Open in app ↗

# Medium     🔍 Search                                    🔔  👤

# Josef Machytka: Speaker Portfolio — PostreSQL, DuckDB, Databases in General

**Josef Machytka**

8 min read  ·  Nov 3, 2024

▶ Listen        ⬆ Share        ••• More

*Expert Talks on PostgreSQL, DuckDB, Databases in general, Data LakeHouse, Data Ingestion and Data Analysis*

## Introduction

With over 30 years of experience in database technologies, I have cultivated a deep understanding of diverse database systems, including extensive work with PostgreSQL. I worked with multiple heterogeneous data ingestion and data processing pipelines, focusing on delivery of valuable business results. The need to solve often very complicated technical problems equipped me with the expertise necessary for insightful and impactful talks. I am passionate about sharing my knowledge, empowering others in the field of databases and data processing and making complex topics accessible.

My talks are living and evolving, I steadily improve them with new content reflecting progress in new versions of the corresponding software. My speaking style combines practical insights with important technical details, ensuring that audiences walk away with really valuable knowledge they can apply. Slides contain a lot of additional information, making them valuable resources even outside the context of my presentations.

*Thank you for considering me as a speaker for your next event. I look forward to the opportunity to share my knowledge and help your audience gain deeper insights into PostgreSQL and database topics.*

## About Me

**Name:** Josef Machytka

**My current job:** Professional Service Consultant, PostgreSQL specialist at NetApp Open Source Systems

**Experience:**
* 30+ years of production experience with different databases:
PostgreSQL 13 years, BigQuery 7y, Oracle 15y, MySQL 12y,
Elasticsearch 5y, MS SQL 5y, Sybase ASE, FoxPro
* 10+ years of experience with high volume and velocity data ingestion pipelines,
Data Analysis, Data Warehouse and Data Lakehouse architecture
* 2+ years of practical experience with different LLMs, their architecture and principles

**Accounts:** LinkedIn, ResearchGate.net, Academia.edu
**Mastodon:** @JosefMachytka@me.dm

## Current Talk Portfolio

- *GIN, BTREE_GIN, GIST, BTREE_GIST, HASH and BTREE indexes on JSONB data*

- *PostgreSQL and DuckDB: Supercharging ad-hoc Data Analysis and ETL*

- *Building a Data Lakehouse with PostgreSQL: Dive into Formats, Tools, Techniques, and Strategies*

- *PostgreSQL Connections Memory Usage: How Much, Why, and When?*

- *Partitioning and Clustering: An Overview of Solutions with a deep dive into PostgreSQL implementation*

See my latest speaking engagements at the end of this article.

# GIN, BTREE_GIN, GIST, BTREE_GIST, HASH and BTREE indexes on JSONB data

**Duration:** 45 minutes

**Target Audience:** Application developers, data analysts

**Overview:** Talk summarizes several months long and still ongoing internal project testing usage and performance of GIN and BTREE_GIN with different operator classes, GIST and BTREE_GIST indexes for GeoJSON data and also standard HASH and BTREE indexes specifically on JSONB data. Tested on several real life datasets with a total size of dozens of GBs. Also, the influence of TOAST compression algorithms, parallelism, memory settings, table statistics on processing JSONB data was tested. Objective of this project was to gather relevant experience to be able to help our customers with their problems, because the majority of articles on the web about JSONB data in PostgreSQL show only trivial examples without any reasonable value for developers solving multiple performance issues related to JSONB data. The talk also discusses practical limitations developers would face if they try to fully decompose JSONB data into relational tables.

**Key Takeaways:**

- Understanding of use cases and performance of different types of indexes for JSONB data

- The impact of system settings like TOAST compression, parallelism, and memory on performance and usage of indexes on JSONB data

- Limitations and considerations for decomposing JSONB data into relational structures

- Insights into internal structure of different types of indexes

**Slides:** academia.edu

**Presented at:**

- Prague PostgreSQL Developer Day 2024 (article on NetApp blog)

- Swiss PG day 2024 (article on NetApp blog)

- Berlin PostgreSQL MeetUp October 2024 (MeetUp entry)

# PostgreSQL and DuckDB: Supercharging ad-hoc Data Analysis and ETL

**Duration:** 45 minutes

**Target Audience:** Data Analysts, Data Scientists, App developers

**Overview:**

DuckDB is a powerful and practical tool, designed by data analysts in academia to address common use cases efficiently. It complements PostgreSQL workflows by streamlining ad-hoc data analysis, ETL processes, cross-database queries, and basic data migrations. DuckDB's lightweight yet highly efficient architecture eliminates the need for extensive setup or specialized tools, making it a valuable asset for simplifying complex tasks.

