

# Analytics-Zoo: 统一的大数据分析+AI平台

Analytics Zoo: A Unified Data Analytics + AI Platform



Intel 软件架构师 大数据分析和人工智能创新院





# Why Analytics-Zoo



# BUDU

#### Distributed, High-Performance Deep Learning Framework for Apache Spark\*

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

Unifying Analytics + Al on Apache Spark\*

## ANALYTICS ZOOO Analytics + AI Platform Distributed TensorFlow\*, Keras\* and BigDL on Apache Spark\*

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





### Real-World ML/DL Applications Are Complex Data Analytics Pipelines



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











#### **Chasm b/w Deep Learning and Big Data Communities**



Deep learning experts

Real-world users (big data users, data scientists, analysts, etc.)







#### Large-Scale Image Recognition



<u>https://software.intel.com/en-us/articles/building-large-scale-image-feature-</u> <u>extraction-with-bigdl-at-jdcom</u>







#### **Standard Spark jobs**

 No changes to the Spark or Hadoop clusters needed **Data parallel** 

• Each Spark task runs the same model on a subset of the data (batch) "Zero" code change

#### • Directly support TensorFlow, Keras and Caffe Model **Seamlessly deployed on production big data clusters**

• Only need to install on driver node.





# What's Analytics-Zoo



#### Distributed, High-Performance Deep Learning Framework for Apache Spark\*

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

#### Unifying Analytics + AI on Apache Spark\*



#### Analytics + AI Platform

#### Distributed TensorFlow\*, Keras\* and **BigDL on Apache Spark\***

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





## Analytics-Zoo: Unified Analytics + AI Platform for BigData

Use case	Recommendation	Anomaly Detection		Text Classification		Text Matching	
Model	Image Classification Object De		ection	Seq2Seq Tran		rmer BERT	
Feature Enginee	ering image	3D image	text	Times	eries		
High Level Pipelines	tfpark: Distribute	Distributed Keras w/ autograd on Spark					
	nnframes: Spark D Pipelines for D	Distributed Model Serving ( batch, streaming & online )					
Backend	TensorFlow* Keras	* BigDL Open	VINO	MKLDNN Apach	ne Spark*	Apache Flink*	
	https://github.com/intel-analytics/analytics-zoo						





#### **Analytics-Zoo: Run as Standard Spark Programs**













## **Distributed Training in Analytics-Zoo**



Peer-2-Peer All-Reduce synchronization





#### **Distributed TF & Keras on Spark**

 Data wrangling and analysis using PySpark

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

 Distributed training / **inference on Spark** 

```
#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
```

Write TensorFlow code inline in PySpark program

optimizer = TFOptimizer.from loss(**loss**, Adam(...)) \ optimizer.optimize(end\_trigger=MaxEpoch(5))





#### **Spark Dataframe & ML Pipeline for DL**

*#Spark dataframe transformations* parquetfile = spark.read.parquet(...) train df = parquetfile.withColumn(...)

#### *#Keras API*

model = Sequential() .add(MaxPooling2D(pool size=(2, 2))) \

*#Spark ML pipeline* Estimater = NNEstimater(**model**, CrossEntropyCriterion()) \ .setFeaturesCol("image") nnModel = estimater.fit(train\_df)

```
.add(Convolution2D(32, 3, 3, activation='relu', input shape=...)) \
.add(Flatten()).add(Dense(10, activation='softmax')))
      .setLearningRate(0.003).setBatchSize(40).setMaxEpoch(5) \
```





## **Distributed Model Serving**



#### Distributed model serving in Web Service, Flink, Kafka, Storm, etc.

Plain Java or Python API, with OpenVINO and DL Boost (VNNI) support

Bolt





# Analytics-Zoo use cases



#### **Computer vision Based Product Defect Detection in Midea**



https://software.intel.com/en-us/articles/industrial-inspection-platform-in-midea-and-kuka-using-distributed-tensorflow-<u>on-analytics</u>









### **Recommender AI Service in MasterCard**



<u>https://software.intel.com/en-us/articles/deep-learning-with-analytic-zoo-optimizes-mastercard-recommender-</u> <u>ai-service</u>





### **Deep Learning Made Easy for BigData**







- software or service activation. Learn more at intel.com, or from the OEM or retailer.
- No computer system can be absolutely secure.
- benchmark results, visit http://www.intel.com/performance.

Intel, the Intel logo, Xeon, Xeon phi, Lake Crest, etc. are trademarks of Intel Corporation in the U.S. and/or other countries. \*Other names and brands may be claimed as the property of others. © 2019 Intel Corporation

### **LEGAL DISCLAIMERS**

•Intel technologies' features and benefits depend on system configuration and may require enabled hardware,

• Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. Consult other sources of information to evaluate performance as you consider your purchase. For more complete information about performance and











#### 扫码加入社群 与志同道合的码友一起 Code Up



该群属于"阿里云ACE"部门群,仅组织内部成员可 以加入,如果组织外部人员收到此分享,需要先申 请加入该组织。





