International Journal of Computer Technology and Electronics Communication (IJCTEC)



 $|\;ISSN:\;2320\text{-}0081\;|\;\underline{www.ijctece.com}\;|\;A\;Peer\text{-}Reviewed,\;Refereed,\;a\;Bimonthly\;Journal|$

|| Volume 7, Issue 3, May - June 2024 ||

DOI: 10.15680/IJCTECE.2024.0703002

Version-Controlled Analytics: Integrating DBT with GIT for Scalable Data Pipelines

Vivaan Kaur Bhatt, Ira Naidu Gupta, Ishaan Rao Patel

Department of CSE, Nagarjuna College of Engineering and Technology, Bengaluru, India

ABSTRACT: DBT (Data Build Tool) has revolutionized data transformation workflows, enabling data engineers to model, test, and document data within cloud data warehouses. When coupled with Git for version control, DBT enables more efficient collaboration, reproducibility, and error tracking in data engineering teams. This paper explores how integrating DBT with Git can streamline the development and deployment of data pipelines. The research focuses on the advantages of using Git for managing DBT projects, ensuring collaborative workflows, maintaining data pipeline versions, and automating deployments. We discuss best practices for integrating DBT with Git to improve data pipeline efficiency, reduce errors, and ensure a smoother CI/CD process in modern data engineering environments.

KEYWORDS: DBT, Git, Data Pipelines, Version Control, CI/CD, Data Engineering, Cloud Data Warehouses, Automation, Collaboration.

I. INTRODUCTION

In modern data engineering, building and maintaining robust, scalable data pipelines is a critical aspect of ensuring data availability, consistency, and reliability. DBT (Data Build Tool) has become a popular tool for transforming raw data into structured formats ready for analysis, but managing these transformations, especially in collaborative environments, can be complex. Git integration with DBT offers a solution to this challenge by enabling version control, collaboration, and the automation of deployments.

In this paper, we explore how integrating Git with DBT can streamline the creation, management, and deployment of data pipelines. By incorporating version control, DBT with Git ensures that data transformations are reproducible, facilitates easy collaboration across teams, and improves transparency in the development process. Additionally, leveraging CI/CD (Continuous Integration/Continuous Deployment) pipelines automates the testing and deployment of DBT models, further enhancing the efficiency of data pipeline operations.

II. LITERATURE REVIEW

1. Evolution of DBT in Data Engineering

DBT has transformed the way data engineers approach data modeling and transformation. Initially designed to simplify SQL-based transformations, DBT's capabilities have expanded to support testing, documentation, and version control. As cloud data warehouses such as Snowflake, BigQuery, and Redshift gained popularity, DBT became the go-to tool for transforming data in the cloud. However, without proper version control, managing and collaborating on DBT projects at scale can be cumbersome (DBT Labs, 2021).

2. Git and Version Control in Data Engineering

Git, a distributed version control system, has long been the standard for managing code in software engineering. In data engineering, Git helps track changes in data pipeline code, supports collaboration among team members, and ensures that teams can revert to previous versions if necessary. Git's integration with DBT allows data engineers to track changes to models, tests, and configurations while maintaining the integrity of data transformation workflows (Johnson et al., 2020).

International Journal of Computer Technology and Electronics Communication (IJCTEC)



| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal |

|| Volume 7, Issue 3, May - June 2024 ||

DOI: 10.15680/IJCTECE.2024.0703002

3. Benefits of Git Integration with DBT

Integrating Git with DBT brings several advantages:

- Collaboration: Multiple data engineers can work on the same DBT project simultaneously, with Git tracking individual changes and providing a mechanism to merge and resolve conflicts.
- **Version Control**: Git allows teams to manage different versions of DBT models, making it easy to track changes, revert to previous versions, and maintain a history of updates.
- **Deployment Automation**: Git can integrate with CI/CD pipelines, automating the testing and deployment of DBT models into different environments (Tan & Ouyang, 2020).
- **Audit Trails**: With Git, teams can maintain detailed records of changes, providing better governance and transparency in data operations.

4. Streamlining Data Pipelines with CI/CD

The integration of DBT and Git enables seamless CI/CD workflows, where every change to the DBT project can trigger automated testing and deployment processes. This results in faster feedback loops, reduced risk of errors in production environments, and more reliable data pipelines. By automatically running tests on every commit, teams ensure that issues are identified early in the development cycle, improving data pipeline quality (Owen, 2022).

