# About.html

- Apache Hadoop PMC, ASF Member

- Yahoo! -> Hortonworks

- 10 years of (only) Hadoop
  - Finally the job-adverts asking for "10 years of Hadoop experience" have validity

- 'Rewritten' the Hadoop processing side – Became Apache Hadoop YARN

- Running compute platform teams at Hortonworks: YARN, MapReduce, Slider, container cloud on YARN

# Agenda

- Introduction

- Past

- Present & Future

**HORTONWORKS®**

# Hadoop Compute Platform – Today and Tomorrow

- It's all about data!

- Layers that enable applications and higher order frameworks that interact with data

- Multi-colored YARN
  - Apps
  - Long running services

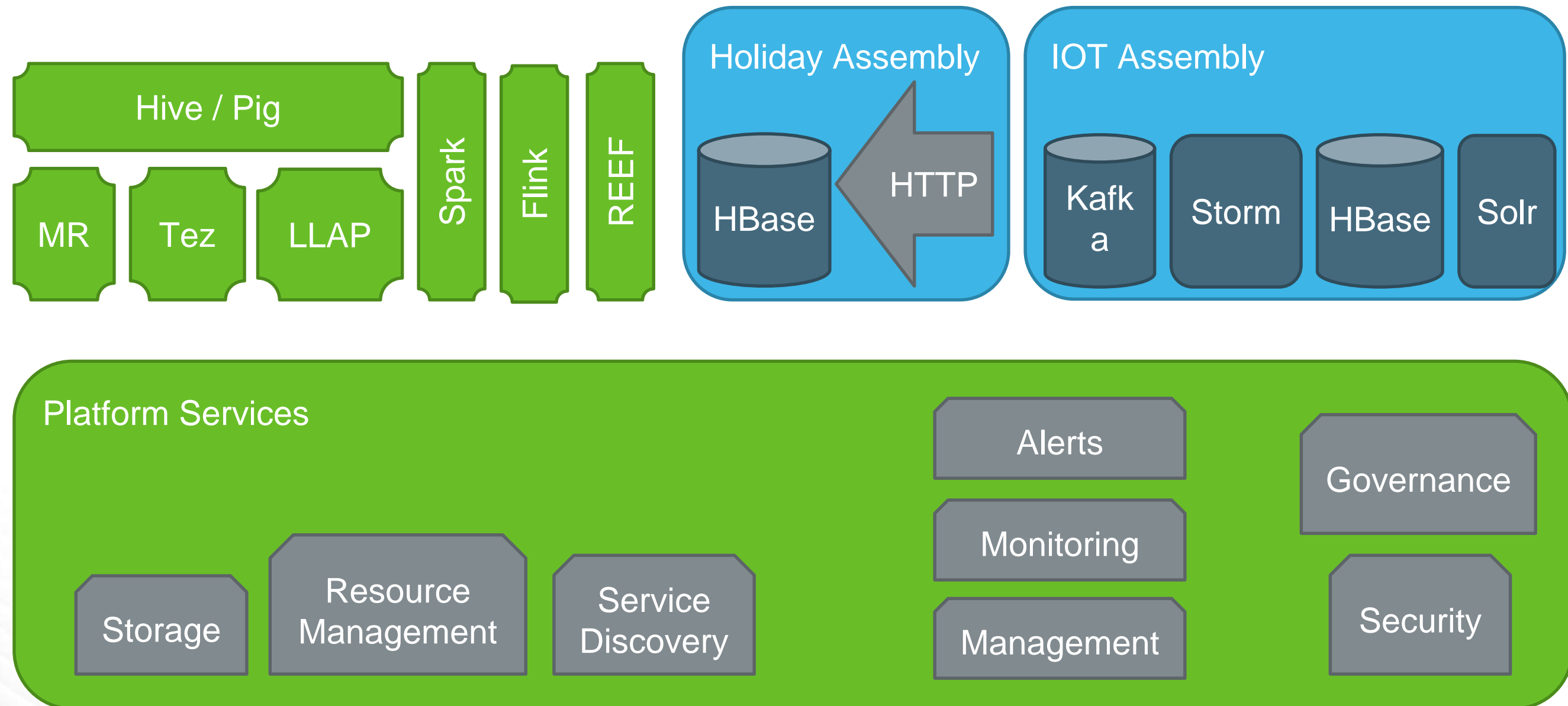- Admins and admin tools (Ambari) for cluster management and monitoring

https://www.flickr.com/photos/happyskrappy/15699919424

HORTONWORKS®

# Why?

- Different asks from different actors
- On isolation, capacity allocations, scheduling

Faster!

More!

SLA!

Everything!

Right now!

Best for my cluster
Throughput
Utilization
Elasticity
Service uptime
Security
ROI

HORTONWORKS®

# Hadoop Compute Platform – Today and Tomorrow

# Past: A quick history

**HORTONWORKS**®

# A brief Timeline: Pre GA

- **Sub-project of Apache Hadoop**

- **Releases tied to Hadoop releases**

- **Alphas and betas**
  - In production at several large sites for MapReduce already by that time

| June-July 2010 | August 2011 | May 2012 | August 2013 |
| --- | --- | --- | --- |

Hortonworks

# A brief Timeline: GA Releases 1/2

- 1st GA
- MR binary compatibility
- YARN API cleanup
- Testing!

