

The AIQ Meta-Testbed: Pragmatically Bridging Academic AI Testing and Industrial Q Needs

Markus Borg

RISE Research Institutes of Sweden AB
Dept. of Computer Science, Lund University
Lund, Sweden
markus.borg@ri.se

ABSTRACT

AI solutions seem to appear in any and all application domains. As AI becomes more pervasive, the importance of quality assurance increases. Unfortunately, there is no consensus on what artificial intelligence means and interpretations range from simple statistical analysis to sentient humanoid robots. On top of that, quality is a notoriously hard concept to pinpoint. What does this mean for AI quality? In this paper, we share our working definition and a pragmatic approach to address the corresponding quality assurance with a focus on testing. Finally, we present our ongoing work on establishing the AIQ Meta-Testbed.

KEYWORDS

artificial intelligence, machine learning, quality assurance, software testing, testbed

ACM Reference format:

Markus Borg. 2016. The AIQ Meta-Testbed: Pragmatically Bridging Academic AI Testing and Industrial Q Needs. In *Proceedings of ACM Conference, Washington, DC, USA, July 2017 (Conference'17)*, 7 pages. DOI: 10.1145/nmnnnnn.nnnnnnn

1 INTRODUCTION

The number of AI applications is constantly growing. Across diverse domains, enterprises want to harness AI technology to explore the lucrative promises expressed by AI advocates. As AI becomes pervasive, there is inevitably a need to build trust in this type of software. Furthermore, critical AI is on the rise, i.e., applications will not be restricted to entertainment and games. AI is already fundamental in many business-critical applications such as ad optimization and recommendation systems. As the technology further evolves, many believe that safety-critical AI will soon become commonplace in the automotive [22] and medical domains [18]. Other examples of critical AI, with other types of quality requirements, will be found in the finance industry and the public sector. Unfortunately, how to best approach Quality Assurance (QA) for AI applications remains an open question.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference'17, Washington, DC, USA

© 2016 ACM. 978-x-xxxx-xxxx-x/YY/MM...\$15.00

DOI: 10.1145/nmnnnnn.nnnnnnn

A fundamental issue originates already in the terminology, i.e., the concept of “AI quality”. First, there are several different definitions of AI, and their interpretations range from simple statistical analysis to the sentient humanoid robotics of the science fiction literature. Furthermore, AI appears to be a moving target, as what was considered AI when the term was coined in the 1950s would hardly qualify as AI today. Second, in the same vein, quality is a notoriously difficult aspect to pinpoint [34]. Quality is a multi-dimensional patchwork of different product aspects that influences the user’s experience. Moreover, quality is highly subjective and largely lies in the eye of the beholder. Taken together, AI quality is a truly challenging concept to approach, i.e., a subjective mish-mash of user experience regarding a type of technology with unclear boundaries that also change over time. There is a need for pragmatic interpretations to help advance research and practice related to AI quality – we provide ours in Section 3.

Contemporary AI solutions are dominated by Machine Learning (ML) and in particular supervised learning. A pragmatic first step would be to initially focus QA accordingly. As development of systems that rely on supervised learning introduces new challenges, QA must inevitably adapt. No longer is all logic expressed by programmers in source code instructions, instead ML models are trained on large sets of annotated data. Andrej Karpathy, AI Director at Tesla, refers to this paradigm of solution development as “Software 2.0” and claims that for many applications that require a mapping from input to output, it is easier to collect and annotate appropriate data than to explicitly write the mapping function.¹ As we embark on the AI quality journey, we argue that methods for QA of “Software 2.0” should evolve first – we refer to this as *MLware*.

The rest of this paper is organized as follows. Section 2 motivates the importance of *MLware* QA, elaborates on the intrinsic challenges, and presents closely related work. Section 3 introduces the working definitions used in our work on establishing the AIQ Meta-Testbed, which is further described in Section 4. Finally, Section 5 concludes our position paper.

2 BACKGROUND AND RELATED WORK

Fueled by Internet-scale data and enabled by massive compute, ML using Deep Neural Networks (DNN), i.e., neural networks with several layers, has revolutionized several application areas. Success stories include computer vision, speech recognition, and machine translation. We will focus the discussion on DNNs, but many of the involved QA issues apply also to other families of ML, e.g.,

¹bit.ly/3dKeUEH

support vector machines, logistic regression, and random forests – software that is not only coded, but also trained.

