**BigDL**

Distributed, High-Performance
Deep Learning Framework
for Apache Spark

https://github.com/intel-analytics/bigdl

**ANALYTICS ZOO**

Unified Analytics + AI Platform
Distributed TensorFlow, Keras, PyTorch and BigDL
on Apache Spark

https://github.com/intel-analytics/analytics-zoo

## Accelerating Data Analytics + AI Solutions At Scale

intel AI

# What's on-going in Spark + AI Community

**– views from a contributor & practitioner**

**Shengsheng Huang**

**Intel AnalyticsZoo team**

# Agenda

- **Efforts for building unified data analytics + AI in production**
- **Efforts to support emerging AI applications**

# Agenda

- **Efforts for building unified data analytics + AI in production**
- **Efforts to support emerging AI applications**

# What's new in spark + ai community

**Spark 3. 0**
- Optimizations on SQL execution (adaptive query execution, dynamic partition pruning )
- DataSourceV2
- Project Hydrogen (Barrier execution mode, Accelerator-aware scheduling, optimized data exchange)
- Spark Graph
- Spark on Kubernetes
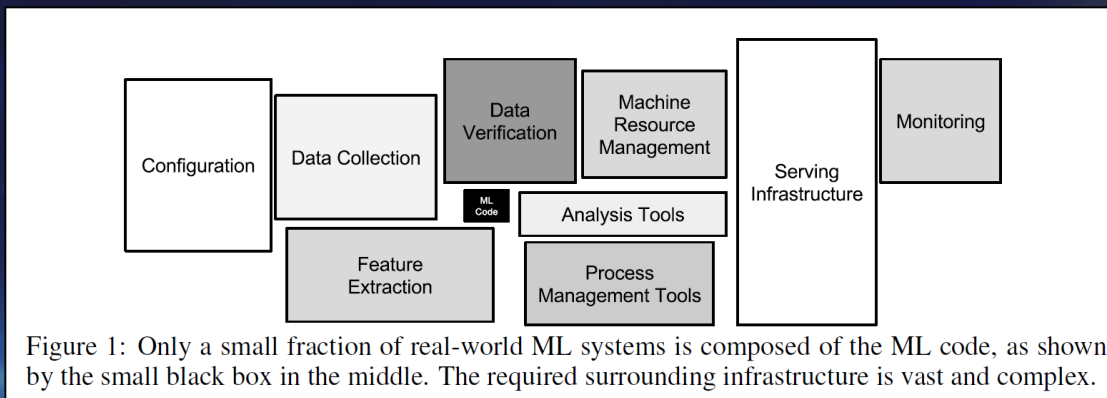- ...

**All for Productivity**

**MLFlow – ML lifecycle management**
- Tracking – log code, data, config, results of experiments, and compare & query
- Projects – code packaging format for reproducible runs on any platform
- Models – model packaging format for sending models to diverse deployment tools.

**Koalas – pandas API on Spark**

**Delta Lake – ACID layer upon data lakes**

# Rationale behind the efforts in community



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

"Hidden Technical Debt in Machine Learning Systems", Sculley et al., Google, NIPS 2015 Paper

- **Integration/Injection of heterogenous data models/sources, computation models, software/hardware components, ... (e.g. DataSourceV2, Project Hydrogen, Spark Graph)**
- **E2E Workflow, ML Lifecycle, Serving, Deployment, Orchestration, ... (e.g. MLFlow, KubeFlow, Seldon, TFX)**
- **Efficiency & Reliability (e.g. SQL-related optimizations, Delta Lake)**
- **Friendly APIs (e.g. Koalas)**