- 1st Post GA
- Bug fixes
- Alpha features

- RM Fail-over
- CS Preemption
- Timeline Service V1

- Writable REST APIs
- Timeline Service V1 security

- Rolling Upgrades
- Services
- Node labels

- Moving to JDK 7+
- Pluggable YARN authentication

15 October 2013

24 February 2014

07 April 2014

11 August 2014

18 November 2014

21 Apr 2015

| 2.2 | 2.3 | 2.4 | 2.5 | 2.6 | 2.7 |

Hortonworks

# A brief Timeline: GA Releases 2/2

- Stabilization

- Stabilization
- All community parties brought together

- First alpha!

- Stabilization

- Second alpha

- Next major release

25 January 2016

25 August 2016

03 September 2016

18 October 2016

25 January 2017

22 March 2017

| 2.7.2 | 2.7.3 | 3.0.0-alpha1 | 2.6.5 | 3.0.0-alpha2 | 2.8.0 |

Hortonworks

# Present & Future

HORTONWORKS

Service workloads

Usability

SLAs

Scale

Much faster scheduling

More powerful scheduling

Containers

GPUs / FPGAs

**Hortonworks**

# Last few Hadoop releases

**Apache Hadoop 2.8.0**

**Apache Hadoop 3.0.x**

**Hortonworks**

# Apache Hadoop 2.8.0

**HORTONWORKS**®

# Application priorities

- YARN-1963

- Within a leaf-queue

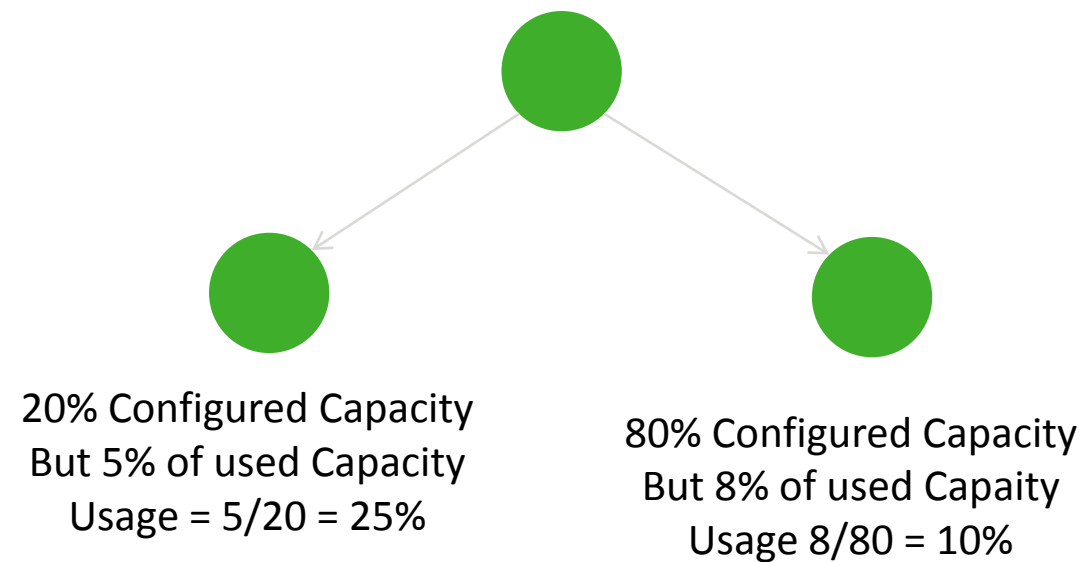| FIFO Policy | App 1 | App 2 | App 3 | App 4 |

