

# Apache Spark: The New Language Of Analytics



Presented By:



In Partnership With:



# About WiseWithData

- Founded in 2015, WiseWithData is Canada's leading data science consulting firm focused on scalable open-source analytics technology
  - Management Consulting
  - Technical Consulting
  - Analytical Products & Solutions
  - Education
- We partner with leading technology firms like Mariner Innovations
- Some of our recent clients include:

**Deloitte.**



MD Financial  
Management  
CMA Companies

**Third  
Bridge.**



Health  
Canada



Immigration, Refugees  
and Citizenship Canada



Employment and  
Social Development Canada



Indigenous and  
Northern Affairs Canada



HOUSE OF COMMONS  
CHAMBRE DES COMMUNES

**Service  
Canada**



# Solving Problems With Data



## Horizontal Solutions

- BI/Data Visualization
- Data Warehousing
- Forecasting
- Fraud & Security Intelligence
- IT Service Management
- Marketing Optimization
- Predictive Maintenance
- Risk Management

## Vertical Solutions

- Finance / Banking
- Financial Accounting & Audit
- Government
- Healthcare
- Insurance
- Logistics
- Retail
- Telecom

# Data Science In A Nutshell

## Descriptive Analytics

- BI and Reporting
- Data Visualization
- Advanced Search
- Graph Processing
- Natural Language Processing

## Predictive Analytics

- Predictive Models
- Forecasting
- Classification
- Clustering
- Recommendation

## Prescriptive Analytics

- Test & Learn Strategies
- Simulations
- Decision Analysis
- Lifetime Valuations / NPV models
- Operations Research / Optimization

## Data Governance

- Data Security and Ownership
- Data Privacy
- Master Data Management
- Metadata Management
- Data Monitoring

## Data Management

- Data Profiling
- Data Quality & Cleansing
- Data Matching
- Data Integration / Transformation
- Data Warehousing

# Why It's Time To Ditch "Big Data"

- The term is no longer relevant
- It's now synonymous with analytics in general, which denigrates the enormous value you can get from "Small Data"
- Data size should not dictate your choice of analytics software or your approach to solving a business problem



# And Start Talking About

## Analytics Software Tools That Are...

- Open and Free
- Fast and Efficient
- Simple To Use
- Scalable and Fault-tolerant
- Comprehensive In Capabilities

# Data Science Platforms

An Ideal Data Science Platform Would Provide

- A full set of capabilities for all analytics uses
- Little need to copy or duplicate data
- Integrated workflows within the same tool and seamless of integration of multiple tools
- The ability to quickly iterate

# A Brief History Of Data Science Platforms



Expensive – Proprietary  
Slow – Unscalable

Inexpensive – Open-Source  
Fast – Scalable

- 1976 - SAS
- 1968/1969+ IBM  
Cognos, SPSS, etc.