From a QA perspective, developing systems based on DNNs constitutes a paradigm shift compared to conventional systems [7]. No longer do human engineers explicitly express all logic in source code, instead DNNs are trained using enormous amounts of historical data. A state-of-the-art DNN might be composed of hundreds of millions of parameter weights that is neither applicable for code review nor code coverage testing [27] – best practices in industry and also mandated by contemporary safety standards. As long as ML applications are restricted to non-critical entertainment applications (e.g., video games and smartphone camera effects) this might not be an issue. However, when ML applications are integrated into critical systems, they must be trustworthy.

The automotive domain is currently spearheading work on dependable ML, reflected by work on the emerging safety standard ISO/PAS 21448. DNNs are key enablers for vehicle environmental perception, which is a prerequisite for autonomous features such as lane departure detection, path planning, and vehicle tracking. While DNNs have been reported to outperform human classification accuracy for specific tasks, they will occasionally misclassify new input. Recent work shows that DNNs trained for perception can drastically change their output if only a few pixels change [4]. The last decade resulted in many beaten ML benchmarks, but as illustrated by this example, there is a pressing need to close the gap between ML application development and its corresponding QA.

There are established approaches to QA for conventional software, i.e., software expressed in source code. Best practices have been captured in numerous textbooks over the years, e.g., by Schumeyer [29], Galin [12], Mistrik *et al.* [24], and Walkinshaw [34]. Developers write source code that can be inspected by others as part of QA. As a complement, static code analysis tools can be used to support source code quality. Unfortunately, the logic encapsulated in a trained ML model cannot be targeted by QA approaches that work on the source code level. ML models in general, and DNN models in particular, are treated as black boxes. While there is growing interest in research on explainable AI [1], interpreting the inner workings of ML is still an open problem. This is a substantial issue when explainability is fundamental, e.g., when safety certification is required [6] or when demonstrating legal compliance [33] (such as GDPR or absence of illegal discrimination in the trained model).

On the other hand, source code inspection and analysis are also not sufficient tools to perform QA of conventional software systems. During development, software solutions rapidly grow into highly complex systems whose QA rarely can be restricted to analysis – although substantial research effort has been dedicated to formal methods [35] including formal verification in model-driven engineering [14]. In practice, software QA revolves around well-defined processes [3, 15] and a backbone of software testing. Software testing, i.e., learning about the system by executing it, is the quintessential approach to software QA [13, 20, 25].

In the software engineering community, there is momentum on evolving practices to replace ad-hoc development of AI-enabled systems by systematic engineering approaches. A textbook by Hulten on “Building Intelligent Systems” [16] is recommended reading in related courses by Kästner at Carnegie Mellon University [21]

and Jamshidi at University of South Carolina. Kästner also provides an annotated bibliography of related academic research², as does the SE4ML group at Leiden Institute of Advanced Computer Science³, recently summarized in an academic paper [31]. Bosch *et al.* recently presented a research agenda for engineering of AI systems [8], sharing what they consider the most important activities to reach production-quality AI systems.

In recent years, numerous papers proposed novel testing techniques tailored for ML. Zhang *et al.* conducted a comprehensive survey of 144 papers on ML testing [36], defined as “any activities designed to reveal ML bugs” where an ML bug is “any imperfection in a machine learning item that causes a discordance between the existing and the required conditions.” Riccio *et al.* conducted another secondary study, analyzing 70 primary studies on functional testing of ML-based systems [32]. The authors do not use the term “bug” for misclassifications, as any ML component will sometimes fail to generalize. We agree with this view, and avoid terms such as ML bugs, model bugs and the like when referring to functional inefficiencies of MLware.

3 AI QUALITY ASSURANCE – WORKING DEFINITIONS

As discussed in Section 1, AI quality is a challenging concept to define. Consequently, QA for AI is at least as hard to specify. Still, we need a working definition to initiate efforts in this direction. In this section, we present the rationale behind our working definition of AI quality and AI quality assurance. Moreover, we introduce several related terms we use in collaborations with industry partners.

The original definition of AI from the 1950s is “*the science and engineering of making intelligent machines*”. Unfortunately, this definition turns AI into a moving target, as expectations on what constitutes an intelligent machine change over time – a computer program for logistics optimization in a warehouse would have been considered intelligent in the 1950s whereas it now could be part of an undergraduate computer science course. Since the term AI was introduced, it has often been used to refer to software solutions of the future, displaying increasingly human-like capabilities. The notation of “intelligence” is still common when referring to the gist of AI/ML applications, as in Hulten’s textbook [16], but ideally we want a definition that remains the same over time.

