# Microsoft ML for Apache Spark

Unifying Machine Learning Ecosystems at Massive Scales

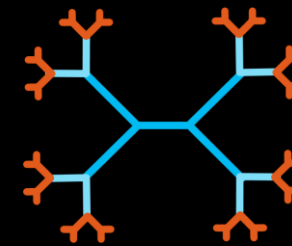**Mark Hamilton**

**Microsoft, MIT**

**marhamil@microsoft.com**

# Overview

- Background
  - Spark + SparkML
  - MMLSpark
- Unifying ML Ecosystems
  - LightGBM, CNTK, Vowpal Wabbit
  - Multilingual Bindings
- Microservice Orchestration
  - Cognitive Services on Spark
- Model Deployment with Spark Serving
- Use Cases
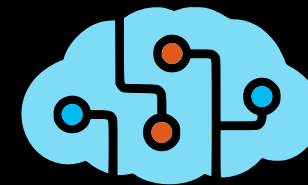  - The Snow Leopard Trust

APACHE Spark™

Vowpal Wabbit

LightGBM

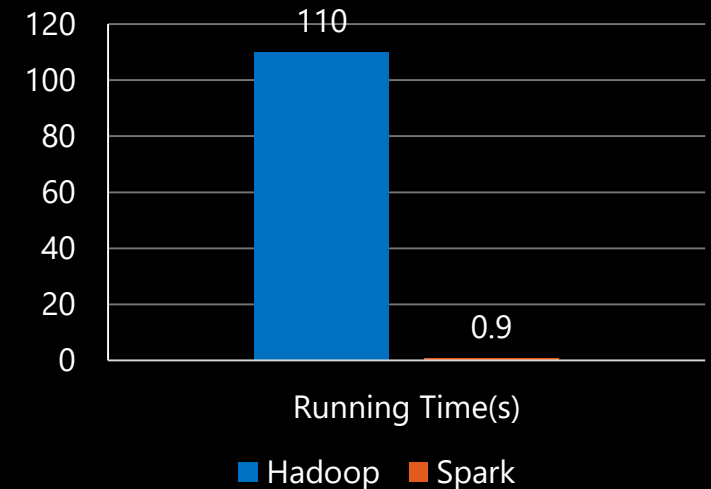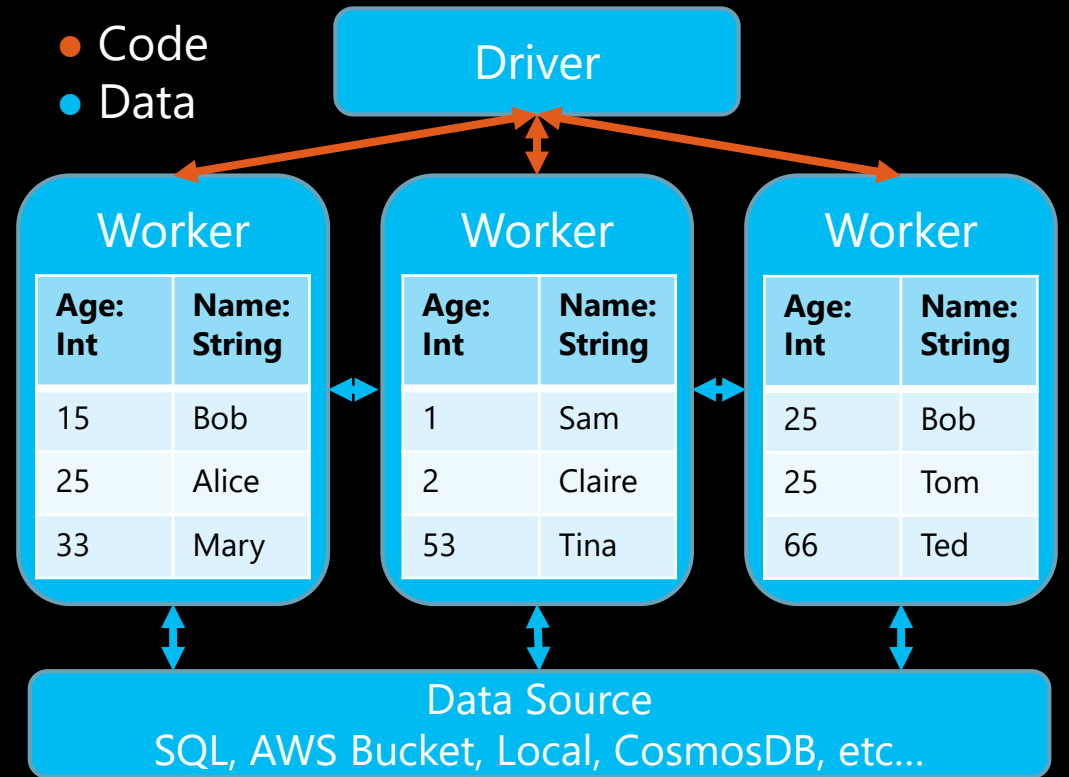CNTK

