# Towards Trustworthy Machine Learning in Production: An Overview of the Robustness in MLOps Approach

FIRAS BAYRAM, Department of Mathematics and Computer Science, Karlstad University, Sweden

BESTOUN S. AHMED, Department of Mathematics and Computer Science, Karlstad University, Sweden

Artificial intelligence (AI), and especially its sub-field of Machine Learning (ML), are impacting the daily lives of everyone with their ubiquitous applications. In recent years, AI researchers and practitioners have introduced principles and guidelines to build systems that make reliable and trustworthy decisions. From a practical perspective, conventional ML systems process historical data to extract the features that are consequently used to train ML models that perform the desired task. However, in practice, a fundamental challenge arises when the system needs to be operationalized and deployed to evolve and operate in real-life environments continuously. To address this challenge, Machine Learning Operations (MLOps) have emerged as a potential recipe for standardizing ML solutions in deployment. Although MLOps demonstrated great success in streamlining ML processes, thoroughly defining the specifications of robust MLOps approaches remains of great interest to researchers and practitioners. In this paper, we provide a comprehensive overview of the trustworthiness property of MLOps systems. Specifically, we highlight technical practices to achieve robust MLOps systems. In addition, we survey the existing research approaches that address the robustness aspects of ML systems in production. We also review the tools and software available to build MLOps systems and summarize their support to handle the robustness aspects. Finally, we present the open challenges and propose possible future directions and opportunities within this emerging field. The aim of this paper is to provide researchers and practitioners working on practical AI applications with a comprehensive view to adopt robust ML solutions in production environments.

CCS Concepts: • **Computing methodologies** → **Artificial intelligence**; **Machine learning**; • **General and reference** → *Surveys and overviews*.

Additional Key Words and Phrases: Artificial intelligence, machine learning, Trustworthy AI, robustness, MLOps systems, DataOps, ModelOps, model performance

## 1 INTRODUCTION

In the current landscape of real-world applications, the exponential generation of data coupled with their inherent complexity and dynamism requires the development of efficient analysis mechanisms. With the increasing difficulty in managing data-driven applications for monitoring and prediction, the deployment of artificial intelligence (AI) systems has played a vital role in real-world applications due to fast-paced advances and innovation in the field [137]. These AI

Authors' addresses: Firas Bayram, firas.bayram@kau.se, Department of Mathematics and Computer Science, Karlstad University, Universitetsgatan 2, 65188, Karlstad, Sweden; Bestoun S. Ahmed, Department of Mathematics and Computer Science, Karlstad University, Universitetsgatan 2, 65188, Karlstad, Sweden, bestoun@kau.se.
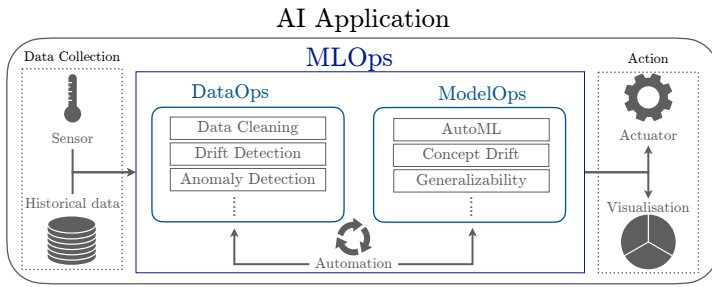
**111**

Fig. 1. MLOps system inside the AI application

systems, equipped with various toolkits, offer a wide range of solutions that can be adopted to help decision makers take rapid action in various applications [3]. Machine learning (ML), the heart of AI, empowers applications to learn and adapt automatically, paving the way for intelligent and autonomous operation. [78].

Traditionally, the development pipeline of ML systems involves the processing of historical data, which represents the raw input, to extract the features to train an ML model(s) that would perform the desired task offline to produce the output [45]. The pipeline includes additional sub-tasks like feature engineering, validation, testing, and performance evaluation. However, this static approach, while useful for initial research and experimentation, often proved inadequate in real-world scenarios. The fundamental challenge arises when the ML solution needs to be operationalized and deployed to function in a real-world environment [134]. In such scenarios, additional constraints and limitations should be specified to perform the tasks that characterize the ML solution, as the offline solution may not meet real-world requirements. For that, the transformation from offline to online ML system implementation has recently been of great importance to shift towards systems that perform in real-time and continuously circulate in production.

Machine Learning Operations (MLOps) have emerged as a potential recipe for designing ML solutions in online deployment scenarios. MLOps is an extension of DevOps, which is a well-known and mature topic in software engineering, to run ML systems in production [135]. MLOps is a set of practices and techniques that aims to automate the delivery of ML models to production [5]. Conceptually, MLOps adopts DevOps principles for building ML systems. In particular, the core principles MLOps uses from DevOps are continuous integration (CI) and continuous delivery (CD). Iteratively and continuously, CI/CD facilitates the automated process of releasing reliable software in short cycles [70]. However, ML systems exhibit different specifications and involve more components than traditional software systems, which require special management and maintenance. MLOps supports the development of ML systems in all phases of the entire ML workflow, from data retrieval to the release and deployment of the solution.

As illustrated in Fig. 1, MLOps serves as a vital companion throughout the lifecycle of AI products, ensuring the continuous monitoring and delivery of ML systems. However, there exists a great need to establish and clarify the elements and standards necessary for achieving optimal system performance. During the last decades, ML research has been highly dependent on quantitative metrics such as accuracy and loss measures to identify system performance and guide the selection of the ML approach [76]. However, at the holistic level of AI application systems, such metrics may not fully capture task performance. For example, a technically accurate solution that compromises ethical standards, such as utilizing private information in datasets or deploying an ML system with weak robustness to changes, poses significant concerns for AI systems. To address these

issues, researchers and practitioners in the AI community have begun to consider these aspects and introduce more principles while building AI systems [93]. The aspects and principles are known as AI trustworthiness [175].

Recently, trustworthy AI has made substantial strides in both academia and industry. Many governing bodies, organizations, and companies have formulated requirements and introduced principles to build trustworthy AI systems. The most popular definition of trustworthy AI was established by the European Union[1], which states that trustworthy AI systems have three key components: lawful, ethical, and robust [30]. Naturally, any MLOps framework designed within a trustworthy AI product must inherit the same principles and meet the requirements. In fact, from a technical standpoint, the MLOps pipeline is the most significant part of the AI system in production, and the implementation of such a trustworthy strategy must be thoroughly considered when devising an operationalized solution for an AI application.

Since MLOps is more concerned with the performance and quality of the system [131] and has less to do with ethical and lawful issues, we will narrow the focus of this paper to the technical robustness aspect of trustworthy AI systems. Robustness, in particular, is of major importance in practical MLOps systems. A robust ML system must be able to retain its effectiveness even when faced with unforeseen or unexpected challenges. The occurrence of system failures or inaccuracies in such settings can potentially cause interruptions or safety risks for users. Therefore, to maintain user trust and confidence in AI solutions, the robustness aspect of MLOps systems must be thoroughly addressed during development and deployment to secure sustained reliability in the long run and ultimately the responsible and effective deployment of ML in the dynamic real world.

This paper formalizes the specifications of robust ML systems in production through an in-depth literature search. In particular, we explore the connection between the robustness property, as a principal component of trustworthy AI systems, and the MLOps approach, as a primary methodology to streamline ML solutions. By formalizing this connection, which has not been exhaustively covered in the literature, we aim to provide researchers and practitioners with a consolidated view of the field. We summarize the contribution of this paper as follows:

- We organize the robustness aspects of ML solutions in production into three components and provide detailed definitions of the related concepts.
- We derive a set of specifications for a robust ML system, which is a seminal part of trustworthy AI systems, and reflect it with the workflow of ML solutions in production.
- We give a comprehensive overview of the specifications of robust ML systems and the constituent elements of MLOps systems.
- We review the literature on existing approaches that address the robustness aspects of ML systems in production.
- We iterate the available tools to build MLOps frameworks and discuss their support to deal with the various robustness aspects.
- Based on our broad analysis, we propose future directions and opportunities for researchers to further investigate in the area.

The remainder of this paper is organized as follows. In Section 2, we provide a general overview of the field and details of the research methodology followed in this work on surveying trustworthy AI and MLOps. In Section 3, we derive specifications to achieve robust MLOps systems and provide definitions related to the specifications of trustworthy AI systems and the elements of MLOps systems. In Section 4, we survey the approaches of academic research that address the robustness aspect of ML systems in production and summarize their pros and cons. We compare the different

---

[1]https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=60419

(a) Trustwothy AI.                                                    (b) MLOps.
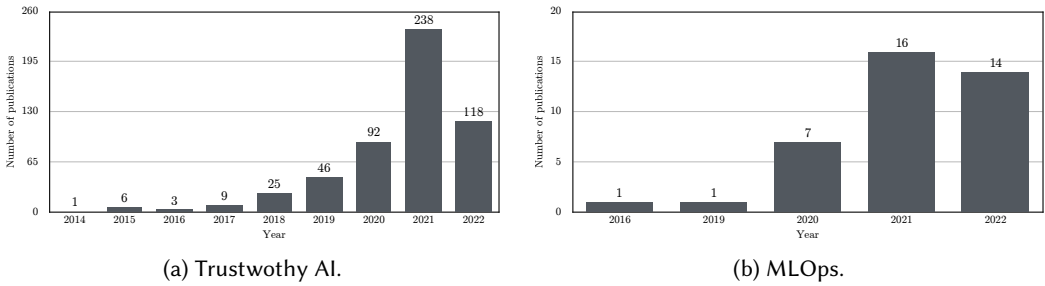
Fig. 2. Number of publications retrieved by searching the Web of Science indexing platform for the terms trustworthy AI and MLOps

technologies used to build MLOps systems in Section 5. We provide open challenges and potential future directions in Section 6. Finally, we conclude the survey with Section 7.

## 2 MOTIVATION AND RESEARCH METHOD

Trustworthy AI and MLOps are relatively new research topics in the maturing field of artificial intelligence (AI) in recent years. When searching for the terms trustworthy AI and MLOps on the Web of Science indexing platform[2], the results reveal that the two topics have begun to attract increasing research attention over the past five years. However, trustworthy AI is a more mature topic than MLOps in the literature. As shown in Fig. 2, trustworthy AI began to emerge in 2018 with 25 publications. The number increased to 238 publications in 2021 and 118 in 2022 (the year of writing this manuscript). Similarly to MLOps, publications have increased in the past two years, from 1 publication in 2019 to 16 in 2021.

### 2.1 Trustworthy AI

Trustworthy AI is a framework to verify that the system is operating according to the expectations of stakeholders [83]. With the rapid growth of AI applications, researchers demonstrated the urge to shift the focus to building trustworthy AI systems [177]. Therefore, trustworthy AI has been introduced in numerous application domains, such as self-driving cars [81], dentistry [103], the architecture, engineering and construction (AEC) industry [50], and education [170]. For example, in the context of self-driving cars, trustworthy AI ensures safe and reliable autonomous vehicle operation, promoting passenger and pedestrian safety. Similarly, the healthcare field also benefits from trustworthy AI systems that help in the precise diagnosis of diseases and the planning of treatments, ultimately improving patient care and treatment outcomes.

Given this, guidelines and specifications have been proposed to frame the limits and set the requirements for a system to be recognized as a trustworthy AI system. These guidelines establish a systematic framework for verifying and confirming the integrity of AI systems in different fields of application. In this section, we summarize the characteristics of trustworthy AI systems and highlight the robustness property and its significance in maintaining the trustworthiness of AI systems.

*2.1.1 Characteristics of Trustworthy AI systems.* The High-Level Expert Group on Artificial Intelligence (AI HLEG) was appointed by the European Commission to develop a guide document that describes the key specifications of trustworthy AI systems [38]. The guidelines aim to provide the principles to which every AI system should adhere. Nevertheless, the guidelines document consists
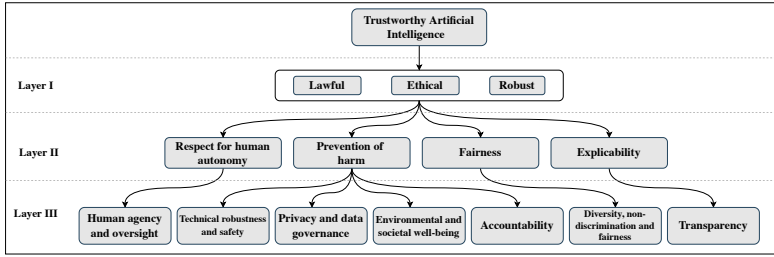
---

[2]https://webofscience.com/

Fig. 3. EU's requirement for trustworthy AI inspired by [38, 83].

of three layers of abstraction, from the most abstract to the most concrete. At the conceptual level, according to EU guidelines, to make an AI system trustworthy, it should align three main components at the top level throughout its evolution; the components are summarized as follows: [82]:

(1) **Lawful**: it ensures that the AI system must comply with all applicable laws and regulations, such as the General Data Protection Regulation (GDPR).
(2) **Ethical**: it ensures that the development of AI systems is in accordance with the ethical values of human society.
(3) **Robust**: it ensures that the AI system is technically and socially reliable so that it does not generate harmful risks that AI systems can cause.

The EU emphasized that the three components should be addressed during the development of AI systems and that the components should interact in deployment to provide a trustworthy service. To follow these guidelines, the document has identified four additional ethical principles for the development, deployment, and use of AI systems, which are: (a) respect for human autonomy, (b) prevention of harm, (c) fairness, and (d) explicability. These four principles are composed of seven fundamental requirements for trustworthy AI systems: (1) human agency and oversight, (2) technical robustness and safety, (3) privacy and data governance, (4) transparency, (5) diversity, non-discrimination, and fairness, (6) environmental and societal well-being, and (7) accountability. The EU requirements for trustworthy AI with different levels of abstraction are summarized in Fig. 3.

*2.1.2   Robustness property of Trustworthy AI systems.* Technical robustness is one of the imperative requirements for trustworthy AI systems, as discussed previously. The EU document related this requirement to the prevention of harm; in other words, the system should behave in a way that minimizes the associated risks and ensures the physical and mental integrity of humans. However, the technical details of the implementation of robust AI systems are not extensively discussed in the EU guidance document. Meanwhile, another detailed guide document has recently been published that surveys the topics relevant to building trustworthy AI systems has recently been published by the International Organization for Standardization (ISO) and the International Electrotechnical Commission (IEC) [75]. The technical report stated that the specification of trustworthy AI systems is out of the scope of the document; however, the report discussed related definitions and terminologies, threats and challenges, and mitigation measures for trustworthy AI systems.

The ISO/IEC report has provided a general definition of robustness as "*the ability of a system to maintain its level of performance under any condition.*" However, this general definition can be more concrete depending on the type of AI system. The robustness property generally ensures that the system can operate continuously according to its design. More specific definitions of robustness
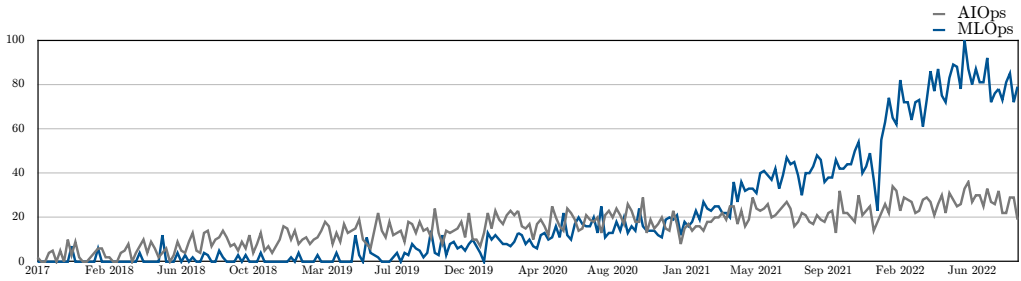
Fig. 4.  Google search trends [3]for AIOps and MLOps terms since 2017.

have been given in the ISO/IEC technical report depending on the functionality of the AI system, such as interpolation, classification, solving tasks or scoring.

