

What is Analytics Zoo



Distributed, High-Performance
Deep Learning Framework
for Apache Spark



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



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





Software

在Flink上使用Analytics Zoo进行实时、分布式深度学习模型推理



目录 Contents

1. 大规模人工智能应用面临的挑战

AI production at scale is facing lots of challenges.

2. 统一的大数据分析及人工智能

Integrated Data Analytics and AI.

3. 跨行业的端到端客户案例实践

Cross-industry End to End Use Cases.

01 大规模人工智能应用面临的挑战

01 AI production at scale is facing lots of challenges

以数据为中心的世界

The Data-Centric World

全球超过 **OVER**
一半 HALF 数据 **OF THE WORLD'S DATA**

创建于过去
WAS CREATED IN THE LAST
两年 2 YEARS

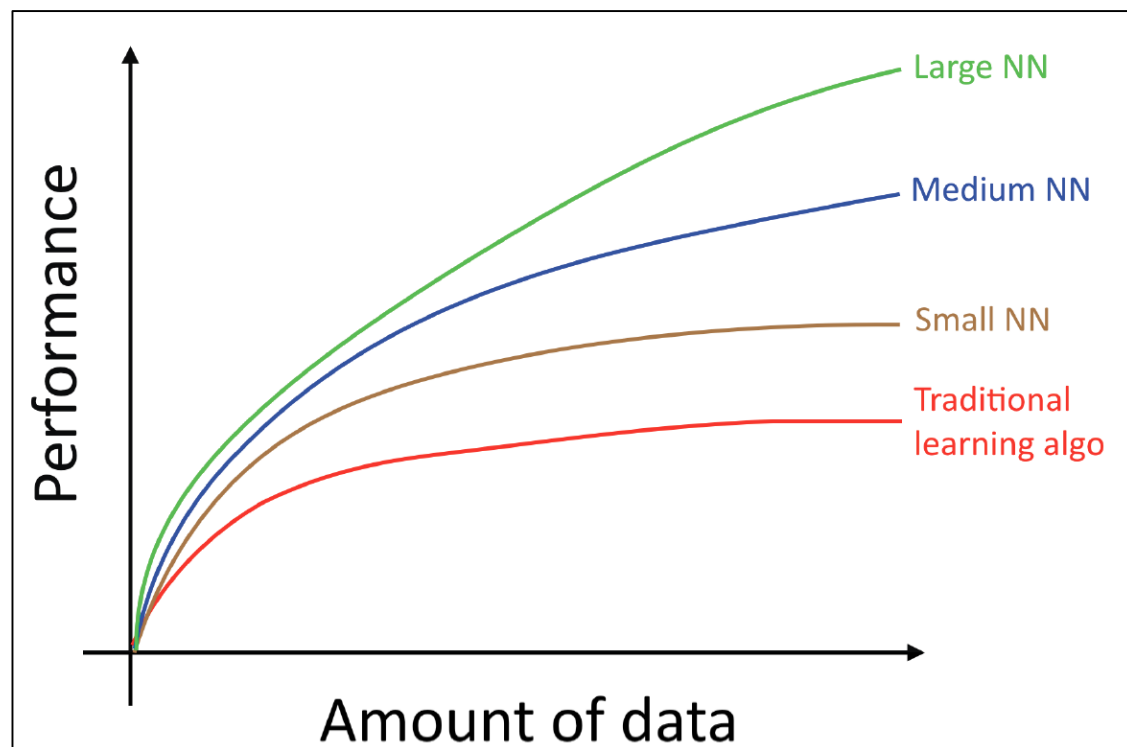
其中只有不到 **LESS THAN**
2% 的数据 **HAS BEEN ANALYZED**
经过了分析

大规模人工智能应用

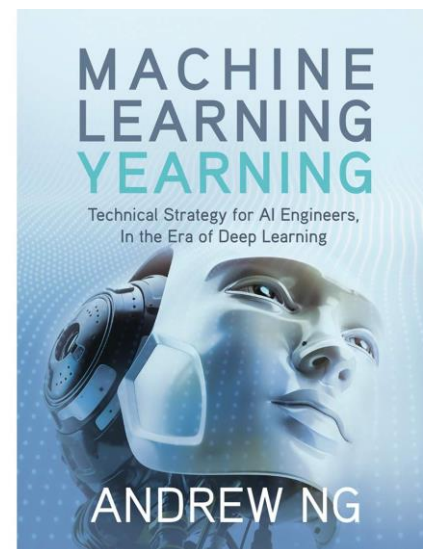
AI Production at Scale

数据驱动深度学习和人工智能应用

Data drives deep learning and AI production



**“Machine Learning Yearning”,
Andrew Ng, 2016**

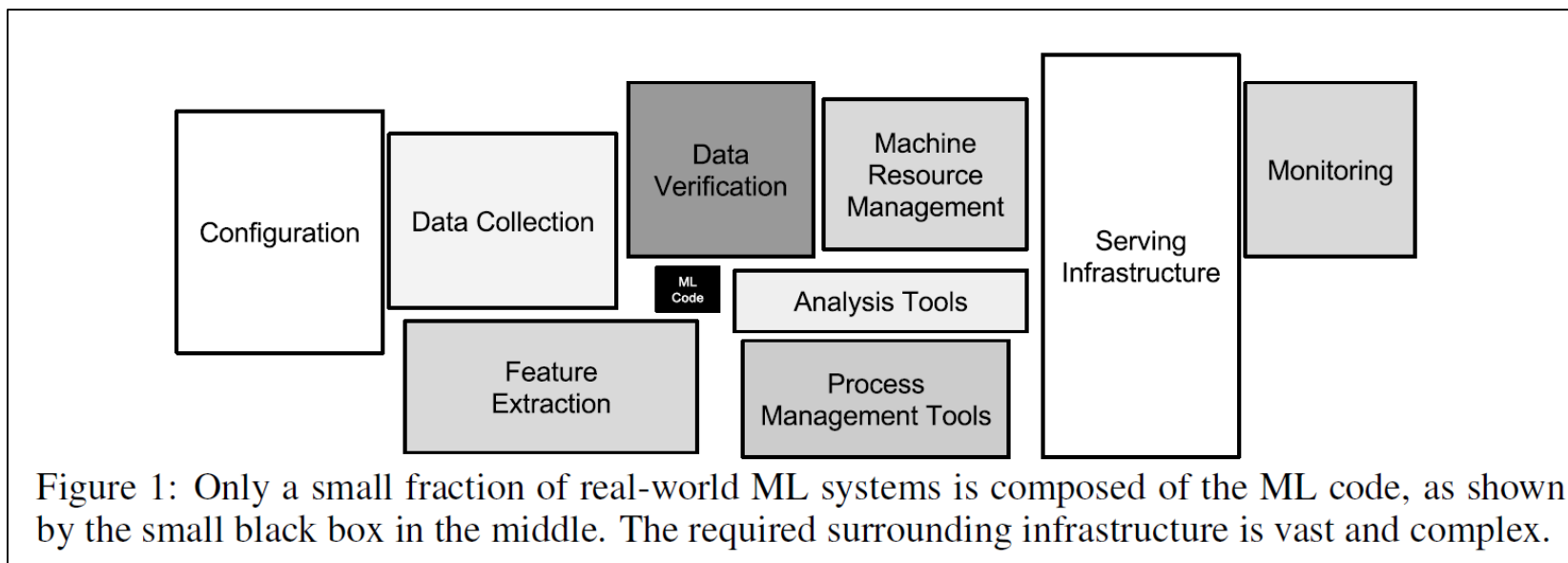


大规模人工智能应用

AI Production at Scale

正面临巨大的挑战

Facing Lots of Challenges



“Hidden Technical Debt in Machine Learning Systems”, Sculley et al., Google, NIPS 2015

