

How Bigtop Leveraged Docker for Build Automation and One-Click Hadoop Provisioning

Evans Ye

Apache Big Data 2015
Budapest

APACHE:
BIG_DATA
EUROPE



Securing Your Journey
to the Cloud

Who am I

- Apache Bigtop PMC member
- Software Engineer at Trend Micro
- Develop Big Data platform
- Develop cyber security solutions using Big Data

Outline

- What is Apache Bigtop?
- Achieving Build Automation
- One-Click Hadoop Provisioning
- Bigtop at Trend Micro

What is Apache Bigtop ?



Bigtop is a project for

Packaging

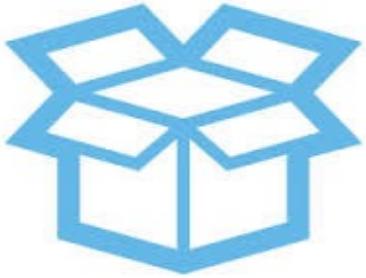
Testing

Deployment

Virtualization



for you to easily build your own Big Data Stack



Packaging

- **Build RPM, DEB packages for Hadoop ecosystem**
 - \$ yum -y install hadoop
 - \$ apt-get upgrade hbase
- **Why not just untar ?**
 - Version control, upgrade/downgrade, dependency management
 - Unified view of configuration, log, binary directories
 - Daemons for better process management

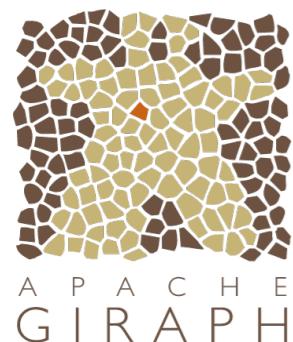
Supported Linux distro.

Bigtop supports all major Linux distributions



- You can find binary repos here:
 - <http://www.apache.org/dist/bigtop/bigtop-1.0.0/repos/>

Supported components





Testing

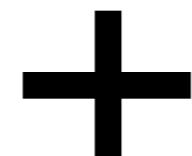
- A fundamental testing problem may happen in your Big Data Stack, too

Hadoop 2.4.1

Tez 0.6.0

Hive 1.1.0

Hadoop 2.6.0



Tez 0.6.2

Hive 1.2.0

Hadoop 2.7.1

Tez 0.7.0

Hive 1.2.1



Testing

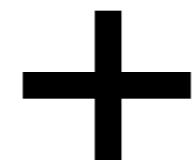
- A fundamental testing problem may happen in your Big Data Stack, too

Hadoop 2.4.1

Tez 0.6.0

Hive 1.1.0

Hadoop 2.6.0



Tez 0.6.2

Hive 1.2.0

Hadoop 2.7.1

Tez 0.7.0

Hive 1.2.1

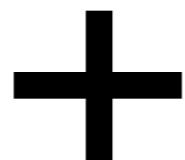


Testing

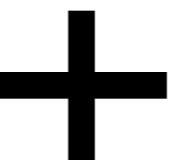
- A fundamental testing problem may happen in your Big Data Stack, too

Does this combination still work?

Hadoop 2.6.0



Tez 0.6.2



Hive 1.2.0

Hadoop 2.7.1

Tez 0.7.0

Hive 1.2.1

There's a need for
integration tests !

Bigtop Tests

- Bigtop has a testing framework that provides APIs
- Built-in tests that Bigtop provides:
 - **smoke-tests**
 - Basic Hadoop ecosystem interoperability
 - **integration tests**
 - Advanced tests runs from jar files

Okay,
then how to run those tests?

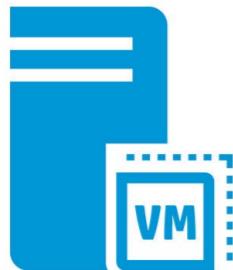


Deployment

- Use Bigtop Puppet to deploy a fully functional, distributed Hadoop cluster
- Integration tests running on a fully distributed cluster
- **Bigtop Puppet**
 - Masterless Puppet
 - Numbers of components are supported
 - Kerberos enabled cluster deployment

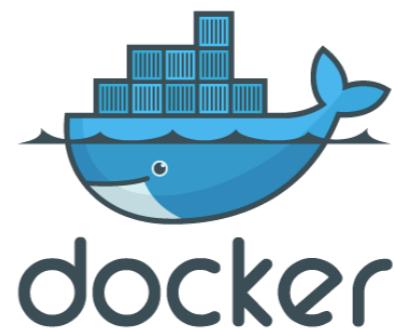


Where to deploy ?



Virtualization

- Automate the infrastructure creation & setup
- **Bigtop Provisioner**
 - Virtualbox VMs auto-provisioning
 - Docker containers auto-provisioning



From user point of view

- **Hadoop App developers**
 - Provision Hadoop cluster to test your code on
- **Cluster administrators**
 - Use Bigtop Puppet to deploy and manage your cluster
 - Run Bigtop tests to ensure your cluster is working
- **Vendors**
 - Build your own Hadoop distribution, customized from Bigtop bits

What is Apache Bigtop?

Achieving Build Automation

One-Click Hadoop Provisioning

Bigtop at Trend Micro

What's the challenge ?

- As you can see,
we support a number of Linux distributions (**m**)

Bigtop supports all major Linux distributions



What's the challenge ?

- Numbers of components (**n**)



We need to build them all
→ $m * n$

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

Preparing build environment

Tool requirements for building Bigtop

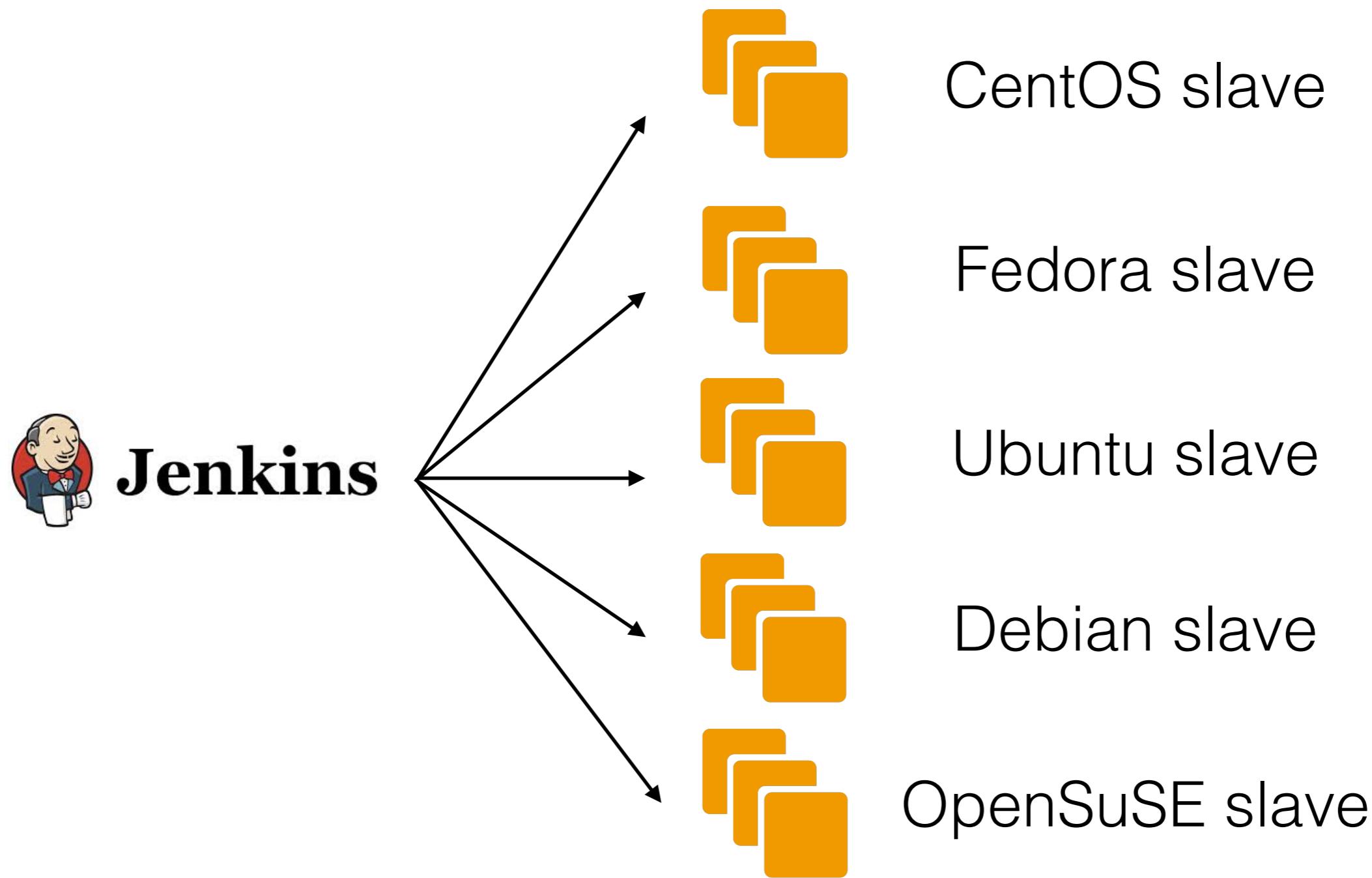
On all systems	Also on RPM-based systems	Also on DEB-based systems
<ul style="list-style-type: none">• Java JDK 1.6• Apache Ant• Apache Maven• wget• tar• git• subversion• gcc• gcc-c++• make• fuse• protobuf-compiler• autoconf• automake• libtool• sharutils• xmllt	<ul style="list-style-type: none">• lzo-devel•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

Seriously ?

