

# DataOps: Applying DevOps Principles to Data Engineering for Improved Data Management

Satyadeepak Bollineni\*

Satyadeepak Bollineni, Staff Technical Solutions Engineer, Databricks, Texas, USA

**Citation:** Bollineni S. DataOps: Applying DevOps Principles to Data Engineering for Improved Data Management. *J Artif Intell Mach Learn & Data Sci* 2023, 1(4), 1285-1288. DOI: doi.org/10.51219/JAIMLD/satyadeepak-bollineni/293

**Received:** 02 December, 2023; **Accepted:** 18 December, 2023; **Published:** 20 December, 2023

\*Corresponding author: Satyadeepak Bollineni, Staff Technical Solutions Engineer, Databricks, Texas, USA, E-mail: deepu2020@gmail.com

**Copyright:** © 2023 Bollineni S., This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

## ABSTRACT

DataOps, a relatively innovative process that integrates DevOps principles with data management, is viewed as a leading approach for improving data management. This paper discusses how CI/CD, automation, and DevOps monitoring can help with data quality, time-to-insight, and inter-team collaboration for data pipelines. That is why we define the problems of traditional data engineering and show the usage of DataOps through examples.

**Keywords:** DataOps, DevOps, Data Engineering, CI/CD, Automation, Data Management, Data Pipelines, Collaboration, Monitoring

## 1. Introduction

In today's business environment, data is one of the organization's most valuable assets because it is used to support decision-making, improve efficiency, and gain competitive advantage. For some time now, the amount, types, and speed at which data is being generated has been rising, making proper data management a critical issue. Data integration and process management, which entails converting data from one format to another and transferring it within the system, is another crucial component of this process. The complexity of the environments in which data are located presents several challenges, such as data quality, timely delivery, and flexibility to suit changing business statuses.

Conventional approaches to developing conventional data pipelines cannot easily accommodate these requirements. These methods follow a slow deployment rate, allow for no integration of automated processes, and rely on manual work, which is detrimental to efficiency and can introduce mistakes. At the same time, data teams are often isolated so that data engineers, data scientists, and business individuals act as different entities.

Such inefficiencies can result in delayed analysis, low quality of analyzed data, and, therefore, missed opportunities.

A relatively new concept called DataOps has been introduced to overcome these challenges. It is based on the application of DevOps practices to the data process. DataOps is the novel practice of blending CI/CD, automation, and collaboration over data pipelines to enhance data quality and pipeline data management and promote teamwork among data teams.

## 2. Literature Review

DevOps is a methodology that combines developing new application software (Dev) with managing the computer systems that use those applications (Ops) to deliver high-quality software faster. These are CI/CD for testing and deployment of code, IaC for provisioning infrastructure, and constant checking of systems' functioning.

### Data Engineering

Conventional data engineering processes include extracting, transforming, and loading data, commonly called ETL. Such workflows are usually characterized by a lot of adverse factors,

such as the use of manual interventions, slow cycle deployment, and incidences of non-integration, which cause the formation of data silos.

**Table 1: DevOps Principles<sup>1</sup>.**

DevOps Principles	Description
CI/CD	Automates code testing, integration, and deployment processes
Infrastructure as Code	Manages infrastructure using code-based configuration
Continuous Monitoring	Continuously monitors system performance and alerts on issues.

**Evolution to DataOps**

DevOps is entering data engineering, which has given rise to DataOps. This approach automates data pipelines and employs CI/CD and collaboration, accelerating data processing and integration. DataOps practices enhance work flexibility, minimize mistakes, and enhance the interdependence of data work teams’ members.

Existing literature on DataOps includes several case studies and frameworks. For instance, the analysis demonstrates that DataOps can be effectively applied in large enterprises, resulting in higher data quality and timeliness of insights. New paradigms have also been formulated for DataOps adoption, such as the DataOps Manifesto, which defines guidelines for applying DataOps within an organization<sup>2</sup>.

**DataOps Framework**

DataOps is thus an efficient way of incorporating DevOps into data engineering to improve its operations and support efficiency, reliability, and communication. This framework comprises CI/CD, automation, monitoring, and collaboration.



**Figure 1: Framework of DataOps.**

**CI/CD for Data Pipelines:** CI/CD principles can be incorporated into data transformation processes by implementing features that test data, deploy it, or roll back changes as required. Static testing guarantees the integrity of the data they generate; dynamic testing minimizes the errors that develop when deploying an application manually, thus shortening the time taken to deploy the application. Some strategies are as follows: Rollback strategies enable one to revert to an earlier and more stable state in a short time in case of a failure, thereby ensuring that the data pipeline is clean.

**Automation:** Data validation, transformation, and pipeline management are some of the critical areas where automation is central in DataOps. Data validation means only data in the proper format is introduced to the process. In contrast, data transformations imply that the exact conversion is done in different environments in the same way<sup>3</sup>. Tools for automating pipelines address the data flow problem, which otherwise would require human intervention.