02 统一的大数据分析及人工智能

02 Integrated Data Analytics and AI

统一的大数据分析及人工智能

Integrated Data Analytics and AI



大数据上的人工智能

AI on Big Data

BigDL

高性能深度学习框架

High-Performance Deep Learning
Framework for Apache Spark*

software.intel.com/bigdl

ANALYTICS ZOO

统一的分析 + 人工智能平台
Integrated Analytics + AI Toolkit

分布式

TensorFlow、PyTorch、Keras 和 BigDL

高级流水线、参考用例、人工智能模型、特征工程等

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

加快数据分析及人工智能大规模应用

Accelerating DATA Analytics + AI Solutions DEPLOYMENT At SCALE



统一的数据分析和AI流水线

End-to-End Big Data Analytics and AI Pipeline

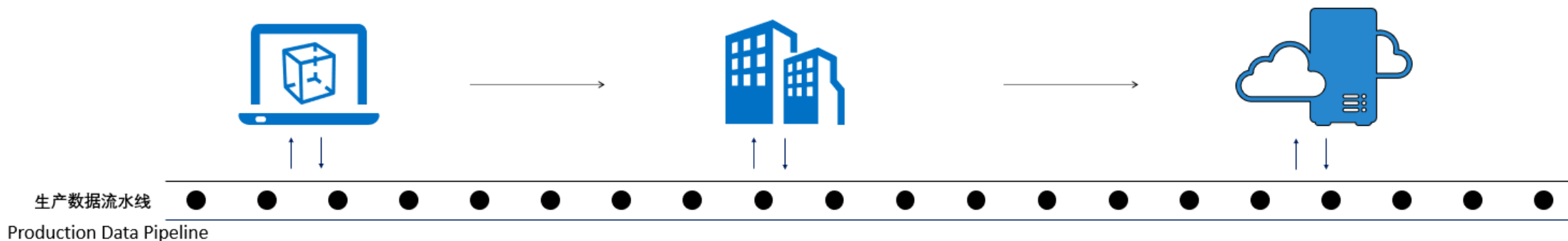
端到端、从原型到生产化部署的无缝扩展

Seamless Scaling from Laptop to Production

在笔记本电脑上使用样本数据构建原型
Prototype on laptop using sample data

在集群上使用历史数据运行模型试验
Experiment on clusters with history data

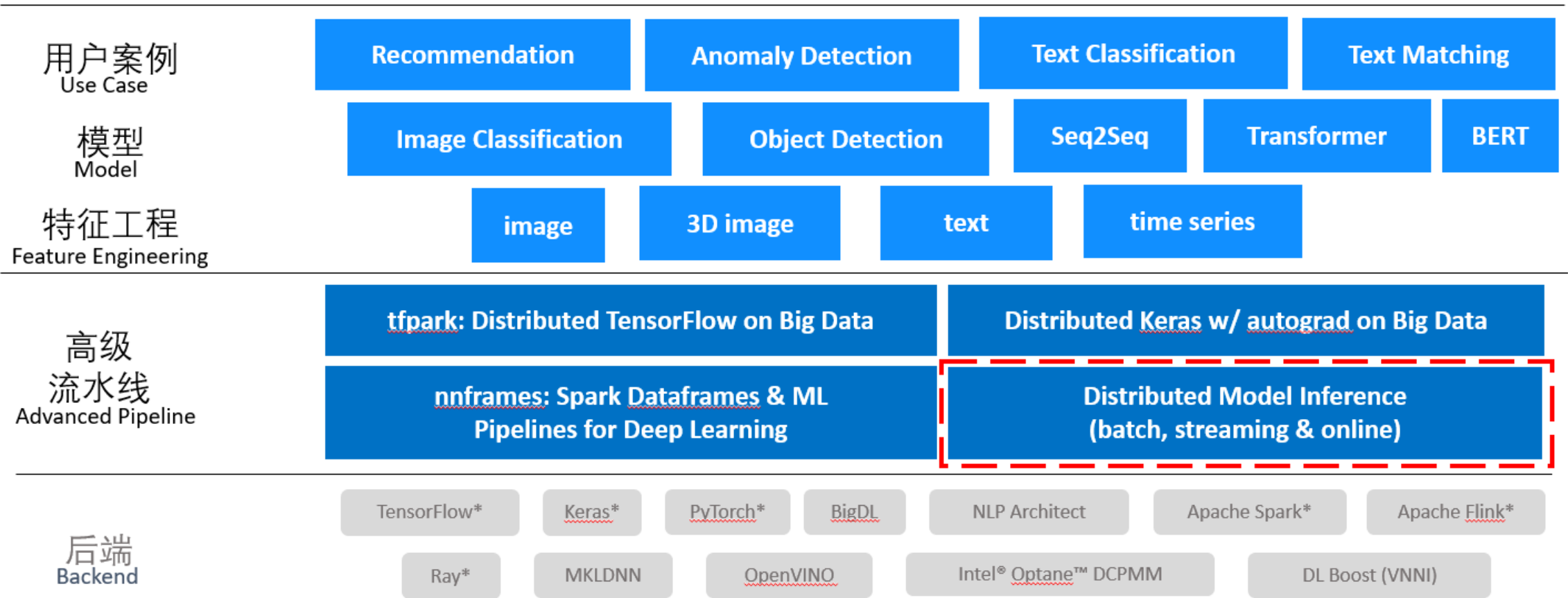
在分布式生产环境中部署
Production deployment w/ distributed data pipeline



- 从笔记本电脑到分布式集群**几乎无需任何代码更改** “Zero” code change from laptop to distributed cluster
- 无需数据拷贝，**直接访问生产大数据系统** Directly access production data without data copy
- 高效构建**端到端**的数据分析+ **AI 流水线原型** Easily prototype the end-to-end pipeline
- 无缝扩展部署到**大数据集群及生产环境** Seamlessly deployed on production big data clusters

统一的大数据分析+人工智能平台

Integrated Big Data Analytics and AI platform



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

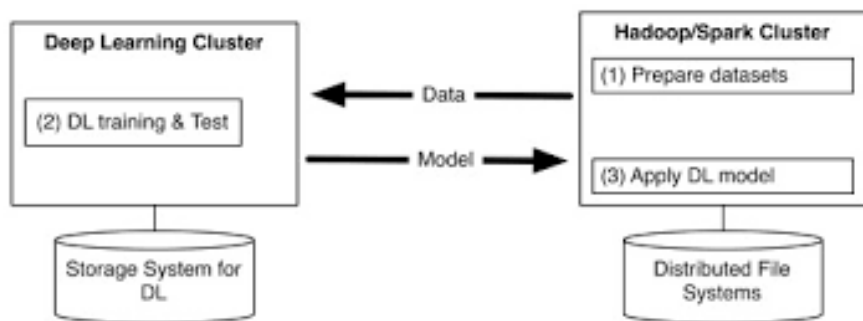


* 文中涉及的其它名称及商标属于各自所有者资产。

分布式 TensorFlow* 流水线

Distributed TensorFlow* Pipeline

- Data loading, processing and feature engineering with Big Data
- Deep learning model development using TensorFlow* or Keras*
- Distributed training / inference on Big Data



#load data

```
train_data = hadoopFile(...).map(...)  
dataset = TFDataset.from_rdd(train_rdd,...)
```

用大数据计算框架载入数据以及
预处理数据或特征工程

#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))
```

用 TensorFlow* 或 Keras* 定义深度
学习模型

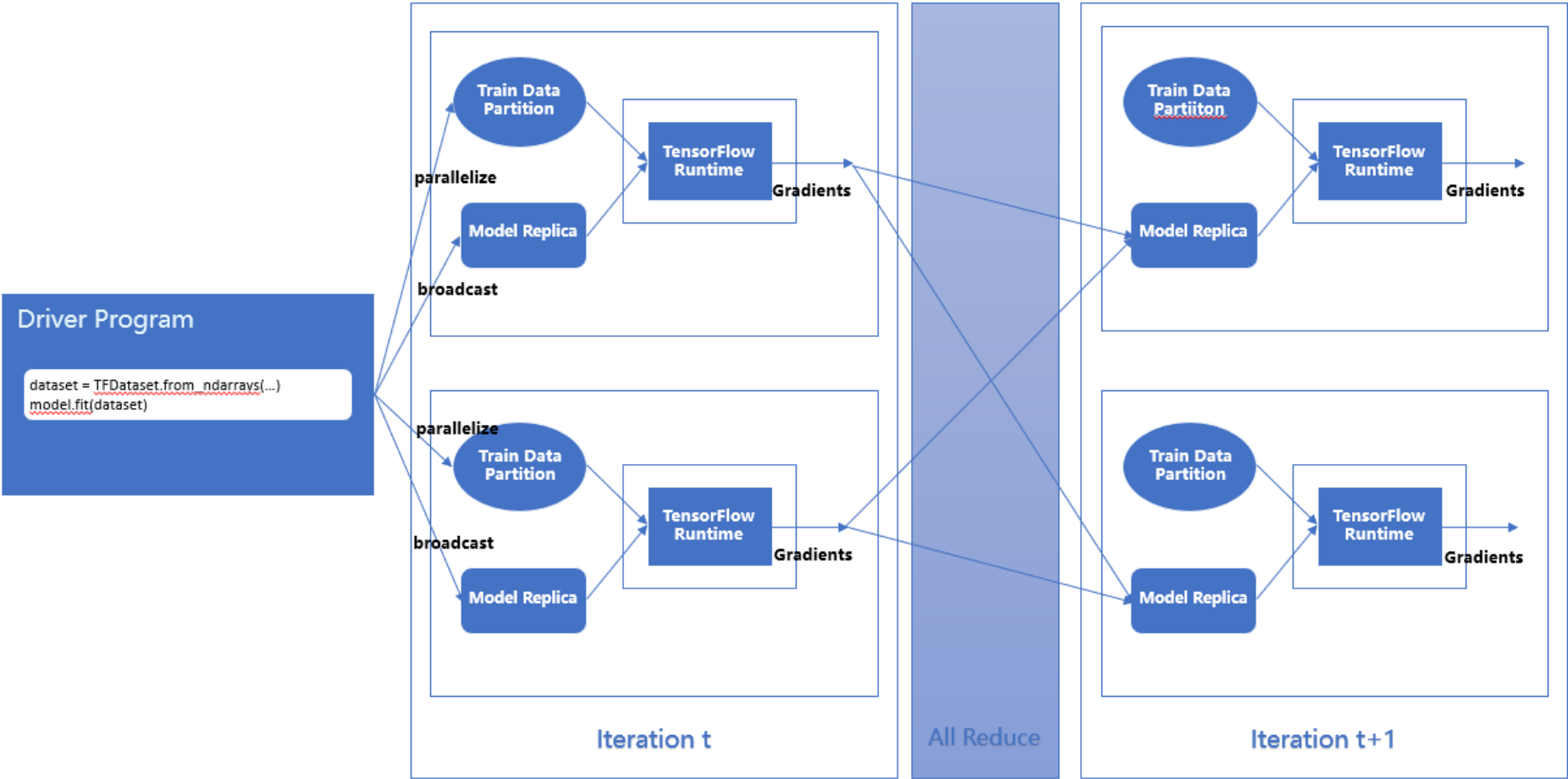
#distributed training