| FIFO Policy With priorities | App 4 P2 | App 3 P3 | App 2 P4 | App 1 P8 |

**HORTONWORKS®**

# Queue priorities

- Today
  - Give to the least satisfied queue first



20% Configured Capacity
But 5% of used Capacity
Usage = 5/20 = 25%

80% Configured Capacity
But 8% of used Capaity
Usage 8/80 = 10%

- With priorities
  - Give to the highest priority queue first

# Preemption within a queue

- **Between apps of different priorities**

App 1 P1
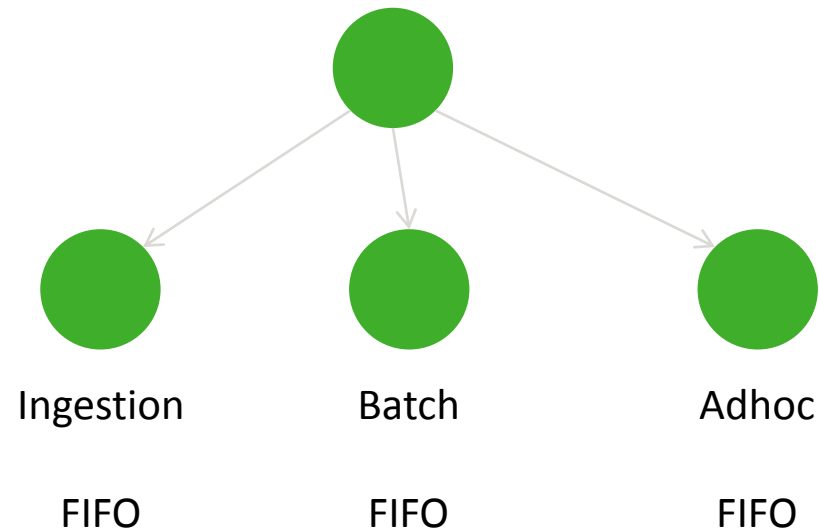
App 1 P2

