

# XLearning : 360深度学习调度平台架构设计

李远策

360大数据基础架构团队负责人

# XLEARNING

“是一款支持多种机器学习、深度学习框架的调度系统。基于Hadoop Yarn完成了对TensorFlow/MXNet/ Caffe/ Theano/ PyTorch/ Keras/ XGBoost等常用框架的集成，同时具备良好的扩展性和兼容性。”

## 为什么要设计XLearning平台

- 服务器资源如何调度（CPU、GPU、MEM）
- 训练数据和训练模型的存储管理
- 深度学习作业管理（状态跟踪、日志查看、Metrics信息）
- 多种深度学习框架的环境部署和多版本管理

## 平台设计考量

- 可否能与现有调度平台融合
- 训练数据的统一管理和高效存取
- 与原生深度学习框架代码兼容、性能一致
- 要支持多租户管理和资源隔离
- 开发和运维工作量

## XLearning架构设计



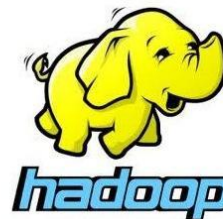
Caffe



PYTORCH

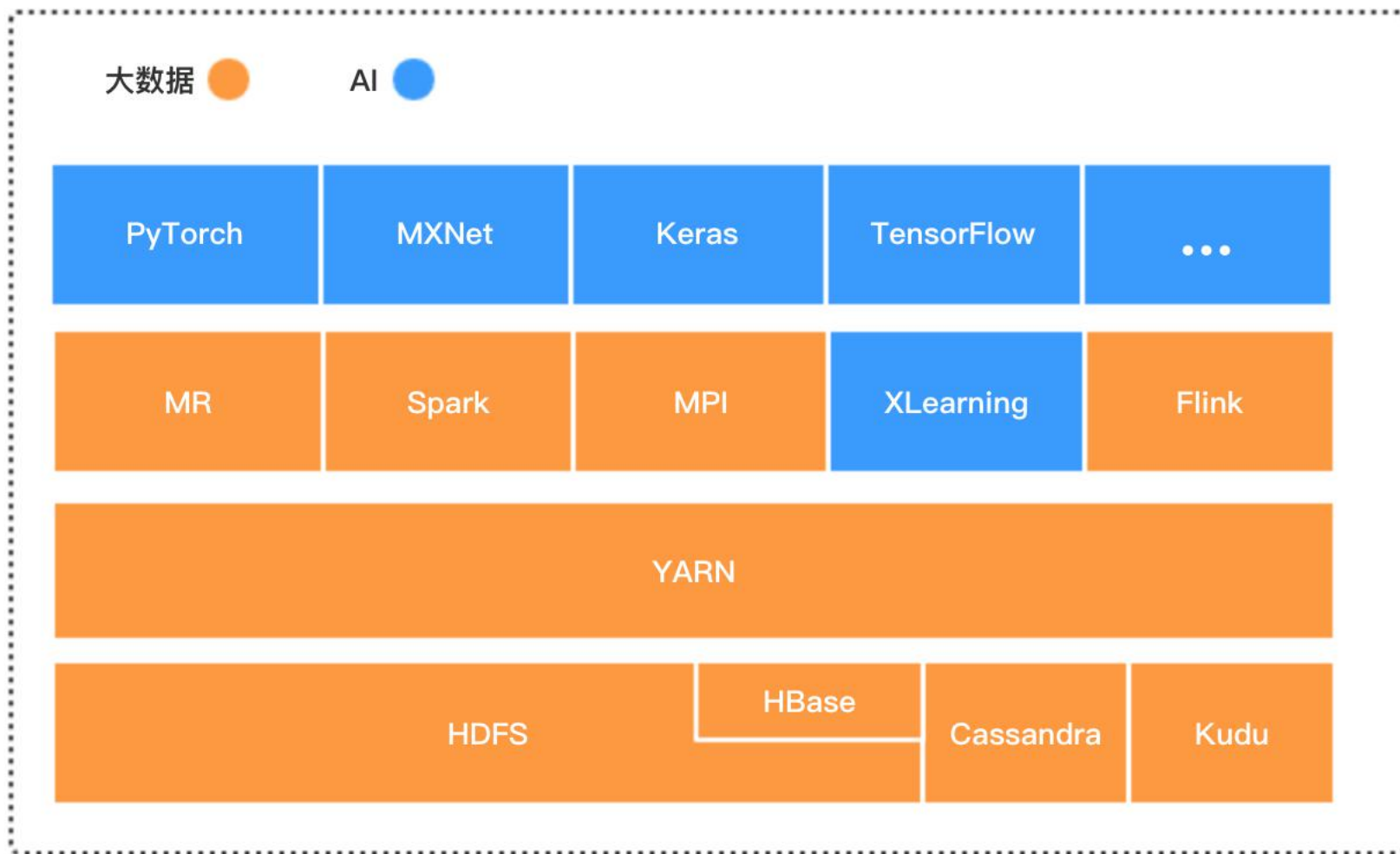
+

Spark

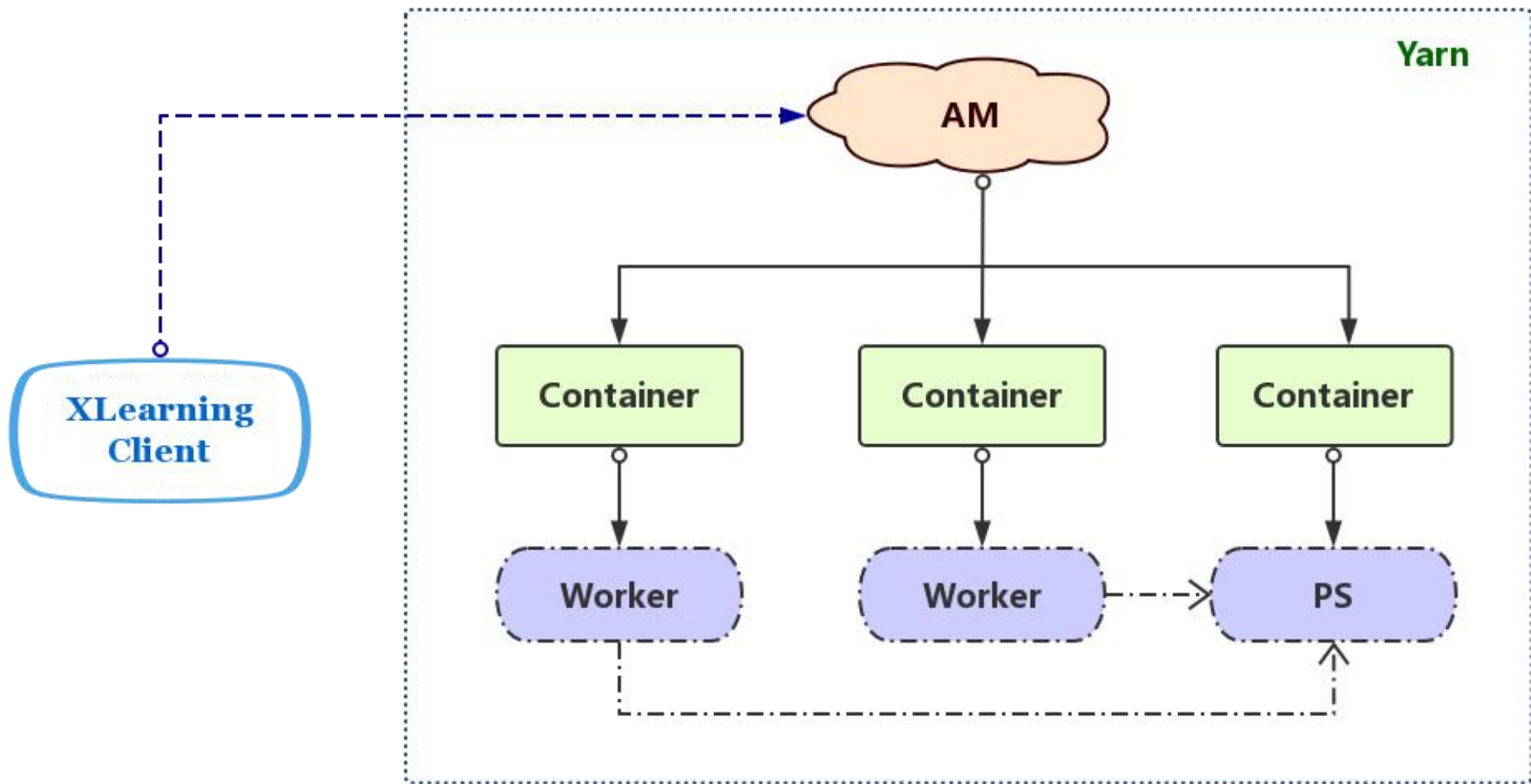


大数据与人工智能的融合平台

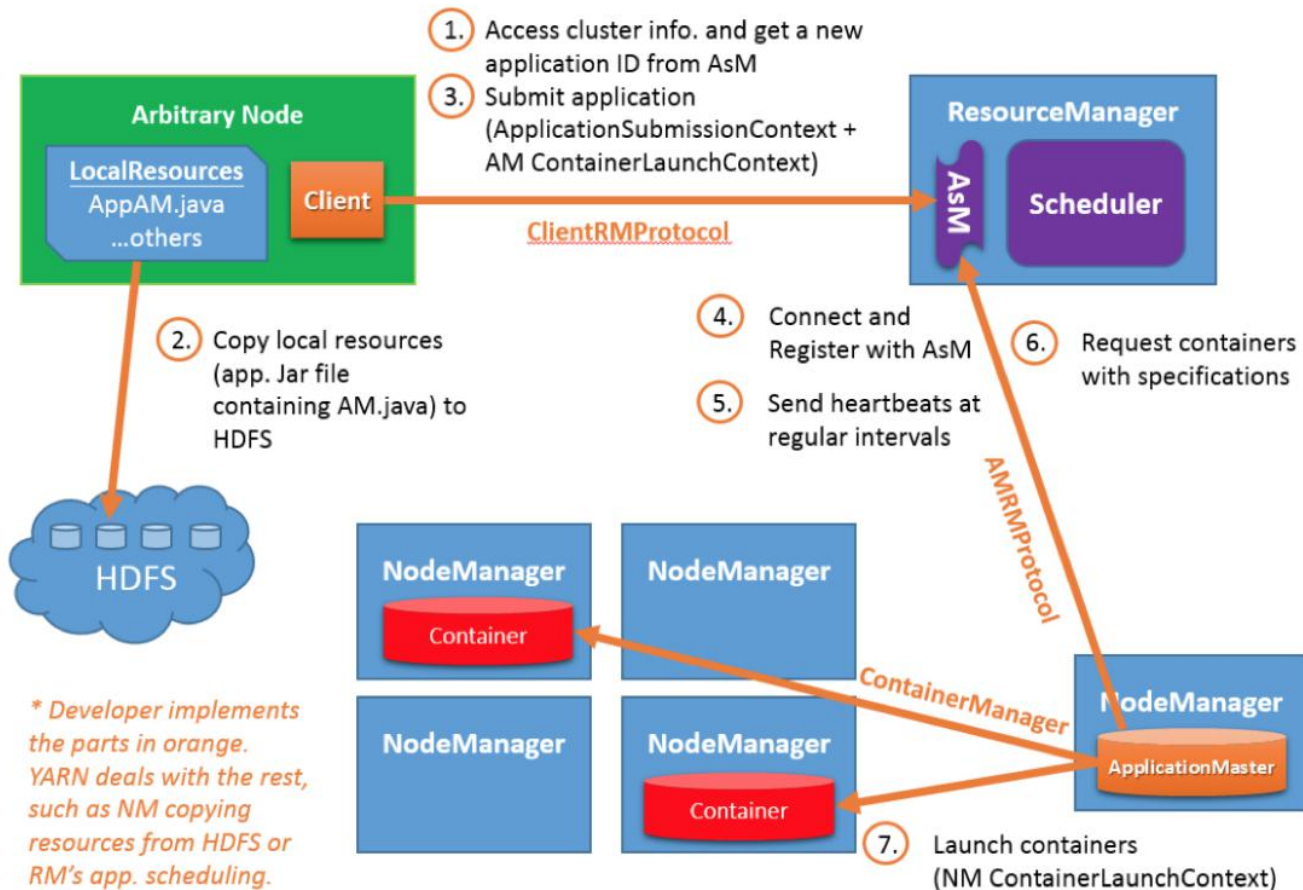
## 大数据+AI平台架构



## XLearning架构简图



## XLearning 执行流程





## XLearning主要功能&特点介绍

- 支持主流的分布式和单机版深度学习框架
- 同时支持同一个深度学习框架的多个版本和用户自定义版本
- 实现了分布式TensorFlow作业ClusterSpec结构的自动构建
- 支持GPU资源调度和隔离，感知GPU设备亲和性（需要Yarn的支持）
- 基于HDFS的数据统一存取，支持多种数据存取模式
- 支持在集群上创建临时GPU虚拟机，解决Debug及个性化GPU需求
- 集成Nvidia Digits系统
- 友好的Web页面，方便查看作业日志、训练进度，硬件资源实时负载信息图表化
- 兼容原生深度学习框架的代码
- 训练效果和性能与原生框架保持一致

