

# Integration of Dbt With Modern Data Stack Technologies

Rameshbabu Lakshmanasamy<sup>1</sup>, Girish Ganachari<sup>2</sup>

<sup>1</sup>Senior Data Engineer, Jewelers Mutual Group

<sup>2</sup>Senior streaming data engineer, Amazon Web Services

## Abstract:

In the current world where information is critical, the timely processing and analysis of data are essential in the decision-making process. Modern data stack is based on the cloud-native, modular, and scalable technology stack that makes the workflows more efficient and flexible. The dbt (Data Build Tool) is one of the critical components in the ecosystem working on SQL level transformations for the cloud data warehouses such as Snowflake and BigQuery. When compared to the other traditional ETL tools, dbt focuses more on the transformation layer, making the tasks easier for data analysts and engineers. This paper discusses how dbt fits into the modern data stack, how it differs from ETL, what tools it complements or interfaces with (like Airflow, Spark, and Kafka), and its implications for data governance and democratization.

**Keywords:** Dbt, Snowflake, Bigquery, Spark, Kafka, ETL

## Introduction:

Traditional ETL (Extract, Load, Transform) tools are those types of tools that extract data from the source systems, clean the data, transform it into the required form, and, at last, load it into the target destination (Zagni, 2023). The most common are Informatica, Talend, and SSIS. These tools are built on elaborated procedures dealing with data processing before it gets to an ultimate storage space. Conventional ETL instruments are still helpful for on-premises information architectures, but they need adaptation when used on cloud data platforms. Loading data after preparation takes some time, which can take more time when working in large data sets, as data has to be shifted from one system to another. Furthermore, the complexity of the ETL structure deems the frameworks brittle and creates a high operation cost concerning changing the settings in cloud-based data landscapes.

On the other hand, dbt (Data Build Tool) uses transformation only mechanism under ELT (Extract, Load, Transform) model. Unlike traditional ETL tools that modify data before it gets to the desired destination, dbt only makes that change when data is already stored in the cloud data warehouses such as Snowflake, Google BigQuery, and Redshift (Cyr & Dorsey, 2022). This shift takes advantage of computing capabilities available in cloud platforms to effect transformations where the data resides, thereby minimizing the movement of clicks of data around the ecosystem of tools: dbt's greatest strength, therefore, is the ability to apply transformations through SQL, which the data analyst can understand without having to be a skilled software developer. Moreover, since dbt transforms the data after loading, the latency issue of traditional ETL tools is minimized, making insights more current.

### **dbt in Cloud Data Warehouses**

Cloud data warehouses like Snowflake, Google BigQuery, Amazon Redshift, and Databricks are highly scalable platforms that enable efficient storage and retrieval of large datasets. They offer fast, nearly real-time processing for heavy analytical workloads and are billed based on usage with elastically scalable capacity. Unlike on-premise databases, these platforms integrate seamlessly with other data tools.

dbt (Data Build Tool) leverages the computing power of cloud data warehouses to perform in-place SQL-based transformations, reducing latency and simplifying operations. Instead of transferring data to external engines, dbt applies transformations directly within the warehouse, enhancing efficiency (Cyr & Dorsey, 2022). For example, in Snowflake, dbt compiles models into optimized SQL queries that Snowflake executes within its multi-cluster, shared-data environment. This allows organizations to manage large-scale data transformations easily, scaling compute resources as needed.

A retail company using dbt with Snowflake can process terabytes of data daily and perform transformations such as aggregations, joins, and filters directly in the warehouse, leading to fast and real-time business insights. This setup optimizes data operations, allowing teams to derive valuable insights efficiently.

