# BDE Architecture

Hajira Jabeen, University of Bonn

## M1-M18 Review Meeting

# Structure

◎Evolution of BDE architecture

◎User of BDE

◎Working

# Platform Description

# Technology assessment

◎Lessons learned:
- A lot of technologies available
- Big Data space moves fast
- High barrier to entry

◎Focus:
- Ease of use
  - ❖ Installation, development, deployment, monitoring
- Flexibility
  - ❖ Keep options open for future
- Reuse effort of the community
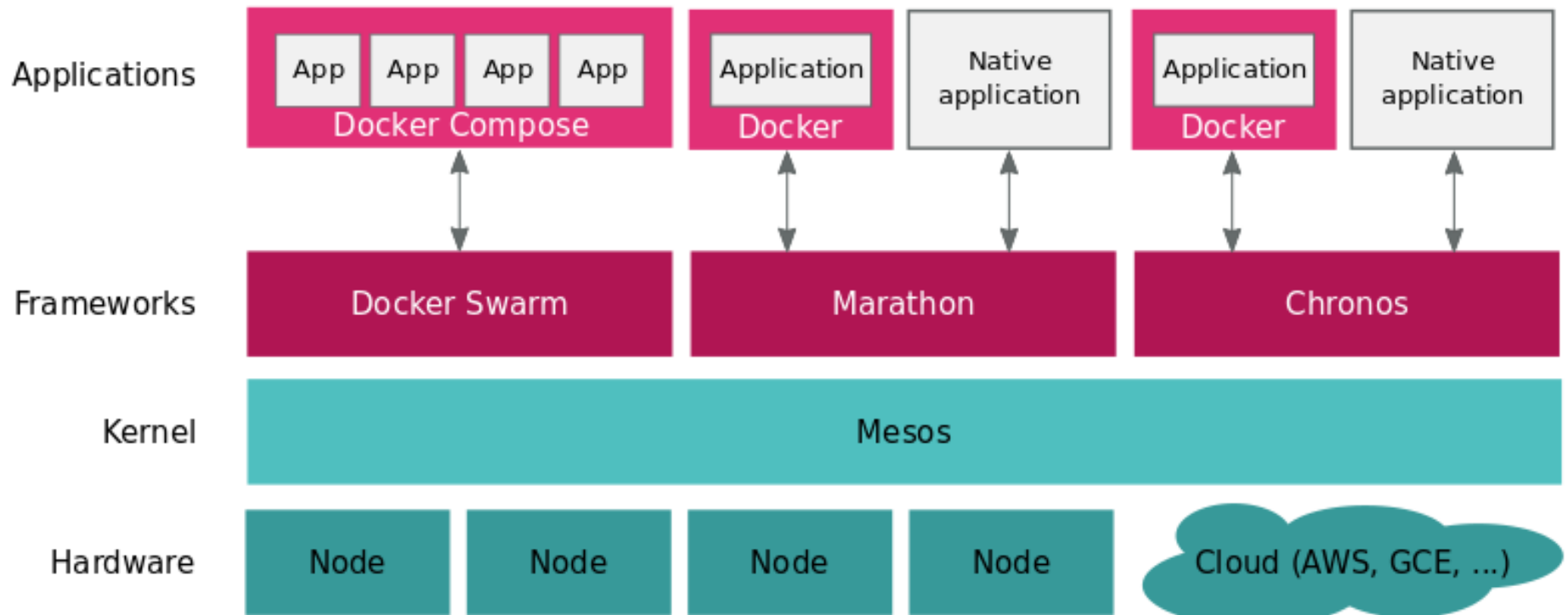  - ❖ Don't reinvent the wheel

# Technical requirements

◎ Input:
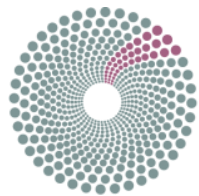- WP2: General requirements elicitation
- WP5: Specific pilot requirements

