

# YAHOO!

## Leveraging Docker for Hadoop Build Automation and Big Data Stack Provisioning

Apache Big Data North America 2017

PRESENTED BY Evans Ye | May 16, 2017

# Who am I

- Software Engineer @ Y! APAC Data Team
- Building data products for...



- Apache Bigtop PMC chair



# Outline

- Quick Intro to Apache Bigtop
- Docker for Bigtop Packaging
- Docker for Bigtop Provisioner
- Docker for Bigtop Sandbox
- Release

# Quick Intro to Apache Bigtop

# Linux Distributions



CentOS



ubuntu  
linux for human beings

fedora 



debian



# Hadoop Distributions





But there're some other great Hadoop ecosystem components..

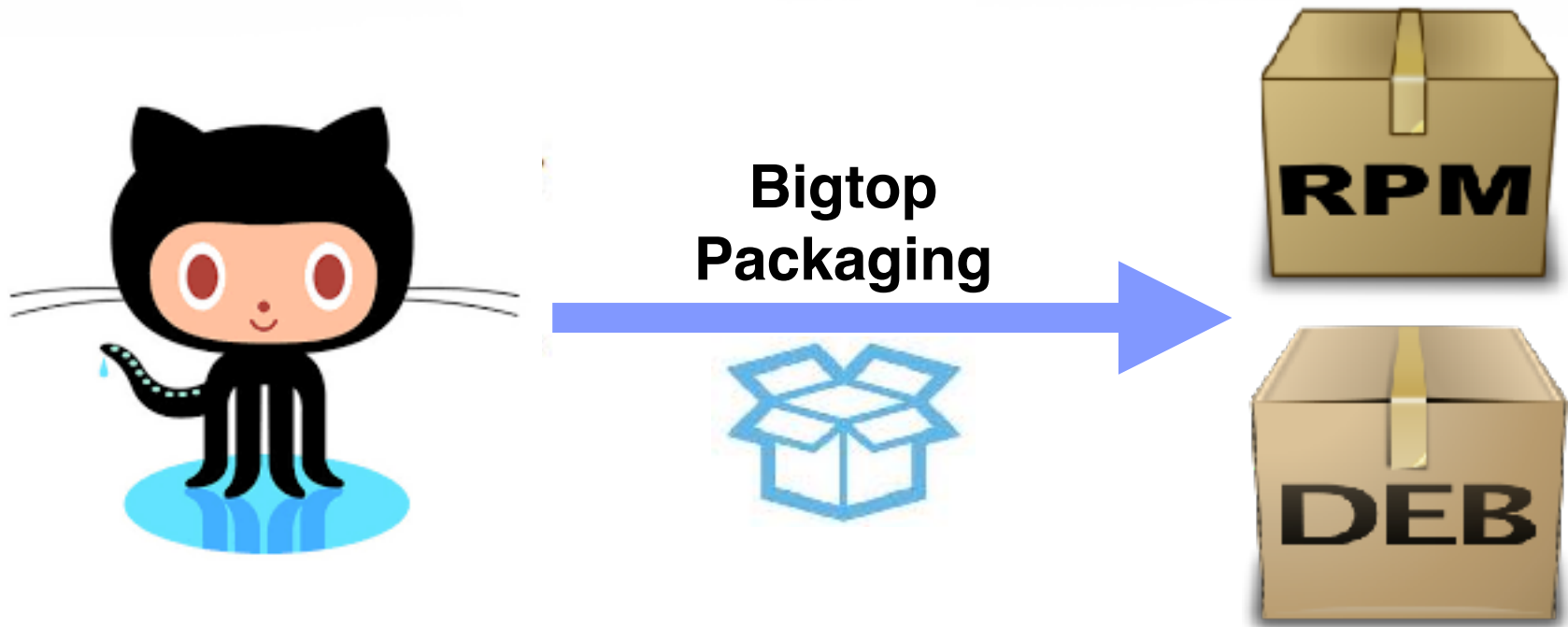


How do I add patches?





# From source code to packages



# Supported components



Spark



Apache

Solr



HUE



APACHE  
HBASE



apache

Ignite



APACHE  
PHOENIX



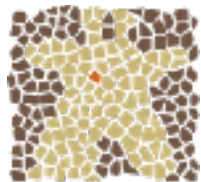
 **hadoop**

OOZIE

 **mahout**

 **HIVE**

TEZ



APACHE  
GIRAPH

  
CRUNCH



YAHOO!

# Bigtop feature set

**Packaging**



**Testing**



**Deployment**



**Virtualization**



**for you to easily build your own Big Data Stack**

# Docker for Bigtop Packaging

# Preparing build environment

## Tool requirements for building Bigtop

On all systems	Also on RPM-based systems	Also on DEB-based systems
<ul style="list-style-type: none"><li>• Java JDK 1.6</li><li>• Apache Ant</li><li>• Apache Maven</li><li>• wget</li><li>• tar</li><li>• git</li><li>• subversion</li><li>• gcc</li><li>• gcc-c++</li><li>• make</li><li>• fuse</li><li>• protobuf-compiler</li><li>• autoconf</li><li>• automake</li><li>• libtool</li><li>• sharutils</li><li>• xslt</li></ul>	<ul style="list-style-type: none"><li>• lzo-devel</li><li>• zlib-devel</li><li>• fuse-devel</li><li>• openssl-devel</li><li>• python-devel</li><li>• libxml2-devel</li><li>• libxslt-devel</li><li>• cyrus-sasl-devel</li><li>• sqlite-devel</li><li>• mysql-devel</li><li>• openldap-devel</li><li>• rpm-build</li><li>• createrepo</li><li>• redhat-rpm-config (RedHat/CentOS only)</li></ul>	<ul style="list-style-type: none"><li>• libxslt1-dev</li><li>• libkrb5-dev</li><li>• libldap2-dev</li><li>• libmysqlclient-dev</li><li>• libsasl2-dev</li><li>• libsqlite3-dev</li><li>• libxml2-dev</li><li>• python-dev</li><li>• python-setuptools</li><li>• liblzo2-dev</li><li>• libzip-dev</li><li>• libfuse-dev</li><li>• libssl-dev</li><li>• build-essential</li><li>• dh-make</li><li>• debhelper</li><li>• devscripts</li><li>• reprepro</li></ul>