### **Exploring Synergies Between dbt and Tools Like Airflow, Spark, Kafka**

#### **dbt and Airflow:**

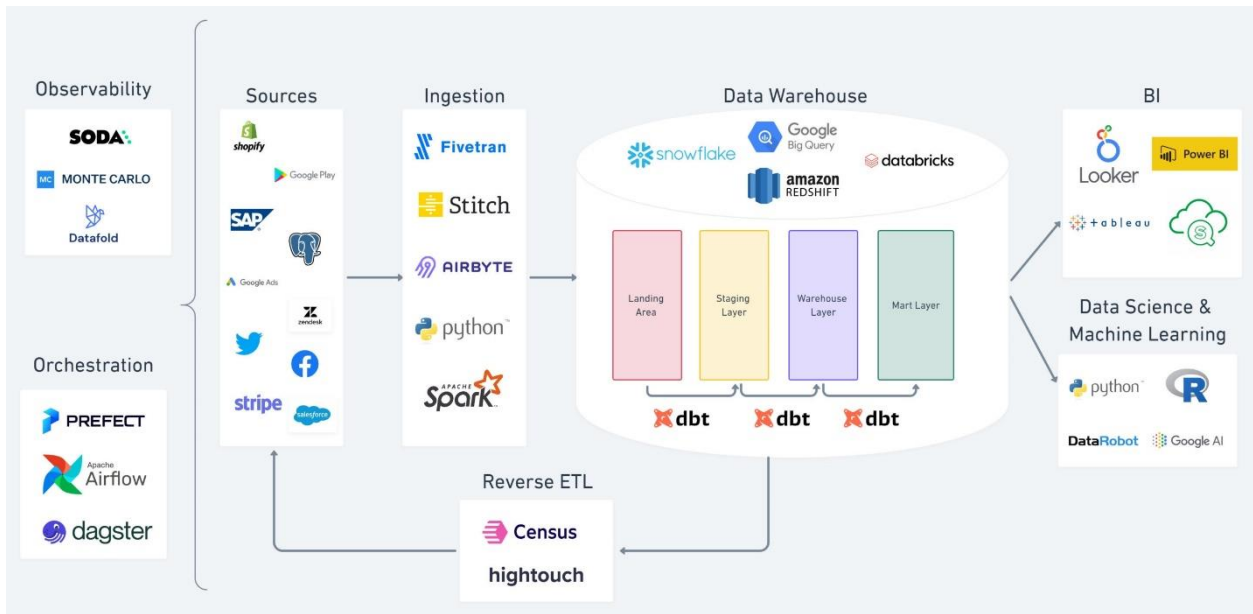
Airflow is a level of coordination that schedules and manages data pipeline tasks or jobs in a workflow. It works alongside dbt by managing the timing and the running of your dbt models in the context of a more extensive data process. According to Kristien et al., (2020) this approach provides an opportunity to integrate dbt runs as tasks in the Airflow DAGs (Directed Acyclic Graphs) to schedule dbt transformations along with other ETL or external processes. For instance, if we have an Airflow DAG sub tag- we can extract data from the source system, load that into the cloud DW- and then call for dbt transformations to execute immediately after load- hence the ETL and T are end-to-end, integrated with Airflow's DAGs.

#### **dbt and Spark:**

While dbt is used for transformations on the SQL level, Apache Spark is used for big data processing. For organizations that are dealing with process-intensive big data in terms of structured data or any other type of data requiring pre-processing prior to processing, Spark may be selected as the data processing tool (Densmore, 2021). Yet, dbt can still be used as the transformation layer, so that dbt models then call Spark to do bulking workloads on data. When both are required, dbt is applied for manageable SQL transformations, while Spark kicks in where the volume or sophistication of the transformation overloads dbt.

#### **dbt and Kafka:**

Kafka also shines in real-time data streaming, where Kafka and similar systems take data and move it through other systems. However, dbt works in batch mode, but it can be connected to downstream Kafka pipelines for near real-time analysis. It is possible for organizations to build hybrid streaming data pipelines to manage real-time ingestion using Kafka streams, and perform transformation and normalization on data stored in cloud warehouses after ingesting using dbt transformations (Cyr & Dorsey, 2023).



(Densmore, 2021).

