

DEV DAY



ML Dev Day Workshop

*Unifying Data Science and
Data Engineering*

WELCOME

- We will begin at 9:05am
- This event will be recorded
- You will receive an email with the Deck and recording tomorrow.
- Please feel free to use the Q&A tool for any questions
- Don't forget to participate in our polls!

Introductions

- Cameron Kashani, cameron.kashani@databricks.com
 - Account Executive
- Lei Pan- lei.pan@databricks.com
 - Solutions Architect
- Aaron Binns - aaron.binns@databricks.com
 - Solutions Architect
- Don Warnes - dong.warnes@databricks.com
 - AWS Alliance Manager

Additional Databricks Attendees: Raise your hand so you may introduce yourselves!

Introductions

- Brad Ito, [Retina](#)
 - CTO & Co-Founder, brad@retina.ai
- Mo Messidi, [Retina](#)
 - DataOps Engineer
- Jonathan Corners - [Pariveda](#)
 - jonathan.corners@parivedasolutions.com
- Brian Edwards - [Pariveda](#)
 - brian.edwards@parivedasolutions.com

Poll

Agenda

- 9:00am Databricks Keynote, Cameron Kashani
- 9:45am Customer Presentation | Retina
- 10:15am Partner Presentation | Pariveda
- 10:30am Break
- 10:40am Data Engineering Interactive Demo & Best Practices:
Preparing Data for Analytics | Aaron Binns
- 11:15am Data Science Interactive Demo & Best Practices: Model
Training and Machine Learning | Lei Pan
- 11:50am Q&A
- 12:00pm Finish

Unified Data Analytics Platform

Accelerating innovation across data science, data engineering, and business analytics

- Global company with 5,000 customers and 450+ partners
- Original creators of popular data and machine learning open source projects
- Empowering data teams to solve the world's toughest problems



Unlocking business value

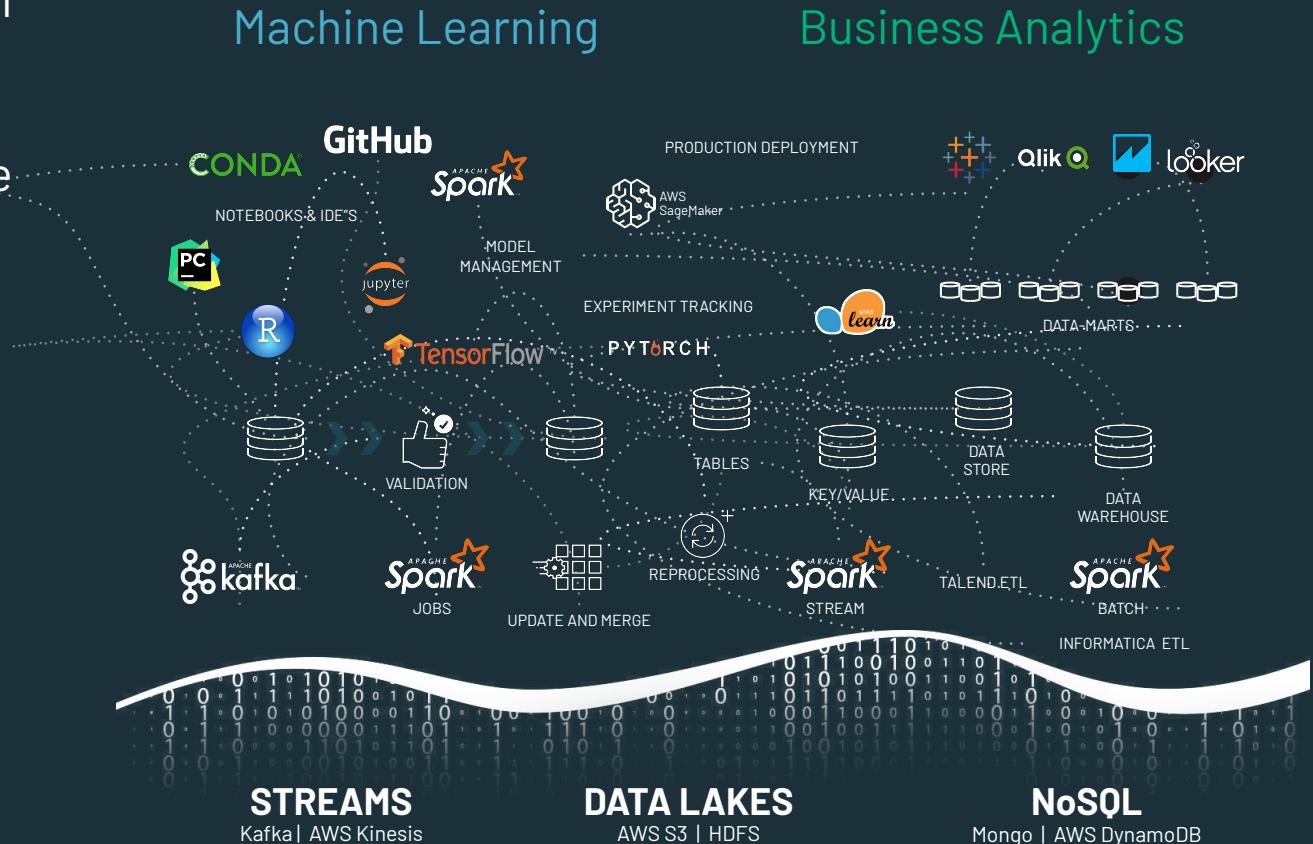
Companies are combining massive data with ML and BI capabilities

Machine Learning

Business Analytics



Most organizations fail to unlock business value due to data, technology and people silos



Unlocking business value: Four challenges

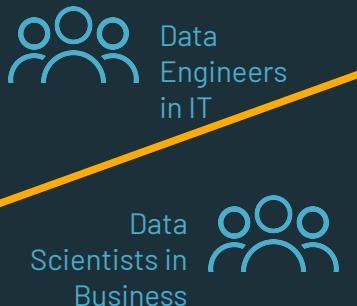
1

Data is messy,
siloed and slow



2

ML is hard,
Production is
harder



3

BI is limited to a
fraction of data

```
110001100011000100010001  
000010111000100101010000  
11110101010101111100111001  
11010100011100110001100011  
00010001000100001011100  
0100101010000111100101010
```