App 1 P3

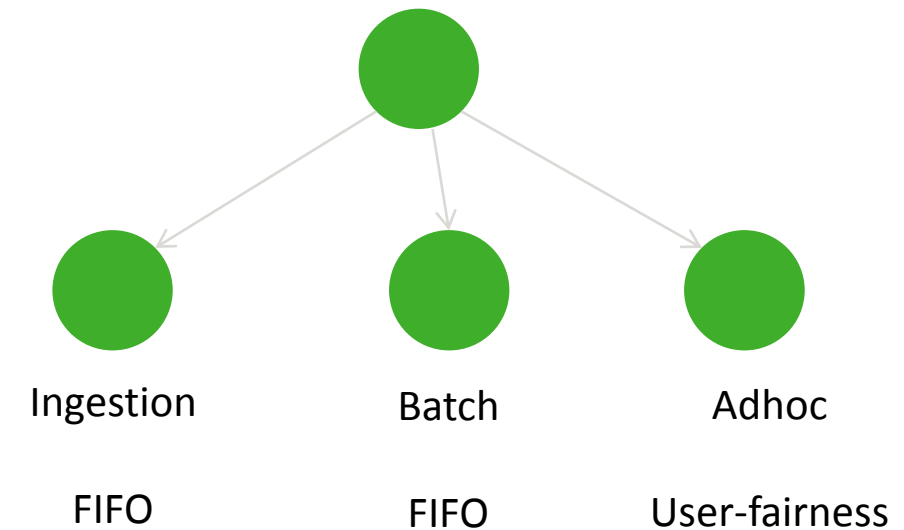- **Between apps of different users**

App 1 U1

App 1 U2

App 1 U3

Hortonworks

# Per-queue Policy-driven scheduling

**Previously**

**Now**



Ingestion     Batch     Adhoc

FIFO      FIFO      FIFO

Ingestion     Batch     Adhoc

FIFO      FIFO     User-fairness

- Coarse policies
- One scheduling algorithm in the cluster
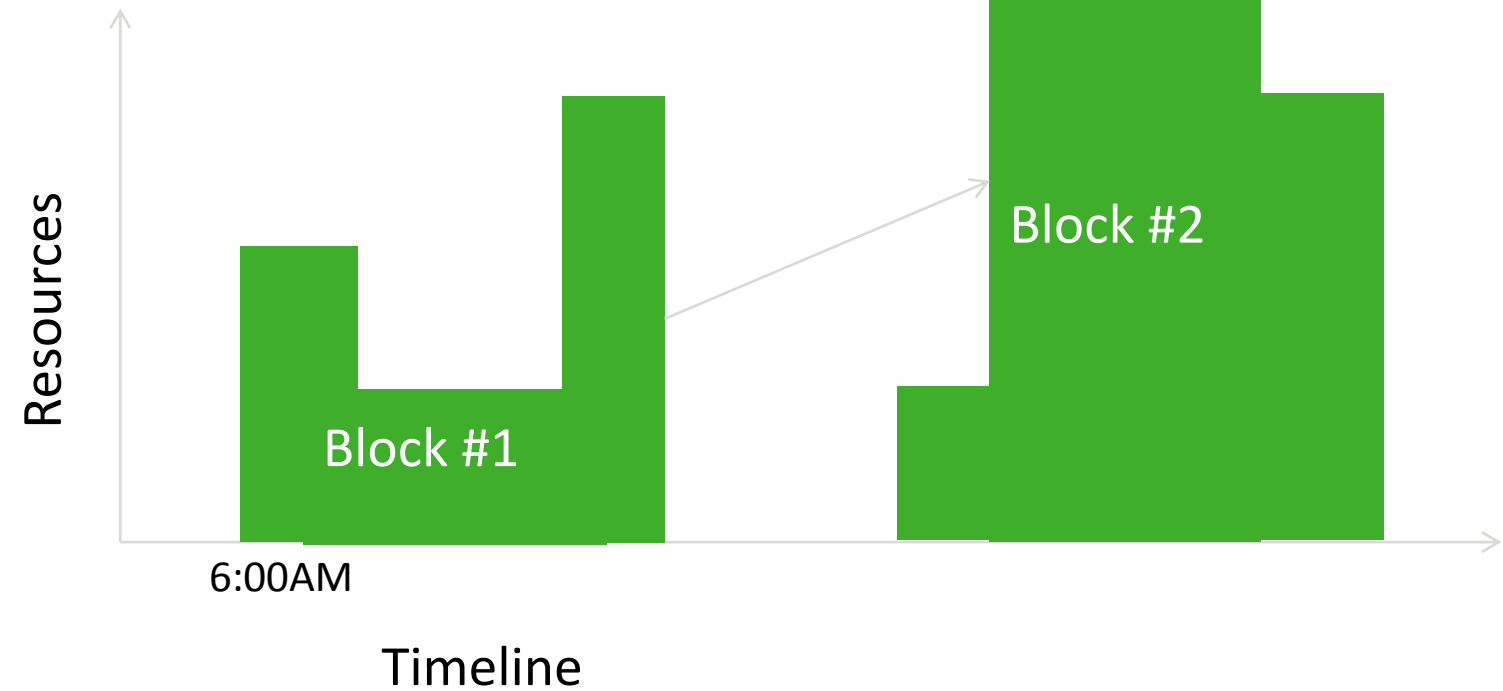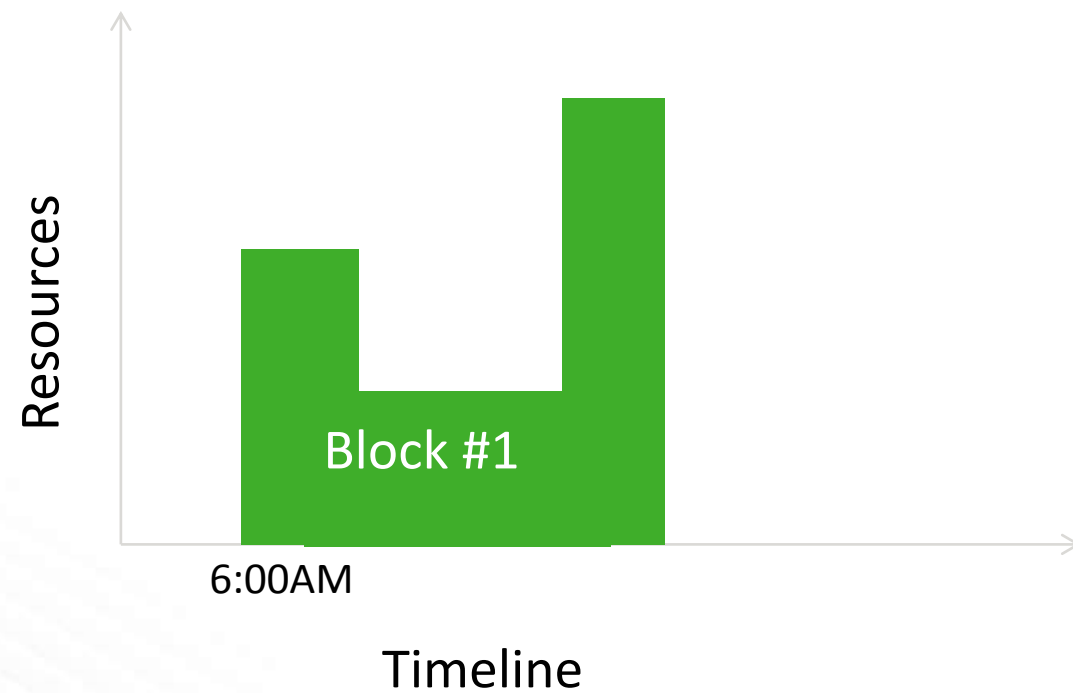- Rigid
- Difficult to experiment

