



# AT&T Billions of Events Processing migration

Praveen Vemulapalli, Director – Technology , AT&T

Akshay Sharma, Sr. Solutions Consultant , Databricks

**June 11, 2024**

# Praveen Vemulapalli

## Things I love to do....

- Love Hiking & Camping
- Love motorcycle riding
- Spend loads of time with my family
- Data & AI Technology evangelism
- Drive change & evolution





# AT&T Background

AT&T started with Bell Patent Association, a legal entity established in 1874 to protect the patent rights of Alexander Graham Bell after he invented the telephone system. Originally a verbal agreement, it was formalized in writing in 1875 as Bell Telephone Company.





By 2024, We're turning to public cloud providers to host our *non-network* workloads. Think traditional IT applications like billing and customer care, and corporate applications like HR and finance (stated in 2019) (source: [https://about.att.com/innovationblog/2019/08/cloud\\_strategy.html](https://about.att.com/innovationblog/2019/08/cloud_strategy.html))

In June 2021, Microsoft and AT&T reached a major milestone when we announced an industry-first collaboration to evolve Microsoft's hybrid cloud technology to support AT&T's 5G core network workloads. (source : <https://azure.microsoft.com/en-us/blog/improving-the-cloud-for-telcos-updates-of-microsoft-s-acquisition-of-att-s-network-cloud/>)

# AT&T's Motivation for Modernizing Hadoop to Databricks



## Change Drivers

### Reduce TCO spend

- Data Centers are Capital Intensive
- Software Cost
- Utility/admin Costs of On-prem Infra

### Nimbleness, Scalability and Innovation

- Models & Data Science Sequenced to Fit Capacity
- Jobs Failed Due to Capacity Constraints
- Adding Data & Analytic Use Cases Required Infrastructure and Increased Sustaining Costs



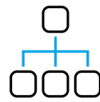
## Future-State Goals



Single Version of Truth



Parallelize, Simplify & Automate



Move Resources up the Value Chain



Free Capital for Growth-Oriented Investments



Enable streaming pipelines & analytics



Empower citizen data scientists & analytics  
+60 BUs



## Success To-Date

- Rationalized +30% of the Data
- Migrated 100% of the User Base
- Accelerate Nimbleness Up to 3x for Key Data Science Activities
- Launched Self-Serve ML Analytics Platform
- Co-located Batch & Streaming Data Products and Analytics
- Streamline Model Recreation/Lineage from Hours to Minutes
- Retired +40% of Servers to Date (100% Q1'23)
- Re-invested Unlocked Resources Improving Effective Cloud Run Rate Value

Source: <https://www.databricks.com/customers/att/migration>



5-year Migration ROI of +300%



A photograph of two hikers in a forest. The hiker on the left is wearing a yellow jacket and a blue beanie, looking at a smartphone. The hiker on the right is wearing a purple jacket and a yellow beanie, holding a trekking pole. The forest has many thin trees and a ground covered in pine needles.

Problem to Solve :  
Large scale event time correlation process

17B+

Events generated by network daily across our apps that do analytics

6400 CPU's

Core Hadoop system was used to manage the daily processing

22-30hrs

Daily batch run times on Proprietary analytics platform for processing



End state :