We argue that the most useful view on AI is to consider it as the next wave of automation in the digital society. Extrapolating from the sequence 1) digitization, 2) digitalization, and 3) digital transformation [28], we consider AI as the next enabling wave in the same direction – allowing automation of more complex tasks than before. Our working definition of AI is “*software that enables automation of tasks that normally would require human intelligence*”. While still imprecise, the definition is good enough for us to later define a delimited subset of AI that deserves our research focus.

Consulting the well-known textbook on AI by Russell and Norvig is one approach to explore the scope of AI [26]. The table of contents lists concepts such as searching, game playing, logic, planning, probabilistic reasoning, natural language processing, perception, robotics, and, of course, learning – all important components

²<https://github.com/ckaestne/seaibib>

³<https://github.com/SE-ML/awesome-sem1>

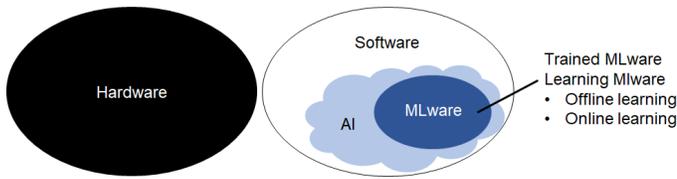


Figure 1: MLware in context.

when mimicking human intelligence. The textbook clearly shows that AI is more than ML. On the other hand, we argue that conventional software QA and testing can be applied to all AI techniques that are implemented in source code. Supervised and unsupervised learning, however, involves a transfer of control from source code to data. Research efforts on QA tailored for this new paradigm are what now would provide the highest return-on-investment. We need to focus on ML-enabled software – we refer to this as *MLware* for short.

Figure 1 illustrates our view on MLware. The future of systems engineering will combine hardware and software components, but the software part needs to be differentiated. A subset of software represents the fuzzy area of AI. We accept that this subset is neither clear-cut nor consistent over time. MLware is a subset of AI that rely on supervised and/or unsupervised learning. All MLware is not made the same. From a QA perspective, we need to distinguish between *trained MLware* that does not learn post deployment and *learning MLware* that keeps improving as new experience is collected post deployment. Learning MLware can be further divided into offline learning (triggered re-training in batches) and online learning (continuous update of trained models).

One might wonder where Reinforcement Learning (RL) fits in our working definition of MLware. Currently, we exclude RL from MLware. The rationale is that in RL, the exploration and exploitation of the learning agent is implemented in source code. RL shares characteristics of both searching and automatic control. We posit that software testing approaches proposed for self-adaptive systems could be generalized to RL [9, 23], and thus the best use of research resources is to focus on supervised and unsupervised learning – the dominating types of ML in practical applications.

A well-cited experience report by Sculley and his Google colleagues presents the vast and complex infrastructure required for successful MLware [30]. The authors describe this in terms of hidden technical debt of ML (cf. the lower part of Figure 2). Building on this discussion, and the expression that “data is the new oil”, our view is that data indeed fuels ML, but conventional source code is still in the driving seat, i.e., MLware is fueled by data and driven by code (cf. the upper part of Figure 2). From this standpoint, it is obvious that conventional approaches to software QA remain essential in the new data-intensive paradigm of MLware. Moreover, just as software QA is dominated by software testing, we expect MLware QA to be dominated by MLware testing.

The phenomenon of software quality has been addressed in plentiful publications. Among other things, this has resulted in standardized software quality models such as ISO/IEC 25010. As MLware still is software, and certainly driven by source code, the existing quality models remain foundational. The sister standard,

ISO/IEC 25012 Data Quality Model, adds a complementary data dimension to the quality discussion. As MLware is fueled by data, this standard is also highly relevant. Our working definition of AI quality is largely an amalgamation of the definitions provided by these two standards in the ISO/IEC 25000 series.

As mentioned in Section 2, there is no consensus in how to refer to issues resulting in MLware misclassifications. Bug is not a suitable term to cover all functional insufficiencies, given its strong connotation to source code defects. Still, we need a new similarly succinct term in the context of MLware. We propose *snag* to refer to the difference between existing and required behaviors of MLware interwoven of data and source code. The root cause of a snag can be a bug either in the learning code or the infrastructure [36], but it is often related to inadequate training data – we call the latter phenomenon a *dug*.