- Fine grained policies
- One scheduling algorithm per queue
- Flexible
- Very easy to experiment!

**HORTONWORKS®**

# Reservations

- *"Run my workload tomorrow at 6AM"*

- **Persistence of the plans with RM failover: YARN-2573**

# Never late again! Job-Level deadline SLOs in YARN

Subru Krishnan
Microsoft

Carlo Curino
Microsoft

https://dataworkssummit.com/san-jose-2017/sessions/never-late-again-job-level-deadline-slos-in-yarn
Wednesday June 14th Room 230A

HORTONWORKS®

# Apache Hadoop 3.x

**HORTONWORKS**

# Apache Hadoop 3.0

Junping Du
Hortonworks

Andrew Wang
Cloudera

https://dataworkssummit.com/san-jose-2017/sessions/apache-hadoop-3-0-community-update
Tuesday June 13th Room 210C

HORTONWORKS®

# Scale!

- Only focusing on sizes of individual clusters

- Tons of sites with clusters made up of multiple thousands of nodes
  - Yahoo!, Twitter, LinkedIn, Microsoft

- Largest clusters the last couple of years
  - 6K-8K

- Roadmap: To 100K thousands and beyond

- Current progress: 40K nodes!

HORTONWORKS®

# Lessons learned from scaling YARN to 40k machines in a multi tenancy environment
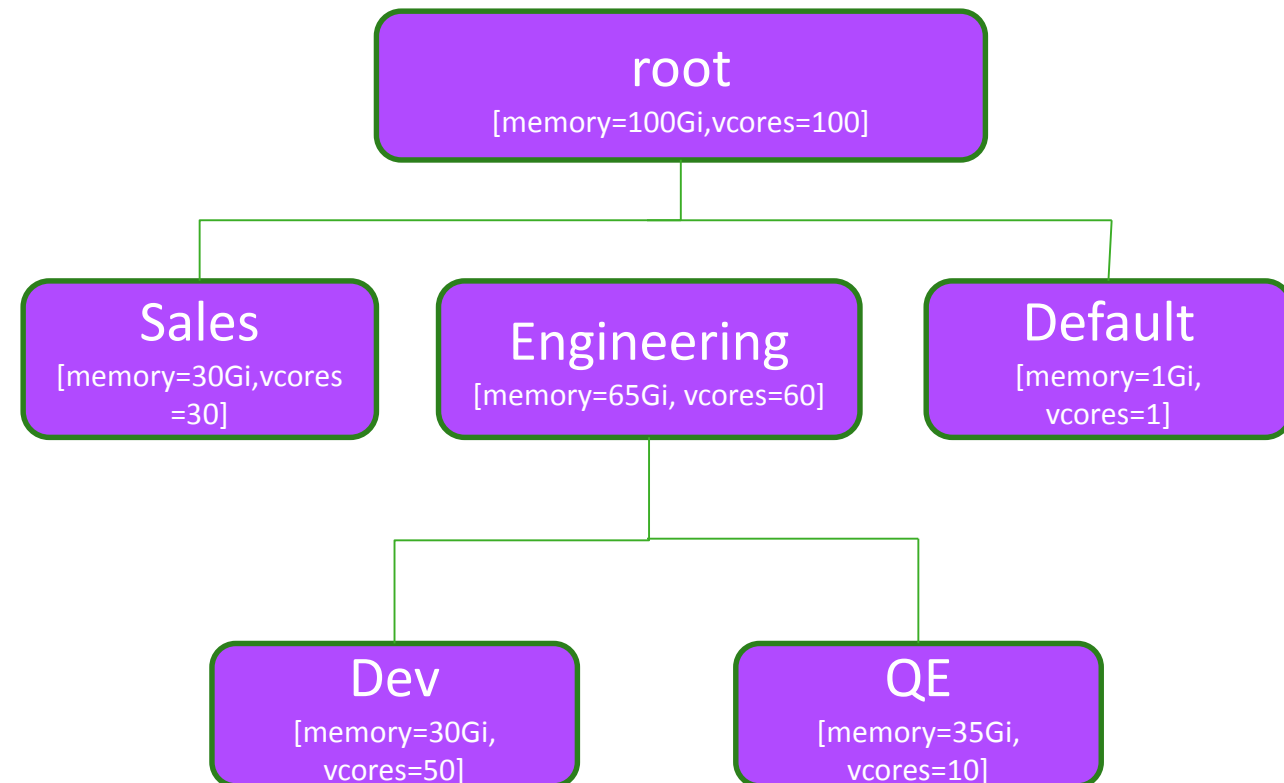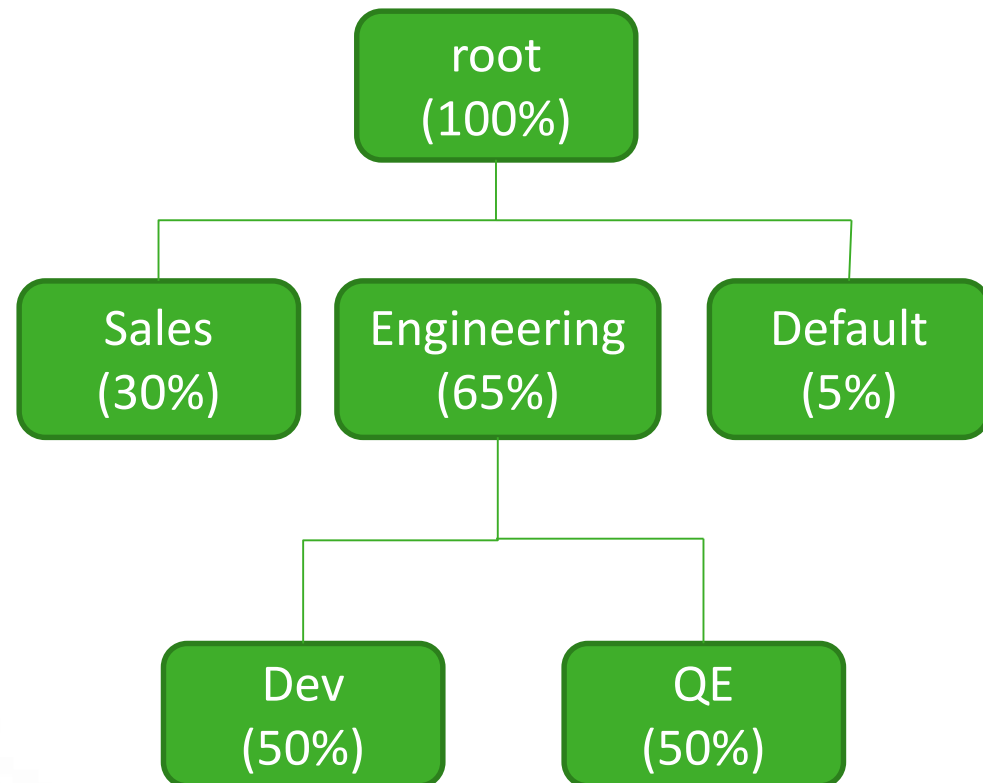
Hitesh Sharma
Microsoft