```
optimizer = TFOptimizer.from_loss(loss, Adam(...))  
optimizer.optimize(end_trigger=MaxEpoch(5))
```

在大数据上分布式训练或者推理

分布式 TensorFlow* 流水线

Distributed TensorFlow* Pipeline

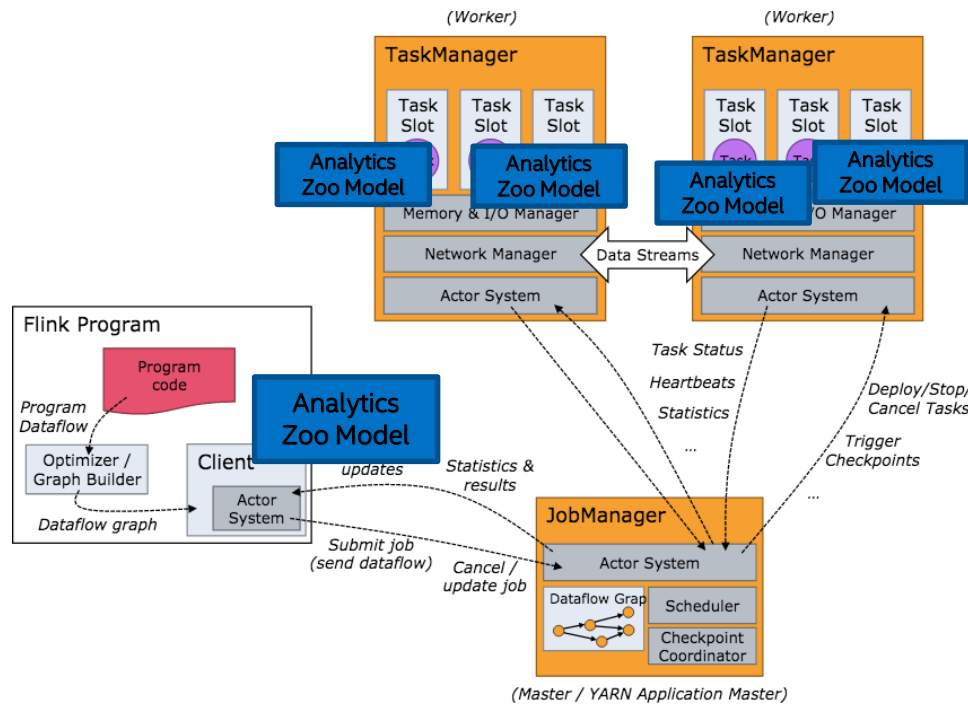


* 文中涉及的其它名称及商标属于各自所有者资产。

分布式、实时 (流式) 模型推理流水线

Distributed and Real time (streaming) Inference Pipeline

- 纯Java或Python API
- 支持Flink*, Spark* Streaming, Storm*, Kafka*等
- 支持Web Services
- 使用OpenVINO和DL Boost(VNNI) 加速



Overview Timeline Exceptions Configuration

Source: Custom Source
Parallelism: 1

Flat Map -> Map
Parallelism: 6

Analytics Zoo Model

Aggregate task statistics by taskmanager

Start Time	End Time	Duration	Name	Bytes received	Records received	Bytes sent	Records sent	Parallelism	Tasks	Status
2019-08-01, 10:09:36	2019-08-01, 10:13:20	3m 44s	Source: Custom Source	0 B	0	537 MB	330	1	0 0 1 0 0 0 0	RUNNING
2019-08-01, 10:09:36	2019-08-01, 10:13:20	3m 44s	Flat Map -> Map	536 MB	326	0 B	320	6	0 0 6 0 0 0 0	RUNNING

Start Time	End Time	Duration	Bytes received	Records received	Bytes sent	Records sent	Attempt	Host	Status
2019-08-01, 10:09:40		3m 40s	83.4 MB	53	0 B	52	1	Almaren-Node-018:41072	RUNNING
2019-08-01, 10:09:40		3m 40s	80.4 MB	54	0 B	53	1	Almaren-Node-018:41072	RUNNING
2019-08-01, 10:09:36		3m 44s	76.2 MB	54	0 B	53	1	Almaren-Node-019:39744	RUNNING
2019-08-01, 10:09:40		3m 40s	96.9 MB	55	0 B	54	1	Almaren-Node-019:39744	RUNNING
2019-08-01, 10:09:40		3m 40s	87.8 MB	55	0 B	54	1	Almaren-Node-020:35343	RUNNING

Analytics Zoo Model

Analytics Zoo Model

Analytics Zoo Model

Analytics Zoo Model

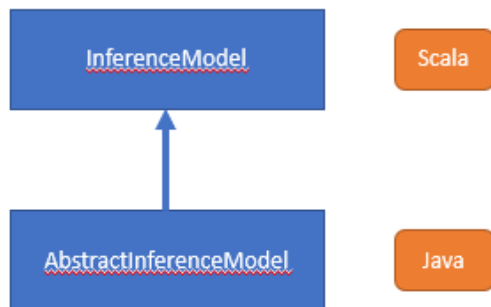
Analytics Zoo Model

Analytics Zoo Model

POJO Style的Inference Model

POJO Style Inference Model

- 纯Java API, 不依赖于任何计算框架, 不需要特别的上下文
- 可使用于单机Java/Scala程序, Web Serving, Cluster Serving包括批处理, 流处理等场景
- 支持Flink*, Spark* Streaming, Storm*, Kafka* 等



```
import com.intel.analytics.zoo.pipeline.inference.AbstractInferenceModel;

public class MyModel extends AbstractInferenceModel {
    public MyModel(int concurrentNum) {
        super(concurrentNum);
    }
    ...
}