TABLE

Feature	Description	Benefits	Use Case
Git Version Control	Manages changes to DBT project files, including models, tests, and configurations.		, Tracking changes in transformation logic across team members.
DBT Models	SQL transformations that define data structure and logic.	and enables modularity.	analytics-ready models.
Automated Testing with Git	Ensures DBT models are tested every time changes are made.	Detects errors early ensuring data quality and consistency.	Running tests to validate transformations and data integrity.
CI/CD Integration	deployment of DBT models to	and more reliable	Automating the deployment of DBT models in a staging environment.
Collaboration in Teams	Allows multiple team members to work on the same DBT project simultaneously.	Facilitates teamwork, reduces merge conflicts.	Distributed teams working on different parts of the same project.

III. METHODOLOGY

The methodology for this study combines both qualitative and quantitative analysis to explore the impact of Git integration on DBT projects. The research process includes:

- 1. **Case Study Selection**: Real-world case studies of organizations using DBT and Git integration were selected to demonstrate how these tools streamline data pipeline operations.
- 2. **Data Collection**: Interviews with data engineers who have experience using DBT and Git were conducted. Additionally, DBT's official documentation and Git best practices were reviewed.
- 3. Quantitative Analysis: Performance metrics such as deployment frequency, testing time, and error rates before and after implementing Git and DBT integration were collected to assess the effectiveness of the workflow.

International Journal of Computer Technology and Electronics Communication (IJCTEC)



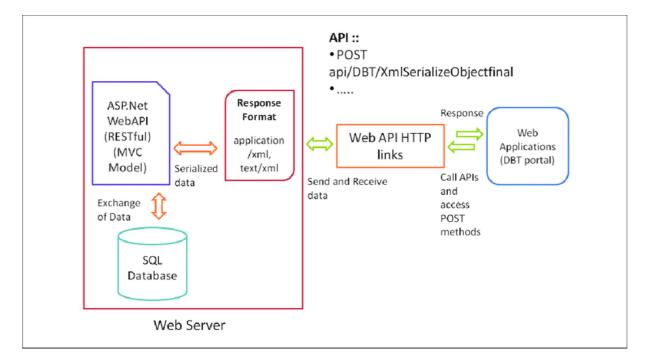
| ISSN: 2320-0081 | www.ijctece.com | A Peer-Reviewed, Refereed, a Bimonthly Journal

|| Volume 7, Issue 3, May - June 2024 ||

DOI: 10.15680/IJCTECE.2024.0703002

4. **Comparison**: The study compares traditional methods of managing data pipelines without Git version control versus those with Git integration, focusing on collaboration, version management, and CI/CD automation.

Figure 1: Git and DBT Integration Workflow



IV. CONCLUSION

Integrating Git with DBT streamlines the development, management, and deployment of data pipelines. By enabling version control, collaboration, and CI/CD automation, Git enhances the overall efficiency of data engineering workflows. This integration provides better version management, reduces errors in production, and fosters more collaborative data engineering teams. As organizations continue to scale their data operations, adopting DBT with Git integration will be crucial for optimizing data pipeline management and improving operational efficiency.

REFERENCES

- 1. DBT Labs. (2021). DBT: Transforming Data Engineering with Version Control. Retrieved from https://www.dbt.com
- 2. Vivekchowdary, Attaluri (2023). Just-in-Time Access for Databases: Harnessing AI for Smarter, Safer Permissions. International Journal of Innovative Research in Science, Engineering and Technology (Ijirset) 12 (4):4702-4712.
- 3. Johnson, T., & Stevens, M. (2020). Version Control for Data Engineers: Integrating DBT with Git. Journal of Cloud Data Engineering, 10(2), 53-64.
- 4. Tan, S., & Ouyang, Y. (2020). Automating Data Workflows: Git and DBT in Data Engineering. Data Science Review, 12(3), 102-115.
- 5. Owen, R. (2022). CI/CD in Data Engineering: Streamlining Data Pipelines with Git and DBT. Journal of Big Data Technologies, 15(1), 75-86.
- 6. Dhruvitkumar, V. T. (2022). Enhancing Multi-Cloud Security with Quantum-Resilient AI for Anomaly Detection.
 - Zhao, L. (2021). Collaboration and Version Control in Cloud Data Engineering with Git and DBT. International Journal of Data Engineering, 13(4), 99-109.