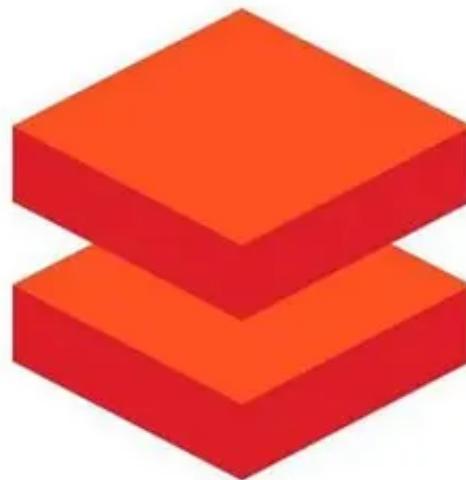


[Open in app ↗](#)[≡ Medium](#) [Search](#) [Write](#)

◆ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#) X



Slash Your Databricks Costs with a Spark Local Setup: A Guide to Smarter Pipelines



Anurag Sharma 4 min read · Jun 7, 2025

24



...

Picture this: your Databricks clusters are humming, churning through terabytes of data, but your cloud bill is starting to look like a phone number.

💡 What if you could slash those costs while building and testing reporting pipelines faster? That's exactly what one data team did by leveraging a **Spark local setup** and a clever code template.

The Problem: Databricks Costs Were Eating Our Lunch

Databricks is best for running Spark workloads, but those clusters don't come cheap. Every time a data engineer spun up a cluster to develop or test a reporting pipeline, the meter was running. 📈 The team needed a way to build and test pipelines without burning through their budget.

Enter the **Spark local setup** — a game-changer that lets them develop on their laptops, not in the cloud.

Fun Fact: Running a medium-sized Databricks cluster for a few hours can cost more than you can think.

Spoiler: This team found a way to cut costs.

The Solution: A Spark Local Setup with a Smart Template

The team crafted a **standardized code template** that every data engineer followed, making pipeline development smooth. This template broke the ETL process into three snappy components: **Reader**, **Transformer**, and **Writer**. (Fancy names for Extract, Transform, Load)

Here is how it worked:

1 Reader: Ingest Data Like a Pro

The Reader component was preloaded with functions to handle every Spark format, including CSV, Avro, Parquet, ORC, Iceberg, and Delta. It forced engineers to apply filters at the read level, cutting down on unnecessary data before it even hit the pipeline.

Pro Move: By filtering early, the team reduced memory usage and processing

time, making local development lightning-fast.

2 Transformer: Business Logic, Step by Step

The Transformer was where the magic happened. It lets engineers write business logic in small, modular segments, chained together using Spark's transform function. This step-by-step approach was like building a LEGO masterpiece — one block at a time, no chaos. 🧱 Each transformation was easy to debug, test, and tweak, all on a local machine.

Why It Rocks: Modular code meant fewer bugs and faster iterations. Plus, running transformations locally saved the team from spinning up costly Databricks clusters for every test.

3 Writer: Output with Precision

The Writer component was picky — it only allowed output in specific columnar formats (like Parquet or Delta) and supported direct loading into Amazon Redshift or RDS. This ensured consistency across pipelines and made downstream integration a breeze.

Bonus: The Writer's strict rules prevented sloppy outputs that could clog up storage or slow down queries.



Beneficial
intelligence

Get The Medium Newsletter
for weekly wisdom.

Subscribe now

The Local Spark Setup: 75% of Work Off the Cloud

The team moved 75% of their query development and testing to a local Spark setup, using tools like IntelliJ or PyCharm. They created sample datasets with a data generator to mimic production data in CSV format. This meant they could build, test, and debug entire pipelines on their laptops without touching Databricks.

Once the pipeline was running smoothly locally, they used the Databricks

CLI to push the code to a Databricks cluster for a final sanity check. If it passed, deployment was a simple matter of hooking the code into their CI/CD pipeline. It was like test-driving a car in your driveway before taking it on the highway. 

Cluster Cost Reduction with Local Development

Monthly cluster costs dropped from \$10,000 to \$2,500 after adopting local development—a 75% reduction!



Impact:

- 💰 **Cost Savings:** By keeping most development local, the team slashed Databricks cluster usage by 75%, saving thousands in compute costs.
- ⚡ **Faster Development:** Local testing meant instant feedback, no waiting for clusters to spin up.
- 🛡️ **Fewer Errors:** The structured template and local testing caught issues early, reducing production failures.



Why This Approach Is a Game-Changer

This Spark local setup with a standardized template isn't just about saving money — it's about working smarter.