Bigtop Toolchain

- Puppet recipes automatically install required libraries, build tools
- To transform a machine into a build environment, simply do
 - \$ gradle toolchain
- **Prerequisites :**
 - Puppet 3.x
 - Java 7

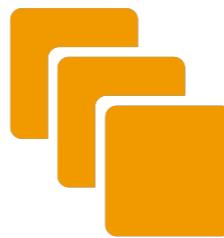
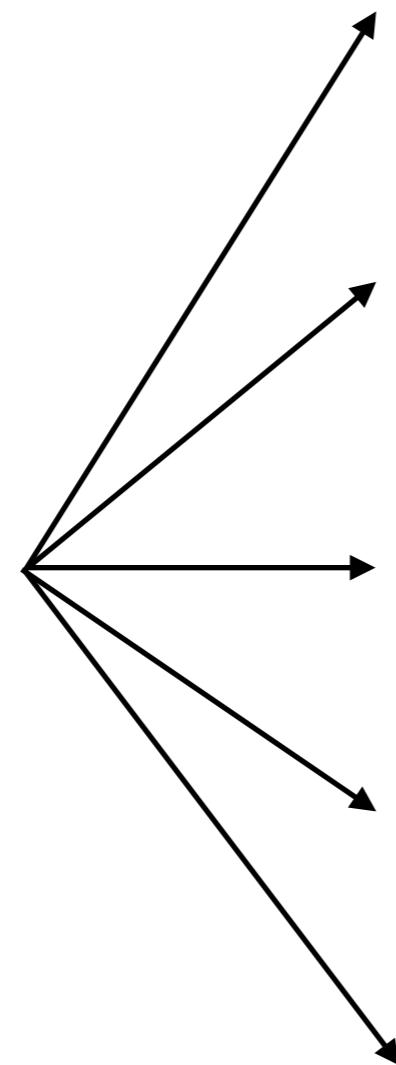
CI infra



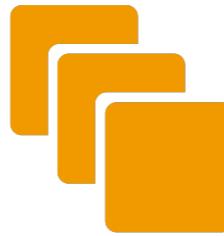
CI infra



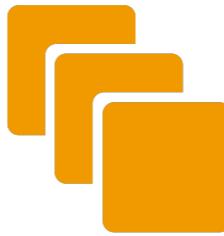
Jenkins



CentOS slave
Bigtop Toolchain



Fedora slave
Bigtop Toolchain



Ubuntu slave
Bigtop Toolchain



Debian slave
Bigtop Toolchain



OpenSuSE slave
Bigtop Toolchain





**Hi Bigtop, How do I test
my packaging code?**



**We got Bigtop Toolchain
for you.**

A contributor

A committer

I need to setup all the
Linux environments on
my laptop?

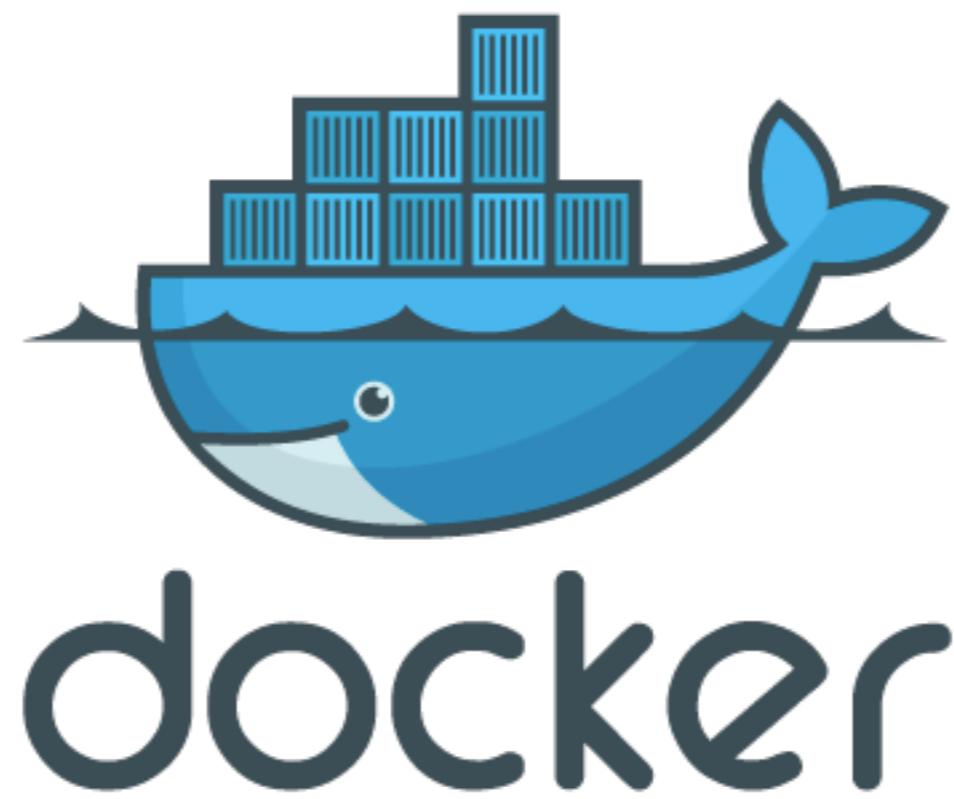


A contributor

(Don't want to answer...)



A committer



Docker - How it works

- Three key techniques in Docker
 - Use Linux Kernel's **Namespaces** to create isolated resources
 - Use Linux Kernel's **cgroups** to constrain resource usage
 - **Union file system** that supports git-like image creation

The nice things about it

- Lightweight
- Fast creation
- Repeatable
- Portability
 - runs on **any** Linux environment