This talk provides actionable insights into integrating DuckDB with PostgreSQL to enhance everyday data workflows. Attendees will explore practical techniques, including the use of DuckDB extensions and Python scripts, to improve agility and productivity in data handling. The focus is on demonstrating DuckDB's ability to deliver simple and powerful solutions for common analytical challenges. Author regularly writes articles on medium.com about abilities of DuckDB.

**Key Takeaways:**

- Small Data Manifesto as a growing trend in data analysis

- Understanding of DuckDB's main features, strengths, and limitations

- Usage of DuckDB as a very efficient enhancement tool for PostgreSQL

- Optimization of analytical workload by combining capabilities of PostgreSQL and DuckDB

**Recording:** PostgreSQL MeetUp for All 2025.01.08

**Slides:** academia.edu

**Presented at:**

- Prague PostgreSQL MeetUp October 2024

- NetApp internal tech talk — 2024.11.26

- PostgreSQL MeetUp for All — 2025.01.08



**Reactions:**

> *"Thank you so much for your talk and participation in our meetup! We're very honored and grateful to have you and your expertise. I've researched DuckDB a lot in the last year and you did a fantastic job of covering everything."*
> *(Elizabeth Garrett Christensen — Crunchy Data)*

## Building a Data Lakehouse with PostgreSQL: Dive into Formats, Tools, Techniques, and Strategies

**Duration:** 45 minutes
**Target Audience:** Application Developers, Database Administrators, Data Analysts

**Overview:**

The evolution of Data Warehouses, Data Lakes, and Data Lakehouses has been marked by many buzzwords, fluctuating trends, and tools that often over-promised but under-delivered. While there are numerous materials on these topics, most of them provide mostly introductory overviews and focus narrowly on a single

technology. And there are even many different opinions about what exactly is Data Lakehouse.

This talk discusses different ways how to understand this topic. It explores data formats and frameworks like Parquet, Apache Iceberg, Delta Lake, Apache Hudi. Discusses different architectures of Data Lakehouse solutions. Also key challenges will be addressed, such as effective Data Governance, compliance with privacy and security standards, and comprehensive data quality checks.

Last part of the talk address current AI hype with its many promises and proposes realistic overview of real capabilities of current Large Language Models and their use cases in Data Lakehouses.

PostgreSQL is extremely well equipped to play a major role in the current Data Lakehouse and AI boom.

**Key Takeaways:**

- A comprehensive overview of Data Lakehouse architecture

- Insights into key data formats and frameworks in modern Data Lakehouses

- Practical ideas for implementing Data Governance practices

- Realistic view of real capabilities of current LLMs in scope of Data Lakehouses

**Slides:** academia.edu

**Presented at:**

- Prague PostgreSQL Developer Day 2025

- NetApp Internal Talk 2025.01.23

## PostgreSQL Connections Memory Usage: How Much, Why, and When?

This talk explores the memory usage of PostgreSQL connections on Debian/Ubuntu running on x86–64 architecture. It gives an overview of memory management concepts, explaining key metrics like virtual, resident, and proportional memory sizes. It also covers various Linux tools for displaying memory usage.

The second part presents practical measurements of PostgreSQL memory usage, based on data from /proc/PID/smaps. It explains why RSS numbers for PostgreSQL connections appear so large after query execution and demonstrates that the actual unique memory usage is only a few dozen megabytes.

Finally, the talk examines how PostgreSQL connections allocate and release additional memory during query processing. It also clarifies where work_mem fits into these numbers and visualizes the process with plots.

