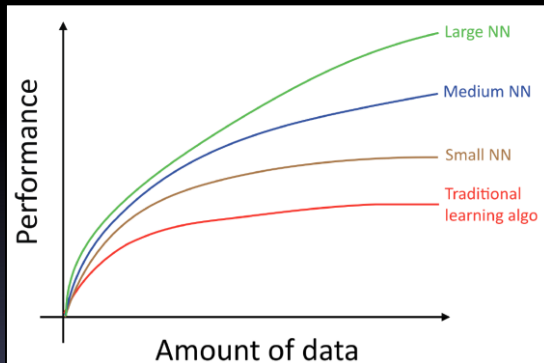# Intel data analytics solution at scale

Zhichao Li (zhichao.li@intel.com)

# Motivations

# Data Scale Driving Deep Learning Process



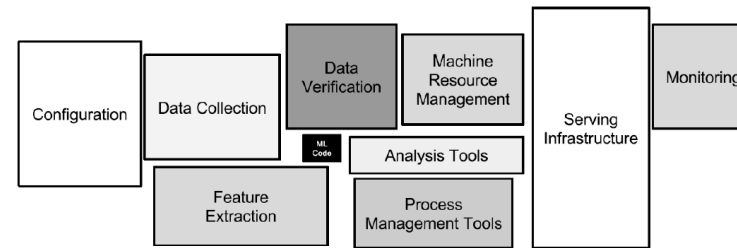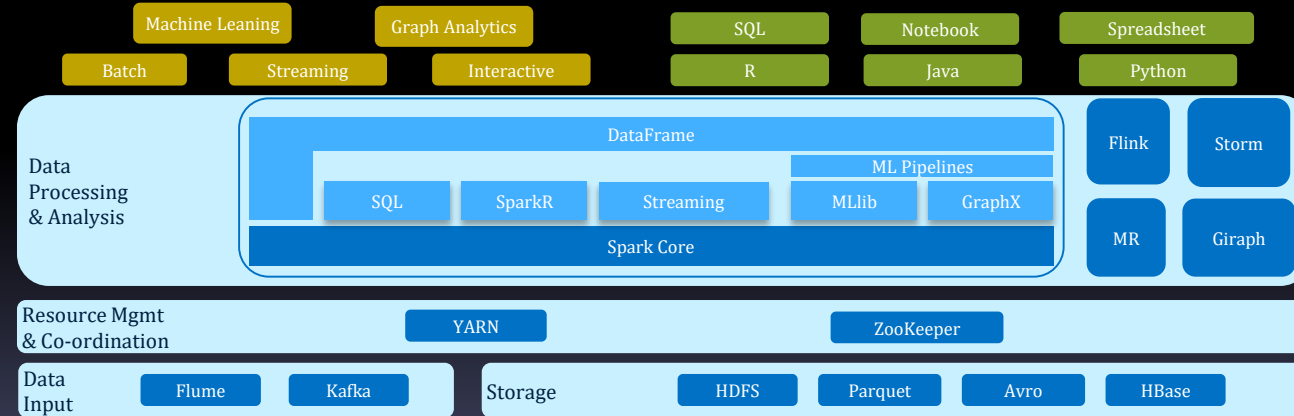"Machine Learning Yearning", Andrew Ng, 2016



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", Google, NIPS 2015 paper

# Hadoop Becoming the Center of Data Gravity

## Hadoop & Spark Platform

| | |
|---|---|
| Machine Leaning | Graph Analytics |
| Batch | Streaming | Interactive |

| | | |
|---|---|---|
| SQL | Notebook | Spreadsheet |
| R | Java | Python |

**Data Processing & Analysis**

- DataFrame
  - ML Pipelines
  - SQL | SparkR | Streaming | MLlib | GraphX
- Spark Core

Flink | Storm

MR | Giraph

**Resource Mgmt & Co-ordination**

YARN | ZooKeeper

**Data Input**

Flume | Kafka
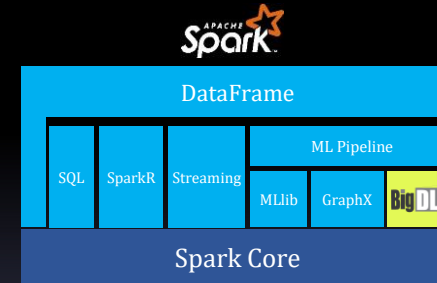
**Storage**

HDFS | Parquet | Avro | HBase

# Overview

# BigDL

## Bringing Deep Learning To Big Data Platform

- **Distributed** deep learning framework for Apache Spark*

- Make deep learning more accessible to **big data users**
  **data scientists**
  - Write deep learning applications as *standard Spark programs*
  - Run on existing Spark/Hadoop clusters *no changes needed*

- Feature parity with popular deep learning frameworks
  - E.g., Caffe, Torch, Tensorflow, etc.

