# Building Data Lakes on AWS with Kafka Connect, Debezium, Apicurio Registry, and Apache Hudi

Learn how to build a near real-time transactional data lake on AWS using a combination of Open Source Software (OSS) and AWS Services
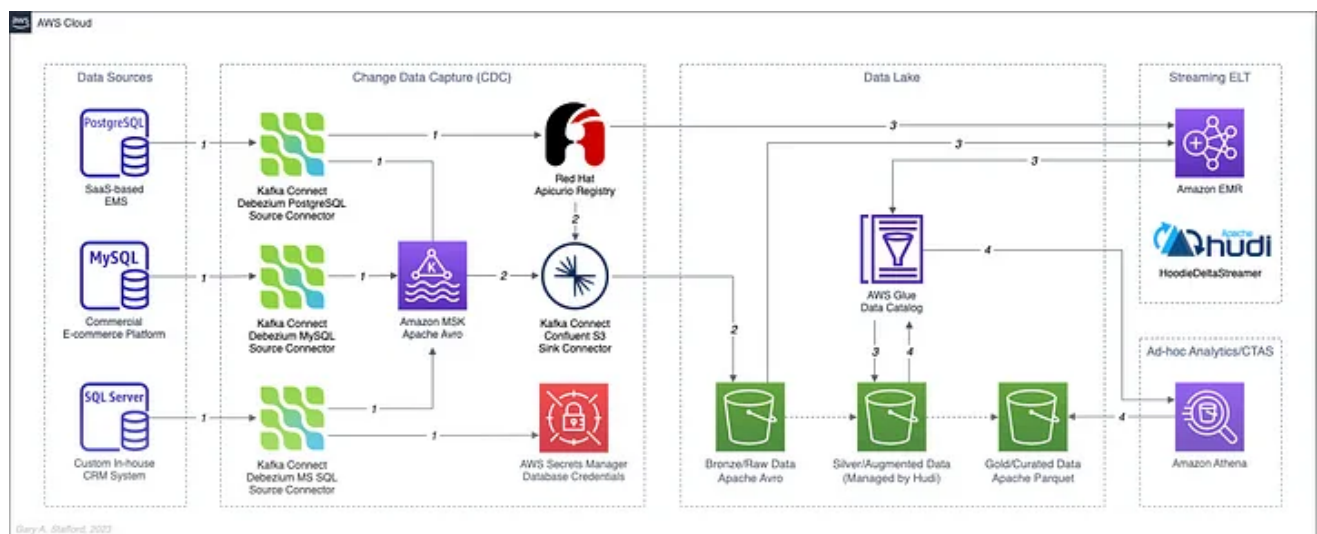
Gary A. Stafford · Follow

Published in ITNEXT

21 min read · Feb 28, 2023

## Introduction

In the following post, we will explore one possible architecture for building a near real-time transactional data lake on AWS. The data lake will be built using a combination of open source software (OSS) and fully-managed AWS services. Red Hat's Debezium, Apache Kafka, and Kafka Connect will be used for change data capture (CDC). In addition, Apache Spark, Apache Hudi, and Hudi's DeltaStreamer will be used to manage the data lake. To complete our architecture, we will use several fully-managed AWS services, including Amazon RDS, Amazon MKS, Amazon EKS, AWS Glue, and Amazon EMR.



The data lake architecture used in this post's demonstration

**Source Code**

The source code, configuration files, EMR Notebook, and a list of commands shown in this post are open-sourced and available on GitHub.

# Create an account to read the full story.

The author made this story available to Medium members only.
If you're new to Medium, create a new account to read this story on us.