Large scale event time correlation process

30%

Cost reduction compared to Hadoop environment – Substantial savings at scale

1000 CPU's

Used dynamically for analytics processing

8Hrs

~60% reduction in data processing time. from 30hrs to 8hrs

Analytics processing moved to Spark & Scala



# Akshay Sharma

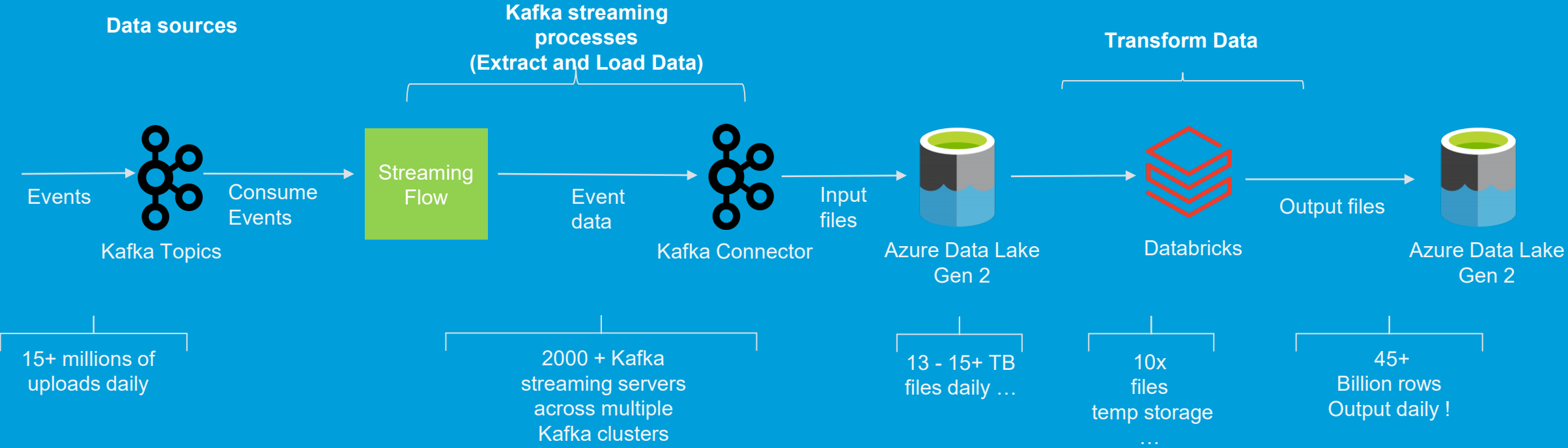
## Things I love to do....

- Listening Music
- Learning new technologies
- Playing PC games
- LeetCode challenges.





# High level Solution Architecture





# Challenges

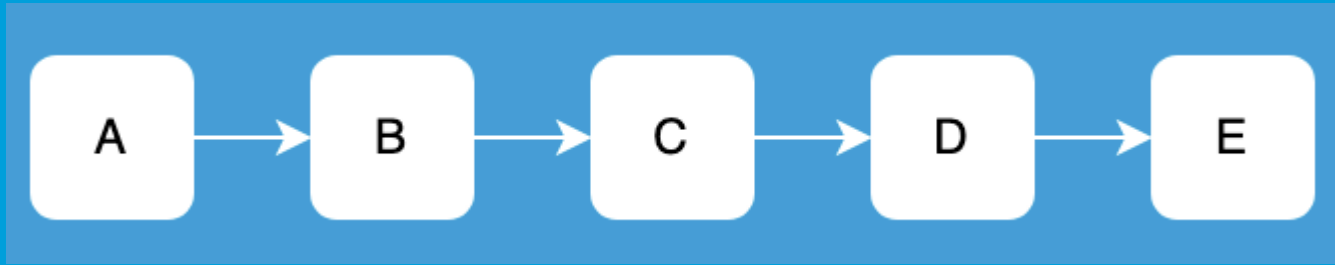
**1. Code Migration** (Loops, Disk IO) MR → RDDs → Dataframes

**2. Tuning Storage account API Rate limits**

**3. Data Quality issues** (DeDuplication, Nulls, DateTime formats)



# Task Orchestration



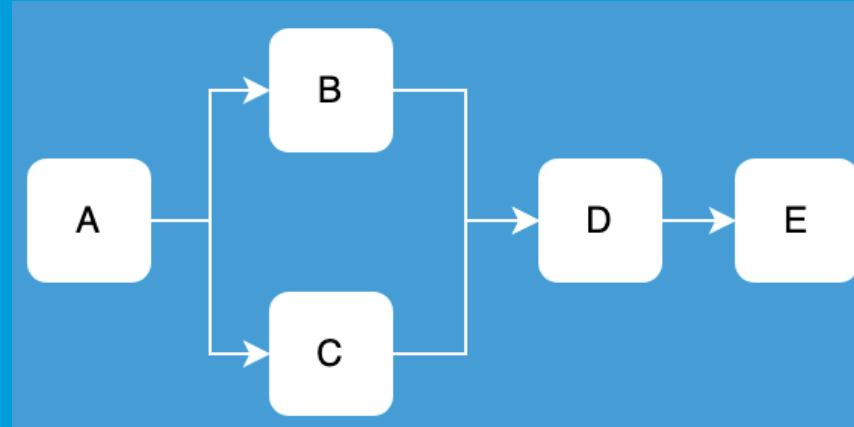
- A = 30 mins ● B = 20 mins
- C = 60 mins ● D = 15 mins
- E = 5 mins

$$30+20+60+15+5 = 130 \text{ mins (2 hrs 10 mins)}$$

Here A, B, C, D, E are individual tasks or let's say *Notebooks* which are going to get executed one after the other.



