



**Databricks**

Unifying Data Science  
and Engineering



# The beginnings of Apache Spark at UC Berkeley

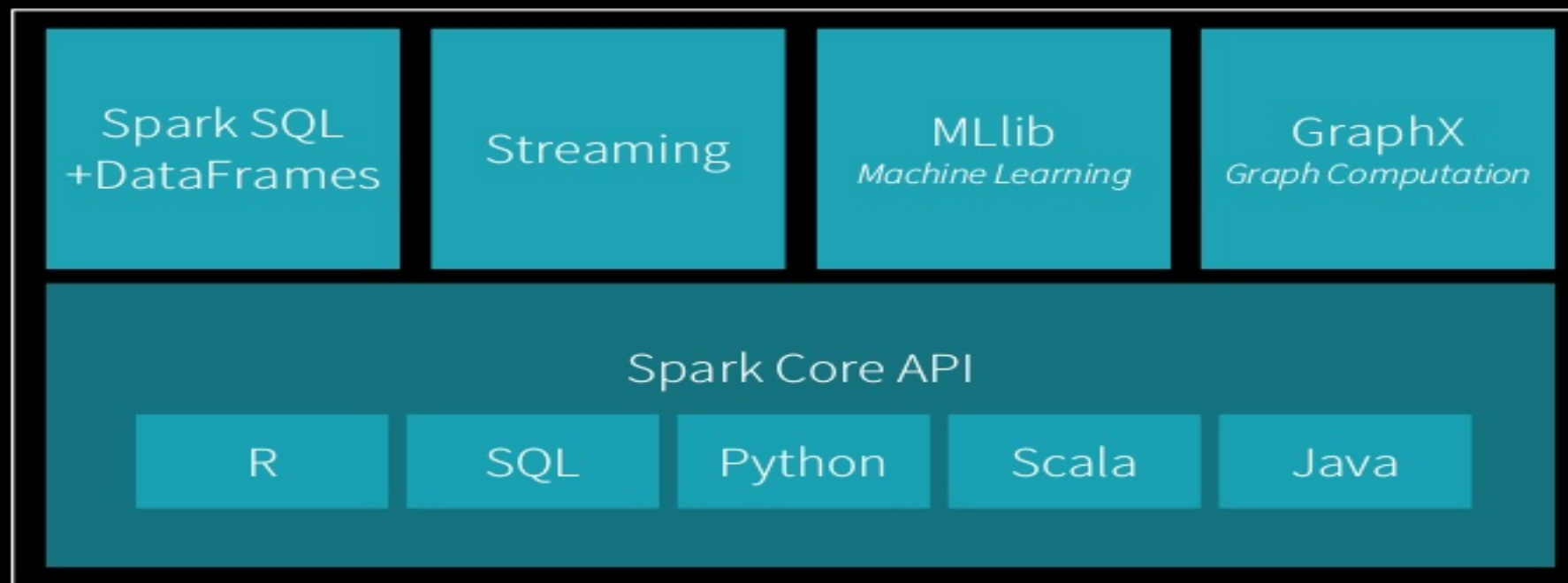
AMPLab funded by tech companies:  
**twitter** **facebook.** **Google**

- Got a glimpse at their most impactful internal projects
- They were leveraging massive amounts of data
- Doing high impact machine learning/AI