- 2012/13 Hadoop 1/2
- 2014 - Spark 1.0
- 2016 - Spark 2.0

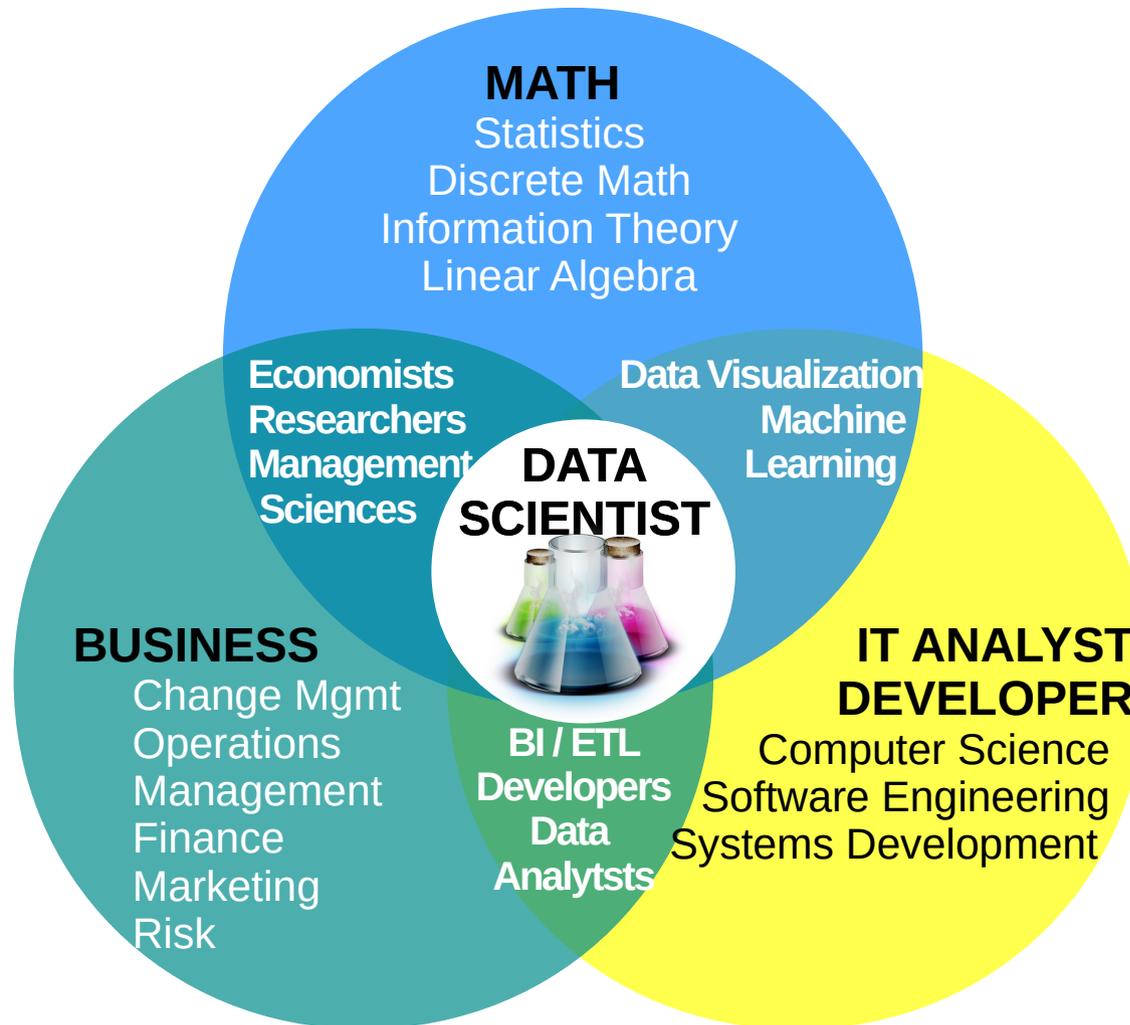
**IBM:** A large proprietary collection of over 17 vendor products purchased by IBM between 2005 and 2011. Core software designed in 1960's (SPSS-1969, Cognos-1968).

**SAS:** A proprietary data science platform. It contains a wide variety of analytic tools and methods, mostly developed in-house.

**Hadoop:** A scalable open-source software platform for scalable, distributed computing. Through a set of ecosystem tools, Hadoop can provide analysis of both structured data and unstructured data.

**Apache Spark:** A open-source high-performance scalable data processing engine. It handles structured & unstructured data, and contains advanced analytical & graph capabilities. An ecosystem of associated tools provide visualization.

# Why Data Scientists Are So Rare



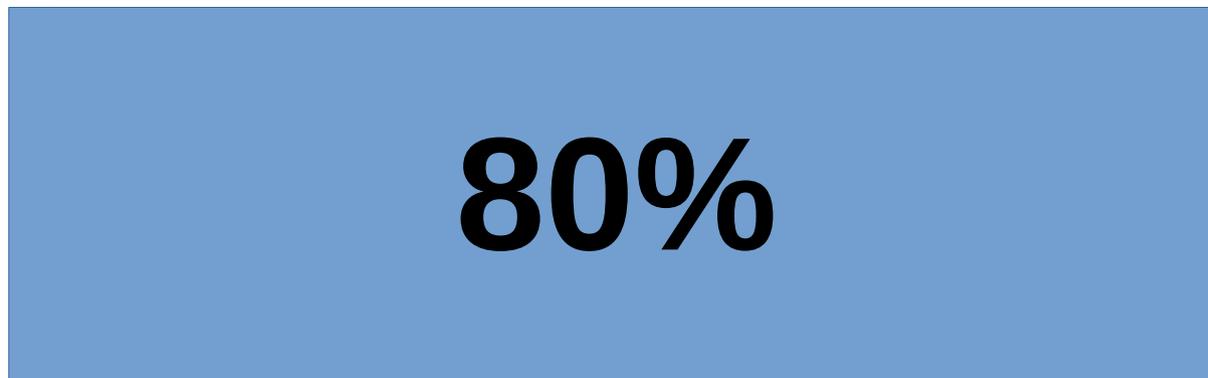
# Data Science Professionals Have The Need For Speed