# Task Orchestration



Here we have enabled parallelism  
By having A FAN-OUT to B and C

Total Time :  $A + \max(B,C) + D + E$

New Time :  $30 + 60 + 15 + 5 = 110$  mins (1 hr 50 mins) (Less by 20 mins)

*Cluster 1 : A, C, D, E*

*Cluster 2 : B*

# Best Practices in Action

Cache and  
Persist

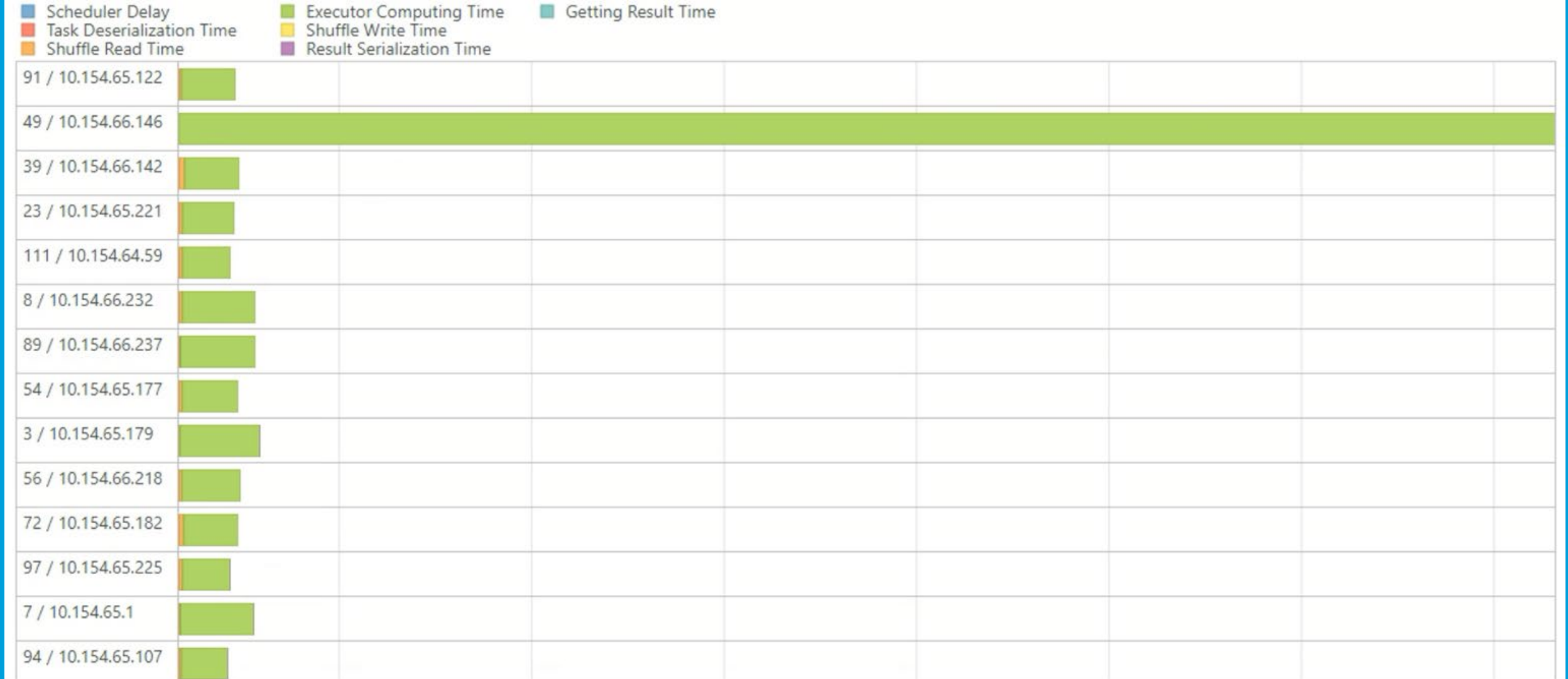
Flexible  
Databricks  
Runtimes

Data  
Distribution

Photon  
Execution



# Data Skew Example

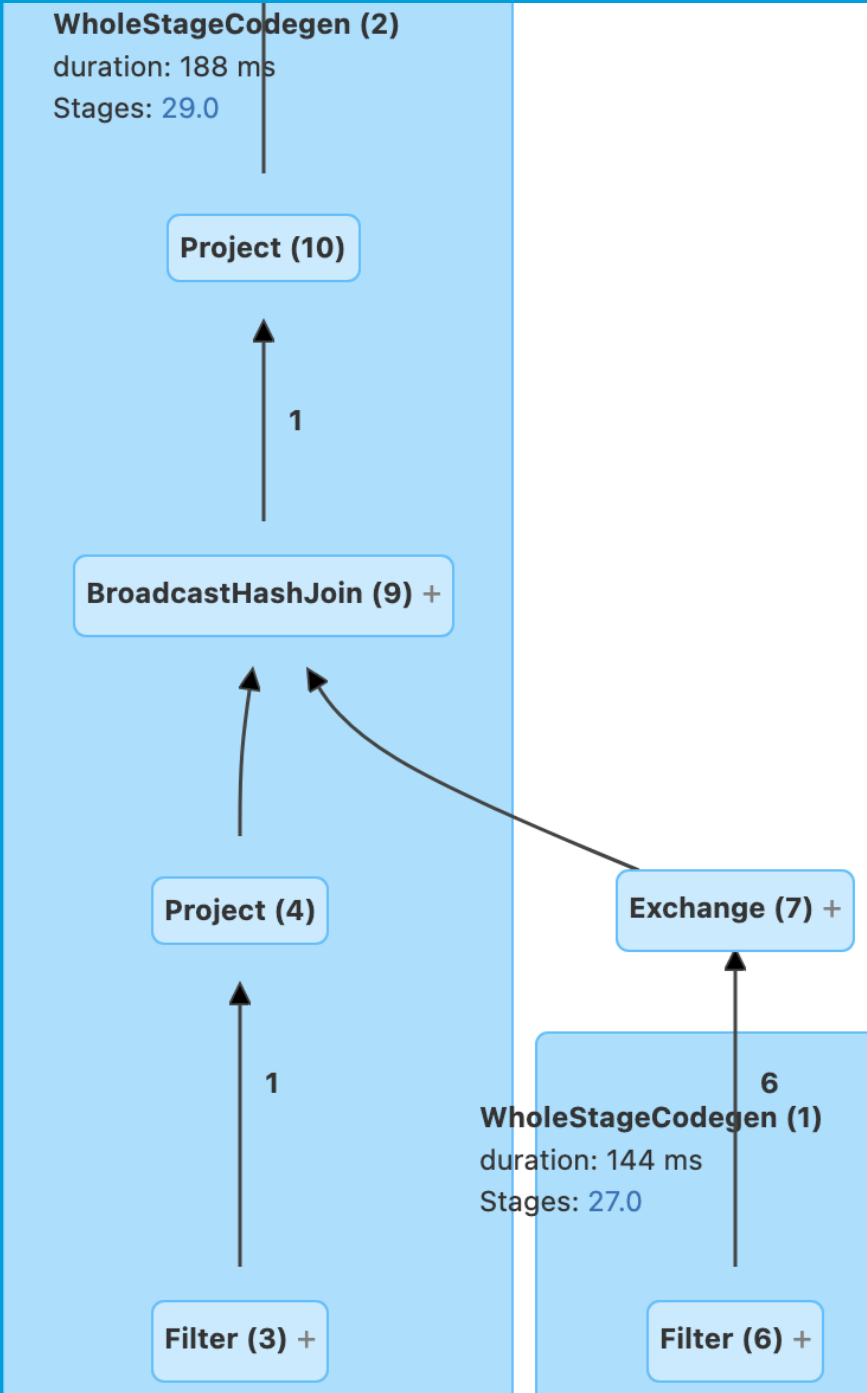
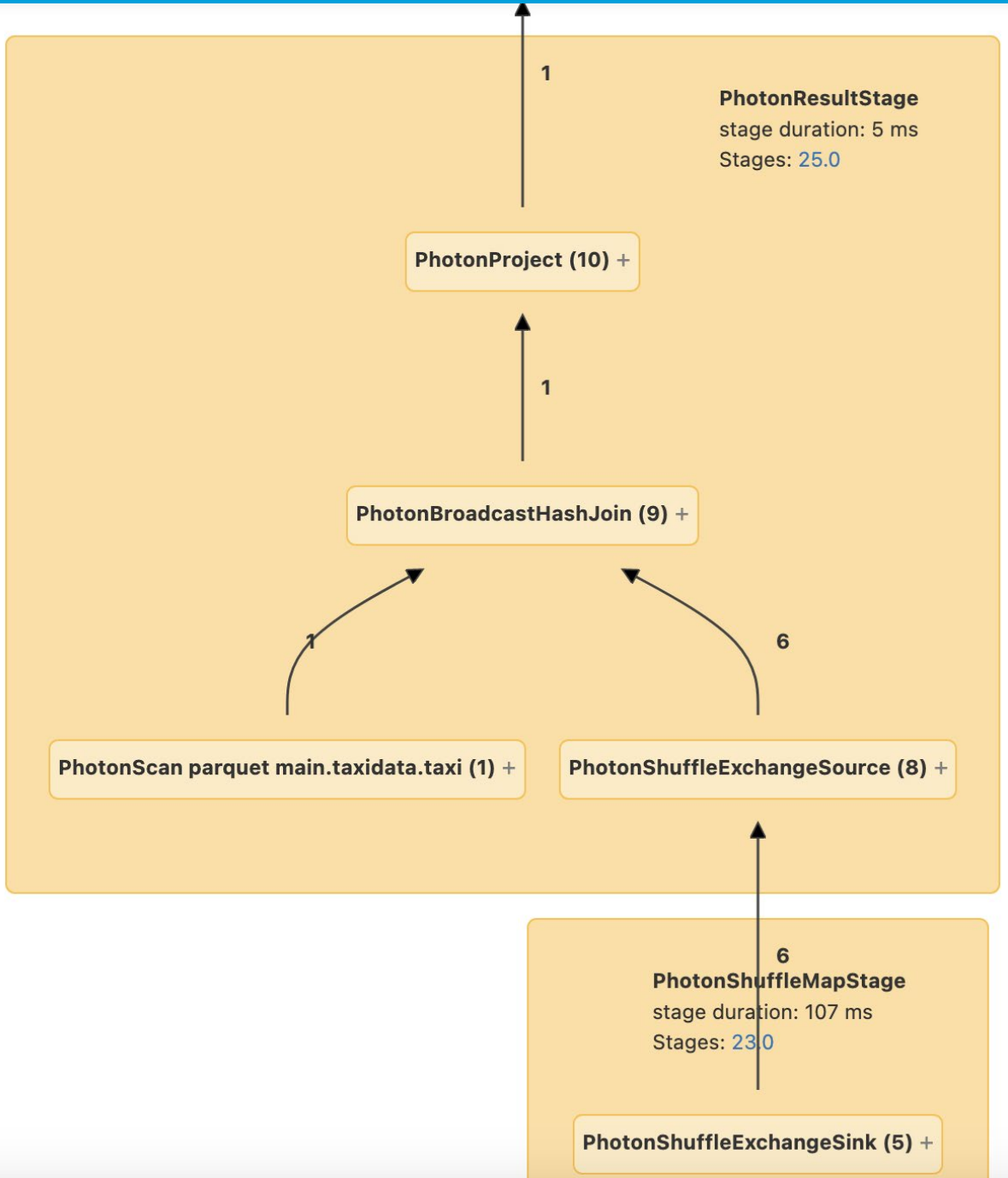


