

STRATEGIC PARTNER



GOLD SPONSOR



SILVER SPONSOR



BRONZE SPONSOR





Azure Databricks 101



Kamil Nowiński

Principal Microsoft Consultant



Kamil Nowiński



Microsoft Data Platform **MVP**
Speaker, blogger, data enthusiast
Principal Microsoft Consultant at Altius
(www.altiusdata.com)
Almost 20 yrs experience as DEV/BI/(DBA)
Member of the Data Community PL
Project member of „SCD Merge Wizard”
Founder of blog SQLPlayer (www.SQLplayer.net)

SQL Server Certificates:
MCITP, MCP, MCTS, MCSA, MCSE Data Platform,
MCSE Data Management & Analytics
Moreover: Bicycle, Running, Digital photography
@NowinskiK, @SQLPlayer

Blog

- Technical posts
- Various skill level
- Cheat sheets
- Recommended books
- Many useful other links
- Interviews (Podcast)
- YouTube Channel



SQL Player

Play with data & have fun!

www.SQLPlayer.net



Scan me

"Ask SQL Family" #podcast



Scan me



YouTube



www.SQLPlayer.net/YouTube

Slides available

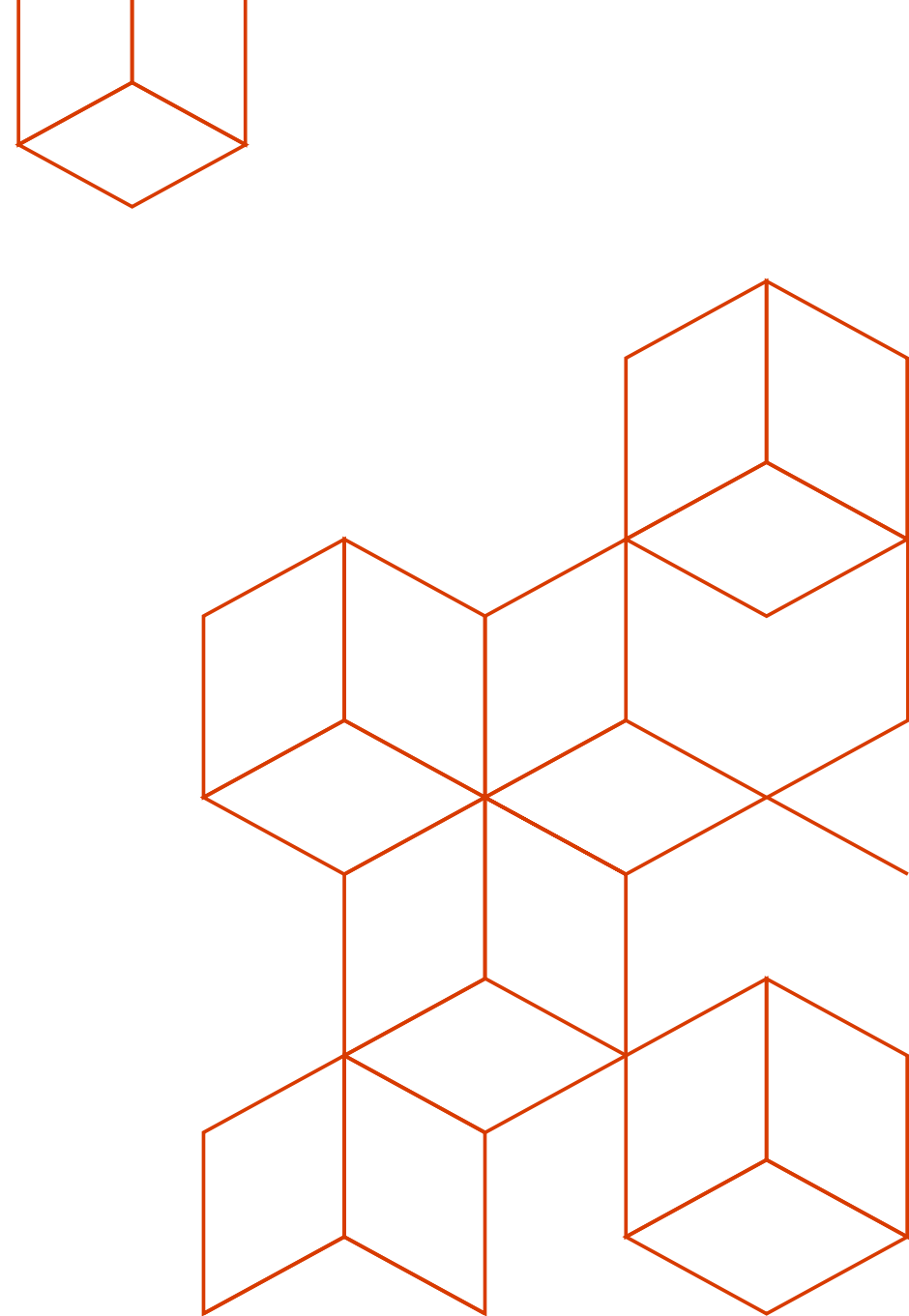


github.com
/NowinskiK/CommunityEvents



- github.com/sqlplayer/
 - SCD-Merge-Wizard
 - DataScriptWriter
- [azure.datafactory.tools](https://github.com/AzureDataFactory/azure.datafactory.tools)
- [azure.datafactory.devops](https://github.com/AzureDataFactory/azure.datafactory.devops)

Azure Databricks



What is Azure Databricks?



A fast, easy and collaborative Apache® Spark™ based analytics platform optimized for Azure



Best of Databricks



Best of Microsoft



Designed in collaboration with the founders of Apache Spark



One-click set up; streamlined workflows



Interactive workspace that enables collaboration between data scientists, data engineers, and business analysts.



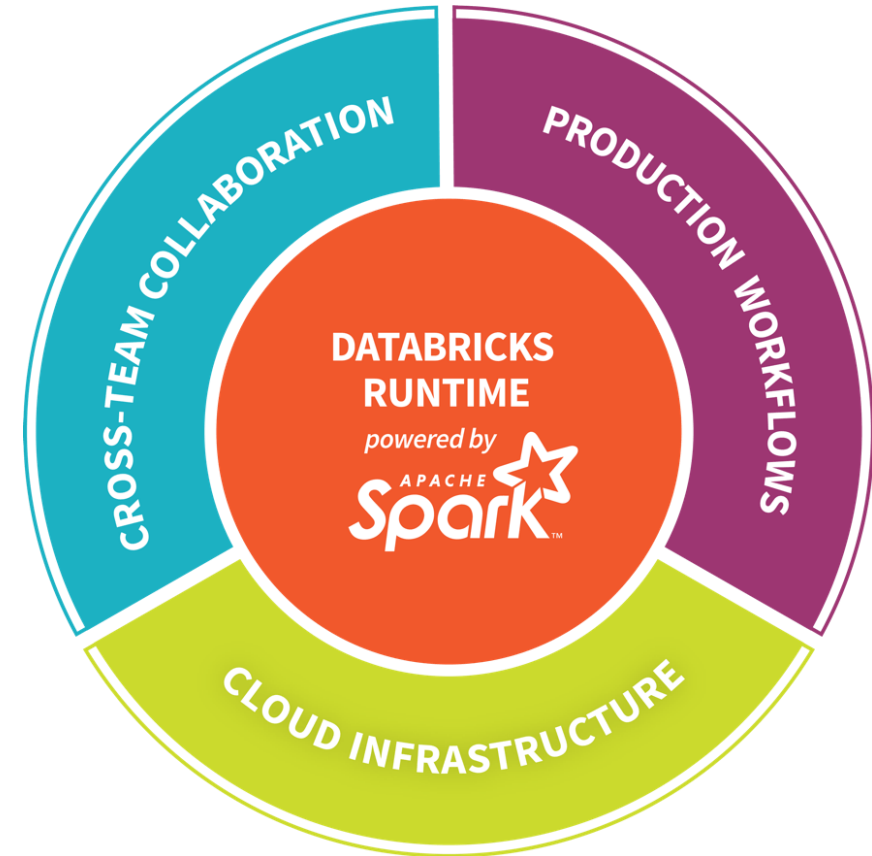
Native integration with Azure services (Power BI, SQL DW, Cosmos DB, Blob Storage)



Enterprise-grade Azure security (Active Directory integration, compliance, enterprise -grade SLAs)

