# Planning BDI Stack for your Big Data Application
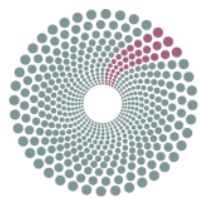
Ivan Ermilov @ ICTCS, Amman, Jordan

# Outline

- Dockerization of Big Data Frameworks (BDF): Why

- What is BDI Stack?

- BDI Stack Lifecycle

- BDI Stack Assembly

- Examples
  - New Spark application
  - Existing Spark application

# Dockerization of BDF: Why

| | Development VM | QA Server | Single Prod Server | Onsite Cluster | Public Cloud | Contributor's laptop | Customer Servers |
|---|---|---|---|---|---|---|---|
| **Static website** | ? | ? | ? | ? | ? | ? | ? |
| **Web frontend** | ? | ? | ? | ? | ? | ? | ? |
| **Background workers** | ? | ? | ? | ? | ? | ? | ? |
| **User DB** | ? | ? | ? | ? | ? | ? | ? |
| **Analytics DB** | ? | ? | ? | ? | ? | ? | ? |
| **Queue** | ? | ? | ? | ? | ? | ? | ? |

docker

# Dockerization of BDF: Why



| | Development VM | QA Server | Single Prod Server | Onsite Cluster | Public Cloud | Contributor's laptop | Customer Servers |
|---|---|---|---|---|---|---|---|
| Static website | | | | | | | |
| Web frontend | | | | | | | |
| Background workers | | | | | | | |
| User DB | | | | | | | |
| Analytics DB | | | | | | | |
| Queue | | | | | | | |

# Dockerization of BDF: Why

◎Development environment

◎Testing environment

◎Staging environment

◎Production environment

**They all the same!**

# Less Duplication = Less Bugs

# What is BDI Stack?

◎Dockerized BDF

◎In one bundle

◎With custom applications

◎docker-compose.yml

# BDI Stack Lifecycle



App templates — Development → Packaging — Ready-made components, best practices

Composition — BDI Stack Builder

Enhancement — BDI Support Components, Instructions

Deployment — Swarm UI, docker-compose

Monitoring — Pipeline/Logging Monitor
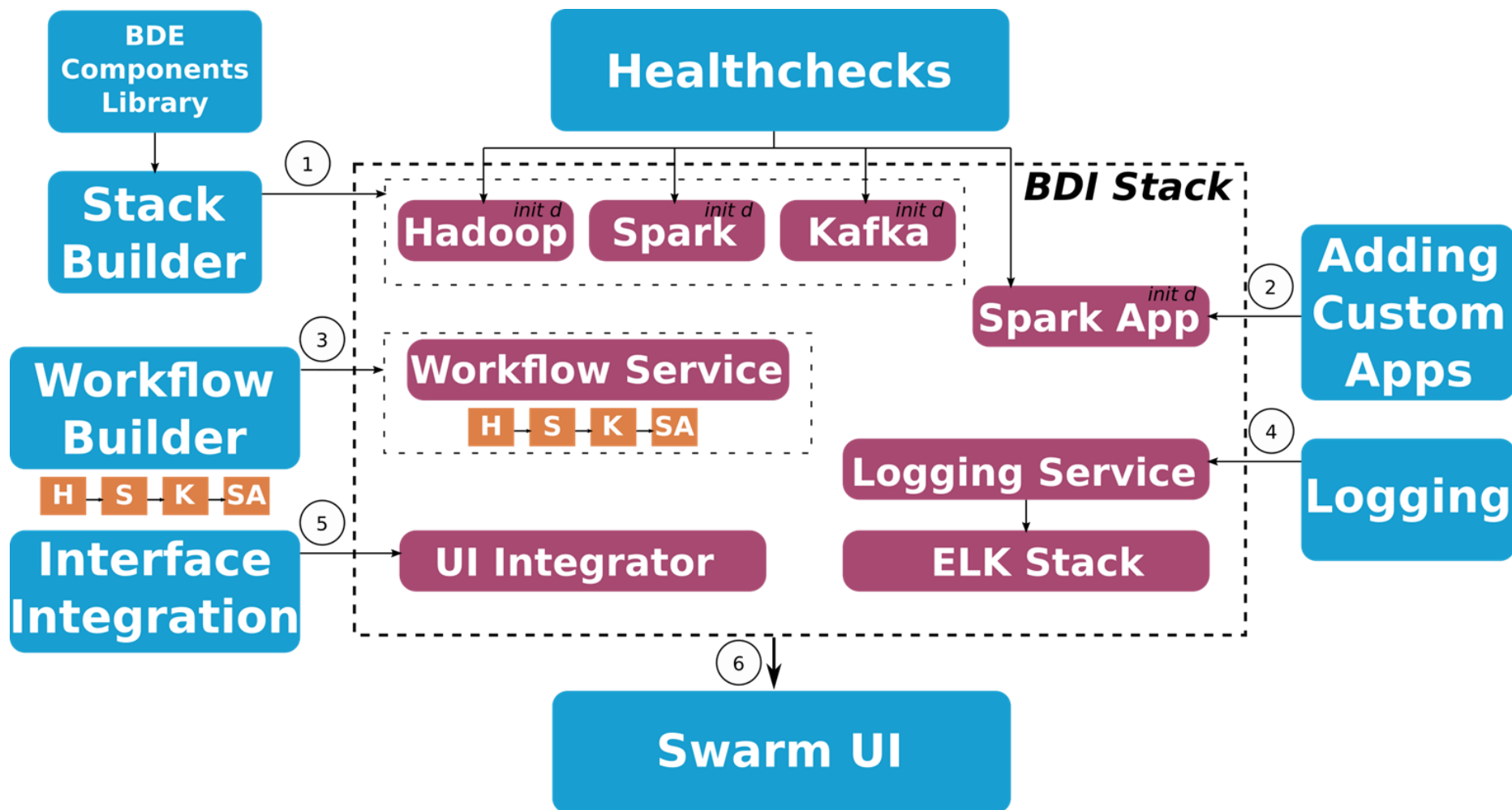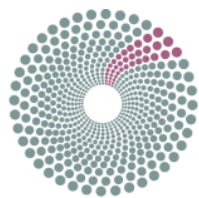
# BDI Stack Assembly