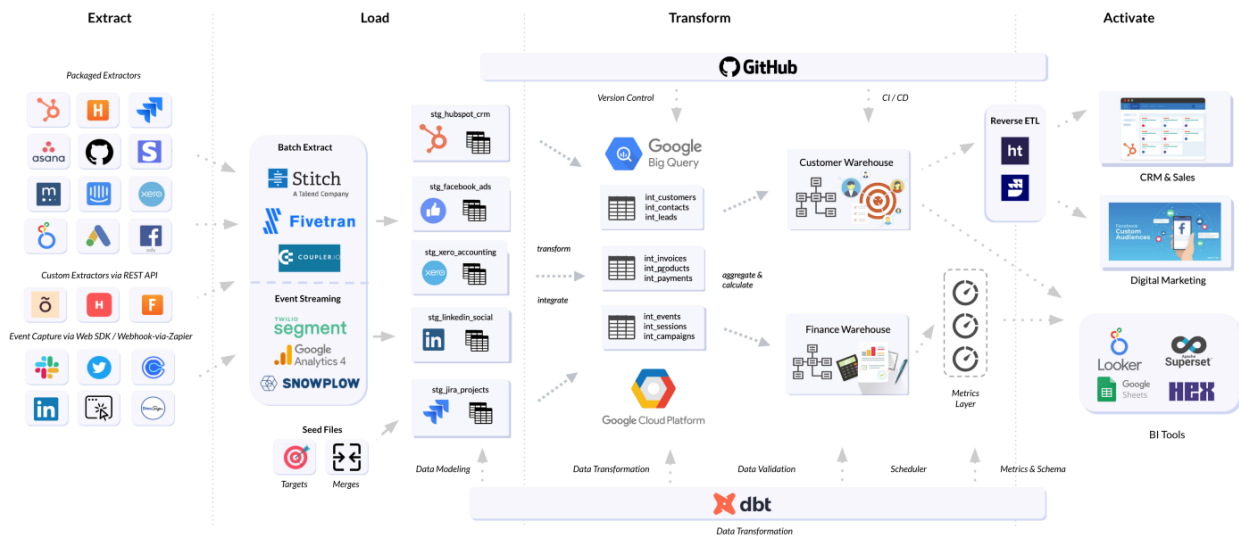
### Best Practices for Building a Cohesive Data Pipeline with dbt

Essentially, working with dbt is all about following the best practices that ensure scalability, maintainability, and reuse while constructing a solid data pipeline. One of the approaches is the use of modular models that entail the segmentation of transformations into smaller manageable parts (Cyr & Dorsey, 2023). With this strategy in place, teams are able to effectively introduce the management of big projects, which helps avoid inconsistency within different projects.

Testing and documentation can be seen as a means of improving data quality: dbt has built-in testing ability, which can be used to define tests that validate the correctness of data transformations; auto-documentation makes sure that data transformations are well-described and easily understandable.

For version control, dbt works with git-based workflows; the teams get a chance to track the changes in the code and put into practice CI/CD principles in their work (Nawaz, 2024). This would make it easy to manage the updates more systematically, hence minimizing some of the errors that are experienced; this, in turn, will significantly add to the overall time required in development cycles.

dbt packages Finally, the usage of iterative dbt packages speeds up the work and provides a kind of universal format for exchanging tasks and procedures, thus reducing considerably the amount of code to write and the number of discrepancies that may arise between projects. Moreover, dbt has the auto-documentation that improves data lineage, this certainly makes it easier for teams to monitor transformations, and also effective for data governance when handling problems. In combination, these practices facilitate a robust, elastic, and adequately documented data stack with dbt while increasing the workflow quality.



(Dash & Swayamsiddha, 2022).

### Adoption Challenges and Strategies, Security Implications of dbt in Data Workflows

Implementing DBT can be quite challenging, mainly for organizations that already have ETL structures in place and need to be updated. The considerable amount of technical debt from traditional ETL tools becomes an issue; often, existing data workflows may become complex and require extensive time and effort to refactor during migration to dbt (Cyr & Dorsey, 2023). Further, migrating from ETL to ELT could be challenging to manage since it alters the scope of teams both in terms of technique and mentality. This puts a need for SQL skills as a barrier for teams who are used to GUI-based ETL tools or even other programming languages.

