# Integrating Data Versioning and Management into CI/CD Pipelines for Machine Learning

## SWAMY PRASADARAO VELAGA

Sr. Program Automation Architect, Department of Information Technology

**Abstract:**

**The rapid evolution and widespread adoption of machine learning (ML) applications have underscored the critical importance of data management practices that ensure reproducibility, reliability, and scalability in model development and deployment. Integrating data versioning and management into Continuous Integration and Continuous Deployment (CI/CD) pipelines for ML represents a pivotal strategy to address these challenges. This survey paper explores the significance of data versioning in CI/CD pipelines, examining key benefits such as enhanced reproducibility of experimental results, effective management of data drift, and compliance with regulatory standards. We delve into the challenges associated with integrating data versioning, including handling large dataset sizes, managing dynamic data sources, and ensuring compatibility across diverse data formats. Moreover, the paper discusses best practices and implementation strategies for adopting data versioning in CI/CD pipelines, emphasizing automation, scalability, and integration with Machine Learning Operations (MLOps). Finally, we outline promising future research directions in data versioning, including advancements in automation, security, and cross-domain collaboration, aimed at further enhancing the reliability and transparency of ML workflows. By addressing these aspects, this paper provides a comprehensive overview of current trends, challenges, and opportunities in leveraging data versioning to optimize CI/CD pipelines for machine learning applications.**

**Keywords: Continuous Deployment, AI Systems, Machine Learning Models, Data Versioning**

## 1. Introduction

In recent years, the proliferation of machine learning (ML) and artificial intelligence (AI) applications across various industries has underscored the need for robust and reliable practices in model development and deployment [1]. Central to the success of these endeavors is the quality and management of data, which serves as the lifeblood of ML algorithms. Data, however, is not static—it evolves over time due to updates, corrections, and additions, posing significant challenges to maintaining consistency and reproducibility in ML workflows.

Data versioning emerges as a critical solution in this landscape, offering a systematic approach to track, manage, and version control datasets throughout the ML lifecycle [1]. Unlike traditional software development, where code changes are meticulously tracked, the dynamic nature of data necessitates specialized tools and methodologies to ensure that ML models trained on specific datasets can be recreated and validated reliably.

In the context of Continuous Integration and Continuous Deployment (CI/CD) pipelines—a cornerstone of agile software development adapted for ML—data versioning becomes indispensable. CI/CD pipelines automate the building, testing, and deployment of ML models, streamlining the development process and accelerating time-to-market

One of the primary challenges in machine learning (ML) research and development is ensuring the reproducibility of experimental results across different environments [2]. Variations in data versions used for model training can introduce discrepancies in model performance, hindering the reliability and comparability of research outcomes. To address this challenge, integrating robust data versioning practices into CI/CD

pipelines for ML models is essential. Data versioning allows researchers to meticulously track and manage the specific datasets used at various stages of model development [3]. By associating precise data versions with corresponding model iterations within CI/CD workflows, researchers can accurately reproduce experimental conditions and validate results consistently. This approach not only facilitates transparent experimentation but also enhances the reliability and trustworthiness of ML research outputs, enabling more informed decision-making and advancing the state-of-the-art in AI technologies [4].

The paper is structured as follows: Section 2 provides a comprehensive literature review, Section 3 covers best practices and implementation strategies, Section 4 discusses challenges in data versioning, and Section 5 concludes with future research directions.

## 2. Literature Review

In this paper, we explore Machine Learning Operations (MLOps), focusing on practices that streamline the development of machine learning projects to achieve rapid time-to-value [5]. The lifecycle of machine learning project development involves various roles, software frameworks, and computing resources, leading to complexity often managed through bundled solutions on commercial cloud platforms, known as vendor lock-in. To address this, we propose an alternative MLOps platform leveraging open-source frameworks on virtual resources. Our approach aligns with the development roles in machine learning, integrating seamlessly into the CI/CD workflow typical of machine learning applications. We illustrate this with a practical example: training and deploying a machine learning model to detect vulnerabilities in software repository code [5].