On the other hand, while accuracy has traditionally been a key metric for the evaluation of AI systems [133], recent studies in ML [163] and deep learning (DL) [185] have emphasized the importance of robustness as another critical measure that highlighted a significant trade-off with the accuracy. Specifically, robustness encompasses a broader set of properties than accuracy as it pertains to the system's ability to maintain functionality amidst various challenges and disturbances. In contrast, accuracy is primarily focused on the predictive performance of the models and serves as a means to achieve robust AI systems [60].

Regarding other system performance metrics, robustness can impact latency through additional checks, but optimization techniques such as parallel processing can help mitigate this effect [178]. However, as robust systems become more complex, interpretability may decrease, although explainable AI can serve to bridge this gap. Furthermore, highly robust systems may pose challenges in terms of maintenance, as they involve complex error-handling mechanisms and configurations. Nevertheless, modular design principles and automation can help streamline maintenance processes and facilitate the management of robust systems over time [151].

Technically, robustness can also be recognized at both the data and algorithm levels [93]. Since the quality of ML solutions is highly dependent on input data, it is natural that data plays a major role in building robust AI systems. Therefore, data quality issues must be carefully investigated and formulated before implementing an AI solution. An important root of data quality issues is that, in real-life applications, data corruptions and inconsistencies are usually unanticipated. Another major issue is that the data will most likely follow changes during the evolution of the system [161]. Likewise, the robustness property must be considered at the algorithmic level, especially in high dimensions, where computations are relatively expensive [40]. In most cases, the robustness of the algorithm in ML is coupled with the quality of the data, since it mainly implies the insusceptibility to data disturbances and the ability to generalize training inputs.

## 2.2  Machine Learning Operations (MLOps)

The field of MLOps is a relatively nascent topic that studies standardized methodologies and strategies for the deployment of ML solutions in real-world AI applications. In the literature, the term AIOps can also be found to describe ML systems in operations [108]. However, according to Google Trends data[3] and as shown in Fig. 4, MLOps has started to be used more frequently than AIOps after 2021. MLOps is a heterogeneous discipline based on the application of DevOps principles to ML software systems [77]. MLOps provides a comprehensive framework for managing the

---

[3]https://trends.google.com/

entire lifecycle of ML models, including data acquisition, data preparation, model training, testing, evaluation, deployment, and continuous monitoring. This approach ensures efficient development and effective maintenance of ML systems in operational environments [63].

The MLOps lifecycle involves several key stages. It begins with ensuring high-quality, reliable data through the collection, cleaning, and preprocessing of raw data to make it suitable for model training. Next, machine learning models are developed and trained by selecting appropriate algorithms, tuning hyperparameters, and iterating to improve performance [136]. Rigorous testing evaluates the model's performance on unseen data to ensure generalization and meet accuracy and reliability standards. The trained model is then deployed into a production environment to make predictions on new data without disrupting existing services. Continuous monitoring of the model's performance in production is essential to detect any degradation, requiring updates and retraining with new data to maintain consistent performance. Managing data versions and model artifacts is crucial for reproducibility and compliance, ensuring all steps in the ML pipeline are well-documented.

In MLOps, testing goes beyond traditional software approaches to specifically address the unique challenges of ML systems. This includes rigorous validation of data quality during collection and preprocessing to optimize model training [191].Unit testing examines specific components like data pipelines and model algorithms to guarantee consistent performance with different datasets. Integration testing validates the seamless interaction between these components throughout the ML lifecycle, from data ingestion to deployment [112]. Additionally, MLOps emphasizes model explainability and robustness testing to enhance transparency and reliability against varying data conditions [125]. By integrating these specialized testing methodologies, MLOps facilitates the deployment of reliable and scalable ML solutions that meet the complex demands of real-world AI applications.

## 3  SPECIFICATIONS OF ROBUSTNESS FOR MLOPS SYSTEMS

Generally speaking, robustness is a broad concept that covers a wide range of aspects and properties [22]. As stated in the ISO document [75], robustness is usually defined for a specific problem in a particular context or domain. In the context of MLOps, robustness and resilience are crucial but distinct concepts [19]. Therefore, researchers from different fields have proposed a set of requirements to provide a robust formulation of their problems [86, 101]. Robustness focuses on the ML system's ability to maintain performance under varying and challenging data conditions, ensuring consistent and accurate predictions even with noisy or adversarial data. Resilience, however, pertains to the ML system's capacity to recover from failures or disruptions, ensuring the entire pipeline can continue operating or quickly resume normal functioning after issues such as hardware failures or software bugs. Our main focus is the robustness property, as it is crucial for guaranteeing the reliability and performance of models in real-world scenarios where data variability is common, ensuring the effectiveness of AI solutions in production environments [152]. In this section, and based on an exhaustive literature search, we dissect the specifications of robust MLOps frameworks, which is a seminal part of building trustworthy AI systems in production environments.

To build an overall robust MLOps system, robustness specifications should be addressed in all system processes. Practically speaking, to build a holistic robust ML system in production, we should address the issues and challenges that affect system robustness at both the data levels and the algorithmic levels, in addition to employing a robust automation process that drives the entire system. However, none of the existing frameworks addresses all aspects of building a robust system. Meanwhile, existing frameworks partially handle the robustness of ML systems; examples are: robustness to outliers [56], robustness to drift [117], or robustness to noise [68]. Therefore, we drill down to the level of granularity of the MLOps system ingredients, i.e., automation, DataOps,
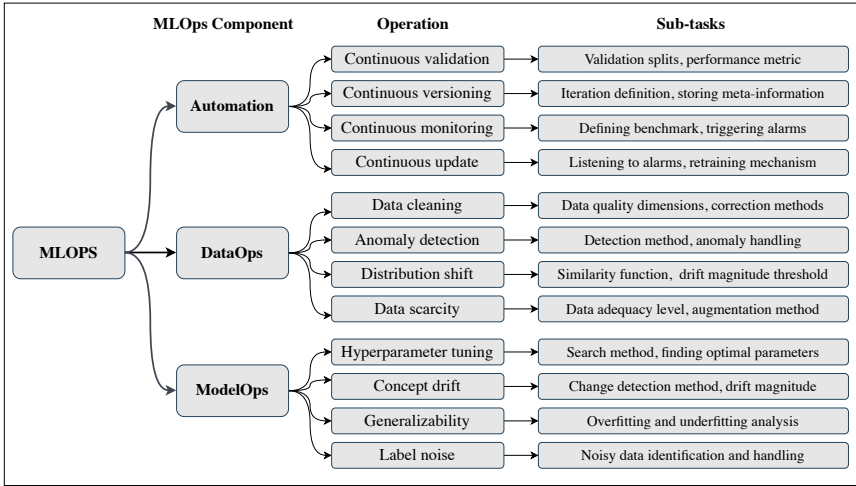
Fig. 5. Summary of MLOps components, operations and the corresponding sub-tasks

and ModelOps. The overall components and operations of the MLOps system, together with the corresponding subtasks, are summarized in Fig. 5.

Illustrating robustness with a concrete example of a traffic prediction system utilized in a navigation application. Robustness in this context refers to the system's ability to adapt to unforeseen situations in road conditions, such as accidents or construction, to deliver accurate and timely route recommendations for users. At the data level, challenges may arise from inconsistencies or inaccuracies in input data, such as missing or incorrect traffic flow information. Employing robust data processing techniques such as data cleansing and anomaly detection can ensure the overall validity of the input data. Furthermore, at the model level, robustness involves developing algorithms that are resistant to variations in traffic patterns and environmental factors. Therefore, robust ML models should generalize well across different road networks and traffic conditions, effectively mitigating the impact of data and algorithmic challenges on the accuracy of prediction. In this way, the system reduces the risk of providing inaccurate navigation guidance, promoting user satisfaction and safety during their journeys.

## 3.1 Robustness in Automation

MLOps systems are based on a continuous automated process that validates, monitors, and updates ML artifacts to maintain high performance during the evolution of the system [150]. To obtain elastic and reliable MLOps frameworks in production, a robust automation system should be implemented to allow rapid updates in the pipeline [79]. To achieve that, MLOps leverages the principles of DevOps, CI/CD, which offer practical techniques to manage the service throughout its development stages [171]. These techniques ensure that elements of the ML system are regularly managed and updated to optimize performance [55]. Since automation involves several steps, ranging from system build-up to system monitoring, the robustness aspect should be addressed at every step.

*3.1.1 Robust Continuous Validation.* ML validation is a diagnostic step that precedes the deployment of the designed solution; it aims to verify the performance under conditions different from those used when building the solution [166]. During ML validation, the robustness of performance could be tested against bias towards training conditions [130]. Therefore, it governs the process

of identifying potential solutions for real deployment environments by estimating performance under a wide range of problem settings [92]. Furthermore, the implementation of a well-defined validation mechanism not only leads to more reproducible solutions [83] but also helps mitigate the risks associated with ML models in general, such as security concerns [111], adversarial attacks [148], and risks related to algorithms and data [158].

The ML validation methodology generally uses several configurations, including dataset splits, performance evaluation metrics, and model hyperparameter settings [172]. The hyperparameter settings of the model are considered an aspect of model robustness [46] and will be further investigated in a separate section. For the other configurations, there are several techniques and settings that can be found in the literature to develop a robust model evaluation methodology. In particular, how to select the validation datasets and performance metrics to evaluate the performance of the tested solution.

- **Selection of Validation Split.** The factor that typically drives the choice of (a) suitable validation split (s) is the robustness aspect that is being tested during the validation phase. Therefore, to increase the robustness of the ML solution, the system should be thoroughly tested against various potential challenges, such as anomalies, drift, or adversarial attacks [27]. Following that design, it is necessary to prepare curated validation splits that contain potential production data issues to thoroughly verify the robustness of the solution. Mostly, validation splits should include outliers, missing data, or distribution shifts. Then, the performance of the system is recorded across the validation splits. Traditionally, many techniques divide the dataset and achieve more robust systems [75]. Cross-validation is the most common statistical method to estimate the performance of an ML model under specific conditions. In the literature, many techniques to perform cross-validation exist. Generally, K-fold cross-validation is the standard procedure for splitting the dataset. The procedure splits the dataset into $K$ subsets that are used to train and validate the model [8]. Some other cross-validation methods are single hold-out random subsampling and leave-one-out cross-validation [18]. However, since building robust solutions requires a systematic way of selecting the splits, random subsampling would not necessarily guarantee the creation of robust splits.
- **Selection of Performance Metric.** After selecting the validation split, the performance of the ML model in the validation split is statistically quantified. Typically, the learning task determines which performance metric to choose. For example, the most popular metrics in regression problems are empirically calculated based on the distance between the actual and predicted values [21]. On the contrary, in classification, most performance metrics are empirically calculated based on the fraction of correctly classified predictions [140]. However, in some problems, a single performance metric does not necessarily provide a good indication of the robustness of the solution. For example, in class imbalance problems, the accuracy measure does not provide a reasonable estimate of the quality of the solution, and the F score, which is a combination of precision and recall measures, is a more explanatory indicator of performance [59]. Another example is Cohen's kappa coefficient [37] for multilabel classification problems. However, there are other popular performance metrics, such as the area under the receiver operating characteristic curve (AUC-ROC), the confusion matrix, and the coefficient of determination $R^2$ [51]. Therefore, it is essential to carefully select several performance metrics in the evaluation phase to capture robustness characteristics from different perspectives.

*3.1.2 Robust Continuous Versioning.* Building ML solutions in production is an iterative process that accompanies the development of the system throughout its life cycle [34]. The ML artifacts will follow changes over time, so it is fundamental to record the historical versions to obtain robust

and tractable MLOps systems. Therefore, incorporating a centralized version control mechanism, a principle derived from continuous integration, is the first step after finalizing a releasable solution [162]. The main benefits of ML versioning are: reproducibility and traceability [110]. On the one hand, using a unique version number allows us to locate and trace the roots of the fault, and on the other hand, it also allows us to reuse a specific state of the system [182]. Continuous versioning requires an effective control mechanism for the governance of ML artifacts. A robust versioning mechanism must ensure that the different ML artifacts, such as configuration and hyperparameters, data and algorithm, have consistent version numbers for a specific iteration [67]. Moreover, the versioning mechanism should clearly define what is considered an iteration. To determine the variations in the ML artifacts to uniquely identify a distinct iteration in the MLOps life cycle.

*3.1.3   Robust Continuous Monitoring.* Lifelong monitoring of MLOps systems is a continuous process to collect feedback on production quality. Monitoring of the post-production model actively checks system performance and diagnoses various elements of the system to alert to abnormal behaviors [23]. Since the construction of MLOps systems involves the integration of several components into the system infrastructure, it is essential to establish a comprehensive monitoring strategy to cover these components. In production, model monitoring includes detecting degradation in both model performance and data quality, such as outliers, distributional shift, or concept drift [87]. Monitoring requires a benchmark to compare the behavior of the current iteration with a pre-defined reference. For example, an evaluation is performed on the basis of the ground-truth labels for the production data to identify the degradation of the ML system's performance [187]. However, the ground-truth might be costly (sometimes impossible) to collect or unavailable immediately [202]. Another example is monitoring the input data drift, which will be examined in more detail later, which involves analyzing the distribution of production data. Thus, a reference distribution should be available for comparison.

*3.1.4   Robust Continuous update.* An update signal is triggered to refresh the learning system to maintain the robustness of the ML model's performance and keep them up to date with the newly arriving data. This step is activated mainly by feedback from model monitoring to extend the existing solution [75]. Typically, the learning algorithm continuously uses the new output of the system to produce new versions of the ML model, known as the feedback loop. To update the model, two key considerations are identified in the literature [138]; the first is when to update the model. The second is how to execute the update. The first consideration focuses on deciding whether we need to update our learning pipeline. This is closely interlinked with model evaluation results; specifically, a threshold should be set, such as the level of performance degradation or the magnitude of distributional shifts. If these thresholds are exceeded, an update signal should be activated to update the ML artifact with the current context of the environment [165]. A typical example is to re-train the model after following a change in the system, for instance, concept drift occurrence. The second consideration has to do with the update procedure of the model. There are many alternatives to update the model, such as performing re-training, adaptation, or replacement of the model [127].

## 3.2   Robustness in DataOps

ML solutions are mainly based on data manipulation to infer patterns in the system [105]. Data are pivotal for most of the phases of MLOps pipelines, as they are used to understand the characteristics of our problem and interpret the event episodes of our observations, in addition to practical use to train and validate the ML model. Due to this, robust data operations should be employed in the data artifact to increase the robustness of the overall MLOps system. DataOps is a new term coined collectively by data engineers, data scientists, and data analysts [118]. DataOps refers to

end-to-end data processing operations in production and is considered an integral part of the entire MLOps system [10]. The challenges of the DataOps pipeline can be roughly divided into two groups. The first group deals with traditional data quality issues that primarily affect the accuracy of the predictive learning task. The second group deals with data operationalization, which is primarily concerned with data handling [26].

The e-commerce industry serves as an example of how robust DataOps practices are implemented in real-world scenarios to maintain operational efficiency and provide customized shopping experiences to customers. For example, within a large online retailer, ML systems analyze customer behavior and preferences to generate personalized product recommendations and targeted marketing strategies. Here, robust DataOps processes are key for managing extensive volumes of transactional data, customer interactions, and inventory records. In such a context, the retailer aggregates and processes data from various sources, including website interactions, purchase histories, and demographic profiles. Robust DataOps practices ensure data quality and reliability by employing data cleansing techniques to eliminate inconsistencies and inaccuracies. Furthermore, ongoing monitoring of data quality helps detect anomalies or discrepancies, facilitating prompt corrective actions.