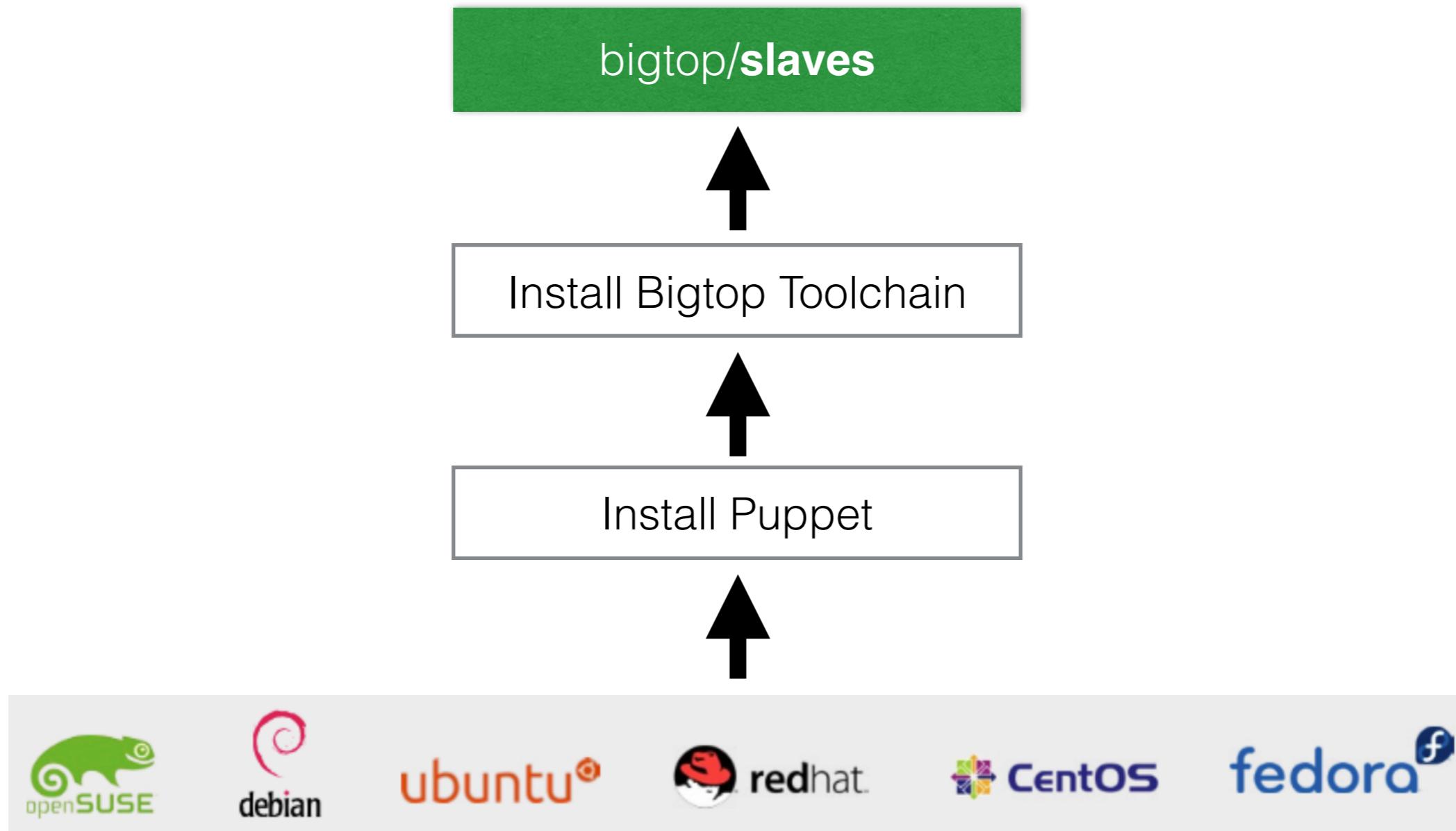
Best fit in our use case



- Run on any kind of Linux distribution in one single machine

Ship build environment in Docker images

Bigtop/slaves images

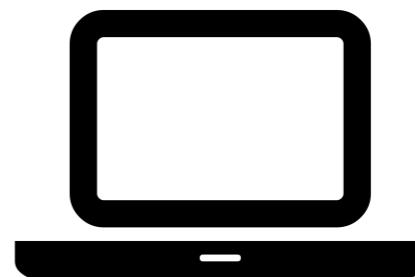
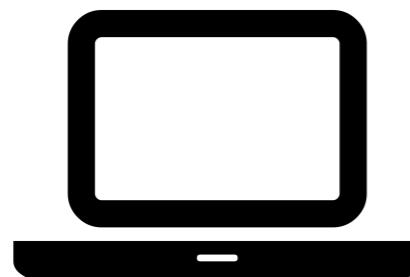
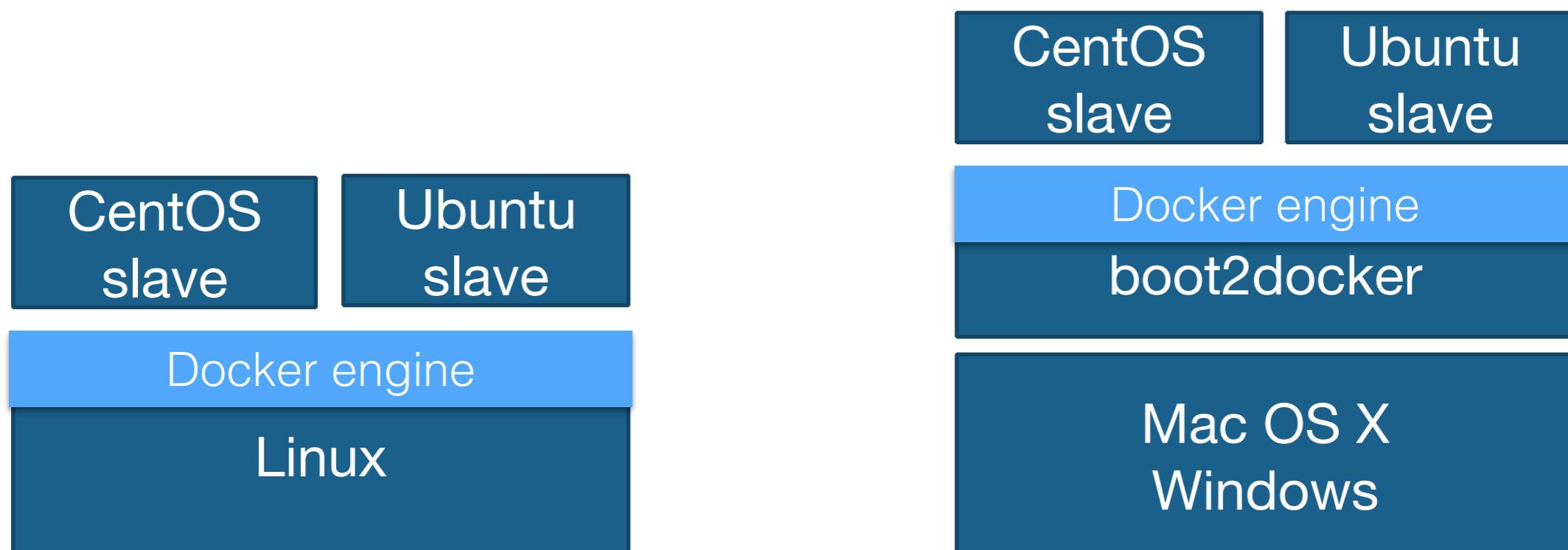


Dockerhub official images

One click to build packages

- \$ git clone <https://github.com/apache/bigtop.git>
- \$ docker run \
--rm \
--volume `pwd`/bigtop:/bigtop \
--workdir /bigtop \
bigtop/slaves:centos-7 \
bash -l -c './gradlew rpm'

Now, easy to build & test on your laptop

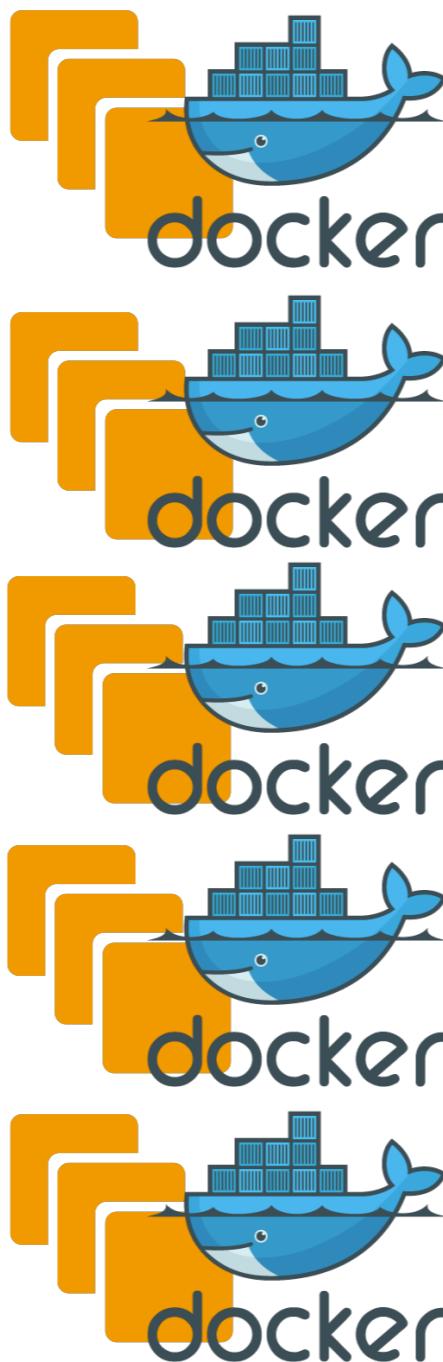
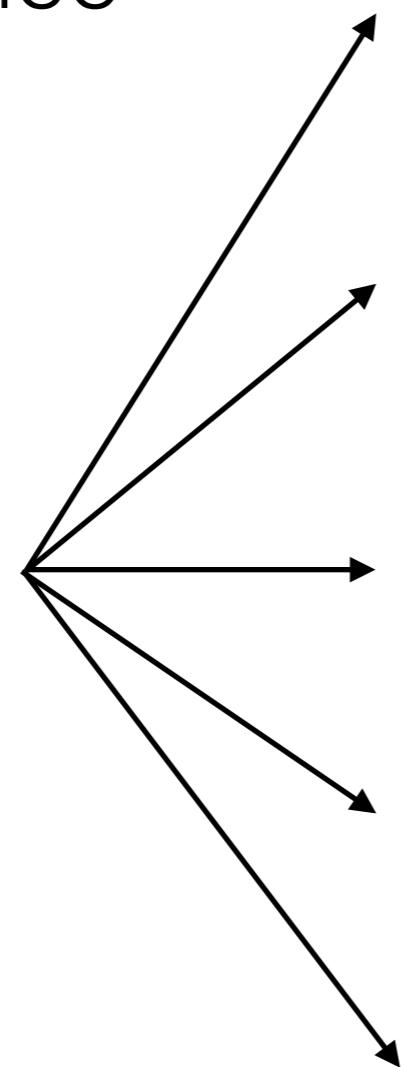