This paper proposes DevOps practices tailored for machine learning applications, aiming to seamlessly integrate development and operational environments [6]. While the initial phases of machine learning development and deployment may appear straightforward, the complexity arises in maintaining these models over time. Without careful design, deploying and using such models can become cumbersome, requiring significant resources for maintenance, enhancement, and monitoring [6]. The paper advocates for applying principles, practices, and tools of continuous integration (CI) and continuous delivery (CD) to mitigate inefficiencies, enable rapid feedback loops, uncover technical debt, enhance value delivery and upkeep, and optimize operational functions in practical machine learning applications [6].

This paper introduces a cloud-based framework and platform designed for comprehensive development and lifecycle management of artificial intelligence (AI) applications [7]. Building upon prior research in platform-level support for cloud-managed deep learning services, the study extends these principles to automate, ensure trust, reliability, traceability, quality control, and reproducibility across AI pipelines. By addressing various use cases and existing challenges, the framework outlines essential components for managing AI application lifecycles. Concrete examples are provided to demonstrate how the framework facilitates the management and execution of model training and continuous learning pipelines, incorporating principles of trusted AI [7].

Microsoft Azure DevOps is a powerful cross-platform automation engine that supports the development, testing, and deployment of both Cloud Native and Non-Cloud Native applications. Its core strength lies in enabling automation through infrastructure as code and seamless integration with Machine Learning Operations (MLOps) frameworks [8]. As software applications increasingly integrate Machine Learning (ML) and Artificial Intelligence (AI) into their development lifecycles, Azure DevOps plays a crucial role in optimizing costs and enhancing customer success [8]. This paper proposes a novel DevOps framework specifically designed for building intelligent Tiny ML dairy agriculture sensors. By leveraging Azure DevOps, this framework aims to streamline product development processes, ensuring high-quality outputs in a cost-effective manner, which is particularly beneficial for small-scale farmers operating within constrained economic environments [8].

The demand for training deep neural networks has surged across academia and industry, but scaling this process remains challenging due to the complexity of tools and hardware involved [9]. This often leads to ad-hoc scripts and custom glue code for managing distributed training setups. Addressing these issues, we introduce

\emph{Alchemist}, an in-house service developed at Apple to streamline distributed training. Designed for simplicity, speed, and scalability, \emph{Alchemist} optimizes the orchestration of distributed training tasks. We detail its architecture, implementation, and provide examples of various distributed training scenarios. Case studies from its deployment in developing autonomous systems highlight significant reductions in training times, enabling efficient handling of large-scale data requirements [9].

This paper offers an analysis of current trends in infrastructure tools supporting modern data-driven application development, deployment, and operation, particularly within the realms of DataOps and MLOps [10]. It explores the application of DevOps methodologies to domains like Data Science, Analytics, Machine Learning, and Artificial Intelligence, emphasizing the need for integrated support across application development and data lifecycle management [10]. Highlighting the importance of data lineage and provenance, the paper discusses recent advancements and tools in cloud-based platforms relevant to DevOps and DataOps, showcasing their diverse applications in different project contexts. Drawing from the development of curricula and courses on DevOps and Cloud-based Software Development at the University of Amsterdam, it reflects on the formulation of the DevOps Body of Knowledge (DevOpsSE BoK) to define essential knowledge areas for SE professionals, facilitating structured curriculum development and profiling. The paper also shares insights from teaching experiences of the DevOps and Cloud-based Software course, illustrating practical applications and educational outcomes [10].

