 2004.
- 4125 [33] S. Beyer and M. Pinzger, "A manual categorization of android app development issues on stack
4126 overflow," in *Proceedings of the 30th International Conference on Software Maintenance and
4127 Evolution*. Victoria, BC, Canada: IEEE, September 2014. DOI 10.1109/ICSME.2014.88.
4128 ISBN 978-0-76-955303-0 pp. 531–535.
- 4129 [34] S. Beyer, C. MacHo, M. Pinzger, and M. Di Penta, "Automatically classifying posts into question
4130 categories on stack overflow," in *Proceedings of the 26th International Conference on Program
4131 Comprehension*. Gothenburg, Sweden: ACM, May 2018. DOI 10.1145/3196321.3196333.
4132 ISBN 978-1-45-035714-2. ISSN 0270-5257 pp. 211–221.
- 4133 [35] J. Biggs and K. Collis, "Evaluating the Quality of Learning: The SOLO Taxonomy (Structure
4134 of the Observed Learning Outcome)," *Management in Education*, vol. 1, no. 4, p. 20, 1987,
4135 DOI 10.1177/089202068700100412. ISBN 0-12-097551-1. ISSN 0892-0206
- 4136 [36] J. J. Blake, L. P. Maguire, T. M. McGinnity, B. Roche, and L. J. McDaid, "The implementation of
4137 fuzzy systems, neural networks and fuzzy neural networks using FPGAs," *Information Sciences*,
4138 vol. 112, no. 1-4, pp. 151–168, 1998, DOI 10.1016/S0020-0255(98)10029-4. ISSN 0020-0255
- 4139 [37] B. S. Bloom, *Taxonomy of Educational Objectives, Handbook 1: Cognitive Domain*, 2nd ed.
4140 Addison-Wesley Longman, 1956. ISBN 978-0-58-228010-6
- 4141 [38] B. W. Boehm, J. R. Brown, and M. Lipow, "Quantitative evaluation of software quality," in
4142 *Proceedings of the 2nd International Conference on Software Engineering*. San Francisco,
4143 California, USA: IEEE, October 1976. ISSN 0270-5257 pp. 592–605.
- 4144 [39] B. Boehm and V. R. Basili, "Software defect reduction top 10 list," *Software Management*, pp.
4145 419–421, 2007, DOI 10.1109/9780470049167.ch12. ISBN 978-0-47-004916-7
- 4146 [40] B. W. Boehm, *Software engineering economics*. Englewood Cliffs, NJ, USA: Prentice-Hall,
4147 1981. ISBN 0-13-822122-7
- 4148 [41] S. Borsci, S. Federici, and M. Lauriola, "On the dimensionality of the System Usability Scale:
4149 A test of alternative measurement models," *Cognitive Processing*, 2009, DOI 10.1007/s10339-
4150 009-0268-9. ISSN 16124782
- 4151 [42] C. Bottomley, "What part writer? What part programmer? A survey of practices
4152 and knowledge used in programmer writing," in *Proceedings of the 2005 IEEE Interna-
4153 tional Professional Communication Conference*. Limerick, Ireland: IEEE, July 2005.
4154 DOI 10.1109/IPCC.2005.1494255, pp. 802–812.
- 4155 [43] M. Boyd and N. Wilson, "Just ask Siri? A pilot study comparing smartphone digital assistants
4156 and laptop Google searches for smoking cessation advice," *PLoS ONE*, vol. 13, no. 3, 2018,
4157 DOI 10.1371/journal.pone.0194811. ISSN 1932-6203
- 4158 [44] O. Boz, "Extracting decision trees from trained neural networks," in *Proceedings of the 8th ACM
4159 SIGKDD International Conference on Knowledge Discovery and Data Mining*. Edmonton,
4160 AB, Canada: ACM, July 2002. DOI 10.1145/775107.775113, pp. 456–461.

- [4161] [45] M. Bramer, *Principles of Data Mining*, ser. Undergraduate Topics in Computer Science. London, England, UK: Springer, 2016, vol. 180, DOI 10.1007/978-1-4471-7307-6. ISBN 978-1-44-717306-9
- [4162]
- [4163]
- [4164] [46] J. Brandt, P. J. Guo, J. Lewenstein, M. Dontcheva, and S. R. Klemmer, “Two studies of opportunistic programming: Interleaving web foraging, learning, and writing code,” in *Proceedings of the SIGCHI Conference on Human Factors in Computing System*. Boston, MA, USA: ACM, April 2009. DOI 10.1145/1518701.1518944. ISBN 978-1-60-558247-4 pp. 1589–1598.
- [4165]
- [4166]
- [4167]
- [4168] [47] L. Bratthall and M. Jørgensen, “Can you trust a single data source exploratory software engineering case study?” *Empirical Software Engineering*, 2002, DOI 10.1023/A:1014866909191. ISSN 1382-3256
- [4169]
- [4170]
- [4171] [48] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and regression trees*. New York, NY, USA: CRC press, 1984. DOI 10.1201/9781315139470. ISBN 978-1-35-146049-1
- [4172]
- [4173]
- [4174] [49] P. Brereton, B. A. Kitchenham, D. Budgen, M. Turner, and M. Khalil, “Lessons from applying the systematic literature review process within the software engineering domain,” *Journal of Systems and Software*, vol. 80, no. 4, pp. 571–583, April 2007, DOI 10.1016/j.jss.2006.07.009. ISSN 0164-1212
- [4175]
- [4176]
- [4177]
- [4178] [50] J. Brooke, “SUS-A quick and dirty usability scale,” *Usability Evaluation in Industry*, pp. 189–194, 1996. ISBN 978-0-74-840460-5
- [4179]
- [4180] [51] ——, “SUS: a retrospective,” *Journal of Usability Studies*, vol. 8, no. 2, pp. 29–40, 2013. ISSN 1931-3357
- [4181]
- [4182] [52] O. Bruna, H. Avetisyan, and J. Holub, “Emotion models for textual emotion classification,” *Journal of Physics: Conference Series*, vol. 772, p. 12063, November 2016, DOI 10.1088/1742-6596/772/1/012063.
- [4183]
- [4184]
- [4185] [53] M. Bunge, “A General Black Box Theory,” *Philosophy of Science*, vol. 30, no. 4, pp. 346–358, October 1963, DOI 10.1086/287954. ISSN 0031-8248
- [4186]
- [4187] [54] BusinessWire, “FileShadow Delivers Machine Learning to End Users with Google Vision API | Business Wire,” [Online] Available: <https://bwnews.pr/2O5qv78>, July 2018, Accessed: 25 January 2019.
- [4188]
- [4189]
- [4190] [55] A. Bussone, S. Stumpf, and D. O’Sullivan, “The role of explanations on trust and reliance in clinical decision support systems,” in *Proceedings of the 2015 IEEE International Conference on Healthcare Informatics*. Dallas, TX, USA: IEEE, October 2015. DOI 10.1109/IChI.2015.26. ISBN 978-1-46-739548-9 pp. 160–169.
- [4191]
- [4192]
- [4193]
- [4194] [56] F. Calefato, F. Lanobile, and N. Novielli, “EmoTxt: a toolkit for emotion recognition from text,” in *Proceedings of the 7th International Conference on Affective Computing and Intelligent Interaction Workshops and Demos*. San Antonio, TX, USA: IEEE, October 2017. DOI 10.1109/ACIIW.2017.8272591, pp. 79–80.
- [4195]
- [4196]
- [4197]
- [4198] [57] F. Calefato, F. Lanobile, F. Maiorano, and N. Novielli, “Sentiment polarity detection for software development,” *Empirical Software Engineering*, vol. 23, no. 3, pp. 1352–1382, 2018, DOI 10.1007/s10664-017-9546-9.
- [4199]
- [4200]
- [4201] [58] G. Canfora, “User-side testing of Web Services,” in *Proceedings of the 9th European Conference on Software Maintenance and Reengineering*. Manchester, England, UK: IEEE, March 2005. DOI 10.1109/csmr.2005.57. ISSN 1534-5351 p. 301.
- [4202]
- [4203]
- [4204] [59] G. Canfora and M. Di Penta, “Testing services and service-centric systems: Challenges and opportunities,” *IT Professional*, vol. 8, no. 2, pp. 10–17, 2006, DOI 10.1109/MITP.2006.51. ISSN 1520-9202
- [4205]
- [4206]
- [4207] [60] R. Caruana, Y. Lou, J. Gehrke, P. Koch, M. Sturm, and N. Elhadad, “Intelligible models for healthcare: Predicting pneumonia risk and hospital 30-day readmission,” in *Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, vol. 2015-Augus. Sydney, Australia: ACM, August 2015. DOI 10.1145/2783258.2788613. ISBN 978-1-45-033664-2 pp. 1721–1730.
- [4208]
- [4209]
- [4210]
- [4211]
- [4212] [61] F. Casati, H. Kuno, G. Alonso, and V. Machiraju, *Web Services-Concepts, Architectures and Applications*, 2004. ISBN 978-3-64-207888-0
- [4213]
- [4214] [62] J. P. Cavano and J. A. McCall, “A framework for the measurement of software quality,” in *Proceedings of the Software Quality Assurance Workshop on Functional and Performance Issues*, vol. 3, no. 5, November 1978, DOI 10.1145/800283.811113, pp. 133–139.
- [4215]
- [4216]

- [63] C. V. Chambers, D. J. Balaban, B. L. Carlson, and D. M. Grasberger, "The effect of microcomputer-generated reminders on influenza vaccination rates in a university-based family practice center." *The Journal of the American Board of Family Practice / American Board of Family Practice*, vol. 4, no. 1, pp. 19–26, 1991, DOI 10.3122/jabfm.4.1.19. ISSN 0893-8652
- [64] J. Cheng and R. Greiner, "Learning bayesian belief network classifiers: Algorithms and system," in *Proceedings of the 14th Biennial Conference of the Canadian Society for Computational Studies of Intelligence*, vol. 2056. Ottawa, ON, Canada: Springer, June 2001. DOI 10.1007/3-540-45153-6_14. ISBN 3-54-042144-0. ISSN 1611-3349 pp. 141–151.
- [65] R. Chillarege, I. S. Bhandari, J. K. Chaar, M. J. Halliday, B. K. Ray, and D. S. Moebus, "Orthogonal Defect Classification—A Concept for In-Process Measurements," *IEEE Transactions on Software Engineering*, vol. 18, no. 11, pp. 943–956, 1992, DOI 10.1109/32.177364. ISSN 0098-5589
- [66] E. Choi, N. Yoshida, R. G. Kula, and K. Inoue, "What do practitioners ask about code clone? a preliminary investigation of stack overflow," in *Proceedings of the 9th International Workshop on Software Clones*, Montreal, QC, Canada, March 2015, DOI 10.1109/IWSC.2015.7069890. ISBN 978-1-46-736914-5 pp. 49–50.
- [67] D. V. Cicchetti, "Guidelines, Criteria, and Rules of Thumb for Evaluating Normed and Standardized Assessment Instruments in Psychology," *Psychological Assessment*, vol. 6, no. 4, pp. 284–290, 1994, DOI 10.1037/1040-3590.6.4.284. ISSN 10403590
- [68] Digital, "Case Study: Finding defects earlier yields enormous savings," [Online] Available: <http://bit.ly/36II2cE>, 2003.
- [69] P. Clark and R. Boswell, "Rule induction with CN2: Some recent improvements," in *Proceedings of the 1991 European Working Session on Learning*. Porto, Portugal: Springer, March 1991. DOI 10.1007/BFb0017011. ISBN 978-3-54-053816-5. ISSN 1611-3349 pp. 151–163.
- [70] J. Cohen, "A Coefficient of Agreement for Nominal Scales," *Educational and Psychological Measurement*, vol. 20, no. 1, pp. 37–46, 1960, DOI 10.1177/001316446002000104. ISSN 1552-3888
- [71] P. Covington, J. Adams, and E. Sargin, "Deep neural networks for youtube recommendations," in *Proceedings of the 10th ACM Conference on Recommender Systems*. Boston, MA, USA: ACM, September 2016. DOI 10.1145/2959100.2959190. ISBN 978-1-45-034035-9 pp. 191–198.
- [72] M. W. Craven and J. W. Shavlik, "Extracting tree-structured representations of trained neural networks," in *Proceedings of the 8th International Conference on Neural Information Processing Systems*, vol. 8. Denver, CO, USA: MIT Press, December 1996. ISBN 978-0-26-220107-0 pp. 24–30.
- [73] J. W. Creswell, *Research design: Qualitative, quantitative, and mixed methods approaches*, 4th ed. SAGE, 2017. ISBN 860-1-40-429618-5
- [74] P. B. Crosby, *Quality is free: The art of making quality certain*. McGraw-Hill, 1979. ISBN 978-0-07-014512-2
- [75] A. Cummaudo, R. Vasa, and J. Grundy, "What should I document? A preliminary systematic mapping study into API documentation knowledge," in *Proceedings of the 13th International Symposium on Empirical Software Engineering and Measurement*. Porto de Galinhas, Recife, Brazil: IEEE, October 2019. DOI 10.1109/ESEM.2019.8870148. ISBN 978-1-72-812968-6. ISSN 1949-3789 pp. 1–6.
- [76] A. Cummaudo, R. Vasa, J. Grundy, M. Abdelrazek, and A. Cain, "Losing Confidence in Quality: Unspoken Evolution of Computer Vision Services," in *Proceedings of the 35th IEEE International Conference on Software Maintenance and Evolution*. Cleveland, OH, USA: IEEE, December 2019. DOI 10.1109/ICSME.2019.00051. ISBN 978-1-72-813094-1 pp. 333–342.
- [77] A. Cummaudo, S. Barnett, R. Vasa, and J. Grundy, "Threshy: Supporting Safe Usage of Intelligent Web Services," 2020, Unpublished.
- [78] A. Cummaudo, S. Barnett, R. Vasa, J. Grundy, and M. Abdelrazek, "Beware the evolving 'intelligent' web service! An integration architecture tactic to guard AI-first components," 2020, Unpublished.
- [79] A. Cummaudo, R. Vasa, S. Barnett, J. Grundy, and M. Abdelrazek, "Interpreting Cloud Computer Vision Pain-Points: A Mining Study of Stack Overflow," in *Proceedings of the 42nd*

- 4273 *International Conference on Software Engineering*. Seoul, Republic of Korea: IEEE, October
4274 2020, In Press.
- 4275 [80] A. Cummaudo, R. Vasa, and J. Grundy, "Assessing API documentation knowledge for computer
4276 vision services," 2020, Unpublished.
- 4277 [81] M. K. Curumsing, A. Cummaudo, U. M. Graestch, S. Barnett, and R. Vasa, "Ranking Computer
4278 Vision Service Issues using Emotion," 2020, Unpublished.
- 4279 [82] M. K. Curumsing, "Emotion-Oriented Requirements Engineering," Ph.D. dissertation, Swin-
4280 burne University of Technology, Hawthorn, VIC, Australia, 2017.
- 4281 [83] H. da Mota Silveira and L. C. Martini, "How the New Approaches on Cloud Computer Vision
4282 can Contribute to Growth of Assistive Technologies to Visually Impaired in the Following
4283 Years?" *Journal of Information Systems Engineering & Management*, vol. 2, no. 2, pp. 1–3,
4284 2017, DOI 10.20897/jisem.201709. ISSN 2468-4376
- 4285 [84] R. M. Davison, M. G. Martinsons, and N. Kock, "Principles of canonical action re-
4286 search," *Information Systems Journal*, vol. 14, no. 1, pp. 65–86, 2004, DOI 10.1111/j.1365-
4287 2575.2004.00162.x. ISSN 1350-1917
- 4288 [85] K. Deb, A. Pratap, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic
4289 algorithm: NSGA-II," *IEEE Transactions on Evolutionary Computation*, vol. 6, no. 2, pp.
4290 182–197, April 2002, DOI 10.1109/4235.996017. ISSN 1089778X
- 4291 [86] K. Dejaeger, F. Goethals, A. Giangreco, L. Mola, and B. Baesens, "Gaining insight into student
4292 satisfaction using comprehensible data mining techniques," *European Journal of Operational
4293 Research*, vol. 218, no. 2, pp. 548–562, 2012, DOI 10.1016/j.ejor.2011.11.022. ISSN 0377-2217
- 4294 [87] I. Dey, *Qualitative Data Analysis: A User-Friendly Guide for Social Scientists*. New York,
4295 NY: Routledge, 1993. DOI 10.4324/9780203412497. ISBN 978-0-41-505852-0
- 4296 [88] V. Dhar, D. Chou, and F. Provost, "Discovering interesting patterns for investment decision
4297 making with GLOWER - A genetic learner overlaid with entropy reduction," *Data Mining and
4298 Knowledge Discovery*, vol. 4, no. 4, pp. 69–80, 2000, DOI 10.1023/A:1009848126475. ISSN
4299 1384-5810
- 4300 [89] V. Dibia, A. Cox, and J. Weisz, "Designing for Democratization: Introducing Novices to
4301 Artificial Intelligence Via Maker Kits," in *Proceedings of the 2017 CHI Conference Extended
4302 Abstracts on Human Factors in Computing Systems*. Denver, CO, USA: ACM, May 2017, pp.
4303 381–384.
- 4304 [90] M. Doderer, K. Yoon, J. Salinas, and S. Kwek, "Protein subcellular localization prediction
4305 using a hybrid of similarity search and Error-Correcting Output Code techniques that produces
4306 interpretable results," *In Silico Biology*, vol. 6, no. 5, pp. 419–433, 2006. ISSN 1386-6338
- 4307 [91] P. Domingos, "Occam's Two Razors: The Sharp and the Blunt," in *Proceedings of the 4th
4308 International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA:
4309 AAAI, August 1998. DOI 10.1.1.40.3278, pp. 37–43.
- 4310 [92] B. Dorn and M. Guzdial, "Learning on the job: Characterizing the programming knowl-
4311 edge and learning strategies of web designers," in *Proceedings of the 28th ACM Conference
4312 on Human Factors in Computing Systems*, vol. 2. Atlanta, GA, USA: ACM, April 2010.
4313 DOI 10.1145/1753326.1753430. ISBN 978-1-60-558929-9 pp. 703–712.
- 4314 [93] F. Doshi-Velez and B. Kim, "Towards A Rigorous Science of Interpretable Machine Learning,"
4315 *arXiv preprint arXiv:1702.08608*, 2017.
- 4316 [94] F. Doshi-Velez, M. Kortz, R. Budish, C. Bavitz, S. J. Gershman, D. O'Brien, S. Shieber,
4317 J. Waldo, D. Weinberger, and A. Wood, "Accountability of AI Under the Law: The Role of
4318 Explanation," *SSRN Electronic Journal*, November 2017, In Press, DOI 10.2139/ssrn.3064761.
- 4319 [95] R. G. Dromey, "A model for software product quality," *IEEE Transactions on Software Engi-
4320 neering*, vol. 21, no. 2, pp. 146–162, 1995, DOI 10.1109/32.345830. ISBN 978-1-11-815666-7.
4321 ISSN 0098-5589
- 4322 [96] C. Drummond and R. C. Holte, "Cost curves: An improved method for visualizing classifier per-
4323 formance," *Machine Learning*, vol. 65, no. 1, pp. 95–130, October 2006, DOI 10.1007/s10994-
4324 006-8199-5. ISSN 0885-6125
- 4325 [97] S. Easterbrook, J. Singer, M.-A. Storey, and D. Damian, "Selecting empirical methods for
4326 software engineering research," in *Guide to Advanced Empirical Software Engineering*, F. Shull,
4327 J. Singer, and D. I. K. Sjøberg, Eds. Springer, November 2007, ch. 11, pp. 285–311. ISBN
4328 978-1-84-800043-8