# Preparing build environment

## Tool requirements for building Bigtop

On all systems	Also on RPM-based systems	Also on DEB-based systems
<ul style="list-style-type: none"><li>• Java JDK 1.6</li><li>• Apache Ant</li><li>• Apache Maven</li><li>• wget</li><li>• tar</li><li>• git</li><li>• subversion</li><li>• gcc</li><li>• gcc-c++</li><li>• make</li><li>• fuse</li><li>• protobuf-compiler</li><li>• autoconf</li><li>• automake</li><li>• libtool</li><li>• sharutils</li><li>• xslt</li></ul>	<ul style="list-style-type: none"><li>• lzo-devel</li><li>• zlib-devel</li><li>• fuse-devel</li><li>• openssl-devel</li><li>• python-devel</li><li>• libxml2-devel</li><li>• libxslt-devel</li><li>• libxslt-devel</li><li>• sqlite-devel</li><li>• mysql-devel</li><li>• openldap-devel</li><li>• rpm-build</li><li>• createrepo</li><li>• redhat-rpm-config (RedHat/CentOS only)</li></ul>	<ul style="list-style-type: none"><li>• libxslt1-dev</li><li>• libkrb5-dev</li><li>• libldap2-dev</li><li>• libmysqlclient-dev</li><li>• libssl2-dev</li><li>• libsqlite3-dev</li><li>• libxml2-dev</li><li>• python-dev</li><li>• python-setuptools</li><li>• liblzo2-dev</li><li>• libzip-dev</li><li>• libfuse-dev</li><li>• libssl-dev</li><li>• build-essential</li><li>• dh-make</li><li>• debhelper</li><li>• devscripts</li><li>• reprepro</li></ul>

...  
**Seriously ?**

# Bigtop Toolchain

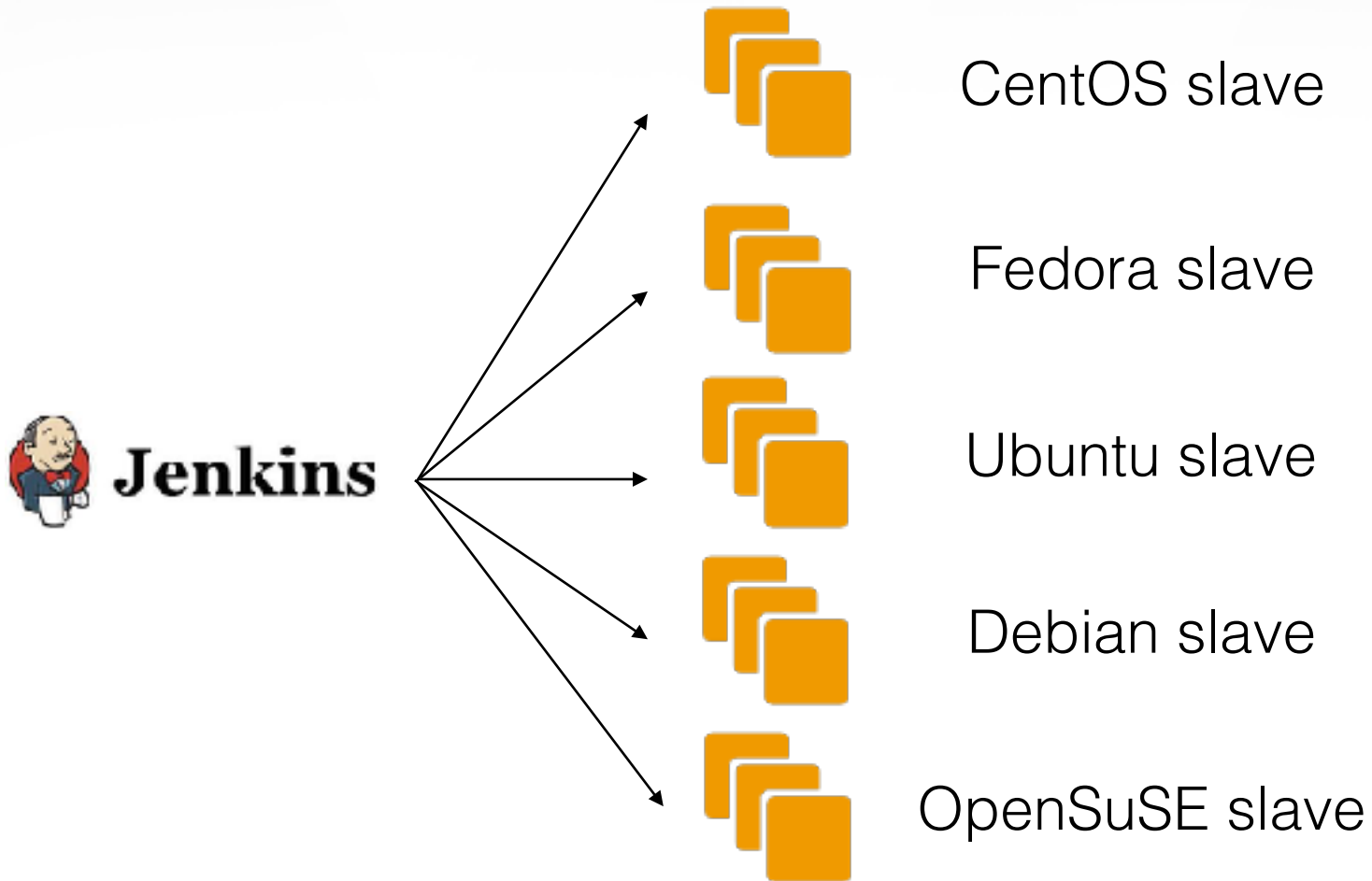
- Puppet recipes to install required libraries, build tools
- To prepare a build environment:

```
git clone https://github.com/apache/bigtop.git
cd bigtop
./bigtop_toolchain/bin/puppetize.sh
./gradlew toolchain
```