## XLearning使用介绍

```
$XLEARNING_HOME/bin/xl-submit \  
  --app-type "tensorflow" \  
  --app-name "tf-demo" \  
  --input /tmp/data/tensorflow#data \  
  --output /tmp/tensorflow_model#model \  
  --files demo.py,dataDeal.py \  
  --launch-cmd "python demo.py" \  
  --worker-num 2 \  
  --worker-memory 10G \  
  --worker-cores 4 \  
  --worker-gcores 2 \  
  --ps-num 1 \  
  --ps-memory 2G \  
  --ps-cores 4\  

```

```
# 作业类型为“TensorFlow”  
# 作业名称为 "tf-demo"  
# 输入数据的HDFS路径  
# 输出模型的HDFS路径  
# 本地TF代码文件  
# TF程序启动指令  
# worker数量为2  
# 每个worker需要10G内存  
# 每个worker需要4个CPU  
# 每个worker需要2个GPU卡  
# PS数量为1  
# 每个PS需要1G内存  
# 每个PS需要4个CPU
```

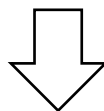
## XLearning数据读取模式介绍

	分片单位	文件格式兼容	读取方式	适用数据量
Download	文件	任意	提前下载	小
Placeholder	文件	任意	直接读取	不限
Stream	Split	受限/用户可定制不同的 <b>inputformat</b>	流式读取	不限

通过 `--input-strategy` 参数选择，默认为Download模式

## XLearning对TensorFlow ClusterSpec的自动构建

```
tf.train.ClusterSpec({  
    "worker": [  
        "worker0.example.com:2222",  
        "worker1.example.com:2222"  
    ],  
    "ps": [  
        "ps0.example.com:2222",  
        "ps1.example.com:2222"  
    ]  
})
```



```
tf.train.ClusterSpec(json.loads(os.environ["TF_CLUSTER_DEF"]))
```

## Yarn首页显示的XLearning作业信息



### All Applications

#### Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	VCores Reserved	GCores Used	GCores Total	GCores Reserved
10284	0	5	10279	15	217 GB	3.43 TB	0 B	34	432	0	7	56	0

#### User Metrics for dr.who

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Containers Pending	Containers Reserved	Memory Used	Memory Pending	Memory Reserved	VCores Used	VCores Reserved
0	0	0	0	0	0	0	0 B	0 B	0 B	0	0

#### Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation
Fair Scheduler	[MEMORY, CPU, GPU]	<memory:1024, vCores:1, gCores:0>

Show 20 entries

ID	User	Name	Application Type	Queue	Application Priority	StartTime	FinishTime	State	FinalStatus	Running Containers	Allocated VCores	Allocated GCores	Allocated Memory(MB)	Reserved Memory
<a href="#">application_1506323180152_10472</a>	xitong	tf-demo	TENSORFLOW	root.default	NORMAL	Mon Dec 18 12:10:34 +0800 2017	N/A	RUNNING	UNDEFINED	4	13	4	24576	0
<a href="#">application_1506323180152_10470</a>	yarn	xunjian_WordCount	MAPREDUCE	root.xunjian	NORMAL	Mon Dec 18 12:08:48 +0800 2017	Mon Dec 18 12:09:00 +0800 2017	FINISHED	SUCCEEDED	N/A	N/A	N/A	N/A	N/A
<a href="#">application_1506323180152_10469</a>	xitong	tf-demo	TENSORFLOW	root.default	NORMAL	Mon Dec 18 12:06:17 +0800 2017	Mon Dec 18 12:09:11 +0800 2017	KILLED	KILLED	N/A	N/A	N/A	N/A	N/A

## XLearning作业首页



### Tensorflow Application application\_1506323180152\_10472

All Containers:

Container ID	Container Host	GPU Device ID	Container Role	Container Status	Start Time	Finish Time	Reporter Progress
container_e10_1506323180152_10472_01_000007	---	0,1	worker	RUNNING	Mon Dec 18 12:10:53 CST 2017	N/A	<div style="width: 100%;"></div>
container_e10_1506323180152_10472_01_000008	---	0,1	worker	RUNNING	Mon Dec 18 12:10:56 CST 2017	N/A	<div style="width: 100%;"></div>
container_e10_1506323180152_10472_01_000002	---	-	ps	RUNNING	Mon Dec 18 12:10:53 CST 2017	N/A	N/A

View TensorBoard:

Tensorboard Info
<a href="http://---.qihoo.net:54201">http://---.qihoo.net:54201</a>

Save Model

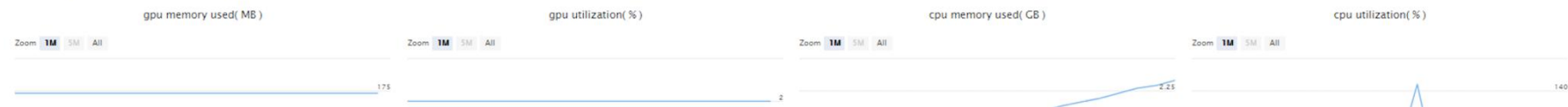
saved the model completed!

Saved timeStamp	Saved path
2017-12-18 12:11:48	/tmp/tensorflow_model/interResult/interResult_2017_12_18_12_11_48

container\_e10\_1506323180152\_10472\_01\_000007 Metrics:



container\_e10\_1506323180152\_10472\_01\_000008 Metrics:



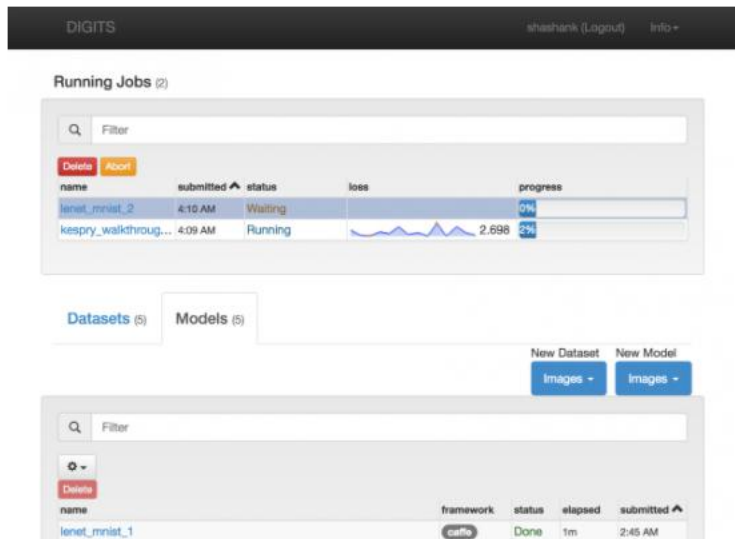
## XLearning GPU虚拟机功能

```
17/09/15 16:27:29 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:30 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:31 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:32 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:33 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:34 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:35 INFO Client: Waiting for vpc login command and password...
17/09/15 16:27:36 INFO Client: Received vpc login command and password from container_e358_1505442218042_0004_01_000002
17/09/15 16:27:36 INFO Client: Login command:ssh root@gpu:~.ssh,ls,cat,ls,cat -p 59136
17/09/15 16:27:36 INFO Client: Password:xLLVIm
```

### 特点&优点:

- 利用集群上GPU的空闲资源创建虚拟机，相比静态分配可以有效提供GPU资源的利用率；
- 采用Docker实现，更轻量，用户方便定制系统环境；
- 随时创建和销毁。

## XLearning集成Nvidia Digits系统



Schedule, monitor, and manage  
neural network training jobs





## XLearning平台设计经验分享

- 提前想清楚平台的价值是什么，能带来什么收益
- 功能设计源于用户需求，不要创造需求，想象出来的功能往往没人用
- 优先解决用户公共需求和最痛点的问题
- 最大程度复用现有的技术工具，尽量避免重复造轮子，劳民伤财
- 重视平台的兼容性（代码、性能、效果等），直接影响平台的推广难度
- WIKI、FAQ能在推广过程中节省大量精力

GIAC | 全球互联网架构大会  
GLOBAL INTERNET ARCHITECTURE CONFERENCE

GIAC

全球互联网架构大会

GLOBAL INTERNET ARCHITECTURE CONFERENCE



扫码关注GIAC公众号

2017.thegiac.com