# Stack Builder

## Stack Builder

Compact list view

**Title**

docker-spark

**Text**

```
version: "2"
services:
  spark-master:
    image: bde2020/spark-master:2.1.0-hadoop2.7
    container_name: spark-master
    ports:
      - "8080:8080"
      - "7077:7077"
    environment:
      - INIT_DAEMON_STEP=setup_spark
      - "constraint:node==<yourmasternode>"
  spark-worker-1:
    image: bde2020/spark-worker:2.1.0-hadoop2.7
    container_name: spark-worker-1
    depends_on:
      - spark-master
    ports:
      - "8081:8081"
    environment:
```

Q Search

To add items from the list to your docker-compose file simply drag and drop them into the text field.

| | |
|---|---|
| bde2020/4store | |
| bde2020/4store-master | |
| flink | 2 |
| hdfs | 6 |

# Adding Custom Apps

```
FROM bde2020/spark-submit:2.1.0-hadoop2.7



ENV ENABLE_INIT_DAEMON=false

ENV SPARK_APPLICATION_PYTHON_LOCATION=

ENV SPARK_MASTER_NAME=sc6-spark-master

ENV SPARK_APPLICATION_ARGS=

ENV SPARK_MASTER_URL=spark://sc6-spark-master:7077

ENV SPARK_MASTER_PORT=7077
```