This paper outlines recommendations for designing and implementing the DevOps and Cloud-based Software Development curricula tailored for Computer Science and Software Engineering master's programs [11]. Central to this approach is the DevOps Body of Knowledge for Software Engineering (DevOpsSE BoK), defining essential Knowledge Areas and Units necessary for SE professionals to excel as DevOps engineers or application developers. The DevOpsSE BoK serves as a foundation for specifying required competences and skills, ensuring consistency in curriculum structure and development. The paper draws from the inaugural course offered during the 2018/2019 academic year at the University of Amsterdam, detailing its structure and instructional methodologies, including project-based learning aimed at enhancing students' team skills in Agile development and knowledge exchange. Additionally, it provides a concise overview of prevalent DevOps definitions, concepts, models, and tools, with a specific emphasis on cloud-based DevOps tools crucial for enabling continuous development and improvement—an essential ethos for contemporary agile and data-driven enterprises [11].

**Table 1: Summary for The Literature Review**

| Reference | Model Used | Application | Highlighted Points |
|---|---|---|---|
| **[5]** | Alchemist | Distributed Training | Streamlined distributed training for deep neural networks; significant reduction in training times; scalable and efficient orchestration of training tasks. |
| **[6]** | MLOps | Machine Learning Development | Integration of CI/CD principles for efficient ML application lifecycle management; focus on reducing maintenance efforts and improving operational efficiency. |
| **[7]** | Cloud-based AI Framework | AI Application Lifecycle Management | Automation, trust, reliability, traceability, and quality control in AI pipelines; management of model training and continuous learning pipelines. |
| **[8]** | Azure DevOps | Cloud Native and Non-Cloud Native Applications | Automation engine supporting infrastructure as code and MLOps frameworks; optimization of development lifecycles integrating ML and AI. |
| **[9]** | Alchemist | Autonomous Systems Development | Simplification, speed, and scalability in distributed training; case studies demonstrating significant training time reductions for large-scale data handling. |

| **[10]** | DevOps and DataOps Tools | Data-driven Application Development | Integration of DevOps methodologies in Data Science, Analytics, ML, and AI domains; emphasis on data lineage and provenance; insights from educational curriculum development. |
| --- | --- | --- | --- |
| **[11]** | DevOpsSE BoK | Software Engineering Curriculum | Definition of essential knowledge areas and units for DevOps engineers and developers; structured curriculum development and educational outcomes from practical applications. |

## 3. Best Practices and Implementation Strategies:

This section provides a structured approach to implementing data versioning in CI/CD pipelines, focusing on best practices that enhance reliability, scalability, and collaboration in ML projects.

Effective integration of data versioning and management into CI/CD pipelines for machine learning requires adherence to best practices and thoughtful implementation strategies [12]. These practices not only ensure data integrity and reproducibility but also streamline collaboration and enhance operational efficiency. Key best practices and strategies include:

1. Versioning Datasets at Source: Start data versioning at the source by implementing robust data management practices [12]. This includes capturing metadata such as creation date, author, and purpose, and using version control systems (VCS) or specialized data versioning tools to manage dataset changes over time.

2. Automated Versioning in Pipelines: Integrate automated data versioning into CI/CD pipelines to ensure consistency and traceability across the ML workflow [13]. Tools like DVC (Data Version Control) or MLflow enable automatic capture and tracking of dataset versions as part of the pipeline execution, minimizing manual errors and ensuring reproducibility of model training and evaluation.

3. Metadata and Lineage Tracking: Implement comprehensive metadata and lineage tracking mechanisms to capture dependencies between datasets, models, and experiments [14]. This includes recording transformations, preprocessing steps, and data splits, providing a clear audit trail for data-driven decisions and compliance with regulatory requirements.

4. Collaborative Data Governance: Establish collaborative data governance frameworks that define roles, responsibilities, and access controls for managing data versions within CI/CD pipelines [15]. Encourage transparency and accountability among team members, ensuring that changes to datasets are documented and validated before integration into model training processes.

5. Regular Validation and Testing: Incorporate validation and testing procedures to verify the integrity and quality of data versions throughout the CI/CD pipeline [16]. This includes automated checks for data consistency, schema evolution, and compatibility with downstream ML tasks, ensuring that only reliable datasets are used for model development and deployment.