Sign up with Google

Sign up with Facebook

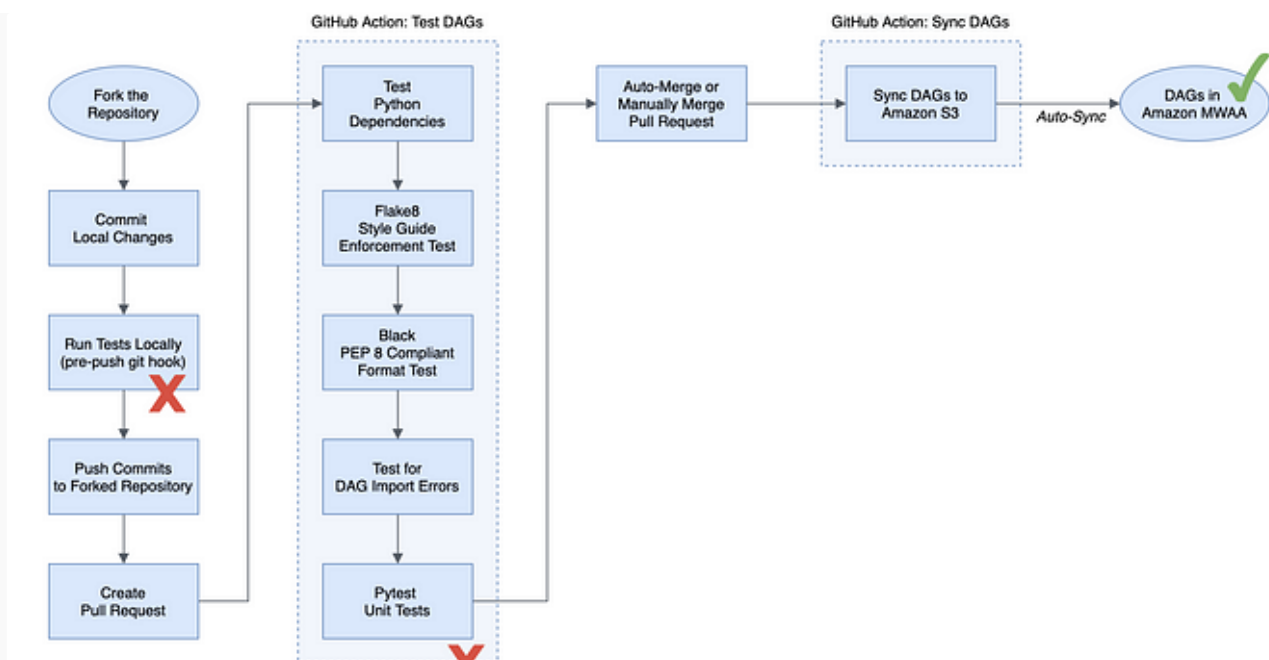Sign up with email

Already have an account? Sign in

## Written by Gary A. Stafford

4.8K Followers · Writer for ITNEXT

Area Principal Solutions Architect @ AWS | 10x AWS Certified Pro | Polyglot Developer | DataOps | GenAI | Technology consultant, writer, and speaker

**More from Gary A. Stafford and ITNEXT**

![Gary A. Stafford avatar] Gary A. Stafford

## DevOps for DataOps: Building a CI/CD Pipeline for Apache Airflow DAGs

Build an effective CI/CD pipeline to test and deploy your Apache Airflow DAGs to Amazon MWAA using GitHub Actions

✦　Dec 13, 2021　👋 165



![Niels Cautaerts avatar] Niels Cautaerts in ITNEXT

## Containers: has the pendulum swung too far?

Containerization has revolutionized the software industry, but using them blindly for everything without considering their drawbacks or...