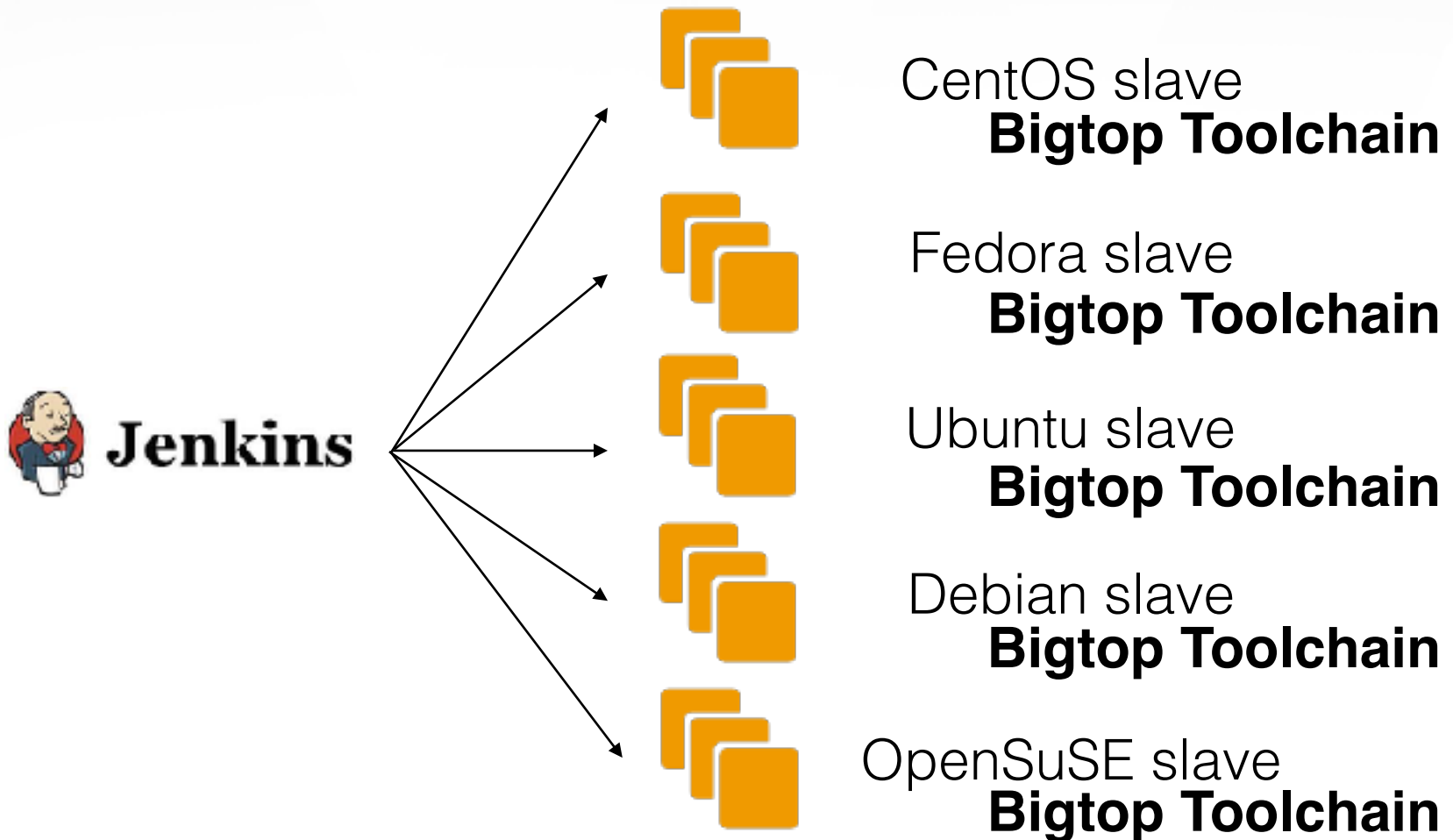
- Prerequisite :
  - Java



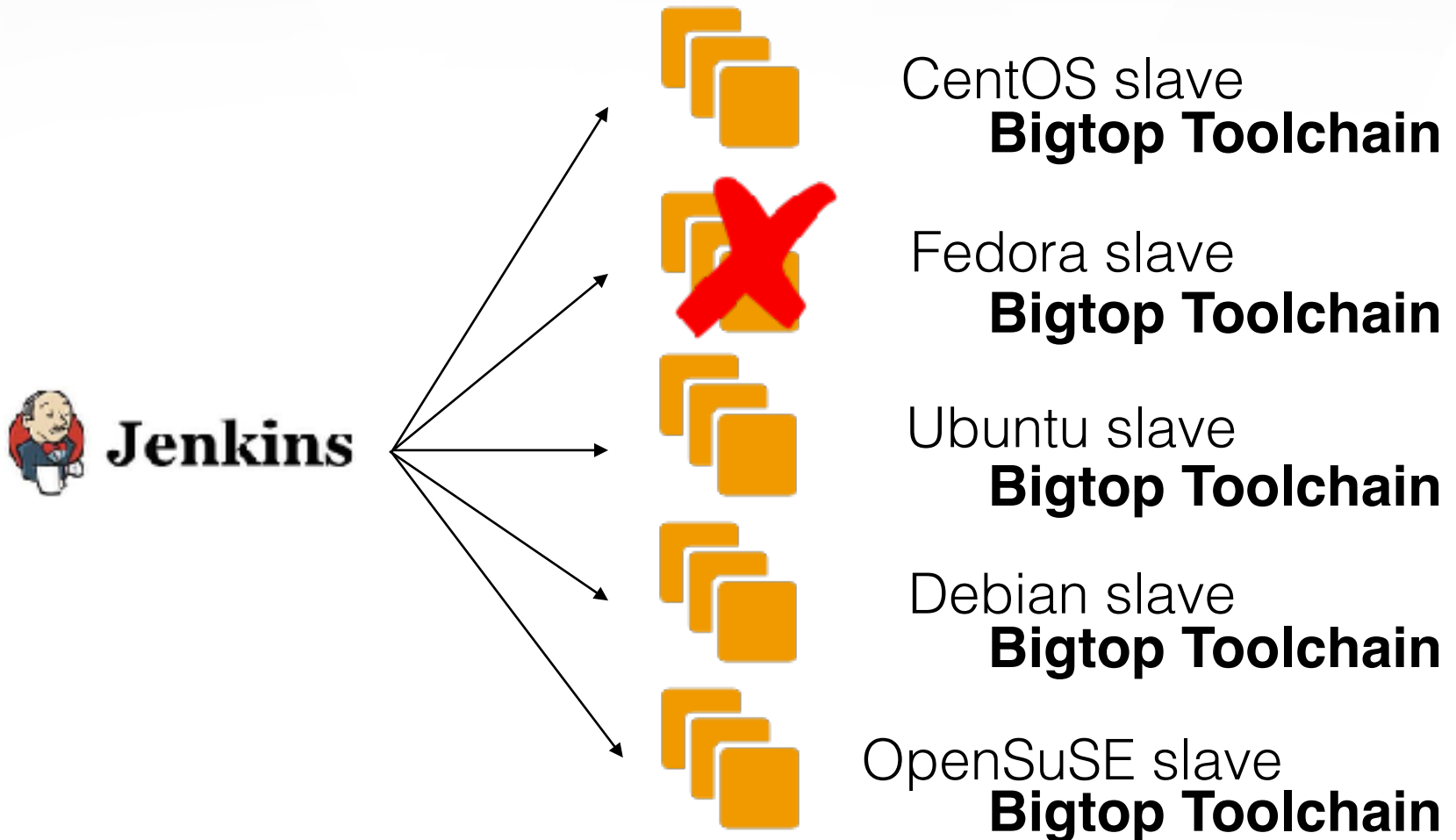
# CI Infrastructure



# CI Infrastructure



# CI Infrastructure

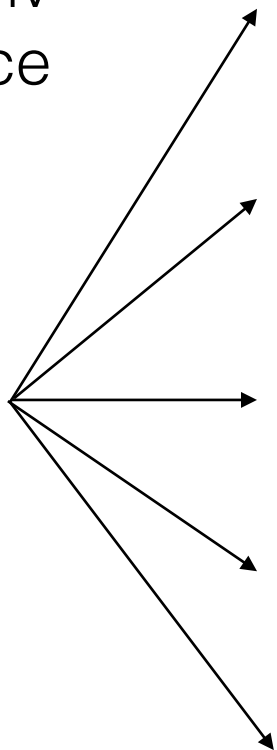


# Dockerized CI Infrastructure

- Immutable env
- Fault tolerance



**Jenkins**



CentOS slave



Fedora slave



Ubuntu slave



Debian slave



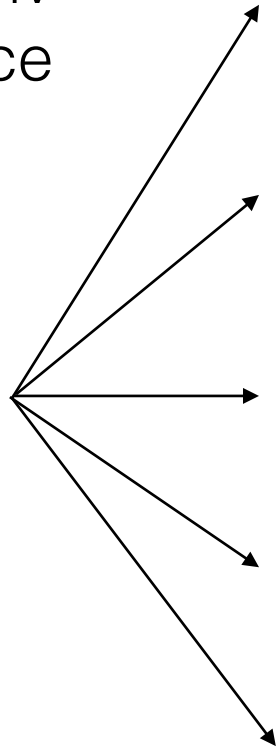
OpenSuSE slave

# Dockerized CI Infrastructure

- Immutable env
- Fault tolerance



**Jenkins**



CentOS slave



Fedora slave  
Ubuntu slave



Debian slave



OpenSuSE slave

# How to build packages

- Execute shell

```
# OS=debian-8  
# COMPONENT=hadoop
```

```
docker run -u jenkins --rm \  
-v `pwd`:/bigtop --workdir /bigtop \  
bigtop/slaves:trunk-$OS \  
bash -l -c "./gradlew allclean $COMPONENT-pkg"
```

- [Bigtop CI Setup Guide](#)

Configuration Matrix	centos-6	centos-7	debian-8	fedora-25	Fedora-25-ppc64le	opensuse-42.1	ubuntu-16.04	ubuntu-16.04-ppc64le