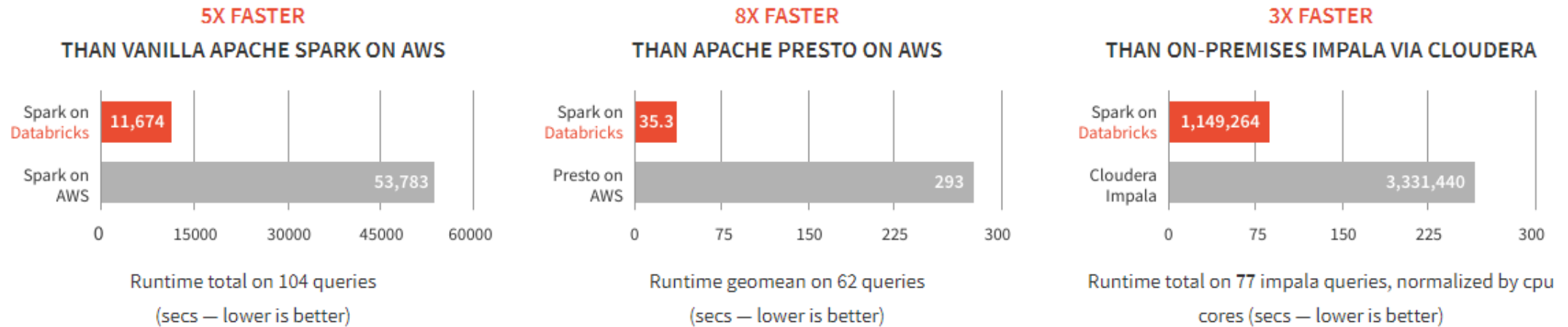
DATABRICKS - COMPANY OVERVIEW

- Founded in late 2013
- By the creators of Apache Spark, original team from UC Berkeley
- Largest code contributor code to Apache Spark
- Main Product: The [Unified Analytics Platform](#)



DATABRICKS SPARK IS FAST

Benchmarks have shown Databricks to often have better performance than alternatives



SOURCE: [Benchmarking Big Data SQL Platforms in the Cloud](#)

Why Spark?



- Open-source data processing engine built around speed, ease of use, and sophisticated analytics
- In memory engine that is up to 100 times faster than Hadoop
- Largest open-source data project with 1000+ contributors
- Highly extensible with support for Scala, Java and Python alongside Spark SQL, GraphX, Streaming and Machine Learning Library