Cognitive Services

Kubernetes

Snow Leopard Trust

THE MET

# APACHE Spark™
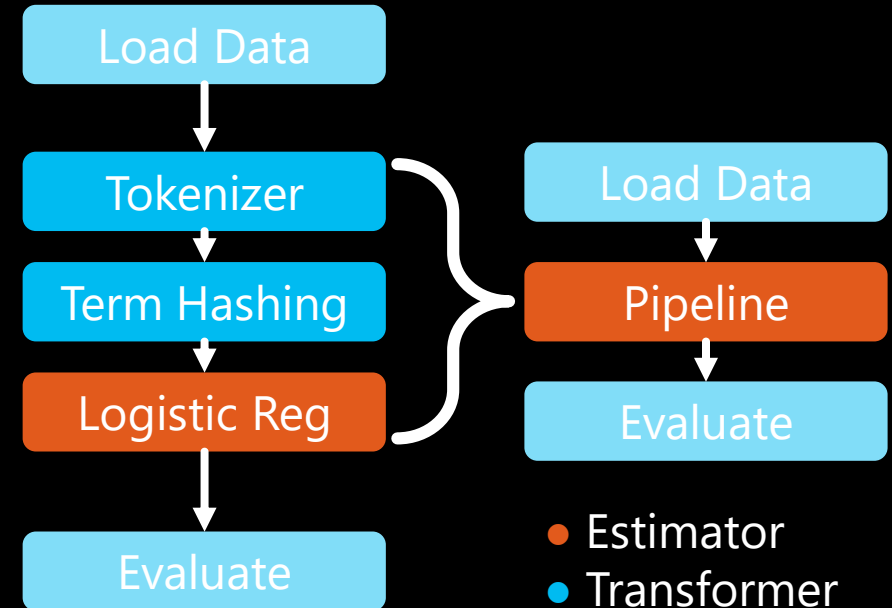
- A **fault-tolerant distributed** computing framework

- Map Reduce + SQL

- Whole program optimization + query pushdown

- Elastic

- Scala, Python, R, Java, Julia

- ML, Graph Processing, Streaming



- Code
- Data

| Driver |

| Worker | | Worker | | Worker |

| Age: Int | Name: String |
|---|---|
| 15 | Bob |
| 25 | Alice |
| 33 | Mary |

| Age: Int | Name: String |
|---|---|
| 1 | Sam |
| 2 | Claire |
| 53 | Tina |

| Age: Int | Name: String |
|---|---|
| 25 | Bob |
| 25 | Tom |
| 66 | Ted |

Data Source
SQL, AWS Bucket, Local, CosmosDB, etc…

Running Time(s)

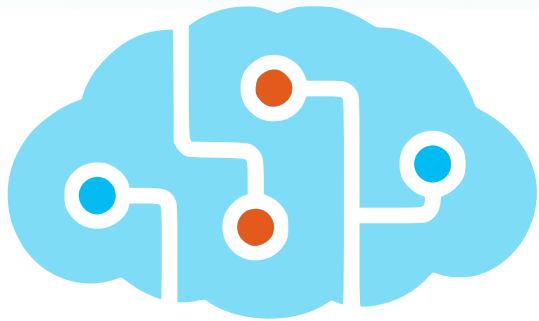- Hadoop  - Spark

110

0.9

# Apache Spark ML

- ▶ High level library for distributed machine learning
- ▶ More general than SciKit-Learn
- ▶ All models have a uniform interface
  - ▶ Can compose models into complex pipelines
  - ▶ Can save, load, and transport models

```
data = spark.read.csv("hdfs://...")
train, test = data.randomSplit([.5,.5])
model = LogisticRegression().fit(train)
predictions = model.transform(test)
```

# Unifying Machine Learning Ecosystems

- Goals
  - Same API
  - Composable
  - Batch, Streaming, Serving
  - Elastically Distributed
  - Fault Tolerant
  - Multi-Language
  - Data Source Agnostic