apex	●	●	●	●	●	●	●	●
ambari	●	●	●	●	●	●	●	●
alluxio	●	●	●	●	●	●	●	●
bigtop-groovy	●	●	●	●	●	●	●	●
bigtop-jsvc	●	●	●	●	●	●	●	●
bigtop-tomcat	●	●	●	●	●	●	●	●
bigtop-utils	●	●	●	●	●	●	●	●
crunch	●	●	●	●	●	●	●	●
datafu	●	●	●	●	●	●	●	●
flume	●	●	●	●	●	●	●	●
flink	●	●	●	●	●	●	●	●
giraph	●	●	●	●	●	●	●	●
gpdb	●	●	●	●	●	●	●	●
hadoop	●	●	●	●	●	●	●	●
hama	●	●	●	●	●	●	●	●
hbase	●	●	●	●	●	●	●	●
hive	●	●	●	●	●	●	●	●
hue	●	●	●	●	●	●	●	●
ignite-hadoop	●	●	●	●	●	●	●	●
kafka	●	●	●	●	●	●	●	●
kite	●	●	●	●	●	●	●	●
mahout	●	●	●	●	●	●	●	●
oozie	●	●	●	●	●	●	●	●
phoenix	●	●	●	●	●	●	●	●
pig	●	●	●	●	●	●	●	●
qfs	●	●	●	●	●	●	●	●
solr	●	●	●	●	●	●	●	●
spark1	●	●	●	●	●	●	●	●
spark	●	●	●	●	●	●	●	●
sqoop	●	●	●	●	●	●	●	●
sqoop2	●	●	●	●	●	●	●	●
tajo	●	●	●	●	●	●	●	●
tez	●	●	●	●	●	●	●	●
yesb	●	●	●	●	●	●	●	●
zookeeper	●	●	●	●	●	●	●	●
zeppelin	●	●	●	●	●	●	●	●

# Bigtop master

# Bigtop early mission accomplished



All major Hadoop distros leverage Bigtop to build its foundation

Leveraged by app providers...