4

Lack of enterprise
readiness

- Fragmented security
- Poor reliability
- Disjointed governance

Make all your data ready for BI and ML

1

Data is messy,
siloed and slow



Unified Data Service

Build open, reliable, fast data lakes with all your data



Big Data



Business Data



Applications



DELTA LAKE™

Open High Quality Fast



hadoop



Amazon S3



Unified
Engine



BI
Reporting



Machine
Learning

Unify data and ML across the full lifecycle

2

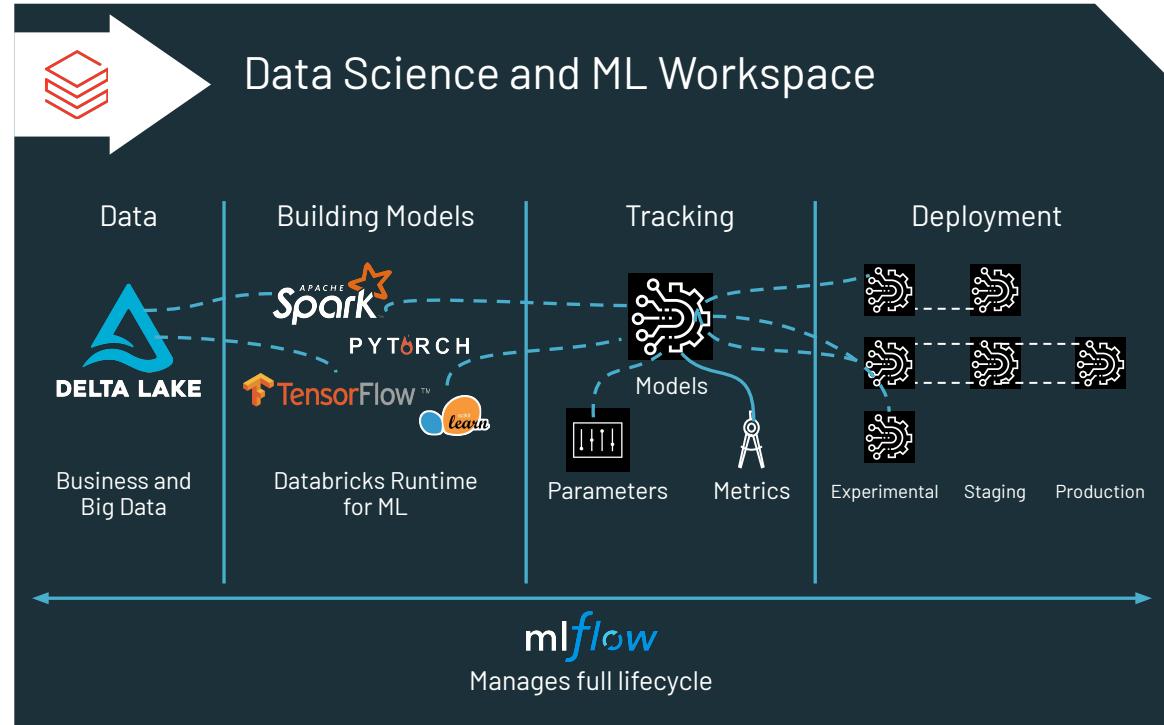
ML is hard,
production is
harder



Data
Engineers
in IT



Data
Scientists
in
Business

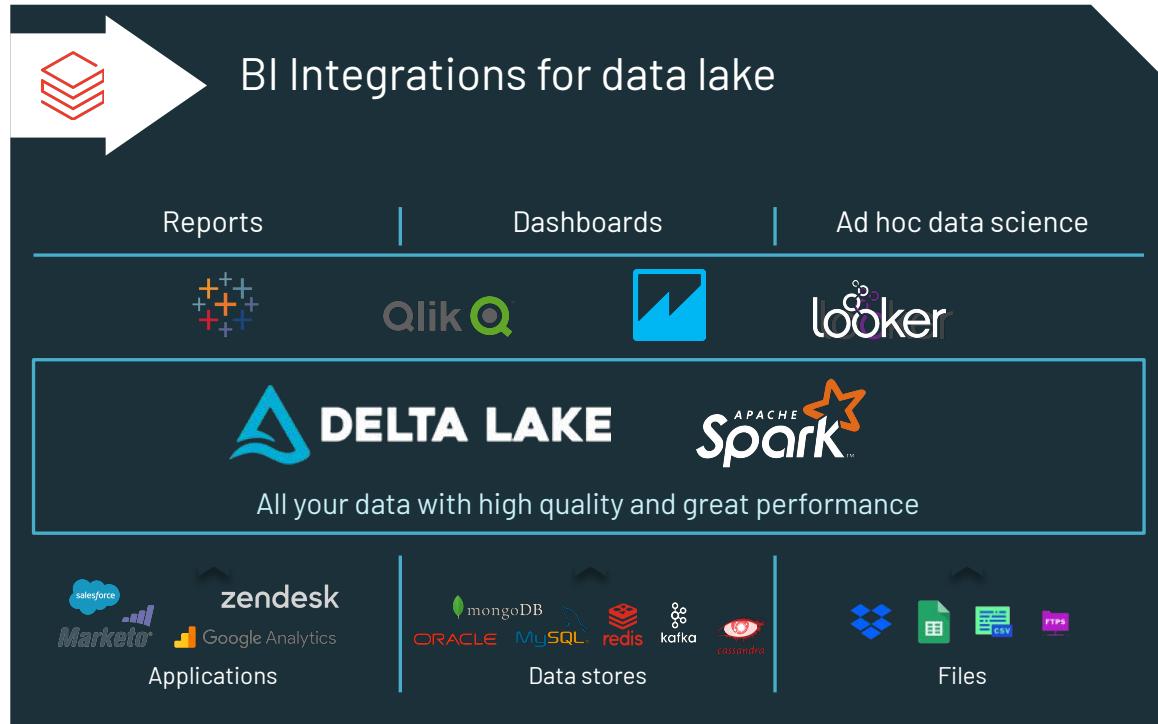


Enable BI directly on all your source data

3

BI is limited to a fraction of data

11000110001100010001000100
00101110001001010100001110
01010100111111001110011101010
001110011000110001100010001
00010000101110001001010100
001111001010