Markus Cozowicz
marcozo@microsoft.com
Data Scientist

Image Processing with Open CV

Gradient Boosting with LightGBM

Deep Learning Pipelines (Databricks)
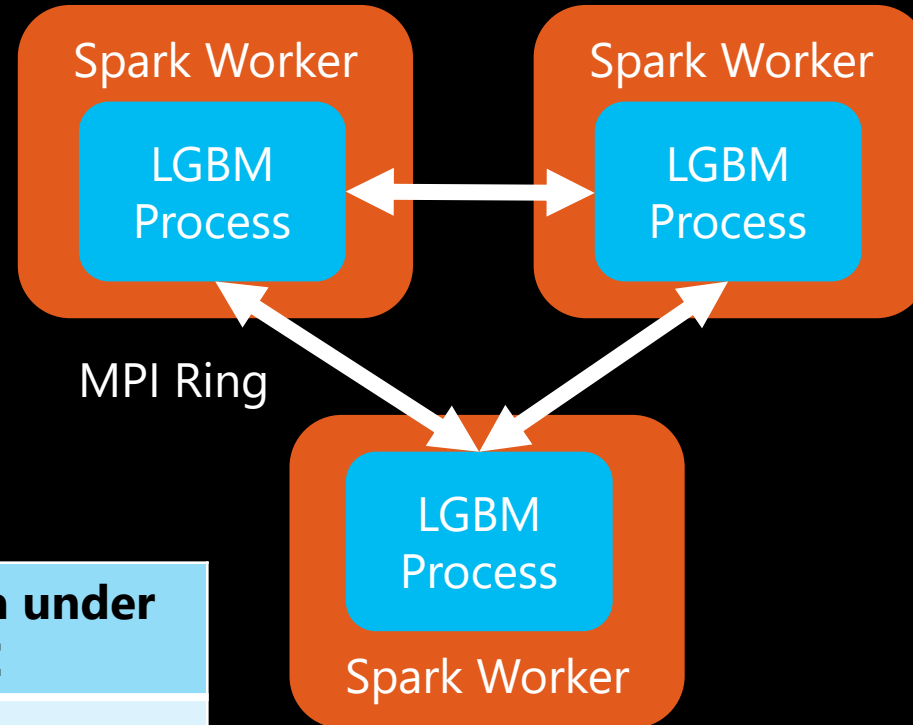
**Text Analytics with Vowpal Wabbit**

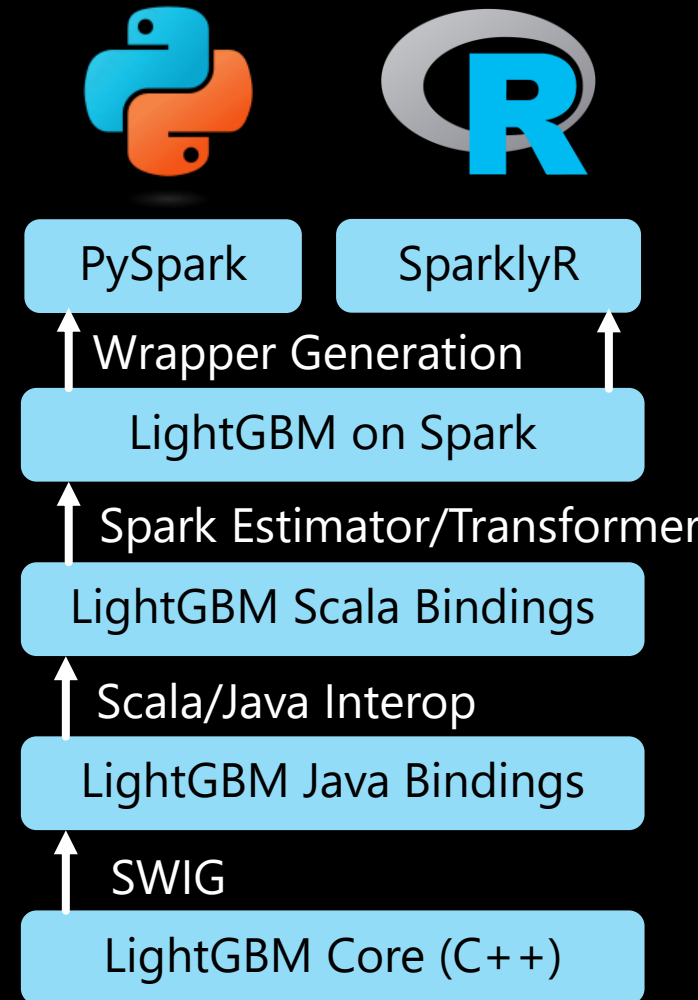Distributed Model Interpretability with LIME