public class ServingExample {
    public static void main(String[] args) throws IOException {
        MyModel model = new MyModel();
        model.load(modelPath, weightPath);
        A data = ...
        List<JTensor> inputs = preProcess(data);
        List<JTensor> outputs = model.predict(inputs);
        B results = postProcess(outputs);
    }
}
```


Inference Model 支持多种深度学习框架的模型

Inference Model supports lots of Deep Learning Frameworks

- 支持多种深度学习框架的模型

- BigDL
- Caffe*
- Tensorflow*
- PyTorch*
- OpenVINO*

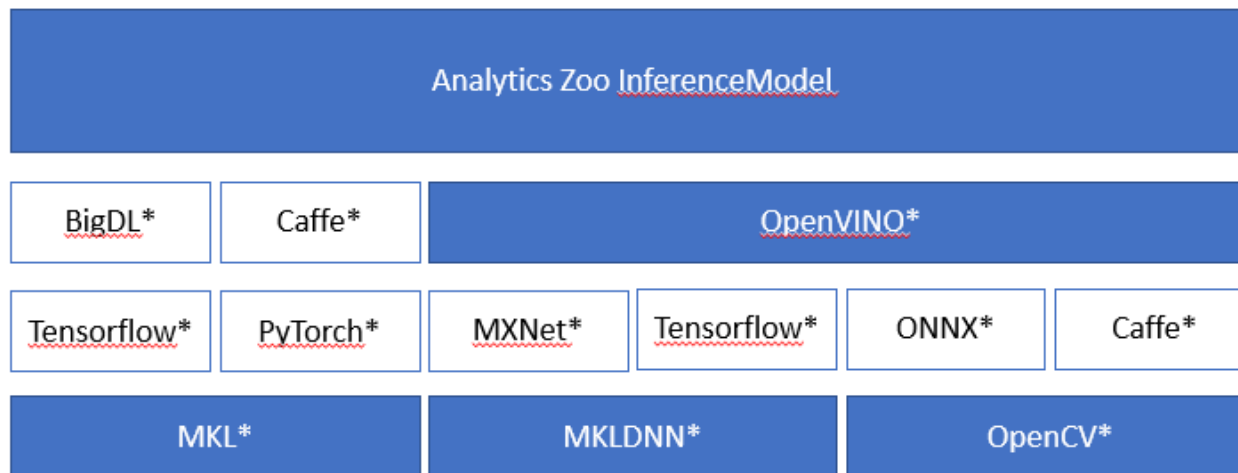
- 简单易用的API

- 加载模型

- Load
- loadCaffe
- loadTF
- loadPyTorch
- loadOpenVINO

- 预测

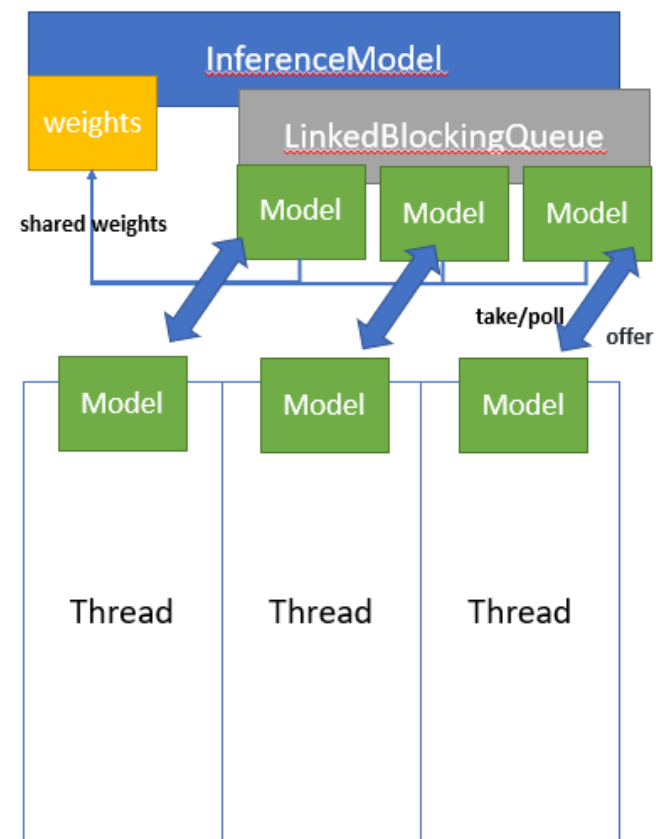
- predict



线程安全的Inference Model

Thread-Safe Inference Model

- 支持线程安全多模型
 - concurrentNum
 - `model = modelQueue.take`
 - autoScalingEnabled
 - `model = modelQueue.poll()`
 - `model = this.originalModel.copy(1)(0)`
- 多模型共享weights



使用OpenVINO*加速模型推理

Model inference accelerating with OpenVINO*

- 支持Image Classification 和Object Detection等
- 支持加载TensorFlow*模型
- 支持模型动态Optimize及Calibrate
- 支持直接加载OpenVINO IR

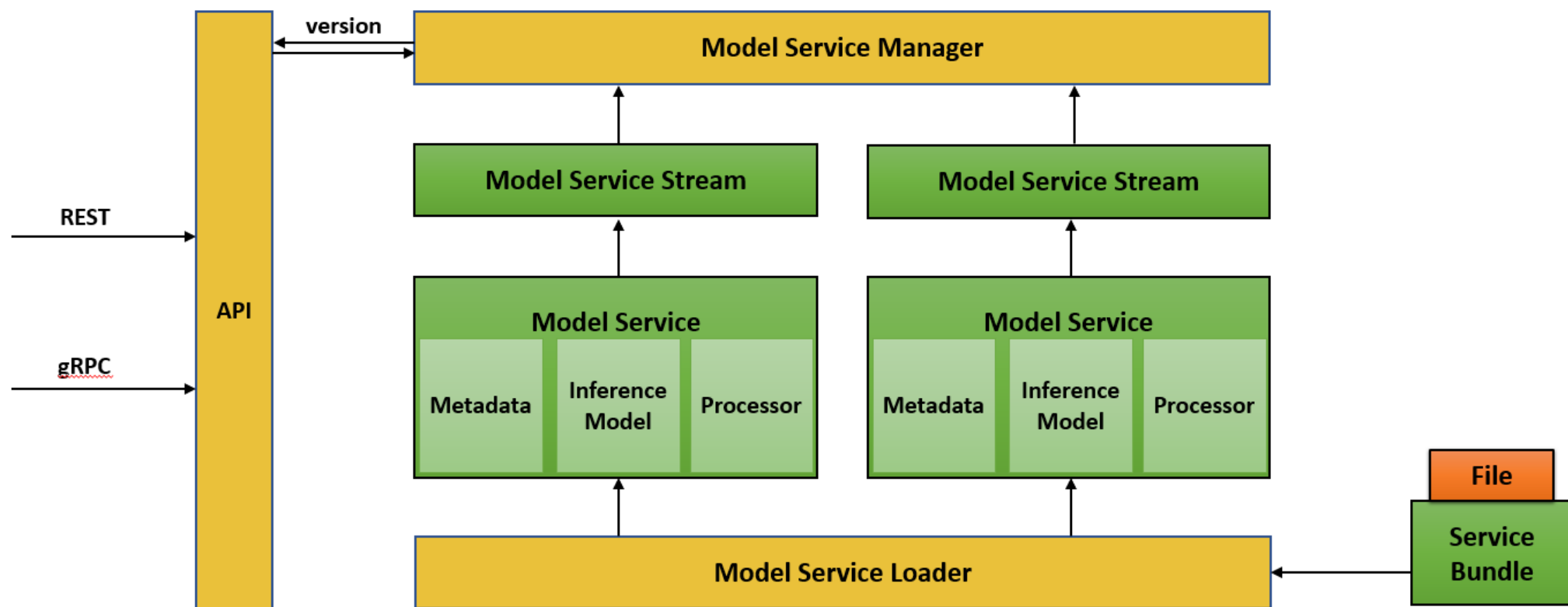
```
from zoo.common.nncontext import init_nncontext
from zoo.feature.image import ImageSet
from zoo.pipeline.inference import InferenceModel
```

```
sc = init_nncontext("OpenVINO Object Detection Inference Example")
images = ImageSet.read(options.img_path, sc,
    resize_height=600, resize_width=600).get_image().collect()
input_data = np.concatenate(
    [image.reshape((1, 1) + image.shape) for image in images], axis=0)
```

```
model = InferenceModel()
model.load_tf(options.model_path, backend="openvino",
    model_type=options.model_type)
predictions = model.predict(input_data)
```