Leverage cloud native platform for enterprise grade solution

4

Lack of enterprise readiness



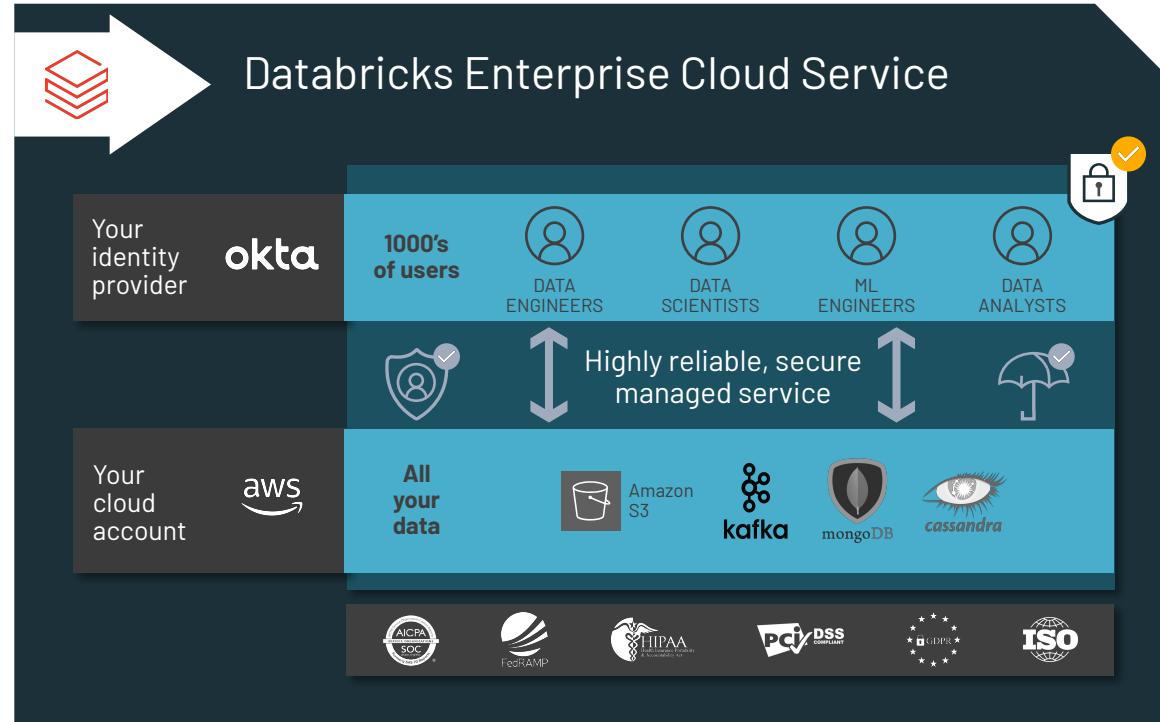
Fragmented security



Poor reliability



Disjointed governance

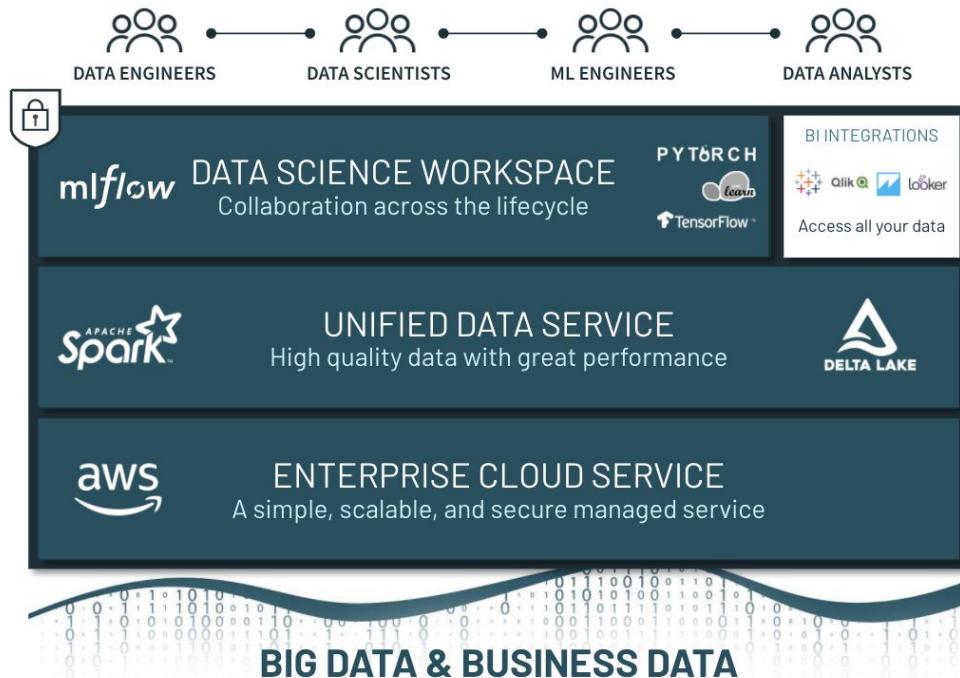


Databricks unified data analytics platform

Data Science, ML, and BI on one cloud platform

Access all business and big data in **open data lake**

Securely integrates with your **cloud ecosystem**



Databricks on AWS

Integrated with your AWS infrastructure

Integrated Data Services

AWS Glue



Amazon S3



Amazon Kinesis



Amazon Redshift



Integrated Management

End-to-End Analytics & ML



Amazon SageMaker



Amazon Quicksight



Amazon DynamoDB



Amazon Athena



AWS Security
Single Sign-On, Identity
Access Management (IAM)



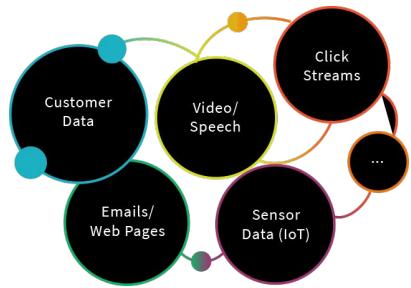
AWS Marketplace
Unified Billing



AWS Directory
Services

Our Business Value at- a-glance

Data engineering, management and transformation



LAKES



STREAMS



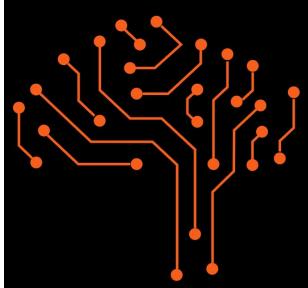
WAREHOUSES



NOSQL



Data Science, Machine Learning and Deep Learning



Data Science and Machine Learning Business Outcomes



- Recommendation & Personalization Engines
- Risk, Fraud, Intrusion Detection & Prevention
- Customer 360 Engagements , Ad Targeting
- Inventory Asset Optimization & Allocation
- Genomics & DNA Sequencing
- Predictive Maintenance, SCM Seasonal Costing
- Sentiment & Customer Churn Analysis
- Security Compliance & Intelligence

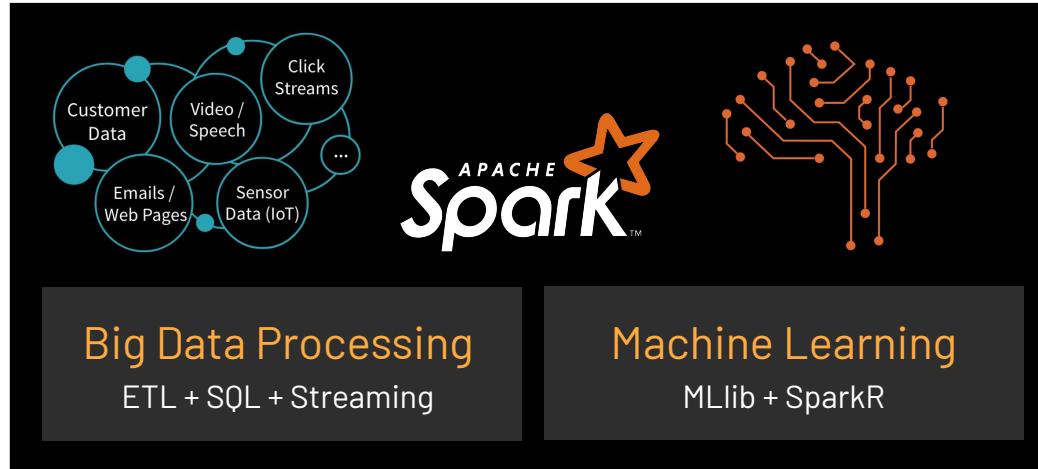
One Unified Data Analytics Platform

// Guaranteed Expertise //

Predictable Business Outcomes

Apache Spark: De-Facto Unified Data Analytics Engine

Uniquely combines Data & AI technologies



facebook.