Deep Learning with CNTK

# Example Backend: LightGBM on Spark

- Barrier Execution for Synchronizing Workers

- Fast Socket/MPI communication

- mapPartitions for Transformer

**Spark Worker**

LGBM Process

**Spark Worker**

LGBM Process

MPI Ring

LGBM Process

**Spark Worker**

| Framework | Time(s) | Area under ROC |
|-----------|---------|----------------|
| XGBoost | 52.60 | .808 |
| SparkML GBT | 82.78 | .788 |
| **LightGBM** | **45.39** | **.812** |

PySpark | SparklyR

Wrapper Generation

LightGBM on Spark

Spark Estimator/Transformer

LightGBM Scala Bindings

Scala/Java Interop

LightGBM Java Bindings

SWIG

LightGBM Core (C++)

Ilya Matiach, ilmat@microsoft.com
Developer, Azure ML

# Cognitive Services

- ▶ High quality pre-built intelligent services

- ▶ No time intensive model training or deployment

- ▶ Leverage Microsoft Research and Azure ML

- ▶ **Available as Docker Containers**

Person

Skateboard

Bing

"Hello World"

Anomaly Finder Result (90 Sensitivity)

Place  Time Range

I had a wonderful trip to Seattle last week and even visited the Space Needle 2 times!

Place

En-US

84% positive

# Vision

Object, scene, and activity detection

Face recognition and identification

Celebrity and landmark recognition

Emotion recognition

Text and handwriting recognition (OCR)

Customizable image recognition

Video metadata, audio, and keyframe extraction and analysis

Explicit or offensive content moderation

# Speech

Speech transcription (speech-to-text)

Custom speech models for unique vocabularies or complex environment

Text-to-speech

Custom Voice

Real-time speech translation

Customizable speech transcription and translation

Speaker identification and verification

# Language

Language detection

Named entity recognition

Key phrase extraction

Text sentiment analysis

Multilingual and contextual spell checking

Explicit or offensive text content moderation

PII detection for text moderation

Text translation

Customizable text translation

Contextual language understanding

# Decision

Q&A extraction from unstructured text

Knowledge base creation from collections of Q&As

Semantic matching for knowledge bases

Customizable content personalization learning

# Search

Ad-free web, news, image, and video search results

Trends for video, news

Image identification, classification and knowledge extraction

Identification of similar images and products

Named entity recognition and classification

Knowledge acquisition for named entities

Search query autosuggest

Ad-free custom search engine creation

# Azure Cognitive Services on Spark

- ▶ Easy to use integration between Spark and the Azure Cognitive Services

- ▶ Composable and pipelinable with all other SparkML models!

- ▶ Exponential Backoffs, Backpressure, Batching, Async Parallelism

- ▶ Fully Fluent API

```
val df = new TextSentiment()
    .setTextCol("text")
    .setOutputCol("sentiment")
    .transform(inputs)
```

| Features | Time (s) | Errors # |
|---|---|---|
| None | 30.8 | 18993 |
| EBO+BP | 1163.0 | 0 |
| EBO+BP+B | 57.1 | 0 |
| **EBO+BP+B+P** | **49.7** | **0** |

# HTTP on Spark

▶ Full Integration between HTTP Protocol and Spark SQL

▶ Spark as a Microservice Orchestrator

▶ Spark + X

▶ Support for all Spark Languages



Web Service

Client    Client    Client

Partition    Partition    Partition

Spark Worker

```
df = SimpleHTTPTransformer()
    .setInputParser(JSONInputParser())
    .setOutputParser(JSONOutputParser()
        .setDataType(schema))
    .setOutputCol("results")
    .setUrl(…)
```

# Deploying on Kubernetes

▶ Works on any k8s cluster

▶ Helm: Package Manager for Kubernetes

```
helm repo add microsoft \
https://microsoft.github.io/charts/repo
helm update

helm install microsoft/spark --version 1.0.0
```

Dalitso Banda, dbanda@microsoft.com
Microsoft AI Development Acceleration Program