Figure 2 presents an overview of our perspective on issues detected in MLware. In the upper left, MLware is illustrated as a type of software that interweaves data (the fuel) and source code (at the helm) to produce output. If a discordance is observed, we call for a snag in the MLware fabric. Assuming that the requirements are valid and the observer interprets them correctly, root causes of snags include bugs and dugs as well as environment issues. The lower part of the figure illustrates the technical debt in machine learning as described by Sculley *et al.* [30]. Bugs can reside in the ML code (the white box), e.g., calling deprecated API methods or incorrect use of tensor shapes [17]. On the other hand, there might also be bugs in the rest of the infrastructure. While the illustrated technical debt revolves around data, all gray boxes will also depend on source code, from small exploratory scripts to mature open source libraries – and the large systems enabling MLware operations [16].

To summarize this section, our position is that research on QA for AI would benefit from adhering to the definitions presented in Table 1.

4 AIQ – AN AI META-TESTBED

Based on the working definitions in Section 3, we plan to support AI QA by establishing an AI meta-testbed. A testbed is a venue that provides a controlled environment to evaluate technical concepts. Under current circumstances, in the middle of the ongoing AI boom⁴, we believe that the establishment of a testbed for testing MLware testing would be the most valuable contribution to AI QA. Assessing the effectiveness of different testing techniques in a controlled setting is not a new idea [5], neither is the concept of testing test cases [37] – but a testbed dedicated to MLware testing is novel. We call it the AIQ Meta-Testbed⁵.

Successful MLware development requires a close connection to the operational environment. The same need has shaped software development at Internet companies, resulting in DevOps – a combination of philosophies, practices, and tools to reduce the time between development and operations while preserving quality [10]. Key enablers are Continuous Integration and Deployment (CI/CD). DevOps that emphasize MLware development is often referred to as MLOps [19], effectively adding Continuous Training (CT) to the

⁴Well aware of the two previous “AI winters”, periods with less interest and funding due to inflated expectations.

⁵metatest.ai

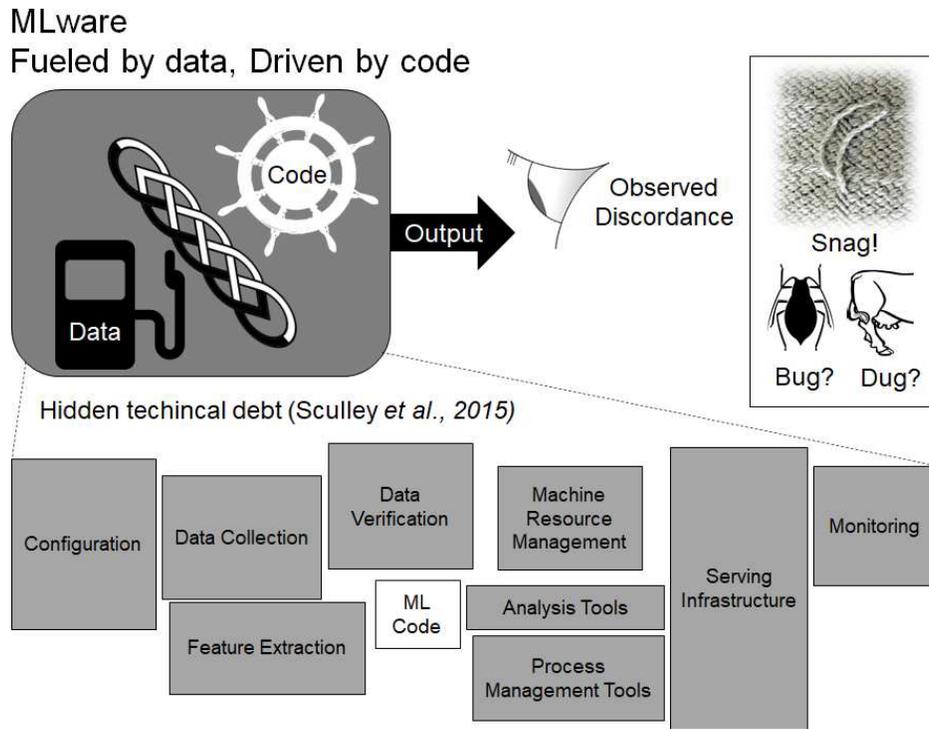


Figure 2: MLware interwoven by data and code. Observed discordances in the output (snags) can originate in source code defects (bugs) or data inadequacies (dugs).

mix. The focus on continuousness is stressed in illustrations by the infinity symbol.

Trust is fundamental for a successful product or service embedding MLware. In 2019, an expert group set up by the European Commission published ethics guidelines for trustworthy AI⁶. As part of the guidelines, seven key requirements are introduced. Table 2 shows a mapping between the EU requirements and the testing properties identified in the survey by Zhang et al. [36]. Our preliminary analysis indicates that all but one requirement has (to some extent) been targeted by academic research. Thus, we believe the time is right for systematic meta-testing in an MLOps context.