LinkedIn

 Nasdaq

Tencent 腾讯
腾讯公司

Bloomberg

Uber



**U.S. Immigration
and Customs
Enforcement**



THE WORLD BANK
IBRD • IDA | WORLD BANK GROUP



amazon.com



 **Baidu 百度**



VIACOM



DEV DAY

 **databricks** | 



500,000
meetup members

In spite of Spark success,
companies are still
struggling with ML

Hardest Part of AI isn't AI, it's Data

"Hidden Technical Debt in Machine Learning Systems," Google NIPS 2015

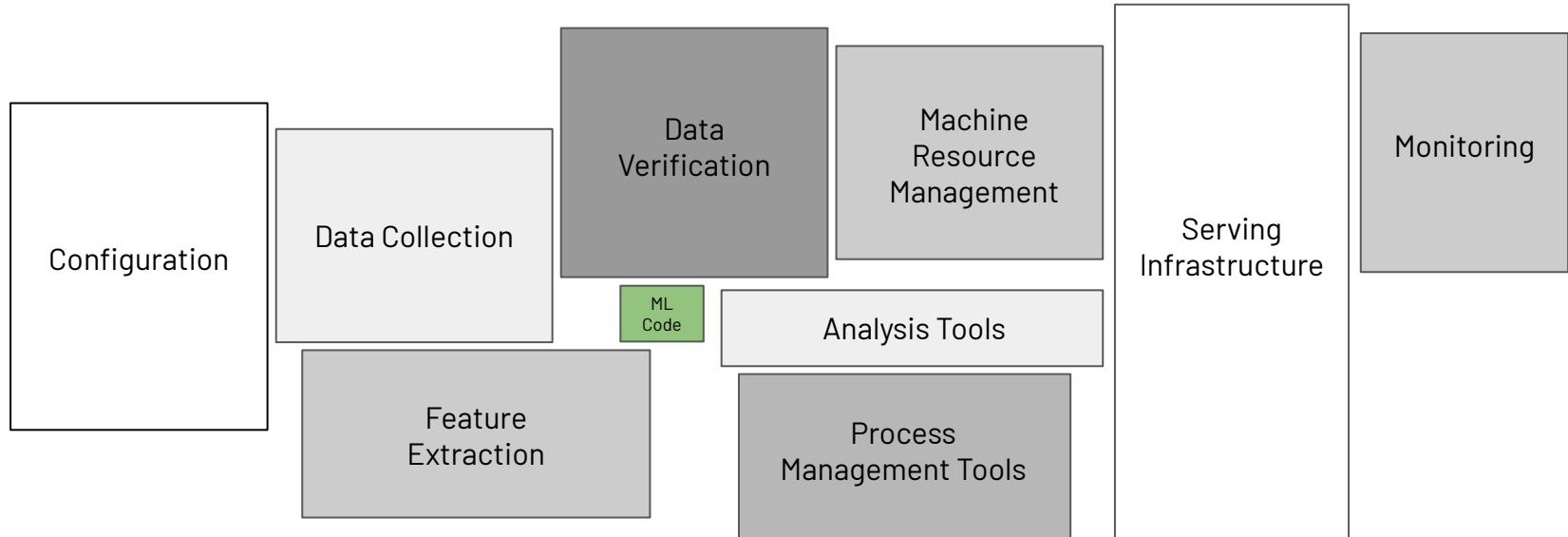
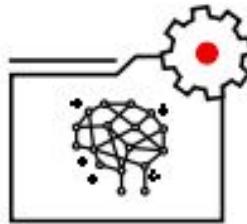


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small green box in the middle. The required surrounding infrastructure is vast and complex.

CIO Survey: AI Promise



CIO Survey: AI Dilemma



Just **17%**

have moved an AI project to
production in a core
business area



It takes an average of
6 months
to bring AI projects to deployment

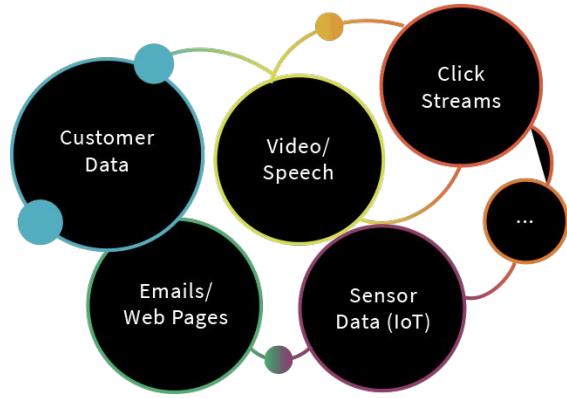
3 major challenges slowing ML Projects:

- 1 Data is not ready for AI
- 2 A Zoo of new ML Frameworks
- 3 Data Science & Engineering silos

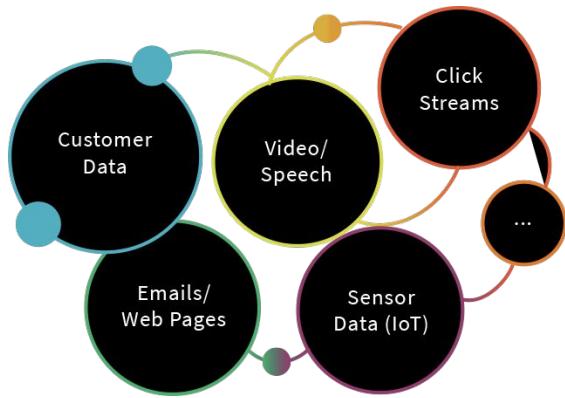
Challenge 1

Data is not ready for AI

Enterprises have been spending millions of dollars getting data into data lakes

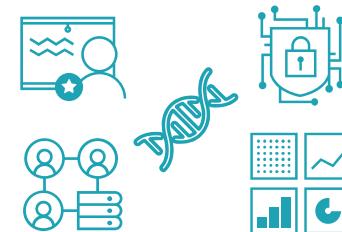


The aspiration is to do analytics & AI on all that data!