Flexible CI infra

- Fault tolerance
- Immutable



Jenkins



CentOS slave

Fedora slave

Ubuntu slave

Debian slave

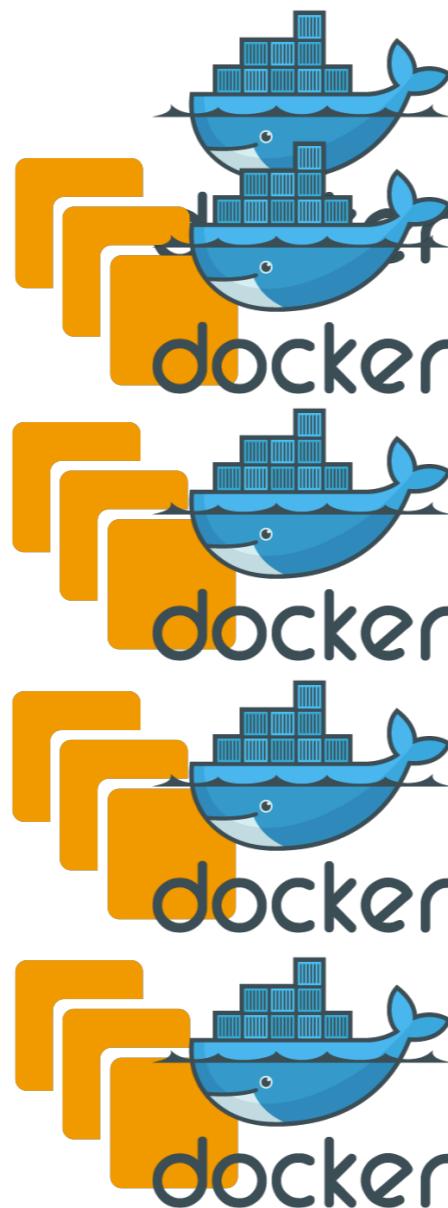
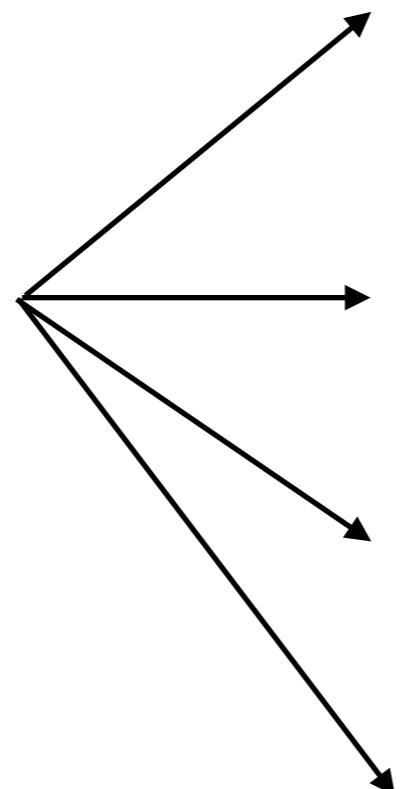
OpenSuSE slave

Flexible CI infra

- Fault tolerance
- Immutable



Jenkins



CentOS slave
Fedora slave

Ubuntu slave

Debian slave

OpenSuSE slave

What is Apache Bigtop?

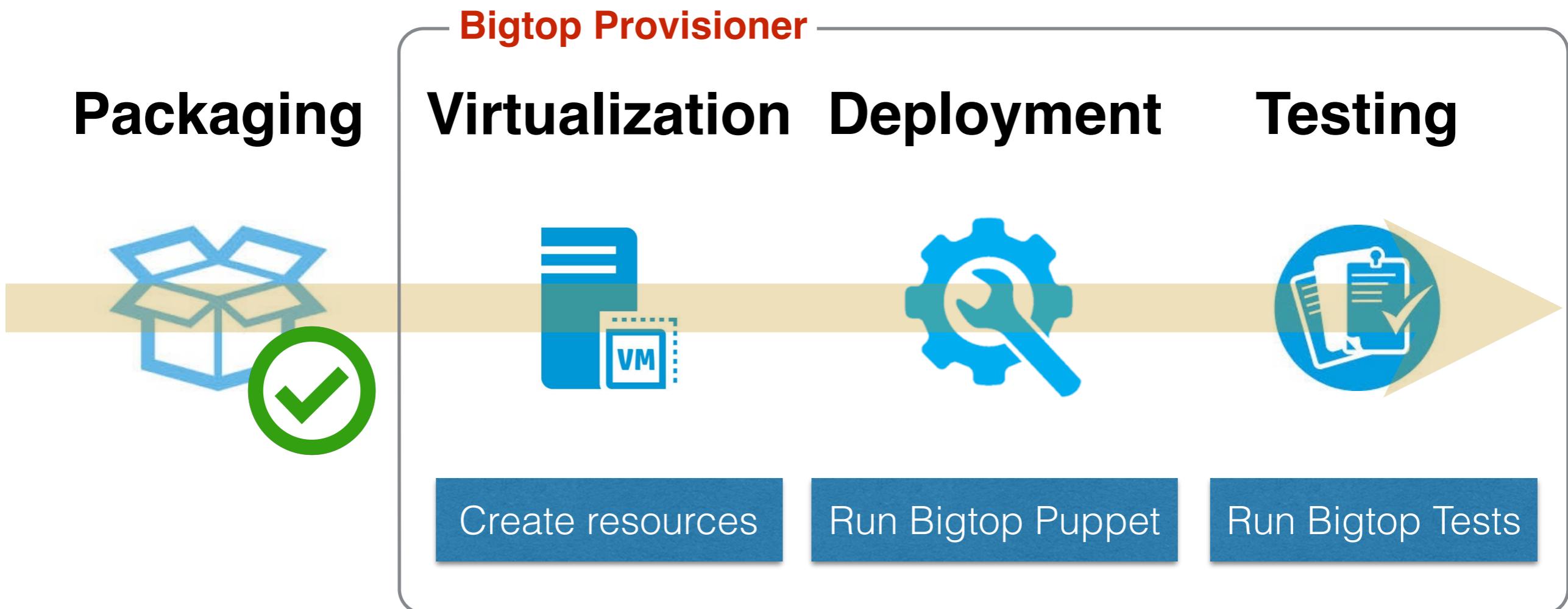
Achieving Build Automation

One-Click Hadoop Provisioning

Bigtop at Trend Micro

Bigtop Provisioner

- A tool to demonstrate the full life cycle of Bigtop



The goal

- **Fast iterative development**
 - Test your code in the cluster, on your laptop, w/o human intervention
- **Flexibility**
 - Choose any combination of components as you want
- **Responsive CI**
 - Integration tests that get you the result in mins
- **A Big Data Stack playground**
 - Spark + Tachyon, Spark + Ignite, etc