# Get out from the Apache dome



# New focus and target user

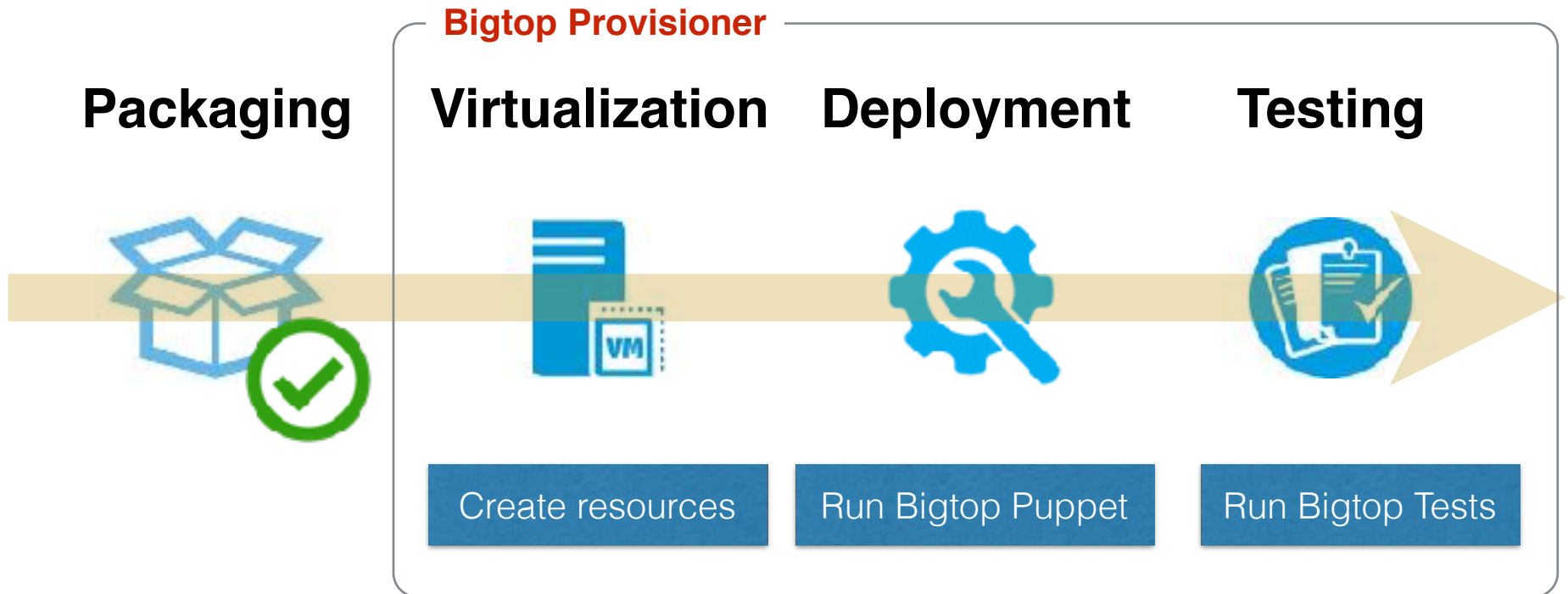
- Data engineers vs Distro. builders
- Solution diversity:
  - Streaming: Flink, Apex
  - In-memory cache: Alluxio, Ignite
  - Non apache: QFS, GPDB
- User/developer tools:
  - Bigtop Provisioner
  - Bigtop Sandbox
- Big data stack references