*3.2.1   Robust Data Cleaning.* Data quality is a cornerstone of any data-driven solution since it may increase/decrease the accuracy of the analytic results [119]. Real-life applications, in which MLOps systems are used primarily, suffer from a high rate of erroneous data recorded from sensors that require cleaning and corrections. Therefore, most of the effort to build an ML solution is devoted to data cleaning to improve data quality [43]. Methods for processing data value errors can be roughly divided into two types: discarding erroneous data or cleaning erroneous data [174]. The latter can also be divided into manual and automatic cleaning. Manual cleaning provides more accurate results, but is more expensive to implement [174]. Data cleaning generally involves performing several activities to foster high levels of data reliability. These activities should be designed coherently to properly handle data corruption and imperfections. The main dimensions of data quality can be summarized as follows [36]:

- **Handling Missing and Duplicate Values.** Missing values and duplicates are very common in real-world applications and handling them is the most basic data cleaning task [48]. In real-world scenarios, missing values can appear for many reasons, such as sensor failures in industrial problems or incomplete laboratory test results in clinical problems [29]. From a robustness perspective, missing data is problematic for several ML models and should be resolved. As an example, the classification and regression tree (CART) models handle the issue by using surrogate splits to mitigate the effect of missing values [153]. Traditionally, missing values are handled by discarding the tuple that includes incomplete data or replacing missing values with a statistically inferred value, known as imputation [100]. Another issue of data quality is data duplication, or redundancy, which is a term that refers to the situation where some data records are repeated more than once in the dataset [196]. Similarly to missing values, the duplication of data samples has a negative implication in learning, as it can skew the input distribution [35], resulting in duplicate bias [7]. To detect duplicate values, Euclidean distance algorithms and Dynamic Time Warping (DTW) are the most widely used techniques, with the latter providing more convenient results in the case of big data [199]. Data fusion is the most widely used approach to correct data duplication [114].
- **Anomaly Detection.** In the ML community, there is a large body of literature dedicated to studying the detection of extreme values in the dataset, known as anomaly detection. Anomalies are classified into three types: point, contextual, or collective anomalies [28]. Extensive prior works in anomaly detection have explored various methodologies and applications

across different fields, providing a foundational understanding that is directly applicable to MLOps. In addition to the negative effect that anomalies have on the performance of predictive systems [72], anomalies can indicate malicious activities or adversarial attacks [74], which can affect the robustness of the overall system. The integration of these established anomaly detection techniques is crucial in MLOps pipelines, where they serve to identify and mitigate abnormal behavior, thereby maintaining system integrity and performance. Thus, anomaly detection is an essential step in any MLOps pipeline to detect abnormal behavior in the system and alleviate associated problems.

- **Domain-based Constraints.** Another step in the data cleaning phase is to deal with data samples that violate the restrictions on attribute values. In classical database management systems, the domain-based constraint is a type of integrity constraint that specifies the semantic rules of the data values that should follow [122]. In almost every application, there are specific rules for the types and values of the different features of the dataset. Simple examples of domain-based constraints are that ages cannot be negative or names cannot be numerical. Usually, domain experts specify the possible ranges and types of attribute values to accurately design a handling mechanism for inconsistent observations in the system [121].

*3.2.2 Robustness to Distribution Shift.* In production environments, the distribution of input data $P(X)$ often demonstrates distinct differences from that of training data due to the dynamic nature of post-deployment data [201]. For lifelong learning systems, the distribution change is an unresolved problem for predictive ML performance, as it tends to decline after the change [192]. Data distribution is a key foundation for many ML algorithms [176]. In deployment, MLOps systems should exhibit high robustness to distributional shifts; this requires systematic detection and adaptation mechanisms to cope with changes [44]. Theoretically, distributional shifts may not be associated with a change in the target concept, known as virtual drift [64]. However, in practice, the model should always be reviewed when there is a change in the data distribution [164]. The data used to train ML algorithms is assumed to be sampled from a distribution defined by an underlying generating function [66]. As discussed in Section 3.1.3, continuous monitoring strategies are deployed for a constant diagnosis of ML artifacts. Regarding distributional shifts, monitoring is generally performed by determining whether the data distribution of the present interval is significantly different from the past distribution using statistical tests [84]. However, some parameters must be defined to establish a robust monitoring strategy, such as the window size to estimate the probability distribution and the change significance threshold, which are selected primarily manually. For adaptation, there exist different techniques that increase the robustness of ML models to distributional shifts, such as importance weighting and uncertainty estimation [197].

*3.2.3 Robustness to Data Scarcity.* Some real-world applications generate a few data points, especially in the medical sector and many industrial production processes [200]. Other related issues caused by data scarcity are applications with underreported data that cause the problem to be sparse [47], or the problem of class imbalance, which is characterized by data-scarce regions of the underrepresented classes [102]. Therefore, the lack of historical data would hinder the development of a learning solution with good generalization performance. Inadequate training samples used to train the model in the training phase would lead to unexpected production behavior or bias towards the overrepresented population [11]. Therefore, data enrichment techniques, such as data augmentation, have been utilized to enhance data-scarce applications by generating more data samples based on the limited available data, or transfer learning by using pre-learned ML models and employing them in a similar task [142]. Additionally, Generative Adversarial Networks (GANs) offer a promising approach to address data scarcity [12]. GANs are a type of deep learning model

that can create entirely synthetic data closely resembling the real data distribution, allowing for significant data expansion. This is especially beneficial for scarce datasets, as demonstrated in applications like diagnosing COVID-19 from lung scans with limited real data [4]. Research has also explored the effectiveness of GANs in the energy sector for power consumption prediction and in defect detection [183].

*3.2.4    Robustness and Data Processing Resources.* In information systems, the resources are not infinite. Any solution developed must be designed within the limits of stringent capacity constraints. In relation to deployment systems, data manifest the bottleneck in the MLOps pipeline [128]. Undoubtedly, performing different data operations requires huge computing and storage resources that should typically be performed in real-time, which causes great difficulty for production systems [41]. As a result, the different stages of the DataOps pipelines should respect the trade-off between the complexity of the operations and the applicability of the solution online. However, not all solutions are suitable for deployment. More robust MLOps systems would require a large number of resources, which may not always be available [6]. Therefore, to determine the pertinence and simplicity of the solution, the MLOps pipeline should be tested in an environment identical to that of deployment.

## 3.3    Robustness in ModelOps

ML models in the wild must continually interact with the surrounding environment to acquire more knowledge and consequently improve their performance. Typically, the environment and its conditions are dynamic and, as the systems evolve, the MLOps systems should perform a set of operations to keep the ML models up-to-date [42]. ModelOps is a general term used to describe operations in any type of model and not specifically for ML models [71]. Researchers used the term ModelOps interchangeably with MLOps when the field emerged in 2018 and 2019 [162]. Subsequently, the MLOps term was mainly used to describe the entire framework that handles ML systems in production. Therefore, to alleviate the confusion, we will refer to the set of operations that are performed on the learning task of the ML model as ModelOps, to follow the naming convention of the MLOps discipline. This clear detachment of data and model operations would facilitate the delineation of the different elements of the MLOps system. To build robust learning algorithms in the long run, the ModelOps pipeline should be designed to be well-performing and address potential challenges and issues.

In our online retail scenario, ModelOps ensures that ML models deployed in production can effectively adapt to seasonal variations and other unforeseen events typical in the retail industry. Alongside addressing seasonal shifts in customer behavior, these models must also handle sudden changes in market trends, promotions, and customer preferences. Robust ModelOps practices involve continuously monitoring model performance and regularly retraining the model to accommodate these dynamic factors. Additionally, robust ModelOps pipelines include strategies for model interpretation, helping stakeholders understand and trust the decisions made by the ML system. Furthermore, ensuring the generalizability of ML models across different customer segments and market conditions is imperative for their success in real-world retail environments. In this section, we thoroughly explore the specifications for building robust ModelOps pipelines; see Section 2.2. We discuss related topics that impact the robustness of learning algorithms as found in the literature.

*3.3.1    Robust Hyperparameter Optimization.* Most learning algorithms require identifying optimal values for a set of hyperparameters that minimize the generalization error of a given function [17]. Thornton *et al.* [160] tested 39 popular ML classifiers on 21 datasets with different hyperparameter settings. The results showed that, on average, the accuracy of the model changes by

46% with the algorithm and the hyperparameter values. However, finding the best combination of hyperparameter values is not straightforward since it involves pre-defining several attributes, such as hyperparameter search methods and value ranges [62]. The challenge of the hyperparameter tuning methods of the ML model in production is the search space, which is very large for most learning algorithms [99], especially complex algorithms such as neural networks [52]. Therefore, researchers have called for robust and automatic hyperparameter optimization methods. This robust and automated mechanism must be addressed when designing the ModelOps pipeline. Robust hyperparameter optimization strategies would select the combination that produces consistent performance when used to evaluate both validation and test splits [124].

*3.3.2  Robustness to Concept Drift.* Another considerable challenge for every lifelong learning system, and MLOps in particular, is the change in the probability distribution that represents the relationship between input data and target variables $P(y|X)$, known as concept drift [149]. It is worth noting that concept drift and distribution shift, are types of the general dataset shift phenomenon [123]. Concept drift is a relevant topic in data stream mining research, wherein the incoming data arrive sequentially with a unique timestamp [1]. Sequential order allows one to segment the learned concepts into periods of time to check for changes. Various reasons can cause concept drift in the system, such as machinery degradation, variations in customer behavior and demands, or adversarial attacks. Similarly to the robustness of the distributional shift presented in Section 3.2.2, the concept drift is first diagnosed and detected to activate an adaptation strategy [97]. Typically, concept drift is detected by a degradation in the predictive performance of the model or a significant deviation in the distribution of input data, or hybrid monitoring for both [15]. However, the model is usually re-trained with the most recent data to provide the adaptation mechanism and more robust performance. Another approach is to utilize ensemble approaches to combine multiple results for the final output [49]. ModelOps should find the balance between the complexity of the concept drift handling solution and the time required to execute the adaptation, as concept drift could have severe consequences for the overall system if not addressed quickly.

*3.3.3  Model Generalizability.* One of the main properties that improve the robustness of an ML solution is the ability to generalize the pattern learned from training data to the entire domain of unseen distributions [173]. In ML research, the term generalizability is tied to underfitting and overfitting, which are generally used to analyze generalization performance. An optimal balance between underfitting and overfitting is required to achieve strong generalizability when building the ML solution [186]. In practice, ML models should generalize well using the available dataset without waiting to obtain large labeled datasets [198]. Therefore, data augmentation procedures have been leveraged to improve the generalization of ML models [198]. A principal approach to validate/improve generalization is to expand the training dataset by adding adversarial examples that can potentially be present in the production data. Other operations that can be found in the literature to handle the generalization of a learning model are cross-validation, regularization, and transfer learning [93]. However, the different operations should be integrated into ModelOps pipelines based on actively evaluating the generalizability of the ML model.

*3.3.4  Robustness to Label Noise.* ML models make decisions based on historical data used in the training phase. Therefore, any corruption in historical data would lead to poor deployment decisions. One type of corruption that may appear in the data is label noise, which affects the interpretation of the data [54]. There are different sources of label noise, including incomplete information, measurement errors, or incorrect labels in the data acquisition phase [88]. The removal of noisy labels is not always feasible, as it is difficult to determine whether the labels are indeed noisy [54]. For example, distinguishing between noisy and correct data records in duplicate records

Table 1. Summary of robust ML approaches in production

| MLOps Component | Operation | References |
|---|---|---|
| **Automation** | Continuous validation | [94, 126, 143] |
| | Continuous versioning | [167, 168, 190] |
| | Continuous monitoring | [39, 126, 145] |
| | Continuous update | [69, 141, 179, 180] |
| **DataOps** | Robust data cleaning | [24, 80, 139, 144, 188] |
| | Robustness to distribution shift | [14, 23, 32, 33, 57, 147] |
| | Robustness to data scarcity | [31, 128, 129] |
| | Resource scheduling | [9, 115] |
| **ModelOps** | Robust hyperparameter tuning | [2, 91, 113] |
| | Robustness to concept drift | [107, 169, 184] |
| | Generalizability | [61, 65, 194] |
| | Robustness to label noise | [58, 109, 132] |

assigned to different target variables in the same dataset is usually challenging or costly [181]. Therefore, the implementation of noise-tolerant strategies in learning algorithms is widely adopted to build solutions that are robust to noise. Researchers have developed model-based methods in noisy environments to obtain noise-tolerant versions of conventional ML learning algorithms. Probabilistic methods are also effective in dealing with label noise, such as Bayesian or clustering methods [53]. These modifications to the learning algorithms should be addressed while building ModelOps pipelines.

A final remark to note is that operations in the different stages of automation management, DataOps, and ModelOps, usually overlap with each other rather than being separate components of linear pipelines. This is why we can view automation processing as the engine of MLOps systems, which mainly distinguishes it from the classical offline frameworks for ML projects. Essentially, MLOps automation management controls and connects different operations in the life cycle of the ML system by activating a specific stage based on the results of another stage. Therefore, it can be run as a standalone service. Using the automated mechanism, the costs of human effort are minimal, and the workflow is more reliable and carried out faster.

## 4   LITERATURE REVIEW ON ROBUST ML APPROACHES

To compile the latest approaches devoted to developing ML systems that address the robustness property in deployment, we conduct an extensive literature search for existing solutions that handle one or more robustness aspects in production that have been proposed in the research community. We begin by reviewing existing surveys that are similar in scope to our work. We then The methods are organized into three levels: automation management, data, and algorithmic. Specifically, we survey the methods that are concerned with the robustness aspects that we presented in Section 3. The retrieved approaches are organized according to the robustness aspect addressed in Table 1.

### 4.1   Related Surveys

Several works have been published to characterize the trustworthiness of AI systems and their different properties. Kaur *et al.* [82] reviewed trustworthy AI approaches based on EU guidelines. The authors dissected some principles and suggestions for building trustworthy AI systems. However, the article mainly discusses the mechanisms of interactions between humans and machines to improve the performance of AI systems. In another article [30], Chatila *et al.* also discussed

the requirements to design trustworthy AI systems. The article discussed the alignment of the deployment and use of such systems with human rights and values. A more detailed article by Kaur *et al.* [83] presents a comprehensive overview of trustworthy AI systems. The article discussed existing strategies to standardize trustworthy AI systems. Furthermore, the article discussed approaches to mitigate AI risks by involving users and society to enrich AI systems. However, the article did not discuss in technical detail how to implement trustworthy AI systems.

A recent article published by Wing [175] supports the trustworthiness of AI systems by benefiting from other communities. The article discusses the advantages of bringing together researchers in different disciplines: trustworthy computing, formal methods, and AI. The author showed that connecting research communities across industry, government, and academia is crucial to establishing trustworthy AI. Similarly, Kumar *et al.* have discussed the requirements of trustworthy AI systems in pervasive computing and big data phenomena [90]. However, the authors reviewed the ethical issues of algorithms and data and provided a real-world application in smart cities.

From a computational perspective, a comprehensive survey has been presented to guide practitioners in achieving trustworthy AI [96]. The survey article provided a broad taxonomy of the field and recent advances. In particular, the article summarized the application of trustworthy AI systems in real-world applications. Furthermore, the article has reviewed the interactions among the different aspects of trustworthy AI. In a relevant work, Li *et al.* [93] have followed a systematic approach to introduce the aspects of trustworthy AI from theoretical and operational perspectives. The authors have also presented guidelines for the different stakeholders in the final product to construct trustworthy AI solutions. However, the article did not merely review existing approaches to build robust AI systems in production.

