# YARN 3.1 and Beyond !

Naganarasimha

ApacheCon

Standard Chartered

Here for good

# Introduction

## Naganarasimha GR

- Apache Hadoop PMC

- Contributing since 5 years

- Senior BigData Architect @ Standard Chartered Bank SG

- Contributed in key YARN features
  - Node attributes
  - ATS V2

# Standard Chartered

## Leading the way in international banking

**86,000**
We employ more than 86,000 people around the world

**125**
Our staff come from 125 different countries

**150**
More than 150 years in business

**>1100**
Branches worldwide

**60**
We operate in 60 countries globally

**Best Consumer Digital Bank**
Global Finance 2017

**100**
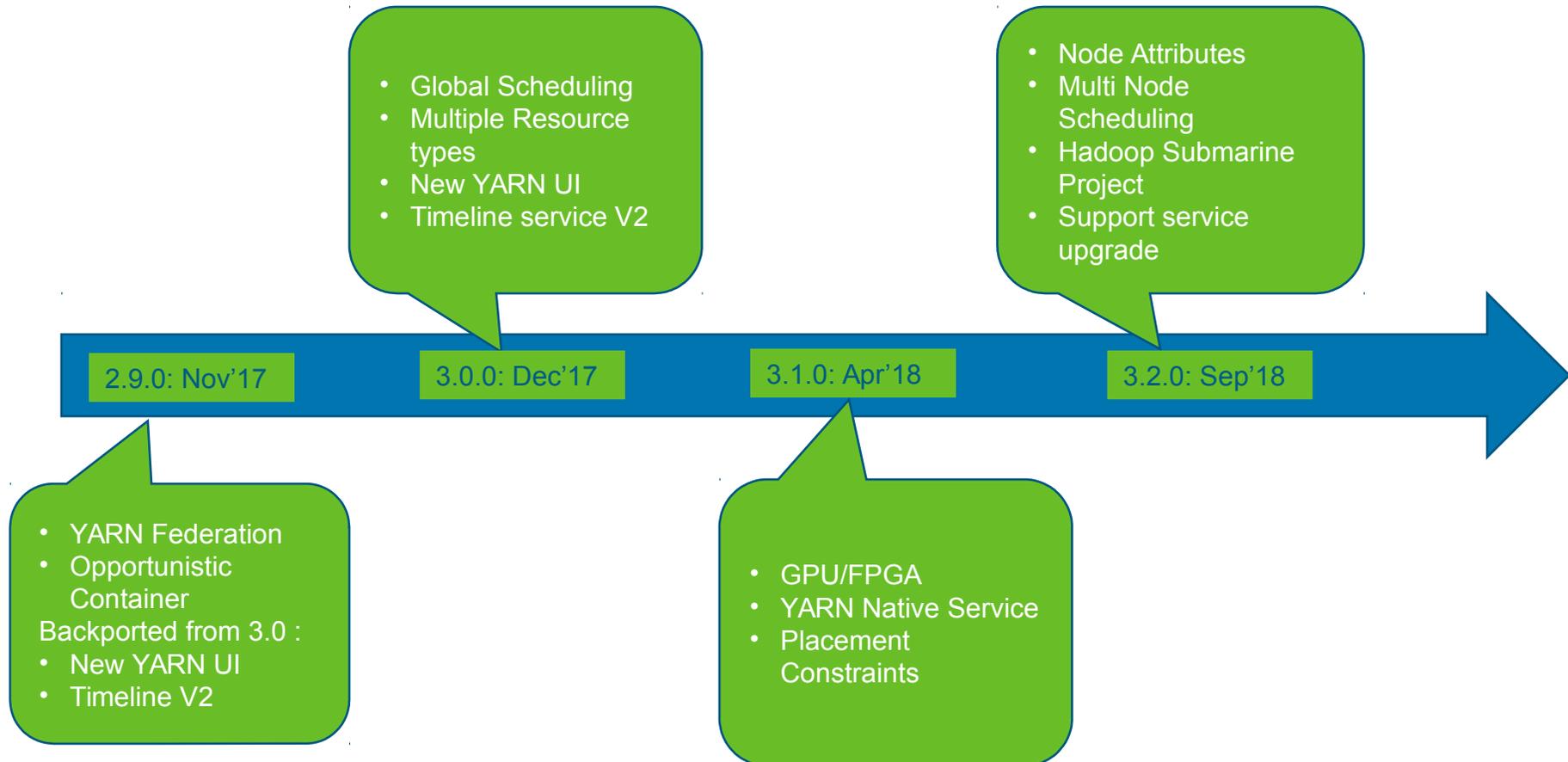Among the top 100 largest companies listed on the London Stock Exchange
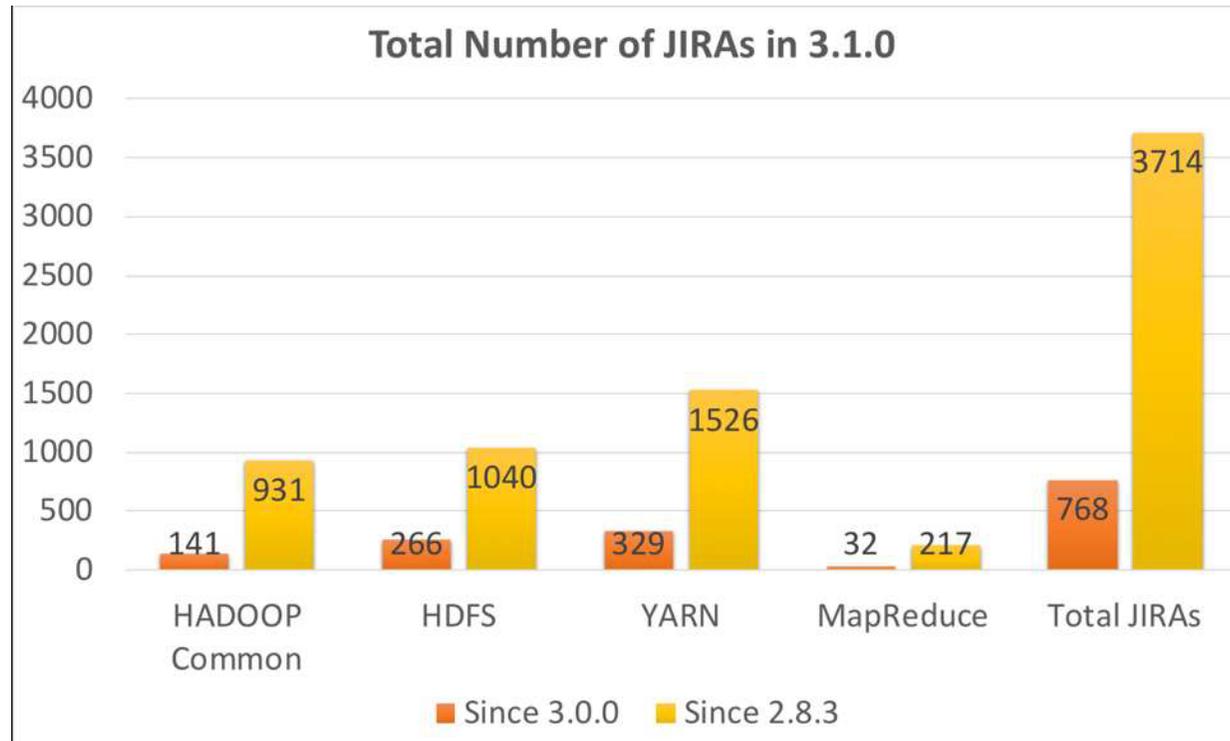
**2**
We're listed on two of Asia's largest stock exchanges
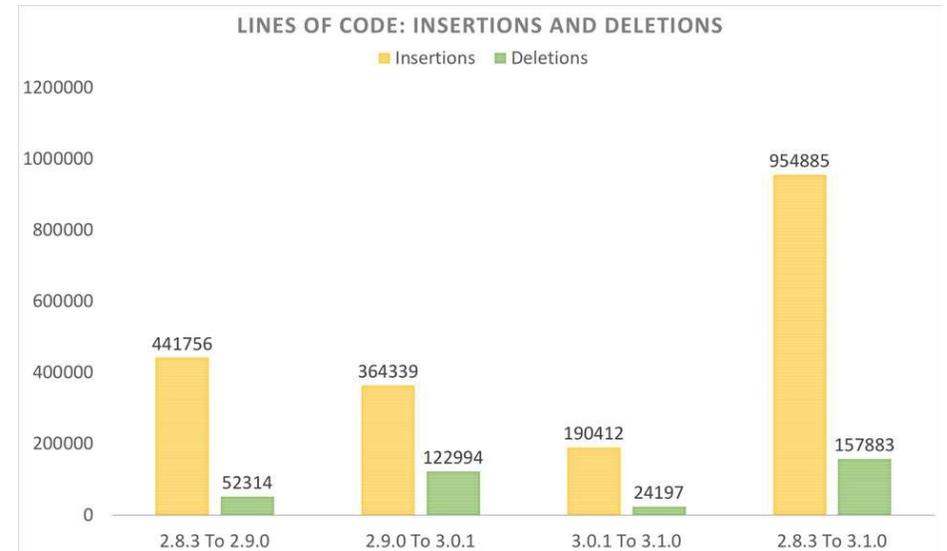
- Introduction
- Present
- Upcoming