**Bigtop
Provisioner**



Vagrant^{w/ provider}

+

Automation Code

+

Bigtop Puppet

||

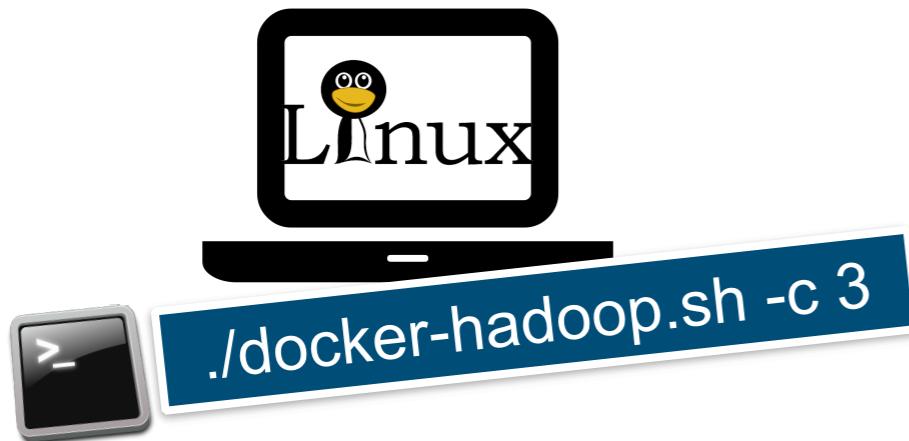
One-click Hadoop Provisioning

Vagrant

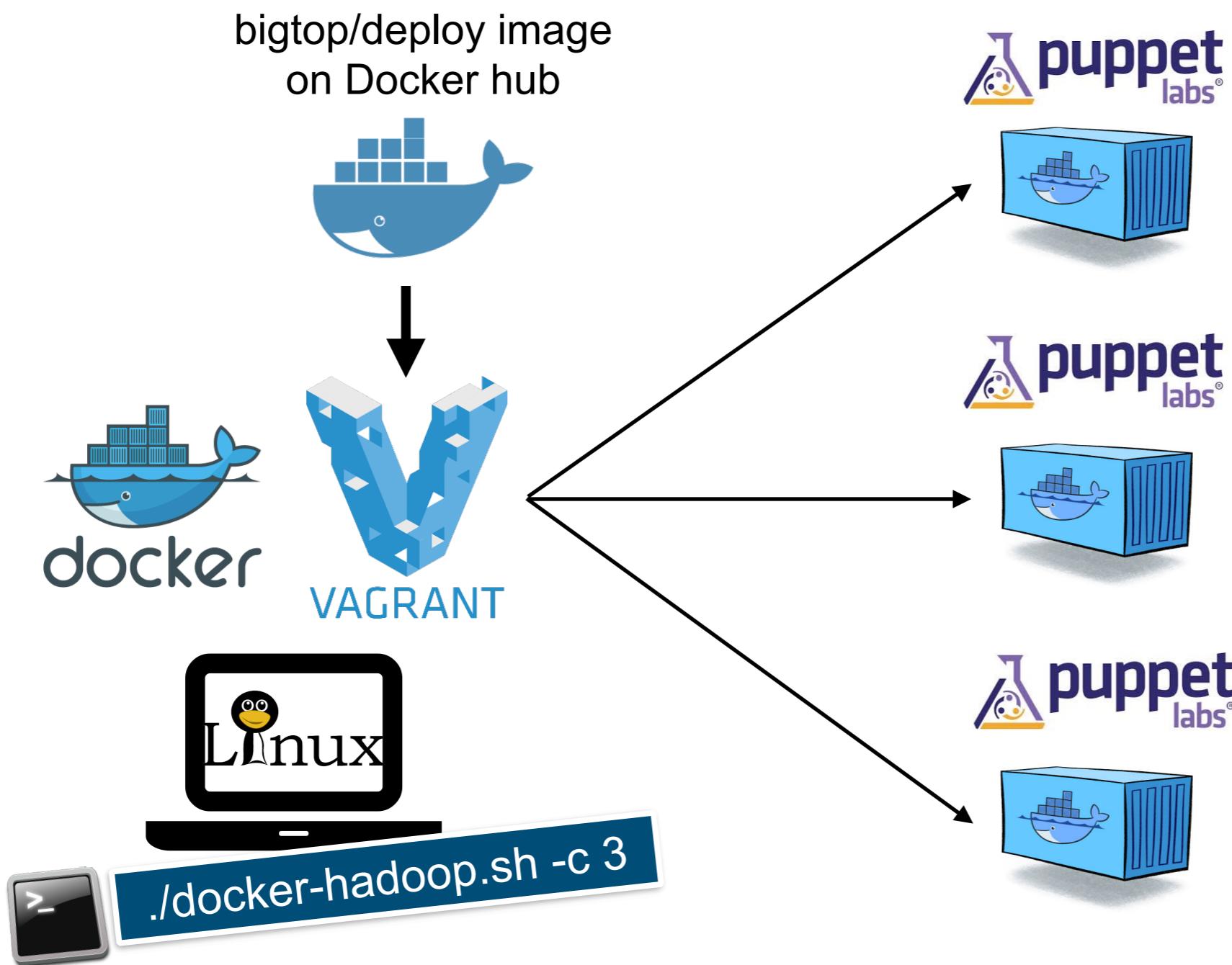
- We use Vagrant as an abstraction layer to support different kind of resource providers