**Business  
Problem**



**Solution**



**80%**



**20%**

**Data preparation is  
80% of the effort**

**Solving the problem  
is 20% of the effort**

## The Need For Speed

**Simple, Powerful API's + Fast Execution = Rapid Iterations  
Rapid Iterations = Agile High Quality Solutions**

# Spark Is A Data Scientist's Dream

Invented by Canadian computer scientist Matei Zaharia, Apache Spark's a revolutionary data science & analytics engine

- **100's of times faster** and more efficient than legacy systems
- **Fully scalable and fault tolerant**, it runs fast and efficiently on everything from laptops to massive clusters of 1,000's of servers
- **Simple easy-to-use coding interfaces** (API) including SQL which can handle streaming or batch applications in the same code
- **Built-in libraries** for data access, streaming, data integration, graph processing, and advanced analytics
- **2<sup>nd</sup> Most active open-source project (after Linux)** - commercially supported by WiseWithData, Databricks, Hortonworks & Others

# Disruptive Spark Adoption

Web/Retail

Industrial/  
High-Tech

Financial

Media/  
Entertainment

Life Sciences

Utilities/NGO/  
Government



tripadvisor

YAHOO!

Tencent 腾讯



amazon



Autotrader



TOYOTA



AUTODESK

ciena

verizon



Adobe



ERICSSON



BARCLAYS



ING



Life Financial

TWO SIGMA



JPMorganChase

nielsen



NETFLIX



ELSEVIER

Bloomberg  
VIACOM

Why?

- ✓ Rapid Low Cost Deployment
- ✓ Powerful Yet Simple Interfaces
- ✓ Exceptional Speed & Productivity



Beth Israel Deaconess  
Medical Center



NOVARTIS

aetna



QuintilesIMS



U.S. Citizenship &  
Immigration  
Services

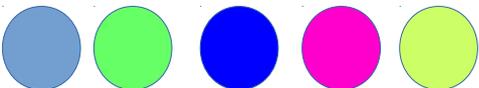
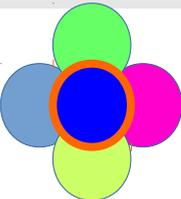
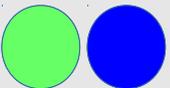


HOUSE OF COMMONS  
CHAMBRE DES COMMUNES



Immigration, Refugees  
and Citizenship Canada

# Why Is Spark So Disruptive?

	<u>Proprietary / Legacy</u>	<u>Spark</u>	<u>Impact</u>
Hardware	\$\$\$\$\$\$\$\$\$\$\$\$	\$	<ul style="list-style-type: none"> <li>• Faster TTV*</li> <li>• Improved ROI</li> </ul>
Software Licence	\$	Free + Optional Support (\$\$\$)	<ul style="list-style-type: none"> <li>• Faster TTV*</li> <li>• Flexibility</li> </ul>
Source Code	Closed	Open	<ul style="list-style-type: none"> <li>• Security</li> <li>• Flexibility</li> <li>• No vendor lock-in</li> </ul>
Tools / Applications			<ul style="list-style-type: none"> <li>• No Integration Mess</li> <li>• Faster TTV*</li> </ul>
Execution Speed	1X	100X+	<ul style="list-style-type: none"> <li>• Faster TTV*</li> <li>• Faster &amp; Better Results</li> </ul>
Productivity	1X	5X+	<ul style="list-style-type: none"> <li>• Better Results</li> <li>• Lower HR costs</li> </ul>
Streaming / Batch			<ul style="list-style-type: none"> <li>• Re-usability</li> <li>• Flexibility</li> </ul>