We wanted to democratize Data + AI



# Apache Spark



# Databricks started in 2013

*Bring Apache Spark to the Enterprise*

Only 1% of enterprises successful with AI

Google facebook. Microsoft

twitter amazon NETFLIX UBER

# Other 99% struggle due to organizational silos

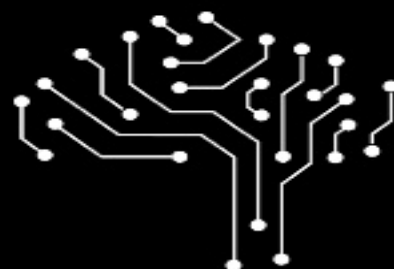
## IT



Data Engineers



## Line of Business



Data Scientists



 databricks

Databricks goal is to unify  
data science & engineering



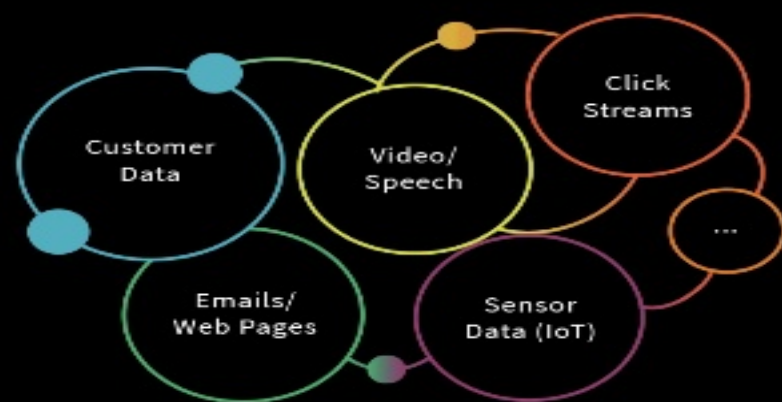
## 3 challenges created by data & AI divide

- ① Data is not ready for AI
- ② Data and AI technology silos
- ③ Data scientists and engineers are in silos