# AI ON BIG DATA

## BigDL

High-Performance
Deep Learning Framework
for Apache Spark

software.intel.com/bigdl

## ANALYTICS ZOO

Unified Analytics + AI Platform
Distributed TensorFlow*, PyTorch*,
Keras* and BigDL on Apache Spark

https://github.com/intel-analytics/analytics-zoo

## ACCELERATING DATA ANALYTICS + AI SOLUTIONS DEPLOYMENT AT SCALE

(intel) AI

# Analytics Zoo
## Unified End-to-End Data Analytics + AI Platform

| | | | | |
|---|---|---|---|---|
| **Use case** | Recommendation | Anomaly Detection | Text Classification | Text Matching |
| **Model** | Image Classification | Object Detection | Seq2Seq | Transformer — BERT |
| **Feature Engineering** | image — 3D image | text | time series | |

**Integrated Analytics/AI Pipelines**

| tfpark: Distributed TF on Spark | Distributed Keras/PyTorch on Spark |
|---|---|
| nnframes: Spark Dataframes & ML Pipelines for Deep Learning | Distributed Model Serving (batch, streaming & online) |

**Backend/Library**

TensorFlow  Keras  PyTorch  BigDL  NLP Architect  Apache Spark  Apache Flink

Ray  MKLDNN  OpenVINO  Intel® Optane™ DCPMM  DL Boost (VNNI)

https://github.com/intel-analytics/analytics-zoo

intel AI

# Distributed TensorFlow on Spark

- **Data wrangling and analysis using PySpark**

- **Deep learning model development using TensorFlow or Keras**

- **Distributed training / inference on Spark**

```python
#pyspark code
train_rdd = spark.hadoopFile(…).map(…)
dataset = TFDataset.from_rdd(train_rdd,…)

#tensorflow code
import tensorflow as tf
slim = tf.contrib.slim
images, labels = dataset.tensors
with slim.arg_scope(lenet.lenet_arg_scope()):
    logits, end_points = lenet.lenet(images, …)
loss = tf.reduce_mean( \
    tf.losses.sparse_softmax_cross_entropy( \
    logits=logits, labels=labels))

#distributed training on Spark
optimizer = TFOptimizer.from_loss(loss, Adam(…))
optimizer.optimize(end_trigger=MaxEpoch(5))
```

**Write TensorFlow code inline in PySpark program**

(intel) AI

# Object Detection and Image Feature Extraction at JD.com*





Image feature extraction pipeline throughput (img/s)

- Reuse existing Hadoop/Spark clusters for deep learning with no changes (image search, IP protection, etc.)
- Efficiently scale out on Spark with superior performance (*3.83x* speed-up vs. GPU severs) as benchmarked by JD

http://mp.weixin.qq.com/s/xUCkzbHK4K06-v5qUsaNQQ
https://software.intel.com/en-us/articles/building-large-scale-image-feature-extraction-with-bigdl-at-jdcom

# Product Recommendations in **Office Depot\***



https://software.intel.com/en-us/articles/real-time-product-recommendations-for-office-depot-using-apache-spark-and-analytics-zoo-on

# Computer Vision Based Product Defect Detection in Midea*

# Particle Classifier for High Energy Physics in CERN*



Deep learning pipeline
for physics data

Model serving using
Apache Kafka and Spark

# Wrap Up

**Community is making efforts to make Spark a unified Analytics + AI platform**

**Analytics Zoo is also working towards similar goal, by**

- Seamless integration various components, e.g. Tensorflow, PyTorch, BigDL, etc.
- Providing full-stack optimizations involving hardware/software (VNNI, MKL-DNN, OpenVINO, etc.)
- Providing ease of use, end-to-end, from laptop to production platform

**We are both contributors and practitioners**. We use, learn, and contribute.

(intel) AI

# Agenda

- **Efforts for building unified data analytics + AI in production**
- **Efforts to support emerging AI applications**

# Towards General AI

Strategies to build AI for game playing, robots, autonomous driving, etc.