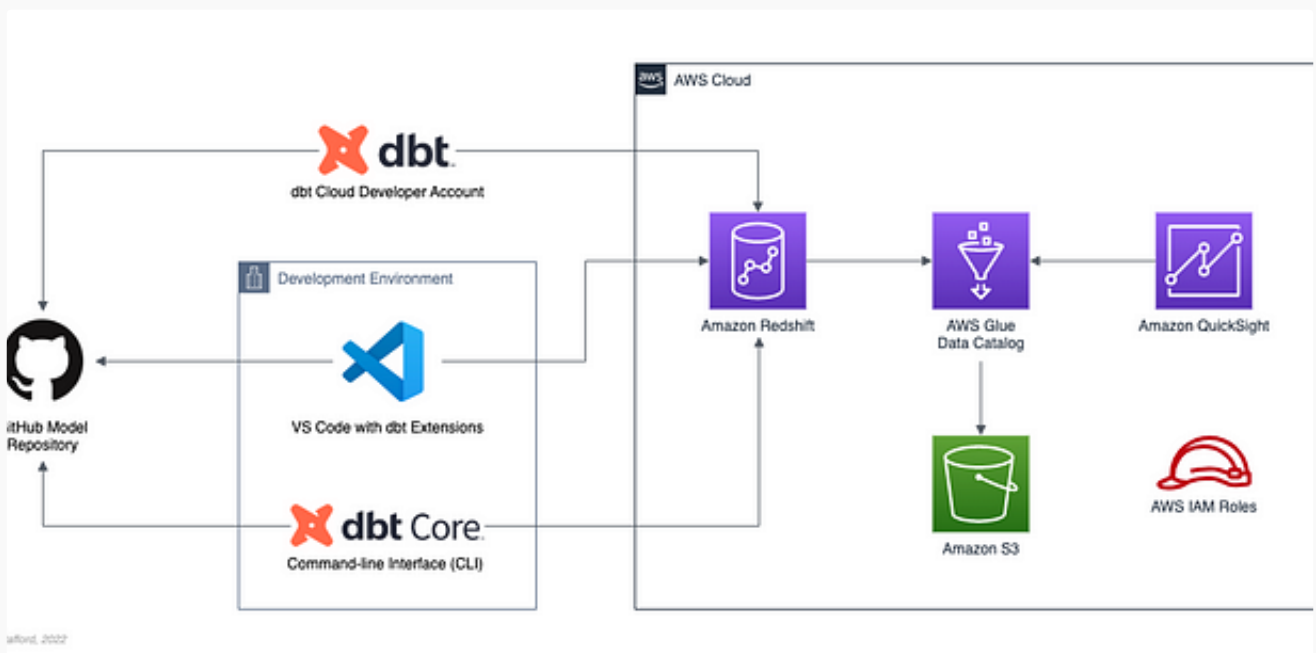👤 Mahdi Mallaki in ITNEXT

## Dockerfile-less and Daemon-less Build

Building a Docker image without requiring a Dockerfile or Docker Daemon

## Lakehouse Data Modeling using dbt, Amazon Redshift, Redshift Spectrum, and AWS Glue

Learn how dbt makes it easy to transform data and materialize models in a modern cloud data lakehouse built on AWS

✦　Aug 19, 2022　👋 197　💬 2　　　　　　　　　　　　　　🔖⁺

See all from Gary A. Stafford

See all from ITNEXT

## Recommended from Medium

## Building a Log Analysis Data Pipeline Using Kafka, Elasticsearch, Logstash, and Kibana — ELK Stack

When it comes to analyzing logs, having a real-time, centralized, and automated solution is a game changer. Instead of sifting through...

👤 Amit Kumar Manjhi

## Handling Schema Evolution in Debezium Kafka Connectors: A Step-by-Step Guide
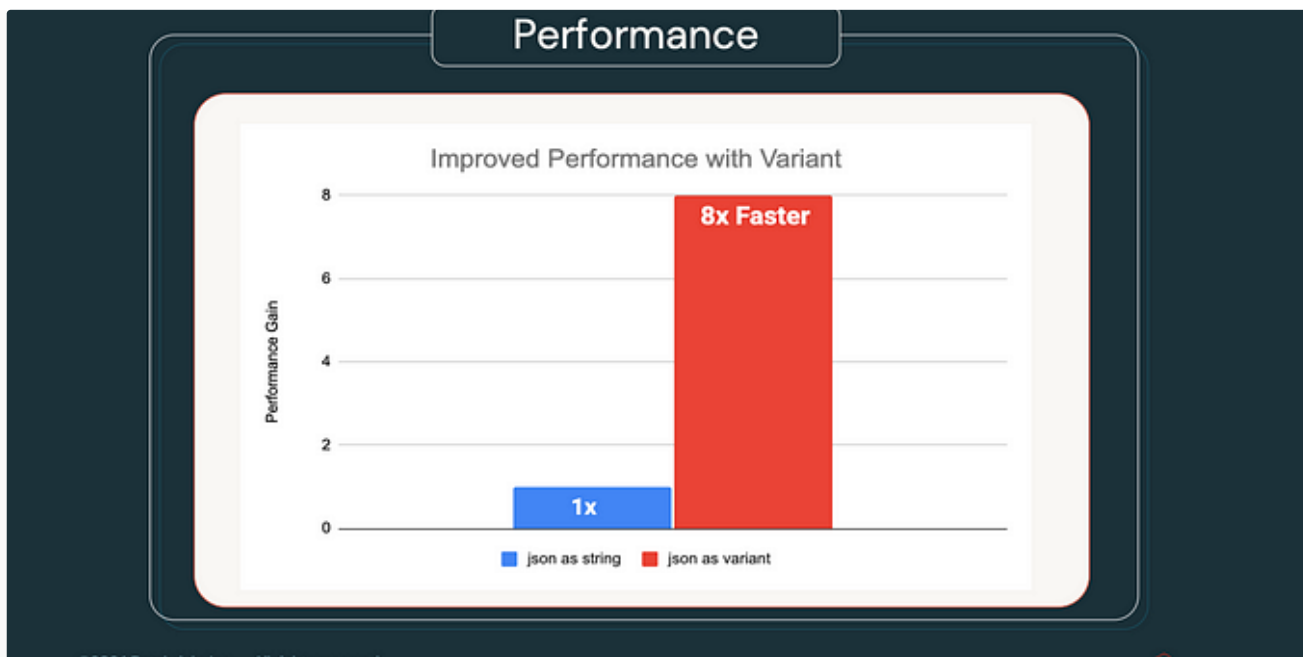
Introduction:

May 4

## Lists



### Natural Language Processing
1715 stories · 1287 saves
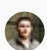
©2024 Databricks Inc. — All rights reserved

Archana Goyal

## What's Next for Apache Spark 4.0: A Comprehensive Overview with Comparisons to Spark 3.x

My articles are open to everyone; non-member readers can read the full article by clicking this link.
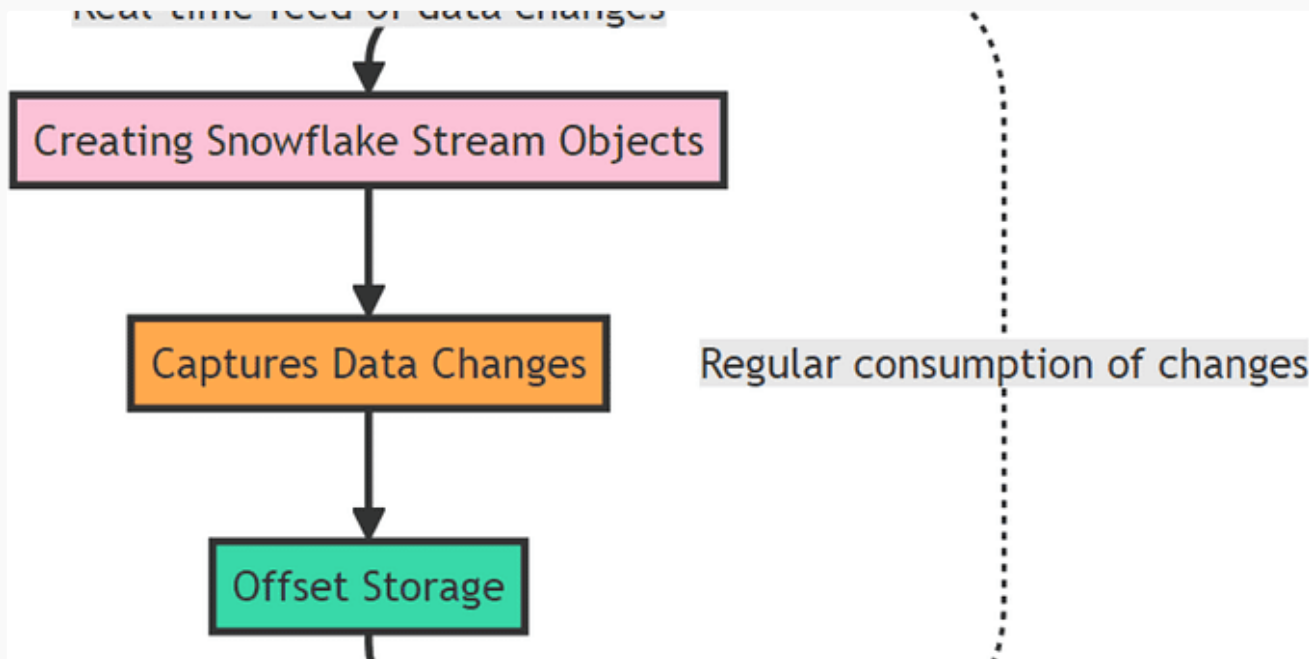
✦ Aug 25   👋 245   💬 2

Hugo Lu

## Snowflake vs. Databricks 2024 (actually useful)

Snowflake vs. Databricks is something we've all heard before, so why not take a different approach

👤 Hanson Olatunde

## Change Data Capture Pipeline Automation in Snowflake

This article offers technical guidance on automating updates and tracking changes of data examples using Snowflake's CDC capabilities. The...

Gavin F.

## Kafka Streams — How to magically join multiple data streams

Seamless Kafka Streams joining just like SQL table joins

✦ Nov 5, 2023   👏 281   💬 3

See more recommendations

Kafka Streams — How to magically join multiple data streams

Seamless Kafka Streams joining just like SQL table joins