Roni Burd
Microsoft

https://dataworkssummit.com/san-jose-2017/sessions/lessons-learned-from-scaling-yarn-to-40k-machines-in-a-multi-tenancy-environment
Wednesday June 14th Room 210A
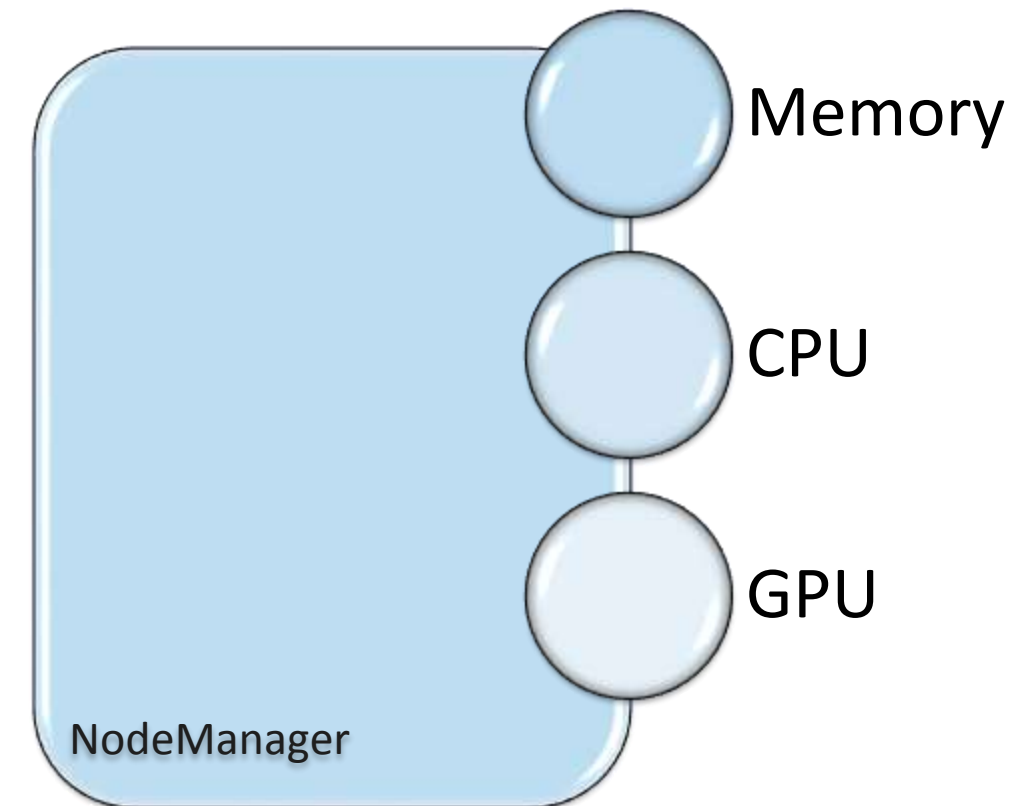
HORTONWORKS®

# Global & Fast Scheduling

- Effort led by Wangda Tan

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

- On trunk
  - 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 with enhancement added recently
  - Opportunities for better placement

HORTONWORKS

# Capacities in numbers vs percentages



root
(100%)

Sales
(30%)

Engineering
(65%)

Default
(5%)

Dev
(50%)

QE
(50%)

root
[memory=100Gi,vcores=100]

Sales
[memory=30Gi,vcores=30]

Engineering
[memory=65Gi, vcores=60]

Default
[memory=1Gi, vcores=1]

Dev
[memory=30Gi, vcores=50]

QE
[memory=35Gi, vcores=10]

HORTONWORKS®

# Resource vectors

- **Till now**
  - Hard coded resources
  - Memory and CPU

- **Now**
  - A generalized vector
  - Admins can create custom Resource Types!

Memory

CPU

GPU

NodeManager

**HORTONWORKS**®

# GPUs on a YARN cluster!

- GPU can speed up compute-intensive applications 10x - 300x times

- Different levels of support
  - Take me to a machine where GPUs are available with Partitions / Node Labels
  - Take me to a machine where GPUs are available
    - give me a **full device only to me** for the **lifetime of my container**
    - give me **multiple full devices only to me** for the **lifetime of my container**
    - give me **full device(s) only to me** for a **portion of the lifetime of my container**
    - give me a **slice of device(s) to me** for a **full / portion of the lifetime of my container**

- More dimensions:
  - CPUs and memory and GPUs and on-GPU memory
  - Topology of multiple GPUs