Companies can help facilitate this by running training sessions and workshops that will allow professionals to learn SQL and dbt, as well as change their ways of thinking about modeling data in a modular fashion and using ELT processes (Densmore, 2021). This is because, in the empirical migration strategy as adopted here, the implementation of dbt starts with relatively less complicated pipelines before progressing to a more sophisticated one.

For security, DBT brings important options that are necessary, for instance, Role-Based Access Control (RBAC) for permission management. It also offers data encryption at rest and row level security to ensure the accomplishment of data security (Densmore, 2021). For compliance with the data governance standards, dbt offers documentation and lineage capabilities to meet the purpose in a structured manner to address compliance nature and improve on it.

### Investigating the use of dbt in Near Real-Time Data Processing

dbt (Data Build Tool) primarily supports batch data processing, suitable for scheduled data transformations based on predefined intervals (e.g., daily or hourly). While effective for cloud data warehouses, dbt's batch processing architecture limits its ability to provide quasi-real-time analysis required by organizations for timely insights. Its reliance on scheduled runs, often managed by tools like Airflow or dbt Cloud's Scheduler, means transformations occur after set intervals, not immediately upon data arrival (Dash & Swayamsiddha, 2022).

To address latency issues, dbt offers features like dbt Cloud for faster transformations and incremental materializations, which process only new data, reducing update time. However, dbt is not designed for

real-time transformations. To bridge this gap, organizations can integrate dbt with streaming technologies like Apache Kafka and Apache Flink. Kafka manages real-time data fetching and streaming, while Flink processes data in near-real-time. Once processed, the data is loaded into a cloud data warehouse for dbt's batch processing, maintaining its transformation, testing, and documentation capabilities (Dash & Swayamsiddha, 2022).

This hybrid approach allows organizations to leverage real-time streaming for immediate analysis while utilizing dbt's robust features for downstream processing, combining the best of both batch and streaming architectures.

### **dbt and Data Governance**

dbt (Data Build Tool) supports data governance by ensuring transparency, accuracy, and accountability in data processes. Its Testing feature allows teams to validate changes using data quality checks (e.g., NULL values, uniqueness, and data integrity) within dbt models (Auliya., 2022). This proactive testing approach helps identify errors before they impact analysis, ensuring data meets required standards.

Data lineage in dbt tracks and visualizes transformations, versioning every change to trace data from raw sources to analytical outputs. This end-to-end tracking is crucial for monitoring data lifecycle, resolving issues, and assigning responsibility, especially in regulated industries like healthcare and finance (Dash & Swayamsiddha, 2022). For example, under HIPAA regulations, dbt can ensure only de-identified and encrypted data passes through transformations, maintaining accuracy and security. Similarly, for GDPR and CCPA compliance, dbt provides detailed documentation and lineage to substantiate data history, access requests, and processing activities.

For broader governance, dbt integrates with cataloging tools like Alation and Collibra, which provide data dictionaries and administrative frameworks. Combining dbt's documentation and lineage capabilities with these tools helps organizations achieve a comprehensive, compliant, and well-governed data environment, supporting transparency and regulatory adherence.

### **The Role of dbt in Promoting Data Democratization**

Data democratization is the process by which data is shared so that everybody in an organization enjoys it and can make decisions for themselves without consulting data experts. In modern organizations, this helps not only data engineers but analysts, product managers, and even business stakeholders collect and analyze data independently, building new innovative processes and making better decisions. DBT (Data Build Tool) is crucial in data democratization since it allows data analysts and engineers to perform data transformation using their favorite SQL (Auliya., 2022). This shift to an analyst's hand helps decentralize and democratize data models and pipelines so that there is less reliance on centralized data engineering teams and more independence and flexibility in data management.