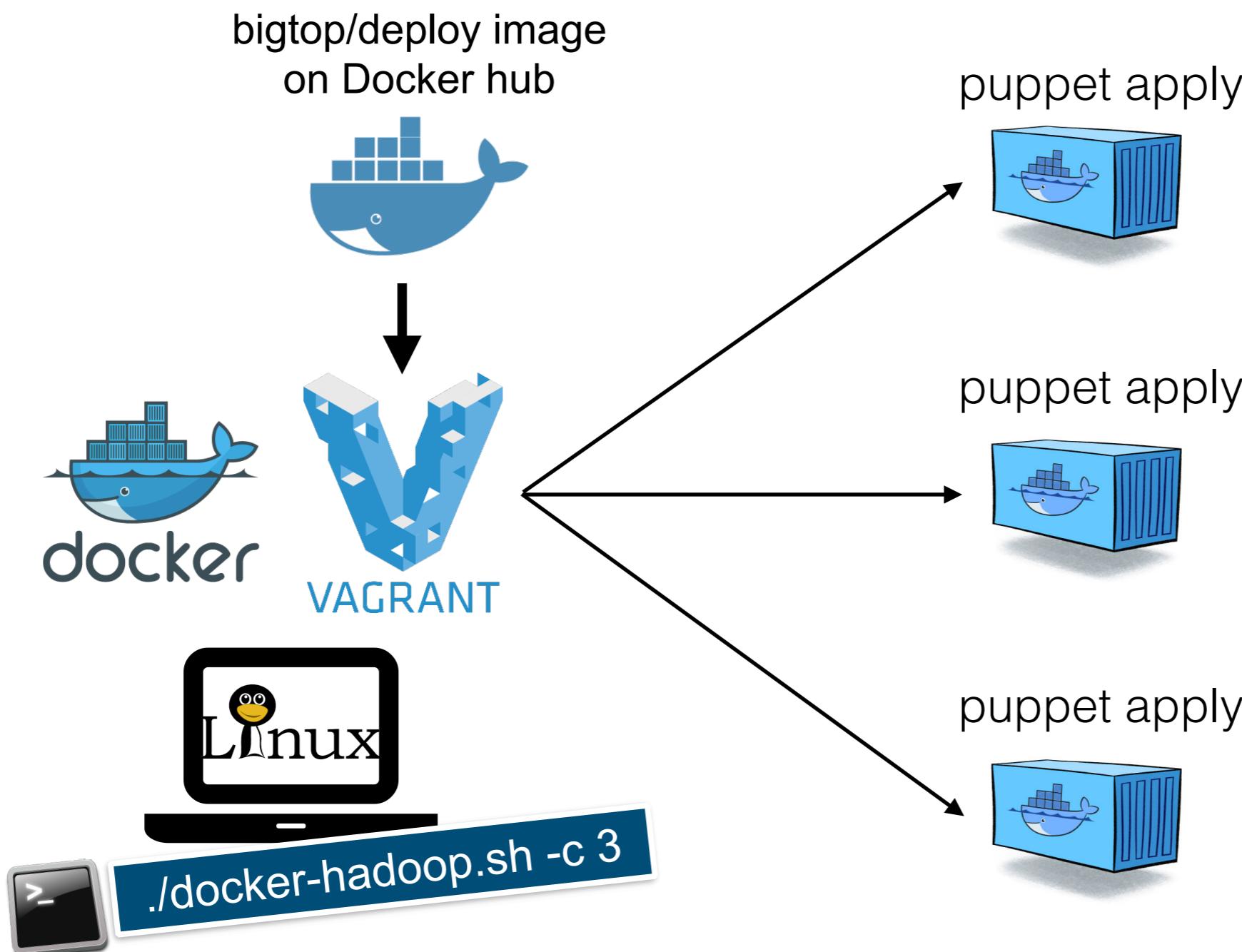
One click Hadoop provisioning



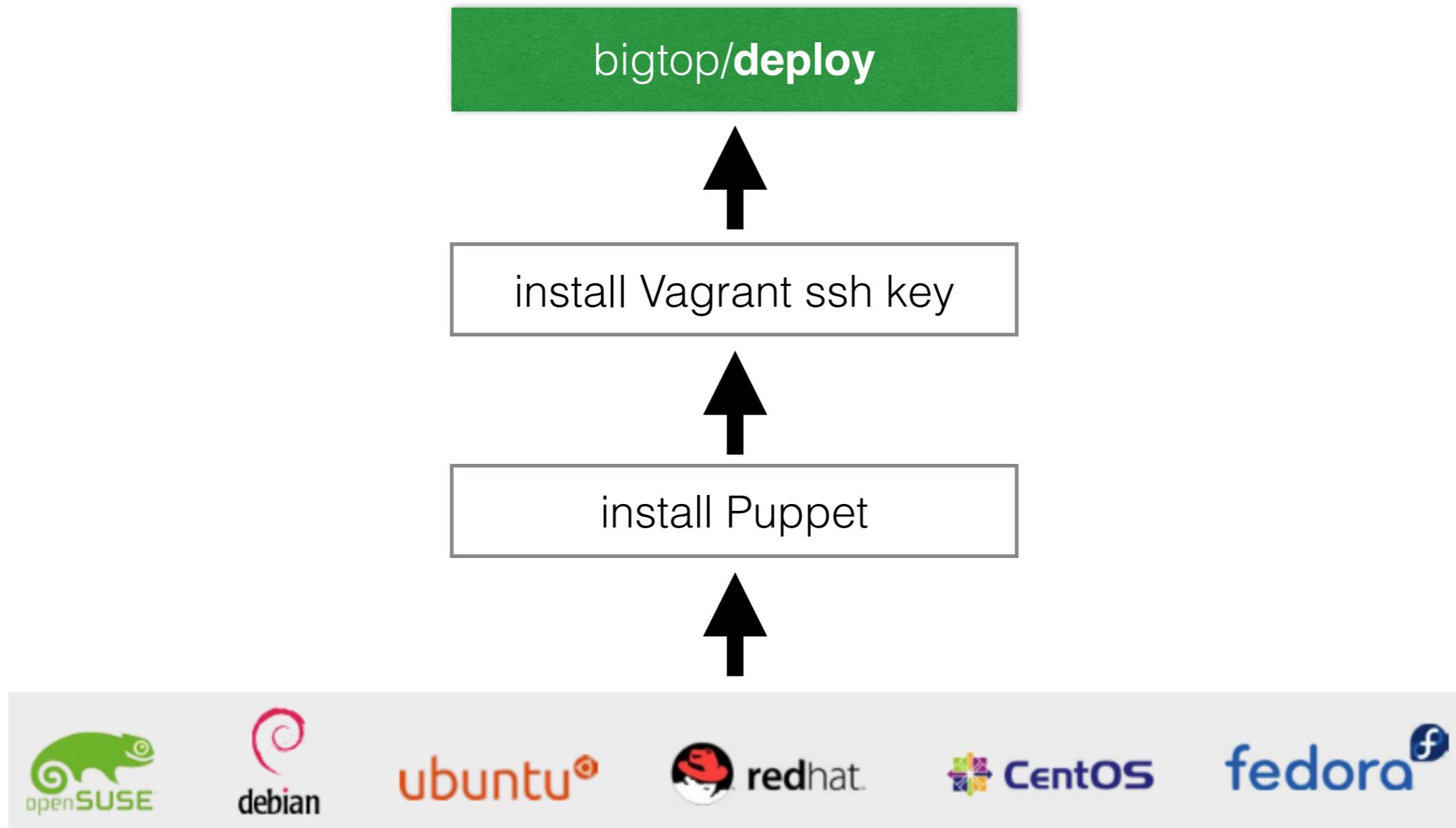
One click Hadoop provisioning



One click Hadoop provisioning

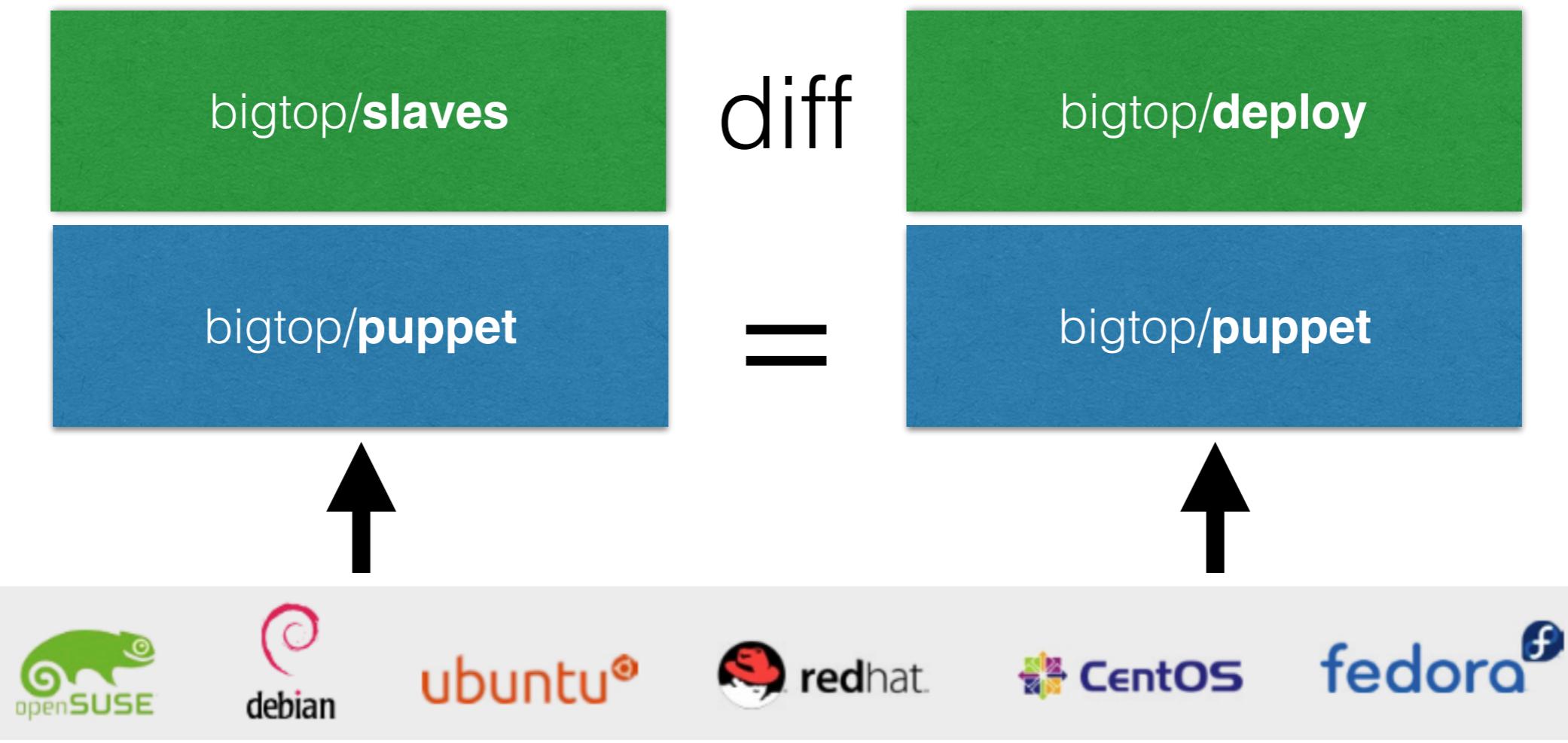


Bigtop/deploy images



Dockerhub official images

Bigtop image hierarchy



All on the Docker hub

Screenshot of the Docker Hub interface showing the Apache Bigtop organization page.

The search bar at the top contains "bigtop".

The sidebar on the left shows a profile picture of a red and white striped circus tent, the name "bigtop", and the text "Apache Bigtop". Below this are links to "The Net", "http://bigtop.apache.org/", and "Joined August 2013".

The main content area displays five Docker repository cards:

Repository	Description	Stars	Pulls	Actions
bigtop/seed	public	0	284	DETAILS
bigtop/puppet	public	1	155	DETAILS
bigtop/slaves	public	2	337	DETAILS
bigtop/jeos	public	0	13	DETAILS
bigtop/deploy	public	1	49	DETAILS

<https://hub.docker.com/u/bigtop/>

A demo is worth 1k words

Oops! demo fail...

<https://asciinema.org/a/55aw3zl2g3dzsfe69198homrl>

Bigtop Provisioner

- Supported providers in Bigtop 1.0.0 release
 - Virtaulbox VM
 - Docker container
- OpenStack support is in master branch now

Use cases

- **For Hadoop app 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 Stack
- **For contributors**
 - Easy to test your packaging, deployment, testing code
- **For vendors**
 - CI out of the box —> patch upstream code made easier

What is Apache Bigtop?

Achieving Build Automation

One Click Hadoop Provisioning

Bigtop at Trend Micro

Trend Micro Hadoop (TMH)

- Use Bigtop as the basis for our internal custom distribution of Hadoop
- Apply **community, internal patches** to upstream projects for business and operational need
- Newest TMH7 is based on Bigtop 1.0 SNAPSHOT

Working with community made our life easier

- Knowing community status made TMH7 release based on Bigtop 1.0 SNAPSHOT possible



Working with community made our life easier

- Knowing community status made TMH7 release based on Bigtop 1.0 SNAPSHOT possible



- Contribute Bigtop Provisioner, packaging code, puppet recipes, bugfixes, CI infra, anything!

Working with community made our life easier

- Leverage Bigtop smoke tests and integration tests with Bigtop Provisioner to evaluate TMH7