# Introduction

Standard Chartered

Here for good

# A brief timeline of Hadoop Releases:

- Global Scheduling
- Multiple Resource types
- New YARN UI
- Timeline service V2

- Node Attributes
- Multi Node Scheduling
- Hadoop Submarine Project
- Support service upgrade

**2.9.0: Nov'17**     **3.0.0: Dec'17**     **3.1.0: Apr'18**     **3.2.0: Sep'18**

- YARN Federation
- Opportunistic Container

Backported from 3.0 :
- New YARN UI
- Timeline V2

- GPU/FPGA
- YARN Native Service
- Placement Constraints

# Community Update: JIRAs in 3.1.0



**Total Number of JIRAs in 3.1.0**

| | HADOOP Common | HDFS | YARN | MapReduce | Total JIRAs |
|---|---|---|---|---|---|
| Since 3.0.0 | 141 | 266 | 329 | 32 | 768 |
| Since 2.8.3 | 931 | 1040 | 1526 | 217 | 3714 |

# Community Update: Source code changes



**FILES CHANGED**

| 2.8.3 To 2.9.0 | 2.9.0 To 3.0.1 | 3.0.1 To 3.1.0 | 2.8.3 To 3.1.0 |
|---|---|---|---|
| 3959 | 3913 | 2613 | 7477 |

**LINES OF CODE: INSERTIONS AND DELETIONS**

Insertions ▪ Deletions

| | 2.8.3 To 2.9.0 | 2.9.0 To 3.0.1 | 3.0.1 To 3.1.0 | 2.8.3 To 3.1.0 |
|---|---|---|---|---|
| Insertions | 441756 | 364339 | 190412 | 954885 |
| Deletions | 52314 | 122994 | 24197 | 157883 |

# Categorization:

**Topics**

**User Centric**

- Right Node Selection
- Ease of use
- Right resource to select

**Admin Centric**

- Faster Allocation
- Ease of monitoring
- Resource Isolation
- Ease of Configuration

# YARN Overview :

# Apache Hadoop 3.1

Here for good

## Moving towards Global & Fast Scheduling

### YARN-5139

- **Problems**
  - Current design of one-node-at-a-time allocation cycle can lead to suboptimal decisions.
  - Several coarse grained locks.

- **With this, we improved to**
  - Look at several nodes at a time
  - Fine grained locks
  - Multiple allocator threads
  - YARN scheduler can allocate 3k+ containers per second ≈ 10 mil allocations / hour!
  - **10X throughput gains**
  - Much better placement decisions

# Traditional scheduling



**Root**

(memory=60G, vcores=40, gpu=10)

**Compliance**

(min=10%, max =30%)

**Retail**

(min=40%, max =50%)

**Engineering**

(min=50%, max =100%)

**Compliance1**

(min=50%, max =50%)

**Compliance2**

(min=50%, max =50%)

**Retail1**

(min=50%, max =50%)

**Retail2**

(min=50%, max =50%)

**Engineering1**

(min=50%, max =50%)

**Enginering2**

(min=50%, max =50%)

# Global Scheduling explained

# Better placement strategies (YARN-6592)

**User Centric**

- **Past**
  - Supported constraints in form of Node Locality.