Data Lake

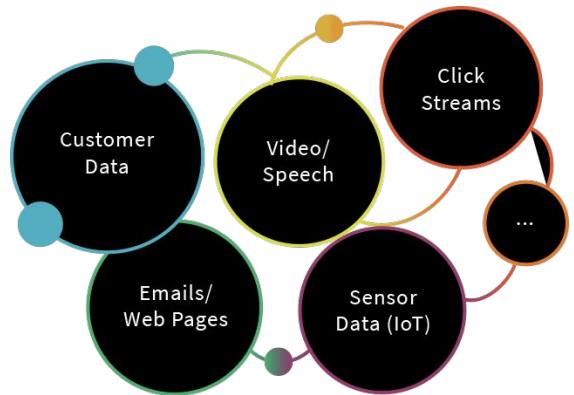
Analytics & AI



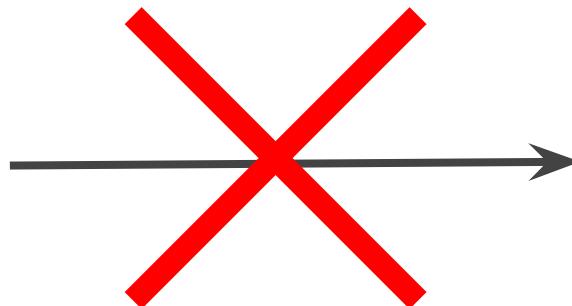
- Recommendation Engines
- Risk, Fraud Detection
- IoT & Predictive Maintenance
- Genomics & DNA Sequencing

But the data is not ready for analytics & AI

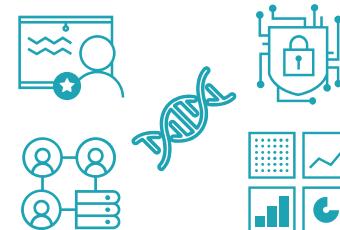
The **majority** of these projects are failing due to
unreliable data!



Data Lake



Analytics & AI



- Recommendation Engines
- Risk, Fraud Detection
- IoT & Predictive Maintenance
- Genomics & DNA Sequencing

Delta Lake: makes data ready for analytics & AI



Reliability

Performance

Analytics & AI



- Recommendation Engines
- Risk, Fraud Detection
- IoT & Predictive Maintenance
- Genomics & DNA Sequencing

A New Standard for Building Data Lakes

Open Source and Open Format



Data Reliability and Quality

Compatible with Spark APIs

Delta Lake: Analytics Ready Data

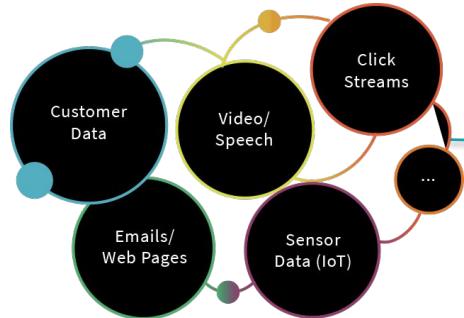
1. Data Reliability

Guarantee High Data Quality
Schema Enforcement & ACID Transactions

2. Query Performance

Very Fast at Scale
Unified Batch and Streaming
Indexing & Caching (10-100x Faster)

LOTS OF NEW DATA



DELTA LAKE

DATA LAKE

Reporting

Dashboards

Alerting

Machine Learning

Challenge 2

A Zoo of new ML Frameworks

Complexity - Zoo of ML frameworks

Machine Learning	Deep Learning	Supporting Libraries	Serving and Monitoring
Scikit-learn , Spark MLlib, H2O, Mlpack, Mahout ...	TensorFlow, Keras, Caffe, PyTorch, Theano, BigDL, SparkDL ...	Python, R, Anaconda, Numpy, Scipy, Pandas, Matplotlib, PyViz ...	MLeap, TF Serving, Cassandra, Redis, TensorBoar d ...

Databricks Runtime for ML

Ready to use clusters with built-in ML Frameworks



XGBoost



GPU support



Challenge 3

Data Scientists & Engineers
are in silos

Data Scientists & Engineers are in Silos

① Data Prep

Hard to make pipelines reliable



② Build Model

Challenging to track and reproduce experiments



③ Deploy Model

Have to ensure reliability, SLAs, and quality



kubernetes



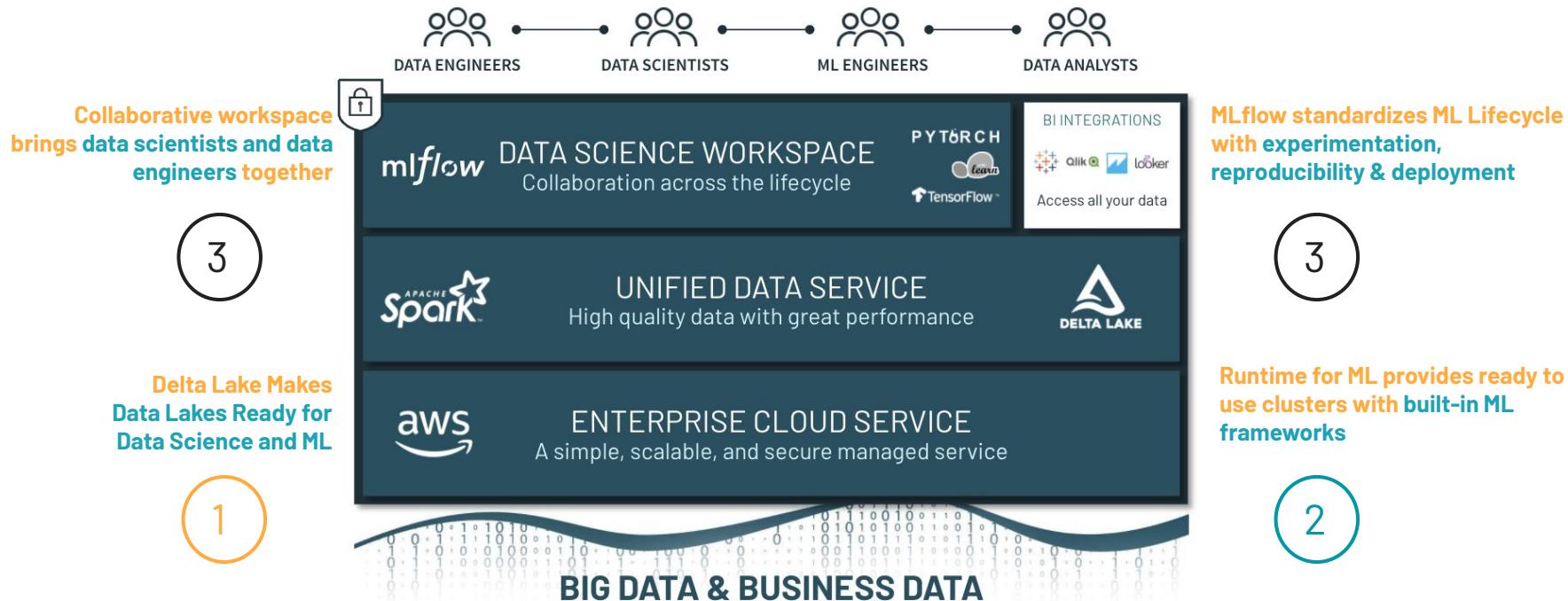
Data Engineers



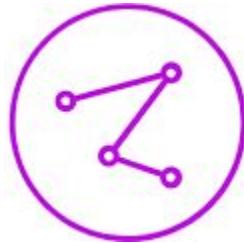
Data Scientists

Databricks Unified Data Analytics Platform

End-to-end ML platform that unifies people, processes and technologies



Customer Presentation



R E T I N A



April 22, 2020



Data Platform Checklist

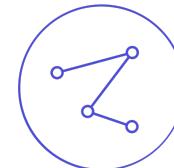
Brad Ito
CTO, Retina AI



Introduction



- 4 startups as CTO
 - Email Marketing Platform
 - Performance Marketing
 - Crowdfunding Platform
 - Retina AI
- 10 Programming Languages
- Numerous data platforms
- Physics, MIT