Working with community made our life easier

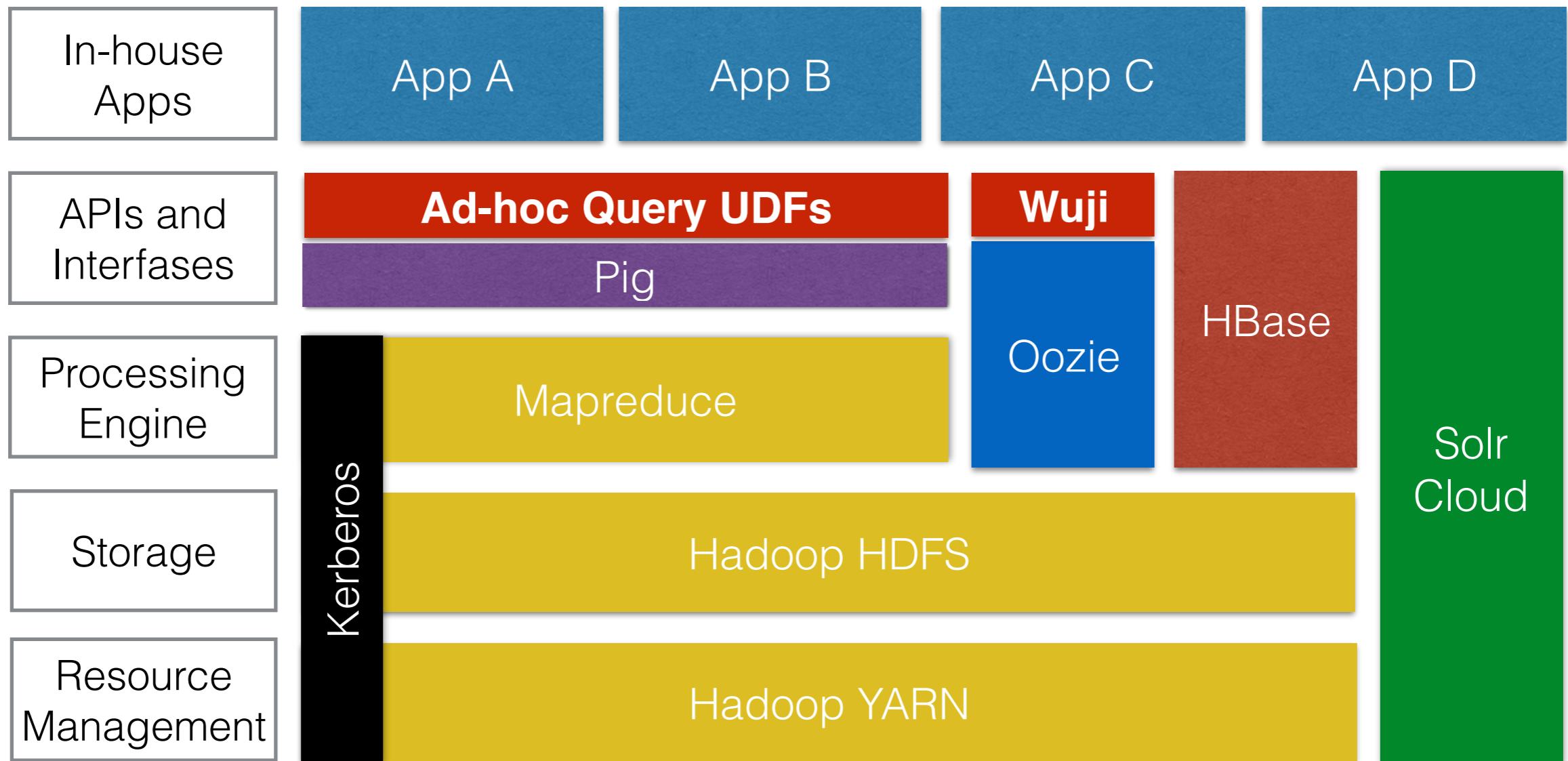
- Leverage Bigtop smoke tests and integration tests with Bigtop Provisioner to evaluate TMH7



- Contribute feedback, evaluation, use case through Production level adoption

Trend Micro Big Data Stack

Powered by Bigtop

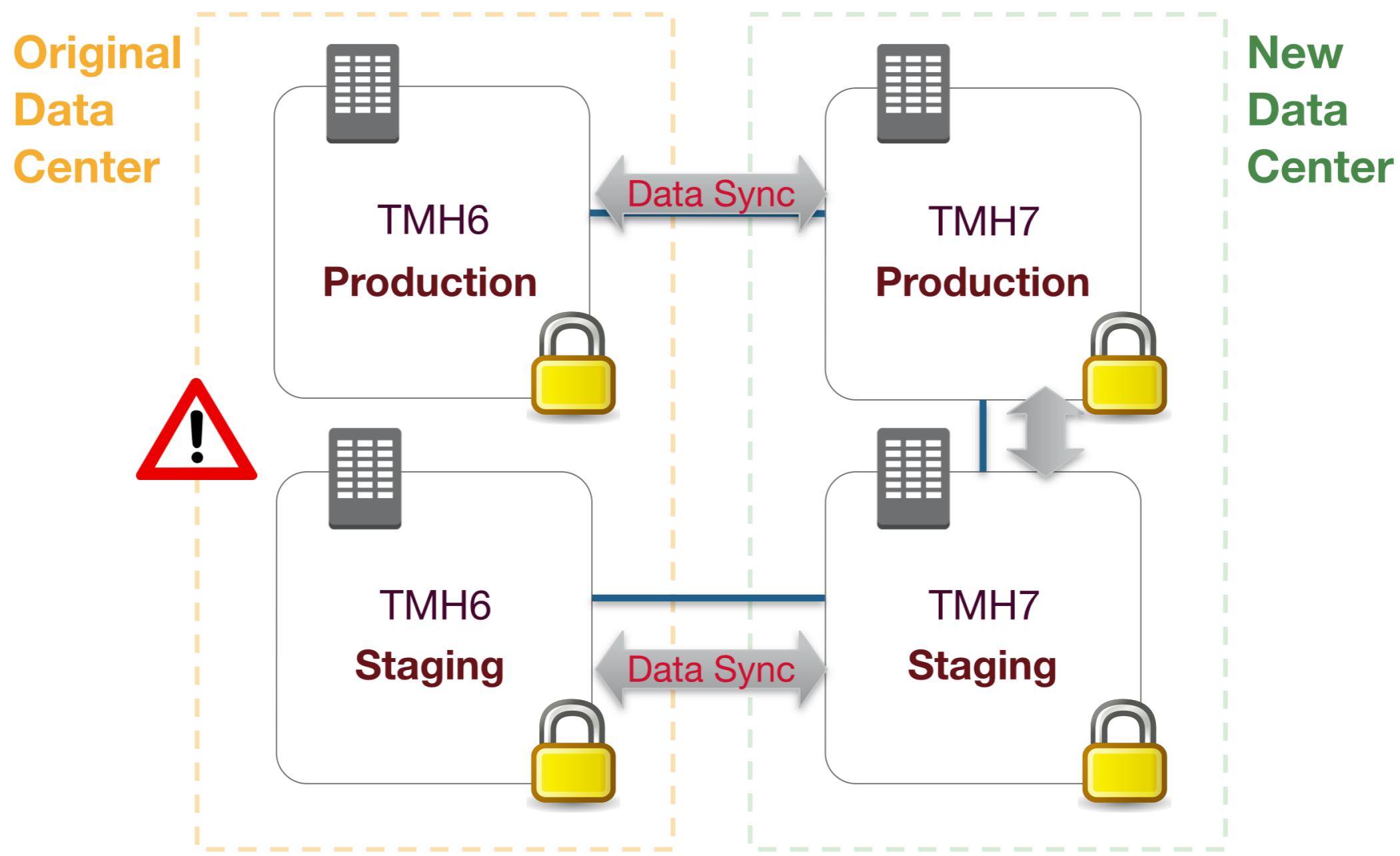


Challenges we're facing !

Migration x Upgrade

- We're moving our cluster to another Data Center
- New cluster will be running **TMH7**
- Need Distcp to sync data from old **TMH6** cluster
- Both are **Kerberos** enabled secure clusters

Sync data between clusters



Distcp between TMH6 and TMH7 ?

Hadoop 2.0.0

Hadoop 2.6.0

Factors that matter

- Hadoop version
- Protocols (hdfs, hftp, webhdfs)
- Where to issue Distcp (Where to run MR job)
- Secure or not

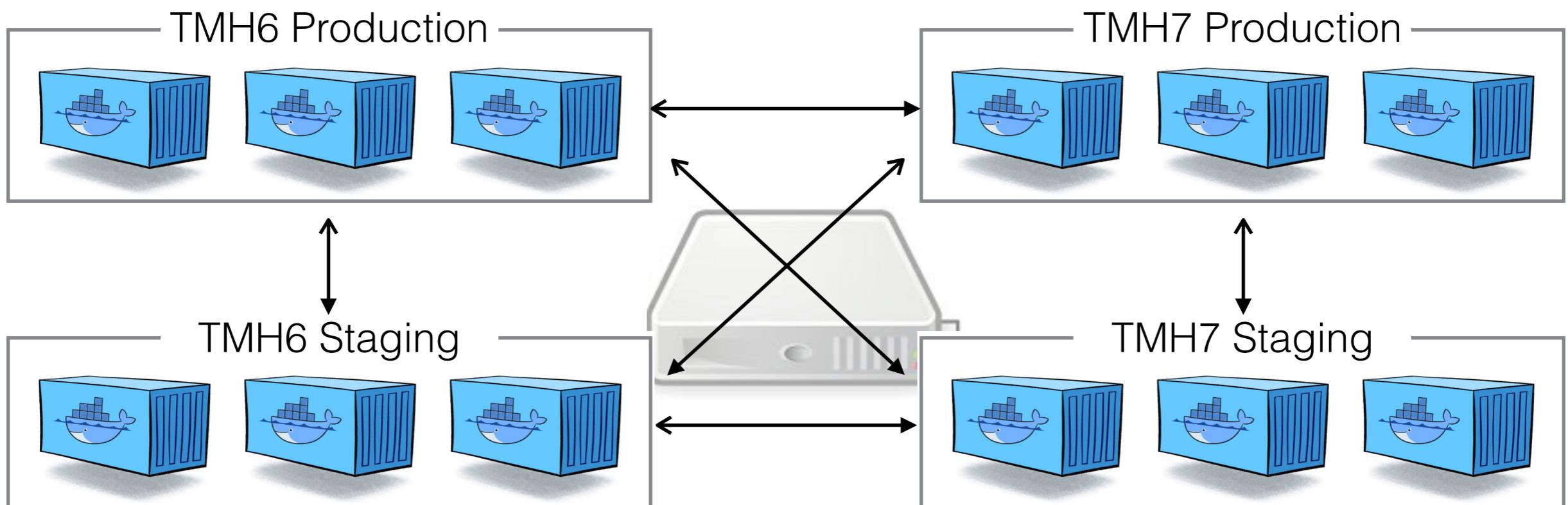
Things always get complicated
with Kerberos...