```
# Print the detection result of the first image.
print(predictions[0])
```

Web Serving



Analytics Zoo *Cluster Serving* 使分布式推理更加简单

Distributed Inference made easy with Analytics Zoo Cluster Serving

部署

- ✓ 一个本地节点或者一个Docker容器
- ✓ 已有的 Flink*/YARN*/Spark*/K8S* 集群

使用



1

- 一条命令:
- 启动Docker容器以及Zoo Cluster Serving

- 此命令指定:
- 输入 和 输出 的队列名字
- 模型 的文件路径
- 预/后处理 的文件路径
- 集群 的访问路径



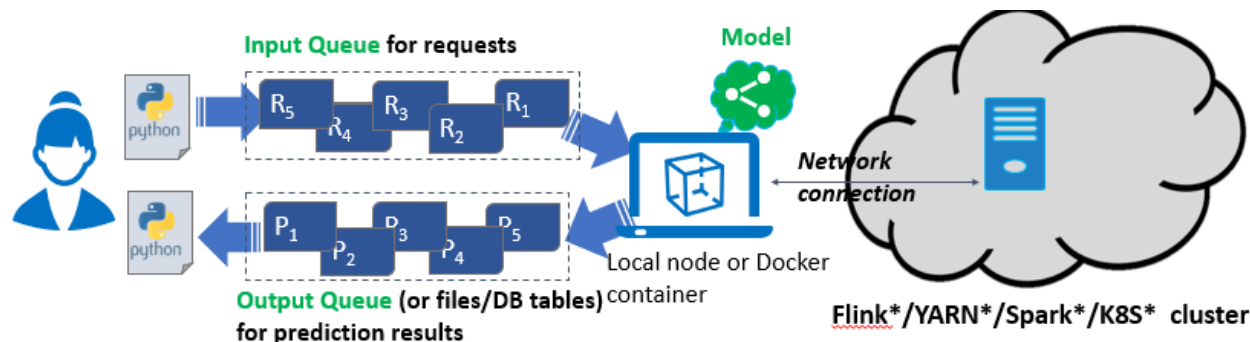
2

- 一个简单的Python脚本:
- 将请求数据发送到 **Input Queue**
- 从 **Output Queue** (或文件/数据库) 获得推理结果



3

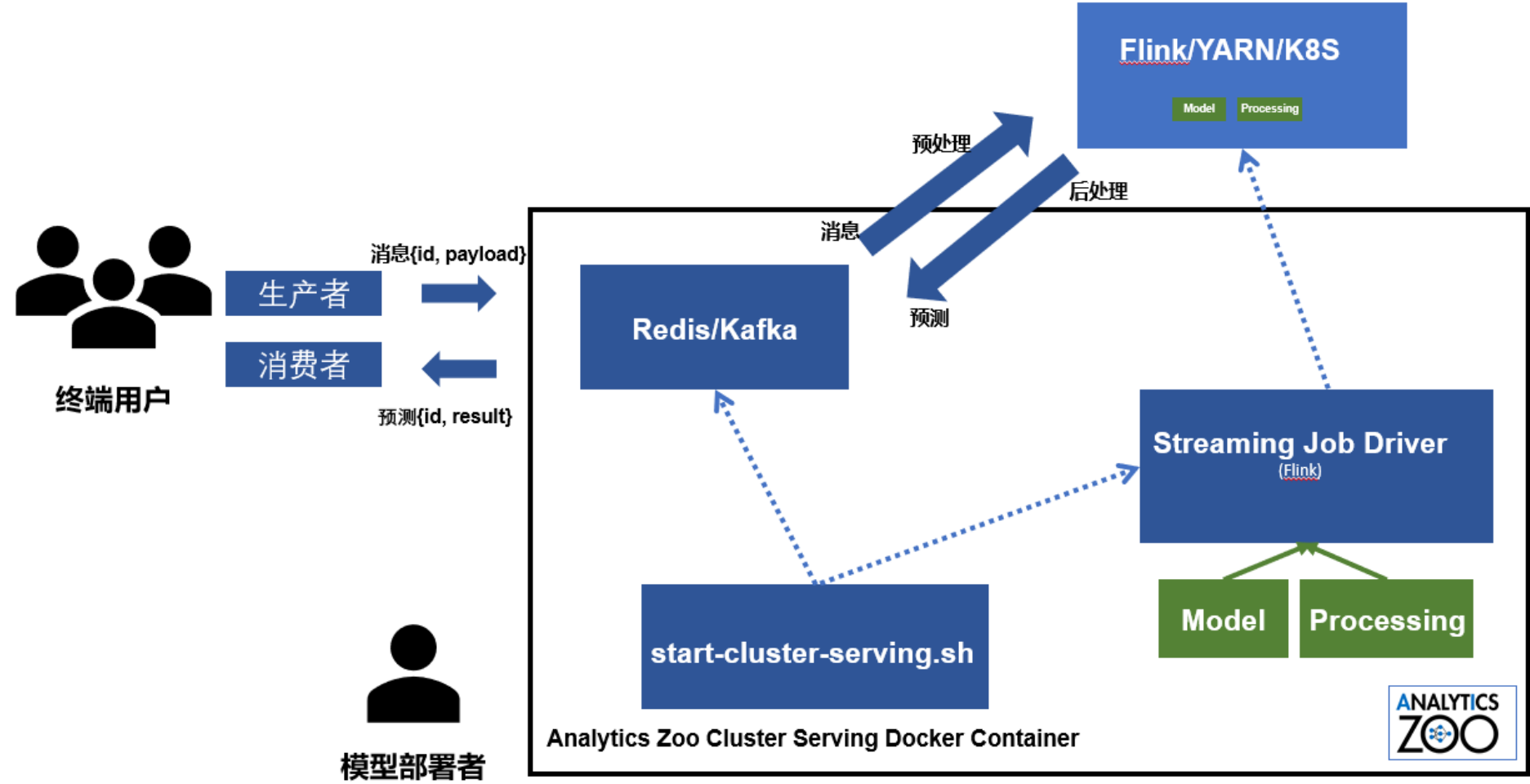
- Analytics Zoo 在集群上自动执行**分布式**、**实时**（流式）模型推理
- 支持 TensorFlow*, Keras*, PyTorch*, Caffe*, BigDL 和 OpenVINO 的模型, 可使用 Int8 加速
- 通过 Flink* 线性扩展



- ✓ 可扩展的分布式推理由Analytics Zoo托管
- ✓ 用户无需为开发和部署复杂的分布式推理方案而费心

Analytics Zoo *Cluster Serving* 使分布式推理更加简单

Distributed Inference made easy with Analytics Zoo Cluster Serving



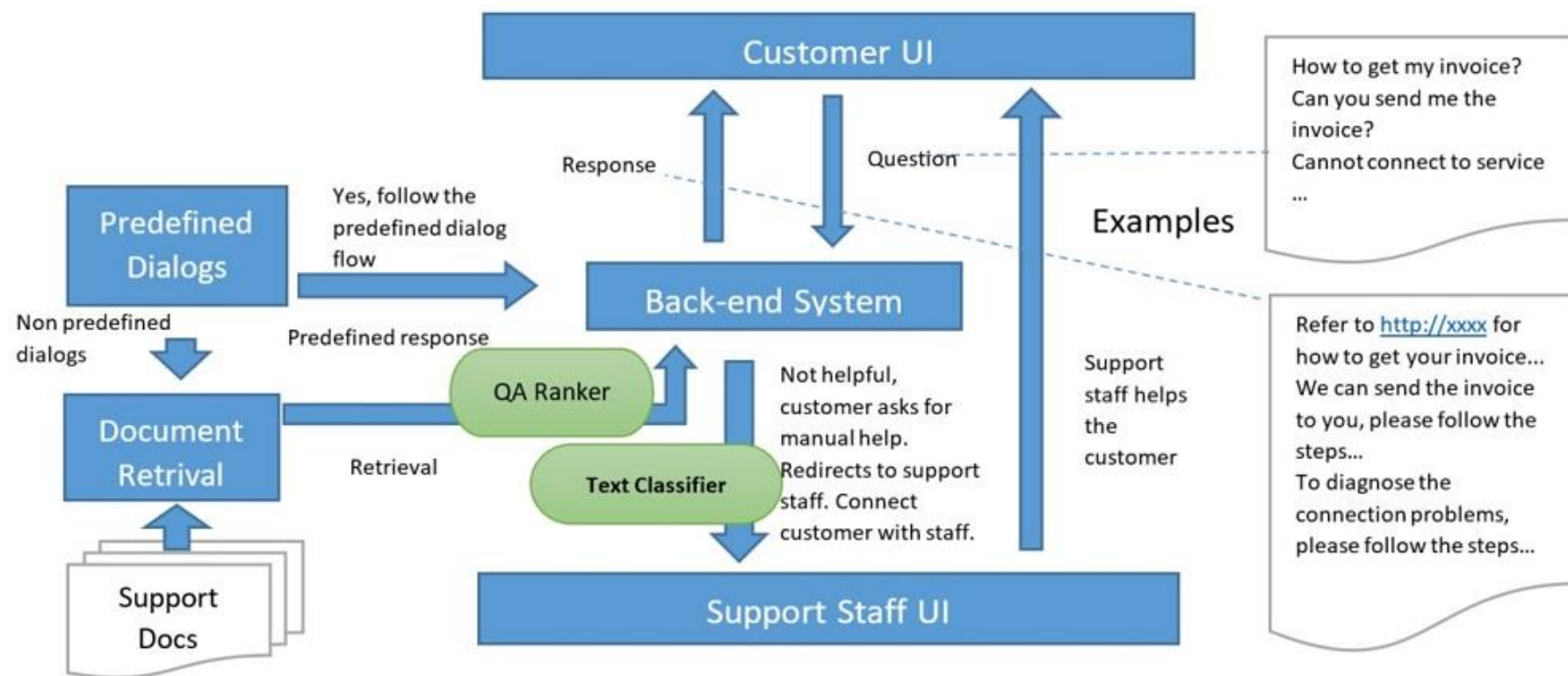
* 文中涉及的其它名称及商标属于各自所有者资产。

03 跨行业的端到端客户案例实践

03 Cross-industry End to End Use Cases

基于NLP的客户服务Chatbot for Microsoft Azure

NLP Based Customer Service Chatbot for Microsoft Azure



<https://software.intel.com/en-us/articles/use-analytics-zoo-to-inject-ai-into-customer-service-platforms-on-microsoft-azure-part-1>
<https://www.infoq.com/articles/analytics-zoo-qa-module/>

云栖社区 > 博客 > 正文

首届！Apache Flink 极客挑战赛强势来袭，重磅奖项等你拿，快来组队报名啦

Ververica

2019-07-24 17:51:26

浏览175

深度学习

大数据

性能优化

机器学习

性能

Apache

钉钉

开源大数据

流计算

大数据分析

ApacheFlink

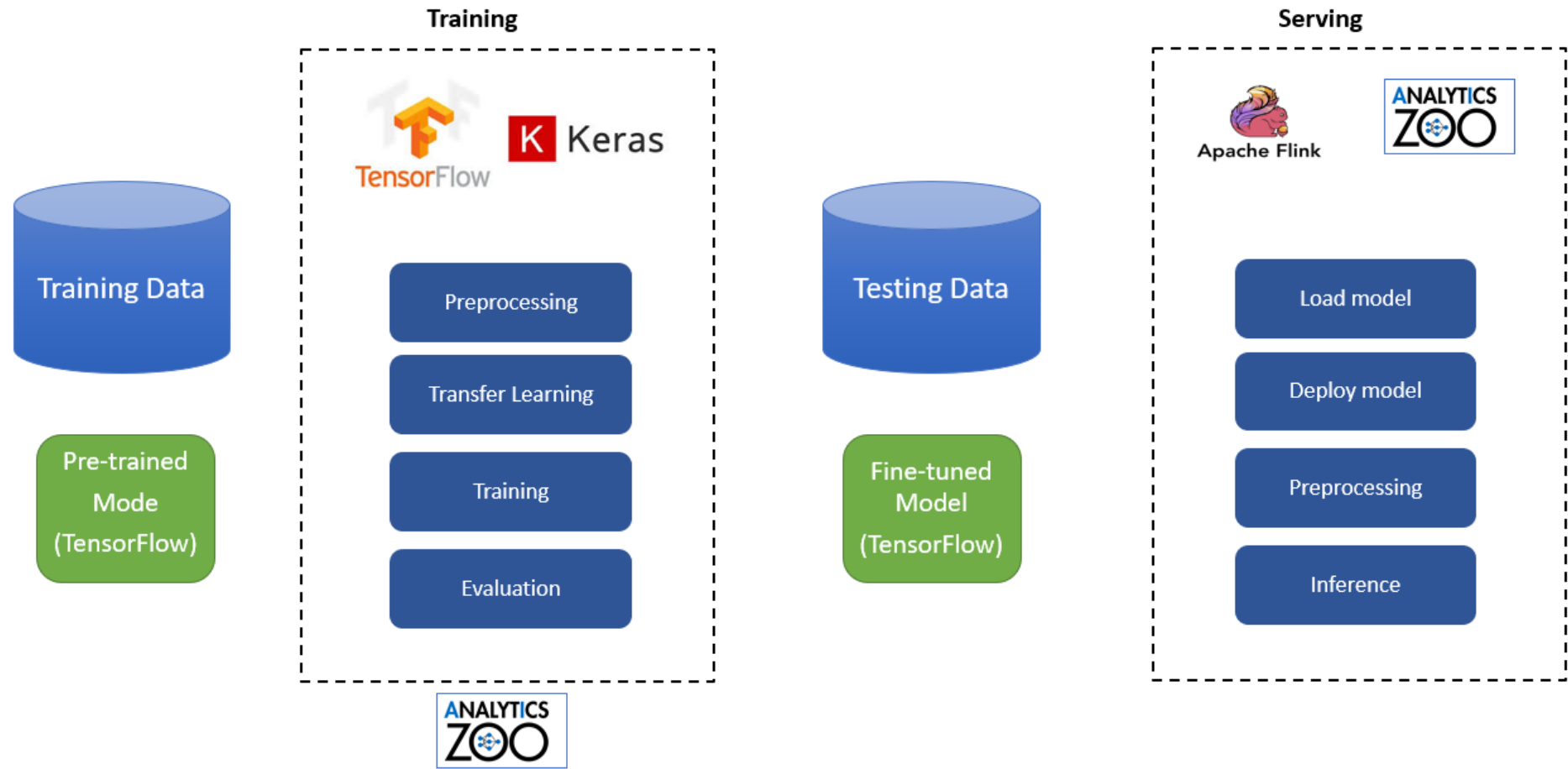
AI及大数据

实时技术

7月24日，阿里云峰会上海开发者大会开源大数据专场，阿里巴巴集团副总裁、计算平台事业部总裁贾扬清与英特尔高级首席工程师、大数据分析和人工智能创新院院长戴金权共同发布首届 Apache Flink 极客挑战赛。



Apache Flink* 极客挑战赛垃圾图片分类



* 文中涉及的其它名称及商标属于各自所有者资产。

使用Analytics Zoo作迁移学习

Transfer Learning with Analytics Zoo

- TFNet load TensorFlow* Saved Model
- Add extra layers
- Training with Estimator