**Key Takeaways:**

- The large RSS values in long-running sessions mostly come from linked shared_buffers

- A newly created connection consumes only up to 10 MB of physical memory, independent of the work_mem setting

- Additional memory is allocated and later released as queries execute

- Work_mem is a "soft maximum limit" — it may not be fully used, but it can also be exceeded

**Slides:** not available online yet

**Presented at:**

- NetApp internal talk 2025.01.27

## Partitioning and Clustering: An Overview of Solutions with a deep dive into PostgreSQL implementation

**Duration:** 45 minutes
**Target Audience:** App developers, system architects

**Overview:** In this presentation, we will examine the implementation of partitioning and clustering in several database systems, such as BigQuery, Snowflake, Oracle and MySQL. Following that, we will discuss a detailed analysis of PostgreSQL's approach to inheritance, partitioning, and clustering. We will check database parameters that

affect performance of these solutions, including rarely used enable_partitionwise_aggregate and enable_partitionwise_join parameters, will compare the results of performance tests between a single large table and partitioned tables on different datasets, look at efficiency of indexes, and discuss the application of multi-level partitioning for different use cases. Additionally, we will delve into details of memory usage, statistics and query optimization. We will also share some tips and tricks — how to subsequently implement partitioning on existing table, how to create child tables with additional columns.

**Key Takeaways:**

- Understanding of partitioning & clustering across different database systems

- In-depth understanding of PostgreSQL implementation

- Practical tips and tricks for performance optimization and implementation of partitions

**Slides:** not available online

**Presented at:**

- NetApp internal workshop 2024.11.12

## My latest speaking engagements

1. **2025.01.29 Prague PostgreSQL Developer Day 2025**
   "Building a Data Lakehouse with PostgreSQL: Dive into Formats, Tools, Techniques, and Strategies"

2. **2025.01.27 NetApp internal talk**
   "PostgreSQL Connections Memory Usage: How Much, Why, and When?"

3. **2025.01.23 NetApp internal talk**
   "Building a Data Lakehouse with PostgreSQL: Dive into Formats, Tools, Techniques, and Strategies"

4. **2025.01.08 PostgreSQL MeetUp for All online event**
   "PostgreSQL and DuckDB: Supercharging ad-hoc Data Analysis and ETL"

5. **2024.11.26 NetApp internal talk**
   "PostgreSQL and DuckDB: Supercharging ad-hoc Data Analysis and ETL"

6. **2024.11.15 NetApp Tech Talk**
   "Many facets of AI hallucinations: factual errors, deep fakes, and creativity"

7. **2024.11.12 Internal Workshop** — 3 talks:
   "PostgreSQL 17: New Features"
   "PostgreSQL Architecture: How PostgreSQL functions under the Hood"
   "PostgreSQL Partitioning: Overview of Solutions, Tips & Tricks, Performance"

8. **2024.10.29 Prague PostgreSQL MeetUp October 2024**
   "PostgreSQL and DuckDB"

9. **2024.10.16 NetApp internal workshop** — "AI Workshop: Exploring Artificial Intelligence" — 3 talks:
   - "The AI Dilemmas: The Philosophical, Ethical, and Legal Implications of AI"
   - "Understanding the Magic of LLMs: A Deep Dive into the Internal Structure of LLMs"
   - "Many facets of AI hallucinations: factual errors, deep fakes, and creativity"

10. **2024.10.09 PostgreSQL MeetUp Berlin October 2024**
    - "GIN, BTREE_GIN, GIST, BTREE_GIST, HASH and BTREE indexes on JSONB data"

11. **2024.09.27 NetApp internal talk**
    - "Beyond the Buzzwords: A pragmatic explanation of LLM terminology"

12. **2024.10.13 NetApp internal talk**
    - "The Art and Science of AI Prompt Engineering"

13. **2024.07.16 PostgreSQL MeetUp Berlin July 2024**
    - Lightning talk: "Can PostgreSQL have a more prominent role in the AI boom?"

14. **2024.06.27 Swiss PG day 2024**
    - "GIN, BTREE_GIN, GIST, BTREE_GIST, HASH and BTREE indexes on JSONB data"

15. **2024.06.05 Prague PostgreSQL Developer Day 2024**
    - "GIN, BTREE_GIN, GIST, BTREE_GIST, HASH and BTREE indexes on JSONB data"

## Photos from my talks



© Tomas Vondra P2D2 — Prague PostgreSQL Developer Day 2024 — talk about indexes on JSONB data

© organizers of Swiss PG day — Swiss PG day 2024 — talk about indexes on JSONB data

© Andreas Scherbaum — PostgreSQL Meetup Berlin July 2024 — lightning talk about PostgreSQL and AI

© Andreas Scherbaum — PostgreSQL Meetup Berlin Octobec 2024 — talk about indexes on JSONB data

© Igor Gavrilov (LinkedIn post) — Prague PostgreSQL Meetup October 2024 — talk about DuckDB

© organizers of Prague PostgreSQL Developer Day 2025

© organizers of Prague PostgreSQL Developer Day 2025

Postgresql          Data Analysis          Data Ingestion          Duckdb

Edit profile

# Written by Josef Machytka

61 Followers  ·  22 Following

I work as PostgreSQL specialist & database reliability engineer at NetApp Deutschland, Open Source Services division.