For MLOps, a few papers have reviewed the current state-of-the-art. The central part of the papers focused on current trends and challenges in building MLOps systems. For example, Kreuzberger *et al.* [89] have highlighted the open challenges in the field and presented an overview of the MLOps architecture. In a review paper, Zhao also discussed the current trends and challenges of MLOps systems [195]. In particular, the paper has focused on applying DevOps principles to ML projects and the additional components needed to satisfy the requirements of MLOps systems. The interplay between AutoML and MLOps has been investigated by Symeonidis *et al.* [156]. In addition to reviewing the current challenges, the authors have proposed an approach based on the incorporation of explainability and sustainability to increase the robustness of MLOps systems. Similarly, Tamburri has described the current challenges in the implementation of sustainability and interpretability in MLOps frameworks [157]. Meanwhile, Testi *et al.* [159] presented a recent taxonomy of MLOps in an overview paper. The paper has defined and standardized the methodologies for building MLOps frameworks. In addition, the paper has also proposed an operationalized approach to building ML systems. In terms of the stages and activities involved in the adoption of MLOps, a systematic review of the literature has described the details of the development of ML systems in operation [77]. Furthermore, the authors have presented a maturity model that should be followed to evolve MLOps practices. The importance of MLOps in data science has been investigated by Mäkinen *et al.* [106]. The authors have collected responses from 331 professionals from 63 countries working in the ML domain. The question was to investigate whether professionals tend to build data-centric or model-centric solutions, or both. The study found that the system should integrate both solutions in the construction of the MLOps pipeline.

## 4.2 Robust Approaches in Automation

Unlike conventional ML schemes in which algorithms work offline, in production, ML solutions are required to leverage the available knowledge to solve new problems continuously and automatically

Table 2. Summary of robust approaches in automation: Pros and cons

| Operation | Reference | Pros | Cons |
|---|---|---|---|
| **Continuous validation** | [94] [143] [126] | No manual intervention needed Configurable user interface Fully automated and scalable, language and library agnostic | Update timing is not dynamic Not fully automated Security and privacy issues are not addressed |
| **Continuous versioning** | [190] [167] [168] | Versioning models and data Versioning multiple pipelines simultaneously Fast iterative ML modeling | Cannot integrate external datasets Manual setting of the code version Still in the initial implementation |
| **Continuous monitoring** | [39] [145] [126] | Scalable and user-friendly access with low latency Front and back ends monitoring Real-time processing | Supports limited learning tasks Limited profiling capabilities Not tested in the long run |
| **Continuous update** | [180] [179] [141] [69] | No accuracy loss Very rapid retraining User-friendly for on-demand usage in edge and cloud systems Collaborative management through a web interface | Only tested on limited ML models Limited to specific types of models Limited investigation of real-world use cases Inference job underutilizes the allocated GPUs |

throughout their lifespan. Automation management involves several operations, including continuous validation, continuous versioning, continuous monitoring, and continuous update. This section surveys current methods that address these operations in the literature. The overall approaches are summarized in terms of their pros and cons in Table 2.

*4.2.1   Continuous Validation Approaches.* Li and Gui proposed a container-based Continuous Machine Learning and Serving (CMS) platform [94]. CMS platform incorporates offline and online stages for model validation to assess the robustness of the ML algorithm before using it in production. While it separates validation stages effectively, it lacks dynamic update timing and relies heavily on a pre-defined Mean Absolute Error (MAE) threshold. In offline validation, 20% of the dataset is used to test the model, the quality of the prediction is evaluated by calculating the Mean Absolute Error (MAE). Subsequently, a new model is introduced to the system if its MAE is less than the current model's error. Once a model passes the offline validation, the online validation is performed to evaluate the model in live data. A pre-defined MAE threshold is then used to evaluate the model's performance during an arbitrary time period. In another approach, Shahoud *et al.* have proposed a framework for managing ML tasks in big data environments [143]. The ML application is based on a microservice-based architecture that supports a scalable solution. Model validation is performed using 20% of the total dataset. To evaluate the performance of the ML model, the mean of the performance metric is calculated for each experiment that is repeated three times. To verify the efficiency of the solution in production, the authors focused on the execution time and overhead metrics of the framework of the different forecasting algorithms rather than prediction accuracy. For artificial intelligence (AIoT) applications, Raj *et al.* [126] have proposed a new robust automated framework for edge machine learning operations (Edge MLOps). The results of the root mean square error (RMSE) were validated by 10-fold cross-validation splits to select the best ML model

of four candidate models. However, integrating additional performance metrics and considering cross-validation techniques could enhance its robustness in model selection.

*4.2.2 Continuous Versioning Approaches.* Knowledge to Environment (K2E) platform was presented to govern data and model catalogues [190]. K2E architecture includes a version control manager that stores the meta-information of the various versions of data and models in the system. The meta-information incorporates the version numbers and pointers to each version's location in the internal datalake. The K2E platform offers efficient version control but lacks integration with external datasets, which limits its versatility. Weide *et al.* [167] have presented another versioning scheme. The scheme uses semantic versioning to distinguish between the types of change that occurred in the ML pipeline, which can offer a very useful interpretation of the evolution of the system. In addition, the scheme allows multiple versions to run simultaneously without conflict, but managing multiple simultaneous versions may pose additional complexity. For iterative model building, the SHERLOCK system was built that enables fast iterative modeling [168]. The system integrates a central repository, ModelDB, that stores all the meta-information about the model in tables. Meta-information includes model specification, model performance, and results of the different models. The meta-information is used later on to compare the models or reuse particular results.

*4.2.3 Continuous Monitoring Approaches.* Velox system has been developed for a scalable model serving pipelines in a low-latency environment [39]. Velox automatically monitors the quality of the model evaluated by selected performance metrics. Additionally, Velox keeps track of the errors associated with each ML model. Monitoring strategies allow the system to continuously verify the generalizability of the model and send re-training signals. However, Velox lacks comprehensive profiling capabilities. In another work, Silva *et al.* presented a monitoring and benchmarking environment that can monitor both the prediction quality and the system resources [145]. In addition to prediction performance, the approach is designed to monitor five metrics of resource consumption: completion time, CPU and GPU usage, memory usage, disk input/output (IO) operations, and network traffic. Edge MLOps [126] also monitors significant RMSE deviations from the performance of the ML model. A fixed threshold value is empirically defined to be used as a benchmark value to compare the current RMSE metric of the deployed ML model. However, relying solely on this fixed threshold may limit its effectiveness in performance management.

*4.2.4 Continuous Update Approaches.* A Provenance-Based Approach for Incrementally Updating Regression Models (PrIU), and its optimized version PrIU-opt, have been proposed to incrementally update the model parameters while preserving prediction performance, but are limited in their testing on specific models [180]. The approaches use gradient descent and its variants to trace the provenance of input data and, consequently, update the model parameters. Similarly, DeltaGrad was proposed for the fast re-training of ML models based on (stochastic) gradient descent [179]. DeltaGrad algorithm is inspired by ideas from the Quasi-Newton methods [25] and can be used to update ML models in production. In addition, in the case of changes in the data, the algorithm can delete the old data and learn from the newly acquired data. Continual-Learning-as-a-Service (CLaaS) paradigm was designed to address the continual learning of ML models in production [141]. Instead of re-training the ML model from scratch, CLaaS integrates several continual learning approaches that update the models upon signals that indicate drastic changes in the distribution of the data and the performance of the ML model. In a related approach, the ModelCI-e (continuous integration and evolution) plugin was designed to easily integrate with existing serving systems [69]. In the online phase, the system automatically checks for significant drifts in the data distribution and invokes an appropriate continual learning method if the drift magnitude exceeds a certain threshold.

Table 3. Summary of robust approaches in DataOps: Pros and cons

| Operation | Reference | Pros | Cons |
|---|---|---|---|
| **Data cleaning** | [80] | Handling missing data efficiently | Requires some manual efforts |
| | [139] | Scalable to large datasets | Not validated in streaming scenarios |
| | [24] | Supports single-batch and inter-batch validation | System's implementation influenced by Google's infrastructure |
| | [188] | User-friendly platform without the need for coding | Lower proficiency in programming may impact the task efficiency |
| | [144] | User-configurable flexibility | Potential information loss |
| **Distribution shift** | [57] | Real-time processing | Prone to false positives |
| | [147] | VAE-based detection of changes without ground-truth labels | Evaluated using use-cases not fully representing real-word scenarios |
| | [14] | Fully automated and scalable | Supports limited ML models |
| | [33] | Can handle complex and high-dimensional features | Only focused on image and textual data |
| | [32] | Low computation cost | Limited to classification tasks |
| | [23] | Rapid, low-latency inference | Only applicable to certain tasks |
| **Data scarcity** | [129] | Support multitask learning and a code-free approach | Does not have support for model versioning |
| | [128] | Reduces the cost and difficulty of training ML models | Requires large amount of unlabeled data |
| | [31] | Uses domain knowledge to generate synthetic data | Was evaluated using simple DL architectures |
| **Resource scheduling** | [115] | Efficiently serves multiple models on edge devices architectures | Limited concurrent access to Coral TPUs from different processors |
| | [9] | Operates on resource-constrained IoT sensors | Long deployment time for new ML models |

Nevertheless, the current approach does not fully utilize the allocated GPUs for inference tasks. Optimizing resource allocation and implementing dynamic load balancing strategies could improve resource utilization efficiency.

## 4.3 Robust Approaches in DataOps

As discussed in Section 3.2, DataOps involves performing several operations, including data cleaning, ensuring the robustness of the entire ML system to data corruptions such as distributional shifts and imbalance datasets, in addition to robustness to data scarcity. In this section, we review existing solutions that address the robustness aspects of DataOps in ML solutions in production. The advantages and disadvantages of the different approaches are presented in Table 3

*4.3.1 Data Cleaning Approaches.* CPClean approach was proposed to handle datasets with missing values to train ML models [80]. The approach uses sequential information maximization to minimize conditional entropy to find the best cleaning strategy. For missing value imputation, CPClean was tested using both simple imputation methods and ML-based imputation methods. However, CPClean requires manual efforts, which could be a bottleneck in large-scale applications. In a different approach, Schelter *et al.* have developed an automated data quality verification platform [139]. The platform integrates an incremental validation methodology of data quality for evolving datasets. In addition, the platform also supports anomaly detection using statistics from historical data quality metrics, , but real-time capabilities should be enhanced. Similarly, Breck *et al.* designed

an automated data validation system in the context of ML [24]. Single-batch validation is performed to analyze deviations from the expected states of data samples within a batch to detect instability in data characteristics. Recently, the DataOps-4G platform was proposed for both specialists and generalists, but it may face limitations in handling advanced data cleaning tasks due to limited programming knowledge. DataOps-4G enables the discovery of data quality issues without the need to write code [188]. The platform handled five data quality issues including missing, non-unique, duplicate, imprecise, and inconsistent values. The platform can be integrated into ML production environments by bundling functionalities and codes into DataOps pipelines. Another toolkit was built to handle data quality issues, namely DQLearn [144].

*4.3.2 Robust Approaches to Distribution Shift.* In a recent study, DRIFT LENS, a real-time unsupervised drift detection approach has been designed to identify and cope with changes in the system using MLOps pipelines [57]. DRIFT LENS drift processing includes offline and online phases, but may be prone to false positives, which could impact decision-making accuracy. In the offline phase, baseline data is extracted to represent the learned concept during training. Baseline data are used to compare the distribution of production data in the online phase to detect changes. Changes are detected on the basis of a comparison with a drift-magnitude threshold value. For medical imaging AI applications, CheXstray has been proposed to achieve unsupervised real-time drift monitoring [147]. For shift detection, the approach employs statistical tests to compare the distribution of a reference dataset with the detection window. In the detection phase, for numerical features, the Kolmogorov-Smirnov (K-S) test is used to measure the distribution shift. In contrast, the chi-square goodness of fit test is used for categorical variables. For scalable industrial processes, an automated and adaptive approach has been proposed to handle distribution changes [14]. The approach uses a covariate shift adaptation technique, namely importance weighting [154], which adjusts the data distribution of the production data samples to match the data distribution of the training datasets. Additionally, the approach has also been used for the rapid integration of new machinery into industrial systems. Yet, the approach may be limited to specific ML models and requires further testing across diverse domains to assess its robustness. Similarly, the MANDOLINE framework has been proposed to evaluate ML models under distribution shift [33]. MANDOLINE uses user-specified slicing functions to capture the distribution shift. The slicing functions are also used to find importance weighting estimates using the Kullback-Leibler importance estimation procedure (KLIEP) [155]. In a different approach, the MASA sampling algorithm has been developed to detect the distribution shift [32]. MASA partitions the dataset into several clusters to estimate changes in ML models using the confusion matrix that is used to evaluate ML performance. This process is repeated until it reaches the stopping rule or the maximum number of iterations. However, it could miss subtle distribution shifts and is limited to classification tasks, potentially restricting its utility in broader contexts. Another system proposed by Breck *et al.* is capable of tracking and identifying distribution shifts between different batches of data [23].

*4.3.3 Robust Approaches to Data Scarcity.* Overton system has been developed to help engineers automate the ML life cycle by providing tools to maintain the AI application [129]. To improve existing data, Overton provides a feature to augment the data by creating synthetic data. The system also supports cold-start use cases where there are no training data available to train the ML model. However, Overton lacks model versioning capabilities. In a different study, the Snorkel system was developed to create labeled datasets using labeling functions [128]. The system can quickly create training data sets by combining sources of weak supervision. Each data point will be annotated with a probabilistic label that will eventually be used by discriminative ML models. For industrial applications with data scarcity, a data augmentation method has been proposed to generate a new set of augmented data points [31]. The method first uses domain knowledge to

decompose the data and then augments the dataset with random and fresh samples, but requires further validation using complex models to ensure its effectiveness across diverse applications.

*4.3.4  Resource Scheduling Approaches.* SensiX++ system has been developed to incorporate MLOps capabilities with edge devices [115]. The system features several components to manage and orchestrate ML solutions in production. The components provide system-wide orchestration and allow the automation and optimization of the data operations to serve multiple models simultaneously. An adaptive scheduler is also incorporated to manage the system resources and inform the input data requirement of the ML models. However, it may require improvements in concurrent resource access. Another new framework used to orchestrate the MLOps system at the edge of IoT systems is Tiny-MLOps [9]. The framework includes several features to deliver adaptation and evolving capabilities to devices with limited resources. Tiny-MLOps is able to orchestrate the entire life cycle of the MLOps system that executes a specific ML-based task, but faces challenges related to long deployment times for new ML models.

## 4.4  Robust Approaches in ModelOps

ML models in production systems continuously evolve and acquire new knowledge. As explained in Section 3.3, various issues should be considered in ML to retain robustness, including hyperparameter optimization, concept drift, generalizability, and label noise. In this section, we examine the literature and survey the approaches that handle the robustness aspects of ModelOps. The pros and cons of the various approaches are outlined in Table 4 to provide a comprehensive summary.

*4.4.1  Hyperparameter Optimization Approaches.* A Bayesian Optimization AutoML Time-series (BOAT) framework was integrated with the MLOps system to automatically select the best hyperparameter configurations [91], but its applicability is limited to univariate time-series data. The framework employs a Monte Carlo Tree Search-based (MCTS) Bayesian optimizer that can efficiently find the optimal hyperparameter configuration in conditional and high-dimensional spaces. Another popular framework for hyperparameter optimization is Optuna [2]. Optuna has the define-by-run feature that allows users to construct the hyperparameter search space dynamically, yet it lacks support for external optimizers. The framework uses a combination of efficient sampling and pruning algorithms to optimize the optimal hyperparameter search. In a different study, the MAGGY framework was introduced for asynchronous parallel hyperparameter search [113]. The framework incorporates several hyperparameter optimization techniques, Bayesian optimization, and Gaussian processes, in addition to user-defined optimizers. MAGGY achieves asynchrony using an early stopping mechanism. However, the framework was evaluated on a single use case, raising concerns about its generalizability.