```
val originalModel = TFNet.fromSavedModel(modelPath, inputs, outputs)
```

```
val model = Sequential[Float]()  
model.add(originalModel)  
model.add(new SpatialAveragePooling[Float](2, 2, globalPooling = true))  
model.add(new Linear[Float](2048, 100))
```

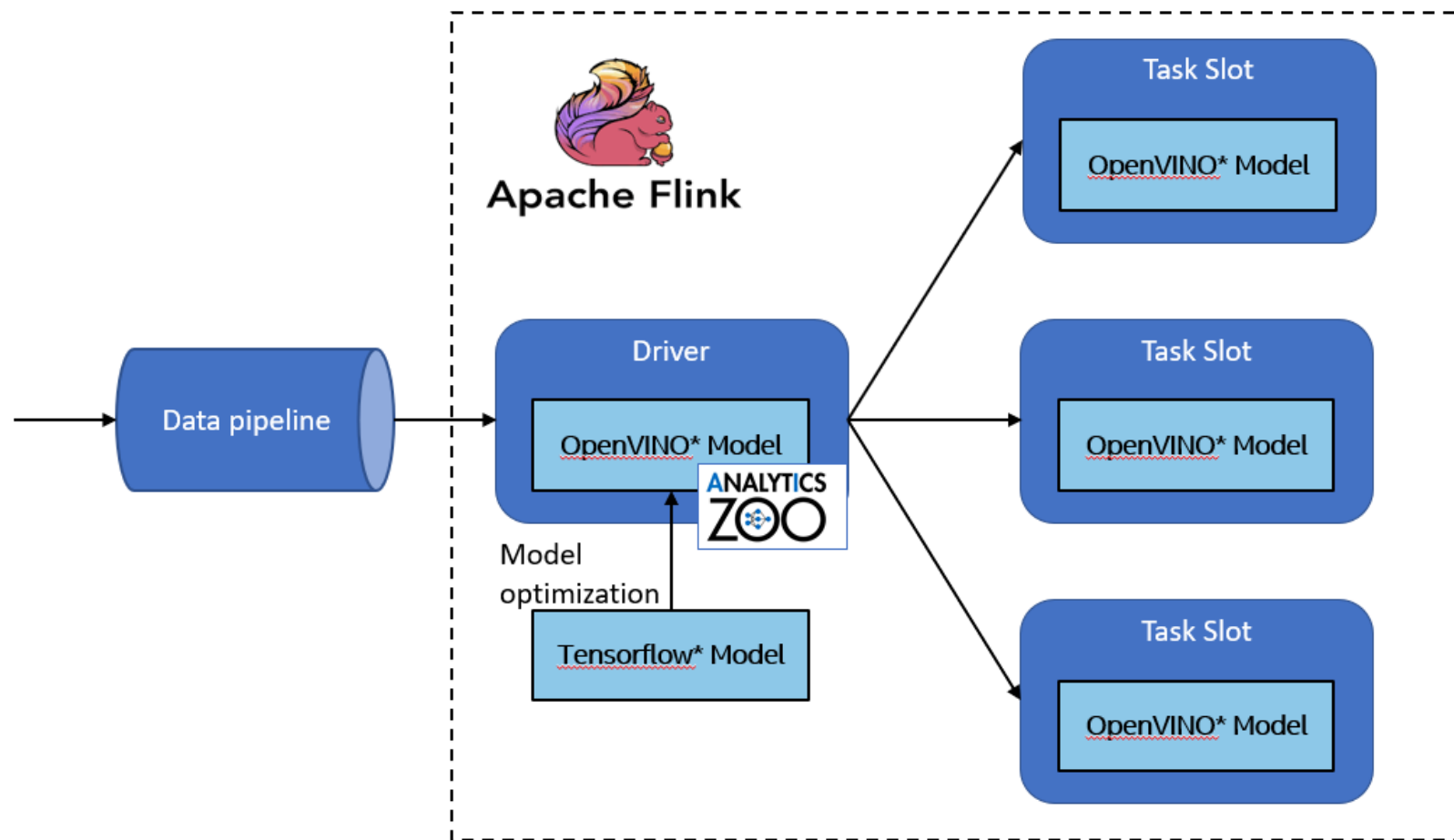
```
val criterion = new CrossEntropyCriterion[Float]()  
val adam = new Adam[Float]()  
val validations = Array(new Top1Accuracy[Float], new Loss[Float])  
val localEstimator = LocalEstimator(model, criterion, adam, validations,  
threadNum)
```