## No responses yet

What are your thoughts?

Respond

## More from Josef Machytka

Josef Machytka

## DuckDB Database File as a New Standard for Sharing Data?

This is not my original idea; I came across it in an excellent article titled "DuckDB Beyond the Hype" by Alireza Sadeghi. However, it...

Dec 30, 2024    👋 18    💬 2



Josef Machytka

## DuckDB Performance Problems with Inappropriate Pivoting Queries on Very Large Datasets

The DuckDB documentation clearly states that this tool is designed for handling datasets fitting into memory. I fully understand that I'm...

Jan 3      👋 6

```
Bob          2100.0    600.0
Charlie      2300.0   1500.0    1100.0
```

```
D pivot pg.sales on (product,year) using sum(sales_amount) group by salesperson order by salesperson;
```

| salesperson<br>varchar | (Laptop, 2022)<br>double | (Laptop, 2023)<br>double | (Phone, 2022)<br>double | (Phone, 2023)<br>double | (Tablet, 2022)<br>double | (Tablet, 2023)<br>double |
|---|---|---|---|---|---|---|
| Alice | 1200.0 | 1400.0 | 800.0 | 900.0 | 300.0 | 400.0 |
| Bob | 1000.0 | 1100.0 | 600.0 | | | |
| Charlie | 1100.0 | 1200.0 | 700.0 | 800.0 | 500.0 | 600.0 |

```
D pivot pg.sales on (year,product) using sum(sales_amount) group by salesperson order by salesperson;
```

| salesperson<br>varchar | (2022, Laptop)<br>double | (2022, Phone)<br>double | (2022, Tablet)<br>double | (2023, Laptop)<br>double | (2023, Phone)<br>double | (2023, Tablet)<br>double |
|---|---|---|---|---|---|---|
| Alice | 1200.0 | 800.0 | 300.0 | 1400.0 | 900.0 | 400.0 |
| Bob | 1000.0 | 600.0 | | 1100.0 | | |
| Charlie | 1100.0 | 700.0 | 500.0 | 1200.0 | 800.0 | 600.0 |

```
D pivot pg.sales on (year) using sum(sales_amount) group by salesperson order by salesperson;
```

```
salesperson    2022    2022
```

Josef Machytka

## Easy and Intelligent Pivot Tables with DuckDB

After exploring the various capabilities of DuckDB in my earlier articles, I want to focus more on its powerful data analytical...

Dec 4, 2024      👋 1



POSTGRES MEETUP FOR ALL
POSTGRES AND DUCKDB
Wednesday, January 8th
10 Pacific | 1 Eastern | 18:00 UTC
JOSEF MACHYTKA