◎ Initial idea: platform profile per V
- Not 1 V that overrules the others per SC
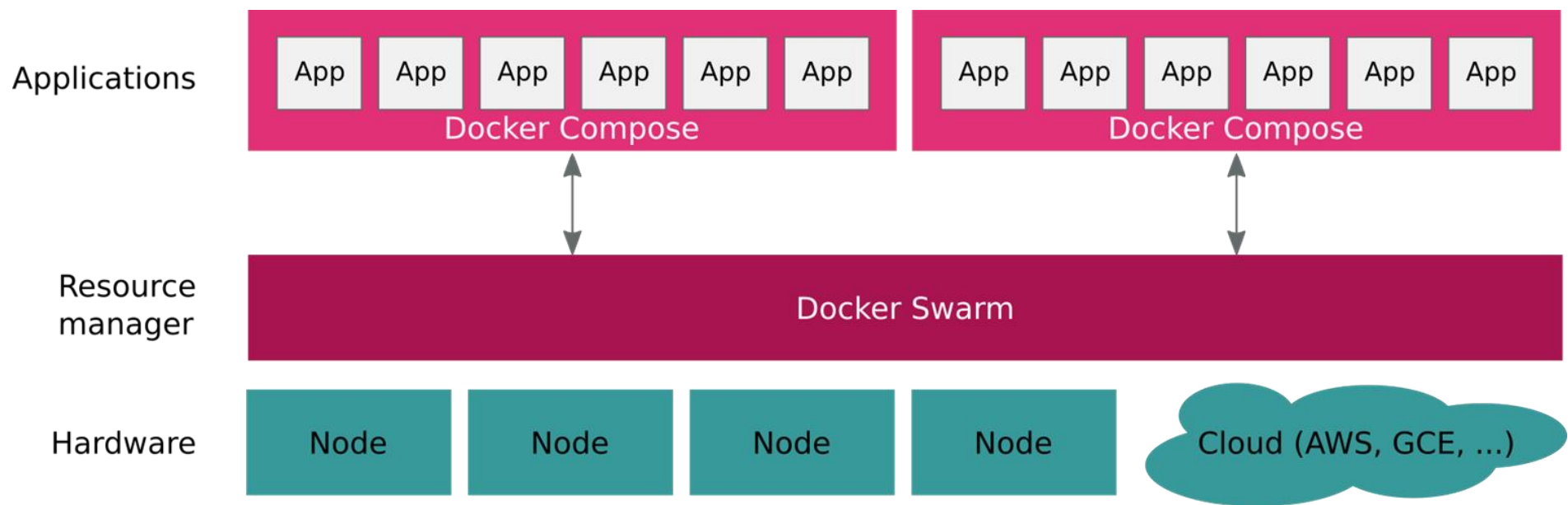  $\Rightarrow$ Provide component suggestions per V