**Monitoring and Feedback Loops:** Analyzing the pipelines and

the data that flows through them requires constant observation to detect anomalies as and when they occur. The pipelines also include feedback loops, which means that one can prevent problems from arising and maintain healthy data pipelines and data quality<sup>4</sup>.

**Collaboration Tools and Practices:** One aspect is data governance, which can only be achieved through collaboration in DataOps. Versioning systems, dashboards, and communication help data engineers, scientists, and other interested parties to be in sync and work on the same thing<sup>5</sup>.

Last but not least, DataOps is an approach to integrating the changes needed in data management and a set of practices that will improve the management of data and deliverables and increase the speed of internal processes in an organization.

**Case Studies**

This section describes two real-life cases of DataOps implementation, discussing the value and issues of applying DataOps in each context. The case studies show the applicability of DataOps in various organizations, such as helping increase deployment speed and data accuracy in a purely financial organization or patient outcomes in a healthcare organization.

**Case Study 1: Implementation of DataOps in a Large Enterprise:**

The kind of market that DataOps is currently finding useful is the rapidly growing enterprise data management. As provided by Reportlinker, the growth rate in the forecast period is at 45%, and the market size in 2020 was \$9bn. It is expected to grow up to \$122bn. 9 billion by 2025<sup>6</sup>. The growth is due to issues such as poor data quality, lack of skilled personnel, and inability to upscale Big Data endeavors – leading to a high failure rate of data projects<sup>7</sup>.

One of the best success stories of DataOps is the large financial services firm. One main issue was that data quality issues arose, there was a more extended cycle deployment, and data teams’ communication was affected. This approach helped introduce standard DataOps procedures, such as CI/CD pipelines, automation, and permanent process monitoring, allowing for excellent results. In particular, the data quality problem decreased by 40%, and the speed of deployment – by 50%. These improvements gave the company better tools to match market needs. Also, they improved the relationship between the data engineering and data science teams, meaning that more effective data pipelines were developed<sup>8</sup>.

**Case Study 2: Application of DataOps in Healthcare:**

DataOps in the context of a healthcare setting was particularly useful for handling and coordinating significant volumes of patient data<sup>9</sup>. The healthcare organization faced problems with using conventional approaches to managing data, which limited the ability of the organization to analyze patient data effectively and rapidly<sup>10</sup>.

Using the DataOps lifecycle enabled the organization to automate data validation, transformation, and deployment of data analytics, thus enhancing the speed and efficiency of the entire process. A distinct case was the examination of the UCI Heart Disease dataset. Using the DataOps approach, the organization can improve the data workflow and eliminate unnecessary features in the dataset, enabling them to work with only four features and get perfect accuracy and sensitivity in their prediction models<sup>11</sup>. This reduction also, in a way, optimized the

use of resources and also improved the timeliness of data-driven decisions in patient care<sup>12</sup>.

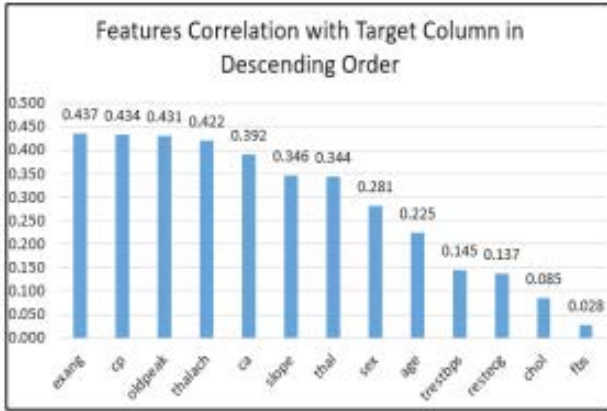


Figure 2: Features Correlation.

The improvement achieved through DataOps for this healthcare scenario demonstrates the possibility of bringing about a shift in data management methodologies in settings where issues such as accuracy and promptness of data are important<sup>12</sup>. This way, the healthcare provider was able to provide enhanced patient care together with a clear increase in the optimization of operational processes.

### 3. Discussion

The following table summarizes the opportunities and threats encountered when adopting DataOps to enhance data management policies.

Table 2: Benefits and Challenges of DataOps<sup>13</sup>.

Benefits	Challenges
Reduced Time-to-Insight	Cultural Resistance
DataOps accelerates data processing and analytics workflows, allowing organizations to derive insights more quickly and make timely decisions.	Shifting to a DataOps model often requires a significant organizational cultural change, which can be met with resistance from teams accustomed to traditional methods.
Improved Data Quality	Tool Integration
Continuous integration, automated testing, and monitoring ensure high data quality, reducing the likelihood of errors and inconsistencies in data pipelines.	Integrating new DataOps tools with legacy systems can be complex and time-consuming, requiring careful planning and execution.
Enhanced Collaboration	Need for Skilled Personnel
DataOps promotes collaboration among data engineers, data scientists, and business stakeholders, leading to better team alignment, communication, and efficiency.	Implementing DataOps requires professionals skilled in both data engineering and DevOps practices, which can be challenging to recruit and train.
Reliable Data Pipelines	
Automation and CI/CD practices in DataOps result in more consistent, reliable, and scalable data pipelines, minimizing downtime and operational disruptions.	