**Deep Reinforcement Learning (DRL)**

**Imitation Learning**

# Parallel Architecture for Deep RL



Massively Parallel Methods for Deep Reinforcement Learning
https://arxiv.org/abs/1507.04296

# Ray On Spark

## Ray

- **https://github.com/ray-project/ray**
- **a distributed framework for emerging AI applications open-sourced by UC Berkeley RISELab**

## RayOnSpark

- **a feature recently added to Analytic Zoo**
- **allows users to directly run Ray programs on Apache Hadoop\*/YARN**
- **Ray applications can be seamlessly integrated into Spark pipeline and operate directly on Spark RDDs or DataFrames.**



riselab

RayOnSpark: Running Emerging AI Applications on Big Data Clusters with Ray and Analytics Zoo

Jason Dai
Jul 29 · 4 min read

*Zhichao Li (zhichao.li@intel.com), Jason Dai (jason.dai@intel.com)*

**https://medium.com/riselab/rayonspark-running-emerging-ai-applications-on-big-data-clusters-with-ray-and-analytics-zoo-923e0136ed6a**

(intel) AI

# Building AI to Play FIFA

**FIFA18\*** – A real-time 3D soccer simulation video game by Electronic Arts\*

**Our Experiment Platform** (collaborations w/ SJTU)

- runs alongside FIFA game in a non-intrusive way
- provides abstraction of game environment (observations, actions, rewards, scores, semantics, etc.)
- Implemented agents: RL, IL, Hybrid (IL + RL)

**Future Work:**

- Transfer between Google Research Football and FIFA?
- Train agents in massive scale w/ Ray & RayOnSpark
- Additional models/scenarios, etc.

https://www.slideshare.net/jason-dai/building-ai-to-play-the-fifa-video-game-using-distributed-tensorflow-on-analytics-zoo

*Other names and brands may be claimed as the property of others.



Results on Shooting Bronze Scenario

| | | Score | Goal Ratio |
|---|---|---|---|
| Human | master | 10112.78 | 92% |
| | demonstrator | 7284.98 | 84.96% |
| Agent | IL | 10345.18 | 92.54% |
| | RL (Policy Gradient) | 5606.31 | 40.25% |
| | Hybrid (RL+IL) | 10514.43 | 95.59% |

(intel) AI

# Scalable AutoML for Time Series Analysis

## AutoML Framework



## Time Series Forecasting w/ AutoML

- **Data processing and feature engineering**
- **Neural network based (hybrid) models**
- **Automated feature selection, model selection, hyper parameter tuning**



https://medium.com/riselab/scalable-automl-for-time-series-prediction-using-ray-and-analytics-zoo-b79a6fd08139

intel AI

# Wrap Up

**We're extending the Spark stack to support emerging AI applications**

- **RayOnSpark**

**We're building emerging AI applications**

- **Building AI to play FIFA**
- **Scalable AutoML for Time Series Analysis**

# More Information on Analytics Zoo

- **Project website**
  - *https://github.com/intel-analytics/analytics-zoo*
- **Tutorials**
  - CVPR 2018: *https://jason-dai.github.io/cvpr2018/*
  - AAAI 2019: *https://jason-dai.github.io/aaai2019/*
- **"BigDL: A Distributed Deep Learning Framework for Big Data"**
  - *In proceedings of ACM Symposium on Cloud Computing 2019 (SOCC'19)*
- **Use cases**
  - *Azure, CERN, MasterCard, Office Depot, Tencent, Midea, etc.*
  - *https://analytics-zoo.github.io/master/#powered-by/*

(intel) AI

# LEGAL NOTICES AND DISCLAIMERS

intel AI

# NLP Based Customer Service Chatbot for **Microsoft Azure**

**BigDL**

Distributed, High-Performance
Deep Learning Framework
for Apache Spark

https://github.com/intel-analytics/bigdl

**ANALYTICS ZOO**

Unified Analytics + AI Platform
Distributed TensorFlow, Keras, PyTorch and BigDL
on Apache Spark

https://github.com/intel-analytics/analytics-zoo

## Accelerating Data Analytics + AI Solutions At Scale

intel AI