- 4329 [98] W. Elazmeh, S. Matwin, D. O’Sullivan, W. Michalowski, and K. Farion, “Insights from pre-
4330 dicting pediatric asthma exacerbations from retrospective clinical data,” in *Proceedings of the*
4331 *22nd Conference on Artificial Intelligence*, vol. WS-07-05. Vancouver, BC, Canada: AAAI,
4332 July 2007. ISBN 978-1-57-735332-4 pp. 10–15.
- 4333 [99] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno,
4334 and D. Song, “Robust Physical-World Attacks on Deep Learning Visual Classification,” in
4335 *Proceedings of the 2017 IEEE Computer Society Conference on Computer Vision and Pattern*
4336 *Recognition*, Honolulu, HI, USA, July 2018, DOI 10.1109/CVPR.2018.00175. ISBN 978-1-
4337 53-866420-9. ISSN 1063-6919 pp. 1625–1634.
- 4338 [100] F. Elder, D. Michie, D. J. Spiegelhalter, and C. C. Taylor, “Machine Learning, Neural, and
4339 Statistical Classification.” *Journal of the American Statistical Association*, vol. 91, no. 433, pp.
4340 436–438, 1996, DOI 10.2307/2291432. ISBN 978-0-13-106360-0. ISSN 0162-1459
- 4341 [101] A. J. Feelders, “Prior knowledge in economic applications of data mining,” in *Proceedings of*
4342 *the 4th European Conference on Principles of Data Mining and Knowledge Discovery*, vol.
4343 1910. Lyon, France: Springer, September 2000. DOI 10.1007/3-540-45372-5_42. ISBN
4344 978-3-54-041066-9. ISSN 1611-3349 pp. 395–400.
- 4345 [102] R. T. Fielding, “Architectural Styles and the Design of Network-based Software Architectures,”
4346 Ph.D. dissertation, University of California, Irvine, 2000.
- 4347 [103] I. Finalyson, “Nondeterministic Finite Automata,” [Online] Available: <http://bit.ly/319GOF9>,
4348 Fredericksburg, VA, USA, 2018.
- 4349 [104] H. Foster, S. Uchitel, J. Magee, and J. Kramer, “Model-based verification of Web service
4350 compositions,” in *Proceedings of the 18th International Conference on Automated Software*
4351 *Engineering*. Linz, Austria: IEEE, September 2004. DOI 10.1109/ase.2003.1240303, pp.
4352 152–161.
- 4353 [105] A. A. Freitas, “A critical review of multi-objective optimization in data mining,” *ACM SIGKDD*
4354 *Explorations Newsletter*, vol. 6, no. 2, p. 77, 2004, DOI 10.1145/1046456.1046467. ISSN 1931-
4355 0145
- 4356 [106] ———, “Comprehensible classification models,” *ACM SIGKDD Explorations Newsletter*, vol. 15,
4357 no. 1, pp. 1–10, March 2014, DOI 10.1145/2594473.2594475. ISSN 1931-0145
- 4358 [107] A. A. Freitas, D. C. Wieser, and R. Apweiler, “On the importance of comprehensible classi-
4359 fication models for protein function prediction,” *IEEE/ACM Transactions on Computational*
4360 *Biology and Bioinformatics*, vol. 7, no. 1, pp. 172–182, 2010, DOI 10.1109/TCBB.2008.47.
4361 ISSN 1545-5963
- 4362 [108] B. J. Frey and D. Dueck, “Clustering by passing messages between data points,” *Science*, vol.
4363 315, no. 5814, pp. 972–976, February 2007, DOI 10.1126/science.1136800. ISSN 0036-8075
- 4364 [109] N. Friedman, D. Geiger, and M. Goldszmidt, “Bayesian Network Classifiers,” *Machine Learn-
4365 ing*, vol. 29, no. 2-3, pp. 131–163, 1997, DOI 10.1002/9780470400531.eorms0099. ISSN
4366 0885-6125
- 4367 [110] G. Fung, S. Sandilya, and R. B. Rao, “Rule extraction from linear support vector machines,”
4368 *Studies in Computational Intelligence*, vol. 80, no. 1, pp. 83–107, 2009, DOI 10.1007/978-3-
4369 540-75390-2_4.
- 4370 [111] D. Gachechiladze, F. Lanubile, N. Novielli, and A. Serebrenik, “Anger and its direction in
4371 collaborative software development,” in *Proceedings of the 39th International Conference on*
4372 *Software Engineering: New Ideas and Emerging Technologies Results Track*, IEEE. Buenos
4373 Aires, Argentina: IEEE, May 2017. DOI 10.1109/ICSE-NIER.2017.18, pp. 11–14.
- 4374 [112] M. Gamer, J. Lemon, I. Fellows, and P. Singh, “Irr: various coefficients of interrater reliability,”
4375 *R package version 0.83*, 2010.
- 4376 [113] V. Garousi and M. Felderer, “Experience-based guidelines for effective and efficient data ex-
4377 traction in systematic reviews in software engineering,” in *Proceedings of the 21st International*
4378 *Conference on Evaluation and Assessment in Software Engineering*, vol. Part F1286. Karl-
4379 skrona, Sweden: ACM, June 2017. DOI 10.1145/3084226.3084238. ISBN 978-1-45-034804-1
4380 pp. 170–179.
- 4381 [114] V. Garousi, M. Felderer, and M. V. Mäntylä, “Guidelines for including grey literature and
4382 conducting multivocal literature reviews in software engineering,” *Information and Software*
4383 *Technology*, vol. 106, pp. 101–121, 2019, DOI 10.1016/j.infsof.2018.09.006. ISSN 0950-5849

- [4384] [115] D. A. Garvin, "What Does 'Product Quality' Really Mean?" *MIT Sloan Management Review*, vol. 26, no. 1, pp. 25–43, 1984. ISSN 0019-848X
- [4385] [116] 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: A Collaborative Ethnography of Documentation Work," *Computer Supported Cooperative Work: CSCW: An International Journal*, vol. 27, no. 3-6, pp. 767–802, May 2018, DOI 10.1007/s10606-018-9333-1. ISSN 15737551
- [4391] [117] GeoSpatial World, "Mapillary and Amazon Rekognition collaborate to build a parking solution for US cities through computer vision," [Online] Available: <http://bit.ly/36AdRmS>, September 2018, Accessed: 25 January 2019.
- [4394] [118] M. Gethsiyal Augasta and T. Kathirvalavakumar, "Reverse engineering the neural networks for rule extraction in classification problems," *Neural Processing Letters*, vol. 35, no. 2, pp. 131–150, 2012, DOI 10.1007/s11063-011-9207-8. ISSN 1370-4621
- [4397] [119] H. L. Gilmore, "Product conformance cost," *Quality progress*, vol. 7, no. 5, pp. 16–19, 1974.
- [4398] [120] R. L. Glass, I. Vessey, and V. Ramesh, "RESRES: The story behind the paper "Research in software engineering: An analysis of the literature"," *Information and Software Technology*, vol. 51, no. 1, pp. 68–70, 2009, DOI 10.1016/j.infsof.2008.09.015. ISSN 0950-5849
- [4401] [121] M. W. Godfrey and D. M. German, "The past, present, and future of software evolution," in *Proceedings of the 2008 Frontiers of Software Maintenance*, Beijing, China, October 2008, DOI 10.1109/FOSM.2008.4659256. ISBN 978-1-42-442655-3 pp. 129–138.
- [4404] [122] M. W. Godfrey and Q. Tu, "Evolution in open source software: a case study," in *Conference on Software Maintenance*. San Jose, CA, USA: IEEE, August 2000. DOI 10.1109/icsm.2000.883030, pp. 131–142.
- [4407] [123] Google LLC, "Classification: Thresholding | Machine Learning Crash Course," [Online] Available: <http://bit.ly/36oMgWb>, 2019, Accessed: 5 February 2020.
- [4409] [124] D. Graziotin, F. Fagerholm, X. Wang, and P. Abrahamsson, "What happens when software developers are (un)happy," *Journal of Systems and Software*, 2018, DOI 10.1016/j.jss.2018.02.041. ISSN 0164-1212
- [4412] [125] P. D. Grünwald, *The Minimum Description Length Principle*. MIT press, 2019. DOI 10.7551/mitpress/4643.001.0001.
- [4414] [126] M. J. Hadley and H. Marc, "Web Application Description Language," [Online] Available: <http://bit.ly/2RXRhQ1>, August 2009.
- [4416] [127] H. A. Haenssle, C. Fink, R. Schneiderbauer, F. Toberer, T. Buhl, A. Blum, A. Kalloo, A. Ben Hadj Hassen, L. Thomas, A. Enk, L. Uhlmann, C. Alt, M. Arenbergerova, R. Bakos, A. Baltzer, I. Bertlich, A. Blum, T. Bokor-Billmann, J. Bowling, N. Braghierioli, R. Braun, K. Budernbachaya, T. Buhl, H. Cabo, L. Cabrijan, N. Cevic, A. Classen, D. Deltgen, C. Fink, I. Georgieva, L. E. Hakim-Meibodi, S. Hanner, F. Hartmann, J. Hartmann, G. Haus, E. Hoxha, R. Karls, H. Koga, J. Kreusch, A. Lallas, P. Majenka, A. Marghoob, C. Massone, L. Mekokishvili, D. Mestel, V. Meyer, A. Neuberger, K. Nielsen, M. Oliviero, R. Pampana, J. Paoli, E. Pawlik, B. Rao, A. Rendon, T. Russo, A. Sadek, K. Samhaber, R. Schneiderbauer, A. Schweizer, F. Toberer, L. Trennheuser, L. Vlahova, A. Wald, J. Winkler, P. Wołbing, and I. Zalaudek, "Man against Machine: Diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists," *Annals of Oncology*, vol. 29, no. 8, pp. 1836–1842, May 2018, DOI 10.1093/annonc/mdy166. ISSN 1569-8041
- [4428] [128] K. A. Hallgren, "Computing Inter-Rater Reliability for Observational Data: An Overview and Tutorial," *Tutorials in Quantitative Methods for Psychology*, vol. 8, no. 1, pp. 23–34, February 2012, DOI 10.20982/tqmp.08.1.p023. ISSN 1913-4126
- [4431] [129] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in *Proceedings of the 30th International Conference on Neural Information Processing Systems*. Barcelona, Spain: Curran Associates Inc., December 2016. DOI 978-1-51-083881-9. ISSN 1049-5258 pp. 3323–3331.
- [4435] [130] S. Haselbock, R. Weinreich, G. Buchgeher, and T. Kriechbaum, "Microservice Design Space Analysis and Decision Documentation: A Case Study on API Management," in *Proceedings of the 11th International Conference on Service-Oriented Computing and Applications, SOCA 2018*, Paris, France, November 2019, DOI 10.1109/SOCA.2018.00008, pp. 1–8.

- 4439 [131] T. Hastie, R. Tibshirani, and J. H. Friedman, *The Elements of Statistical Learning*, 2nd ed., ser.
4440 Data Mining, Inference, and Prediction. Springer, January 2001.
- 4441 [132] B. Hayete and J. R. Bienkowska, “Gotrees: Predicting go associations from protein domain
4442 composition using decision trees,” in *Proceedings of the Pacific Symposium on Biocomputing 2005, PSB 2005*. Hawaii, USA: World Scientific Publishing Company, January 2005.
4443 DOI 10.1142/9789812702456_0013. ISBN 9-81-256046-7 pp. 127–138.
- 4444 [133] A. Head, C. Sadowski, E. Murphy-Hill, and A. Knight, “When not to comment: Questions and
4445 tradeoffs with API documentation for C++ projects,” in *Proceedings of the 40th International
4446 Conference on Software Engineering*, ser. questions and tradeoffs with API documentation for
4447 C++ projects. Gothenburg, Sweden: ACM, May 2018. DOI 10.1145/3180155.3180176.
4448 ISSN 0270-5257 pp. 643–653.
- 4449 [134] R. Heckel and M. Lohmann, “Towards Contract-based Testing of Web Services,” *Electronic Notes in Theoretical Computer Science*, vol. 116, pp. 145–156, January 2005,
4450 DOI 10.1016/j.entcs.2004.02.073. ISSN 1571-0661
- 4451 [135] D. Heckerman, D. M. Chickering, C. Meek, R. Rounthwaite, and C. Kadie, “Dependency
4452 networks for inference, collaborative filtering, and data visualization,” *Journal of Machine
4453 Learning Research*, vol. 1, no. 1, pp. 49–75, 2001, DOI 10.1162/153244301753344614. ISSN
4454 1532-4435
- 4455 [136] M. Henning, “API design matters,” *Communications of the ACM*, vol. 52, no. 5, pp. 46–56,
4456 2009, DOI 10.1145/1506409.1506424. ISSN 0001-0782
- 4457 [137] F. Hohman, A. Head, R. Caruana, R. DeLine, and S. M. Drucker, “Gamut: A design probe
4458 to understand how data scientists understand machine learning models,” in *Proceedings of the
4459 2019 CHI Conference on Human Factors in Computing Systems*. Glasgow, Scotland, UK:
4460 ACM, May 2019. DOI 10.1145/3290605.3300809. ISBN 978-1-45-035970-2
- 4461 [138] F. Hohman, M. Kahng, R. Pienta, and D. H. Chau, “Visual Analytics in Deep Learning: An
4462 Interrogative Survey for the Next Frontiers,” *IEEE Transactions on Visualization and Computer
4463 Graphics*, vol. 25, no. 8, pp. 2674–2693, 2019, DOI 10.1109/TVCG.2018.2843369. ISSN
4464 1941-0506
- 4465 [139] J. W. Horch, *Practical Guide To Software Quality Management*. Artech House, 2003. ISBN
4466 978-1-58-053604-2
- 4467 [140] H. Hosseini, B. Xiao, and R. Poovendran, “Google’s cloud vision API is not robust to noise,” in
4468 *Proceedings of the 16th IEEE International Conference on Machine Learning and Applications*,
4469 vol. 2017-Decem. Cancun, Mexico: IEEE, December 2017. DOI 10.1109/ICMLA.2017.0-
4470 172. ISBN 978-1-53-861417-4 pp. 101–105.
- 4471 [141] D. Hou and L. Mo, “Content categorization of API discussions,” in *Proceedings of the 29th In-
4472 ternational Conference on Software Maintenance*. Eindhoven, Netherlands: IEEE, September
4473 2013. DOI 10.1109/ICSM.2013.17, pp. 60–69.
- 4474 [142] C. Howard, “Introducing Google AI,” [Online] Available: <http://bit.ly/2uI6vAr>, May 2018,
4475 Accessed: 28 August 2018.
- 4476 [143] J. Huang and C. X. Ling, “Using AUC and accuracy in evaluating learning algorithms,”
4477 *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 3, pp. 299–310, 2005,
4478 DOI 10.1109/TKDE.2005.50. ISSN 1041-4347
- 4479 [144] J. Huysmans, K. Dejaeger, C. Mues, J. Vanthienen, and B. Baesens, “An empirical evaluation
4480 of the comprehensibility of decision table, tree and rule based predictive models,” *Decision
4481 Support Systems*, vol. 51, no. 1, pp. 141–154, April 2011, DOI 10.1016/j.dss.2010.12.003.
4482 ISSN 0167-9236
- 4483 [145] IEEE, “IEEE Standard Glossary of Software Engineering Terminology,” 1990.
- 4484 [146] International Organization for Standardization, “ISO25010:2011 - Systems and software engi-
4485 neering – Systems and software Quality Requirements and Evaluation (SQuaRE) – System and
4486 software quality models,” 2011.
- 4487 [147] ———, “ISO 8402:1986 Information Technology - Software Product Evaluation - Quality Char-
4488 acteristics and Guidelines for Their Use,” [Online] Available: <http://bit.ly/37SK4HP>, 1986.
- 4489 [148] ———, “ISO 9000:2015 Quality management systems – Fundamentals and vocabulary,” [Online]
4490 Available: <http://bit.ly/37O4oKo>, 2015.
- 4491