# Architectural design

# Architectural design



Applications

Resource manager

Hardware

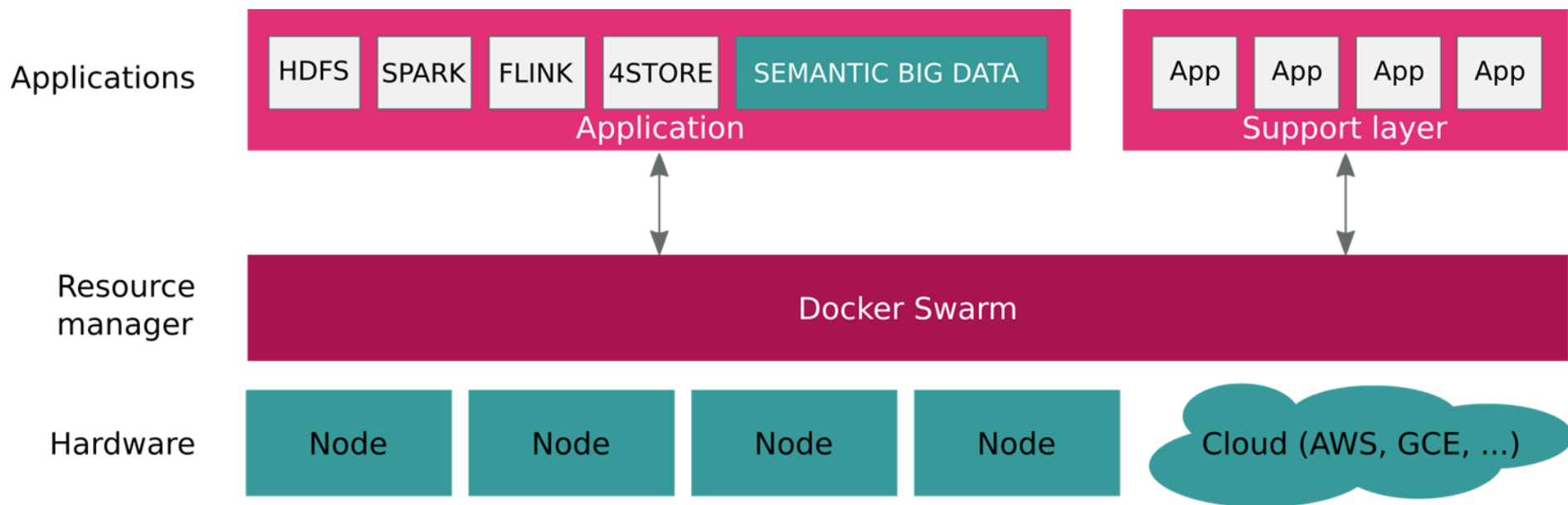

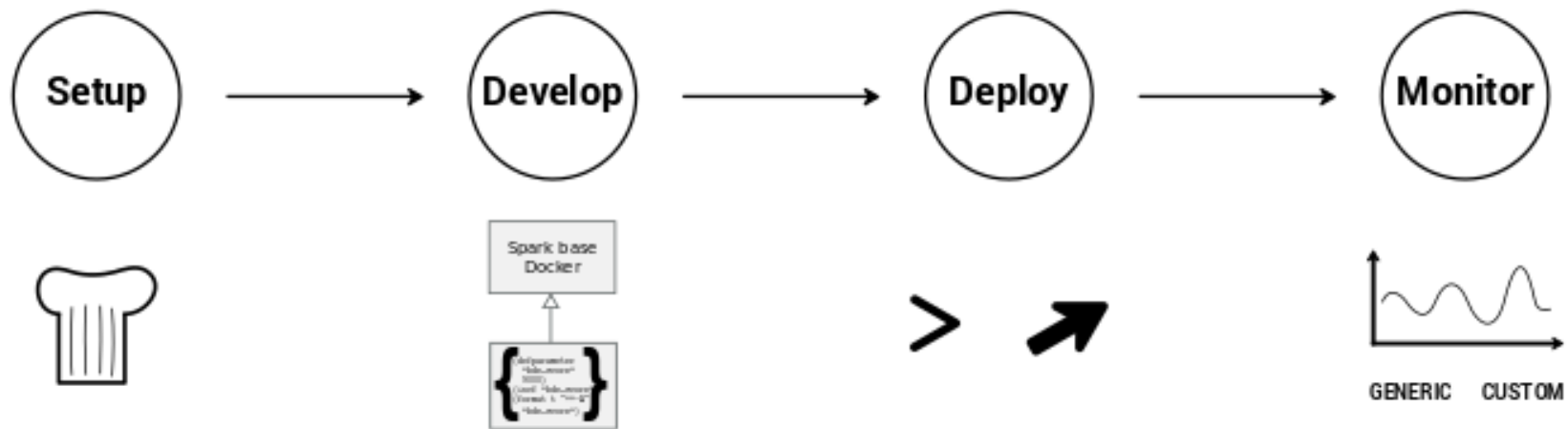

# Architectural design

The minimum knowledge requirements for the BDE user are:

- Ability to write programs for his particular use case
- Inter connectivity of components, if he wants to create a pipeline of different components
- Basics of distributed systems and web-services
- However, this does not exclude experienced users or data scientists from using the platform with ease.