![Josef Machytka] Josef Machytka

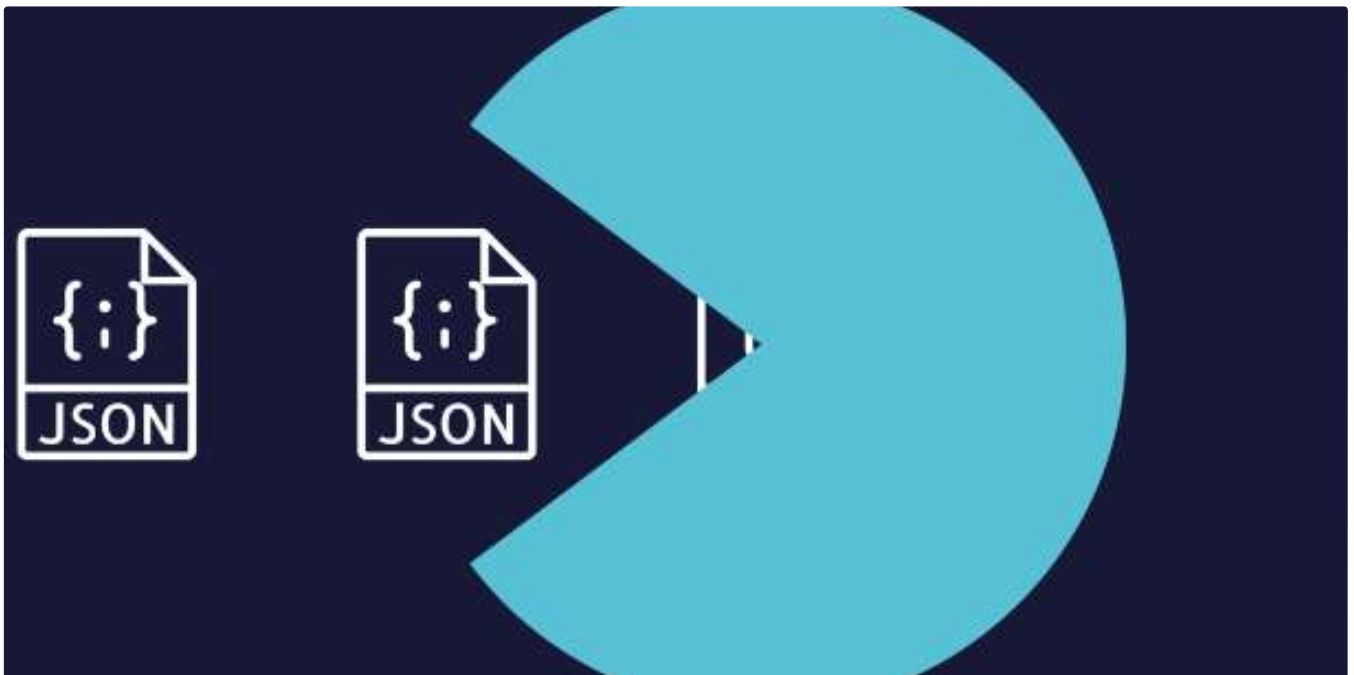# PostgreSQL and DuckDB: Supercharging Ad-Hoc Data Analysis and ETL

On Wednesday, January 8th, 2025, I had the amazing opportunity to present my talk about the DuckDB database online at the "Postgres MeetUp...

Jan 9    👏 14    💬 1                                                          🔖    •••