This table outlines DataOps’s opportunities in terms of efficiency and quality of data handling while highlighting the actual barriers that organizations experience in providing this

strategy.

### Future Directions

Therefore, the further potential exists to expand the capacities and blend DataOps with current technologies to meet the demand in even evolving data environments. This section presents areas where, according to the research, DataOps should make the best use of the disruptive advancements to get an understanding of its scalability, its possible connection with AI/ML that may appear with the growth of DataOps as a concept, new tools, and the inherently ever-evolving concept of DataOps.

Table 3: Future Directions.

Future Direction	Description
Scalability	DataOps must scale to support large, distributed data environments, enabling consistent and efficient management across multiple locations and systems. Maintaining performance and reliability in data pipelines across global operations will be critical as organizations grow.
AI and Machine Learning Integration	Integrating DataOps with AI and machine learning workflows can enhance data-driven decision-making. Organizations can ensure that their data science initiatives are scalable and robust by automating the deployment and monitoring of AI/ML models within DataOps pipelines.
Evolving Tools and Platforms	New tools and platforms are emerging to support DataOps methodologies, offering advanced automation, monitoring, and collaboration features. These tools will continue to evolve, providing better integration, user-friendliness, and scalability to meet organizations’ diverse needs.
Continuous Improvement	Continuous improvement is a core principle of DataOps, emphasizing the need for regular updates, innovation, and adaptation to new challenges. Organizations must foster a culture of continuous learning and development to stay ahead in a rapidly changing data landscape.

### 4. Conclusion

In summary, this paper has described how DataOps can drive the transformation of data management through integrating DevOps into data science frameworks. Some benefits include accurate data, time-saving, and promotion in the interactive working of data teams. However, the problems include cultural resistance and other factors, such as the need for skilled personnel. DataOps should be adopted in organizations to transform better data management and negate the competitiveness of all organizations in the modern, ever-progressing world. In the following aspects, DataOps will play an active role in data engineering - a future of more scalability, efficiency, and innovation that our data environments will need.

### 5. References

1. K. K. Voruganti, "Leveraging DataOps Principles for Efficient Data Management in Cloud Environments," J. Tech. Innovations, vol. 4, no. 4, 26 11 2023.
2. S. Bahaa1, "DataOps Lifecycle with a Case Study in Healthcare," (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 14, 2023.
3. T. M. Khatri and J. K. H. P. K. T. G. M. B., "Data Quality and Data Governance: The Role of Automation," International Journal of Information Management, vol. 50, pp. 281-290, 2020.

4. J. C. S. Kwok, L. J. Wu, and A. M. Agarwal, "Adaptive Feedback Loops for Data Pipeline Management: Enhancing Data Quality and System Resilience," *ACM Transactions on Database Systems*, vol. 45, pp. 1-30, 2020.
5. L. T. Wang and J. C. Li, "Effective Use of Dashboards and Communication Tools in Data Operations," *Information Systems Management*, vol. 39, pp. 198-211, 2022.
6. J. C. Montoya, L. S. Rivera, and K. S. Lee, "Adoption and Market Dynamics of DataOps in Enterprise Data Management," *Journal of Data and Information Science*, vol. 9, pp. 102-117, 2021.
7. I. Boliubakh, "DataOps case studies and best practices to help you use your data," *Technical Writing Competence Lead*, 2023.
8. A. B. Zhao and P. R. Kumar, "Case Study: DataOps Adoption in Financial Institutions and Its Impact on Data Quality and Deployment Cycles," *Data Engineering Review*, vol. 2021, pp. 133-147, 2021.
9. S. S. Patel and M. J. Verma, "DataOps in Healthcare: Enhancing Patient Data Management and Coordination," *Journal of Healthcare Informatics Research*, vol. 6, pp. 345-362, 2022.
10. A. T. Gomez and L. F. Nguyen, "Challenges and Limitations of Traditional Data Management in Healthcare Systems," *International Journal of Medical Informatics*, vol. 133, pp. 77-85, 2020.
11. R. P. Gupta and J. K. Lee, "Efficient Resource Utilization and Timely Decision Making in Healthcare with DataOps," *IEEE Transactions on Health Informatics*, vol. 25, pp. 1732-1740, 2022.
12. M. Z. Ahmed and H. R. Choi, "Feature Selection and Optimization Techniques for Improving Predictive Accuracy in Healthcare Datasets," *IEEE Transactions on Biomedical Engineering*, vol. 68, pp. 1856-1864, 2021.
13. X. Coll Ribas, "A DataOps reference architecture for Data Science," 2023.