Figure 3 presents an overview of the AIQ Meta-Testbed in the MLOps context. We will set up a contemporary MLOps pipeline to allow controlled experiments in the lab while still providing an environment relevant to industry practice. Test automation is the backbone of MLOps, and MLware testing occurs in several phases during the MLware engineering lifecycle [36] (cf. the textboxes in Figure 3). First, the standard practice during model training is to split data into training, validation, and test subsets. We refer to this type of ML model testing as evaluation. Second, offline MLware testing occurs prior to deployment – conducted on different testing levels (input data, ML model, integration, system) and with varying access levels of the MLware under test (white-box, data-box, black-box) as defined by Riccio et al. [32]. Third, online MLware testing occurs after deployment. Common examples include A/B testing and runtime monitoring to detect distributional shifts.

⁶ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai

The AIQ Meta-Testbed will primarily focus on offline MLware testing (the solid-border textbox in Figure 3). We plan to enable meta-testing by providing a *control panel* for toggling testing techniques (C) in Figure 3) corresponding to the testing properties in Table 2, controlled *fault-injection* (A) (e.g., bug/dug injection, hyperparameter changes, mutation operators) and state-of-the-art *test input generation* (B) (e.g., search-based testing, GAN-based synthesis, metamorphic relations, and adequacy-driven generation). The results from both MLware testing and meta-testing will be presented in dashboards (D).

Extrapolating from the publication trends reported in the recent secondary studies [32, 36], there will be an avalanche of MLware testing papers in the next years. Staying on top of the research will become a considerable challenge and for practitioners with limited experience in reading academic papers, the challenge will be insurmountable – motivating the need to create an overview and shortlisting the most promising techniques.

Activities at the AIQ Meta-Testbed will include external replications of studies on MLware testing. By performing controlled meta-testing of the shortlisted techniques, we will be able to provide evidence-based recommendations on what techniques to use and in which contexts. The controlled environment of the AIQ Meta-Testbed will enable exploration of applied research questions, such as:

- Which contextual factors influence the MLware test effectiveness the most?

Table 1: Working definitions of key terms related to the AIQ Meta-Testbed.

Term	Definition	Comments
AI	A subset of software that automates tasks that normally would require human intelligence.	MLware, interwoven by data and source code, is the most precise term to describe our research interest. On the other hand, AI is a dominant term in industry and news media. We propose a pragmatic sacrifice of scientific preciseness in favour of industrial and societal relevance. In practice, we treat AI as synonymous with MLware in discussions with clients.
MLware	A subset of AI that, fueled by data, realizes functionality through supervised and/or unsupervised learning.	
MLware Testing	Any activity that aims to learn about MLware by executing it.	The typical goal of testing is detecting differences between existing and required behavior [2]. Other possible testing goals include exploratory testing and compliance testing.
AI Quality	The capability of MLware to satisfy stated and implied needs under specified conditions while the underlying data satisfy the requirements specific to the application and its context.	MLware combines data and conventional source code, thus we propose the amalgamation of corresponding quality definitions from the IEC/ISO 25000 series. Our proposal is in line with discussions by Felderer <i>et al.</i> in the context of testing data-intensive systems [11].
AI Quality Assurance	Any systematic process to provide confidence that the desired AI Quality is maintained.	QA encompasses many activities throughout the product lifecycle. However, in current AI discussions with clients, we primarily interpret it as MLware testing.
Snag	Any imperfection in MLware that causes a discordance between the existing and the required conditions.	There is an ongoing discussion in the research community about how to refer to MLware misclassifications [32]. We argue against using the term bug whenever there is unexpected output. Instead, we propose calling it a snag in the MLware fabric.
Bug	A source code defect that causes a discordance between the existing and the required conditions.	The term bug has a firmly established meaning, thus we suggest restricting its use to source code. As MLware is driven by code, bugs can cause snags.
Dug	A data inadequacy that causes a discordance between the existing and the required conditions.	With bugs reserved for source code defects, we need a novel expression for the data counterpart. The new term must be a worthy match for the succinct “bug”. Currently, we call them “dugs”.

- Which proposed MLware testing techniques scale to very large DNNs?
- How to best integrate MLware testing in an MLOps pipeline?
- What should be done to limit test maintenance in an MLware testing context?
- After observing a snag, how to support the subsequent root cause analysis?