RETINA

- Customer Lifetime Value experts
 - Early & accurate models
- Customer Intelligence Partner
- Clients include: Capital One, Chegg, Madison Reed, Dollar Shave Club



The Problem

Business Data Pain

- Unable to use data to drive business change
- Data more chore than delight
- Untrusted data
- Expensive and complex systems
 - Being "Big Data" is not enough



Solution:

Data Platform Checklist

- See how you measure up with simple yes / no
- Broadly applicable
- Move from buzzwords to value
- Provide a "North Star" for data initiatives



1. Is your data in one place?

- Accessible via a single platform to ease data discovery
- Single source of truth - trusted data insights
- Can be queried
 - quickly - at the speed of a question
 - without adverse effects (not OLTP, scaled resources)



2. Is your data self-service?

- Technical (Python / R / SQL) users can find / access data
 - can easily share transformed data
- Analysts can find / access data via BI Dashboard
- Applies: DataOps >> Data Engineering Requests
- Meta data available
 - Business context
 - Data lineage

3. Is your data checked for validity?

- Manual checks with business context
- Automated checks
 - Developed from real data quality problems
 - Manual intervention
- Data ownership



4. Is your data resilient?

- Can you correct past data?
- Can you easily make small changes and undo them?
- Are your data transformations idempotent?
 - Event logging over state tracking
 - Universally unique ids
- Are data relationships explicit?
 - example: customer identity resolution



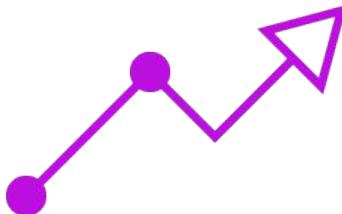
5. Are you optimized for cost / performance?

- Separate cloud compute from storage
- Schemas with appropriate data types
- Data structures
 - In-Memory
 - Columnar
- Data Storage
 - Compressed
 - De-duplicated



6. Is your data at the speed that you need it?

- Understand your data latency requirements
- Use streaming technologies for streaming data
- Schedule batch jobs appropriately
 - At least as fast as data sources
 - No faster than ability to consume data



7. Is your data platform connectable?

- Easily ingest from new data sources
 - example: [Databricks partner integrations](#)
- Easily use data with external systems
 - Queries -> Marketing Automation
 - Dashboards and Reports
 - Data Science and Machine Learning partners



Data Platform Checklist

1. Is your data in one place?
2. Is your data self-service?
3. Is your data checked for validity?
4. Is your data resilient?
5. Are you optimized for cost / performance?
6. Is your data available at the speed that you need it?
7. Is your data platform connectable?



Brad Ito
CTO, Retina AI
brad@retina.ai
www.retina.ai - retina.ai/blog

<https://www.linkedin.com/in/brad-ito/>

Partner Case Study

PARIVEDA

Modern Data Enterprise

Presented By:
Jonathan Corners and Brian Edwards

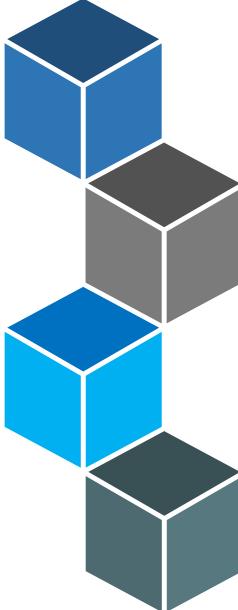
Hear Our Perspectives - <https://parivedaperspectives.com/>



Pariveda Solutions Guiding Principles

TRANSPARENCY IS KEY TO A GOOD RELATIONSHIP

Be frank, give options, and enable our customers to do what makes sense in their context.



CO-CREATE VALUE WITH OUR CUSTOMERS

Giving and sharing is the best way to uncover potential.



TECHNOLOGY IS AN ENABLER, NOT A SOLUTION

Apply the human picture when observing with a technical lens.



MAKE A DIFFERENCE

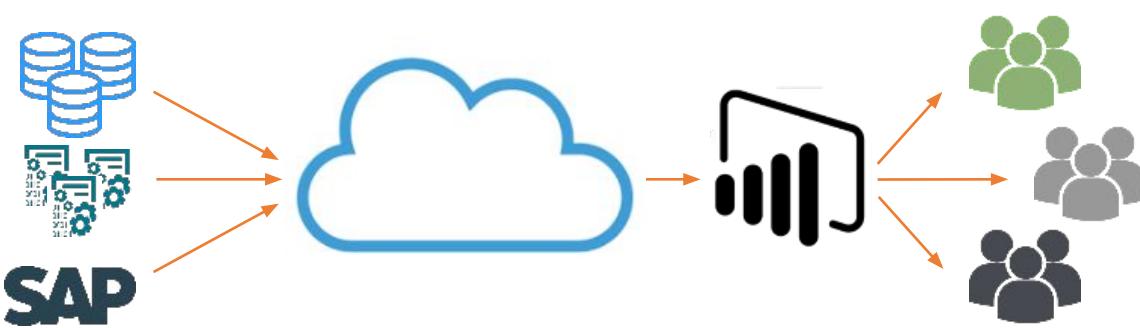
Take pride in the process and the product, but most of all strive for positive impact.

Pariveda's Modern Data Enterprise Maturity Model

	 Uncontrolled	 Reactive	 Proactive	 Resilient
Value	No organization-wide plan or vision guiding efforts and consequently, ROI is unknown or poorly justified	Current state process defined, focused on supporting current state efforts and baseline financial estimates of projected costs and initial value forecasts	Vision of future state defined, focused on incremental improvements inclusive of estimates that consider operations	Vision and guiding principles, cost and value considerations together drive all technology efforts to reach future state
Platform	Data Warehouse – Analytics are descriptive and immature	Data Swamp – Descriptive analytics, reporting, and visualization	Data Lake – Experimentation and beginning of predictive analytics	Data Platform – Optimization for increased scale of use-cases
Governance	Lack of well-defined data stewardship and management	Data stewardship defined, manage requests and issues reactively	Proactive data stewardship and management	Data stewardship and management is operationalized

The modern data enterprise is proactive and resilient with data driven insight!

Case Study: Data Driven Insights for Marketer of Consumer and Commercial Products



Situation

A Fortune 500 company that sells consumer packaged products is looking to accelerate insights by providing analysts with increased interaction and high levels of granularity on their datasets.

Data quality and data scale challenges made it difficult for the business to accurately forecast future sales and inventory needs.

The current approach was high cost for IT to maintain and difficult to integrate new data sources. A new scalable approach was required.

Solution

We worked with the client to developed a scalable, cloud based Digital Data Platform.

Using the cloud, we implemented an ingest, model, enhance, transform, deliver approach to processing large datasets and delivering data at its more granular level to 10,000 users. Automation of repeatable manual processes enabled the team to reallocate operational effort to building components that improved data insights.

As the platform matured, Databricks was introduced to enable advanced analytic scenarios at scale.