*4.4.2  Robust Approaches to Concept Drift.* A novel multi-agent system framework, Driftage, was proposed to handle concept drift in deployment [169]. The framework was built based on Monitor-Analyse-Plan-Execute over a shared Knowledge (MAPE-K) architecture for self-adaptation [73]. For concept drift handling, Driftage follows two flows, Monitor–Analyser and Planner–Executor, for drift detection and analysis. Furthermore, the Analyser and Planner components can be executed using various drift handling techniques. For security applications, the Contrastive Autoencoder for Drifting Detection and Explanation (CADE) method was proposed to detect and understand concept drift occurrence [184]. For drift detection, the method defines the distance function from the training dataset based on contrastive learning. For drift explanation, the method uses distance changes to find the features that drive the significant deviation between the drifting data and its nearest class. For industrial processes, an Online Autoregression with Deep Long-Short-Term Memory (OAR-DLSTM) method was developed to handle hybrid recurring concept drifts. The

Table 4. Summary of robust approaches in ModelOps: Pros and cons

| Operation | Reference | Pros | Cons |
|---|---|---|---|
| **Hyperparameter optimization** | [91] | Customizable Python API with centralized user interface | Supports only univariate time-series data |
| | [2] | Dynamic construction of search space for parameters | No interface to support deploying external optimizers |
| | [113] | Use of parallel compute resources and global optimization | Was only evaluated on a single use-case |
| **Concept drift** | [169] | Flexibility and interpretability | Limited to Python |
| | [184] | Explainable and low computational complexity | Mainly focused on detecting new concepts |
| | [107] | Scalable, adaptable, and flexible for large-scale systems | Cannot handle unseen drifts leading to accuracy drops |
| **Generalization** | [65] | Handles high-dimensional data | Limited to specific tasks |
| | [194] | Efficient in unseen domains | Handles supervised tasks only |
| | [61] | Reduces observational overfitting | Limited to specific tasks |
| **Label noise** | [58] | Robust to train dataset distortions | High number of hyperparameters |
| | [109] | Support real-time systems | Cost and infrastructural limitations |
| | [132] | User-friendly interface | Limited to classification |

method adopts an Online Autoregression (OAR) algorithm to alleviate the challenge of transfer and overlap between concepts. For making predictions, the method defines a sub-model based on the types of recurring concepts. For large-scale systems, Matchmaker was introduced to deal with covariate shift and concept drift problems [107]. The method ranks batches taken from the training dataset based on spatial proximity to the test sample and the accuracy of the trained ML model. After that, the rankings are combined and the model with the best performance is selected for deployment.

*4.4.3 Generalization Approaches.* External data sources have been used by Ho *et al.* to improve the generalizability of ML models [65]. The authors have shown that external validation outperforms internal validation in improving the generalizability of ML models. In addition, considerations for acquiring a proper external validation dataset and deploying the procedures were presented in the article. Finally, the external validation procedure was extended using convergent and divergent validations. For medical imaging, the BigAug method was proposed to generalize 3D medical image segmentation models to unseen domains [194]. The method used data augmentation techniques in the source domain to train neural networks on a large augmented dataset that accommodates several domain-shift situations. However, its effectiveness can be constrained to tasks within the medical imaging domain. Similarly, in reinforcement learning, the SOft Data Augmentation (SODA) method was proposed to improve the generalizability of models using data augmentation and domain randomization [61]. The method incorporates a novel DMControl generalization benchmark that can achieve high performance in test-time generalization to unseen environments.

*4.4.4 Robust Approaches to Label Noise.* An automated dataset optimization (AutoDO) model has been proposed to handle datasets with label noise [58]. AutoDO employs a soft-label sub-model to estimate hyperparameters for each label in the training dataset based on a noise-free validation dataset. These soft labels are then generated using a softmax function. While this method effectively deals with label noise, it may require tuning a large number of hyperparameters, which can be labor-intensive and time-consuming. Furthermore, the effectiveness of AutoDO may be influenced by the quality and representativeness of the noise-free validation dataset. On the other hand, DeFraudNet,

as outlined in [109], utilizes auto-encoder reconstruction error to reduce label noise. By taking both an unlabeled training dataset and a golden label generation dataset as input, DeFraudNet outputs a denoised, weak-labeled training dataset. This approach integrates label generation with model deployment, offering a comprehensive solution to address label noise. However, it may face challenges related to cost and infrastructure limitations, particularly in scenarios where extensive computational resources are required for training and deploying the model. In a more recent study, the Snoopy system was presented to facilitate the execution of automatic feasibility studies before the deployment of ML solutions [132]. The study has focused on label noise situations, which are sources for the non-zero irreducible error, but its applicability is limited to classification tasks. The core component of Snoopy is the Bayes error rate (BER), which estimates the error of a given task.

## 5 REVIEW OF MLOPS TOOLS AND TECHNOLOGIES

Currently, several tools and platforms are available to operationalize ML solutions in production. We will explore the most popular MLOps tools on the basis of their official documentation. Specifically, we review the support of each tool for the various robustness aspects that were surveyed in the previous sections. This can help practitioners in selecting the most useful tool(s) for their applications and in deploying a robust solution. From a technological perspective, the most popular public tools that manage the ML life cycle in deployments are [159]: Amazon SageMaker by Amazon Web Service (AWS) [95], Microsoft Azure [146], Vertex AI (formerly known as AI Platform) [20], MLFlow [189] and TensorFlow Extended (TFX) [13]. The platforms can be categorized into cloud-based platforms: AWS, Azure, and Google Cloud, and non-cloud-based platforms: MLFlow and TFX. Cloud-based platforms are known to be more elastic and easier to use with their interactive interface. However, they are not as customizable as the non-cloud platforms. We provide a brief description of each of these MLOps platforms:

(1) **Amazon SageMaker:** Amazon's interface to productize ML solutions in AWS. The platform offers a free trial to start using its services. The platform provides a wide range of functionalities to manage ML solutions in the cloud through interactive interfaces. Moreover, SageMaker features many built-in algorithms and utilities, in addition to the ability to customize them by users. SageMaker also includes the ability to estimate training costs before running the jobs. To store ML artifacts, SageMaker uses Amazon Simple Storage Service (Amazon S3) buckets.

(2) **Microsoft Azure:** Microsoft's cloud-based platform that supports MLOps pipelines through a set of services. The platform includes an Azure Machine Learning Studio environment to operationalize ML projects. Azure also includes a free trial option to get acquainted with its services. One of the main services provided by Microsoft Azure is Azure DevOps, which manages and automates ML project pipelines. Additionally, Azure SDKs include a diverse range of libraries that support many programming languages to build ML projects. In addition to the built-in packages, Azure also supports ML solutions customization.

(3) **Vertex AI:** Google's cloud-based framework to orchestrate the MLOps pipeline in Google Cloud. The AI platform provides several services for building, evaluating and deploying ML solutions. To start with the services, free credit is offered that can be used to test the platform. Many functionalities and components are integrated into the platform. Dataflow is the component that is concerned with managing the data life cycle and supports batch and stream data processing. Another key component is Cloud Build, which enables smooth automation of ML application operations with a serverless CI/CD platform.

(4) **MLFlow:** An open source platform built based on *open interface* philosophy. MLFlow allows the management of the entire ML life cycle through four primary components: MLflow Tracking, MLflow Projects, MLflow Models, and MLflow Models. MLFlow can deploy ML

Table 5. Summary of the built-in features of MLOps tools

|            | Operation | MLOps tool | | | | |
|------------|-----------|------|-------|-----------|--------|-----|
|            |           | AWS  | Azure | Vertex AI | MLFlow | TFX |
| **Automation** | **Continuous validation** | ✓ | ✓ | ✓ | ✓ | ✓ |
|            | **Continuous versioning** | ✓ | ✓ | ✓ | ✓ | ✓ |
|            | **Continuous monitoring** | ✓ | ✓ | ✓ | ✓ | ✓ |
|            | **Continuous update** | ✓ | ✓ | ✓ | ✓ | ✓ |
| **DataOps** | **Data cleaning** | ✓ | ✓ | ✓ | ✓ | ✓ |
|            | **Anomaly detection** | ✓ | ✓ | ✓ | - | ✓ |
|            | **Distribution shift** | ✓ | ✓ | ✓ | - | ✓ |
|            | **Data augmentation** | - | - | - | - | - |
| **ModelOps** | **Hyperparameter tuning** | ✓ | ✓ | ✓ | ✓ | ✓ |
|            | **Concept drift** | ✓ | - | - | - | - |
|            | **Generalizability** | - | ✓ | - | - | - |
|            | **Label noise** | ✓ | - | ✓ | - | - |

projects and libraries in any programming language. It can also be integrated as a web service with cloud platforms to deploy ML pipelines. MLFlow offers a useful feature in practice that allows users to track performance using interactive plots using a custom-defined metric.

(5) **TFX:** An open source general-purpose platform for production-scale deployments of ML pipelines. TFX was developed on the basis of the TensorFlow [4] software library for ML and AI. TFX supports the most popular programming languages. To implement an ML pipeline, TFX allows users to configure a sequence of components that can handle data processing and model deployment. For quality monitoring, TFX uses validation checks for every stage of the ML pipeline. These validations can evaluate the data quality and performance of the ML models.

To gain insight into the support for the robustness aspect of MLOpsm tools, we compared the built-in features offered by each tool. The comparison results are summarized in Table 5. We can see that all tools support automation management of the ML life cycle. These continuous operations are the cornerstone of any ML project deployed in real-world applications, and they distinguish them from offline ML projects. For DataOps, all platforms provide data management tools, such as data analysis or visualization, to facilitate data quality tracking. On the contrary, other operations are still missing from the tools, such as data augmentation for DataOps. For the rest of the operations, most of the tools support distribution shift detection by employing statistical tests to monitor the data distribution. For ModelOps, all tools incorporate hyperparameter tuning methods to tweak the ML models. Detecting concept drift is only included in AWS SageMake with the CloudWatch metrics repository. CloudWatch monitors the deviation in model performance by monitoring the defined performance metric and activates the model re-training procedure automatically. RobustDG toolkit [104] was developed by Microsoft to enhance the generalizability of the ML model to unseen domains. The toolkit can be easily extended to write problem-specific domain generalization algorithms. Label noise is handled by the data labeling functionality offered by AWS and Vertext AI. Lastly, an important remark is that we have surveyed the built-in features that are integrated by the tools. However, the majority of platforms exhibit flexibility by allowing developers to write their own solutions and deploy them in the framework.

---

[4]https://www.tensorflow.org/

## 6   INSIGHTS AND OPEN CHALLENGES

In this paper, we highlight the lessons learned from the surveyed methods, open challenges and potential future directions to build robust MLOps systems. Although the existing literature provides information on key practices for improving the technical robustness of MLOps approaches, significant gaps still exist, necessitating further research and exploration. Through an extensive literature search, we identified the key insights and several challenges in achieving robustness in MLOps systems and outline opportunities for enhancing the maturity of the field. Our aim is to provide a roadmap to address these challenges and accelerate the widespread adoption of robust machine learning systems in AI applications.

### 6.1   Lessons Learned from Robust MLOps Methods

Several lessons and insights can be inferred from the literature based on the surveyed methods.

- **Monitoring and Performance Optimization:** Implementing advanced monitoring practices in MLOps is crucial for maintaining the reliability and performance of ML systems in production environments. Beyond merely tracking prediction accuracy, real-time monitoring provides deep insights into model behavior and system health, enabling proactive identification of potential issues before they escalate into significant problems.
- **Scalable Infrastructure and Resource Management:** The scalability of infrastructure is essential for accommodating varying workloads in MLOps. Leveraging cloud-native solutions such as AWS, Azure, or GCP provides scalable infrastructure and efficient resource management. This approach supports dynamic workloads and large-scale deployments, enhancing operational flexibility and cost optimization.
- **Customization:** Customization plays a pivotal role in maximizing the effectiveness of MLOps tools and practices. While off-the-shelf MLOps tools offer a broad range of functionalities, customizing these tools to fit specific project requirements can significantly improve outcomes. Integrating tailored solutions and leveraging developer expertise enables organizations to address unique challenges effectively.

### 6.2   Challenges in Implementing Robustness

Implementing robustness in MLOps systems presents multifaceted challenges that require strategic approaches and careful consideration across various dimensions. These challenges span from quantifying overall system robustness to managing critical aspects such as data quality, model interpretability, concept drift, and continuous adaptation.

- **Quantifying System Robustness:** A key challenge in MLOps is quantifying system robustness effectively. This involves developing rigorous methodologies to assess resilience across various dimensions simultaneously, such as handling data perturbations, ensuring model performance in diverse conditions, and maintaining system reliability in dynamic operational environments.
- **Integration and Orchestration:** Deploying and managing ML models at scale within production environments presents technical challenges related to integration and versioning integrity. Ensuring seamless integration with existing infrastructure, efficient resource management, and maintaining version control across deployments are critical concerns. Adopting containerization technologies, orchestration frameworks such as Kubernetes[5], and advanced DevOps practices can streamline deployment processes and enhance scalability, reliability, and maintainability of MLOps systems.

---

[5]https://kubernetes.io/

- **Cost Management:** Optimizing operational costs associated with running ML workloads, such as infrastructure and resource utilization, is critical. Achieving a balance between performance and cost efficiency necessitates advanced resource allocation and management strategies.
- **Human-in-the-Loop:** Effectively integrating human feedback and intervention into automated MLOps pipelines remains a significant challenge. Incorporating dynamic domain expertise and user feedback to enhance model performance and adaptability is an active area of research [116].
- **Heterogeneous Data Sources and Generalization:** Distributed systems, such as federated learning, involve training models across multiple devices or institutions with varying datasets that may differ significantly in distribution, quality, and size [193]. Ensuring that the federated model can generalize well across these heterogeneous data sources without overfitting to any particular client's data remains a substantial challenge.

## 6.3 Future Opportunities

Based on the open challenges identified, we suggest the following directions to explore in the future.

*6.3.1 Trade-offs in Robustness Aspects of MLOps.* The underlying effects of the various robustness aspects are not investigated for ML systems in production. Due to its high impact on the system design, the trade-off between the robustness aspects of MLOps subtasks must be closely examined. This trade-off determines the balance between the various aspects since addressing a particular aspect can improve/deteriorate the system robustness to other aspects. Many studies have reported side effects between robustness aspects [85, 98]. Improving one aspect of robustness may come at the expense of another, highlighting the need for careful consideration and strategic trade-offs. For example, optimizing model accuracy through complex ensemble techniques might sacrifice interpretability, making it challenging to understand the model's decisions. Therefore, it would be valuable to explore how this interplay can be managed in real-world systems. This would help mitigate the potentially reflecting effects of aspects and emphasize specific issues that appear more frequently in certain problems.

*6.3.2 Quantification of Overall Robustness.* The current evaluation methods primarily focus on quantitative performance metrics such as accuracy and precision. However, robustness encompasses broader aspects such as resilience to adversarial attacks, model interpretability, and generalization across diverse datasets. For example, robustness evaluations may involve stress testing models with adversarial examples, analyzing model behavior under different environmental conditions, and assessing the impact of data distribution shifts on model performance. In general, the issue arises in relation to the interpretation of the results. For example, in the literature, concept drift is evaluated mainly using performance degradation [15]. Nevertheless, other factors could cause performance to degrade, and it does not necessarily indicate a concept drift occurrence. Therefore, developing comprehensive and aspect-specific methodologies for evaluating robustness can provide deeper insights into system performance and identify potential vulnerabilities. Additionally, the use of explainable AI techniques can improve our understanding of system behavior and observations, facilitating the construction of quantifiable solutions [16].