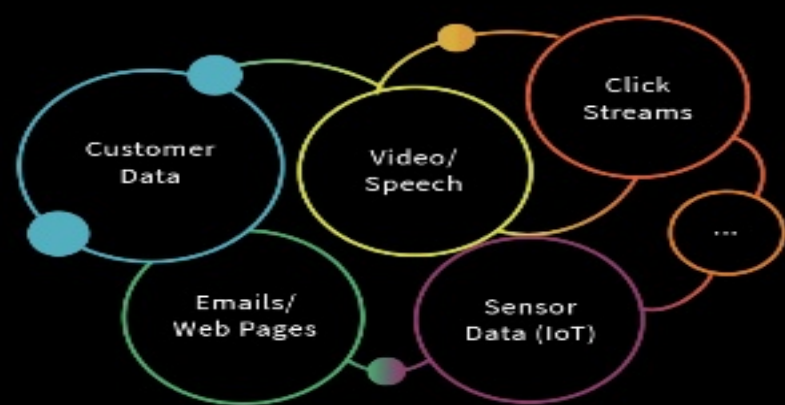
① Data is not ready for AI

# Massive data in data lakes

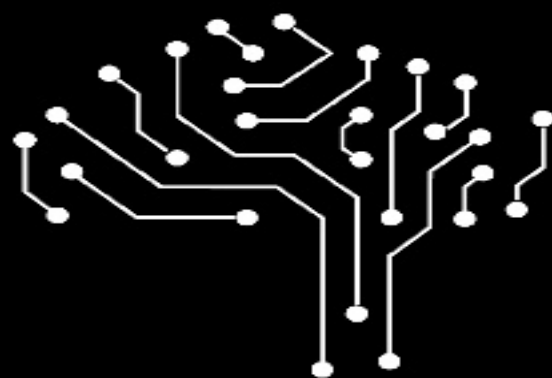


**Data Lake**

# Vision to do AI on that data

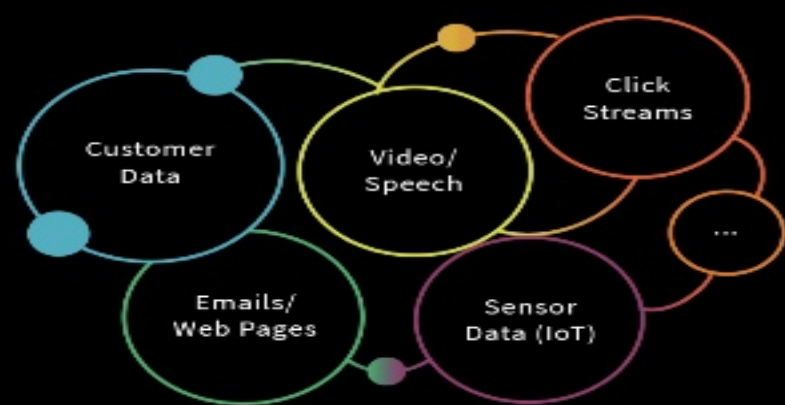


**Data Lake**



**AI**

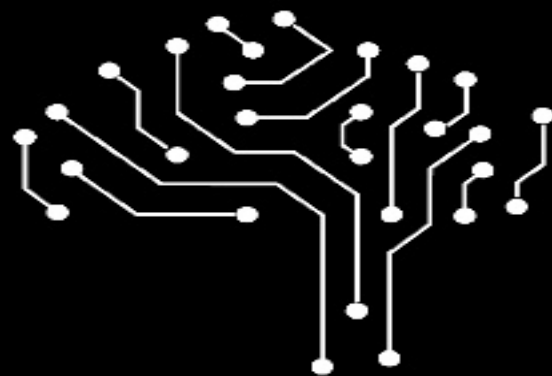