6. Scalability and Performance Optimization: Optimize data versioning practices for scalability and performance by leveraging distributed storage solutions and efficient indexing mechanisms [17]. This allows handling large-scale datasets and accommodating growing computational demands without compromising pipeline efficiency or reliability.

7. Continuous Improvement and Feedback Loop: Foster a culture of continuous improvement by soliciting feedback from stakeholders and iterating on data versioning processes. Regularly review and update versioning strategies based on evolving business requirements, technological advancements, and feedback from end-users to maintain alignment with organizational goals and objectives [17]. Figure 1 presents the best practice for CI/CD pipeline design.
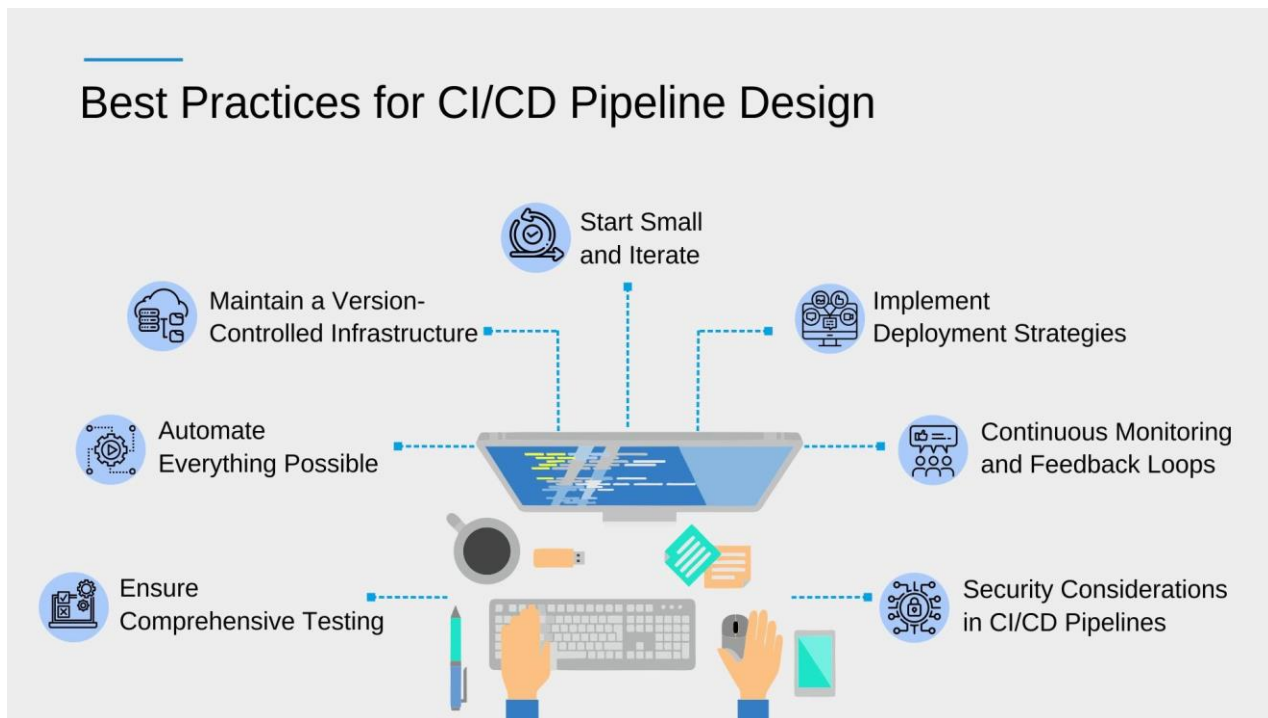
**Figure 1: Best Practices for CI/CD Pipeline Design[1]**

By adopting these best practices and implementation strategies, organizations can effectively integrate data versioning and management into CI/CD pipelines for machine learning. This approach not only enhances the reproducibility and reliability of ML workflows but also promotes collaboration, compliance, and efficiency in delivering AI-driven solutions.