Why it's a data engineer's dream:

- **Consistency:** The template ensured every pipeline followed the same structure, making code reviews and collaboration smooth.
- **Scalability:** The same code that ran locally worked seamlessly in Databricks, so scaling up was effortless.
- **Developer Happiness:** No one likes waiting for clusters or debugging in the cloud. Local development lets the team iterate faster. 😎

How to Steal This Playbook

Ready to slash your Databricks bill? Here's how to get started:

1. **Build a Template:** Create a Reader-Transformer-Writer structure for your pipelines. Predefine functions for common formats and enforce early filtering. Create a flow of code structure that the developer needs to follow to

keep the code simple and easy.

- 2. Set Up Local Spark:** Install Spark locally and use IntelliJ/PyCharm to run and debug your code. Generate sample data to mimic production.
- 3. Test Locally, Deploy Smart:** Use the Databricks CLI to validate your code in the cloud only when it's ready. Then deploy via your CI/CD pipeline.
- 4. Track Savings:** Monitor your Databricks usage to see the cost drop.

The Bottom Line

- **Problem:** Sky-high Databricks costs from constant cluster usage for development and testing.
- **Solution:** A Spark local setup with a standardized Reader-Transformer-Writer template, moving 75% of the work off the cloud.
- **Impact:** Massive cost savings, faster development, and happier engineers.

[Databricks](#)[Spark](#)[Automation](#)[Data Engineer](#)[Data Engineering](#)

Written by Anurag Sharma

82 followers · 3 following

[Edit profile](#)

Data Engineering Specialist with 10+ exp. Passionate about optimizing pipelines, data lineage, and Spark performance and sharing insights to empower data pros!

No responses yet



...



Anurag Sharma him/he

What are your thoughts?

More from Anurag Sharma



 Anurag Sharma

Mastering Code Reviews in Data Engineering: A Step-by-Step Guide

Hey there, data enthusiasts! If you are knee-deep in the world of data pipelines, ETL jobs,...

Feb 2  29



...



 Anurag Sharma

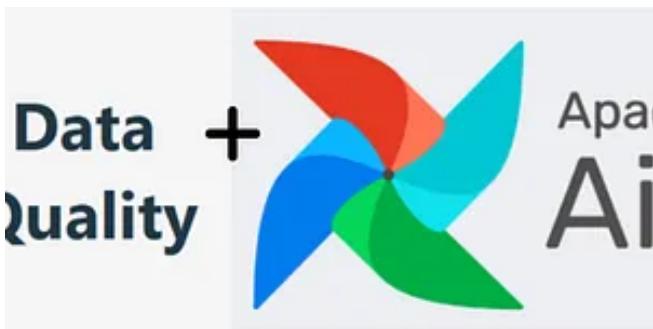
How We Built a Generic Data Cleaning Framework for Every...

Hey Data Folks! Clean data is the foundation of every reliable dashboard, model, and...

Feb 9  2



...



 Anurag Sharma



 Anurag Sharma

Automating Data Quality in Data Engineering: A Game-Changer fo...

As a data engineer, ensuring data quality is like keeping the engine of a car running...

Jun 3, 2025 99



...

Modern Data Platform for ELT Data into Snowflake + AWS

This blog describes a production-ready data platform built for a Data Team that handles...

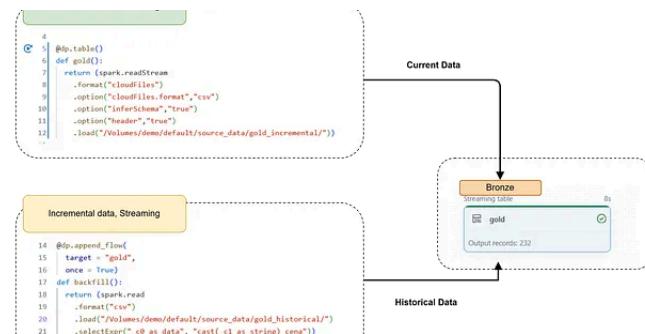
Nov 4, 2025



...

See all from Anurag Sharma

Recommended from Medium



Mariusz Kujawski

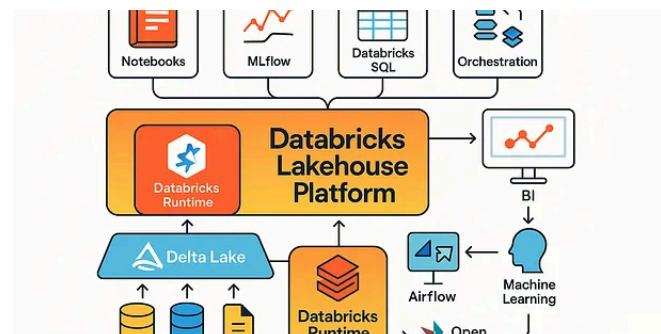
Databricks Declarative Pipelines: Backfill, Time Travel, Delete, and ...