# Data is not ready for AI



**Data Lake**



Inconsistent Data  
Lack of Schema  
Slow Performance and Costly



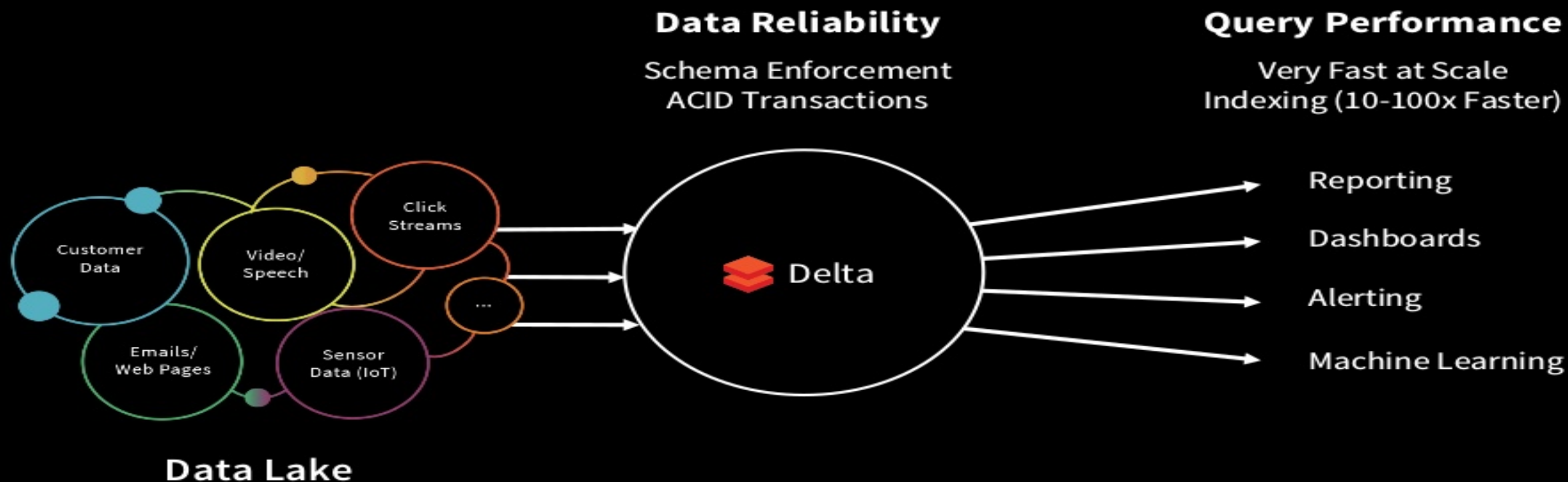
**AI**

# Databricks Delta

*Brings data reliability and performance to data lakes*



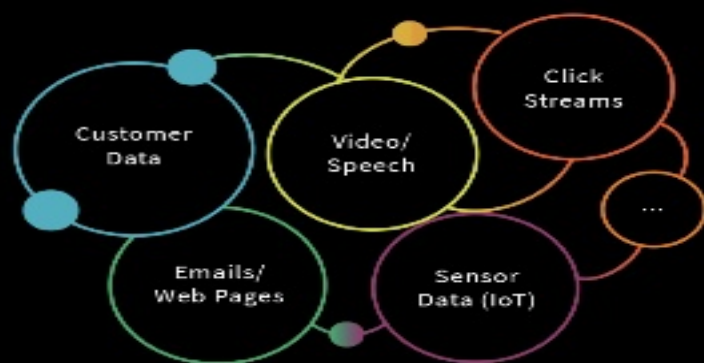
# Databricks Delta: makes data ready for AI



