# Disrupting Big Data with Apache Spark in the Cloud

Ali Ghodsi

databricks™

# The Dawn of Advanced Analytics

**Self-driving cars**
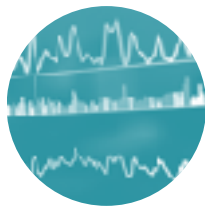
**SIRI/assistants**

**Watson**
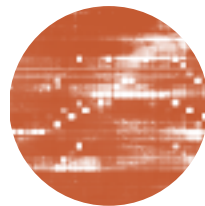
**Not just sci-fi, important applications for businesses**

databricks

# Analytics Transforming Industries

### Predictive analytics

### Anomaly Detection

Predict Product Revenue
Customer Assessment
Targeted Advertising

Fraud Detection
Risk Assessment
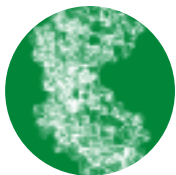Equipment Failure

**Data-Driven Real-time Analytics Applications**

# Today's Data Reality

DATA WAREHOUSES
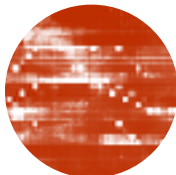
HADOOP DATA LAKES DATA HUBS

CLOUD STORAGE

## Siloed, Fast-Growing Size, Cost
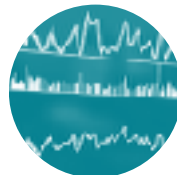
databricks

# The Analytics Gap

## Real-time Data-Driven Analytics Applications

Pharma

Media

Industrial

DATA WAREHOUSES

HADOOP DATA LAKES DATA HUBS

CLOUD STORAGE

## Siloed, Fast-Growing Size, Cost

databricks

# Why is there a gap?

## Real-time Data-Driven Analytics Applications

**(1) Manage Data infrastructure**

- Create, tune, monitor **compute clusters**.
- **Securely** access silos of **disparate data sources**.
- Enforce **proper data governance**.

**(2) Empower teams to be productive**

- Securely share big data clusters among analysts.
- **Interactively explore** data and prototype ideas.
- Debug, troubleshoot, version-control big data applications.

**(3) Establish Production-Ready Applications**

- Setup **robust data pipelines** for ETL/ELT.
- **Productionize real-time** applications with HA, FT.
- Build, serve, maintain advanced machine learning models.

## Siloed, Fast-Growing Size, Cost

# Databricks Cloud-Hosted Platform

**1** **Just-in-Time Data Platform**

**2** **Integrated Workspace**

**3** **Automated Apache Spark Management**

- Separate compute & storage

- Integrate existing data stores

- Efficient cache on first access

- Interactive notebooks, dashboards, reports

- Real-time exploration, machine learning, graph use cases

- Workflow scheduler for ML, streaming, SQL, ETL

- High availability, fault-tolerant, performance-optimized

**Agile**

**Democratize Big Data**

**Production-Ready**

databricks

# The Challenge of Securing Analytics

**End-to-end security a challenge for enterprises**

| Securing file management | Secure table management | Secure cluster management | Secure job workflows | Secure dashboards, report, notebook management |

**Today there are piecemeal solutions, but no comprehensive solution**

# Databricks Enterprise Security (DBES)

**Holistic end-to-end security for Data Analytics**

Files      Tables      Clusters      Workflows      Notebooks, Dashboards, Reports

**DBES provides**
- **Role-based access control**
- **Auditing and governance**
- **Integrated identity-management**
- **Encryption on-disk and on-the-wire**

**The First End-to-End Security Solution for Apache Spark**

databricks

# Enterprise use-cases

Preventing credit card fraud

Predict energy demand based on massive weather data

Predict player churn, predicting network outages

Natural language processing to extract author graph

Generating tailored programs based on big data

Thank you.

databricks™

# Try Apache Spark with Databricks

Try latest version of Apache Spark and preview of Spark 2.0

[http://databricks.com/try](http://databricks.com/try)