- **Now YARN can support a lot more use cases**
  - Co-locate the allocations of a job on the same rack (**affinity**)
  - Spread allocations across machines (**anti-affinity**) to minimize resource interference
  - Allow up to a specific number of allocations in a node group (**cardinality**)

# Better placement strategies (YARN-6592)

- **Affinity**



- **Anti-affinity**

# Absolute Resources Configuration in CS – YARN-5881

**Admin Centric**

- Gives ability to configure Queue resources as below
  <memory=24GB, vcores=20, yarn.io/gpu=2>

- Enables admins to assign different quotas of different resource-types

- No more **"Single percentage value limitation for all resource-types"**

**Root**
(memory=60G, vcores=40, gpu=10)

**Compliance**
(memory=60G, vcores=40, gpu=10)

**Retail**
(memory=60G, vcores=40, gpu=0)

**Engineering**
(memory=20G, vcores=40, gpu=10)

## Auto Creation of Leaf Queues - YARN-7117

- **Easily map a queue** explicitly to user or group with out additional configs
  - For e.g, User X comes in, automatically create a queue for user X with a templated capacity requirements

- **Auto created Queues will be**
  - created runtime based on user mapping
  - cleaned up after use
  - adhering to ACLs

# Usability : UI 1/2

**Admin Centric**

# Usability : UI 2/2

## Timeline  Service 2.0 Improvements

Understanding and Monitoring a Hadoop cluster
itself is a BigData problem

- Using **HBase as backend** for better
  scalability for read/write
- More robust storage fault tolerance
- Migration and compatibility with v.1.5

# Resource profiles and custom resource types

- **YARN** supported only **Memory** and **CPU**

- **Now**
  - A generalized vector for all resources
  - Admin could add arbitrary resource types!

Ease of resource requesting model using profiles for apps



| Profile | Memory | CPU | GPU |
|---------|--------|-----|-----|
| Small | 2 GB | 4 Cores | 0 Cores |
| Medium | 4 GB | 8 Cores | 0 Cores |
| Large | 16 GB | 16 Cores | 4 Cores |

## GPU support on YARN

- **Why?**
  - No need to setup separate clusters
  - Leverage shared compute!

- **Why need isolation?**
  -  Multiple processes use the single GPU will be:
    - Serialized.
    - Cause OOM easily.

- **GPU isolation on YARN:**
  - Granularity is for per-GPU device.
  - Use cgroups / docker to enforce isolation.

Tensorflow 1.2

CUDA Library 5.0

Ubuntu 14:04

docker

**Volume Mount**

GPU Base Lib v1

Host OS

# FPGA on YARN

**User Centric**

- **FPGA isolation on YARN:**
  - Granularity is for per-FPGA device
  - Use Cgroups to enforce the isolation
- Currently, only Intel OpenCL SDK for FPGA is supported.
- Implementation is extensible to other FPGA SDK.

# Services support in YARN

- **A native YARN services framework (YARN-5079)**
  - Native Yarn support to Services
  - Apache Slider retired from Incubator – lessons and key code carried over to YARN

- **Simplified discovery of services via DNS mechanisms: YARN-4757**
  - *regionserver-0.hbase-app-3.hadoop.yarn.site*

- **Application & Services upgrades:  YARN-4726**
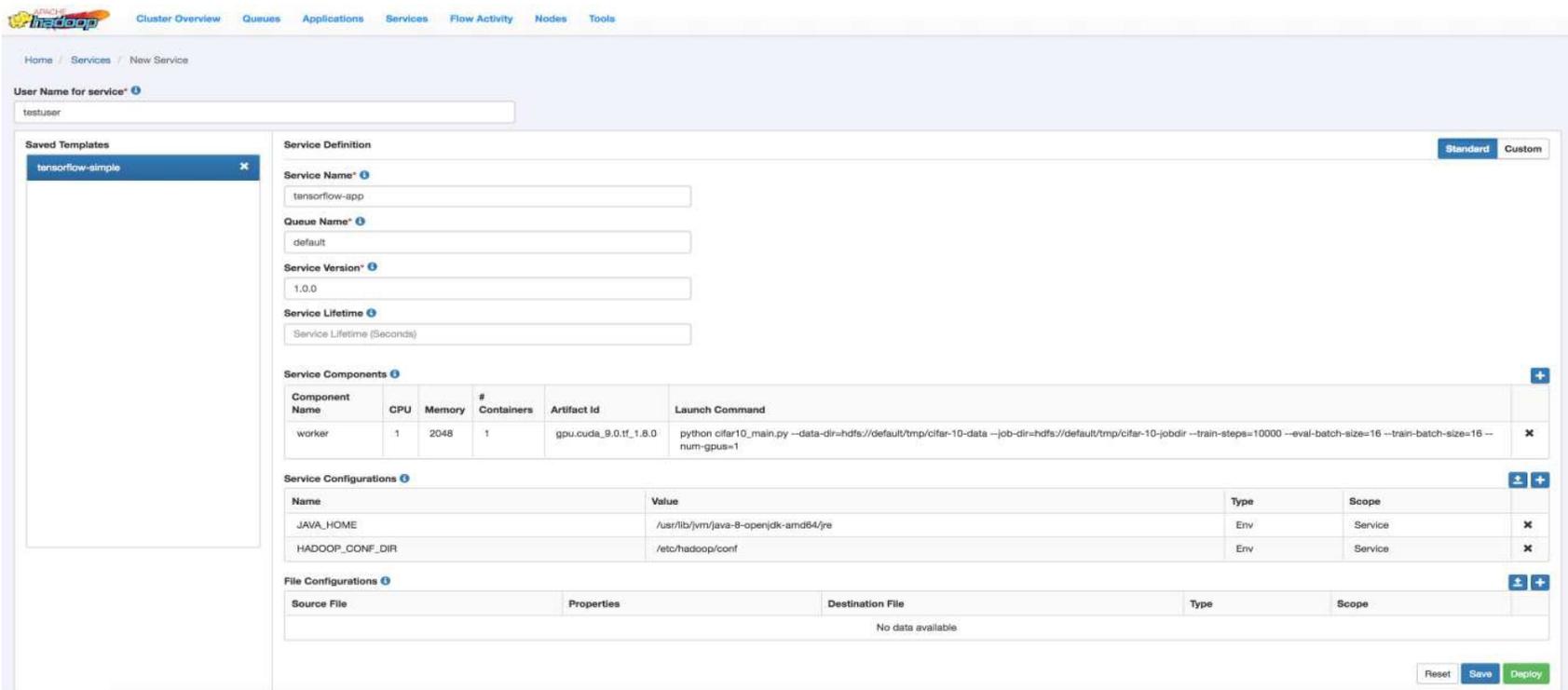  - "Do an upgrade of my HBase app with minimal impact to end-users".