## ② Data and AI technology silos

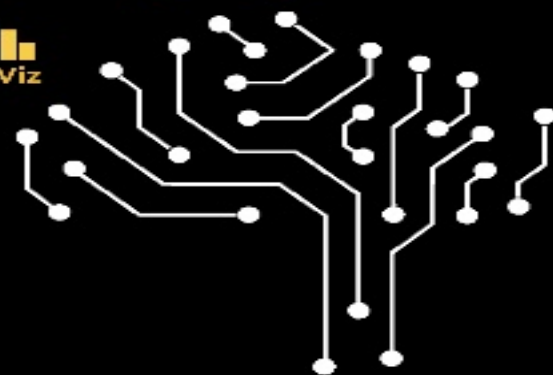


# Data & AI technology silos



Great for Data, but not AI

## Supporting and Deployment Libraries



Great for AI, but not for data

# Data & AI technology silos

Supporting and Deployment Libraries

## Databricks Runtime for ML

Ready to use clusters with built-in ML Frameworks

TensorFlow™

learn

K

XGBoost

PYTORCH



NumPy

GPU support

Azure

aws

Great for Data, but not AI

Great for AI, but not for data

③ Data scientists and engineers are in silos

# Data scientists & engineers are in silos

## ① Data Prep

Hard to make pipelines reliable



## ③ Deploy Model

Have to ensure reliability, SLAs, and quality



kubernetes



docker

Data Engineers



## ② Build Model

Challenging to track and reproduce experiments

GitHub

CONDA



TensorFlow



Excel



Jupyter

XGBoost

Data Scientists



# Databricks MLflow: unifies data scientists & engineers



Data Engineers



Data Scientists



# Databricks MLflow: unifies data scientists & engineers

## ① Data Prep

Build reliable data pipelines  
Track the datasets

**Databricks Delta**

## ③ Deploy Model

Deploy models in production,  
track their quality

**MLflow Serving**

**Data Engineers**



## ② Build Model

Track Experiments  
Reproduce experiments

**MLflow Project & Tracker**  
**Databricks Runtime for ML**

**Data Scientists**



# Databricks MLflow: unifies data scientists & engineers

## ① Data Prep

Build reliable data pipelines  
Track the datasets

**Databricks Delta**

## ③ Deploy Model

Deploy models in production,  
track their quality

**MLflow Serving**

Data Engineers

## ② Build Model

Track Experiments  
Reproduce experiments

**Databricks Runtime for ML**  
**MLflow Project & Tracker**

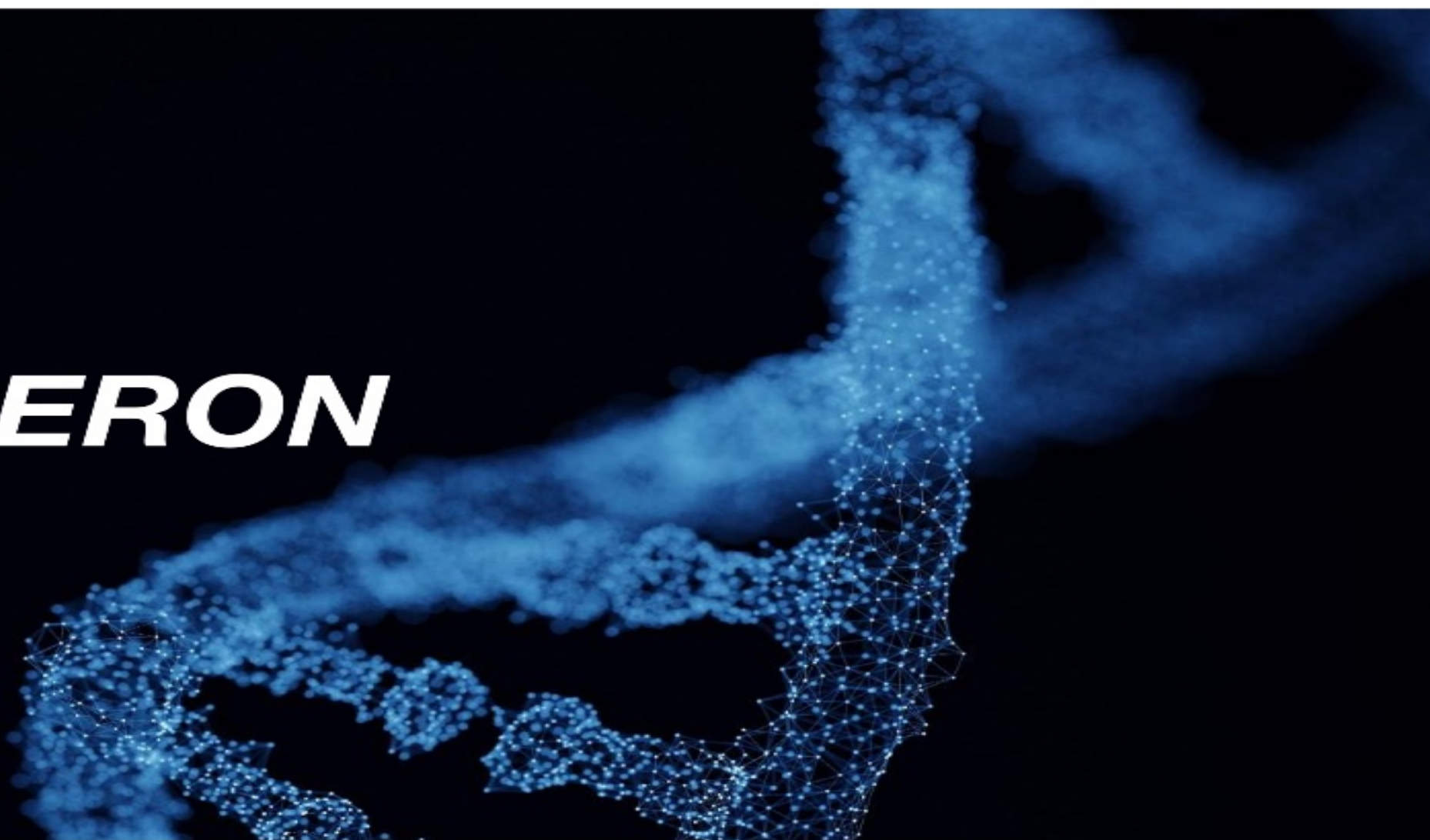
Data Scientists

Announcing: **time travel**





***REGENERON***





Hotels.com<sup>TM</sup>



Databricks makes AI possible  
for the other 99% by  
unifying data science and engineering





# Databricks Unified Analytics demo by **Michael Armbrust**