Initially, data engineering groups were accountable for every aspect of the pipeline. However, with dbt, there is the distributed ownership concept because business teams that are in touch with the data demand can ownership correct the models. This promotes alternatives and quicker response to calls from businesses. dbt also enables self-service analytics where querying and visualizations are collaborative (Dash & Swayamsiddha, 2022). Extensive testing, versioning, and documentation are supported with additional features to prove that data pipelines are correctly maintained and easily explainable. In a broader perspective, it helps establish the actual use of data, which helps create a culture of data democratization.



**Conclusion:**

Dbt (Data Build Tool) is currently one of the most popular instruments for working with data in the modern data stack. This SQL approach empowers data analysts and engineers to assume control over the transformation process as it decentralizes data transformation from multiple centralized data transformation teams (Auliya., 2022). In doing so, dbt aligns well with the current data requirements of fast-paced organizations since it enhances post-load transformations in cloud data warehouses for data pipeline scalability.

With support for cloud data platforms such as Snowflake, BigQuery, and Redshift and support for orchestration systems such as Airflow and real-time systems like Kafka, dbt provides the means for organizations to create scalable and integrated data systems. Due to the materialized testing, version control, and documentation options to guarantee data quality, governance, and transparency, it enhances operational execution and compliance with regulations.

**References:**

1. Zagni, R. (2023). Data Engineering with dbt: A practical guide to building a cloud-based, pragmatic, and dependable data platform with SQL. Packt Publishing Ltd. [https://books.google.com/books?hl=en&lr=&id=okLJEAAQBAJ&oi=fnd&pg=PP1&dq=Integration+of+dbt+with+Modern+Data+Stack+Technologies&ots=3eDxg7wAYS&sig=rCdF-SJFHQED\\_FVEp1NjnLb61E4](https://books.google.com/books?hl=en&lr=&id=okLJEAAQBAJ&oi=fnd&pg=PP1&dq=Integration+of+dbt+with+Modern+Data+Stack+Technologies&ots=3eDxg7wAYS&sig=rCdF-SJFHQED_FVEp1NjnLb61E4)
2. Cyr, C., & Dorsey, D. (2022). Unlocking dbt. <https://link.springer.com/content/pdf/10.1007/978-1-4842-9703-2.pdf>
3. Cyr, C., & Dorsey, D. (2023). Introduction to dbt. In *Unlocking dbt: Design and Deploy Transformations in Your Cloud Data Warehouse* (pp. 1-40). Berkeley, CA: Apress. [https://link.springer.com/chapter/10.1007/978-1-4842-9703-2\\_1](https://link.springer.com/chapter/10.1007/978-1-4842-9703-2_1)
4. Auliya, M., Aziz, A., & Nurharjadmo, W. (2022, December). Big Data Modern Stack for District Government. In *1st International Conference on Demographics and Civil-registration (INCODEC 2021)* (pp. 103-109). Atlantis Press. <https://www.atlantis-press.com/proceedings/incodec-21/125977931>
5. Kristien, M., Spink, T., Campbell, B., Sarkar, S., Stark, I., Franke, B., ... & Topham, N. (2020). Fast and correct load-link/store-conditional instruction handling in DBT systems. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 39(11), 3544-3554. <https://ieeexplore.ieee.org/abstract/document/9211550>
6. Dash, B., & Swayamsiddha, S. (2022, December). Reverse etl for improved scalability, observability, and performance of modern operational analytics-a comparative review. In *2022 OITS International Conference on Information Technology (OCIT)* (pp. 491-494). IEEE. <https://ieeexplore.ieee.org/abstract/document/10053738>
7. Densmore, J. (2021). Data pipelines pocket reference. O'Reilly Media. <https://books.google.com/books?hl=en&lr=&id=SxgcEAAQBAJ&oi=fnd&pg=PR2&dq=Integration+of+dbt+with+Modern+Data+Stack+Technologies+and+traditional+ETL&ots=vFxEWdrADv&sig=Hs48SfE8BdVqtYKpoo2chk33CKo>