---

See all from Josef Machytka

---

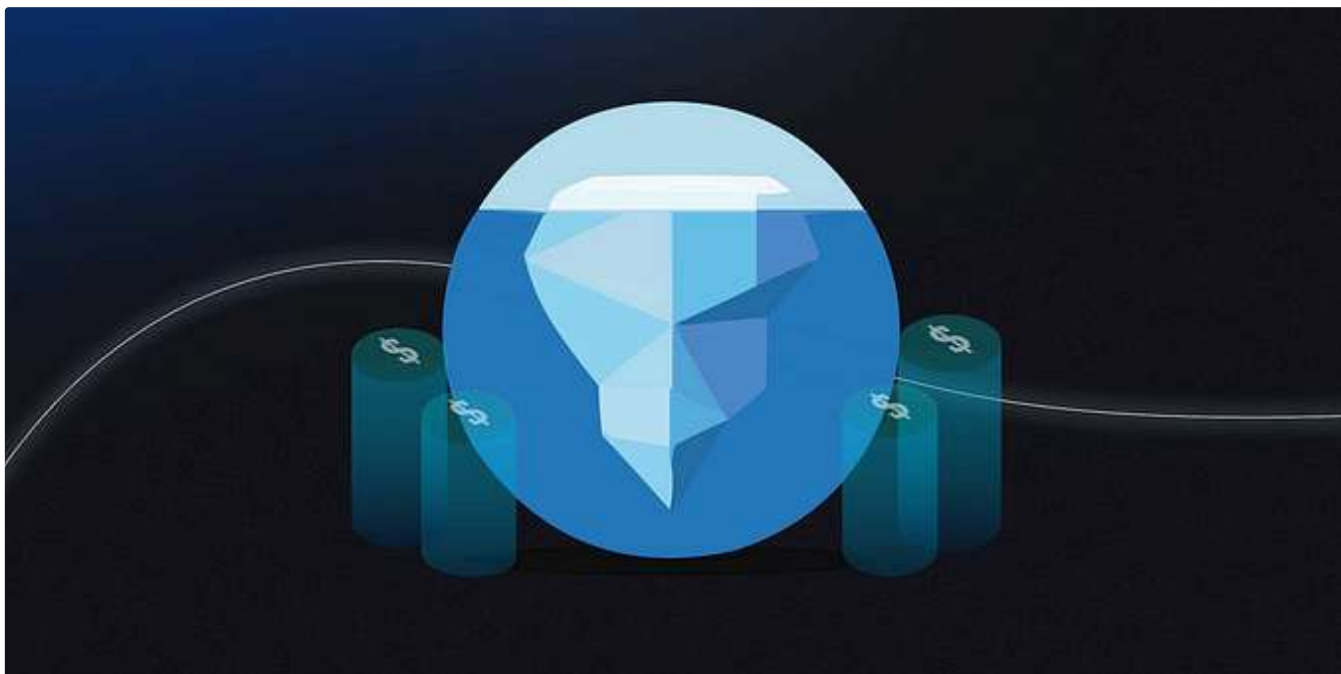# Recommended from Medium



![Jordi Puig] Jordi Puig

## DLT is Great!

Learn the basics of dlt with this simple tutorial and forget about the data engineering pains of extracting and loading data.

Oct 12, 2024    👏 85    💬 2                                                   🔖    •••

---

In Towards Dev by RisingWave Labs

## The True Power of Apache Iceberg: Revolutionizing Modern Data Architectures

Apache Iceberg is more than just a new technology; it's a paradigm shift in how we manage and utilize data.

4d ago

## Lists



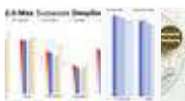**Practical Guides to Machine Learning**

10 stories · 2182 saves



**ChatGPT prompts**

51 stories · 2522 saves



**Staff picks**

806 stories · 1601 saves



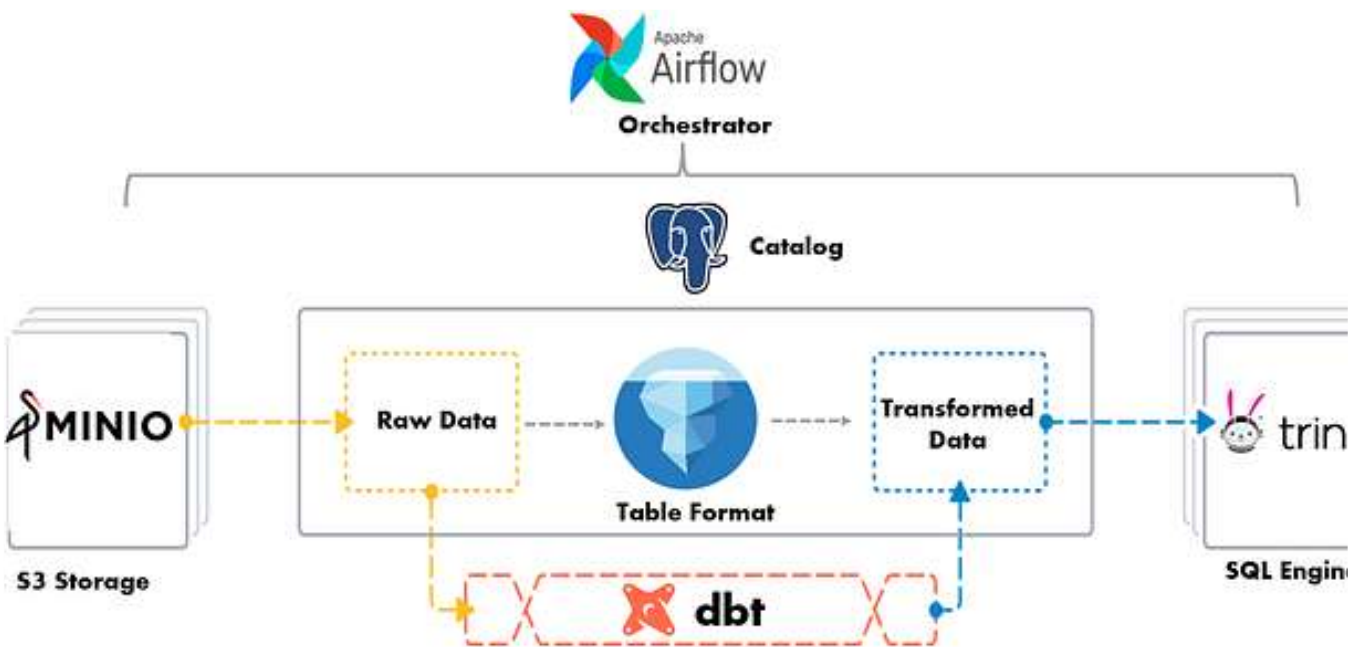**Natural Language Processing**

1908 stories · 1563 saves

In Level Up Coding by Anna Geller

## 2025 Data Engineering & AI Trends

How GenAI, new data regulations, Postgres, DuckDB, and open table formats affect data engineering in 2025 and beyond
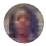
Jan 24    👏 216    💬 2