## Simplified APIs for service definitions

**User Centric**

- Applications need simple APIs

- Need to be deployable "easily"

- Simple REST API layer (YARN-4793)

- Spawn services & Manage them

```
"name": "kafka-app-1",
"lifetime": "3600",
"components": [
    {
        "name": "KAFKABROKER",
        "number_of_containers": 3,
        "unique_component_support" : "true",
        "artifact": {
            "id": "registry.eng.hortonworks.com/hwx-assemblies/kafka:0.10.1",
            "type": "DOCKER"
        },
        "launch_command": "sleep 60; /usr/hdp/current/kafka-broker/bin/kafka-server-start.sh /etc/kafka/conf/server.properties",
        "resource": {
            "cpus": 1,
            "memory": "512"
        },
        "configuration": {
            "files": [
                {
                    "type": "PROPERTIES",
                    "dest_file": "/etc/kafka/conf/server.properties",
                    "props": {
                        "broker.id": "${COMPONENT_ID}",
                        "zookeeper.connect": "${CLUSTER_ZK_QUORUM}${SERVICE_ZK_PATH}",
                        "listeners": "PLAINTEXT://kafkabroker${COMPONENT_ID}.${SERVICE_NAME}.${USER}.${DOMAIN}:9092",
                        "zookeeper.session.timeout.ms": "80000",
                        "zookeeper.connection.timeout.ms": "80000"
                    }
                }
            ]
        }
    }
]
}
```

# How to run a new service in YARN ?

# Apache Hadoop 3.2

# Node Attributes (YARN-3409)

- "Take me to a node with JDK 10"

- Node Partition vs. Node Attribute
- **Node Partition**
  - One partition for one node
  - ACL
  - Shares between queues
  - Pre emption enforced
- **Attribute**
  - For Container placement
  - No ACL's and Shares
  - First come first serve



Partition - 1 (15 nodes)
Queue-A (40%)
Queue-B (60%)

Partition - 2 (15 nodes)
A (25%)
B (25%)
C (25%)
D (25%)

**Node 1**
os.type=ubuntu
os.version=14.10
glibc.version=2.20
JDK.version=8u20

**Node 2**
os.type=RHEL
os.version=5.1
GPU.type=x86_64
JDK.version=7u20

**Node 16**
os.type=windows
os.version=7
JDK.version=8u20

**Node 17**
os.type=SUSE
os.version=12
GPU.type=i686
JDK.version=7u20

# Node Attributes (YARN-3409)

- **Distributed Node Attributes**
    - NM can detect its attributes
    - Script based and Config based detection.
    - Attribute prefix : *yarn.nm.io*

- **Centralised Node Attributes**
    - Configured through CLI and REST
    - Admin ACL's to configure
    - Attribute prefix : *yarn.rm.io*