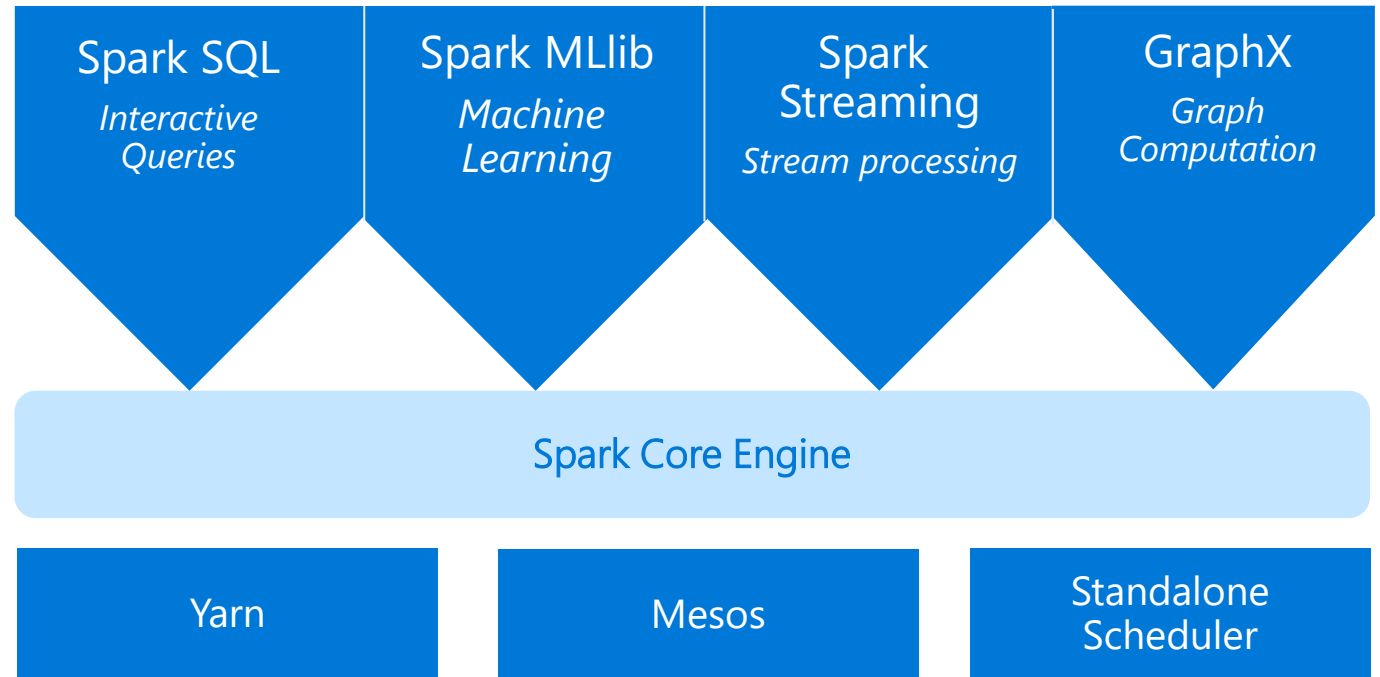
APACHE SPARK

A unified, distributed, open source engine for large-scale data processing

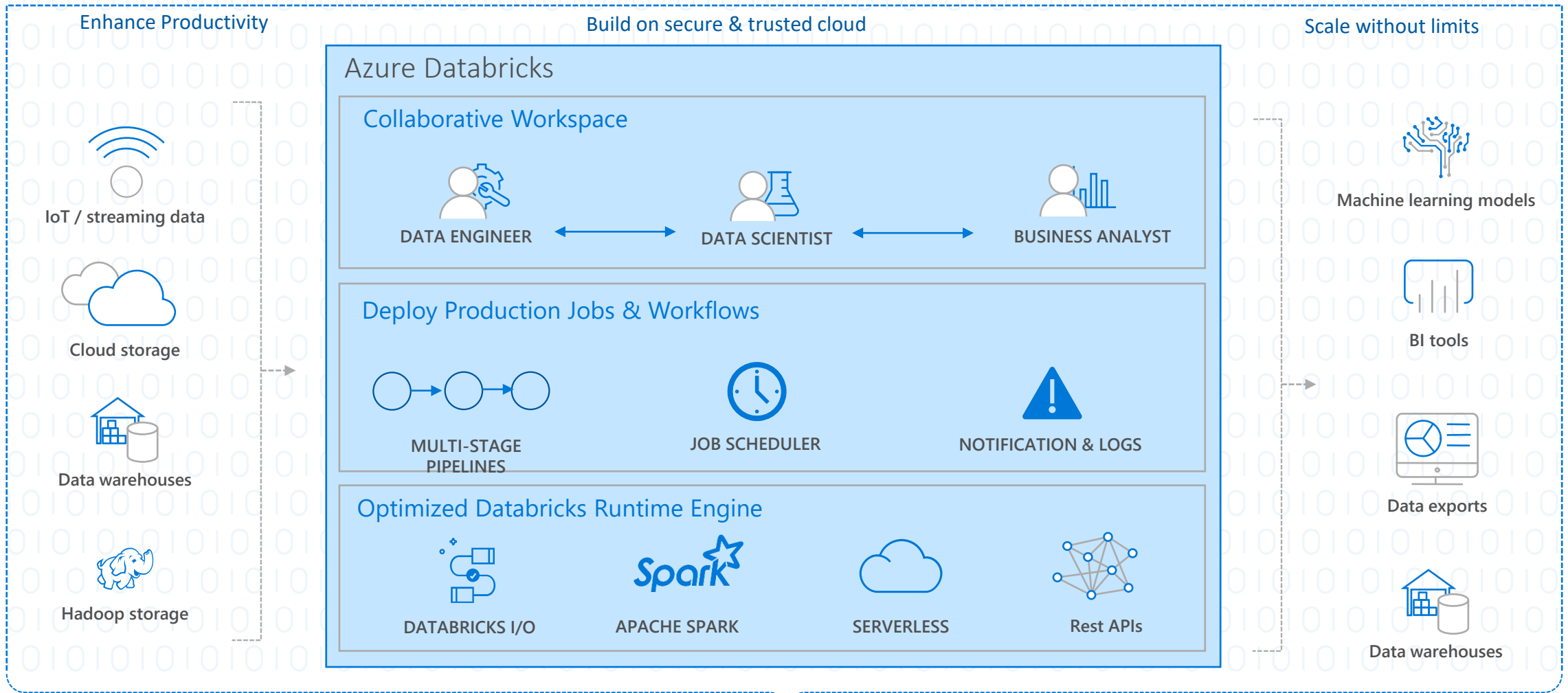


Spark Unifies:

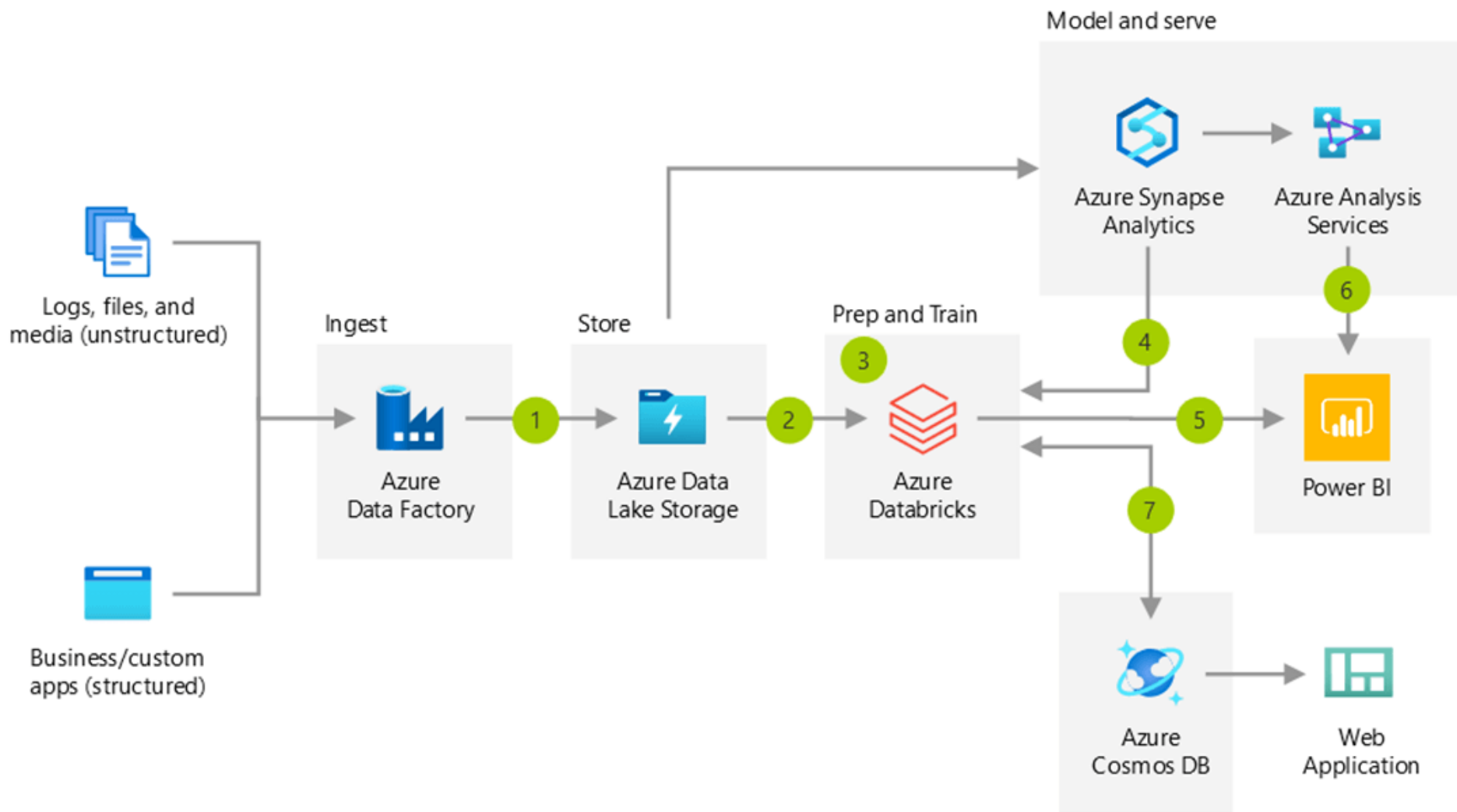
- Batch Processing
- Interactive SQL
- Real-time processing
- Machine Learning
- Deep Learning
- Graph Processing