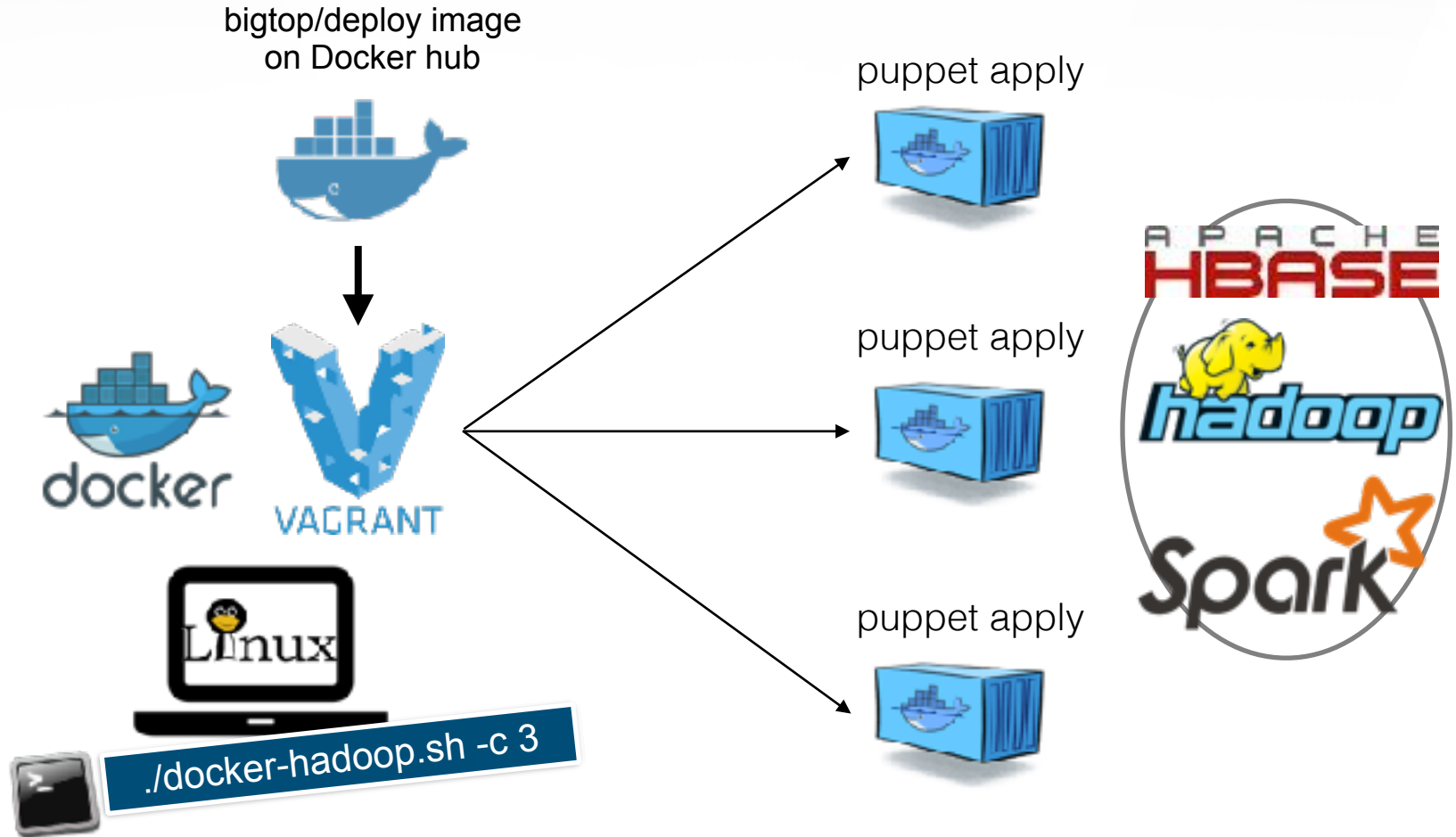
# Docker for Bigtop Provisioner

# Bigtop Provisioner

- A tool to demonstrate full life cycle of Bigtop



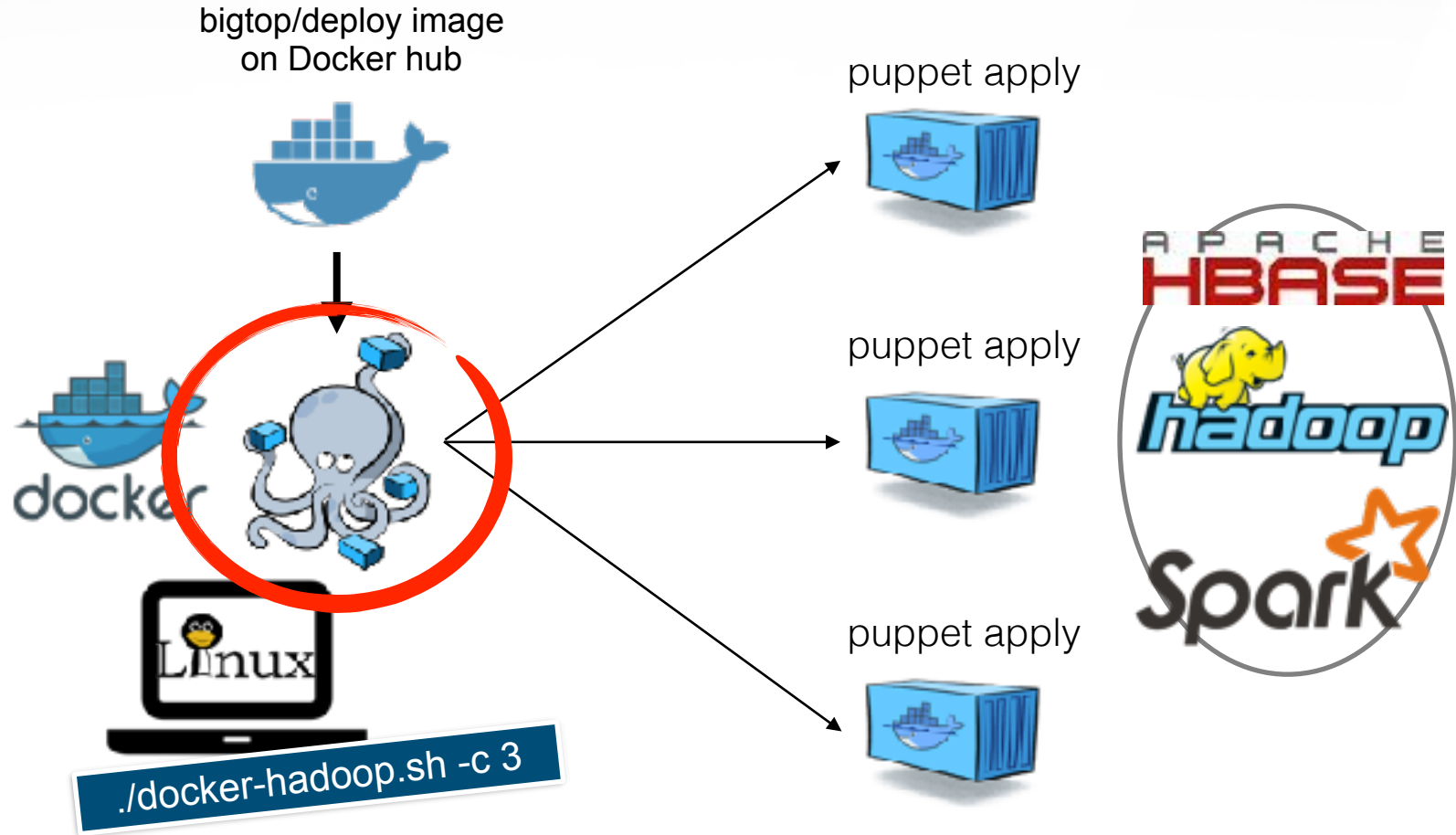
# One click Hadoop provisioning (Bigtop 1.0.0)



# What's the problem with Vagrant's Docker Provider?

- Need to add vagrant public key into docker images
- Too many issues with auto-created boot2docker VM
- A bug for docker provider keep opening for 2ys
  - Waiting for machine to boot' hangs infinitely
- Can not share same code for different providers anyway
- Not all the docker options supported in Vagrantfile
- ^#?& slow

# Replaced by docker-compose (Bigtop 1.2.0)



# Advantages

- No need to create customized image beforehand
- Better compatibility with Docker's native solutions
- Clear, simple yaml file for orchestration settings
- Supports new features such as overlay network
- Leverage Swarm for multi-node cluster deployment
- Fast —> better user experience



# How to run Docker Provisioner

- Execute shell

```
# See bigtop/provisioner/docker/*.yaml  
CONFIG=YOUR_CUSTOM_CONF.yaml
```

```
# provision  
./gradlew -Pconfig=${CONFIG} -Pnum_instances=1 \  
docker-provisioner
```

```
# destroy provisioned cluster  
./gradlew docker-provisioner-destroy
```

- [Bigtop CI Setup Guide](#)

Configuration Matrix	centos6	centos7	debian8	ubuntu_xenial
alluxio	●	●	●	●
apex	●	●	●	●
crunch	●	●	●	●
flink	●	●	●	●
flume	●	●	●	●
hbase	●	●	●	●
giraph	●	●	●	●
hive	●	●	●	●
httpfs	●	●	●	●
hue	●	●	●	●
ignite_hadoop	●	●	●	●
karbon	●	●	●	●
mahout	●	●	●	●
oozie	●	●	●	●
pig	●	●	●	●
qfs	●	●	●	●
spark	●	●	●	●
sqoop2	●	●	●	●
tez	●	●	●	●
yarn	●	●	●	●
ycsb	●	●	●	●
zeppelin	●	●	●	●
zookeeper	●	●	●	●

# Visibility for deployments

# Use Cases

- **For application developers, cluster admins, users**
  - Run a Hadoop cluster to test your code on
  - Try & test configurations before applying to Production
  - Play around with Bigtop Big Data Stacks
- **For contributors**
  - Easy to test your packaging, deployment, testing code
- **For Distro. builders**
  - CI matrix → patch upstream code made easier

# Docker for Bigtop Sandbox

# Introducing Bigtop Sandbox

- Easiest way to get started
- Docker images that has Bigtop stacks installed and configured
- Pseudo cluster up & running w/ zero installation
- Command-line tool for you to build your own stack

# Docker Image layer Interface

Customized big data stack

Deploy & management tool

Base image (OS)