# Node Attributes (YARN-3409)

**Use cases :**

- **Hardware Constraints** : To identify specific kind of resources like
    - GPU, FPGA,
    - SSD, # of disks,
    - InfiniBand,
    - (dual) network cards,
- **Task Constraints:** Task or container specific constraints like
    - To run on specific Operating system versions.
    - Processor architecture
    - software library versions
- **Experimental** : Based on dynamic attributes like
    - Load average,
    - disk usage
    - Network

# Container overcommit (YARN-1011)

- Every user says "Give me 16GB for my task", even though it's only needed at peak
- Each node has some allocated but unutilized capacity. Use such capacity to run opportunistic tasks
- Preempt such tasks when needed

# Auto-spawning of system services (YARN-8048)

**Admin Centric**

- "Start this service when YARN starts"

- "initd for YARN"
- Services are started during the yarn bootstrap
  - For example YARN ATSv2 needs Hbase, so Hbase is system service of YARN.
  - Only Admin can configure
  - Started along with ResourceManager
  - Place spec files under yarn.service.systemservice.dir FS path

```
SYSTEM_SERVICE_DIR_PATH
|---- sync
|    |--- user1
|    |    |---- service1.yarnfile
|    |    |---- service2.yarnfile
|    |--- user2
|    |    |---- service3.yarnfile
|    |    ....
|    |
|---- async
|    |--- user3
|    |    |---- service1.yarnfile
|    |    |---- service2.yarnfile
|    |--- user4
|    |    |---- service3.yarnfile
|    |    ....
|    |
```

## TensorFlow on YARN (YARN-8220)

- Run deep learning workloads on the same cluster as analytics, stream processing etc!

- Integrated with latest TensorFlow 1.8 and has **GPU** support
  - Use simple command to run TensorFlow app by using Native Service spec file (Yarnfile)

    *yarn app -launch distributed-tf <path-to-saved-yarnfile>*

  - A simple python command line utility also could be used to auto-create Yarnfile
    ```
    python submit_tf_job.py
    --remote_conf_path hdfs:///tf-job-conf
    --input_spec example_tf_job_spec.json
    --docker_image gpu.cuda_9.0.tf_1.8.0
    --job_name distributed-tf-gpu
    --user tf-user
    --domain tensorflow.site
    --distributed --kerberos
    ```

# TensorFlow on YARN (YARN-8220)

- Run deep learning workloads on the same cluster as analytics, stream processing etc!

- Integrated with latest TensorFlow 1.8 and has **GPU** support
  - Use simple command to run TensorFlow app by using Native Service spec file (Yarnfile)

    *yarn app -launch distributed-tf <path-to-saved-yarnfile>*

  - A simple python command line utility also could be used to auto-create Yarnfile
    python submit_tf_job.py
    --**remote_conf_path** hdfs:///tf-job-conf
    --**input_spec** example_tf_job_spec.json
    --**docker_image** gpu.cuda_9.0.tf_1.8.0
    --**job_name** distributed-tf-gpu
    --**user** tf-user
    --**domain** tensorflow.site
    --**distributed** --**kerberos**

# TensorFlow on YARN (YARN-8220)

- **Sample Yarnfile for TensorFlow job**

```json
{
  "name": "distributed-tf",
  "version": "1.0.0",
  "components": [
    {
      "name": "worker",
      "dependencies": [],
      "resource": {
        "cpus": 1,
        "memory": "4096",
        "additional" : {
          "yarn.io/gpu" : {
            "value" : 1
          }
        }
      },
      "launch_command": "cd /test/models/tutorials/image/cifar10_estimator && python cifar10_main.py --data-dir=hdfs://default/tmp
      "number_of_containers": 1
    }
  ],
  "kerberos_principal" : {
    "principal_name" : "test-user@EXAMPLE.COM",
    "keytab" : "file:///etc/security/keytabs/test-user.headless.keytab"
  }
}
```

## Other related talks :

- **Deep learning on YARN - Running distributed Tensorflow / MXNet / Caffe / XGBoost on Hadoop clusters**
  - **Speakers** : Wangda Tan
  - Thursday, 27th Sep, 11:20
  - Ballroom

- **Running distributed TensorFlow in production: challenges and solutions on YARN 3.0**
  - **Speakers** : Wangda Tan & Yanbo Liang
  - Thursday, 27th Sep, 15:40
  - Ballroom

Queries ?