- 4493 [149] ——, “ISO/IEC 9126 Information Technology - Software Product Evaluation - Quality Charac-
4494 teristics and Guidelines for Their Use,” [Online] Available: <http://bit.ly/2tgMHUE>, November
4495 1999.
- 4496 [150] S. Inzunza, R. Juárez-Ramírez, and S. Jiménez, “API Documentation,” in *Proceedings of the*
4497 *6th World Conference on Information Systems and Technologies*. Naples, Italy: Springer,
4498 March 2018. DOI 10.1007/978-3-319-77712-2_22, pp. 229–239.
- 4499 [151] A. Iyengar, “Supporting Data Analytics Applications Which Utilize Cognitive Services,” in
4500 *Proceedings of the 37th International Conference on Distributed Computing Systems*. Atlanta,
4501 GA, USA: IEEE, June 2017. DOI 10.1109/ICDCS.2017.172. ISBN 978-1-53-861791-5 pp.
4502 1856–1864.
- 4503 [152] N. Japkowicz and M. Shah, *Evaluating learning algorithms: A classification perspective*.
4504 Cambridge University Press, 2011, vol. 9780521196, DOI 10.1017/CBO9780511921803. ISBN
4505 978-0-51-192180-3
- 4506 [153] M. W. M. Jaspers, M. Smeulers, H. Vermeulen, and L. W. Peute, “Effects of clinical decision-
4507 support systems on practitioner performance and patient outcomes: A synthesis of high-quality
4508 systematic review findings,” *Journal of the American Medical Informatics Association*, vol. 18,
4509 no. 3, pp. 327–334, 2011, DOI 10.1136/amiajnl-2011-000094. ISSN 1067-5027
- 4510 [154] S. Y. Jeong, Y. Xie, J. Beaton, B. A. Myers, J. Stylos, R. Ehret, J. Karstens, A. Efeoglu, and
4511 D. K. Busse, “Improving documentation for eSOA APIs through user studies,” in *Proceedings*
4512 *of the First International Symposium on End User Development*, vol. 5435 LNCS. Siegen,
4513 Germany: Springer, March 2009. DOI 10.1007/978-3-642-00427-8_6. ISSN 0302-9743 pp.
4514 86–105.
- 4515 [155] T. Jiang and A. E. Keating, “AVID: An integrative framework for discovering functional rela-
4516 tionship among proteins,” *BMC Bioinformatics*, vol. 6, no. 1, p. 136, 2005, DOI 10.1186/1471-
4517 2105-6-136. ISSN 1471-2105
- 4518 [156] B. Jimerson and B. Gregory, “Pivotal Cloud Foundry, Google ML, and Spring,” [Online]
4519 Available: <http://bit.ly/2RUBIIL>, San Francisco, CA, USA, December 2017.
- 4520 [157] Y. Jin, *Multi-Objective Machine Learning*, ser. Studies in Computational Intelligence. Berlin,
4521 Heidelberg: Springer, 2006. DOI 10.1007/3-540-33019-4. ISBN 978-3-54-030676-4
- 4522 [158] U. Johansson and L. Niklasson, “Evolving decision trees using oracle guides,” in *Proceedings*
4523 *of the 2009 IEEE Symposium on Computational Intelligence and Data Mining*. Nashville,
4524 TN, USA: IEEE, May 2009. DOI 10.1109/CIDM.2009.4938655. ISBN 978-1-42-442765-9
4525 pp. 238–244.
- 4526 [159] M. Jørgensen, T. Dybå, K. Liestøl, and D. I. K. Sjøberg, “Incorrect results in software engineer-
4527 ing experiments: How to improve research practices,” *Journal of Systems and Software*, vol.
4528 116, pp. 133–145, 2016, DOI 10.1016/j.jss.2015.03.065. ISSN 0164-1212
- 4529 [160] J. M. Juran, *Juran on Planning for Quality*. New York, NY, USA: The Free Press, 1988. ISBN
4530 978-0-02-916681-9
- 4531 [161] N. Juristo and O. S. Gómez, “Replication of software engineering experiments,” in *Proceedings*
4532 *of the LASER Summer School on Software Engineering*. Elba Island, Italy: Springer, 2011.
4533 DOI 10.1007/978-3-642-25231-0_2. ISBN 978-3-64-225230-3. ISSN 0302-9743 pp. 60–88.
- 4534 [162] N. Juristo and A. M. Moreno, *Basics of Software Engineering Experimentation*. Boston, MA,
4535 USA: Springer, March 2001. DOI 10.1007/978-1-4757-3304-4.
- 4536 [163] A. Karwath and R. D. King, “Homology induction: The use of machine learning to improve se-
4537 quence similarity searches,” *BMC Bioinformatics*, vol. 3, no. 1, p. 11, 2002, DOI 10.1186/1471-
4538 2105-3-11. ISSN 1471-2105
- 4539 [164] K. A. Kaufman and R. S. Michalski, “Learning from inconsistent and noisy data: The AQ18
4540 approach,” in *Proceedings of the 11th European Conference on Principles and Practice of*
4541 *Knowledge Discovery in Databases*, vol. 1609. Warsaw, Poland: Springer, September 1999.
4542 DOI 10.1007/BFb0095128. ISBN 3-540-65965-X. ISSN 1611-3349 pp. 411–419.
- 4543 [165] D. Kavaler, D. Posnett, C. Gibler, H. Chen, P. Devanbu, and V. Filkov, “Using and asking:
4544 APIs used in the Android market and asked about in StackOverflow,” in *Proceedings of the*
4545 *5th International Conference on Social Infomatics*. Kyoto, Japan: Springer, November 2013.
4546 DOI 10.1007/978-3-319-03260-3_35. ISBN 978-3-31-903259-7. ISSN 0302-9743 pp. 405–
4547 418.

- 4548 [166] B. Kim, "Interactive and Interpretable Machine Learning Models for Human Machine Collaboration," Ph.D. dissertation, Massachusetts Institute of Technology, 2015.
- 4549
- 4550 [167] B. Kim, C. Rudin, and J. Shah, "The Bayesian case model: A generative approach for case-based reasoning and prototype classification," in *Proceedings of the 28th Conference on Neural Information Processing Systems*, Montreal, QC, Canada, December 2014. ISSN 1049-5258 pp. 1952–1960.
- 4551
- 4552
- 4553
- 4554 [168] B. Kitchenham and S. Charters, "Guidelines for performing Systematic Literature Reviews in Software Engineering," Software Engineering Group, Keele University and Department of Computer Science, University of Durham, Keele, UK, Tech. Rep., 2007.
- 4555
- 4556
- 4557 [169] B. A. Kitchenham and S. L. Pfleeger, "Personal opinion surveys," in *Guide to Advanced Empirical Software Engineering*, F. Shull, J. Singer, and D. I. K. Sjøberg, Eds. Springer, November 2007, ch. 3, pp. 63–92. ISBN 978-1-84-800043-8
- 4558
- 4559
- 4560 [170] B. A. Kitchenham, T. Dybå, and M. Jorgensen, "Evidence-Based Software Engineering," in *Proceedings of the 26th International Conference on Software Engineering*. Edinburgh, Scotland, UK: IEEE, May 2004. ISBN 978-0-76-952163-3 pp. 273–281.
- 4561
- 4562
- 4563 [171] H. K. Klein and M. D. Myers, "A set of principles for conducting and evaluating interpretive field studies in information systems," *MIS Quarterly: Management Information Systems*, vol. 23, no. 1, pp. 67–94, 1999, DOI 10.2307/249410. ISSN 0276-7783
- 4564
- 4565
- 4566 [172] A. J. Ko and Y. Riche, "The role of conceptual knowledge in API usability," in *Proceedings of the 2011 IEEE Symposium on Visual Languages and Human Centric Computing*. Pittsburgh, PA, USA: IEEE, September 2011. DOI 10.1109/VLHCC.2011.6070395. ISBN 978-1-45771245-6 pp. 173–176.
- 4567
- 4568
- 4569
- 4570 [173] A. J. Ko, B. A. Myers, and H. H. Aung, "Six learning barriers in end-user programming systems," in *Proceedings of the 2004 IEEE Symposium on Visual Languages and Human Centric Computing*. Rome, Italy: IEEE, September 2004. DOI 10.1109/vlhcc.2004.47. ISBN 0-78-038696-5 pp. 199–206.
- 4571
- 4572
- 4573
- 4574 [174] I. Kononenko, "Inductive and bayesian learning in medical diagnosis," *Applied Artificial Intelligence*, vol. 7, no. 4, pp. 317–337, 1993, DOI 10.1080/08839519308949993. ISSN 1087-6545
- 4575
- 4576 [175] J. Kotula, "Using patterns to create component documentation," *IEEE Software*, vol. 15, no. 2, pp. 84–92, 1998, DOI 10.1109/52.663791. ISSN 0740-7459
- 4577
- 4578 [176] K. Krippendorff, *Content Analysis*, ser. An Introduction to Its Methodology. SAGE, 1980. ISBN 978-1-50-639566-1
- 4579
- 4580 [177] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017, DOI 10.1145/3065386. ISSN 1557-7317
- 4581
- 4582
- 4583 [178] J. A. Krosnick, "Survey Research," *Annual Review of Psychology*, vol. 50, no. 1, pp. 537–567, February 1999, DOI 10.1146/annurev.psych.50.1.537. ISSN 0066-4308
- 4584
- 4585 [179] A. Kurakin, I. J. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," in *Proceedings of the 5th International Conference on Learning Representations*, Toulon, France, April 2017.
- 4586
- 4587
- 4588 [180] G. Laforge, "Machine Intelligence at Google Scale," in *QCon*, London, England, UK, June 2018.
- 4589
- 4590 [181] H. Lakkaraju, S. H. Bach, and J. Leskovec, "Interpretable decision sets: A joint framework for description and prediction," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. San Francisco, CA, USA: ACM, August 2016. DOI 10.1145/2939672.2939874. ISBN 978-1-45-034232-2 pp. 1675–1684.
- 4591
- 4592
- 4593
- 4594 [182] J. R. Landis and G. G. Koch, "The Measurement of Observer Agreement for Categorical Data," *Biometrics*, vol. 33, no. 1, p. 159, March 1977, DOI 10.2307/2529310. ISSN 0006341X
- 4595
- 4596 [183] N. Lavrač, "Selected techniques for data mining in medicine," *Artificial Intelligence in Medicine*, vol. 16, no. 1, pp. 3–23, 1999, DOI 10.1016/S0933-3657(98)00062-1. ISSN 0933-3657
- 4597
- 4598 [184] T. Lei, R. Barzilay, and T. Jaakkola, "Rationalizing neural predictions," in *Proceedings of the 9th International Joint Conference on Natural Language Processing and Conference on Empirical Methods in Natural Language Processing*. Austin, TX, USA: Association for Computational Linguistics, November 2016. DOI 10.18653/v1/d16-1011. ISBN 978-1-94-562625-8 pp. 107–117.
- 4599
- 4600
- 4601
- 4602

- [4603] [185] T. C. Lethbridge, S. E. Sim, and J. Singer, "Studying software engineers: Data collection
[4604] techniques for software field studies," *Empirical Software Engineering*, vol. 10, no. 3, pp.
[4605] 311–341, July 2005, DOI 10.1007/s10664-005-1290-x. ISSN 1382-3256
- [4606] [186] R. J. Light, "Measures of response agreement for qualitative data: Some generalizations and
[4607] alternatives," *Psychological Bulletin*, vol. 76, no. 5, pp. 365–377, 1971, DOI 10.1037/h0031643.
[4608] ISSN 0033-2909
- [4609] [187] E. Lima, C. Mues, and B. Baesens, "Domain knowledge integration in data mining using
[4610] decision tables: Case studies in churn prediction," *Journal of the Operational Research Society*,
[4611] vol. 60, no. 8, pp. 1096–1106, 2009, DOI 10.1057/jors.2008.161. ISSN 0160-5682
- [4612] [188] B. Lin, F. Zampetti, G. Bavota, M. Di Penta, M. Lanza, and R. Oliveto, "Sentiment anal-
[4613] ysis for software engineering: How far can we go?" in *Proceedings of the 40th Inter-*
[4614] *national Conference on Software Engineering*. Gothenburg, Sweden: ACM, May 2018.
[4615] DOI 10.1145/3180155.3180195, pp. 94–104.
- [4616] [189] M. Linares-Vásquez, G. Bavota, M. Di Penta, R. Oliveto, and D. Poshyvanyk, "How do API
[4617] changes trigger stack overflow discussions? A study on the android SDK," in *Proceedings of*
[4618] *the 22nd International Conference on Program Comprehension*. Hyderabad, India: ACM,
[4619] June 2014. DOI 10.1145/2597008.2597155. ISBN 978-1-4503-2879-1 pp. 83–94.
- [4620] [190] Z. C. Lipton, "The mythos of model interpretability," *Communications of the ACM*, vol. 61,
[4621] no. 10, pp. 35–43, 2018, DOI 10.1145/3233231. ISSN 1557-7317
- [4622] [191] M. Litwin, *How to Measure Survey Reliability and Validity*. Thousand Oaks, CA, USA: SAGE,
[4623] 1995, vol. 7, DOI 10.4135/9781483348957. ISBN 978-0-80-395704-6
- [4624] [192] Y. Liu, T. Kohlberger, M. Norouzi, G. E. Dahl, J. L. Smith, A. Mohtashamian, N. Olson,
[4625] L. H. Peng, J. D. Hipp, and M. C. Stumpe, "Artificial Intelligence-Based Breast Cancer Nodal
[4626] Metastasis Detection." *Archives of Pathology & Laboratory Medicine*, vol. 143, no. 7, pp.
[4627] 859–868, July 2017, DOI 10.5858/arpa.2018-0147-OA. ISSN 1543-2165
- [4628] [193] D. Lo Giudice, C. Mines, A. LeClair, R. Curran, and A. Homan, "How AI Will Change
[4629] Software Development And Applications," [Online] Available: <http://bit.ly/38RiAIN>, Forrester
[4630] Research, Inc., Tech. Rep., November 2016.
- [4631] [194] R. Lori and M. Oded, *Data mining with decision trees*. World Scientific Publishing Company,
[4632] 2008, vol. 69. ISBN 978-9-81-277171-1
- [4633] [195] R. E. Lyons and W. Vanderkulk, "The Use of Triple-Modular Redundancy to Improve Computer
[4634] Reliability," *IBM Journal of Research and Development*, vol. 6, no. 2, pp. 200–209, April 2010,
[4635] DOI 10.1147/rd.62.0200. ISSN 0018-8646
- [4636] [196] W. Maalej and M. P. Robillard, "Patterns of knowledge in API reference documentation," *IEEE
[4637] Transactions on Software Engineering*, 2013, DOI 10.1109/TSE.2013.12. ISSN 0098-5589
- [4638] [197] G. Malgieri and G. Comandé, "Why a right to legibility of automated decision-making exists
[4639] in the general data protection regulation," *International Data Privacy Law*, vol. 7, no. 4, pp.
[4640] 243–265, June 2017, DOI 10.1093/idpl/ixp019. ISSN 2044-4001
- [4641] [198] L. Mandel, "Describe REST Web services with WSDL 2.0," [Online] Available: <https://ibm.co/313RoNV>, May 2008, Accessed: 28 August 2018.
- [4642] [199] T. E. Marshall and S. L. Lambert, "Cloud-based intelligent accounting applications: Accounting
[4643] task automation using IBM watson cognitive computing," *Journal of Emerging Technologies
[4644] in Accounting*, vol. 15, no. 1, pp. 199–215, 2018, DOI 10.2308/jeta-52095. ISSN 1558-7940
- [4645] [200] D. Martens, J. Vanthienen, W. Verbeke, and B. Baesens, "Performance of classification mod-
[4646] els from a user perspective," *Decision Support Systems*, vol. 51, no. 4, pp. 782–793, 2011,
[4647] DOI 10.1016/j.dss.2011.01.013. ISSN 0167-9236
- [4648] [201] A. I. Martins, A. F. Rosa, A. Queirós, A. Silva, and N. P. Rocha, "European Portuguese
[4649] Validation of the System Usability Scale (SUS)," in *Procedia Computer Science*, 2015,
[4650] DOI 10.1016/j.procs.2015.09.273. ISSN 18770509
- [4651] [202] P. Mayring, "Mixing Qualitative and Quantitative Methods," in *Mixed Methodology in Psycho-
[4652] logical Research*. Sense Publishers, 2007, ch. 6, pp. 27–36. ISBN 978-9-07-787473-8
- [4653] [203] J. A. McCall, P. K. Richards, and G. F. Walters, "Factors in Software Quality: Concept and
[4654] Definitions of Software Quality," General Electric Company, Griffiss Air Force Base, NY, USA,
[4655] Tech. Rep. RADC-TR-77-369, November 1977.
- [4656] [204] J. McCarthy, "Programs with common sense," in *Proceedings of the Symposium on the Mech-
[4657] anization of Thought Processes*, Cambridge, MA, USA, 1963, pp. 1–15.
- [4658]