5 SUMMARY AND CONCLUDING REMARKS

AI is becoming a pervasive subset of software, thus the elusive concepts of AI quality and QA are increasingly important. We argue

that pragmatic interpretations are needed to advance the field, and introduce a working definition of MLware as a subset of software within AI that realizes functionality through machine learning by interweaving data and source code. Furthermore, we define AI quality as “the capability of MLware to satisfy stated and implied needs under specified conditions while the underlying data satisfy the requirements specific to the application and its context”. We recommend that AI QA first and foremost should be interpreted as MLware testing and that the term bug shall be reserved for source

Table 2: Mapping the EU requirements for trustworthy AI and the testing properties targeted by publications on MLware testing as identified by Zhang *et al.* [36]. Gray cells show functional testing, i.e., the scope of Riccio *et al.*'s secondary study [32]

EU Key Requirements for Trustworthy AI	Testing Properties (Zhang <i>et al.</i> , 2020)
Human agency and oversight support (human autonomy and decision-making)	Correctness (probability to "get things right")
Technical robustness and safety (accuracy, reliability, reproducibility, resilience to attack)	Overfitting degree (model's capacity to fit the data)
Privacy and data governance (respect for privacy and integrity)	Robustness (resilience towards perturbations)
Transparency (traceability, explainability, communication)	Security (resilience against antagonistic attacks)
Diversity, non-discrimination and fairness (avoidance of unfair bias, accessibility)	Data privacy (ability to preserve private data)
Societal and environmental wellbeing (sustainability and social impact)	Efficiency (training and inference time)
Accountability (auditability, responsibility before and after use)	Fairness (freedom from discrimination)
	Interpretability (understandable decision rationales)

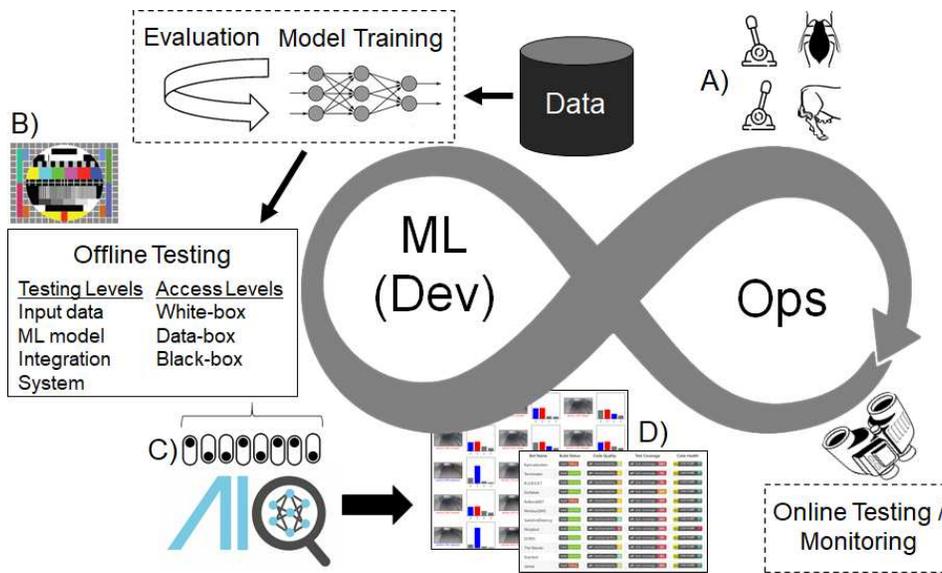


Figure 3: The AIQ Meta-Testbed in the MLOps context. We will focus on providing A) fault-injection, B) test input generation for offline testing, C) a control panel for toggling offline testing techniques, and D) presenting the results in dashboards.

code defects – instead we propose “snag” to refer to observed discordances in the MLware fabric. Finally, we present the AIQ Meta-Testbed – bridging academic research on MLware testing and industrial needs for quality by providing evidence-based recommendations based on replication studies in a controlled environment.

ACKNOWLEDGEMENTS

This work was funded by Plattformen at Campus Helsingborg, Lund University.

REFERENCES

- [1] Amina Adadi and Mohammed Berrada. 2018. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). *IEEE Access* 6 (2018), 52138–52160. <https://doi.org/10.1109/ACCESS.2018.2870052>
- [2] Paul Ammann and Jeff Offutt. 2016. *Introduction to Software Testing*. Cambridge University Press.
- [3] Noushin Ashrafi. 2003. The Impact of Software Process Improvement on Quality: In Theory and Practice. *Information & Management* 40, 7 (2003), 677–690. [https://doi.org/10.1016/S0378-7206\(02\)00096-4](https://doi.org/10.1016/S0378-7206(02)00096-4)
- [4] Aharon Azulay and Yair Weiss. 2019. Why Do Deep Convolutional Networks Generalize so Poorly to Small Image Transformations? *J. of Machine Learning Research* 20 (2019), 25.