- High performance
  - Powered by Intel MKL and multi-threaded programming

- Efficient scale-out
  - Leveraging Spark for distributed training & inference



https://github.com/intel-analytics/BigDL
https://bigdl-project.github.io/

# Analytics

## Build and Productionize Deep Learning Apps for Big Data at Scale
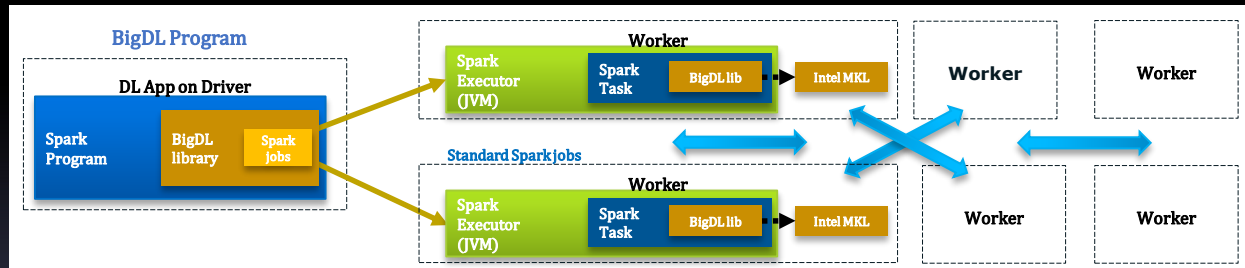
| | |
|---|---|
| **Reference Use Cases** | • Anomaly detection<br>• Sentiment analysis<br>• Fraud detection<br>• Chatbot, sequence prediction, etc. |
| **Built-In Deep Learning Models** | • Image classification<br>• Object detection<br>• Text classification<br>• Recommendations<br>• Sequence-to-sequence, GAN, etc. |
| **Feature Engineering** | Feature transformations for<br>• Image, text, 3D imaging, time series, speech, etc. |
| **High-Level Pipeline APIs** | • Native deep learning support in Spark DataFrames and ML Pipelines<br>• Autograd, Keras and transfer learning APIs for model definition<br>• Support for model serving/inference pipelines |
| **Backbends** | Spark, BigDL, TensorFlow, etc. |

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

# BigDL Run as Standard Spark Programs

**Standard Spark jobs**
- No changes to the Spark or Hadoop clusters needed

**Iterative**
- Each iteration of the training runs as a Spark job

**Data parallel**
- Each Spark task runs the same model on a subset of the data (batch)

# Models Interoperability Support
### (e.g., between TensorFlow, Keras, Caffe, Torch, BigDL models)

## Load existing TensorFlow, Keras, Caffe, Torch Model

Useful for inference and model fine-tuning

Allows for transition from single-node for distributed application deployment

- Allows for model sharing between data scientists and production engineers

# Use Cases

# Cloud & Big Data Platforms

Running BigDL, Deep Learning for Apache Spark, on AWS*
(Amazon* Web Service)
https://aws.amazon.com/blogs/ai/running-bigdl-deep-learning-for-apache-spark-on-aws/

BigDL on Alibaba* Cloud
E-MapReduce*

https://yq.aliyun.com/articles/73347

BigDL on CDH* and Cloudera* Data Science Workbench*
http://blog.cloudera.com/blog/2017/04/bigdl-on-cdh-and-cloudera-data-science-workbench/

BigDL Spark deep learning library VM now available on Microsoft* Azure* Marketplace https://azure.microsoft.com/en-us/blog/bigdl-spark-deep-learning-library-vm-now-available-on-microsoft-azure-marketplace/

BigDL in KMR* Service of Kingsoft* Cloud

https://docs.ksyun.com/read/latest/33/_book/bigDL.html

Using BigDL in IBM* Data Science Experience

https://medium.com/ibm-data-science-experience/using-bigdl-in-data-science-experience-for-deep-learning-on-spark-f1cf30ad6ca0

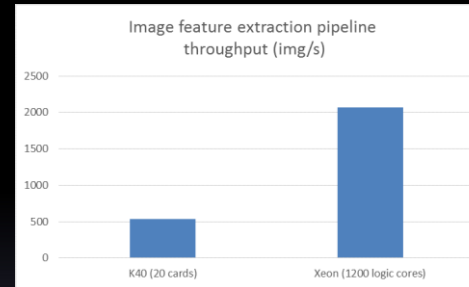Using BigDL for deep learning with Apache Spark and Google* Cloud Dataproc*
https://cloud.google.com/blog/big-data/2018/04/using-bigdl-for-deep-learning-with-apache-spark-and-google-cloud-dataproc

Intel's BigDL on Databricks*

https://databricks.com/blog/2017/02/09/intels-bigdl-databricks.html

BigDL Shipped in Cray* Urika-XC* Analytics Software Suite

https://www.cray.com/blog/scalable-deep-learning-bigdl-urika-xc-software-suite/

# Object Detection and Image Feature Extraction in JD

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

For more complete information about performance and benchmark results, visit www.intel.com/benchmarks.

# User-Merchant Propensity Modeling in MasterCard



Implementation : run BigDL & ALS over Spark on Hadoop

https://www.ai-expo.net/northamerica/talk/using-deep-learning-intel-bigdl-optimized-personalized-card-linked-offer/
https://conferences.oreilly.com/strata/strata-ca/public/schedule/detail/63897

# Neural Recommendation Engine in China Life

Realize re-discovery of life insurance business, accurately and effectively recommend products.

# Image Similarity Search for **MLSListings**



MLSlistings built image-similarity based house recommendations using BigDL on Microsoft Azure

https://software.intel.com/en-us/articles/using-bigdl-to-build-image-similarity-based-house-recommendations

# NLP Based Call Center Routing in GigaSpaces

# 3D Medical Image Analysis in UCSF

# Partner With Us

- Use Analytics-Zoo & Share your Experience

- Use Intel Optimized Libraries & Frameworks

- Leverage Intel Developer Zone Resources

Source code: https://github.com/intel-analytics/analytics-zoo/
Documents: https://analytics-zoo.github.io/

# Legal Disclaimers