- 4659 [205] B. McGowen, "Machine learning with Google APIs," [Online] Available: <http://bit.ly/3aUQpo2>, January 2019.
- 4660 [206] M. L. McHugh, "Interrater reliability: The kappa statistic," *Biochimia Medica*, vol. 22, no. 3, pp. 276–282, 2012, DOI 10.11613/bm.2012.031. ISSN 1330-0962
- 4661 [207] S. G. McLellan, A. W. Roesler, J. T. Tempest, and C. I. Spinuzzi, "Building more usable APIs," *IEEE Software*, vol. 15, no. 3, pp. 78–86, 1998, DOI 10.1109/52.676963. ISSN 0740-7459
- 4662 [208] L. McLeod and S. G. MacDonell, "Factors that affect software systems development project outcomes: A survey of research," *ACM Computing Surveys*, vol. 43, no. 4, p. 24, 2011, DOI 10.1145/1978802.1978803. ISSN 0360-0300
- 4663 [209] J. Meltzoff and H. Cooper, *Critical thinking about research: Psychology and related fields*, 2nd ed. American Psychological Association, 2018. DOI 10.1037/0000052-000.
- 4664 [210] 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, August 2018, DOI 10.1177/0047281617721853. ISSN 1541-3780
- 4665 [211] T. Mens and S. Demeyer, *Software Evolution*. Berlin, Heidelberg: Springer, 2008. DOI 10.1007/978-3-540-76440-3. ISBN 978-3-54-076439-7
- 4666 [212] T. Mens, S. Demeyer, M. Wermelinger, R. Hirschfeld, S. Ducasse, and M. Jazayeri, "Challenges in software evolution," in *Proceedings of the 8th International Workshop on Principles of Software Evolution*, vol. 2005. Lisbon, Portugal: IEEE, September 2005. DOI 10.1109/I-WPSE.2005.7. ISBN 0-76-952349-8. ISSN 1550-4077 pp. 13–22.
- 4667 [213] A. C. Michalos and H. A. Simon, *The Sciences of the Artificial*. MIT press, 1970, vol. 11, no. 1, DOI 10.2307/3102825.
- 4668 [214] D. Michie, "Machine learning in the next five years," in *Proceedings of the 3rd European Conference on European Working Session on Learning*. Glasgow, Scotland, UK: Pitman Publishing, Inc., October 1988. ISBN 978-0-27-308800-4 pp. 107–122.
- 4669 [215] G. A. Miller, "WordNet: A Lexical Database for English," *Communications of the ACM*, vol. 38, no. 11, pp. 39–41, November 1995, DOI 10.1145/219717.219748. ISSN 1557-7317
- 4670 [216] D. Moody, "The physics of notations: Toward a scientific basis for constructing visual notations in software engineering," *IEEE Transactions on Software Engineering*, vol. 35, no. 6, pp. 756–779, 2009, DOI 10.1109/TSE.2009.67. ISSN 0098-5589
- 4671 [217] A. Murgia, P. Tourani, B. Adams, and M. Ortú, "Do developers feel emotions? an exploratory analysis of emotions in software artifacts," in *Proceedings of the 11th Working Conference on Mining Software Repositories*. Hyderabad, India: ACM, May 2014. DOI 10.1145/2597073.2597086, pp. 262–271.
- 4672 [218] B. A. Myers, S. Y. Jeong, Y. Xie, J. Beaton, J. Stylos, R. Ehret, J. Karstens, A. Efeoglu, and D. K. Busse, "Studying the Documentation of an API for Enterprise Service-Oriented Architecture," *Journal of Organizational and End User Computing*, vol. 22, no. 1, pp. 23–51, January 2010, DOI 10.4018/joeuc.2010101903. ISSN 1546-2234
- 4673 [219] C. Myers, A. Furqan, J. Nebolsky, K. Caro, and J. Zhu, "Patterns for how users overcome obstacles in Voice User Interfaces," in *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, vol. 2018-April. Montreal, QC, Canada: ACM, April 2018. DOI 10.1145/3173574.3173580. ISBN 978-1-45-035620-6 p. 6.
- 4674 [220] S. Nakajima, "Model-Checking Verification for Reliable Web Service," in *Proceedings of the First International Symposium on Cyber World*. Montreal, QC, Canada: IEEE, November 2002. ISBN 978-0-76-951862-6 pp. 378–385.
- 4675 [221] M. Narayanan, E. Chen, J. He, B. Kim, S. Gershman, and F. Doshi-Velez, "How do Humans Understand Explanations from Machine Learning Systems? An Evaluation of the Human-Interpretability of Explanation," *IEEE Transactions on Evolutionary Computation*, 2018, In Press.
- 4676 [222] S. Narayanan and S. A. McIlraith, "Simulation, verification and automated composition of web services," in *Proceedings of the 11th International Conference on World Wide Web*. Honolulu, HI, USA: ACM, May 2002. DOI 10.1145/511446.511457. ISBN 1-58-113449-5 pp. 77–88.
- 4677 [223] B. J. Nelson, "Remote Procedure Call," Ph.D. dissertation, Carnegie Mellon University, 1981.
- 4678 [224] H. F. Niemeyer and A. C. Niemeyer, "Apportionment methods," *Mathematical Social Sciences*, vol. 56, no. 2, pp. 240–253, 2008. ISSN 0165-4896

- [225] N. Novielli, F. Calefato, and F. Lanubile, "The challenges of sentiment detection in the social programmer ecosystem," in *Proceedings of the 7th International Workshop on Social Software Engineering*, Bergamo, Italy, August 2015, DOI 10.1145/2804381.2804387. ISBN 978-1-45-033818-9 pp. 33–40.
- [226] ——, "A gold standard for emotion annotation in Stack Overflow," in *Proceedings of the 15th International Conference on Mining Software Repositories*, IEEE. Gothenburg, Sweden: ACM, May 2018. DOI 10.1145/3196398.3196453, pp. 14–17.
- [227] K. Nybom, A. Ashraf, and I. Porres, "A systematic mapping study on API documentation generation approaches," in *Proceedings of the 44th Euromicro Conference on Software Engineering and Advanced Applications*. Prague, Czech Republic: IEEE, August 2018. DOI 10.1109/SEAA.2018.00081. ISBN 978-1-53-867382-9 pp. 462–469.
- [228] 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," in *Proceedings of the 20th Annual International Conference on Computer Documentation*. Toronto, ON, Canada: ACM, October 2002. DOI 10.1145/584955.584976, pp. 133–141.
- [229] T. Ohtake, A. Cummaudo, M. Abdelrazek, R. Vasa, and J. Grundy, "Merging intelligent API responses using a proportional representation approach," in *Proceedings of the 19th International Conference on Web Engineering*. Daejeon, Republic of Korea: Springer, June 2019. DOI 10.1007/978-3-030-19274-7_28. ISBN 978-3-03-019273-0. ISSN 1611-3349 pp. 391–406.
- [230] Open Software Foundation, "Part 3: DCE Remote Procedure Call (RPC)," in *OSF DCE application development guide: revision 1.0*. Prentice Hall, December 1991.
- [231] N. Oreskes, K. Shrader-Frechette, and K. Belitz, "Verification, validation, and confirmation of numerical models in the earth sciences," *Science*, vol. 263, no. 5147, pp. 641–646, 1994, DOI 10.1126/science.263.5147.641. ISSN 0036-8075
- [232] A. L. M. Ortiz, "Curating Content with Google Machine Learning Application Programming Interfaces," in *EIAPortugal*, July 2017.
- [233] M. Ortú, A. Murgia, G. Destefanis, P. Tourani, R. Tonelli, M. Marchesi, and B. Adams, "The emotional side of software developers in JIRA," in *Proceedings of the 13th International Conference on Mining Software Repositories*, ACM. Austin, TX, USA: ACM, May 2016. DOI 10.1145/2901739.2903505, pp. 480–483.
- [234] F. E. B. Otero and A. A. Freitas, "Improving the interpretability of classification rules discovered by an ant colony algorithm: Extended results," in *Evolutionary Computation*, vol. 24, no. 3. ACM, 2016. DOI 10.1162/EVCO_a_00155. ISSN 1530-9304 pp. 385–409.
- [235] A. Pal, S. Chang, and J. A. Konstan, "Evolution of experts in question answering communities," in *Proceedings of the 6th International AAAI Conference on Weblogs and Social Media*. Dublin, Ireland: AAAI, June 2012. ISBN 978-1-57-735556-4 pp. 274–281.
- [236] R. Parekh, "Designing AI at Scale to Power Everyday Life," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. Halifax, NS, Canada: ACM, August 2017. DOI 10.1145/3097983.3105815, p. 27.
- [237] D. L. Parnas and S. A. Vilkomir, "Precise documentation of critical software," in *Proceedings of 10th IEEE International Symposium on High Assurance Systems Engineering*. Plano, TX, USA: IEEE, November 2007. DOI 10.1109/HASE.2007.63. ISSN 1530-2059 pp. 237–244.
- [238] K. Patel, J. Fogarty, J. A. Landay, and B. Harrison, "Investigating statistical machine learning as a tool for software development," in *Proceedings of the 26th SIGCHI Conference on Human Factors in Computing Systems*, ser. CHI '08. Florence, Italy: ACM, April 2008. DOI 10.1145/1357054.1357160. ISBN 978-1-60-558011-1 pp. 667–676.
- [239] C. Pautasso, O. Zimmermann, and F. Leymann, "RESTful web services vs. "Big" web services: Making the right architectural decision," in *Proceedings of the 17th International Conference on World Wide Web*. Beijing, China: ACM, April 2008. DOI 10.1145/1367497.1367606. ISBN 978-1-60-558085-2
- [240] M. Pazzani, "Comprehensible knowledge discovery: gaining insight from data," in *Proceedings of the First Federal Data Mining Conference and Exposition*, Washington, DC, USA, 1997, pp. 73–82.

- 4768 [241] M. J. Pazzani, S. Mani, and W. R. Shankle, "Acceptance of rules generated by machine learning
4769 among medical experts," *Methods of Information in Medicine*, vol. 40, no. 5, pp. 380–385,
4770 2001, DOI 10.1055/s-0038-1634196. ISSN 0026-1270
- 4771 [242] J. Pearl, "The seven tools of causal inference, with reflections on machine learning," *Communications of the ACM*, vol. 62, no. 3, pp. 54–60, 2019, DOI 10.1145/3241036. ISSN 1557-7317
- 4772 [243] K. Petersen and C. Gencel, "Worldviews, research methods, and their relationship to validity
4773 in empirical software engineering research," in *Proceedings of the Joint Conference of the
4774 23rd International Workshop on Software Measurement and the 8th International Conference
4775 on Software Process and Product Measurement*. Ankara, Turkey: IEEE, October 2013.
4776 DOI 10.1109/IWSM-Mensura.2013.22. ISBN 978-0-76-955078-7 pp. 81–89.
- 4777 [244] K. Petersen, R. Feldt, S. Mujtaba, and M. Mattsson, "Systematic mapping studies in software en-
4778 gineering," in *Proceedings of the 12th International Conference on Evaluation and Assessment
4779 in Software Engineering, EASE 2008*, 2008, DOI 10.14236/ewic/ease2008.8, pp. 68–77.
- 4780 [245] Z. Pezzementi, T. Tabor, S. Yim, J. K. Chang, B. Drozd, D. Guttendorf, M. Wagner, and
4781 P. Koopman, "Putting Image Manipulations in Context: Robustness Testing for Safe Perception,"
4782 in *Proceedings of the 15th IEEE International Symposium on Safety, Security, and Rescue
4783 Robotics*. Philadelphia, PA, USA: IEEE, August 2018. DOI 10.1109/SSRR.2018.8468619.
4784 ISBN 978-1-53-865572-6 pp. 1–8.
- 4785 [246] H. Pham, *System Software Reliability*, 1st ed. Springer, 2000. ISBN 978-1-84-628295-9
- 4786 [247] M. Piccioni, C. A. Furia, and B. Meyer, "An empirical study of API usability," in *Proceedings
4787 of the 13th International Symposium on Empirical Software Engineering and Measurement*.
4788 Baltimore, MD, USA: IEEE, October 2013. DOI 10.1109/ESEM.2013.14. ISSN 1949-3770
4789 pp. 5–14.
- 4790 [248] R. M. Pirsig, *Zen and the art of motorcycle maintenance: An inquiry into values*, 1st ed.
4791 HarperTorch, 1974. ISBN 9-780-06-058946-2
- 4792 [249] R. S. Pressman, *Software Engineering: A Practitioner's Approach*, 8th ed. McGraw-Hill,
4793 2005. ISBN 978-0-07-802212-8
- 4794 [250] J. R. Quinlan, "Some elements of machine learning," in *Proceedings of the 9th International
4795 Workshop on Inductive Logic Programming*, vol. 1634. Bled, Slovenia: Springer, June 1999.
4796 DOI 10.1007/3-540-48751-4_3. ISBN 3-54-066109-3. ISSN 1611-3349 pp. 15–18.
- 4797 [251] ———, *C4.5: Programs for machine learning*. San Francisco, CA, USA: Morgan Kauffmann,
4798 1993. ISBN 978-1-55-860238-0
- 4799 [252] N. Rama Suri, V. S. Srinivas, and M. Narasimha Murty, "A cooperative game theoretic approach
4800 to prototype selection," in *Proceedings of the 11th European Conference on Principles and
4801 Practice of Knowledge Discovery in Databases*. Warsaw, Poland: Springer, September 2007.
4802 DOI 10.1007/978-3-540-74976-9_58. ISBN 978-3-54-074975-2. ISSN 0302-9743 pp. 556–
4803 564.
- 4804 [253] M. Reboucas, G. Pinto, F. Ebert, W. Torres, A. Serebrenik, and F. Castor, "An Empirical Study
4805 on the Usage of the Swift Programming Language," in *Proceedings of the 23rd International
4806 Conference on Software Analysis, Evolution, and Reengineering*. Suita, Japan: IEEE, March
4807 2016. DOI 10.1109/saner.2016.66, pp. 634–638.
- 4808 [254] A. Reis, D. Paulino, V. Filipe, and J. Barroso, "Using online artificial vision services to assist
4809 the blind - An assessment of Microsoft Cognitive Services and Google Cloud Vision," *Advances
4810 in Intelligent Systems and Computing*, vol. 746, no. 12, pp. 174–184, 2018, DOI 10.1007/978-
4811 3-31-77712-2_17. ISBN 978-3-31-977711-5. ISSN 2194-5357
- 4812 [255] M. T. Ribeiro, S. Singh, and C. Guestrin, "'Why Should I Trust You?': Explaining the Predic-
4813 tions of Any Classifier," in *Proceedings of the 22nd ACM SIGKDD International Conference
4814 on Knowledge Discovery and Data Mining*. San Francisco, CA, USA: ACM, August 2016.
4815 DOI 2939672.2939778, pp. 1135–1144.
- 4816 [256] M. Ribeiro, K. Grolinger, and M. A. M. Capretz, "MLaaS: Machine learning as a service,"
4817 in *Proceedings of the 14th International Conference on Machine Learning and Applications*.
4818 Miami, FL, USA: IEEE, December 2015. DOI 10.1109/ICMLA.2015.152. ISBN 978-1-50-
4819 900287-0 pp. 896–902.
- 4820 [257] G. Richards, V. J. Rayward-Smith, P. H. Sönksen, S. Carey, and C. Weng, "Data mining for
4821 indicators of early mortality in a database of clinical records," *Artificial Intelligence in Medicine*,
4822 vol. 22, no. 3, pp. 215–231, 2001, DOI 10.1016/S0933-3657(00)00110-X. ISSN 0933-3657
- 4823