```
val trainData = Cifar10DataLoader.loadTrainData(imageDirPath)  
.filter(_.label() <= 100).slice(0, 10 * batchSize)  
val testData = Cifar10DataLoader.loadTestData(imageDirPath)  
.filter(_.label() <= 100).slice(0, 10 * batchSize)
```

```
localEstimator.fit(trainData, testData,  
ImageProcessing.labeledBGRIImageToMiniBatchTransformer,  
batchSize, epoch)
```

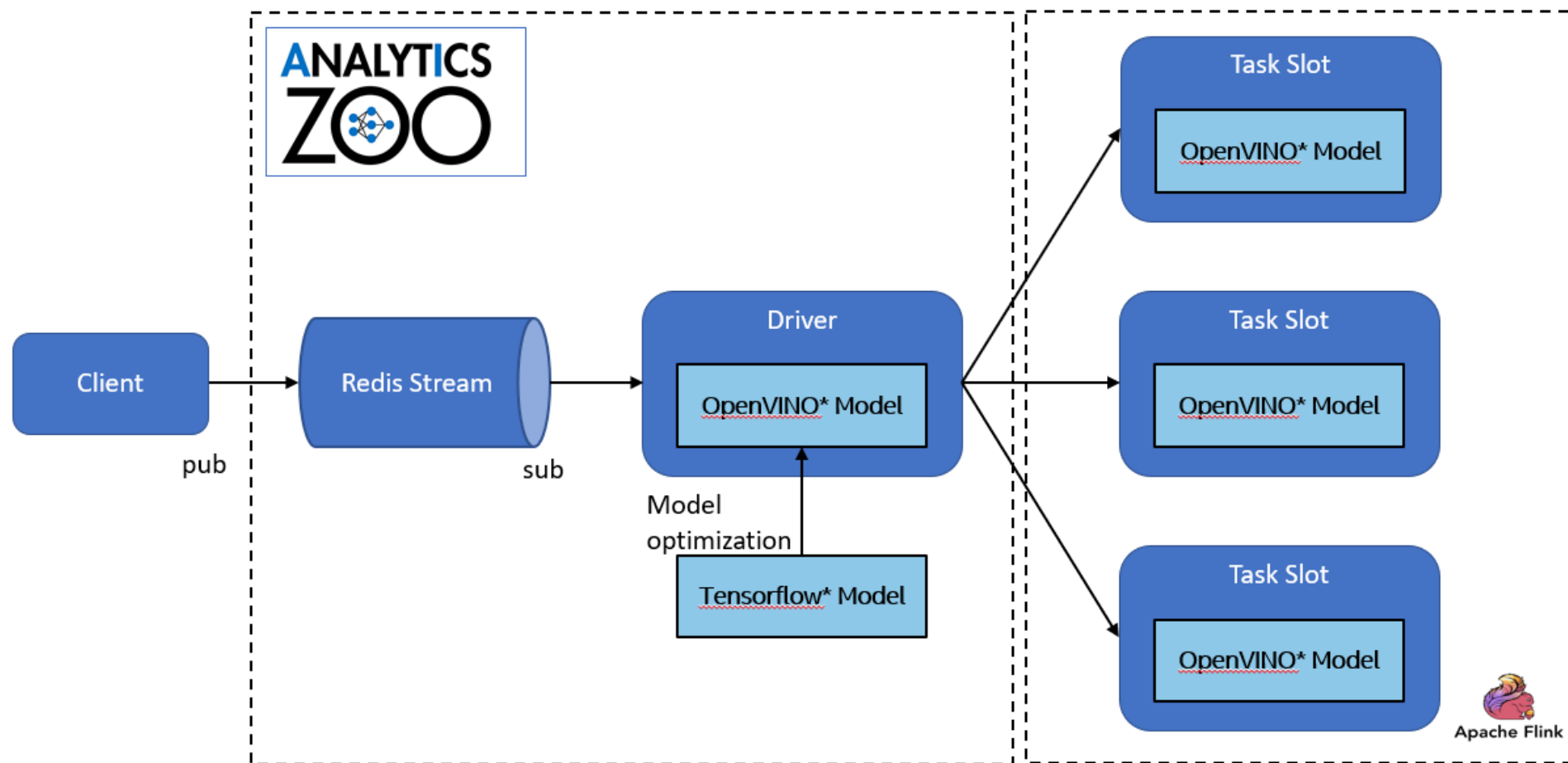
在Apache Flink*中使用Analytics Zoo进行分布式模型推理

Distributed Model Serving with Analytics Zoo in Apache Flink*



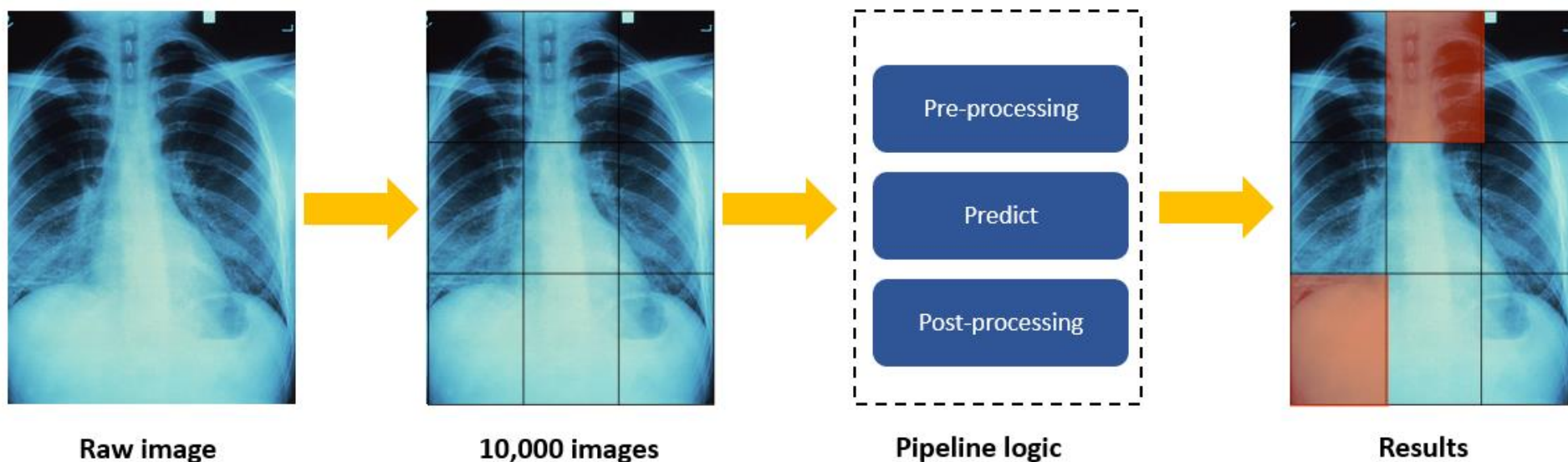
使用Analytics Zoo Cluster Serving进行分布式模型推理

Distributed Model Serving with Analytics Zoo Cluster Serving



使用Analytics Zoo Cluster Serving加速医疗影像分析

Accelerate medical image analysis with Analytics Zoo Cluster Serving



- 结果正确，但性能不可接受，每张原始图片需要1-2小时的处理与预测时间，很难扩展
- 性能瓶颈：预处理（split, crop, resize and normalization），推理

Unacceptable performance, 1-2 hours of processing and prediction, hard to scale
Performance bottlenecks: preprocessing (split, crop, resize and normalization)

<https://en.wikipedia.org/wiki/X-ray>

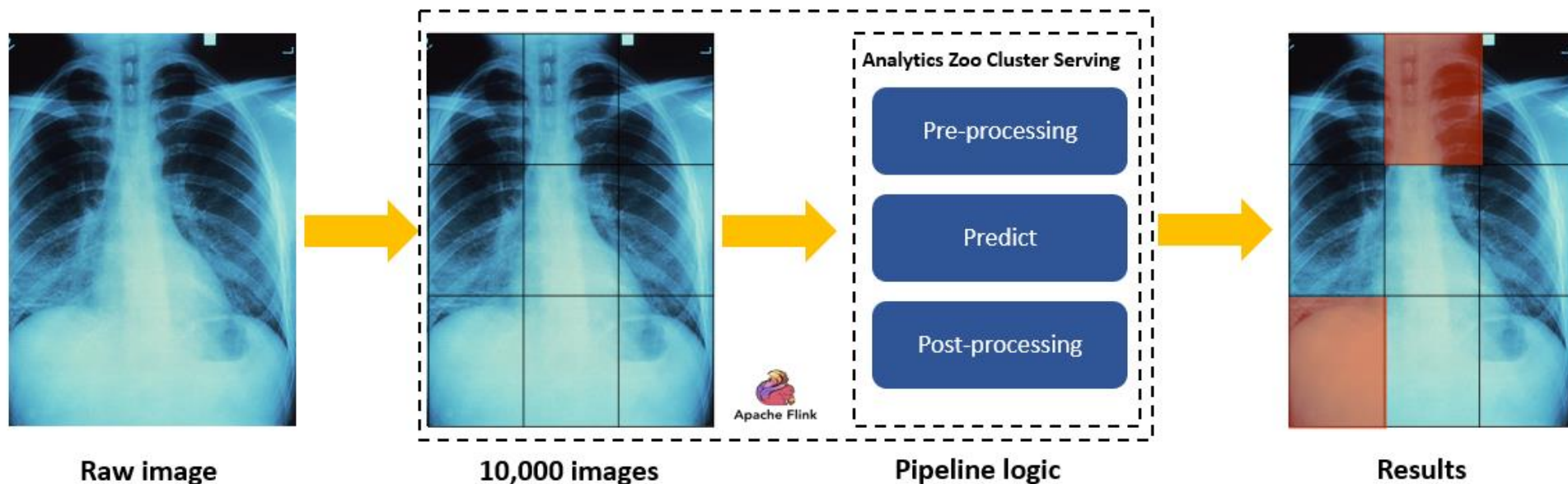
* 文中涉及的其它名称及商标属于各自所有者资产。



Software

使用Analytics Zoo Cluster Serving加速医疗影像分析

Accelerate medical image analysis with Analytics Zoo Cluster Serving



- 使用Analytics Zoo Cluster Serving, 处理时间: 1-2小时→秒级
- 并发式图像处理使用Apache Flink*与Analytics Zoo(OpenCV*)
- 并发式模型推理使用Analytics Zoo(Caffe*, MKLDNN)
- 易于实现及扩展

With Analytics Zoo cluster serving, 1-2 hours → seconds