---



In Dev Genius by Haq Nawaz

## Building an End-to-End Data Lake ELT Pipeline using Modern Data Stack
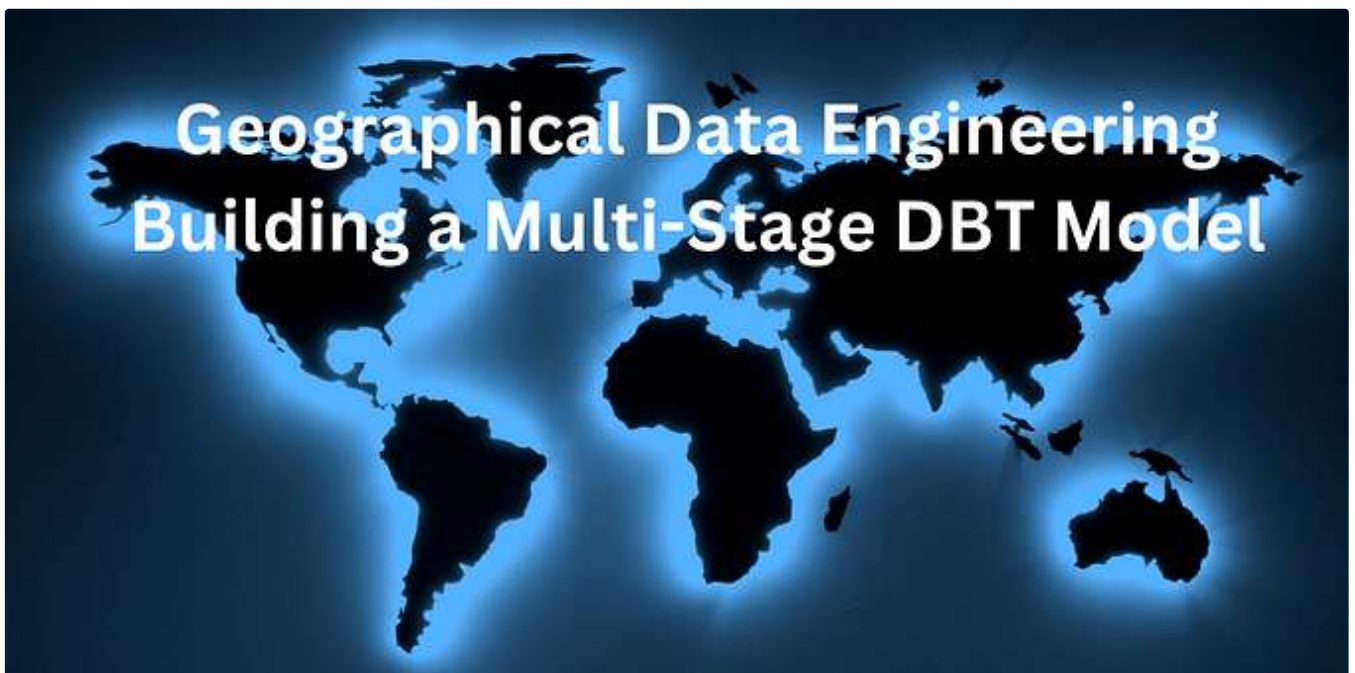
Using Airflow, dbt, Iceberg, MinIO, and Trino

👤 Nkem Onyemachi

## Data Load Tool (dlt) Series 2: Creating Pipelines from Different Sources

In the first article of this series,

Geographical Data Engineering
Building a Multi-Stage DBT Model

🗺️ Data Dev Backyard

## Geographical Data Engineering: Building a Multi-Stage DBT Model and Querying Using DuckDB

In the previous tutorials, we have introduced DBT, and discussed the core concepts behind DBT. In this tutorial, we would like to create...

✦       5d ago      👋 1                                                                                       🔖+            •••

─────────────────────────────────────────────────────────

( See more recommendations )