- [258] G. Ridgeway, D. Madigan, T. Richardson, and J. O’Kane, “Interpretable Boosted Naïve Bayes Classification,” in *Proceedings of the 4th International Conference on Knowledge Discovery and Data Mining*. New York, NY, USA: AAAI, 1998, pp. 101–104.
- [259] RightScale Inc., “State of the Cloud Report: DevOps Trends,” Tech. Rep., 2016.
- [260] G. Ritzer and E. Guba, “The Paradigm Dialog,” *Canadian Journal of Sociology*, vol. 16, no. 4, p. 446, 1991, DOI 10.2307/3340973. ISSN 0318-6431
- [261] M. P. Robillard, “What makes APIs hard to learn? Answers from developers,” *IEEE Software*, vol. 26, no. 6, pp. 27–34, 2009, DOI 10.1109/MS.2009.193. ISSN 0740-7459
- [262] M. P. Robillard and R. Deline, “A field study of API learning obstacles,” *Empirical Software Engineering*, vol. 16, no. 6, pp. 703–732, 2011, DOI 10.1007/s10664-010-9150-8. ISSN 1382-3256
- [263] M. P. Robillard, A. Marcus, C. Treude, G. Bavota, O. Chaparro, N. Ernst, M. A. Gerostall, M. Godfrey, M. Lanza, M. Linares-Vásquez, G. C. Murphy, L. Moreno, D. Shepherd, and E. Wong, “On-demand developer documentation,” in *Proceedings of the 33rd IEEE International Conference on Software Maintenance and Evolution*. Shanghai, China: IEEE, September 2017. DOI 10.1109/ICSME.2017.17, pp. 479–483.
- [264] H. Robinson, J. Segal, and H. Sharp, “Ethnographically-informed empirical studies of software practice,” *Information and Software Technology*, vol. 49, no. 6, pp. 540–551, 2007, DOI 10.1016/j.infsof.2007.02.007. ISSN 0950-5849
- [265] C. Rosen and E. Shihab, “What are mobile developers asking about? A large scale study using stack overflow,” *Empirical Software Engineering*, vol. 21, no. 3, pp. 1192–1223, 2016, DOI 10.1007/s10664-015-9379-3. ISSN 1573-7616
- [266] W. Rosenberry, D. Kenney, and G. Fisher, *Understanding DCE*. O’Reilly & Associates, Inc., 1992. ISBN 978-1-56-592005-7
- [267] A. Rosenfeld, R. Zemel, and J. K. Tsotsos, “The Elephant in the Room,” *arXiv preprint arXiv:1808.03305*, 2018.
- [268] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, “Right for the right reasons: Training differentiable models by constraining their explanations,” in *Proceedings of the 26th International Joint Conferences on Artificial Intelligence*, Melbourne, Australia, August 2017, DOI 10.24963/ijcai.2017/371. ISBN 978-0-99-924110-3. ISSN 1045-0823 pp. 2662–2670.
- [269] R. J. Rubey and R. D. Hartwick, “Quantitative measurement of program quality,” in *Proceedings of the 1968 23rd ACM National Conference*. Las Vegas, NV, USA: ACM, August 1968. DOI 10.1145/800186.810631. ISBN 978-1-45-037486-6 pp. 671–677.
- [270] J. Rutten, M. Kwiatkowska, G. Norman, and D. Parker, *Mathematical techniques for analyzing concurrent and probabilistic systems*, ser. CRM Monograph Series, P. Panangaden and F. van Breugel, Eds. American Mathematical Society, 2004, vol. 23.
- [271] M. E. C. Santos, J. Polvi, T. Taketomi, G. Yamamoto, C. Sandor, and H. Kato, “Usability scale for handheld augmented reality,” in *Proceedings of the ACM Symposium on Virtual Reality Software and Technology, VRST*, 2014, DOI 10.1145/2671015.2671019. ISBN 9781450332538
- [272] J. Sauro and J. R. Lewis, “When designing usability questionnaires, does it hurt to be positive?” in *Proceedings of the 2011 SIGCHI Conference on Human Factors in Computing Systems*, Vancouver, BC, Canada, May 2011, DOI 10.1145/1978942.1979266, pp. 2215–2223.
- [273] M. Schwabacher and P. Langley, “Discovering communicable scientific knowledge from spatio-temporal data,” in *Proceedings of the 18th International Conference on Machine Learning*. Williamstown, MA, USA: Morgan Kaufmann, June 2001. ISBN 978-1-55-860778-1 pp. 489–496.
- [274] A. Schwaighofer and N. D. Lawrence, *Dataset shift in machine learning*, J. Quiñonero-Candela and M. Sugiyama, Eds. Cambridge, MA, USA: The MIT Press, 2008. ISBN 978-0-26-217005-5
- [275] T. A. Schwandt, “Qualitative data analysis: An expanded sourcebook,” *Evaluation and Program Planning*, vol. 19, no. 1, pp. 106–107, 1996, DOI 10.1016/0149-7189(96)88232-2. ISSN 0149-7189
- [276] D. Sculley, M. E. Otey, M. Pohl, B. Spitznagel, J. Hainsworth, and Y. Zhou, “Detecting adversarial advertisements in the wild,” in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM. San Diego, CA, USA: ACM, August 2011. DOI 10.1145/2020408.2020455, pp. 274–282.

- 4880 [277] D. Sculley, G. Holt, D. Golovin, E. Davydov, T. Phillips, D. Ebner, V. Chaudhary, M. Young, J. F.
4881 Crespo, and D. Dennison, “Hidden technical debt in machine learning systems,” in *Proceedings*
4882 of the 29th Conference on Neural Information Processing Systems, Montreal, QC, Canada,
4883 December 2015. ISBN 0262017091, 9780262017091. ISSN 1049-5258 pp. 2503–2511.
- 4884 [278] C. B. Seaman, “Qualitative methods,” in *Guide to Advanced Empirical Software Engineering*,
4885 F. Shull, J. Singer, and D. I. K. Sjøberg, Eds. Springer, November 2007, ch. 2, pp. 35–62.
4886 ISBN 978-1-84-800043-8
- 4887 [279] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, “Grad-CAM:
4888 Visual Explanations from Deep Networks via Gradient-Based Localization,” *International*
4889 *Journal of Computer Vision*, pp. 618–626, 2019, DOI 10.1007/s11263-019-01228-7. ISSN
4890 1573-1405
- 4891 [280] S. Sen and L. Knight, “A genetic prototype learner,” in *Proceedings of the International Joint*
4892 *Conference on Artificial Intelligence*. Montreal, QC, Canada: Morgan Kaufmann, August
4893 1995, pp. 725–733.
- 4894 [281] C. E. Shannon and W. Weaver, “The mathematical theory of communication,” *The Bell*
4895 *System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948, DOI 10.1002/j.1538-
4896 7305.1948.tb01338.x.
- 4897 [282] P. Shaver, J. Schwartz, D. Kirson, and C. O’Connor, “Emotion knowledge: Further exploration
4898 of a prototype approach,” *Journal of Personality and Social Psychology*, vol. 52, no. 6, pp.
4899 1061–1086, 1987, DOI 10.1037/0022-3514.52.6.1061.
- 4900 [283] M. Shaw, “Writing good software engineering research papers,” in *Proceedings of the 25th*
4901 *International Conference on Software Engineering*. Portland, OR, USA: IEEE, May 2003.
4902 ISBN 978-0-76-951877-0 pp. 726–736.
- 4903 [284] M. Shepperd, “Replication studies considered harmful,” in *Proceedings of the 40th Interna-*
4904 *tional Conference on Software Engineering*. Gothenburg, Sweden: ACM, May 2018.
4905 DOI 10.1145/3183399.3183423. ISBN 978-1-45-035662-6. ISSN 0270-5257 pp. 73–76.
- 4906 [285] D. J. Sheskin, *Handbook of Parametric and Nonparametric Statistical Procedures*. New York,
4907 NY, USA: Chapman and Hall/CRC, 2004. DOI 10.4324/9780203489536.
- 4908 [286] L. Si and J. Callan, “A semisupervised learning method to merge search engine results,”
4909 *ACM Transactions on Information Systems*, vol. 21, no. 4, pp. 457–491, October 2003,
4910 DOI 10.1145/944012.944017. ISSN 1046-8188
- 4911 [287] J. Singer, S. E. Sim, and T. C. Lethbridge, “Software engineering data collection for field
4912 studies,” in *Guide to Advanced Empirical Software Engineering*, F. Shull, J. Singer, and D. I. K.
4913 Sjøberg, Eds. Springer, November 2007, ch. 1, pp. 9–34. ISBN 978-1-84-800043-8
- 4914 [288] S. Singh, M. T. Ribeiro, and C. Guestrin, “Programs as Black-Box Explanations,” *arXiv preprint*
4915 *arXiv:1611.07579*, November 2016.
- 4916 [289] V. S. Sinha, S. Mani, and M. Gupta, “Exploring activeness of users in QA forums,” in *Proce-*
4917 *edings of the 10th Working Conference on Mining Software Repositories*. San Francisco, CA,
4918 USA: IEEE, May 2013. DOI 10.1109/MSR.2013.6624010. ISBN 978-1-46-732936-1. ISSN
4919 2160-1852 pp. 77–80.
- 4920 [290] M. Sokolova and G. Lapalme, “A systematic analysis of performance measures for classifi-
4921 cation tasks,” *Information Processing and Management*, vol. 45, no. 4, pp. 427–437, 2009,
4922 DOI 10.1016/j.ipm.2009.03.002. ISSN 0306-4573
- 4923 [291] I. Sommerville, *Software Engineering*, 9th ed. Boston, MA, USA: Addison-Wesley, 2011.
4924 ISBN 978-0-13-703515-1
- 4925 [292] P. Spector, *Summated Rating Scale Construction*. Newbury Park, CA, USA: SAGE, 1992.
4926 DOI 10.4135/9781412986038. ISBN 978-0-80-394341-4
- 4927 [293] R. Stevens, J. Ganz, V. Filkov, P. Devanbu, and H. Chen, “Asking for (and about) permissions
4928 used by Android apps,” in *Proceedings of the 10th Working Conference on Mining Software*
4929 *Repositories*. San Francisco, CA, USA: IEEE, May 2013. ISBN 978-1-46-732936-1 pp. 31–40.
- 4930 [294] J. Su, D. V. Vargas, and K. Sakurai, “One Pixel Attack for Fooling Deep Neural
4931 Networks,” *IEEE Transactions on Evolutionary Computation*, vol. 23, no. 5, pp. 828–841, 2019,
4932 DOI 10.1109/TEVC.2019.2890858. ISSN 1941-0026
- 4933 [295] G. H. Subramanian, J. Nosek, S. P. Raghunathan, and S. S. Kanitkar, “A comparison of
4934 the decision table and tree,” *Communications of the ACM*, vol. 35, no. 1, pp. 89–94, 1992,
4935 DOI 10.1145/129617.129621. ISSN 1557-7317

- [296] S. Subramanian, L. Inozemtseva, and R. Holmes, "Live API documentation," in *Proceedings of the 36th International Conference on Software Engineering*. Hyderabad, India: ACM, May 2014. DOI 10.1145/2568225.2568313. ISSN 0270-5257 pp. 643–652.
- [297] S. Sun, W. Pan, and L. L. Wang, "A Comprehensive Review of Effect Size Reporting and Interpreting Practices in Academic Journals in Education and Psychology," *Journal of Educational Psychology*, vol. 102, no. 4, pp. 989–1004, 2010, DOI 10.1037/a0019507. ISSN 0022-0663
- [298] D. Szafron, P. Lu, R. Greiner, D. S. Wishart, B. Poulin, R. Eisner, Z. Lu, J. Anvik, C. Macdonell, A. Fyshe, and D. Meeuwis, "Proteome Analyst: Custom predictions with explanations in a web-based tool for high-throughput proteome annotations," *Nucleic Acids Research*, vol. 32, 2004, DOI 10.1093/nar/gkh485. ISSN 0305-1048
- [299] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," in *Proceedings of the 2nd International Conference on Learning Representations*. Banff, AB, Canada: ACM, April 2014.
- [300] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the Inception Architecture for Computer Vision," in *Proceedings of the 2016 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. Las Vegas, NV, USA: IEEE, June 2016. DOI 10.1109/CVPR.2016.308. ISBN 978-1-46-738850-4. ISSN 1063-6919 pp. 2818–2826.
- [301] M. B. W. Tabor, "Student Proves That S.A.T. Can Be: (D) Wrong," [Online] Available: <https://nyti.ms/2UiKrrd>, New York, NY, USA, February 1997.
- [302] A. Tahir, A. Yamashita, S. Licorish, J. Dietrich, and S. Counsell, "Can you tell me if it smells? A study on how developers discuss code smells and anti-patterns in Stack Overflow," in *Proceedings of the 22nd International Conference on Evaluation and Assessment in Software Engineering*. Christchurch, New Zealand: ACM, June 2018. DOI 10.1145/3210459.3210466. ISBN 978-1-45-036403-4 pp. 68–78.
- [303] G. Tassey, *The economic impacts of inadequate infrastructure for software testing*. National Institute of Standards and Technology, September 2002. DOI 10.1080/10438590500197315. ISBN 978-0-75-672618-8
- [304] 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, June 2004, DOI 10.1016/j.infsof.2003.10.004. ISSN 0950-5849
- [305] R. S. Taylor, "Question-Negotiation and Information Seeking in Libraries," *College and Research Libraries*, vol. 29, no. 3, 1968, DOI 10.5860/crl_29_03_178.
- [306] S. W. Thomas, B. Adams, A. E. Hassan, and D. Blostein, "Studying software evolution using topic models," *Science of Computer Programming*, vol. 80, pp. 457–479, 2014, DOI 10.1016/j.scico.2012.08.003. ISSN 0167-6423
- [307] S. Thrun, "Is Learning The n-th Thing Any Easier Than Learning The First?" in *Proceedings of the 8th International Conference on Neural Information Processing Systems*. Denver, CO, USA: MIT Press, November 1996. ISSN 1049-5258 p. 7.
- [308] C. Treude, O. Barzilay, and M. A. Storey, "How do programmers ask and answer questions on the web?" in *Proceedings of the 33rd International Conference on Software Engineering*. Honolulu, HI, USA: ACM, May 2011. DOI 10.1145/1985793.1985907. ISBN 978-1-45-030445-0. ISSN 0270-5257 pp. 804–807.
- [309] G. Uddin and F. Khomh, "Automatic Mining of Opinions Expressed About APIs in Stack Overflow," *IEEE Transactions on Software Engineering*, February 2019, In Press, DOI 10.1109/TSE.2019.2900245. ISSN 1939-3520
- [310] G. Uddin and M. P. Robillard, "How API Documentation Fails," *IEEE Software*, vol. 32, no. 4, pp. 68–75, June 2015, DOI 10.1109/MS.2014.80. ISSN 0740-7459
- [311] M. Usman, R. Britto, J. Börstler, and E. Mendes, "Taxonomies in software engineering: A Systematic mapping study and a revised taxonomy development method," *Information and Software Technology*, vol. 85, pp. 43–59, May 2017, DOI 10.1016/j.infsof.2017.01.006. ISSN 0950-5849
- [312] A. Van Assche and H. Blockeel, "Seeing the forest through the trees learning a comprehensible model from a first order ensemble," in *Proceedings of the 17th International Conference on Inductive Logic Programming*. Corvallis, OR, USA: Springer, June 2007. DOI 10.1007/978-3-540-78469-2_26. ISBN 3-540-078468-3. ISSN 0302-9743 pp. 269–279.