## 4. Challenges in Data Versioning

Integrating data versioning into Continuous Integration and Continuous Deployment (CI/CD) pipelines for machine learning introduces several complex challenges that impact data management, model reproducibility, and operational efficiency. Addressing these challenges is crucial to ensuring the reliability and scalability of ML workflows. Key challenges include:

1. Large Dataset Sizes: Machine learning models often require large-scale datasets for training, testing, and validation. Managing versions of these large datasets can lead to significant storage and computational overheads within CI/CD pipelines [15]. Efficient storage solutions and optimized version control mechanisms are essential to handle the volume and velocity of data updates without compromising pipeline performance or scalability.

2. Data Drift and Distribution Changes: Real-world datasets are dynamic and subject to changes in distribution over time (data drift). Integrating data versioning must account for these changes to maintain model accuracy and reliability [16]. Tracking and comparing data distributions across different versions require sophisticated monitoring and alerting systems within CI/CD pipelines, enabling proactive detection and mitigation of data drift effects on model performance.

3. Compatibility with Different Data Formats and Sources: ML projects often involve heterogeneous data sources and formats, ranging from structured databases to unstructured text or image data. Ensuring compatibility and consistency in versioning across diverse data types poses technical challenges [17]. Data versioning tools and frameworks must support seamless integration with various data formats and sources,

---

[1] https://www.neebal.com/blog/seamless-software-releases_-a-deep-dive-into-ci_cd-pipelines-1

providing unified versioning capabilities without imposing constraints on data representation or preprocessing workflows.

4. Versioning Granularity and Efficiency: Determining the appropriate granularity of data versioning poses a trade-off between granularity and storage efficiency. While fine-grained versioning captures detailed changes at the dataset level, it can increase storage requirements and complicate version management [15]. Balancing granularity with practical storage and operational considerations requires careful design of version control policies and pruning strategies within CI/CD pipelines.

5. Collaborative Data Governance and Access Control: Effective data versioning necessitates collaborative data governance frameworks that define roles, responsibilities, and access controls across multidisciplinary teams [17]. Ensuring data integrity and preventing unauthorized modifications or access requires robust authentication mechanisms and audit trails within CI/CD pipelines. Establishing clear protocols for data sharing, review, and validation enhances transparency and accountability in versioning practices.

6. Regulatory Compliance and Data Privacy: Compliance with regulatory standards (e.g., GDPR, HIPAA) mandates rigorous data governance practices, including data versioning and auditability [17]. CI/CD pipelines must enforce data privacy measures and ensure secure handling of sensitive information across versioned datasets. Implementing encryption, anonymization techniques, and policy-driven access controls are essential to safeguarding data integrity and maintaining regulatory compliance in ML workflows.

Navigating these challenges requires a comprehensive understanding of data lifecycle management, version control principles, and integration strategies tailored to the unique requirements of machine learning applications. By addressing these challenges effectively, organizations can leverage data versioning within CI/CD pipelines to enhance the reproducibility, reliability, and scalability of ML models in production environments.

## 5. Conclusion

The integration of data versioning and management into Continuous Integration and Continuous Deployment (CI/CD) pipelines represents a pivotal advancement in enhancing the reliability, reproducibility, and efficiency of machine learning (ML) workflows. Throughout this survey paper, we have explored the importance of data versioning in addressing critical challenges such as data reproducibility across environments, management of data drift, and compatibility with diverse data formats. By adopting best practices and implementation strategies, organizations can establish robust frameworks for versioning datasets, tracking lineage, and ensuring compliance with regulatory standards.

Data versioning enables researchers and data scientists to maintain a transparent audit trail of dataset changes, facilitating reproducible experiments and informed decision-making in model development and deployment. The automation and integration of versioning processes within CI/CD pipelines streamline collaboration, accelerate experimentation cycles, and mitigate risks associated with data inconsistency and model degradation over time.