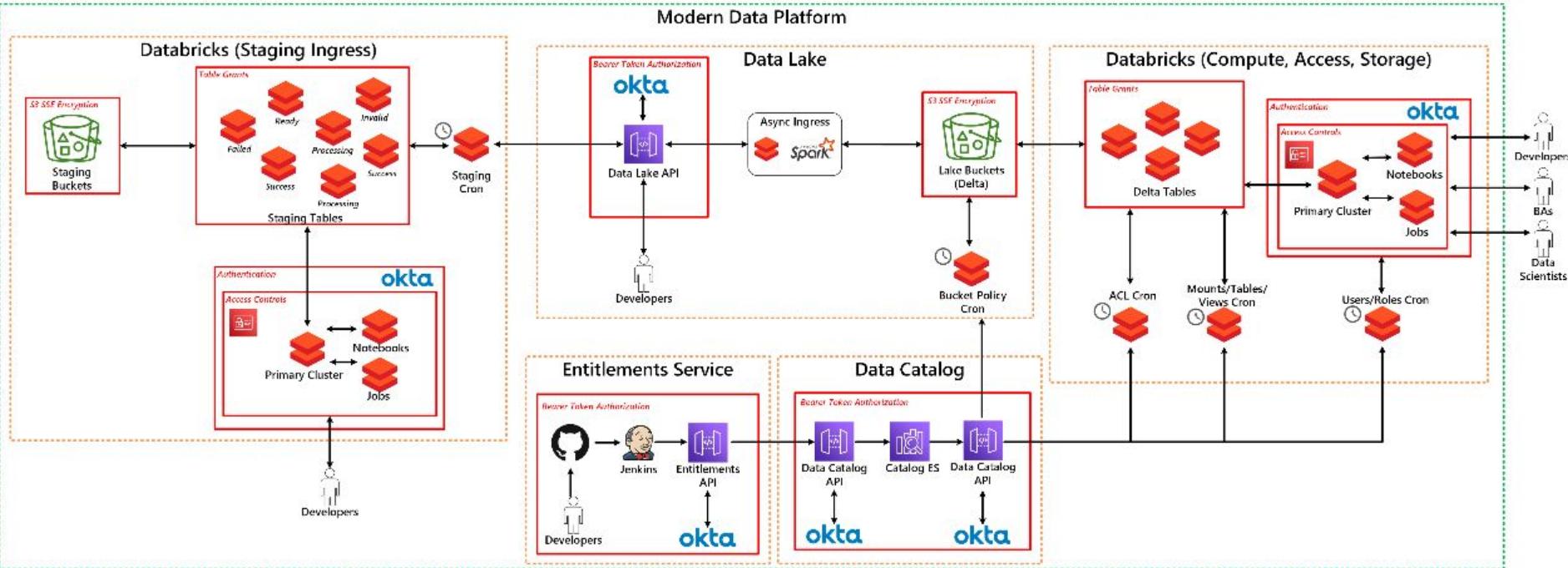
Value Provided

"Speed to insight" enabled executive decision makers to respond faster and more accurately to market changes.

Reduced IT effort on data quality issues by 25%. Reduced implementation and operational run costs by leveraging a patterned based approach for adding and processing datasets.

Data Science teams could identify and create new reports and models to make accurate predictions. Same day updates are now possible, allowing the ability to report on information across the entire enterprise.

Databricks simplifies and upskills multiple areas of the Modern Data Platform



DataForge - Component-Driven Data Platform



Ingest: Capture the data from the source (especially for ephemeral sources)

Model: Handle differences in input format for common data sets

Enhance: Add additional context, make each data record more usable

Transform: Combine and aggregate data sets to produce new data products that are ready for consumption

Deliver: Retrieve and send data sets / records to support data based applications and analytics

How Pariveda Can Help

How do I convince my boss
of the value of modern data
tools & practices?

Where should we start?
What use cases have the
best bang for the buck?

How do I turn my data
swamp into a data lake

AI/ML APPLICATION: PARTICIPANTS & OUTCOMES



Product Leads
Business Unit Heads

CTOs & CDOs
BI/Analytics Leads

IDEATION	ML PROOF OF CONCEPT	PRODUCTION INTEGRATION
<p>What high-value decisions does your business make that ML-based predictions can help with?</p> <p>What data sources are available to support models? Can AI techniques encode unstructured data sources to enable new applications?</p> <p>How complex will the ML model development be for each idea?</p> <p><i>Prioritized concept cards capturing high-value systems that can be enabled by Machine Learning</i></p>	<p>Includes 'ideation' plus an implementation of:</p> <p>What is the predictive accuracy of the ML model in aiding in high value decisions?</p> <p>How can the model, running on the AWS Platform, be executed to make manual predictions?</p> <p>How will the model's predictions be integrated into business processes?</p> <p><i>ML model accuracy report alongside plan and estimates for production application build</i></p>	<p>Includes model from PoC integrated into systems and processes:</p> <p>How do we build/maintain/operate APIs & data pipelines for ML apps?</p> <p>How will business processes change to leverage predictions that the ML POC has made easily available?</p> <p>How can I measure the business value generated by the new predictive capability?</p> <p><i>Realized business value from POC alongside plan to build Data Lake Platform and enable ML @ Scale</i></p>

PRICING

1. Best results achieved when workshop delivered over 2 days. Day one focused on ideation, day 2 on prioritization, with complexity analysis done offline between the days.

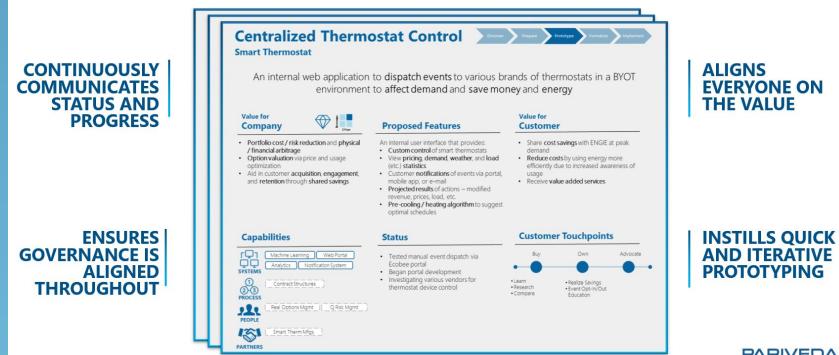
Free 1-2 Day Workshop¹

Price based on complexity of POC

Price based on POC outcomes

Concept Cards

Pariveda's Concept Card minimizes the barriers and positions your organization to have meaningful conversations around new business ideas



For more information: jonathan.corners@parivedasolutions.com



Questions

Brian.Edwards@parivedasolutions.com

Jonathan.Corners@parivedasolutions.com

SPARK+AI SUMMIT

JUNE 22-26, 2020 | ORGANIZED BY  databricks®

THE VIRTUAL EVENT FOR DATA TEAMS

- Extended to 5 days with over 200 sessions
- 4x the pre-conference training
- Keynotes by visionaries and thought leaders

NOW FREE



Migrating Hadoop to Databricks Webinar

The banner features the Databricks logo (a red cube icon followed by the word "databricks") and the AWS logo (the orange arrow icon) side-by-side. Below the logos, the word "WEBINAR" is written in red capital letters. The main title of the webinar, "Migrating On-premises Hadoop to a Cloud Data Lake", is displayed in large black text. To the right of the text, there are two black and white portrait photos of men: Denis Dubeau on the left and Brian Dirking on the right. At the bottom of the banner, there is a dark horizontal bar containing the names and titles of the speakers: Denis Dubeau (AWS Partner Solution Architect Lead, Databricks) and Brian Dirking (Sr. Director Partner Marketing, Databricks).