- 4991 [313] B. Venners, "Design by Contract: A Conversation with Bertrand Meyer," *Artima Developer*,
4992 2003.
- 4993 [314] W. Verbeke, D. Martens, C. Mues, and B. Baesens, "Building comprehensible customer churn
4994 prediction models with advanced rule induction techniques," *Expert Systems with Applications*,
4995 vol. 38, no. 3, pp. 2354–2364, 2011, DOI 10.1016/j.eswa.2010.08.023. ISSN 0957-4174
- 4996 [315] F. Wachter, Mitterstadt, "EU regulations on algorithmic decision-making and a "right to explanation"," in *Proceedings of the 2016 ICML Workshop on Human Interpretability in Machine*
4997 *Learning*, New York, NY, USA, June 2016, pp. 26–30.
- 4998 [316] B. Wang, Y. Yao, B. Viswanath, H. Zheng, and B. Y. Zhao, "With great training comes great
5000 vulnerability: Practical attacks against transfer learning," in *Proceedings of the 27th USENIX*
5001 *Security Symposium*. Baltimore, MD, USA: USENIX Association, July 2018. ISBN 978-1-
5002 93-913304-5 pp. 1281–1297.
- 5003 [317] K. Wang, *Cloud Computing for Machine Learning and Cognitive Applications: A Machine*
5004 *Learning Approach*. Cambridge, MA, USA: MIT Press, 2017. ISBN 978-0-26-203641-2
- 5005 [318] S. Wang, D. Lo, and L. Jiang, "An empirical study on developer interactions in StackOverflow,"
5006 in *Proceedings of the 28th Annual ACM Symposium on Applied Computing*. Coimbra, Portugal:
5007 ACM, March 2013. DOI 10.1145/2480362.2480557, pp. 1019–1024.
- 5008 [319] W. Wang and M. W. Godfrey, "Detecting API usage obstacles: A study of iOS and android
5009 developer questions," in *Proceedings of the 10th Working Conference on Mining Software*
5010 *Repositories*. San Francisco, CA, USA: IEEE, May 2013. DOI 10.1109/MSR.2013.6624006.
5011 ISBN 978-1-46-732936-1. ISSN 2160-1852 pp. 61–64.
- 5012 [320] R. Watson, "Development and application of a heuristic to assess trends in API documentation,"
5013 in *Proceedings of the 30th ACM International Conference on Design of Communication*.
5014 Seattle, WA, USA: ACM, October 2012. DOI 10.1145/2379057.2379112. ISBN 978-1-45-
5015 031497-8 pp. 295–302.
- 5016 [321] R. Watson, M. Mark Starnes, J. Jeannot-Schroeder, and J. H. Spyridakis, "API documentation
5017 and software community values: A survey of open-source API documentation," in *Proceedings*
5018 *of the 31st ACM International Conference on Design of Communication*. Greenville, SC,
5019 USA: ACM, September 2013. DOI 10.1145/2507065.2507076, pp. 165–174.
- 5020 [322] S. Weerawarana, F. Curbera, F. Leymann, T. Storey, and D. F. Ferguson, *Web Services Platform*
5021 *Architecture*. Crawfordsville, IN, USA: Prentice-Hall, 2005. ISBN 0-13-148874-0
- 5022 [323] D. Wettschereck, D. W. Aha, and T. Mohri, "A Review and Empirical Evaluation of Feature
5023 Weighting Methods for a Class of Lazy Learning Algorithms," *Artificial Intelligence Review*,
5024 vol. 11, no. 1-5, pp. 273–314, 1997, DOI 10.1007/978-94-017-2053-3_11. ISSN 0269-2821
- 5025 [324] J. Wexler, M. Pushkarna, T. Bolukbasi, M. Wattenberg, F. Viegas, and J. Wilson, "The What-If
5026 Tool: Interactive Probing of Machine Learning Models," *IEEE Transactions on Visualization*
5027 *and Computer Graphics*, vol. 26, no. 1, pp. 56–65, 2019, DOI 10.1109/tvcg.2019.2934619.
5028 ISSN 1077-2626
- 5029 [325] H. Wickham, "A Layered grammar of graphics," *Journal of Computational and Graphical*
5030 *Statistics*, vol. 19, no. 1, pp. 3–28, January 2010, DOI 10.1198/jcgs.2009.07098. ISSN 1061-
5031 8600
- 5032 [326] R. Wieringa, N. Maiden, N. Mead, and C. Rolland, "Requirements engineering paper classification
5033 and evaluation criteria: a proposal and a discussion," *Requirements Engineering*, vol. 11,
5034 no. 1, pp. 102–107, March 2006, DOI 10.1007/s00766-005-0021-6. ISSN 0947-3602
- 5035 [327] I. H. Witten, E. Frank, M. A. Hall, and C. J. Pal, *Data Mining: Practical Machine Learning*
5036 *Tools and Techniques*. Morgan Kaufmann, 2016. DOI 10.1016/c2009-0-19715-5. ISBN
5037 978-0-12-804291-5
- 5038 [328] C. Wohlin and A. Aurum, "Towards a decision-making structure for selecting a research design
5039 in empirical software engineering," *Empirical Software Engineering*, vol. 20, no. 6, pp. 1427–
5040 1455, May 2015, DOI 10.1007/s10664-014-9319-7. ISSN 1573-7616
- 5041 [329] C. Wohlin, P. Runeson, M. Höst, M. C. Ohlsson, B. Regnell, and A. Wesslén, *Experimentation*
5042 *in Software Engineering*. Berlin, Heidelberg: Springer, 2012. DOI 10.1007/978-3-642-
5043 29044-2. ISBN 978-3-64-229044-2
- 5044 [330] M. L. Wong and K. S. Leung, *Data Mining Using Grammar Based Genetic Programming and*
5045 *Applications*. Springer, 2002. DOI 10.1007/b116131. ISBN 978-0-79-237746-7

- 5046 [331] M. R. Wrobel, "Emotions in the software development process," in *2013 6th International*
5047 *Conference on Human System Interactions (HSI)*. IEEE, 2013, pp. 518–523.
- 5048 [332] X. Yi and K. J. Kochut, "A CP-nets-based design and verification framework for web services
5049 composition," in *Proceedings of the 2004 IEEE International Conference on Web Services*. San
5050 Diego, CA, USA: IEEE, July 2004. DOI 10.1109/icws.2004.1314810. ISBN 0-76-952167-3
5051 pp. 756–760.
- 5052 [333] R. K. Yin, *Case study research and applications: Design and methods*, 6th ed. Los Angeles,
5053 CA, USA: SAGE, 2017. ISBN 978-1-50-633616-9
- 5054 [334] J. Zahálka and F. Železný, "An experimental test of Occam's razor in classification," *Machine*
5055 *Learning*, vol. 82, no. 3, pp. 475–481, 2011, DOI 10.1007/s10994-010-5227-2. ISSN 0885-6125
- 5056 [335] J. Zhang and R. Kasturi, "Extraction of Text Objects in Video Documents: Recent Progress," in
5057 *Proceedings of the 8th International Workshop on Document Analysis Systems*. Nara, Japan:
5058 IEEE, September 2008. DOI 10.1109/das.2008.49, pp. 5–17.
- 5059 [336] B. Zupan, J. Demšar, M. W. Kattan, J. R. Beck, and I. Bratko, "Machine learning for survival
5060 analysis: a case study on recurrence of prostate cancer," *Artificial intelligence in medicine*,
5061 vol. 20, no. 1, pp. 59–75, 2000.
- 5062 [337] M. Zur Muehlen, J. V. Nickerson, and K. D. Swenson, "Developing web services choreography
5063 standards - The case of REST vs. SOAP," *Decision Support Systems*, vol. 40, no. 1, pp. 9–29,
5064 July 2005, DOI 10.1016/j.dss.2004.04.008. ISSN 0167-9236

5065

5066

List of Online Artefacts

5067

5068 The online artefacts listed below have been downloaded and stored on the Deakin
5069 Research Data Store (RDS) for archival purposes at the following location:

5070 RDS29448-Alex-Cummaudo-PhD/datasets/webrefs

- 5071 [338] Affectiva, Inc., “Home - Affectiva : Affectiva,” <http://bit.ly/36sgbMM>, 2018, accessed: 15
5072 October 2018.
- 5073 [339] Amazon Web Services, Inc., “Detecting Labels in an Image,” <https://amzn.to/2TBNtTa>, 2018,
5074 accessed: 28 August 2018.
- 5075 [340] ———, “Detecting Objects and Scenes,” <https://amzn.to/2TDed5V>, 2018, accessed: 28 August
5076 2018.
- 5077 [341] ———, “Amazon Rekognition,” <https://amzn.to/2TyT2BL>, 2018, accessed: 13 September 2018.
- 5078 [342] ———, “Aws release notes,” <https://go.aws/2v0RYjr>, 2019, accessed: 18 March 2019.
- 5079 [343] ———, “Actions - amazon rekognition,” <https://amzn.to/392p3dH>, 2019, accessed: 18 March
5080 2019.
- 5081 [344] ———, “Amazon rekognition | aws machine learning blog,” <https://go.aws/37Q7lKc>, 2019, ac-
5082 cessed: 18 March 2019.
- 5083 [345] ———, “Amazon rekognition image,” <https://go.aws/2ubB6qc>, 2019, accessed: 18 March 2019.
- 5084 [346] ———, “Best practices for sensors, input images, and videos - amazon rekognition,” <https://amzn.to/2uZlWo0>, 2019, accessed: 18 March 2019.
- 5085 [347] ———, “Exercise 1: Detect objects and scenes in an image (console) - amazon rekognition,” <https://amzn.to/36TlNm>, 2019, accessed: 18 March 2019.
- 5086 [348] ———, “Java (sdk v1) code samples for amazon rekognition - aws code sample,” <https://amzn.to/2ugTle3>, 2019, accessed: 18 March 2019.
- 5087 [349] ———, “Limits in amazon rekognition - amazon rekognition,” <https://amzn.to/2On6n0h>, 2019,
5088 accessed: 18 March 2019.
- 5089 [350] ———, “Step 1: Set up an aws account and create an iam user - amazon rekognition,” <https://amzn.to/2tqW4kI>, 2019, accessed: 18 March 2019.
- 5090 [351] ———, “Troubleshooting amazon rekognition video - amazon rekognition,” <https://amzn.to/3b763fS>, 2019, accessed: 18 March 2019.
- 5091 [352] Beijing Geling Shentong Information Technology Co., Ltd., “DeepGlint,” <http://bit.ly/2uHHdPS>, 2018, accessed: 3 April 2019.
- 5092 [353] Beijing Kuangshi Technology Co., Ltd., “Megvii,” <http://bit.ly/2WJYFzk>, 2018, accessed: 3
5093 April 2019.

- 5100 [354] Clarifai, Inc., “Enterprise AI Powered Computer Vision Solutions | Clarifai,” <http://bit.ly/2TB3kSa>, 2018, accessed: 13 September 2018.
- 5101 [355] CloudSight, Inc., “Image Recognition API & Visual Search Results | CloudSight AI,” <http://bit.ly/2UmNPCw>, 2018, accessed: 13 September 2018.
- 5102 [356] Cognitec Systems GmbH, “The face recognition company - Cognitec,” <http://bit.ly/38VguBB>, 2018, accessed: 15 October 2018.
- 5103 [357] A. Cummaudo, “ICSE 2020 Submission #564 Supplementary Materials,” <http://bit.ly/2Z8zOKW>, 2019.
- 5104 [358] Deep AI, Inc., “DeepAI: The front page of A.I. | DeepAI,” <http://bit.ly/2TBNYgf>, 2018, accessed: 26 September 2018.
- 5105 [359] Google LLC, “Best practices for enterprise organizations | documentation | google cloud,” <http://bit.ly/2v0RSs5>, 2019, accessed: 18 March 2019.
- 5106 [360] ——, “Detect Labels | Google Cloud Vision API Documentation | Google Cloud,” <http://bit.ly/2TD5key>, 2018, accessed: 28 August 2018.
- 5107 [361] ——, “Class EntityAnnotation | Google.Cloud.Vision.V1,” <http://bit.ly/2TD5fpg>, 2018, accessed: 28 August 2018.
- 5108 [362] ——, “Vision API - Image Content Analysis | Cloud Vision API | Google Cloud,” <http://bit.ly/2TD9mBs>, 2018, accessed: 13 September 2018.
- 5109 [363] ——, “Machine learning glossary | google developers,” <http://bit.ly/3b38VdL>, 2019, accessed: 18 March 2019.
- 5110 [364] ——, “Open Images Dataset V4,” <http://bit.ly/2Ry2vvF>, 2019, accessed: 9 November 2018.
- 5111 [365] ——, “Quickstart: Using client libraries | cloud vision api documentation | google cloud,” <http://bit.ly/2RRMQHG>, 2019, accessed: 18 March 2019.
- 5112 [366] ——, “Release notes | cloud vision api documentation | google cloud,” <http://bit.ly/2UipY5J>, 2019, accessed: 18 March 2019.
- 5113 [367] ——, “Sample applications | cloud vision api documentation | google cloud,” <http://bit.ly/2SdoB5r>, 2019, accessed: 18 March 2019.
- 5114 [368] ——, “Tips & tricks | cloud functions documentation | google cloud,” <http://bit.ly/2GZNc8Z>, 2019, accessed: 18 March 2019.
- 5115 [369] ——, “Vision ai | derive image insights via ml | cloud vision api | google cloud,” <http://bit.ly/31nWoNx>, 2019, accessed: 18 March 2019.
- 5116 [370] Guangzhou Tup Network Technology, “TupuTech,” <http://bit.ly/2uF4IsN>, 2018, accessed: 3 April 2019.
- 5117 [371] Imappa Technologies, “Imappa - powerful image recognition APIs for automated categorization & tagging in the cloud and on-premises,” <http://bit.ly/2TxsyRe>, 2018, accessed: 13 September 2018.
- 5118 [372] International Business Machines Corporation, “Watson Visual Recognition - Overview | IBM,” <https://ibm.co/2TBNIO4>, 2018, accessed: 13 September 2018.
- 5119 [373] ——, “Watson Tone Analyzer,” <https://ibm.co/37w3y4A>, 2019, accessed: 25 January 2019.
- 5120 [374] Kairos AR, Inc., “Kairos: Serving Businesses with Face Recognition,” <http://bit.ly/30WHGNs>, 2018, accessed: 15 October 2018.
- 5121 [375] Microsoft Corporation, “azure-sdk-for-java/ImageTag.java,” <http://bit.ly/38IDPWU>, 2018, accessed: 28 August 2018.
- 5122 [376] ——, “Image Processing with the Computer Vision API | Microsoft Azure,” <http://bit.ly/2YqhkS6>, 2018, accessed: 13 September 2018.
- 5123 [377] ——, “How to call the Computer Vision API,” <http://bit.ly/2TD5oJk>, 2018, accessed: 28 August 2018.
- 5124 [378] ——, “Call the computer vision api - azure cognitive services | microsoft docs,” <http://bit.ly/2vHSdjT>, 2019, accessed: 18 March 2019.
- 5125 [379] ——, “Content tags - computer vision - azure cognitive services | microsoft docs,” <http://bit.ly/2vESzHX>, 2019, accessed: 18 March 2019.
- 5126 [380] ——, “Github - azure-samples/cognitive-services-java-computer-vision-tutorial: This tutorial shows the features of the microsoft cognitive services computer vision rest api.” <http://bit.ly/37N1yoN>, 2019, accessed: 18 March 2019.
- 5127 [381] ——, “Improving your classifier - custom vision service - azure cognitive services | microsoft docs,” <http://bit.ly/37SBkRQ>, 2019, accessed: 18 March 2019.

- 5156 [382] ——, “Quickstart: Computer vision client library for .net - azure cognitive services | microsoft
5157 docs,” <http://bit.ly/2vF3wJC>, 2019, accessed: 18 March 2019.
- 5158 [383] ——, “Release notes - custom vision service - azure cognitive services | microsoft docs,”
5159 <http://bit.ly/2UIPiaW>, 2019, accessed: 18 March 2019.
- 5160 [384] ——, “Sample: Explore an image processing app in c# - azure cognitive services | microsoft
5161 docs,” <http://bit.ly/2u4mPMh>, 2019, accessed: 18 March 2019.
- 5162 [385] ——, “Tutorial: Generate metadata for azure images - azure cognitive services | microsoft
5163 docs,” <http://bit.ly/2RRnARK>, 2019, accessed: 18 March 2019.
- 5164 [386] ——, “Tutorial: Use custom logo detector to recognize azure services - custom vision - azure
5165 cognitive services | microsoft docs,” <http://bit.ly/2RUGwPH>, 2019, accessed: 18 March 2019.
- 5166 [387] ——, “What is computer vision? - computer vision - azure cognitive services | microsoft docs,”
5167 <http://bit.ly/37SomDx>, 2019, accessed: 18 March 2019.
- 5168 [388] SenseTime, “SenseTime,” <http://bit.ly/2WH6RjF>, 2018, accessed: 3 April 2019.
- 5169 [389] Shanghai Yitu Technology Co., Ltd., “Yitu Technology,” <http://bit.ly/2uGvxgf>, 2018, accessed:
5170 3 April 2019.
- 5171 [390] Stack Overflow User #1008563 ‘samiles’, “AWS Rekognition PHP SDK gives invalid image
5172 encoding error,” <http://bit.ly/31Sgpec>, 2019, accessed: 22 June 2019.
- 5173 [391] Stack Overflow User #10318601 ‘reza naderii’, “google cloud vision category detecting,”
5174 <http://bit.ly/31Uf32t>, 2019, accessed: 22 June 2019.
- 5175 [392] Stack Overflow User #10729564 ‘gabgob’, “Multiple Google Vision OCR requests at once?”
5176 <http://bit.ly/31P09dU>, 2019, accessed: 22 June 2019.
- 5177 [393] Stack Overflow User #1453704 ‘deoptimancode’, “Human body part detection in Android,”
5178 <http://bit.ly/31T5pxd>, 2019, accessed: 22 June 2019.
- 5179 [394] Stack Overflow User #174602 ‘geekyaleks’, “aws Rekognition not initializing on iOS,” <http://bit.ly/31UeqG9>, 2019, accessed: 22 June 2019.
- 5180 [395] Stack Overflow User #2251258 ‘James Dorfman’, “All GoogleVision label possibilities?”
5181 <http://bit.ly/31R4FZi>, 2019, accessed: 22 June 2019.
- 5182 [396] Stack Overflow User #2521469 ‘Hillary Sanders’, “Is there a full list of potential labels that
5183 Google’s Vision API will return?” <http://bit.ly/2KNnJSB>, 2019, accessed: 22 June 2019.
- 5184 [397] Stack Overflow User #2604150 ‘user2604150’, “Google Vision Accent Character Set NodeJs,”
5185 <http://bit.ly/31TsVdp>, 2019, accessed: 22 June 2019.
- 5186 [398] Stack Overflow User #3092947 ‘Mark Bench’, “Google Cloud Vision OCR API returning
5187 incorrect values for bounding box/vertices,” <http://bit.ly/31UeZjf>, 2019, accessed: 22 June
5188 2019.
- 5189 [399] Stack Overflow User #3565255 ‘CSquare’, “Vision API topicality and score always the same,”
5190 <http://bit.ly/2TD5As2>, 2019, accessed: 22 June 2019.
- 5191 [400] Stack Overflow User #4748115 ‘Latifa Al-jiffry’, “similar face recognition using google cloud
5192 vision API in android studio,” <http://bit.ly/31WhMZY>, 2019, accessed: 22 June 2019.
- 5193 [401] Stack Overflow User #4852910 ‘Gaurav Mathur’, “Amazon Rekognition Image caption,” <http://bit.ly/31P08qm>, 2019, accessed: 22 June 2019.
- 5194 [402] Stack Overflow User #5294761 ‘Eury Pérez Beltré’, “Specify language for response in Google
5195 Cloud Vision API,” <http://bit.ly/31SsUGG>, 2019, accessed: 22 June 2019.
- 5196 [403] Stack Overflow User #549312 ‘GroovyDotCom’, “Image Selection for Training Visual Recog-
5197 nition,” <http://bit.ly/31W8lcw>, 2019, accessed: 22 June 2019.
- 5198 [404] Stack Overflow User #5809351 ‘J.Doe’, “How to confidently validate object detection results
5199 returned from Google Cloud Vision.” <http://bit.ly/31UcCNy>, 2019, accessed: 22 June 2019.
- 5200 [405] Stack Overflow User #5844927 ‘Amit Pawar’, “Google cloud Vision and Clarifai doesn’t Support
5201 tagging for 360 degree images and videos,” <http://bit.ly/31StuEm>, 2019, accessed: 22 June 2019.
- 5202 [406] Stack Overflow User #5924523 ‘Akash Dathan’, “Can i give aspect ratio in Google Vision api?”
5203 <http://bit.ly/2KSJwsp>, 2019, accessed: 22 June 2019.
- 5204 [407] Stack Overflow User #6210900 ‘Mike Grommet’, “Are the Cloud Vision API limits in docu-
5205 mentation correct?” <http://bit.ly/31SsNLg>, 2019, accessed: 22 June 2019.
- 5206 [408] Stack Overflow User #6649145 ‘I. Sokolyk’, “How to get a position of custom object on image
5207 using vision recognition api,” <http://bit.ly/3210Q49>, 2019, accessed: 22 June 2019.
- 5208 [409] Stack Overflow User #6841211 ‘NigelJL’, “Google Cloud Vision - Numbers and Numerals
5209 OCR,” <http://bit.ly/31P07mi>, 2019, accessed: 22 June 2019.