\*TTV=Time To Value

Time To See Spark In Action



# The Secrets Of Spark's Success



- **Be Fault Tolerant, Not Wasteful** – Replicate lineage, not work



- **Use Your Memory** – Reduce storage I/O usage by leveraging memory



- **Be Lazy** – Use a lazy execution model to plan and optimize work in advance



- **Go Columnar!** - Use partitioned columnar data structures like Parquet



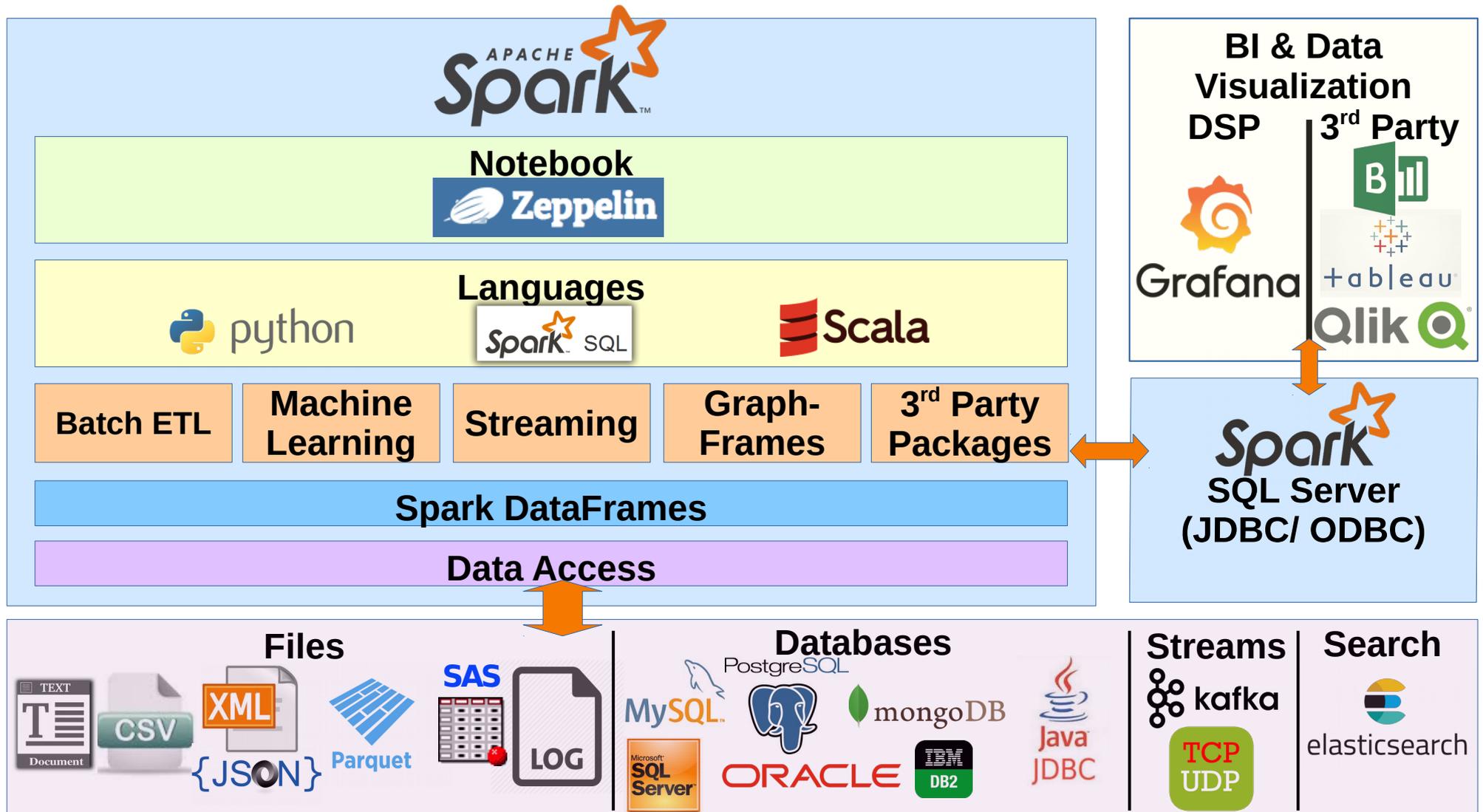
- **Keep It Simple Stupid (KISS)** – DataFrames are powerful, simple, and easy to learn



- **Speak Your Language** – You can choose between Python, SQL, Scala, R, or Java