*6.3.3 Holistic Robust MLOps Approach.* Building a holistic, robust ML system for production requires comprehensive frameworks that address all aspects of robustness simultaneously. While existing approaches in the literature partially handle robustness, there is a lack of clarity on how these approaches collectively address the multitude of factors influencing robustness. To bridge this

gap, it is imperative to adopt a unified approach that considers the interplay between various factors that influence robustness. The plethora of MLOps tools and platforms available today presents a unique opportunity to construct a holistic, robust MLOps approach. Using the customizable features and integrative capabilities of these tools, a holistic robust MLOps system can be achieved ultimately in real-world production environments.

*6.3.4    Flexible MLOps Architecture.* The dynamicity of real-world applications requires adaptation in the architectural design of solution systems [120]. Traditional monolithic architectures are often rigid and difficult to modify, hindering agility and innovation in MLOps systems. An interesting line of research is to develop modular and extensible architectures that enable flexible integration of new data sources, algorithms, and computational resources. One promising approach is the adoption of microservices-based architectures, which facilitates the development of scalable and decoupled MLOps systems, allowing practitioners to iterate and experiment with different components independently. Additionally, the implementation of containerization and orchestration technologies can enhance scalability and resource utilization, enabling the deployment and management of ML workloads more efficiently. Therefore, future research should focus on designing extensible architectures that promote the development of elastic solutions. These solutions should possess the ability to adapt and evolve in response to changing business requirements and technological advancements.

## 7    CONCLUSION

In this survey, we start with the EU's guidelines for the requirements of AI systems to be trustworthy. One of the key principles suggested in the guidelines is that the system should be technically robust. However, the robustness property may have different definitions depending on the domain and the application. However the common point of all definitions is that a robust system can withstand unforeseen conditions during its evolution. AI solutions and its most widespread sub-field of ML are employed in a vast spectrum of applications. In real-world scenarios, the deployment of ML solutions in production requires a set of operations to be performed. MLOps approach has emerged to standardize ML workflows in deployment settings. Therefore, the robustness property should be addressed in MLOps towards building trustworthy ML in production. Through a thorough literature search, we draw out the current challenges and devise the specifications and practices for robust-driven ML solutions in production. By embracing these practices to deliver MLOps pipelines, the primary attention is devoted to building robust ML systems in the long run rather than focusing on short-term gains. In the literature, a key observation is that the robustness property has been used loosely to refer to an approach that partially addresses robustness aspects. Therefore, to build an overall robust ML framework in production, robustness aspects should be fully covered in the constituent ingredients of MLOps, specifically, in the three interdependent components of MLOps: automation management, which is the pillar of MLOps, DataOps, and ModelOps. For each component, we provide related definitions and concepts. Additionally, we organize existing academic research that handles the robustness aspects of production ML according to MLOps components. We also review the most prevalent MLOps tools and their default support of robustness aspects to enable researchers and practitioners to select the suitable tool for their problem.

# REFERENCES

[1] Supriya Agrahari and Anil Kumar Singh. 2021. Concept drift detection in data stream mining: A literature review. *Journal of King Saud University-Computer and Information Sciences* (2021).

[2] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 2623–2631.

[3] Ismaeel Al Ridhawi, Safa Otoum, Moayad Aloqaily, and Azzedine Boukerche. 2020. Generalizing AI: challenges and opportunities for plug and play AI solutions. *IEEE Network* 35, 1 (2020), 372–379.

[4] Hazrat Ali, Christer Grönlund, and Zubair Shah. 2023. Leveraging GANs for data scarcity of COVID-19: Beyond the hype. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 659–667.

[5] Sridhar Alla and Suman Kalyan Adari. 2021. *Beginning MLOps with MLFlow*. Springer.

[6] Sridhar Alla and Suman Kalyan Adari. 2021. What is mlops? In *Beginning MLOps with MLFlow*. Springer, 79–124.

[7] Miltiadis Allamanis. 2019. The adverse effects of code duplication in machine learning models of code. In *Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*. 143–153.

[8] Davide Anguita, Luca Ghelardoni, Alessandro Ghio, Luca Oneto, and Sandro Ridella. 2012. The 'K'in K-fold cross validation. In *20th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning (ESANN)*. i6doc. com publ, 441–446.

[9] Mattia Antonini, Miguel Pincheira, Massimo Vecchio, and Fabio Antonelli. 2022. Tiny-MLOps: a framework for orchestrating ML applications at the far edge of IoT systems. In *2022 IEEE International Conference on Evolving and Adaptive Intelligent Systems (EAIS)*. IEEE, 1–8.

[10] Harvinder Atwal. 2020. Dataops technology. In *Practical DataOps*. Springer, 215–247.

[11] Rohit Babbar and Bernhard Schölkopf. 2019. Data scarcity, robustness and extreme multi-label classification. *Machine Learning* 108, 8 (2019), 1329–1351.

[12] Ms Aayushi Bansal, Dr Rewa Sharma, and Dr Mamta Kathuria. 2022. A systematic review on data scarcity problem in deep learning: solution and applications. *ACM Computing Surveys (CSUR)* 54, 10s (2022), 1–29.

[13] Denis Baylor, Eric Breck, Heng-Tze Cheng, Noah Fiedel, Chuan Yu Foo, Zakaria Haque, Salem Haykal, Mustafa Ispir, Vihan Jain, Levent Koc, et al. 2017. TFX: A tensorflow-based production-scale machine learning platform. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. 1387–1395.

[14] Firas Bayram, Bestoun S Ahmed, Erik Hallin, and Anton Engman. 2022. A Drift Handling Approach for Self-Adaptive ML Software in Scalable Industrial Processes. In *2022 37th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE.

[15] Firas Bayram, Bestoun S. Ahmed, and Andreas Kassler. 2022. From concept drift to model degradation: An overview on performance-aware drift detectors. *Knowledge-Based Systems* 245 (2022), 108632.

[16] Vaishak Belle and Ioannis Papantonis. 2021. Principles and practice of explainable machine learning. *Frontiers in big Data* (2021), 39.

[17] James Bergstra and Yoshua Bengio. 2012. Random search for hyper-parameter optimization. *Journal of machine learning research* 13, 2 (2012).

[18] Daniel Berrar. 2019. Cross-Validation.

[19] Arjun Nitin Bhagoji, Daniel Cullina, Chawin Sitawarin, and Prateek Mittal. 2018. Enhancing robustness of machine learning systems via data transformations. In *2018 52nd Annual Conference on Information Sciences and Systems (CISS)*. IEEE, 1–5.

[20] Ekaba Bisong. 2019. An overview of google cloud platform services. *Building Machine Learning and Deep Learning Models on Google Cloud Platform* (2019), 7–10.

[21] Alexei Botchkarev. 2018. Performance metrics (error measures) in machine learning regression, forecasting and prognostics: Properties and typology. *arXiv preprint arXiv:1809.03006* (2018).

[22] James V Bradley. 1978. Robustness? *Brit. J. Math. Statist. Psych.* 31, 2 (1978), 144–152.

[23] Eric Breck, Shanqing Cai, Eric Nielsen, Michael Salib, and D Sculley. 2017. The ML test score: A rubric for ML production readiness and technical debt reduction. In *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 1123–1132.

[24] Eric Breck, Neoklis Polyzotis, Sudip Roy, Steven Whang, and Martin Zinkevich. 2019. Data Validation for Machine Learning. In *MLSys*.

[25] Richard H Byrd, Jorge Nocedal, and Robert B Schnabel. 1994. Representations of quasi-Newton matrices and their use in limited memory methods. *Mathematical Programming* 63, 1 (1994), 129–156.

[26] Antonio Capizzi, Salvatore Distefano, and Manuel Mazzara. 2019. From devops to devdataops: Data management in devops processes. In *International Workshop on Software Engineering Aspects of Continuous Development and New Paradigms of Software Production and Deployment*. Springer, 52–62.

[27] Nicholas Carlini and David Wagner. 2017. Towards evaluating the robustness of neural networks. In *2017 ieee symposium on security and privacy (sp)*. Ieee, 39–57.

[28] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 3 (2009), 1–58.

[29] Guoqing Chao, Yuan Luo, and Weiping Ding. 2019. Recent advances in supervised dimension reduction: A survey. *Machine learning and knowledge extraction* 1, 1 (2019), 341–358.

[30] Raja Chatila, Virginia Dignum, Michael Fisher, Fosca Giannotti, Katharina Morik, Stuart Russell, and Karen Yeung. 2021. Trustworthy ai. In *Reflections on Artificial Intelligence for Humanity*. Springer, 13–39.

[31] Ayan Chatterjee, Bestoun S Ahmed, Erik Hallin, and Anton Engman. 2022. Testing of machine learning models with limited samples: an industrial vacuum pumping application. In *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*. 1280–1290.

[32] Lingjiao Chen, Tracy Cai, Matei Zaharia, and James Zou. 2021. Did the model change? Efficiently assessing machine learning API shifts. *arXiv preprint arXiv:2107.14203* (2021).

[33] Mayee Chen, Karan Goel, Nimit S Sohoni, Fait Poms, Kayvon Fatahalian, and Christopher Ré. 2021. Mandoline: Model Evaluation under Distribution Shift. In *International Conference on Machine Learning*. PMLR, 1617–1629.

[34] Mandepudi Nobel Chowdary, Bussa Sankeerth, Chennupati Kumar Chowdary, and Manu Gupta. 2022. Accelerating the Machine Learning Model Deployment using MLOps. In *Journal of Physics: Conference Series*, Vol. 2327. IOP Publishing, 012027.

[35] Abdur Chowdhury and Joshua Alspector. 2003. Data duplication: an imbalance problem?. In *ICML'2003 Workshop on Learning from Imbalanced Data Sets (II), Washington, DC*. Citeseer.

[36] Xu Chu, Ihab F Ilyas, Sanjay Krishnan, and Jiannan Wang. 2016. Data cleaning: Overview and emerging challenges. In *Proceedings of the 2016 international conference on management of data*. 2201–2206.

[37] Jacob Cohen. 1968. Weighted kappa: nominal scale agreement provision for scaled disagreement or partial credit. *Psychological bulletin* 70, 4 (1968), 213.

[38] European Commission, Content Directorate-General for Communications Networks, and Technology. 2019. *Ethics guidelines for trustworthy AI.* Publications Office.

[39] Daniel Crankshaw, Peter Bailis, Joseph E Gonzalez, Haoyuan Li, Zhao Zhang, Michael J Franklin, Ali Ghodsi, and Michael I Jordan. 2014. The missing piece in complex analytics: Low latency, scalable model management and serving with velox. *arXiv preprint arXiv:1409.3809* (2014).

[40] Ilias Diakonikolas, Gautam Kamath, Daniel M Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. 2021. Robustness meets algorithms. *Commun. ACM* 64, 5 (2021), 107–115.

[41] Ousmane Diallo, Joel JPC Rodrigues, and Mbaye Sene. 2012. Real-time data management on wireless sensor networks: A survey. *Journal of Network and Computer Applications* 35, 3 (2012), 1013–1021.

[42] Alican Dogan and Derya Birant. 2021. Machine learning and data mining in manufacturing. *Expert Systems with Applications* 166 (2021), 114060.

[43] Pedro Domingos. 2012. A few useful things to know about machine learning. *Commun. ACM* 55, 10 (2012), 78–87.

[44] Anton Dries and Ulrich Rückert. 2009. Adaptive concept drift detection. *Statistical Analysis and Data Mining* 2, 5-6 (2009), 311–327.

[45] Mengnan Du, Ninghao Liu, and Xia Hu. 2019. Techniques for interpretable machine learning. *Commun. ACM* 63, 1 (2019), 68–77.

[46] Xianping Du, Hongyi Xu, and Feng Zhu. 2021. Understanding the effect of hyperparameter optimization on machine learning models for structure design problems. *Computer-Aided Design* 135 (2021), 103013.

[47] Michaela Dvorzak and Helga Wagner. 2016. Sparse Bayesian modelling of underreported count data. *Statistical Modelling* 16, 1 (2016), 24–46.

[48] Lisa Ehrlinger, Thomas Grubinger, Bence Varga, Mario Pichler, Thomas Natschläger, and Jürgen Zeindl. 2018. Treating missing data in industrial data analytics. In *2018 Thirteenth International Conference on Digital Information Management (ICDIM)*. IEEE, 148–155.

[49] Ryan Elwell and Robi Polikar. 2011. Incremental learning of concept drift in nonstationary environments. *IEEE Transactions on Neural Networks* 22, 10 (2011), 1517–1531.

[50] Newsha Emaminejad and Reza Akhavian. 2022. Trustworthy AI and robotics: Implications for the AEC industry. *Automation in Construction* 139 (2022), 104298.

[51] Bradley J Erickson and Felipe Kitamura. 2021. Magician's corner: 9. Performance metrics for machine learning models. *Radiology: Artificial Intelligence* 3, 3 (2021).

[52] Matthias Feurer and Frank Hutter. 2019. Hyperparameter optimization. In *Automated machine learning*. Springer, Cham, 3–33.

[53] Benoît Frénay and Michel Verleysen. 2013. Classification in the presence of label noise: a survey. *IEEE transactions on neural networks and learning systems* 25, 5 (2013), 845–869.

[54] Luís PF Garcia, André CPLF de Carvalho, and Ana C Lorena. 2015. Effect of label noise in the complexity of classification problems. *Neurocomputing* 160 (2015), 108–119.

[55] Satvik Garg, Pradyumn Pundir, Geetanjali Rathee, PK Gupta, Somya Garg, and Saransh Ahlawat. 2021. On continuous integration/continuous delivery for automated deployment of machine learning models using mlops. In *2021 IEEE fourth international conference on artificial intelligence and knowledge engineering (AIKE)*. IEEE, 25–28.

[56] Jian Ge, Han Li, Hongpeng Wang, Haobin Dong, Huan Liu, Wenjie Wang, Zhiwen Yuan, Jun Zhu, and Haiyang Zhang. 2019. Aeromagnetic compensation algorithm robust to outliers of magnetic sensor based on Huber loss method. *IEEE Sensors Journal* 19, 14 (2019), 5499–5505.

[57] Salvatore Greco and Tania Cerquitelli. 2021. Drift Lens: Real-time unsupervised Concept Drift detection by evaluating per-label embedding distributions. In *2021 International Conference on Data Mining Workshops (ICDMW)*. IEEE, 341–349.

[58] Denis Gudovskiy, Luca Rigazio, Shun Ishizaka, Kazuki Kozuka, and Sotaro Tsukizawa. 2021. Autodo: Robust autoaugment for biased data with label noise via scalable probabilistic implicit differentiation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16601–16610.

[59] Xinjian Guo, Yilong Yin, Cailing Dong, Gongping Yang, and Guangtong Zhou. 2008. On the class imbalance problem. In *2008 Fourth international conference on natural computation*, Vol. 4. IEEE, 192–201.

[60] Ronan Hamon, Henrik Junklewitz, and Ignacio Sanchez. 2020. Robustness and explainability of artificial intelligence. *Publications Office of the European Union* (2020).

[61] Nicklas Hansen and Xiaolong Wang. 2021. Generalization in reinforcement learning by soft data augmentation. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 13611–13617.

[62] Lars Hertel, Julian Collado, Peter Sadowski, Jordan Ott, and Pierre Baldi. 2020. Sherpa: Robust hyperparameter optimization for machine learning. *SoftwareX* 12 (2020), 100591.

[63] Nipuni Hewage and Dulani Meedeniya. 2022. Machine learning operations: A survey on MLOps tool support. *arXiv preprint arXiv:2202.10169* (2022).

[64] Fabian Hinder, André Artelt, and Barbara Hammer. 2020. Towards non-parametric drift detection via dynamic adapting window independence drift detection (DAWIDD). In *International Conference on Machine Learning*. PMLR, 4249–4259.

[65] Sung Yang Ho, Kimberly Phua, Limsoon Wong, and Wilson Wen Bin Goh. 2020. Extensions of the external validation for checking learned model interpretability and generalizability. *Patterns* 1, 8 (2020), 100129.