# Docker Image layer

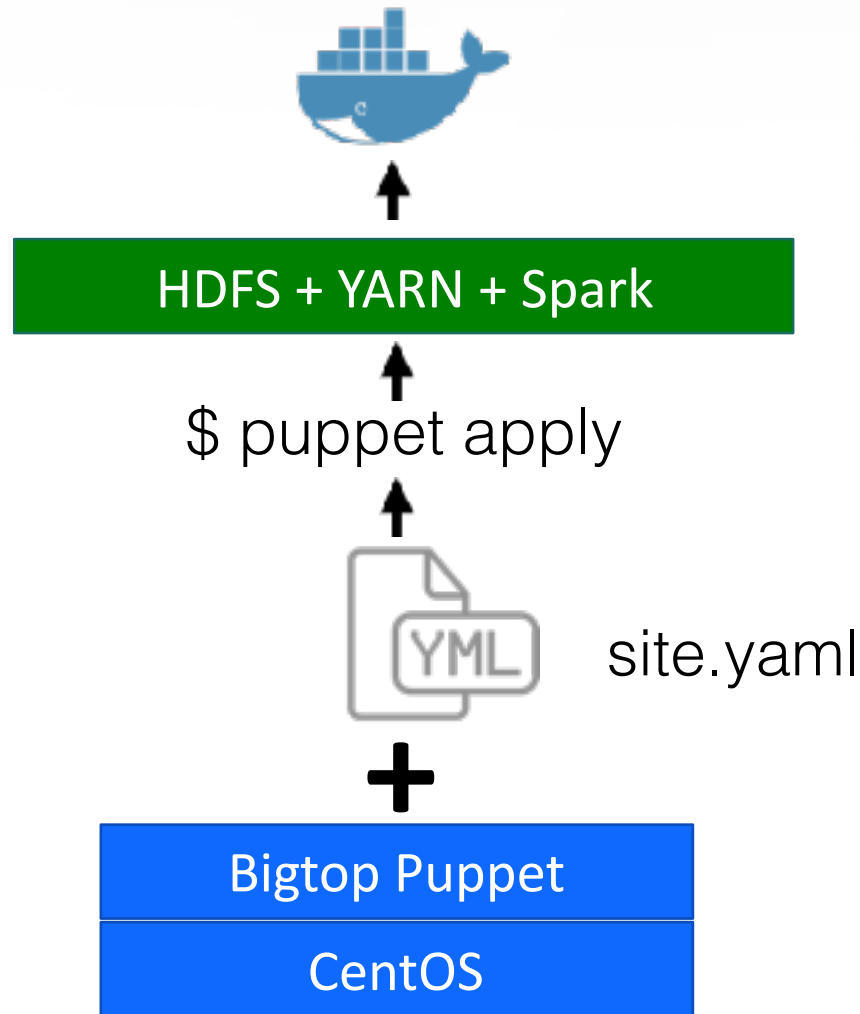
## Concrete implementation

HDFS + YARN + Spark

Bigtop Puppet

`bigtop/puppet:ubuntu-16.04`

# Building images





# How to build

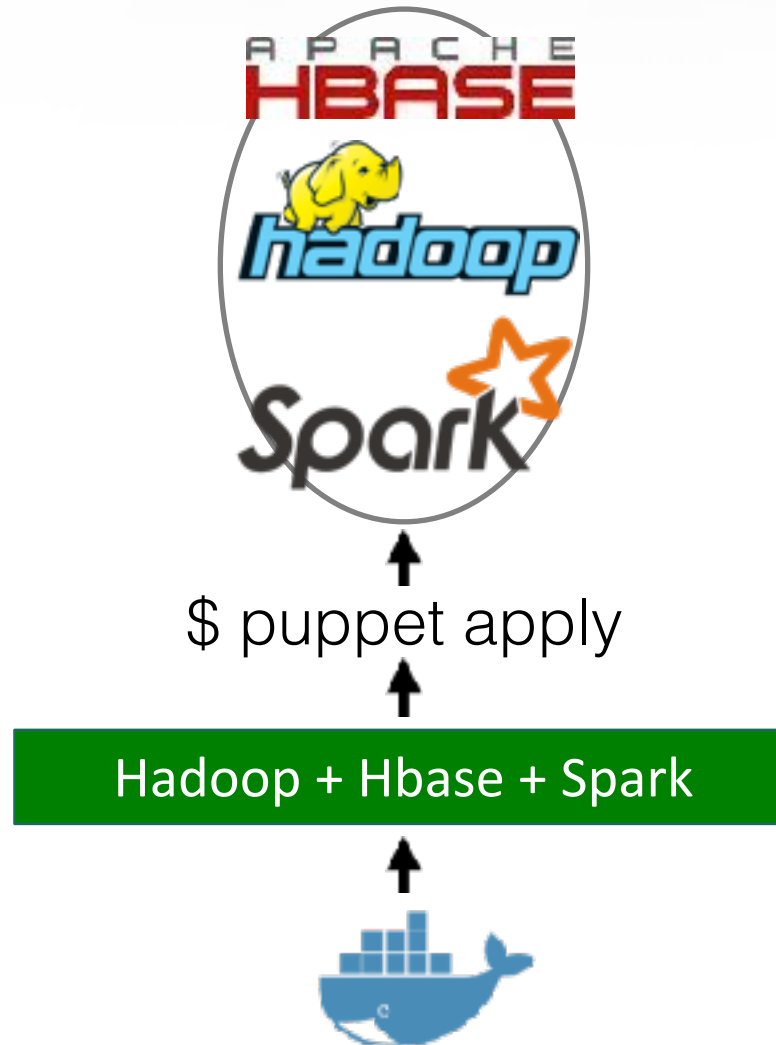
```
git clone https://github.com/apache/bigtop.git  
cd bigtop/docker/sandbox
```

```
./build.sh -a evansye -o ubuntu-16.04 \  
-c hdfs,yarn,spark
```

- Specify custom conf:

```
./build.sh -a evansye -o ubuntu-16.04 \  
-f site.yaml -t apache_big_data_2017_miami
```

# Running images



# How to run

```
docker run --name sandbox -d \  
-p 50070:50070 -p 8088:8088 \  
bigtop/sandbox:apache_big_data_2017_miami
```