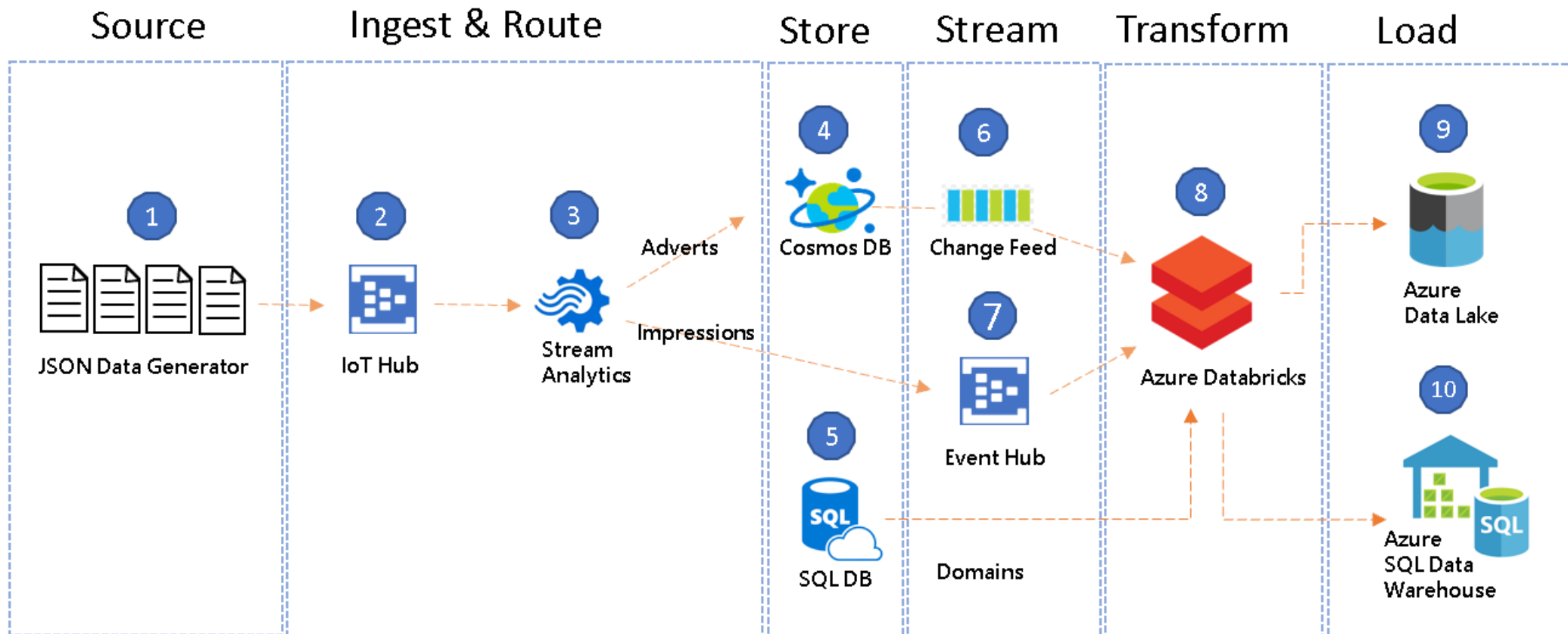
A Z U R E D A T A B R I C K S



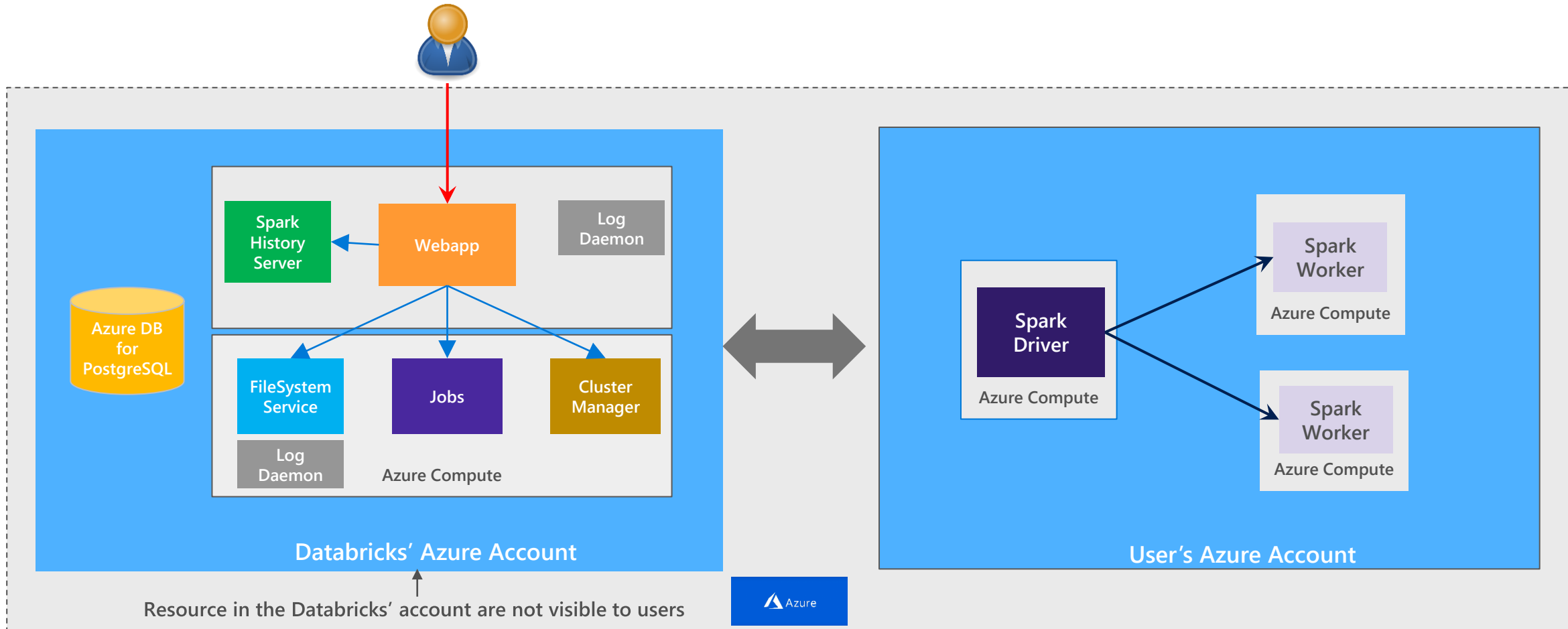
ARCHITECTURE EXAMPLE



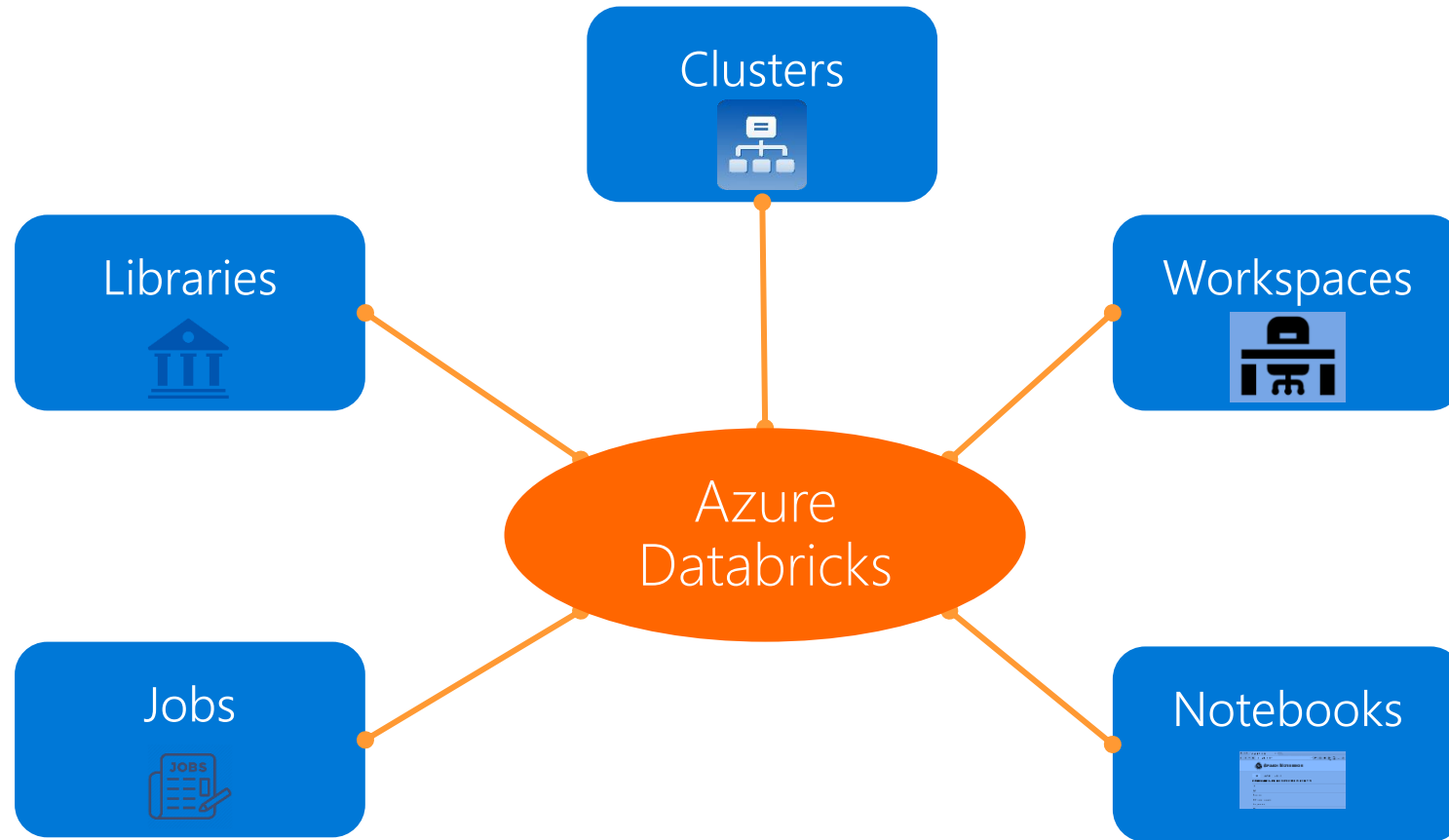
ARCHITECTURE EXAMPLE



AZURE DATABRICKS CLUSTER ARCHITECTURE



AZURE DATABRICKS CORE ARTIFACTS



Azure Databricks

Demo

Azure Databricks

Delta Lake

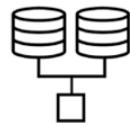
DELTA LAKE - Reliable Data Lakes at Scale on Azure



Data Versioning



ACID Transactions



Optimized Layouts



Fast Streaming



Efficient Upserts



Schema Enforcement