HORTONWORKS®

# Hadoop ecosystem boosts Tensorflow and machine learning technologies
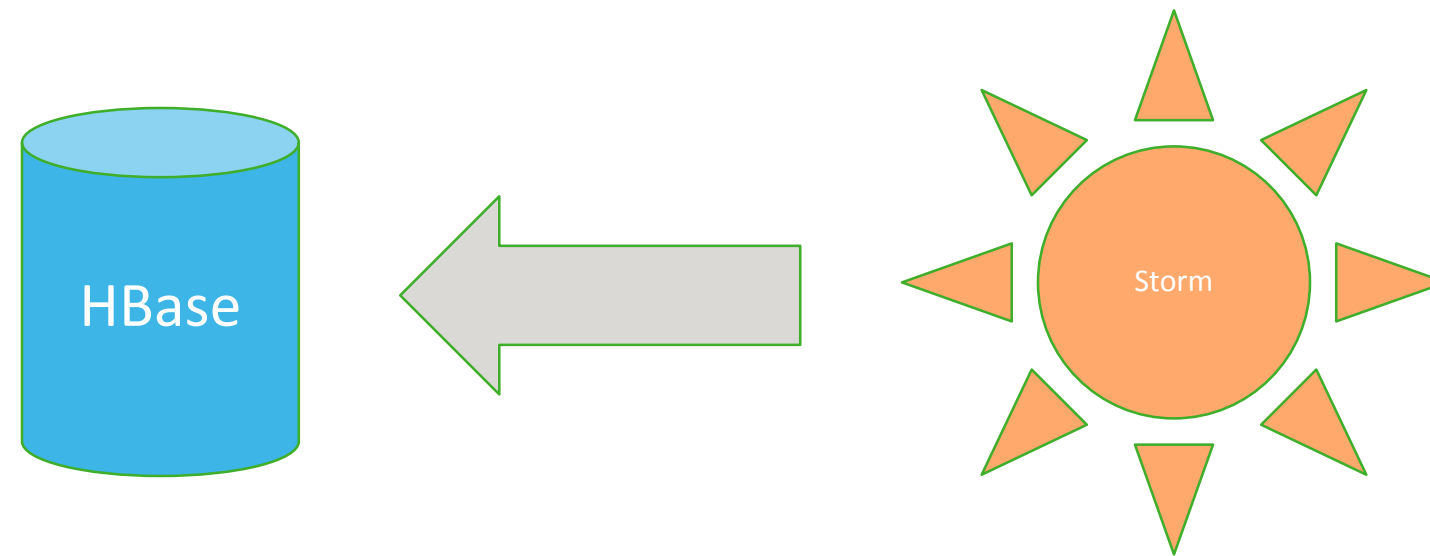
Wangda Tan
Hortonworks

Yanbo Liang
Hortonworks

https://dataworkssummit.com/san-jose-2017/sessions/hadoop-ecosystem-boosts-tensorflow-and-machine-learning-technologies

Wednesday June 14th Ballroom B

**HORTONWORKS**®

# Better placement strategies

- Affinity

- Anti-affinity
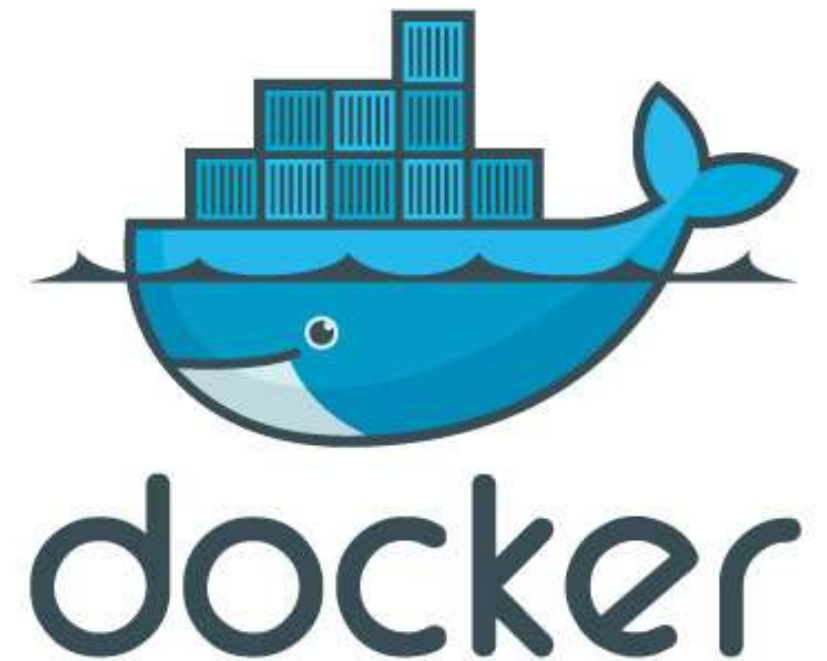
# Packaging



- Containers
  - Lightweight mechanism for packaging and resource isolation
  - Popularized and made accessible by Docker
  - Can replace VMs in some cases
  - Or more accurately, VMs got used in places where they didn't need to be
- Native integration ++ in YARN
  - Support for "Container Runtimes" in LCE: YARN-3611
  - Process runtime
  - Docker runtime

HORTONWORKS®

# Simplified APIs for service definitions

- Applications need simple APIs

- Need to be deployable "easily"

- Simple REST API layer fronting YARN
  - https://issues.apache.org/jira/browse/YARN-4793
  - [Umbrella] Simplified API layer for services and beyond

- Spawn services & Manage them

```
1  {
2      "name": "nginx",
3      "lifetime": "3600",
4      "queue": "default-developers",
5      "components" :
6      [
7          {
8              "name": "NGINX",
9              "dependencies": [ ],
10             "number_of_containers": 1,
11             "artifact": {
12                 "id": "nginx:latest",
13                 "type": "DOCKER"
14             },
15             "launch_command": "nginx -d daemon off",
16             "resource": {
17                 "cpus": 1,
18                 "memory": "1024"
19             }
20         }
21     ]
22  }
```

HORTONWORKS®

# Services support



- Application & Services upgrades
  - "Do an upgrade of my Spark / HBase apps with minimal impact to end-users"
  - YARN-4726

- Simplified discovery of services via DNS mechanisms: YARN-4757
  - *regionserver30.hbase-app-3.0.vinodkv.yarn.site*

**HORTONWORKS®**

# Services Framework

- Platform is only as good as the tools

- A native YARN services framework
  - https://issues.apache.org/jira/browse/YARN-4692
  - [Umbrella] Native YARN framework layer for services and beyond