- [5] V.R. Basili and R.W. Selby. 1987. Comparing the Effectiveness of Software Testing Strategies. *IEEE Transactions on Software Engineering* SE-13, 12 (1987), 1278–1296. <https://doi.org/10.1109/TSE.1987.232881>
- [6] Markus Borg. 2019. Explainability First! Cousteauing the Depths of Neural Networks to Argue Safety. In *Explainable Software for Cyber-Physical Systems (ES4CPS): Report from the GI Dagstuhl Seminar 19023*, Joel Greenyer, Malte Lochau, and Thomas Vogel (Eds.), 26–27. <http://arxiv.org/abs/1904.11851>
- [7] Markus Borg et al. 2019. Safely Entering the Deep: A Review of Verification and Validation for Machine Learning and a Challenge Elicitation in the Automotive Industry. *J. of Automotive Software Engineering* 1, 1 (2019), 1–19. <https://doi.org/10.2991/jase.d.190131.001>
- [8] Jan Bosch, Ivica Crnkovic, and Helena Holmstrom Olsson. 2020. Engineering AI Systems: A Research Agenda. *arXiv:2001.07522 [cs]* (Jan. 2020). <http://arxiv.org/abs/2001.07522>
- [9] Kai-Yuan Cai. 2002. Optimal Software Testing and Adaptive Software Testing in the Context of Software Cybernetics. *Information and Software Tech.* 44, 14 (2002), 841–855. [https://doi.org/10.1016/S0950-5849\(02\)00108-8](https://doi.org/10.1016/S0950-5849(02)00108-8)
- [10] FMA Erich, Chintan Amrit, and Maya Daneva. 2017. A Qualitative Study of DevOps Usage in Practice. *Journal of Software: Evolution and Process* 29, 6 (2017), e1885.
- [11] Michael Felderer, Barbara Russo, and Florian Auer. 2019. On Testing Data-Intensive Software Systems. In *Security and Quality in Cyber-Physical Systems Engineering*, Stefan Biffl, Matthias Eckhart, Arndt Lüder, and Edgar Weippl (Eds.). Springer International Publishing, 129–148. https://doi.org/10.1007/978-3-030-25312-7_6
- [12] Daniel Galin. 2003. *Software Quality Assurance: From Theory to Implementation*. Pearson, Harlow, England ; New York.
- [13] D. Gelperin and B. Hetzel. 1988. The Growth of Software Testing. *Commun. ACM* 31, 6 (1988), 687–695. <https://doi.org/10.1145/62959.62965>
- [14] Carlos A. Gonzalez and Jordi Cabot. 2014. Formal Verification of Static Software Models in MDE: A Systematic Review. *Information and Software Tech.* 56, 8 (2014), 821–838. <https://doi.org/10.1016/j.infsof.2014.03.003>
- [15] J. Herbsleb et al. 1997. Software Quality and the Capability Maturity Model. *Commun. ACM* 40, 6 (1997), 30–40.
- [16] Geoff Hulten. 2018. *Building Intelligent Systems: A Guide to Machine Learning Engineering* (1 ed.). Apress, New York, NY.
- [17] Nargiz Humbatova, Gunel Jahangirova, Gabriele Bavota, Vincenzo Riccio, Andrea Stocco, and Paolo Tonella. 2020. Taxonomy of Real Faults in Deep Learning Systems. In *Proc. of the 42nd Int'l Conf. on Software Engineering*.
- [18] Fei Jiang et al. 2017. Artificial Intelligence in Healthcare: Past, Present and Future. *Stroke and Vascular Neurology* 2, 4 (2017), 230–243. <https://doi.org/10.1136/svn-2017-000101> arXiv:<https://svn.bmj.com/content/2/4/230.full.pdf>
- [19] Ioannis Karamitsos, Saeed Albarhami, and Charalampos Apostolopoulos. 2020. Applying DevOps Practices of Continuous Automation for Machine Learning. *Information* 11, 7 (2020), 363.
- [20] Mohamad Kassab, Joanna F. DeFranco, and Phillip A. Laplante. 2017. Software Testing: The State of the Practice. *IEEE Software* 34, 5 (2017), 46–52. <https://doi.org/10.1109/MS.2017.3571582>
- [21] Christian Kästner and Eunsuk Kang. 2020. Teaching Software Engineering for AI-Enabled Systems. *arXiv:2001.06691 [cs]* (Jan. 2020). <http://arxiv.org/abs/2001.06691>