[66] T. R. Hoens, R. Polikar, and N. Chawla. 2011. Learning from streaming data with concept drift and imbalance: an overview. *Progress in Artificial Intelligence* 1 (2011), 89–101.

[67] Fred Hohman, Kanit Wongsuphasawat, Mary Beth Kery, and Kayur Patel. 2020. Understanding and visualizing data iteration in machine learning. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–13.

[68] Hossein Hosseini, Baicen Xiao, and Radha Poovendran. 2017. Google's cloud vision api is not robust to noise. In *2017 16th IEEE international conference on machine learning and applications (ICMLA)*. IEEE, 101–105.

[69] Yizheng Huang, Huaizheng Zhang, Yonggang Wen, Peng Sun, and Nguyen Binh Duong Ta. 2021. Modelci-e: Enabling continual learning in deep learning serving systems. *arXiv preprint arXiv:2106.03122* (2021).

[70] Jez Humble and David Farley. 2010. *Continuous delivery: reliable software releases through build, test, and deployment automation*. Pearson Education.

[71] Waldemar Hummer, Vinod Muthusamy, Thomas Rausch, Parijat Dube, Kaoutar El Maghraoui, Anupama Murthi, and Punleuk Oum. 2019. Modelops: Cloud-based lifecycle management for reliable and trusted ai. In *2019 IEEE International Conference on Cloud Engineering (IC2E)*. IEEE, 113–120.

[72] Olumuyiwa Ibidunmoye, Francisco Hernández-Rodriguez, and Erik Elmroth. 2015. Performance anomaly detection and bottleneck identification. *ACM Computing Surveys (CSUR)* 48, 1 (2015), 1–35.

[73] Didac Gil De La Iglesia and Danny Weyns. 2015. MAPE-K formal templates to rigorously design behaviors for self-adaptive systems. *ACM Transactions on Autonomous and Adaptive Systems (TAAS)* 10, 3 (2015), 1–31.

[74] Ihab F Ilyas and Xu Chu. 2019. *Data cleaning*. Morgan & Claypool.

[75] ISO/IEC. 2022. *Information technology — Artificial intelligence — Overview of trustworthiness in artificial intelligence*. Technical Report 24028:2020. Geneva, Switzerland.

[76] Anna Jobin, Marcello Ienca, and Effy Vayena. 2019. The global landscape of AI ethics guidelines. *Nature Machine Intelligence* 1, 9 (2019), 389–399.

[77] Meenu Mary John, Helena Holmström Olsson, and Jan Bosch. 2021. Towards mlops: A framework and maturity model. In *2021 47th Euromicro Conference on Software Engineering and Advanced Applications (SEAA)*. IEEE, 1–8.

[78] Michael I Jordan and Tom M Mitchell. 2015. Machine learning: Trends, perspectives, and prospects. *Science* 349, 6245 (2015), 255–260.

[79] Ioannis Karamitsos, Saeed Albarhami, and Charalampos Apostolopoulos. 2020. Applying DevOps practices of continuous automation for machine learning. *Information* 11, 7 (2020), 363.

[80] Bojan Karlaš, Peng Li, Renzhi Wu, Nezihe Merve Gürel, Xu Chu, Wentao Wu, and Ce Zhang. 2020. Nearest Neighbor Classifiers over Incomplete Information: From Certain Answers to Certain Predictions. *Proc. VLDB Endow.* 14, 3 (nov 2020), 255–267.

[81] Gour Karmakar, Abdullahi Chowdhury, Rajkumar Das, Joarder Kamruzzaman, and Syed Islam. 2021. Assessing trust level of a driverless car using deep learning. *IEEE Transactions on Intelligent Transportation Systems* 22, 7 (2021), 4457–4466.

[82] Davinder Kaur, Suleyman Uslu, and Arjan Durresi. 2020. Requirements for trustworthy artificial intelligence–a review. In *International Conference on Network-Based Information Systems.* Springer, 105–115.

[83] Davinder Kaur, Suleyman Uslu, Kaley J Rittichier, and Arjan Durresi. 2022. Trustworthy artificial intelligence: a review. *ACM Computing Surveys (CSUR)* 55, 2 (2022), 1–38.

[84] Yoshinobu Kawahara and Masashi Sugiyama. 2012. Sequential change-point detection based on direct density-ratio estimation. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 5, 2 (2012), 114–127.

[85] Nitish Shirish Keskar and Richard Socher. 2017. Improving generalization performance by switching from adam to sgd. *arXiv preprint arXiv:1712.07628* (2017).

[86] Hiroaki Kitano. 2007. Towards a theory of biological robustness. , 137 pages.

[87] Janis Klaise, Arnaud Van Looveren, Clive Cox, Giovanni Vacanti, and Alexandru Coca. 2020. Monitoring and explainability of models in production. *arXiv preprint arXiv:2007.06299* (2020).

[88] Michał Koziarski, Michał Woźniak, and Bartosz Krawczyk. 2020. Combined cleaning and resampling algorithm for multi-class imbalanced data with label noise. *Knowledge-Based Systems* 204 (2020), 106223.

[89] Dominik Kreuzberger, Niklas Kühl, and Sebastian Hirschl. 2023. Machine learning operations (mlops): Overview, definition, and architecture. *IEEE access* (2023).

[90] Abhishek Kumar, Tristan Braud, Sasu Tarkoma, and Pan Hui. 2020. Trustworthy AI in the age of pervasive computing and big data. In *2020 IEEE International Conference on Pervasive Computing and Communications Workshops (PerCom Workshops).* IEEE, 1–6.

[91] John Joy Kurian, Marcel Dix, Ido Amihai, Glenn Ceusters, and Ajinkya Prabhune. 2021. BOAT: A bayesian optimization automl time-series framework for industrial applications. In *2021 IEEE Seventh International Conference on Big Data Computing Service and Applications (BigDataService).* IEEE, 17–24.

[92] Robyn Larracy, Angkoon Phinyomark, and Erik Scheme. 2021. Machine learning model validation for early stage studies with small sample sizes. In *2021 43rd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC).* IEEE, 2314–2319.

[93] Bo Li, Peng Qi, Bo Liu, Shuai Di, Jingen Liu, Jiquan Pei, Jinfeng Yi, and Bowen Zhou. 2022. Trustworthy AI: From Principles to Practices. *ACM Comput. Surv.* (aug 2022). https://doi.org/10.1145/3555803

[94] KeDi Li and Ning Gui. 2020. CMS: A continuous machine-learning and serving platform for industrial big data. *Future Internet* 12, 6 (2020), 102.

[95] Edo Liberty, Zohar Karnin, Bing Xiang, Laurence Rouesnel, Baris Coskun, Ramesh Nallapati, Julio Delgado, Amir Sadoughi, Yury Astashonok, Piali Das, et al. 2020. Elastic machine learning algorithms in amazon sagemaker. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data.* 731–737.

[96] Haochen Liu, Yiqi Wang, Wenqi Fan, Xiaorui Liu, Yaxin Li, Shaili Jain, Yunhao Liu, Anil K. Jain, and Jiliang Tang. 2022. Trustworthy AI: A Computational Perspective. *ACM Trans. Intell. Syst. Technol.* (jun 2022).

[97] Jie Lu, Anjin Liu, Fan Dong, Feng Gu, João Gama, and Guangquan Zhang. 2019. Learning under Concept Drift: A Review. *IEEE Transactions on Knowledge and Data Engineering* 31, 12 (2019), 2346–2363.

[98] Yang Lu, Yiu-ming Cheung, and Yuan Yan Tang. 2017. Dynamic Weighted Majority for Incremental Learning of Imbalanced Data Streams with Concept Drift.. In *IJCAI.* 2393–2399.

[99] Gang Luo. 2016. A review of automatic selection methods for machine learning algorithms and hyper-parameter values. *Network Modeling Analysis in Health Informatics and Bioinformatics* 5, 1 (2016), 1–16.

[100] Yuan Luo, Peter Szolovits, Anand S Dighe, and Jason M Baron. 2018. 3D-MICE: integration of cross-sectional and longitudinal imputation for multi-analyte longitudinal clinical data. *Journal of the American Medical Informatics Association* 25, 6 (2018), 645–653.

[101] Richard M Lusby, Jesper Larsen, and Simon Bull. 2018. A survey on robustness in railway planning. *European Journal of Operational Research* 266, 1 (2018), 1–15.

[102] Yuting Lyu, Junghui Chen, and Zhihuan Song. 2021. Synthesizing labeled data to enhance soft sensor performance in data-scarce regions. *Control Engineering Practice* 115 (2021), 104903.

[103] J Ma, L Schneider, S Lapuschkin, R Achtibat, M Duchrau, J Krois, F Schwendicke, and W Samek. 2022. Towards Trustworthy AI in Dentistry. *Journal of Dental Research* (2022), 00220345221106086.

[104] Divyat Mahajan, Shruti Tople, and Amit Sharma. 2021. The Connection between Out-of-Distribution Generalization and Privacy of ML Models. *arXiv preprint arXiv:2110.03369* (2021).

[105] Mohammad Saeid Mahdavinejad, Mohammadreza Rezvan, Mohammadamin Barekatain, Peyman Adibi, Payam Barnaghi, and Amit P Sheth. 2018. Machine learning for Internet of Things data analysis: A survey. *Digital Communications and Networks* 4, 3 (2018), 161–175.

[106] Sasu Mäkinen, Henrik Skogström, Eero Laaksonen, and Tommi Mikkonen. 2021. Who needs MLOps: What data scientists seek to accomplish and how can MLOps help?. In *2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN)*. IEEE, 109–112.

[107] Ankur Mallick, Kevin Hsieh, Behnaz Arzani, and Gauri Joshi. 2022. Matchmaker: Data Drift Mitigation in Machine Learning for Large-Scale Systems. *Proceedings of Machine Learning and Systems* 4 (2022), 77–94.

[108] Adnan Masood and Adnan Hashmi. 2019. AIOps: predictive analytics & machine learning in operations. In *Cognitive Computing Recipes*. Springer, 359–382.

[109] Jose Mathew, Meghana Negi, Rutvik Vijjali, and Jairaj Sathyanarayana. 2021. DeFraudNet: An End-to-End Weak Supervision Framework to Detect Fraud in Online Food Delivery. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 85–99.

[110] Beatriz MA Matsui and Denise H Goya. 2022. MLOps: A Guide to its Adoption in the Context of Responsible AI. In *2022 IEEE/ACM 1st International Workshop on Software Engineering for Responsible Artificial Intelligence (SE4RAI)*. IEEE, 45–49.

[111] Gary McGraw, Richie Bonett, Victor Shepardson, and Harold Figueroa. 2020. The Top 10 Risks of Machine Learning Security. *Computer* 53, 6 (2020), 57–61.

[112] Noha Medhat, Sherin M Moussa, Nagwa Lotfy Badr, and Mohamed F Tolba. 2020. A framework for continuous regression and integration testing in IoT systems based on deep learning and search-based techniques. *IEEE Access* 8 (2020), 215716–215726.

[113] Moritz Meister, Sina Sheikholeslami, Amir H Payberah, Vladimir Vlassov, and Jim Dowling. 2020. Maggy: Scalable asynchronous parallel hyperparameter search. In *Proceedings of the 1st Workshop on Distributed Machine Learning*. 28–33.

[114] Tong Meng, Xuyang Jing, Zheng Yan, and Witold Pedrycz. 2020. A survey on machine learning for data fusion. *Information Fusion* 57 (2020), 115–129.

[115] Chulhong Min, Akhil Mathur, Utku Gunay Acer, Alessandro Montanari, and Fahim Kawsar. 2021. SensiX++: Bringing mlops and multi-tenant model serving to sensory edge devices. *arXiv preprint arXiv:2109.03947* (2021).

[116] Eduardo Mosqueira-Rey, Elena Hernández-Pereira, David Alonso-Ríos, José Bobes-Bascarán, and Ángel Fernández-Leal. 2023. Human-in-the-loop machine learning: a state of the art. *Artificial Intelligence Review* 56, 4 (2023), 3005–3054.

[117] Richard Hugh Moulton, Herna L Viktor, Nathalie Japkowicz, and João Gama. 2018. Clustering in the presence of concept drift. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 339–355.

[118] Aiswarya Raj Munappy, David Issa Mattos, Jan Bosch, Helena Holmström Olsson, and Anas Dakkak. 2020. From ad-hoc data analytics to dataops. In *Proceedings of the International Conference on Software and System Processes*. 165–174.

[119] Jack E Olson. 2003. *Data quality: the accuracy dimension*. Elsevier.

[120] Peyman Oreizy, Michael M Gorlick, Richard N Taylor, Dennis Heimhigner, Gregory Johnson, Nenad Medvidovic, Alex Quilici, David S Rosenblum, and Alexander L Wolf. 1999. An architecture-based approach to self-adaptive software. *IEEE Intelligent Systems and Their Applications* 14, 3 (1999), 54–62.

[121] Soya Park, April Yi Wang, Ban Kawas, Q Vera Liao, David Piorkowski, and Marina Danilevsky. 2021. Facilitating knowledge sharing from domain experts to data scientists for building nlp models. In *26th International Conference on Intelligent User Interfaces*. 585–596.

[122] Xiaolei Qian and Gio Wiederhold. 1986. Knowledge-based Integrity Constraint Validation.. In *VLDB*, Vol. 86. Citeseer, 25–28.

[123] Joaquin Quionero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. 2009. *Dataset Shift in Machine Learning*. The MIT Press.

[124] Alberto Racca and Luca Magri. 2021. Robust Optimization and Validation of Echo State Networks for learning chaotic dynamics. *Neural Networks* 142 (2021), 252–268.

[125] Emmanuel Raj. 2021. *Engineering MLOps: Rapidly build, test, and manage production-ready machine learning life cycles at scale*. Packt Publishing Ltd.

[126] Emmanuel Raj, David Buffoni, Magnus Westerlund, and Kimmo Ahola. 2021. Edge MLOps: An automation framework for aiot applications. In *2021 IEEE International Conference on Cloud Engineering (IC2E)*. IEEE, 191–200.

[127] Sebastian Raschka and Vahid Mirjalili. 2019. *Python machine learning: Machine learning and deep learning with Python, scikit-learn, and TensorFlow 2*. Packt Publishing Ltd.

[128] Alexander Ratner, Stephen H Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. 2017. Snorkel: Rapid training data creation with weak supervision. In *Proceedings of the VLDB Endowment. International Conference on Very Large Data Bases*, Vol. 11. NIH Public Access, 269.

[129] Christopher Ré, Feng Niu, Pallavi Gudipati, and Charles Srisuwananukorn. 2019. Overton: A data system for monitoring and improving machine-learned products. *arXiv preprint arXiv:1909.05372* (2019).

[130] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. 2018. Learning to reweight examples for robust deep learning. In *International conference on machine learning*. PMLR, 4334–4343.

[131] Cedric Renggli, Luka Rimanic, Nezihe Merve Gürel, Bojan Karlaš, Wentao Wu, and Ce Zhang. 2021. A data quality-driven view of mlops. *arXiv preprint arXiv:2102.07750* (2021).

[132] Cedric Renggli, Luka Rimanic, Luka Kolar, Nora Hollenstein, Wentao Wu, and Ce Zhang. 2020. Ease. ml/snoopy: Towards Automatic Feasibility Study for Machine Learning Applications. *arXiv preprint arXiv:2010.08410* (2020).

[133] Andras Rozsa, Manuel Günther, and Terrance E Boult. 2016. Are accuracy and robustness correlated. In *2016 15th IEEE international conference on machine learning and applications (ICMLA)*. IEEE, 227–232.

[134] Cynthia Rudin and Kiri L Wagstaff. 2014. Machine learning for science and society. , 9 pages.