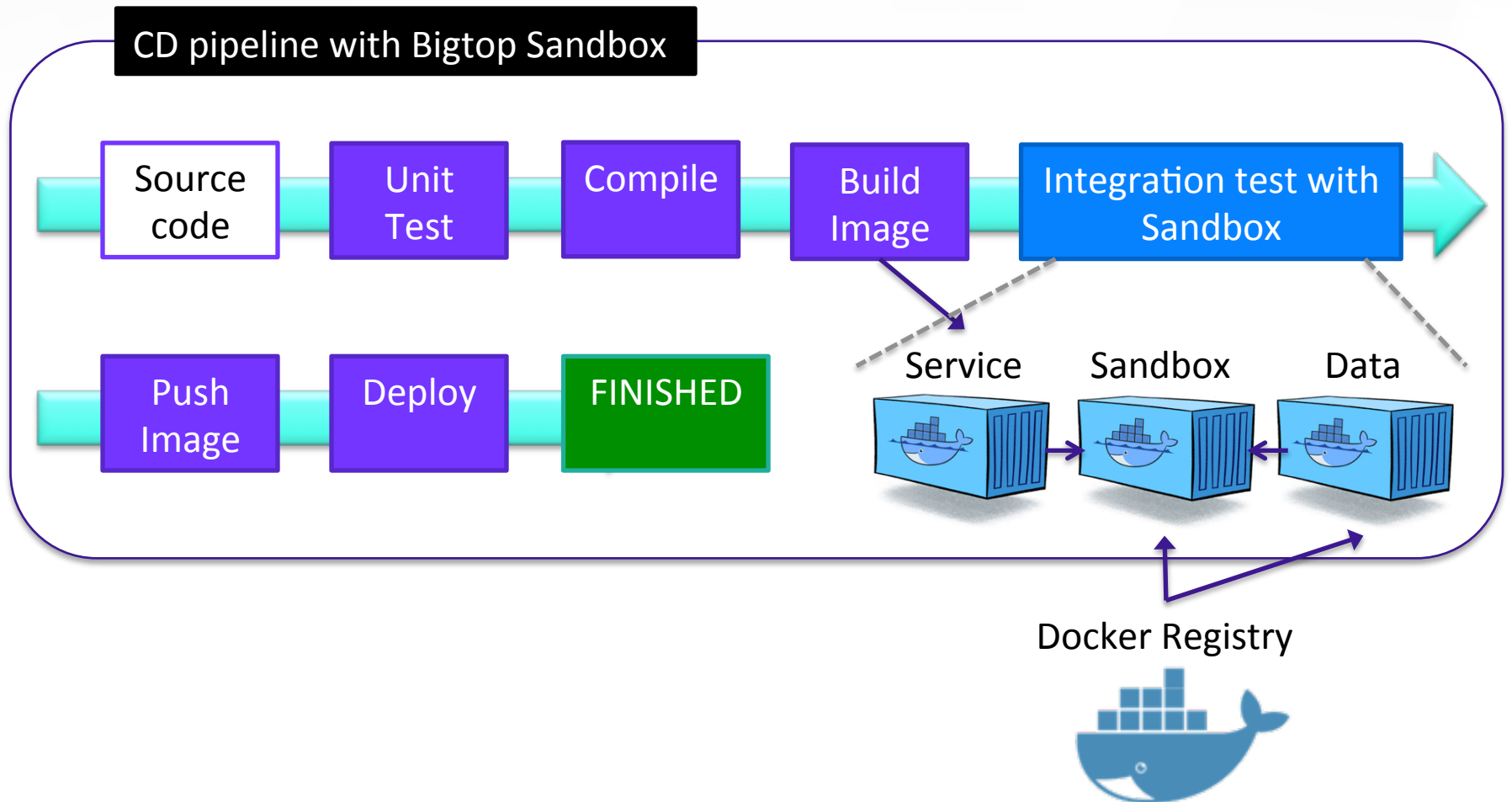
```
docker logs -f sandbox
```

```
docker exec sandbox spark-example SparkPi
```

	Bigtop Provisioner	Bigtop Sandbox
Scalable	V	X
Portable	X	V
Flexibility	High	Medium
Speed	> 2 mins	> 15 secs
Requires Network	V	X

	<b>Bigtop Provisioner</b>	<b>Bigtop Sandbox</b>
<b>Data engineers</b>	Multi-node cluster testing	Build/use sandboxes for dev & test
<b>Ops</b>	Multi-node cluster testing	Single node testing
<b>Contributors</b>	Test packages, puppet recipes, test cases	Test packages, puppet recipes, test cases
<b>Distro. Builders</b>	Test packages, puppet recipes, test cases	Provide Sandboxes

# Integration test in CI/CD pipeline



# Future

- Production deployment using Sandbox image
  - --net host or SDN
  - External volumes for fsimage, data, logs, etc
  - Cluster orchestration
    - Kubernetes?

**Release**



# Bigtop 1.2.0 Released Apr., 2017

- New components:
  - Ambari 2.5.0
  - GPDB 5.0.0-alpha.0  
(Greenplum)
- Featured upgrade:
  - Hadoop 2.7.3
  - Spark 2.1.0
  - Kafka 0.10.1.1
  - HBase 1.1.3
  - and more

# What's new in Bigtop 1.2.0?

- New features:
  - Juju bigtop charms
  - Bigtop Sandbox (alpha)
- Improvement:
  - Bigtop Docker Provisioner made faster

# Juju Cloud Weather Report

Bundle	Date	Latest Test Results			
cwr_bundle_bigdata_dev_hadoop_processing	Apr 13, 2017 at 03:01	AWS ✔	Azure ○	GCE ✔	Icd ▲
cwr_bundle_bigdata_dev_spark_processing	Apr 13, 2017 at 06:41	AWS ✔	Azure ○	GCE ▲	Icd ✘
cwr_bundle_hadoop_nbase	May 12, 2017 at 20:17	AWS ✔	Azure ✔	GCE ▲	Icd ○
cwr_bundle_hadoop_kafka	May 12, 2017 at 22:57	AWS ✔	Azure ✔	GCE ▲	Icd ○
cwr_bundle_hadoop_processing	May 12, 2017 at 20:37	AWS ✔	Azure ✔	GCE ▲	Icd ○
cwr_bundle_hadoop_spark	May 12, 2017 at 17:04	AWS ✔	Azure ✔	GCE ✔	Icd ○
cwr_bundle_spark_processing	May 12, 2017 at 23:26	AWS ✔	Azure ✔	GCE ▲	Icd ○
cwr_charm_release_ubuntu_devenv_in_cs_kwmanroe_java_devenv	May 10, 2017 at 20:05	AWS ✔	Azure ○	GCE ○	Icd ○

✔ Test Passed ✘ Test Failed ▲ Infrastructure Failure ○ No Test Result

# Road ahead

- AARCH 64 support
- Enhance support set in Bigtop Puppet
- Extend the CI matrix to Bigtop Tests
- Ambari Bigtop integration
- Big data stack references

# We want you!

- Join mailing list, ask questions, suggest features, etc
- Contribute (components, tutorials, docs)
- Report bugs
- Reference
  - Home page: <http://bigtop.apache.org/>
  - mailing list: <http://bigtop.apache.org/mail-lists.html>
  - Document: <https://cwiki.apache.org/confluence/display/BIGTOP/Index>
  - Source code: <https://github.com/apache/bigtop>
  - Packages: <https://www.apache.org/dist/bigtop/bigtop-1.2.0/repos/>
  - JIRA: <https://issues.apache.org/jira/browse/BIGTOP>



Thank you !



Questions?