# User profiles

# Platform installation

› Manual installation guide

› Using Docker Machine
- On local machine (VirtualBox)
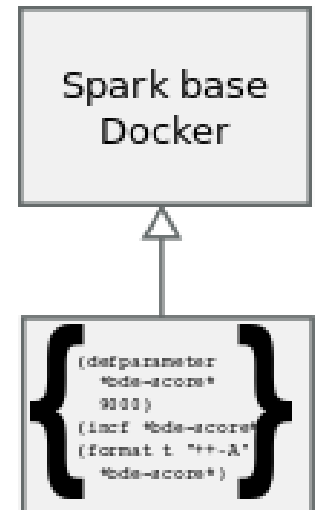- In cloud (AWS, DigitalOcean, Azure)
- Bare metal

◎ Screencast

# Developing a component

◎ Base Docker images
  - ○ Serve as a template for a (Big Data) technology
  - ○ Easily extendable custom algorithm/data

◎ Published components
  - ○ Responsibilities divided b/w partners
  - ○ Image repositories on GitHub
  - ○ Automated builds on DockerHub
  - ○ Documentation on BDE Wiki
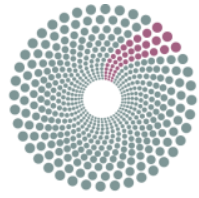
Spark base Docker

# Deploying a Big Data pipeline

◉ Pipeline:
  collection of communicating components
  to solve a specific problem

◉ Described in Docker Compose
  ○ Component configuration
  ○ Application topology

◉ Orchestrator required for initialization process
  ○ Components may depend on each other
  ○ Components may require manual intervention

# Scalability of BDE

◎ 1000 Nodes

◎ 3000 Containers

◎ 1 Swarm Manager

◎ Docker swarm V 1.0

# BDE vs Hadoop distributions

# BDE vs Hadoop distributions

| | Hortonworks | Cloudera | MapR | Bigtop | BDE |
|---|---|---|---|---|---|
| *File System* | HDFS | HDFS | NFS | HDFS | **HDFS** |
| *Installation* | Native | Native | Native | Native | **lightweight virtualization** |
| *Plug & play components (no rigid schema)* | no | no | no | no | **yes** |
| *High Availability* | Single failure recovery (yarn) | Single failure recovery (yarn) | Self healing, mult. failure rec. | Single failure recovery (yarn) | **Multiple Failure recovery** |
| *Cost* | Commercial | Commercial | Commercial | Free | **Free** |
| *Scaling* | Freemium | Freemium | Freemium | Free | **Free** |
| *Addition of custom components* | Not easy | No | No | No | **Yes** |
| *Integration testing* | yes | yes | yes | yes | **--** |
| *Operating systems* | Linux | Linux | Linux | Linux | **All** |
| *Management tool* | Ambari | Cloudera manager | MapR Control system | - | **Docker swarm UI+ Custom** |

# BDE vs Hadoop distributions

BDE is:

- Not built on top of existing distributions
- Targets
  - Communities
  - Research institutions
- Bridges scientists and open data
- Multi Tier research efforts towards Smart Data

# User interfaces

◎ Target: facilitate use of the platform

◎ Available interfaces
- Workflow UIs
  - ❖ Workflow Builder
  - ❖ Workflow Monitor
- Swarm UI
- Integrator UI

# BDE Workflow builder



BDE Workflow Builder       Workflows

## k-means demo

| k-means Spark demo app

## Steps

⇅
### Setup HDFS
Booting of the HDFS cluster.

`setup_hdfs`

DELETE

⇅
### Setup Spark
Starts the Spark master and workers.

`setup_spark`

DELETE

⇅
### Populate HDFS with core data
Please upload the location data to the HDFS filesystem. This is a manual step. Press finish when you're done

`populate_hdfs`

DELETE

# BDE Workflow monitor