# Making Spark Simple To Deploy & Use WWD Data Science Platform (DSP)

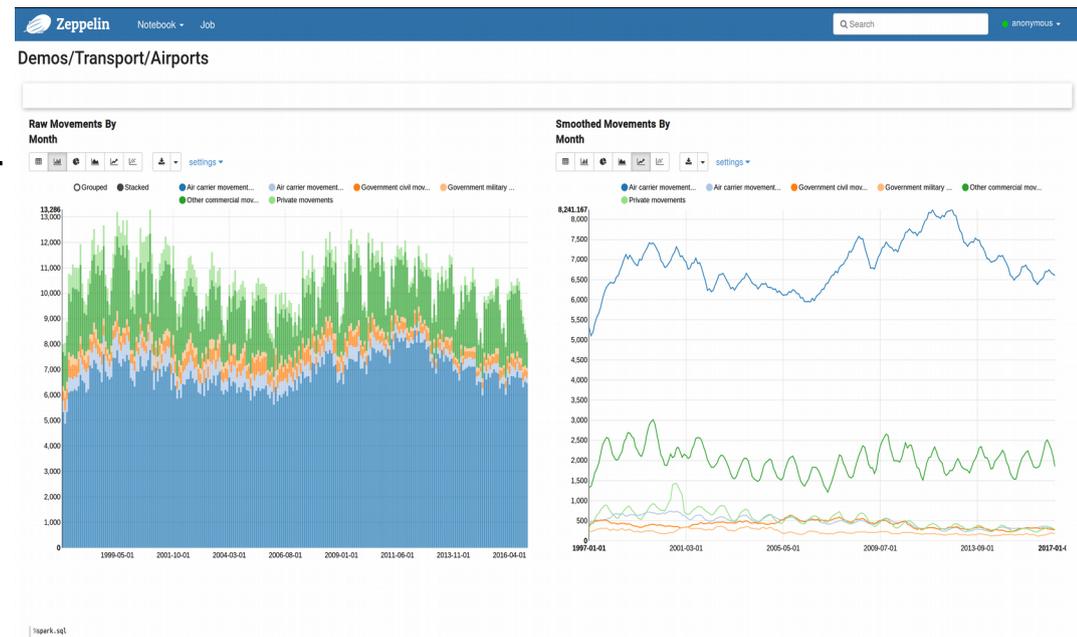


# DSP Spotlight - Apache Zeppelin

## Simple & Powerful Spark Notebook

### Apache Zeppelin

- A multi-purpose Notebook interface runs on top of Spark for data science use cases such as:
  - Documentation
  - Data Ingestion
  - Data Discovery
  - Advanced Analytics
  - Data Visualization
- Makes using Spark simple and painless
- Create simple dynamic forms for beautiful analytic applications



# DSP Is The Foundation Of Our Leading Tools & Solutions

Industry Solutions	<p style="text-align: center;"><b>GuardDog</b></p> <p>Full featured security intelligence suite. It includes business rules, anomaly detection, predictive models, network analysis and investigative capabilities for:</p> <ul style="list-style-type: none"> <li>• Fraud, Waste and Abuse</li> <li>• Public Safety and Border Security</li> <li>• Anti-Money Laundering</li> <li>• IT Operations Management</li> </ul>		<p style="text-align: center;"><b>LogHaus</b></p> <p>Log ingestion and data warehousing solution. It efficiently transforms log files or streams into highly structured data for reporting. Supported log formats include:</p> <ul style="list-style-type: none"> <li>• Firewall Logs</li> <li>• Web Application Logs</li> <li>• Windows Event Logs</li> <li>• Generic SYSLOG formatted logs</li> </ul>	
	Analytic Tools	<p style="text-align: center;"><b>MatchBox</b></p> <p>Comprehensive data quality tool for matching, standardization and entity resolution. It can process address, personal identity and business name data.</p>	<p style="text-align: center;"><b>PreViso</b></p> <p>Business forecasting tool for developing large scale hierarchical forecast models. It includes model reconciliation and monitoring capabilities.</p>	<p style="text-align: center;"><b>ModelEyes</b></p> <p>Tool for model monitoring and governance. Provides standardized reports for assessing the stability and validity production models.</p>
Foundation		<b>WWD Data Science Platform</b>		
	<b>Data Management &amp; Governance</b>	<b>Data Integration &amp; Transformation</b>	<b>Data Visualization &amp; Reporting</b>	<b>Advanced Analytics</b>

# Fast, Simple, Scalable Analytics

## Your Spark Journey Starts Today

### **WiseWithData & Mariner Innovations**

- Spark based analytical products and solutions
- Spark advisory, deployment, development, education, and support services
- WiseWithData.com / MarinerInnovations.com
- Contact Us:
  - Ian Ghent (Ian.Ghent@WiseWithData.com)
  - Karen Fegarty (Karen.fegarty@marinerinnovations.com)

### **Spark+AI Summits**

- San Francisco, USA - June 4-6, 2018
- London, UK – October 2-4, 2018

### **Spark Online Resources**

- spark.apache.org
- Apache Spark YouTube Channel
- Apache Spark StackOverload list