WEBINAR

Migrating On-premises
Hadoop to a
Cloud Data Lake

Denis Dubeau
AWS Partner Solution Architect Lead, Databricks

Brian Dirking
Sr. Director Partner Marketing, Databricks

May 21, 2020 | 10am PT

<https://bit.ly/HadoopToDB>

Agenda

- 9:00am Databricks Keynote, Cameron Kashani
- 9:45am Customer Presentation | Retina
- 10:15am Partner Presentation | Pariveda
- 10:30am Break
- 10:40am Data Engineering Interactive Demo & Best Practices:
Preparing Data for Analytics | Aaron Binns
- 11:15am Data Science Interactive Demo & Best Practices: Model
Training and Machine Learning | Lei Pan
- 11:50am Q&A
- 12:00pm Finish

Break

To follow-along with our technical demos after the break, please have your Databricks Community Edition account created. No credit card needed.

<https://databricks.com/signup/signup-community>

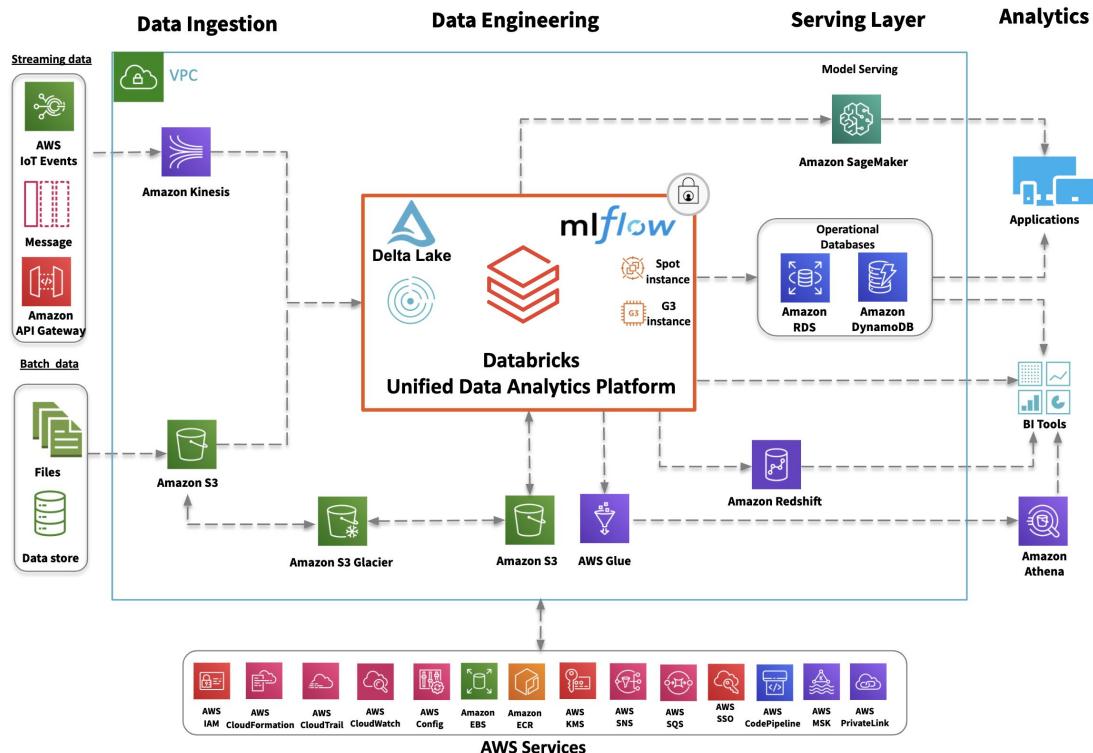
Demo Notebook links are in your reminder email.

AWS | Databricks Dev Day Quick Guide

<https://dbricks.co/aws-devdays>

Databricks Unified Data Analytics Technical Interactive Demo

Databricks & AWS data lake implementation



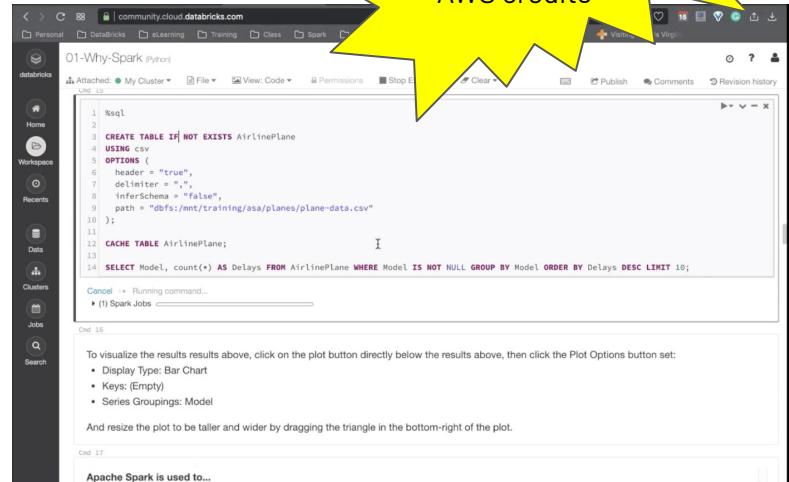
Free Training

Three virtual 2-hour training sessions using Databricks on AWS:

- Getting Started with Apache Spark
- Data Engineering and Streaming Analytics
- Machine Learning

Sign up now:

<http://bit.ly/TrainingAWS>



Speak to your Databricks team to get \$50 in free AWS credits

```
1 #sql
2
3 CREATE TABLE IF NOT EXISTS AirlinePlane
4 USING csv
5 OPTIONS (
6   header = "true",
7   inferSchema = "true",
8   inferSchema = "false",
9   path = "dbfs:/mnt/training/asa/planes/plane-data.csv"
10 );
11
12 CATCH TABLE AirlinePlane;
13
14 SELECT Model, count(*) AS Delays FROM AirlinePlane WHERE Model IS NOT NULL GROUP BY Model ORDER BY Delays DESC LIMIT 10;
```

To visualize the results results above, click on the plot button directly below the results above, then click the Plot Options button set:

- Display Type: Bar Chart
- Keys: (Empty)
- Series Groupings: Model

And resize the plot to be taller and wider by dragging the triangle in the bottom-right of the plot.

Cmd 16

Cmd 17

Apache Spark is used to...

Post Event Survey

<https://dbricks.co/mldevday>

Thank you!