Parallel image processing using Apache Flink * and Analytics Zoo (OpenCV *)
Parallel model inference using Analytics Zoo (Caffe *, MKLDNN)
Easy to implement and scale

其他跨行业的端到端客户案例实践

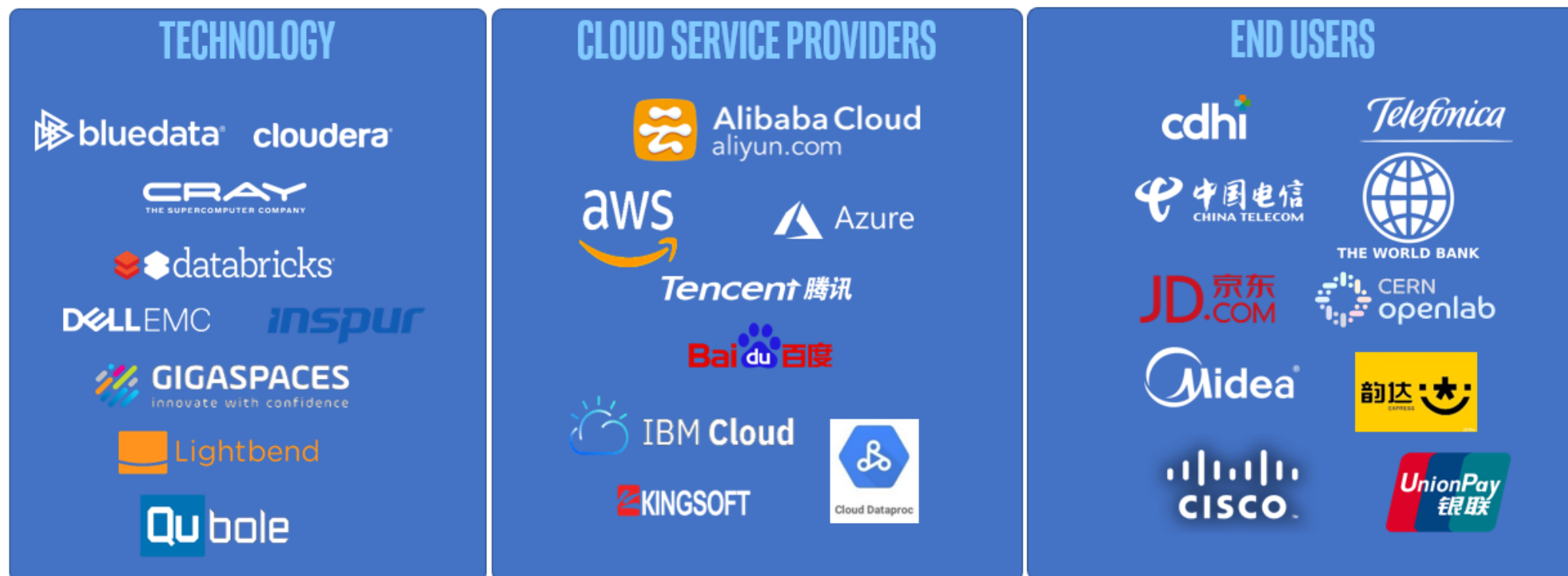
Other End to END Use Cases Examples

- Office Depot*: 基于用户 Session 行为的产品推荐
 - <https://software.intel.com/en-us/articles/real-time-product-recommendations-for-office-depot-using-apache-spark-and-analytics-zoo-on>
 - <https://conferences.oreilly.com/strata/strata-ca-2019/public/schedule/detail/73079>
- 美的*: 工业视觉检测云平台
 - <https://software.intel.com/en-us/articles/industrial-inspection-platform-in-midea-and-kuka-using-distributed-tensorflow-on-analytics>
 - <https://www.intel.cn/content/www/cn/zh/analytics/artificial-intelligence/midea-case-study.html>
- CERN*: 基于深度学习的高能物理粒子事件分类
 - <https://db-blog.web.cern.ch/blog/luca-canali/machine-learning-pipelines-high-energy-physics-using-apache-spark-bigdl>
 - <https://databricks.com/session/deep-learning-on-apache-spark-at-cerns-large-hadron-collider-with-intel-technologies>

更多的案例实践

And Many More

Not a full list



software.intel.com/AlonBigData

Thanks



Distributed, High-Performance
Deep Learning Framework
for Apache Spark



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



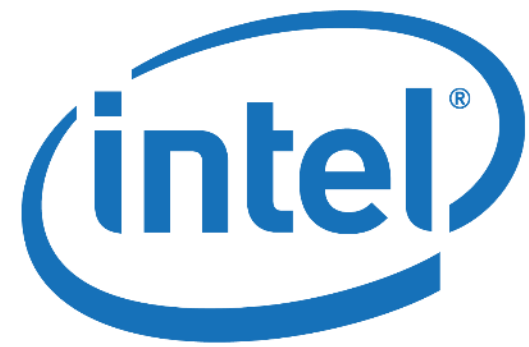
Unified Analytics + AI Platform
Distributed TensorFlow, Keras, PyTorch and BigDL on
Apache Spark



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

Accelerating Data Analytics + AI Solutions At Scale





Software