Spark Declarative Pipelines (SDP) are becoming a more powerful tool with the...

Nov 13, 2025 28 2



...



In Tech with Abhishek by Abhishek Kumar Gupta

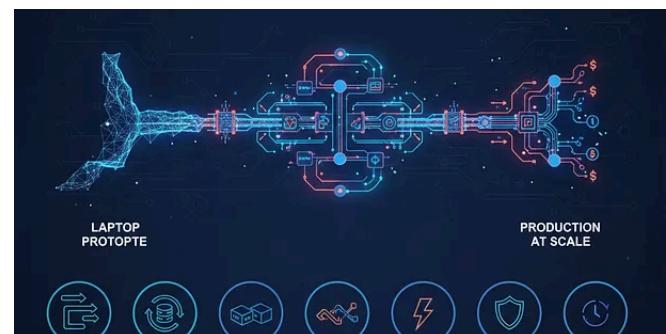
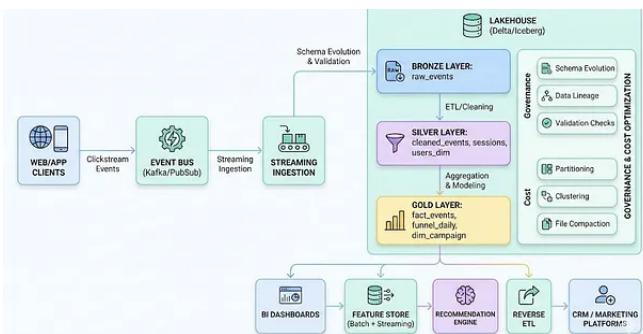
🔥 50 Databricks Interview Questions & Answers: The Ultima...

Unlock your Databricks interview success with expert questions and practical, detailed...

Aug 23, 2025 140 3



...



Sarah Sagi

Data Engineering System Design: Clickstream Data Into a Modern...

A complete end-to-end architecture for building scalable clickstream analytics and...

Dec 12, 2025

12



...



Reliable Data Engineering

15 Data Pipeline Architecture Patterns Every Engineer Should...

Your pipeline works on your laptop. Here's how to make it work in production — at scale...

Feb 15

140

2



...

```

: line magics:
%alias_magic %autoawait %autocall %automagic %autosave %bookmark %cd %clear %
%`s %conda %config %connect_info %cp %debug %dhist %dirs %doctest_mode %edit %
%edit %killbgscripts %ldir %less %lf %k %ll %load %load_ext %loadpy %logoff %
%logstate %logstop %ls %lsmagic %lx %macro %magic %mamba %man %matplotlib %
%more %new %notebook %profile %page %pastebin %pdb %pdef %pdoc %pfile %pinfo %
%pipeline_codegen %popd %pprint %precision %prun %psearch %psource %pushd %pwd %pyc %
%isole %quickref %recall %rehashx %reload_ext %rep %rerun %reset %reset_selective %
%rm %rmdir %run %save %sc %set_cell_max_output_size_in_mb %set_env %store %sx %synt %
%timeit %top %unalias %unload_ext %uv %who %who_ls %whos %xdel %xmode
%is ON, % prefix IS NOT needed for line magics.

: cell magics:
%%ML %%%SVG %%bash %%capture %%code_wrap %%debug %%file %%html %%javascript %%js %%
%%memprofile %%koprofile %%perl %%pipeline_codegen %%profile %%prun %%pypy %%pyth%
%%python3 %%ruby %%script %%sh %%sql %%svg %%sx %%system %%time %%timeit %%writefile

: is ON, % prefix IS NOT needed for line magics.

```



Hubert Dudek

Hidden Magic Commands in Databricks Notebooks

If you're working with Databricks notebooks, you're probably familiar with basic magic...

Jan 10

120

5



...

	<code>i₂ id</code>	<code>i₂ name</code>	<code>i₂ email</code>
1	Customer 1	customer1@example.com	c*****@example.com
2	Customer 2	customer2@example.com	c*****@example.com
3	Customer 3	customer3@example.com	c*****@example.com
4	Customer 4	customer4@example.com	c*****@example.com
5	Customer 5	customer5@example.com	c*****@example.com
6	Customer 6	customer6@example.com	c*****@example.com
7	Customer 7	customer7@example.com	c*****@example.com
8	Customer 8	customer8@example.com	c*****@example.com
9	Customer 9	customer9@example.com	c*****@example.com
10	Customer 10	customer10@example.com	c*****@example.com



Christian Hansen

A Step-by-Step Guide to Setting Up ABAC in Databricks (Unity Catalog)

How to use governed tags, dynamic policies, and UDFs to implement scalable attribute-...

Dec 1, 2025

16



...

[See more recommendations](#)