- 5212 [410] Stack Overflow User #7064840 ‘Josh’, “Google Cloud Vision fails at batch annotate images.
5213 Getting Netty Shaded ClosedChannelException error,” <http://bit.ly/31UrBH9>, 2019, accessed:
5214 22 June 2019.
- 5215 [411] Stack Overflow User #7187987 ‘tuanars10’, “Adding a local path to Microsoft Face API by
5216 Python,” <http://bit.ly/2KLeMt3>, 2019, accessed: 22 June 2019.
- 5217 [412] Stack Overflow User #7219743 ‘Davide Biraghi’, “Google Vision API does not recognize
5218 single digits,” <http://bit.ly/31Ws1Nj>, 2019, accessed: 22 June 2019.
- 5219 [413] Stack Overflow User #738248 ‘lavuy’, “Meaning of score in Microsoft Cognitive Service’s
5220 Entity Linking API,” <http://bit.ly/2TD9vWw>, 2019, accessed: 22 June 2019.
- 5221 [414] Stack Overflow User #7604576 ‘Alagappan Narayanan’, “Text extraction - line-by-line,” <http://bit.ly/31Yc21s>, 2019, accessed: 22 June 2019.
- 5222 [415] Stack Overflow User #7692297 ‘lucas’, “Can Google Cloud Vision generate labels in Spanish
5223 via its API?” <http://bit.ly/31UcBsY>, 2019, accessed: 22 June 2019.
- 5224 [416] Stack Overflow User #7896427 ‘David mark’, “Google Api Vision, ““before_request”” error,”
5225 <http://bit.ly/31Z27Zt>, 2019, accessed: 22 June 2019.
- 5226 [417] Stack Overflow User #8210103 ‘Cosmin-Ioan Leferman’, “Google Vision API text detection
5227 strange behaviour - Javascript,” <http://bit.ly/31Ucyxi>, 2019, accessed: 22 June 2019.
- 5228 [418] Stack Overflow User #8411506 ‘AsSportac’, “How can we find an exhaustive list (or graph)
5229 of all logos which are effectively recognized using Google Vision logo detection feature?”
5230 <http://bit.ly/31Z27IX>, 2019, accessed: 22 June 2019.
- 5231 [419] Stack Overflow User #8594124 ‘God Himself’, “How to set up AWS mobile SDK in iOS project
5232 in Xcode,” <http://bit.ly/31St2pE>, 2019, accessed: 22 June 2019.
- 5233 [420] Stack Overflow User #9006896 ‘Dexter Intelligence’, “Getting wrong text sequence when image
5234 scanned by offline google mobile vision API,” <http://bit.ly/31Sgr5O>, 2019, accessed: 22 June
5235 2019.
- 5236 [421] Stack Overflow User #9913535 ‘Sahil Mehra’, “Google Vision API: ModuleNotFoundError:
5237 module not found ‘google.oauth2’,” <http://bit.ly/31VIZfU>, 2019, accessed: 22 June 2019.
- 5238 [422] Symisc Systems, S.U.A.R.L, “Computer Vision & Media Processing APIs | PixLab,” <http://bit.ly/2UlkW9K>, 2018, accessed: 13 September 2018.
- 5239 [423] Talkwalker Inc., “Image Recognition - Talkwalker,” <http://bit.ly/2TyT7W5>, 2018, accessed: 13
5240 September 2018.
- 5241 [424] TheySay Limited, “Sentiment Analysis API | TheySay,” <http://bit.ly/37AzTHI>, 2019, accessed:
5242 25 January 2019.
- 5243
- 5244

Part IV

Appendices

APPENDIX A

Additional Materials

A.1 On the Development, Documentation and Usage of Web APIs

The development of web application programming interfaces (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 [266]. This framework used the synchronous communication paradigm remote procedure calls (RPCs) first introduced by Nelson [223] that allows procedures to be called in a remote address space as if it were local. Its communication paradigm, DCE/RPC [230], 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 [230]). 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) [61] 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 (SLA) whereby the consumer and provider agree upon a set level of service, which in some cases are legally binding [26]. 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 [317, 322]). These attributes are included within SOAP headers; thus, QoS aspects are independent from the transport layer and instead exist at the application layer [239]. 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.

While it is rich in metadata and verbosity, discussions on whether this was a benefit or drawback came about the mid-2000s [239, 337] 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).

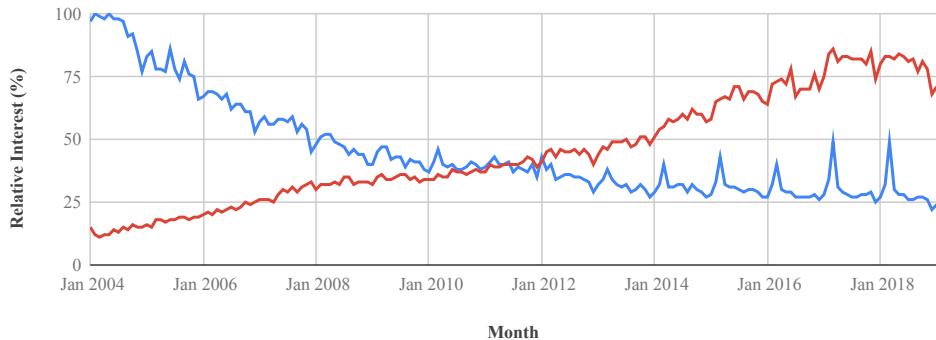


Figure A.1: Worldwide search interest for SOAP (blue) and REST (red) since 2004. Source: Google Trends.

Listing A.1: A SOAP HTTP POST consumer request to retrieve customer record #43456 from a web service provider. Source: [18].

```

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: [18].