ACID Transaction Guarantees

- Atomic, Consistent, Isolated, Durable

Versioned parquet files

- Delta transaction log keeps track of all operations

Efficient Upserts

- *MERGE, DELETE, UPDATE*

Small file compaction w/ no interrupt to availability

- *OPTIMIZE* and *VACUUM*

Z-Order partitioning w/ up to 100x perf

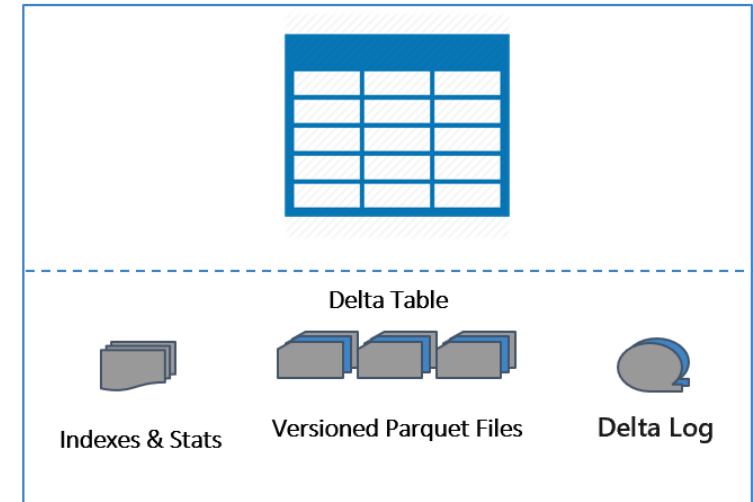
- New multidimensional partitioning enables data skipping

Time Travel

- Audit history, Pipeline Debugging, Data Reproducibility

Delta Table =

Parquet + Transaction Log + Indexes/Stats



Open Source



THE
APACHE
SOFTWARE FOUNDATION

+

THE
LINUX
FOUNDATION

Resources



[Azure Databricks – strona główna](#)

[Apache Spark](#)

[Databricks Community Edition](#)

[Microsoft Learn: Databricks](#)

[PowerShell module: azure.databricks.cicd.tools](#)

[Databricks CLI](#)

[Spark by {Examples}](#)

Questions?



Thank you!



kamil@nowinski.net



@NowinskiK

@SQLPlayer



SQLPlayer.net



<https://github.com/NowinskiK/CommunityEvents>



Kamil Nowinski

Microsoft Data Platform MVP

MCSE Data Platform & MCSE Data Management and Analytics