[135] Philipp Ruf, Manav Madan, Christoph Reich, and Djaffar Ould-Abdeslam. 2021. Demystifying mlops and presenting a recipe for the selection of open-source tools. *Applied Sciences* 11, 19 (2021), 8861.

[136] Iqbal H Sarker. 2021. Machine learning: Algorithms, real-world applications and research directions. *SN Computer Science* 2, 3 (2021), 1–21.

[137] Iqbal H Sarker. 2022. AI-based modeling: Techniques, applications and research issues towards automation, intelligent and smart systems. *SN Computer Science* 3, 2 (2022), 1–20.

[138] Sebastian Schelter, Felix Biessmann, Tim Januschowski, David Salinas, Stephan Seufert, and Gyuri Szarvas. 2018. On challenges in machine learning model management. (2018).

[139] Sebastian Schelter, Dustin Lange, Philipp Schmidt, Meltem Celikel, Felix Biessmann, and Andreas Grafberger. 2018. Automating large-scale data quality verification. *Proceedings of the VLDB Endowment* 11, 12 (2018), 1781–1794.

[140] Naeem Seliya, Taghi M Khoshgoftaar, and Jason Van Hulse. 2009. A study on the relationships of classifier performance metrics. In *2009 21st IEEE international conference on tools with artificial intelligence*. IEEE, 59–66.

[141] Rudy Semola, Vincenzo Lomonaco, and Davide Bacciu. 2022. Continual-Learning-as-a-Service (CLaaS): On-Demand Efficient Adaptation of Predictive Models. *arXiv preprint arXiv:2206.06957* (2022).

[142] Ali Shafahi, Parsa Saadatpanah, Chen Zhu, Amin Ghiasi, Christoph Studer, David Jacobs, and Tom Goldstein. 2019. Adversarially robust transfer learning. *arXiv preprint arXiv:1905.08232* (2019).

[143] Shadi Shahoud, Sonja Gunnarsdottir, Hatem Khalloof, Clemens Duepmeier, and Veit Hagenmeyer. 2019. Facilitating and managing machine learning and data analysis tasks in Big Data environments using web and microservice technologies. In *Proceedings of the 11th International Conference on Management of Digital EcoSystems*. 80–87.

[144] Shrey Shrivastava, Dhaval Patel, Nianjun Zhou, Arun Iyengar, and Anuradha Bhamidipaty. 2020. DQLearn: a toolkit for structured data quality learning. In *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 1644–1653.

[145] Lucas Cardoso Silva, Fernando Rezende Zagatti, Bruno Silva Sette, Lucas Nildaimon dos Santos Silva, Daniel Lucrédio, Diego Furtado Silva, and Helena de Medeiros Caseli. 2020. Benchmarking machine learning solutions in production. In *2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 626–633.

[146] Julian Soh and Priyanshi Singh. 2020. Machine learning operations. In *Data Science Solutions on Azure*. Springer, 259–279.

[147] Arjun Soin, Jameson Merkow, Jin Long, Joesph Paul Cohen, Smitha Saligrama, Stephen Kaiser, Steven Borg, Ivan Tarapov, and Matthew P Lungren. 2022. CheXstray: Real-time Multi-Modal Data Concordance for Drift Detection in Medical Imaging AI. *arXiv preprint arXiv:2202.02833* (2022).

[148] Liwei Song, Reza Shokri, and Prateek Mittal. 2019. Privacy risks of securing machine learning models against adversarial examples. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 241–257.

[149] Yiliao Song, Jie Lu, Anjin Liu, Haiyan Lu, and Guangquan Zhang. 2021. A Segment-Based Drift Adaptation Method for Data Streams. *IEEE Transactions on Neural Networks and Learning Systems* (2021).

[150] Ola Spjuth, Jens Frid, and Andreas Hellander. 2021. The machine learning life cycle and the cloud: implications for drug discovery. *Expert opinion on drug discovery* 16, 9 (2021), 1071–1079.

[151] John Spray, Roopak Sinha, Arnab Sen, and Xingbin Cheng. 2021. Building Maintainable Software Using Abstraction Layering. *IEEE Transactions on Software Engineering* 48, 11 (2021), 4397–4410.

[152] Monika Steidl, Michael Felderer, and Rudolf Ramler. 2023. The pipeline for the continuous development of artificial intelligence models—Current state of research and practice. *Journal of Systems and Software* 199 (2023), 111615.

[153] Masahiro Sugimoto, Masahiro Takada, and Masakazu Toi. 2013. Comparison of robustness against missing values of alternative decision tree and multiple logistic regression for predicting clinical data in primary breast cancer. In *2013 35th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*. IEEE,

3054–3057.

[154] Masashi Sugiyama and Motoaki Kawanabe. 2012. *Machine Learning in Non-Stationary Environments: Introduction to Covariate Shift Adaptation*. The MIT Press.

[155] Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul Buenau, and Motoaki Kawanabe. 2007. Direct importance estimation with model selection and its application to covariate shift adaptation. *Advances in neural information processing systems* 20 (2007).

[156] Georgios Symeonidis, Evangelos Nerantzis, Apostolos Kazakis, and George A Papakostas. 2022. MLOps-Definitions, Tools and Challenges. In *2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC)*. IEEE, 0453–0460.

[157] Damian A Tamburri. 2020. Sustainable mlops: Trends and challenges. In *2020 22nd international symposium on symbolic and numeric algorithms for scientific computing (SYNASC)*. IEEE, 17–23.

[158] Samson Tan, Araz Taeihagh, and Kathy Baxter. 2022. The Risks of Machine Learning Systems. *arXiv preprint arXiv:2204.09852* (2022).

[159] Matteo Testi, Matteo Ballabio, Emanuele Frontoni, Giulio Iannello, Sara Moccia, Paolo Soda, and Gennaro Vessio. 2022. MLOps: A Taxonomy and a Methodology. *IEEE Access* (2022).

[160] Chris Thornton, Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. 2013. Auto-WEKA: Combined selection and hyperparameter optimization of classification algorithms. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*. 847–855.

[161] Sunil Thulasidasan, Sushil Thapa, Sayera Dhaubhadel, Gopinath Chennupati, Tanmoy Bhattacharya, and Jeff Bilmes. 2021. An effective baseline for robustness to distributional shift. In *2021 20th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 278–285.

[162] Mark Treveil, Nicolas Omont, Clément Stenac, Kenji Lefevre, Du Phan, Joachim Zentici, Adrien Lavoillotte, Makoto Miyazaki, and Lynn Heidmann. 2020. *Introducing MLOps*. O'Reilly Media.

[163] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. 2018. Robustness may be at odds with accuracy. *arXiv preprint arXiv:1805.12152* (2018).

[164] Alexey Tsymbal, Mykola Pechenizkiy, Pádraig Cunningham, and Seppo Puuronen. 2008. Dynamic integration of classifiers for handling concept drift. *Information Fusion* 9, 1 (2008), 56–68. Special Issue on Applications of Ensemble Methods.

[165] Juan Pablo Usuga Cadavid, Samir Lamouri, Bernard Grabot, Robert Pellerin, and Arnaud Fortin. 2020. Machine learning applied in production planning and control: a state-of-the-art in the era of industry 4.0. *Journal of Intelligent Manufacturing* 31, 6 (2020), 1531–1558.

[166] Andrius Vabalas, Emma Gowen, Ellen Poliakoff, and Alexander J Casson. 2019. Machine learning algorithm validation with a limited sample size. *PloS one* 14, 11 (2019), e0224365.

[167] Tom Van Der Weide, Dimitris Papadopoulos, Oleg Smirnov, Michal Zielinski, and Tim Van Kasteren. 2017. Versioning for end-to-end machine learning pipelines. In *Proceedings of the 1st Workshop on Data Management for End-to-End Machine Learning*. 1–9.

[168] Manasi Vartak, Pablo Ortiz, Kathryn Siegel, Harihar Subramanyam, Samuel Madden, and Matei Zaharia. 2015. Supporting fast iteration in model building. In *NIPS Workshop LearningSys*. 1–6.

[169] Diogo Munaro Vieira, Chrystinne Fernandes, Carlos Lucena, and Sérgio Lifschitz. 2021. Driftage: a multi-agent system framework for concept drift detection. *GigaScience* 10, 6 (2021), giab030.

[170] Stéphan Vincent-Lancrin and Reyer van der Vlies. 2020. Trustworthy artificial intelligence (AI) in education: Promises and challenges. (2020).

[171] Manish Virmani. 2015. Understanding DevOps & bridging the gap from continuous integration to continuous delivery. In *Fifth international conference on the innovative computing technology (intech 2015)*. IEEE, 78–82.

[172] Kiri Wagstaff. 2012. Machine learning that matters. *arXiv preprint arXiv:1206.4656* (2012).

[173] Jindong Wang, Cuiling Lan, Chang Liu, Yidong Ouyang, Tao Qin, Wang Lu, Yiqiang Chen, Wenjun Zeng, and Philip Yu. 2022. Generalizing to unseen domains: A survey on domain generalization. *IEEE Transactions on Knowledge and Data Engineering* (2022).

[174] Xi Wang and Chen Wang. 2019. Time series data cleaning: A survey. *Ieee Access* 8 (2019), 1866–1881.

[175] Jeannette M Wing. 2021. Trustworthy AI. *Commun. ACM* 64, 10 (2021), 64–71.

[176] Ian H. Witten, Eibe Frank, Mark A. Hall, and Christopher J. Pal. 2017. Chapter 9 - Probabilistic methods. In *Data Mining: Practical Machine Learning Tools and Techniques* (4th ed.), Ian H. Witten, Eibe Frank, Mark A. Hall, and Christopher J. Pal (Eds.). Morgan Kaufmann, 335–416.

[177] Robert Woitsch, Wilfrid Utz, Anna Sumereder, Bernhard Dieber, Benjamin Breiling, Laura Crompton, Michael Funk, Karin Bruckmüller, and Stefan Schumann. 2021. Collaborative model-based process assessment for trustworthy AI in robotic platforms. In *International Conference on Society 5.0*. Springer, 163–174.

[178] Huaming Wu, Ziru Zhang, Chang Guan, Katinka Wolter, and Minxian Xu. 2020. Collaborate edge and cloud computing with distributed deep learning for smart city internet of things. *IEEE Internet of Things Journal* 7, 9 (2020), 8099–8110.

[179] Yinjun Wu, Edgar Dobriban, and Susan Davidson. 2020. Deltagrad: Rapid retraining of machine learning models. In *International Conference on Machine Learning*. PMLR, 10355–10366.

[180] Yinjun Wu, Val Tannen, and Susan B Davidson. 2020. Priu: A provenance-based approach for incrementally updating regression models. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 447–462.

[181] Tong Xiao, Tian Xia, Yi Yang, Chang Huang, and Xiaogang Wang. 2015. Learning from massive noisy labeled data for image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2691–2699.

[182] Doris Xin, Litian Ma, Jialin Liu, Stephen Macke, Shuchen Song, and Aditya Parameswaran. 2018. Accelerating human-in-the-loop machine learning: Challenges and opportunities. In *Proceedings of the second workshop on data management for end-to-end machine learning*. 1–4.

[183] Chaobin Xu, Wei Li, Xiaohui Cui, Zhenyu Wang, Fengling Zheng, Xiaowu Zhang, and Bin Chen. 2024. Scarcity-GAN: Scarce data augmentation for defect detection via generative adversarial nets. *Neurocomputing* 566 (2024), 127061.

[184] Limin Yang, Wenbo Guo, Qingying Hao, Arridhana Ciptadi, Ali Ahmadzadeh, Xinyu Xing, and Gang Wang. 2021. CADE: Detecting and explaining concept drift samples for security applications. In *30th USENIX Security Symposium (USENIX Security 21)*. 2327–2344.

[185] Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Russ R Salakhutdinov, and Kamalika Chaudhuri. 2020. A closer look at accuracy vs. robustness. *Advances in neural information processing systems* 33 (2020), 8588–8601.

[186] Xue Ying. 2019. An overview of overfitting and its solutions. In *Journal of physics: Conference series*, Vol. 1168. IOP Publishing, 022022.

[187] Zachary Young and Robert Steele. 2022. Empirical evaluation of performance degradation of machine learning-based predictive models–A case study in healthcare information systems. *International Journal of Information Management Data Insights* 2, 1 (2022), 100070.

[188] Shaochen Yu, Tianwa Chen, Lei Han, Gianluca Demartini, and Shazia Sadiq. 2022. DataOps-4G: On Supporting Generalists in Data Quality Discovery. *IEEE Transactions on Knowledge and Data Engineering* (2022).

[189] Matei Zaharia, Andrew Chen, Aaron Davidson, Ali Ghodsi, Sue Ann Hong, Andy Konwinski, Siddharth Murching, Tomas Nykodym, Paul Ogilvie, Mani Parkhe, et al. 2018. Accelerating the machine learning lifecycle with MLflow. *IEEE Data Eng. Bull.* 41, 4 (2018), 39–45.

[190] Gorka Zárate, Raúl Miñón, Josu Díaz-de Arcaya, and Ana I Torre-Bastida. 2022. K2E: Building MLOps Environments for Governing Data and Models Catalogues while Tracking Versions. In *2022 IEEE 19th International Conference on Software Architecture Companion (ICSA-C)*. IEEE, 206–209.

[191] Marc Zeller, Thomas Waschulzik, Reiner Schmid, and Claus Bahlmann. 2024. Toward a safe MLOps process for the continuous development and safety assurance of ML-based systems in the railway domain. *AI and Ethics* (2024), 1–8.

[192] Friedemann Zenke, Ben Poole, and Surya Ganguli. 2017. Continual learning through synaptic intelligence. In *International Conference on Machine Learning*. PMLR, 3987–3995.

[193] Chen Zhang, Yu Xie, Hang Bai, Bin Yu, Weihong Li, and Yuan Gao. 2021. A survey on federated learning. *Knowledge-Based Systems* 216 (2021), 106775.

[194] Ling Zhang, Xiaosong Wang, Dong Yang, Thomas Sanford, Stephanie Harmon, Baris Turkbey, Bradford J Wood, Holger Roth, Andriy Myronenko, Daguang Xu, et al. 2020. Generalizing deep learning for medical image segmentation to unseen domains via deep stacked transformation. *IEEE transactions on medical imaging* 39, 7 (2020), 2531–2540.

[195] Yizhen Zhao. 2021. Machine Learning in Production: A Literature Review.

[196] Yanjie Zhao, Li Li, Haoyu Wang, Haipeng Cai, Tegawendé F Bissyandé, Jacques Klein, and John Grundy. 2021. On the impact of sample duplication in machine-learning-based android malware detection. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 30, 3 (2021), 1–38.

[197] Aurick Zhou and Sergey Levine. 2021. Bayesian Adaptation for Covariate Shift. *Advances in Neural Information Processing Systems* 34 (2021), 914–927.

[198] Kaiyang Zhou, Ziwei Liu, Yu Qiao, Tao Xiang, and Chen Change Loy. 2021. Domain generalization in vision: A survey. *arXiv preprint arXiv:2103.02503* (2021).

[199] Lina Zhou, Shimei Pan, Jianwu Wang, and Athanasios V Vasilakos. 2017. Machine learning on big data: Opportunities and challenges. *Neurocomputing* 237 (2017), 350–361.

[200] Xiaofeng Zhou, Naiju Zhai, Shuai Li, and Haibo Shi. 2022. Time Series Prediction Method of Industrial Process with Limited Data Based on Transfer Learning. *IEEE Transactions on Industrial Informatics* (2022).

[201] Zhi-Hua Zhou. 2022. Open-environment machine learning. *National Science Review* 9, 8 (2022), nwac123.

[202] Indre Žliobaite. 2010. Change with delayed labeling: When is it detectable?. In *2010 IEEE International Conference on Data Mining Workshops*. IEEE, 843–850.