Looking ahead, the field of data versioning in ML CI/CD pipelines presents several promising avenues for future research and development. One critical area is the enhancement of automation capabilities to streamline data versioning processes across large-scale and distributed ML environments. This includes automating data lineage tracking to capture dependencies and transformations across datasets, implementing efficient version pruning strategies to manage storage costs, and optimizing computational resources for scalability and performance. Additionally, integrating data versioning practices with MLOps frameworks is crucial for fostering end-to-end automation, monitoring, and governance of ML models throughout their lifecycle. This integration can enhance operational efficiency, facilitate model deployment, and ensure consistency in versioning practices across diverse ML workflows. Another key challenge lies in effectively managing dynamic data sources and streaming data streams, necessitating the development of techniques to ensure real-

time responsiveness and adaptability in CI/CD pipelines. Furthermore, advancing security measures and privacy-preserving techniques in data versioning operations is essential to safeguard sensitive data and ensure compliance with stringent regulatory requirements such as GDPR and HIPAA.

**References**
1. Vadavalasa, Ram Mohan. "End to end CI/CD pipeline for Machine Learning." International Journal of Advance Research, Ideas and Innovation in Technology 6 (2020): 906-913.
2. Zhou, Yue, Yue Yu, and Bo Ding. "Towards mlops: A case study of ml pipeline platform." 2020 International conference on artificial intelligence and computer engineering (ICAICE). IEEE, 2020.
3. Dhaliwal, Neha. "VALIDATING SOFTWARE UPGRADES WITH AI: ENSURING DEVOPS, DATA INTEGRITY AND ACCURACY USING CI/CD PIPELINES." JOURNAL OF BASIC SCIENCE AND ENGINEERING 17.1 (2020).
4. Körner, Christoph, and Kaijisse Waaijer. Mastering Azure Machine Learning: Perform large-scale end-to-end advanced machine learning in the cloud with Microsoft Azure Machine Learning. Packt Publishing Ltd, 2020.
5. Liu, Yan, et al. "Building a platform for machine learning operations from open source frameworks." IFAC-PapersOnLine 53.5 (2020): 704-709.
6. Karamitsos, Ioannis, Saeed Albarhami, and Charalampos Apostolopoulos. "Applying DevOps practices of continuous automation for machine learning." Information 11.7 (2020): 363.
7. Hummer, Waldemar, et al. "Modelops: Cloud-based lifecycle management for reliable and trusted ai." 2019 IEEE International Conference on Cloud Engineering (IC2E). IEEE, 2019.
8. Vuppalapati, Chandrasekar, et al. "Automating tiny ml intelligent sensors devops using microsoft azure." 2020 ieee international conference on big data (big data). IEEE, 2020.
9. Ma, Minghuang, et al. "Democratizing production-scale distributed deep learning." arXiv preprint arXiv:1811.00143 (2018).
10. Demchenko, Yuri. "From DevOps to DataOps: Cloud based Software Development and Deployment." Proc. The International Conference on High Performance Computing and Simulation (HPCS 2020). 2020.
11. Demchenko, Yuri, et al. "Teaching DevOps and cloud based software engineering in university curricula." 2019 15th International Conference on eScience (eScience). IEEE, 2019.
12. Bai, Haishi. Zen of Cloud: Learning Cloud Computing by Examples. CRC Press, 2019.
13. Belmont, Jean-Marcel. Hands-On Continuous Integration and Delivery: Build and release quality software at scale with Jenkins, Travis CI, and CircleCI. Packt Publishing Ltd, 2018.
14. Arundel, John, and Justin Domingus. Cloud Native DevOps with Kubernetes: building, deploying, and scaling modern applications in the Cloud. O'Reilly Media, 2019.
15. Atwal, Harvinder, and Harvinder Atwal. "DevOps for DataOps." Practical DataOps: Delivering Agile Data Science at Scale (2020): 161-189.
16. Jokinen, Oskari. "Software development using DevOps tools and CD pipelines, A case study." Helsingin yliopisto (2020): 54.
17. Viitasuo, Ella. "Adding security testing in DevOps software development with continuous integration and continuous delivery practices." (2020).