Kerberos cross realm authentication

- We need to do Distcp across 4 secure clusters
—> Kerberos cross realm auth across 4 clusters
- Our Hadoop management tool needs to support this through auto-configuring
- Developing the management tool is challenging !

Docker comes to the rescue

- Run multiple secure HA clusters on Docker
- Dev & test the management tool iteratively

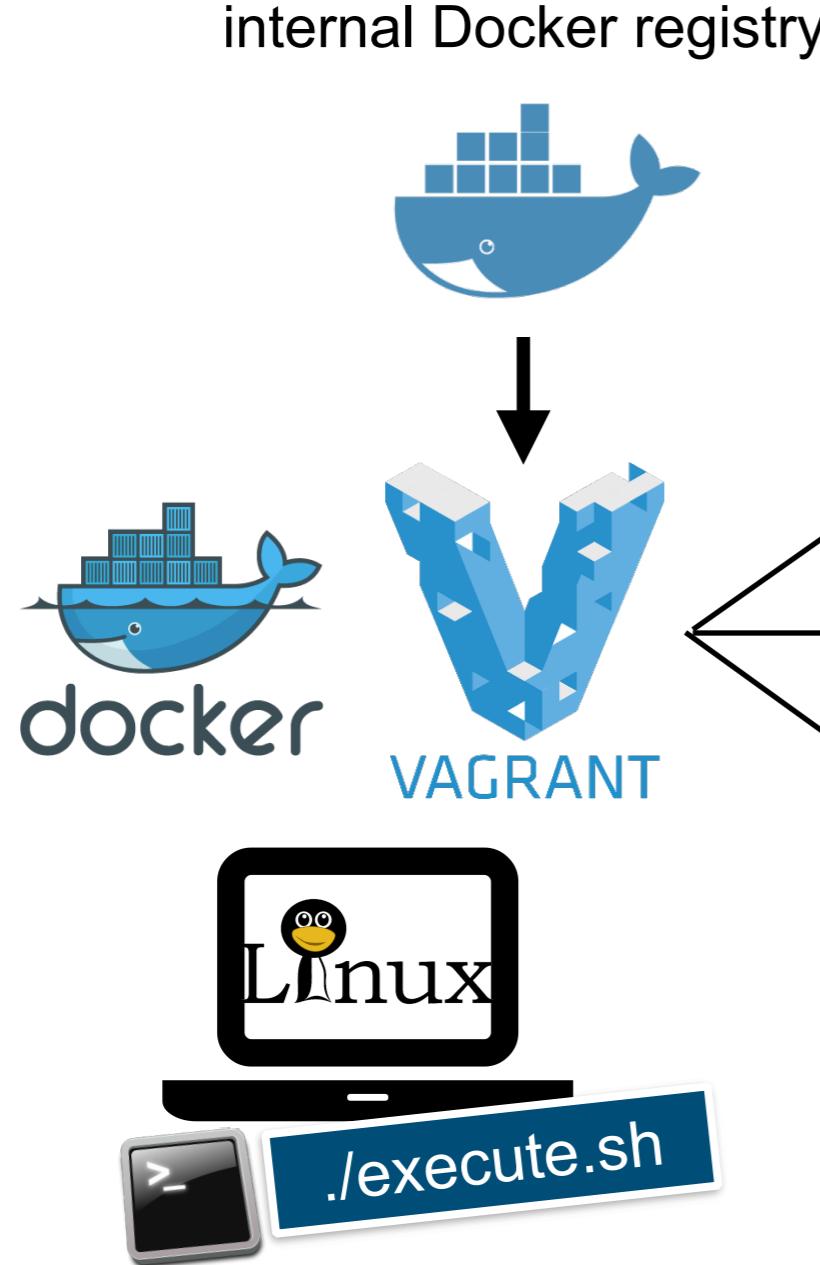


Hadoop apps dev & test on Docker

Hadoocker

- A Devops toolkit for Hadoop app developer to develop and test its code on
- Provides fixed Big Data Stack preload images
 - > dev & test env up and running w/o deployment
 - > support end-to-end CI test for apps
- A Hadoop env for apps to test against our new Hadoop distribution
- Same features are also delivered in Vagrant boxes
- <https://github.com/evans-ye/hadoocker>

Docker based dev & test env



Hadoop server

Hadoop client

data

TMH7



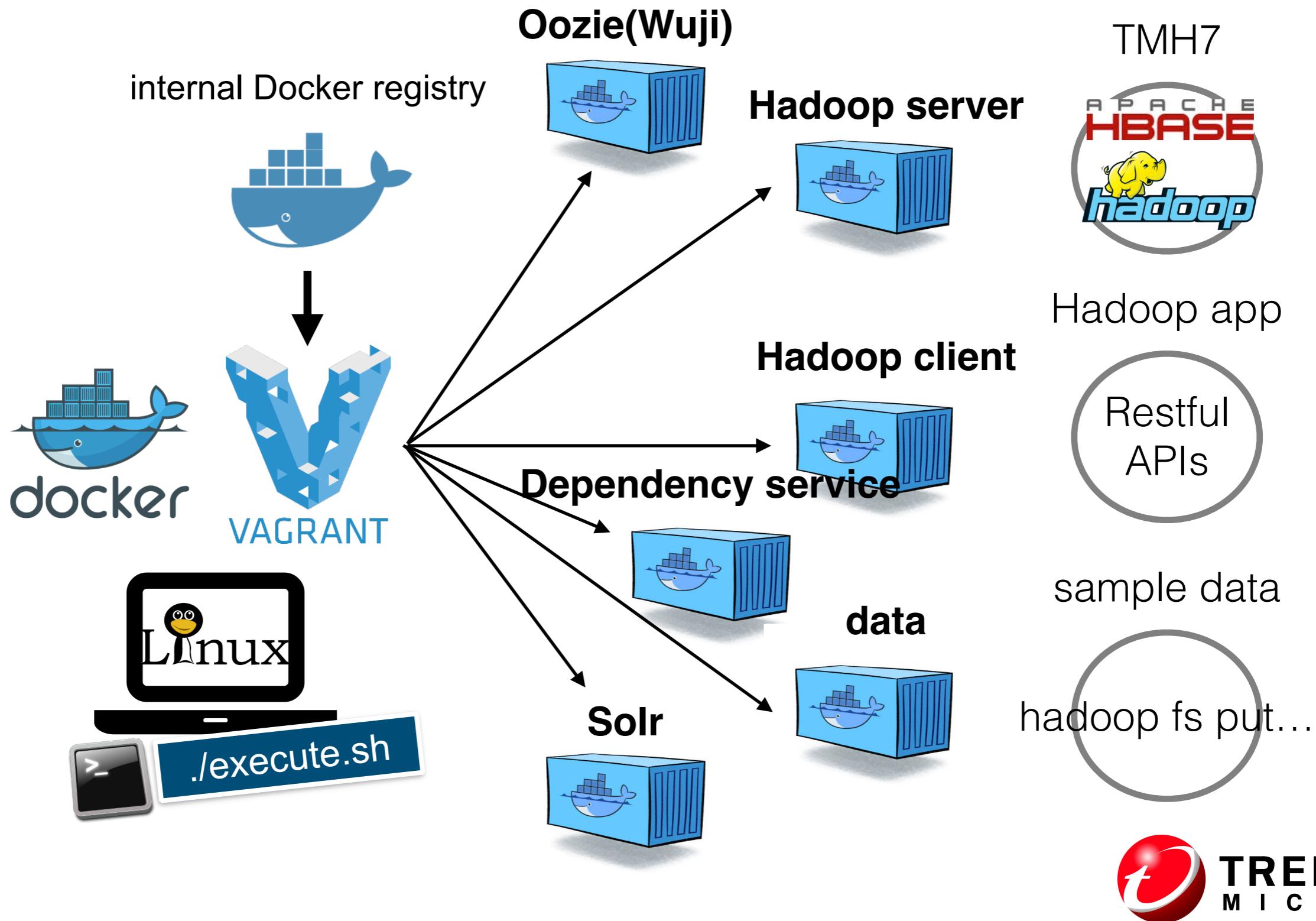
Hadoop app

Restful APIs

sample data

hadoop fs put

Docker based dev & test env



Summary

Summary

- Bigtop helps you to build your own Big Data Stack
- Building packages made easier by offering immutable build environment through Docker images
- Bigtop Provisioner creates you a Hadoop cluster by one click with flexibility
- Use Docker to simulate complex environment to ease development and testing efforts
- Ship apps in Docker images to dev & test anywhere

Wait,
how's the demo ?



Thank you !

Questions ?