```

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 [239]. That is, the HTTP is considered as a transport protocol solely. In 2000, Fielding [102] 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 of 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 [30].
- **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) [198]). Therefore, no RESTful service uses a standardised response document or invocation syntax. While there are proposals, such as WADL [126], 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: [18].

```
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.

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

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.

A.2 Additional Figures

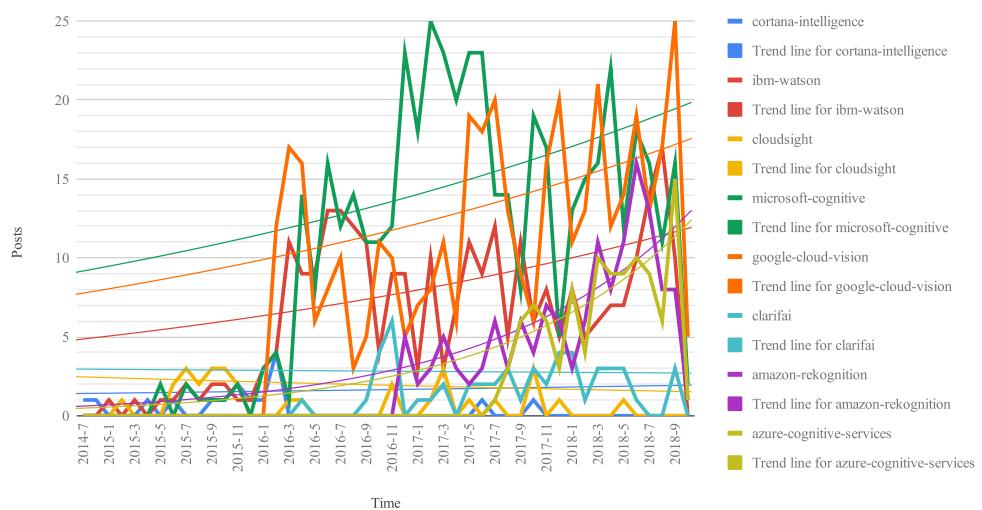
The following figures are listed in this section:

- List
- of
- every
- additional
- figure (and page ref)
- and
- why
- it
- is
- relevant.

*< todo: Include causal model: draft1 & 4 >
< todo: Include technical domain model >
< todo: Include threshy: decision boundary > < todo: Include threshy: domain
model > < todo: Include threshy: sequence diagrams >
< todo: Include ICSE workflow >
< todo: Include architectural model: new brc > < todo: Include architectural
model: evolution > < todo: Include architectural model: creation of request >
< todo: Include architectural model: evolution > < todo: Include architectural
model: class diagram > < todo: Include architectural model: evolution > < todo:
Include architectural model: overall state diagram > < todo: Include architectural
model: page6 > < todo: Include architectural model: page7 >*

Figure A.2: A Broad Range of artificial intelligence (AI)-Based Products And Services Is Already Visible. (From [193].)

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 computer vision services (CVSs).

A.3 Reference Architecture Source Code

(TODO: ICVS benchmarker code)

APPENDIX B

Supplementary Materials to Chapter 7

B.1 Detailed Overview of Our Proposed Taxonomy

An overview of the 5 dimensions and categories (sub-dimensions) within our proposed taxonomy. ILS = In-Literature Score, calculated as a percentage of the number of papers that make the recommendation of all 21 primary sources. IPS = In-Practice Score, calculated as the average compliance to the SUS. Colour scales indicate relevancy weight within ILS or IPS values for comparative purposes, where red = *lowest* and green = *highest*. GCV, AWS, ACV = Presence of category in Google Cloud Vision, Amazon Rekognition, and Azure Cloud Vision documentation. Presence indicated as *fully present* (●), *partially present* (◐), and *not present* (○).

Key	Description	Primary Sources	ILS	IPS	GCV	AWS	ACV
A1	Quick-start guides to rapidly get started using the API in a specific programming language.	S4, S9, S10	0.14	0.88	●	○	●
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	0.52	0.56	●	●	●
A3	Explanations of the API's high-level architecture to better understand intent and context.	S1, S2, S4, S11, S14, S16, S19, S20	0.38	0.70	●	●	●
A4	Source code implementation and code comments (where applicable) to understand the API author's mindset.	S1, S4, S7, S12, S13, S17, S20	0.33	0.47	○	○	○
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	0.71	0.89	●	●	●
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	0.57	0.54	○	●	●
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	0.24	0.66	○	○	●
A8	Best-practices of implementation to assist with debugging and efficient use of the API.	S1, S2, S4, S5, S7, S8, S9, S14	0.38	0.68	○	●	○
A9	An exhaustive list of all major components that exist within the API.	S4, S16, S19	0.14	0.69	○	●	●
A10	Minimum system requirements and dependencies to use the API.	S4, S7, S13, S17, S19	0.24	0.71	○	○	●
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	0.38	0.86	○	●	○
A12	Error definitions that describe how to address a specific problem.	S1, S2, S4, S5, S9, S11, S13	0.33	0.84	○	○	○

Continued on next page...

Key	Description	Primary Sources	ILS	IPS	GCV	AWS	ACV
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	0.48	0.80	●	●	●
B2	Descriptions of the types of applications the API can develop.	S2, S4, S9, S11, S15, S18	0.29	0.57	○	○	●
B3	Descriptions of the types of users who should use the API.	S4, S9	0.10	0.44	○	○	○
B4	Descriptions of the types of users who will use the product the API creates.	S4	0.05	0.42	○	○	○
B5	Success stories about the API used in production.	S4	0.05	0.49	○	●	●
B6	Documentation to compare similar APIs within the context to this API.	S2, S6, S13, S18	0.19	0.47	○	○	●
B7	Limitations on what the API can and cannot provide.	S4, S5, S8, S9, S14, S16	0.29	0.94	○	●	●
C1	Descriptions of the relationship between API components and domain concepts.	S3, S10	0.10	0.51	○	○	●
C2	Definitions of domain-terminology and concepts, with synonyms if applicable.	S2, S3, S4, S6, S7, S10, S14, S16	0.38	0.74	○	○	○
C3	Generalised documentation for non-technical audiences regarding the API and its domain.	S4, S8, S16	0.14	0.55	●	●	●
D1	A list of FAQs.	S4, S7	0.10	0.75	●	●	●
D2	Troubleshooting suggestions.	S4, S8	0.10	0.56	○	○	○
D3	Diagrammatically representing API components using visual architectural representations.	S6, S13, S20	0.14	0.63	○	○	○
D4	Contact information for technical support.	S4, S8, S19	0.14	0.21	●	●	●
D5	A printed/printable resource for assistance.	S4, S6, S7, S9, S16	0.24	0.56	○	●	●

Continued on next page...

Key	Description	Primary Sources	ILS	IPS	GCV	AWS	ACV
D6	Licensing information.	S7	0.05	0.66	○	○	◐
E1	Searchable knowledge base.	S3, S4, S6, S10, S14, S17, S18	0.33	0.81	●	●	●
E2	Context-specific discussion forum.	S4, S10, S11	0.14	0.58	●	●	◐
E3	Quick-links to other relevant documentation frequently viewed by developers.	S6, S16, S20	0.14	0.63	○	○	○
E4	Structured navigational style (e.g., breadcrumbs).	S6, S10, S20	0.14	0.58	●	●	●
E5	Visualised map of navigational paths to certain API components in the website.	S6, S14, S20	0.14	0.50	○	○	○
E6	Consistent look and feel of documentation.	S1, S2, S3, S5, S6, S8, S10, S15, S20	0.43	0.70	●	●	●

B.2 Sources of Documentation

Sources of documentation used for the validation of the taxonomy. For clarity, exact webpages are not referenced for each category, but can be found in supplementary materials which can be downloaded from the URL listed in the paper.

Service	Document Sources
Google Cloud Vision	https://cloud.google.com/vision/docs/quickstart-client-libraries https://googleapis.github.io/google-cloud-java/google-cloud-clients/apidocs/index.html https://cloud.google.com/vision/#cloud-vision-use-cases https://cloud.google.com/vision/docs/quickstart-client-libraries#using_the_client_library https://cloud.google.com/vision/docs/tutorials https://cloud.google.com/community/tutorials?q=vision https://cloud.google.com/vision/docs/samples#mobile_platform_examples https://cloud.google.com/docs/enterprise/best-practices-for-enterprise-organizations https://cloud.google.com/functions/docs/bestpractices/tips https://cloud.google.com/vision/#derive-insight-from-images-with-our-powerful-cloud-vision-api https://cloud.google.com/vision/docs/quickstart-client-libraries https://cloud.google.com/vision/docs/release-notes https://cloud.google.com/vision/docs/reference/rpc/google.rpc#google.rpc.Code https://cloud.google.com/vision/#derive-insight-from-your-images-with-our-powerful-----pretrained-api-models-or-easily-train-custom-vision-models-with-automl-----vision-beta https://cloud.google.com/vision/#insight-from-your-images https://developers.google.com/machine-learning/glossary/ https://cloud.google.com/vision/docs/resources https://cloud.google.com/vision/sla https://cloud.google.com/vision/docs/data-usage https://cloud.google.com/vision/docs/support#searchbox https://cloud.google.com/vision/docs/support

Continued on next page...

Service	Document Sources
Amazon Rekognition	<p>https://docs.aws.amazon.com/rekognition/latest/dg/getting-started.html</p> <p>https://docs.aws.amazon.com/AWSJavaSDK/latest/javadoc/index.html</p> <p>https://aws.amazon.com/rekognition/#Rekognition_Image_Use_Cases</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/labels-detect-labels-image.html</p> <p>https://aws.amazon.com/rekognition/getting-started/#Tutorials</p> <p>https://aws.amazon.com/blogs/machine-learning/category/artificial-intelligence/amazon-rekognition/</p> <p>https://docs.aws.amazon.com/code-samples/latest/catalog/code-catalog-java-example_code-rekognition.html</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/best-practices.html</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/API_Operations.html</p> <p>https://aws.amazon.com/rekognition/image-features/</p> <p>https://aws.amazon.com/releasenotes/?tag=releasenotes%23keywords%23amazon-rekognition</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/setting-up.html</p> <p>https://aws.amazon.com/rekognition/</p> <p>https://aws.amazon.com/rekognition/</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/limits.html</p> <p>https://aws.amazon.com/rekognition/pricing/</p> <p>https://aws.amazon.com/rekognition/sla/</p> <p>https://aws.amazon.com/rekognition/faqs/</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/video-troubleshooting.html</p> <p>https://docs.aws.amazon.com/rekognition/latest/dg/rekognition-dg.pdf</p> <p>https://github.com/awsdocs/amazon-rekognition-developer-guide/issues</p> <p>https://forums.aws.amazon.com/thread.jspa?threadID=285910</p>

Continued on next page...

Service	Document Sources
Azure Computer Vision	https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/quickstarts-sdk/csharp-analyze-sdk https://docs.microsoft.com/en-us/java/api/overview/azure/cognitiveservices/client/computervision?view=azure-java-stable https://docs.microsoft.com/en-us/azure/architecture/example-scenario/ai/intelligent-apps-image-processing https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/tutorials/java-tutorial https://docs.microsoft.com/en-us/azure/cognitive-services/custom-vision-service/logo-detector-mobile https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/tutorials/storage-lab-tutorial https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/tutorials/csharpTutorial https://docs.microsoft.com/en-us/azure/cognitive-services/custom-vision-service/getting-started-improving-your-classifier https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/home#analyze-images-for-insight https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/vision-api-how-to-topics/howtocallvisionapi https://docs.microsoft.com/en-us/azure/cognitive-services/custom-vision-service/release-notes https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/ https://azure.microsoft.com/en-au/services/cognitive-services/computer-vision/ https://azure.microsoft.com/en-us/pricing/details/cognitive-services/computer-vision/ https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/concept-tagging-images https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/home https://azure.microsoft.com/en-us/support/legal/sla/cognitive-services/v1_1/ https://docs.microsoft.com/en-au/azure/cognitive-services/computer-vision/faq https://azure.microsoft.com/en-us/support/legal/

B.3 List of Primary Sources

Below lists the primary sources identified in our systematic mapping study. They are listed in order of assignment to the taxonomy described in Appendix B.1.

- S1. M. P. Robillard, "What makes APIs hard to learn? Answers from developers," *IEEE Software*, vol. 26, no. 6, pp. 27–34, 2009, DOI 10.1109/MS.2009.193. ISSN 0740-7459
- S2. M. P. Robillard and R. Deline, "A field study of API learning obstacles," *Empirical Software Engineering*, vol. 16, no. 6, pp. 703–732, 2011, DOI 10.1007/s10664-010-9150-8. ISSN 1382-3256
- S3. A. J. Ko and Y. Riche, "The role of conceptual knowledge in API usability," in *Proceedings of the 2011 IEEE Symposium on Visual Languages and Human Centric Computing*. Pittsburg, PA, USA: IEEE, September 2011. DOI 10.1109/VLHCC.2011.6070395. ISBN 978-1-45771245-6 pp. 173–176
- S4. 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," in *Proceedings of the 20th Annual International Conference on Computer Documentation*. Toronto, ON, Canada: ACM, October 2002. DOI 10.1145/584955.584976, pp. 133–141
- S5. R. Watson, M. Mark Stammes, J. Jeannot-Schroeder, and J. H. Spyridakis, "API documentation and software community values: A survey of open-source API documentation," in *Proceedings of the 31st ACM International Conference on Design of Communication*. Greenville, SC, USA: ACM, September 2013. DOI 10.1145/2507065.2507076, pp. 165–174
- S6. 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 *Proceedings of the First International Symposium on End User Development*, vol. 5435 LNCS. Siegen, Germany: Springer, March 2009. DOI 10.1007/978-3-642-00427-8_6. ISSN 0302-9743 pp. 86–105
- S7. E. Aghajani, C. Nagy, O. L. Vega-Marquez, M. Linares-Vasquez, L. Moreno, G. Bavota, and M. Lanza, "Software Documentation Issues Unveiled," in *Proceedings of the 41st International Conference on Software Engineering*. Montreal, QC, Canada: IEEE, May 2019. DOI 10.1109/ICSE.2019.00122. ISBN 978-1-72-810869-8. ISSN 0270-5257 pp. 1199–1210
- S8. S. Haselbock, R. Weinreich, G. Buchgeher, and T. Krichbaum, "Microservice Design Space Analysis and Decision Documentation: A Case Study on API Management," in *Proceedings of the 11th International Conference on Service-Oriented Computing and Applications, SOCA 2018*, Paris, France, November 2019, DOI 10.1109/SOCA.2018.00008, pp. 1–8
- S9. S. Inzunza, R. Juárez-Ramírez, and S. Jiménez, "API Documentation," in *Proceedings of the 6th World Conference on Information Systems and Technologies*. Naples, Italy: Springer, March 2018. DOI 10.1007/978-3-319-77712-2_2, pp. 229–239
- S10. 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, August 2018, DOI 10.1177/0047281617721853. ISSN 1541-3780
- S11. 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: A Collaborative Ethnography of Documentation Work," *Computer Supported Cooperative Work: CSCW: An International Journal*, vol. 27, no. 3-6, pp. 767–802, May 2018, DOI 10.1007/s10606-018-9333-1. ISSN 15737551
- S12. A. Head, C. Sadowski, E. Murphy-Hill, and A. Knight, "When not to comment: Questions and tradeoffs with API documentation for C++ projects," in *Proceedings of the 40th International Conference on Software Engineering*, ser. questions and tradeoffs with API documentation for C++ projects. Gothenburg, Sweden: ACM, May 2018. DOI 10.1145/3180155.3180176. ISSN 0270-5257 pp. 643–653

- S13. L. Aversano, D. Guardabascio, and M. Tortorella, “Analysis of the Documentation of ERP Software Projects,” *Procedia Computer Science*, vol. 121, pp. 423–430, January 2017, DOI 10.1016/j.procs.2017.11.057. ISSN 1877-0509
- S14. M. P. Robillard, A. Marcus, C. Treude, G. Bavota, O. Chaparro, N. Ernst, M. A. Gerosall, M. Godfrey, M. Lanza, M. Linares-Vásquez, G. C. Murphy, L. Moreno, D. Shepherd, and E. Wong, “On-demand developer documentation,” in *Proceedings of the 33rd IEEE International Conference on Software Maintenance and Evolution*. Shanghai, China: IEEE, September 2017. DOI 10.1109/ICSME.2017.17, pp. 479–483
- S15. R. Watson, “Development and application of a heuristic to assess trends in API documentation,” in *Proceedings of the 30th ACM International Conference on Design of Communication*. Seattle, WA, USA: ACM, October 2012. DOI 10.1145/2379057.2379112. ISBN 978-1-45031497-8 pp. 295–302
- S16. W. Maalej and M. P. Robillard, “Patterns of knowledge in API reference documentation,” *IEEE Transactions on Software Engineering*, 2013, DOI 10.1109/TSE.2013.12. ISSN 0098-5589
- S17. D. L. Parnas and S. A. Vilkomir, “Precise documentation of critical software,” in *Proceedings of 10th IEEE International Symposium on High Assurance Systems Engineering*. Plano, TX, USA: IEEE, November 2007. DOI 10.1109/HASE.2007.63. ISSN 1530-2059 pp. 237–244
- S18. C. Bottomley, “What part writer? What part programmer? A survey of practices and knowledge used in programmer writing,” in *Proceedings of the 2005 IEEE International Professional Communication Conference*. Limerick, Ireland: IEEE, July 2005. DOI 10.1109/IPCC.2005.1494255, pp. 802–812
- S19. 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, June 2004, DOI 10.1016/j.infsof.2003.10.004. ISSN 0950-5849
- S20. J. Kotula, “Using patterns to create component documentation,” *IEEE Software*, vol. 15, no. 2, pp. 84–92, 1998, DOI 10.1109/52.663791. ISSN 0740-7459
- S21. S. G. McLellan, A. W. Roesler, J. T. Tempest, and C. I. Spinuzzi, “Building more usable APIs,” *IEEE Software*, vol. 15, no. 3, pp. 78–86, 1998, DOI 10.1109/52.676963. ISSN 0740-7459

B.4 Survey Questions

This section contains the exact text of the survey described in Section 7.5.1. Our instrument also included questions where answers were not included in the research reported in this article, e.g. questions 1 and 2 regarding consent and ensuring participants have had development experience. Images used within the survey have been removed.

Developer opinions towards the importance of web API documentation recommendations

In this study, we are finding out how important recommendations of web API documentation are to developers. From this, we will improve AI-powered APIs. While there are screenshots of example APIs in the questions, think of an API that you have used based on **your own prior experience** when answering these questions. Thanks for taking the time to answer these questions; it should only take you about **10–20 minutes** to complete.

Attribution Notice

Portions of this questionnaire are reproduced from work created and shared by Google and used according to terms described in the Creative Commons 3.0 Attribution License.

Implementation-specific documentation of web APIs

When answering these questions please answer with respect to **your own experience** in learning web APIs (if applicable). Any examples provided exist solely to help illustrate the statement. For each question, please nominate how much you agree with the following statements: [*Strongly agree, Somewhat agree, Neither agree nor disagree, Somewhat disagree, Strongly disagree*]

- Q3a. I think quick-start guides with code that help me get started with an API's client library are important. e.g., quick-start guides that show how to get started and interact with the API and its responses.
- Q3b. I don't find low-level documentation of all classes and methods particularly helpful. e.g., a generated online reference manual from Javadoc comments.
- Q3c. I would imagine that explanations of the API's high-level architecture, context and rationale would be important to better understand how to consume the API. e.g., a graphic showing how the API could fit into the wider context of an application.
- Q3d. If I want to understand why an API did something that I didn't expect, the source code comments generally don't help me. e.g., an example from the Lodash API that describes why `set.add` isn't directly returned.
- Q3e. I find small code snippets with comments to demonstrate a single component's basic functionality within the API a useful way to learn. e.g., 10-30 lines of code to demonstrating various how-tos of a computer vision API.
- Q3f. I think it's cumbersome to read through step-by-step tutorials that show how to build something non-trivial with multiple components using the API. e.g., a ten-step tutorial documenting how to combine face recognition, face analysis, scene description, and landmark detection API components to generate descriptions of photos.

- Q3g. I think it's useful to download source code of production-ready applications that demonstrate the use of multiple facets of the API. e.g., a downloadable iOS app that demonstrates how to perform image analysis on an iPhone/iPad.
- Q3h. I think official documentation describing the 'best-practices' of how to use the API to assist with debugging and efficiency is not helpful. e.g., an article describing the correct ways of doing things, the best tools to use, and how to write well-performing code.
- Q3i. I believe an exhaustive list of all major components in the API without excessive detail would be useful when learning an API. e.g., a computer vision web API might list object detection, object localisation, facial recognition, and facial comparison as its 4 components.
- Q3j. I believe minimum system requirements and/or dependencies to use the API do not always need to be part of official documentation. e.g., I can find descriptions of how to get started with a Python environment for a cloud platform on community forums instead of the API's website.
- Q3k. I think instructions on how to install or access the API, update it, and the frequency of its release cycle is all useful information to know about. e.g., a list showing the latest releases, what was added and how to update your application to make use of it.
- Q3l. Error codes describing specific problems with an API are not helpful. e.g., a list of canonical HTTP error codes and how to interpret them.

Rationale-specific documentation of web APIs

When answering these questions please answer with respect to **your own experience** in learning web APIs (if applicable). Any examples provided exist solely to help illustrate the statement. For each question, please nominate how much you agree with the following statements: [*Strongly agree, Somewhat agree, Neither agree nor disagree, Somewhat disagree, Strongly disagree*]

- Q4a. I think that, as a starting point when beginning to learn about an API, I would like to read about descriptions of the API's purpose and overview.
- Q4b. I don't find descriptions of the types of applications the API can develop helpful.
- Q4c. I believe that descriptions of the types of developers who should and shouldn't use the API is important to know.
- Q4d. I don't think that descriptions of the types of end-users who will use the product built using the API is important to know in advance.
- Q4e. I think that if I read success stories about when the API was previously used in production, I would have a better indicator of how I could use that API.
- Q4f. I think that documentation that compares an API to other, similar APIs confusing and not important.
- Q4g. I believe it is important to know about what the limitations are on what the API can and cannot provide.

Conceptual-specific documentation of web APIs

When answering these questions please answer with respect to **your own experience** in learning web APIs (if applicable). Any examples provided exist solely to help illustrate the

statement. For each question, please nominate how much you agree with the following statements: *[Strongly agree, Somewhat agree, Neither agree nor disagree, Somewhat disagree, Strongly disagree]*

- Q5a. I wouldn't read through theory about the API's domain that relates theoretical concepts to API components and how both work together.
 - Q5b. I think it is important to know the definitions of the API's domain-specific terminology and concepts (with synonyms where needed). e.g., a computer vision API that uses machine learning should list machine learning concepts.
 - Q5c. It's not really important to document information about the API to non-technical audiences, such as managers and other stakeholders. e.g., pricing information, uptime information, QoS metrics/SLAs etc.
-

General-support documentation of web APIs

When answering these questions please answer with respect to **your own experience** in learning web APIs (if applicable). Any examples provided exist solely to help illustrate the statement. For each question, please nominate how much you agree with the following statements: *[Strongly agree, Somewhat agree, Neither agree nor disagree, Somewhat disagree, Strongly disagree]*

- Q6a. I find lists of Frequently Asked Questions (FAQs) helpful.
 - Q6b. When something goes wrong, I don't read through troubleshooting suggestions for specific problems straight away as I like to solve it myself.
 - Q6c. I like to see diagrammatic representations of an API's components using visual architectural visualisations. e.g., UML class diagram, sequence diagram.
 - Q6d. I wouldn't look for email addresses and/or phone number for technical support in an API's documentation.
 - Q6e. I generally refer to a programmer's reference guide or textbook about the API when I need to.
 - Q6f. I don't think it's important to read about the licensing information about the API.
-

The effect of structure and tooling on web API documentation

When answering these questions please answer with respect to **your own experience** in learning web APIs (if applicable). Any examples provided exist solely to help illustrate the statement. For each question, please nominate how much you agree with the following statements: *[Strongly agree, Somewhat agree, Neither agree nor disagree, Somewhat disagree, Strongly disagree]*

- Q7a. I would like to use a searchable knowledge base to find information.
- Q7b. I think a context-specific discussion forum between developers isn't very helpful as it just introduces noise. e.g., issue trackers, Slack group.
- Q7c. I think links to other similar documentation frequently viewed by other developers would be useful. e.g., 'people who viewed this also viewed...'
- Q7d. If I get lost within the API's documentation, a 'breadcrumbs'-style of navigation isn't very useful to me.

- Q7e. A visualised map of navigational paths to common API components in the website would be useful to have. e.g., a large and complex API for Enterprise Service-Oriented Architecture where I could click into various boxes to read about components and arrows to read about how they are related.
- Q7f. I believe ensuring consistent look and feel of all documentation isn't necessary to a good API documentation.
-

Demographics

- Q8a. Are you, or do you aspire to be, a professional programmer? Or would you consider programming a hobby?
[Professional, Hobbyist]
- Q8b. How many years have you been programming?
[1–5 years, 6–10 years, 11–15 years, 16–20 years, 21–30 years, 31–40 years, 41+ years]
- Q8c. In what type of role would you say your current job falls into?
[Back-end developer, Data or business analyst, Data scientist or machine learning specialist, Database administrator, Designer, Desktop or enterprise applications developer, DevOps specialist, Educator or academic researcher, Embedded applications or devices developer, Engineering manager, Front-end developer, Full-stack developer, Game or graphics developer, Marketing or sales professional, Mobile developer, Product manager, QA or test developer, Student, System administration]
- Q8d. What level of seniority would you say this role falls into?
[Intern Role, Graduate Role, Junior Role, Mid-Tier Role, Senior Role, Lead Role, Principal Role, Management, N/A (e.g., I am a student), Other]
- Q8e. What industry would you say you work in?
[Cloud-based solutions or services, Consulting, Data and analytics, Financial technology or services, Healthcare technology or services, Information technology, Media, advertising, publishing, or entertainment, Other software development, Retail or eCommerce, Software as a service (SaaS) development, Web development or design, N/A (e.g., I am a student), Other industry not listed here]
-

*** End of Survey ***

APPENDIX C

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 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 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

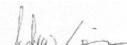
Author 4

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	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. <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. <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: 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

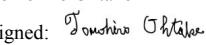
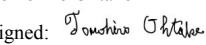
I agree to be named as one of the authors of this work, and confirm:

- 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 Abdelrazeck

 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 D

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 Abdelrazeq, 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