# Model Deployment with Spark Serving

▶ Sub-millisecond RESTful Model Deployment on Spark Clusters

**Batch API:**
```
spark.read.parquet.load(…)
    .select(…)
```

**Streaming API:**
```
spark.readStream.kafka.load(…)
    .select(…)
```

**Serving API:**
```
spark.readStream.server("0.0.0.0", 5000).load(…)
    .select(…)
```

AI for Earth

Snow Leopard Trust

# Endangered Status Matters



BBC NEWS

Snow leopard no longer 'endangered'

14 September 2017

Statement on IUCN Red List Status Change of the Snow Leopard

*The Snow Leopard Trust, one the leading conservation organizations working to protect this cat, opposes the IUCN's decision to change the snow leopard's Red List status from 'Endangered' to 'Vulnerable'.*

# Remote Camera Trapping

# Creating a labelled Training Dataset

# Creating a labelled Training Dataset

# Transfer Learning with ResNet 50



Filters from Zeiler + Fergus 2013

# Performance

## Without Deep Featurization

## With Deep Featurization, Augmentation, and Temporal Ensembling



**Accuracy 65.6%**

**Accuracy 94.7%**

# Goal: Identify Individual Leopards



Source: HotSpotter - Patterned Species Instance Recognition

# Automating Detection with LIME on Spark

# LIME on Spark

# End to End Architecture

# Results

**Human Labels**

**Unsupervised FRCNN Outputs**

**Human Labels**

**Unsupervised FRCNN Outputs**

# Microsoft Machine Learning for Apache Spark

v0.18

**Microsoft's Open Source Contributions to Apache Spark**

Distributed Machine Learning

Fast Model Deployment

Microservice Orchestration

Multilingual Binding Generation

www.aka.ms/spark

Azure/mmlspark

# Thanks to

- You all!
- **Ilya Matiach: LightGBM on Spark**
- **Markus Cozowicz: VW on Spark**
- Sudarshan Raghunathan, Christina Lee, Daniel Ciborowski, Eli Barzilay, Tong Wen, Pablo Castro, Chris Hoder, Ryan Gaspar, Henrik Neilsen, Andrew Schonhoffer, Joseph Sirosh
- Microsoft NERD Garage Team + MIT Externship Program
- Snow Leopard Trust: Koustubh Sharma, Rhetick Sengupta, Jeff Brown, Michael Despines
- Microsoft Development Acceleration Team:
  - Dalitso Banda, Casey Hong, Karthik Rajendran, Manon Knoertzer, Tayo Amuneke, Alejandro Buendia
- Azure CAT, AzureML, and Azure Search Teams

# Get in Touch

- Support: mmlspark-support@microsoft.com
- Me: marhamil@microsoft.com
- Github ⬤ : Azure/mmlspark
- Website: www.aka.ms/spark
- Paper:  www.aka.ms/spark-paper
- Contributions Welcome!
- Check out our MSR Podcast on Oct 2

# Backup Slides

AI for Cultural Institutions

Microsoft    THE MET    MIT

# Celebrating 2 years of Open Access at The MET

▶ In 2016 The MET Released 400k images under open access

▶ This past winter the MET released a new subject-keyword dataset of image annotations

▶ MIT, The MET, and Microsoft participated in a 3-day hackathon to create intelligent experiences using the new collection

**OA** OpenAccess

Goals:

Create new works of art

Use new work to explore existing art

Explore further with intelligent search

Needed Technologies:

Generative Adversarial Networks

Reverse image search

Elasticsearch with Cognitive Services

# Generative Adversarial Art

# Custom Reverse Image Search

Query Image   ResNet Featurizer   Deep Features   Fast Nearest Neighbor Lookup   Closest Match

MMLSpark

SparkML LSH or Annoy

Filters from Zeiler + Fergus 2013

# Example Nearest Neighbors

# Intelligent Search Index



▶ Pipe images through Computer Vision API to annotate image for searching

▶ Stream images and intelligent annotations to Azure Search

**Query Image:**

**Describe Image Output:**

A picture containing a person

A picture containing a glass, cup

A fish swimming underwater

**Deep Feature Nearest Neighbors:**

# End Application: Gen Studio

AI for Accessibility

Seeing AI

# Currency Identification

# A Familiar Architecture...



| Queries |
|---------|
| 1 Dollar |
| 5 Dollars |
| 10 Dollars |
| 20 Dollars |

Labelled Images

$1

$5

$20

Deep Features

Bing Image Search → Union + Distinct → MobileNet → Logistic Regression → Azure Machine Learning + Spark Serving

Prep Data | Train | Deploy