- [22] Hod Lipson and Melba Kurman. 2016. *Driverless: Intelligent Cars and the Road Ahead*. MIT Press.
- [23] Sara Mahdavi-Hezavehi et al. 2017. A Systematic Literature Review on Methods That Handle Multiple Quality Attributes in Architecture-Based Self-Adaptive Systems. *Information and Software Tech.* 90 (2017), 1–26. <https://doi.org/10.1016/j.infsof.2017.03.013>
- [24] Ivan Mistrik, Richard M. Soley, Nour Ali, John Grundy, and Bedir Tekinerdogan (Eds.). 2016. *Software Quality Assurance: In Large Scale and Complex Software-intensive Systems*. Morgan Kaufmann, Waltham, Mass.
- [25] Alessandro Orso and Gregg Rothermel. 2014. Software Testing: A Research Travelogue (2000ff?2014). In *Future of Software Engineering Proceedings*. 117–132. <https://doi.org/10.1145/2593882.2593885>
- [26] Stuart Russell and Peter Norvig. 2009. *Artificial Intelligence: A Modern Approach* (3 ed.). Pearson, Upper Saddle River, NJ.
- [27] Rick Salay, Rodrigo Queiroz, and Krzysztof Czarnecki. 2018. *An Analysis of ISO 26262: Machine Learning and Safety in Automotive Software*. SAE Technical Paper 2018-01-1075. <https://www.sae.org/publications/technical-papers/content/2018-01-1075/>
- [28] Daniel R. A. Schallmo and Christopher A. Williams. 2018. History of Digital Transformation. *Digital Transformation Now!* (2018), 3–8. https://doi.org/10.1007/978-3-319-72844-5_2
- [29] Gordon Schulmeyer. 1987. *Handbook Of Software Quality Assurance* (1 ed.). Prentice Hall, Lebanon, Indiana, USA.
- [30] D. Sculley et al. 2015. Hidden Technical Debt in Machine Learning Systems. In *Proc. of the 28th Int'l Conf. on Neural Information Proc. Systems*. 2503–2511.
- [31] Alex Serban, Koen van der Blom, Holger Hoos, and Joost Visser. 2020. Adoption and Effects of Software Engineering Best Practices in Machine Learning. In *Proc. of the 14th International Symposium on Empirical Software Engineering and Measurement*.
- [32] Riccio Vincenzo, Gunel Jahangirova, Andrea Stocco, Nargiz Humbatova, Michael Weiss, and Paolo Tonella. 2020. Testing Machine Learning based Systems: A Systematic Mapping. *Empirical Software Engineering (To appear)* (2020).
- [33] Andreas Vogelsang and Markus Borg. 2019. Requirements Engineering for Machine Learning: Perspectives from Data Scientists. In *Proc. of the 27th Int'l Requirements Engineering Conf. Workshops*. 245–251. <https://doi.org/10.1109/REW.2019.00050>
- [34] Neil Walkinshaw. 2017. *Software Quality Assurance: Consistency in the Face of Complexity and Change*. Springer.
- [35] Danny Weyns et al. 2012. A Survey of Formal Methods in Self-Adaptive Systems. In *Proc. of the 5th Int'l C* Conf. on Comp. Science and Software Eng.* 67–79. <https://doi.org/10.1145/2347583.2347592>
- [36] Jie M. Zhang et al. 2020. Machine Learning Testing: Survey, Landscapes and Horizons. *IEEE Transactions on Software Engineering* (2020), (Early Access). <https://doi.org/10.1109/TSE.2019.2962027>
- [37] Qianqian Zhu, Annibale Panichella, and Andy Zaidman. 2018. A Systematic Literature Review of How Mutation Testing Supports Quality Assurance Processes. *Software Testing, Verification and Reliability* 28, 6 (2018), e1675. <https://doi.org/10.1002/stvr.1675>

This figure "EUmapping.png" is available in "png" format from:

<http://arxiv.org/ps/2009.05260v1>

This figure "aiq.png" is available in "png" format from:

<http://arxiv.org/ps/2009.05260v1>

This figure "issues.png" is available in "png" format from:

<http://arxiv.org/ps/2009.05260v1>

This figure "mlware.png" is available in "png" format from:

<http://arxiv.org/ps/2009.05260v1>