- Assembly: Supporting a DAG of apps:
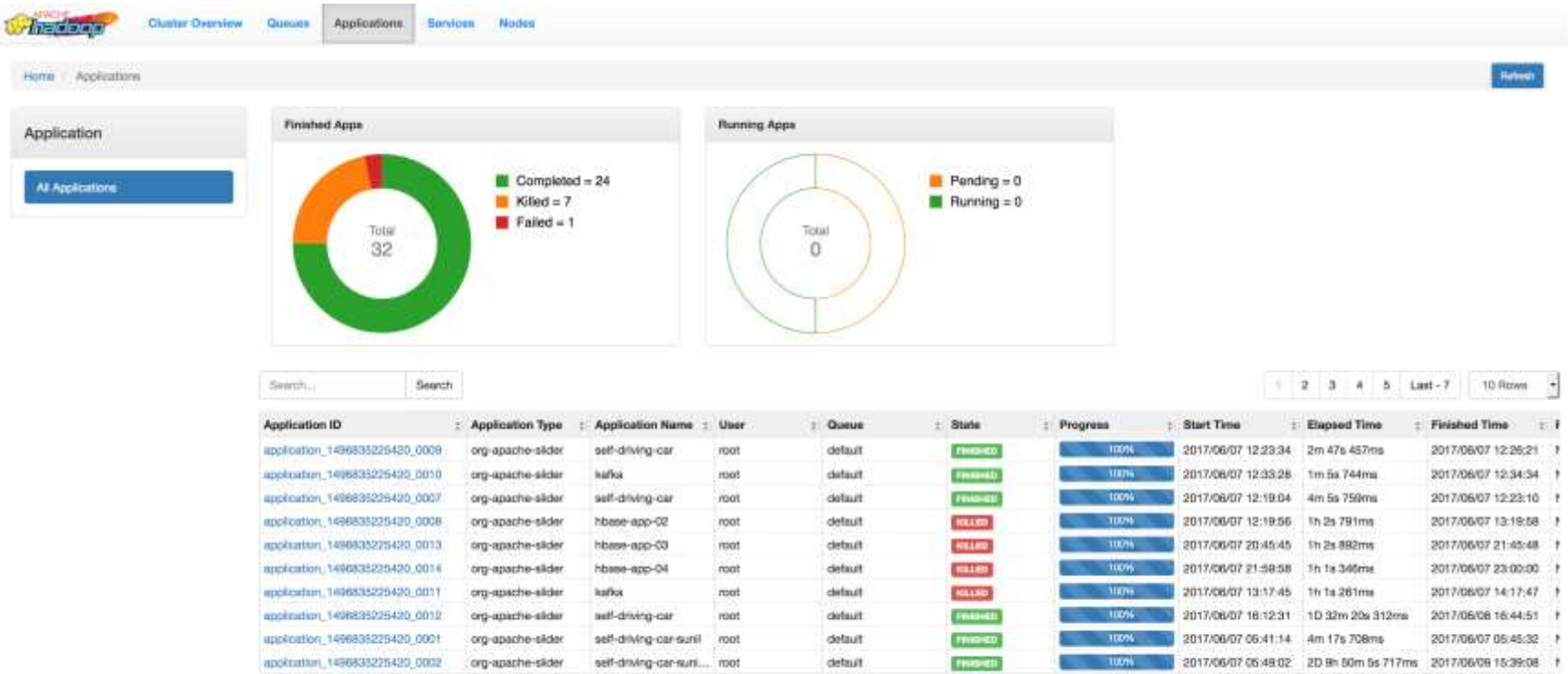  - https://issues.apache.org/jira/browse/SLIDER-875

# Running a container cloud on YARN

Shane Kumpf
Hortonworks

Jian He
Hortonworks

https://dataworkssummit.com/san-jose-2017/sessions/running-a-container-cloud-on-yarn
Thursday June 15th Room 210C

**HORTONWORKS®**

# User experience



API based queue management
Decentralized

New web UI

Improved logs management
Live application logs

HORTONWORKS®

# Timeline Service

- Application History
  - "Where did my containers run?"
  - "Why is my application slow?"
  - "Is it really slow?"
  - "Why is my application failing?"
  - "What happened with my application? Succeeded?"

- Cluster History
  - Run analytics on historical apps!
  - "User with most resource utilization"
  - "Largest application run"
  - "Why is my cluster slow?"
  - "Why is my cluster down?"
  - "What happened in my clusters?"

- Collect and use past data
  - To schedule "my application" better
  - To do better capacity planning

# Timeline Service 2.0

- ## Next generation
  - Today's solution helped us understand the space
  - Limited scalability and availability

- ## *"Analyzing Hadoop Clusters is becoming a big-data problem"*
  - Don't want to throw away the Hadoop application metadata
  - Large scale
  - Enable near real-time analysis: *"Find me the user who is hammering the FileSystem with rouge applications. Now."*

- ## Timeline data stored in HBase and accessible to queries

**HORTONWORKS®**

# Building a modern end-to-end open source Big Data reference application

Edgar Orendain

UC Berkeley / Hortonworks

https://dataworkssummit.com/san-jose-2017/sessions/building-a-modern-end-to-end-open-source-big-data-reference-application

Wednesday June 14th Ballroom C

# Thank you!