# Photon

The next-generation engine for the lakehouse







# Key Takeaways

1. Stick with Dataframes and it's supported features

2. Consider your Storage Account.

3. Data quality impacts parallel processing.



# Databricks Workflows

1 Ease Of Use

2 Detailed Monitoring

3 Scalable

# Databricks Workflows

**Jobs** consist of one or more **Tasks**



Databricks Notebooks



Python Scripts



Python Wheels



SQL Files/Queries



Delta Live Tables Pipeline



dbt



Java JAR file



Spark Submit

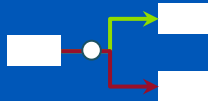
**Control flows** can be established between **Tasks**.



Sequential



Parallel



Conditionals (Run If)



Jobs-as-a-Task (Modular)



For-Each Loop

**Jobs** supports different **Triggers**



Manual Trigger



Scheduled (Cron)



API Trigger



File Arrival Triggers



Continuous (Streaming)

# Databricks Workflows

## Task Dependencies

When a task is **Done**, it can be in a **Success**, **Failure**, or **Excluded** state.

### All Succeeded

*Default behaviour*



### At Least 1 Succeeded

*e.g. Fan in with at least some success*



### None Failed

*e.g. Run task(s) at the end of DAG if nothing fails*



### All Done

*e.g. Perform clean up even if tasks have failed or excluded*



### All Least 1 Failed

*e.g. Perform clean-up with observability or specific actions*

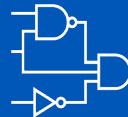


### All Failed

*e.g. Perform clean-up with observability or specific actions*



## Parameterisation



### Job Parameters

Passed into each Task with behaviour based on the type  
*e.g. additional options for JARs, spark-submit, Python Args*



### Job Contexts

Special set of templated variables that provide introspective metadata about job and task  
*e.g. run\_id, job\_id, start\_time*



### Task Values

Custom parameters that can be shared between Tasks in a Job  
*e.g. anything that can be programmatically set or retrieved!*

## Webhooks

Allows customers to build **event-driven integrations** with Databricks.

Supported destinations are **Slack** and **Webhooks**, with the below **notification events**:

For example, you can send a message to a Slack #channel when:



**On start:** Send a message to a when a job or a parent run is started



**On success:** when a job or a parent run finished without any errors



**On failure:** when a job fails or a parent run is terminated with one of the children in a failed state.





**THANK  
YOU**