# WorkFlow Builder

**BDE Workflow Builder**

Workflows

# My Workflow

My Workflow

## Steps

⇅

Init Hadoop

Initialization of Hadoop

init-hadoop

DELETE

⇅

Init Spark

Initialization of Spark

init-spark
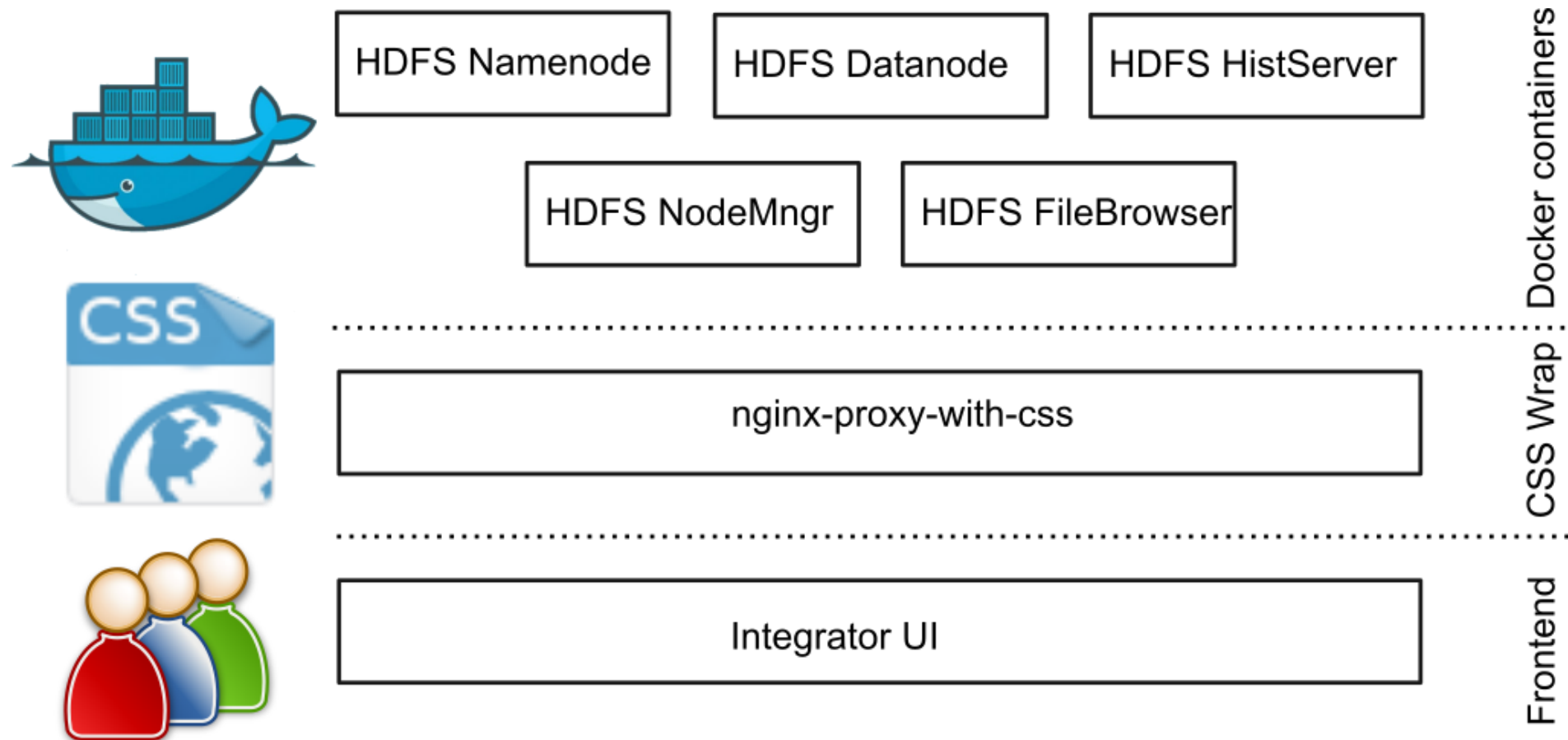
DELETE

*step_code*

# Logging Monitor

◎ Network logging for HTTP

- Capture network interface as PCAP
- Convert to HAR (json)
- Expand HAR
- Dump into ELK stack

HDFS Namenode

HDFS Datanode

HDFS HistServer

HDFS NodeMngr

HDFS FileBrowser

nginx-proxy-with-css

Integrator UI

Docker containers

CSS Wrap

Frontend

https://www.big-data-europe.eu/user-interface-integration-in-bdi-platform-integrator-ui-application/

# UI Integrator

# Reverse Proxy/CSSWrapper

◎Simple injection of custom CSS

```
strabon:

  image: bde2020/strabon

  links:

    - csswrapper

  expose:

    - "8080"

  environment:

    VIRTUAL_HOST: "strabon.big-data-europe.aksw.org"

    VIRTUAL_PORT: "8080"

    CSS_SOURCE: "strabon"
```

https://www.big-data-europe.eu/using-reverse-proxy-inside-bde-platform-jwildernginx-setup-for-docker-swarm/
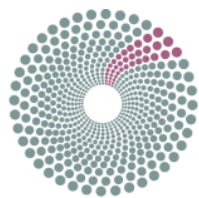https://github.com/big-data-europe/demo-integrator-ui

# Swarm UI

**Swarm UI**

Repositories     Pipelines

## Repositories

test

Located at https://github.com/big-data-europe/demo-spark-sensor-data. Has 1 connected pipelines.

**EDIT**     **LAUNCH**

Create new repository

# Swarm UI
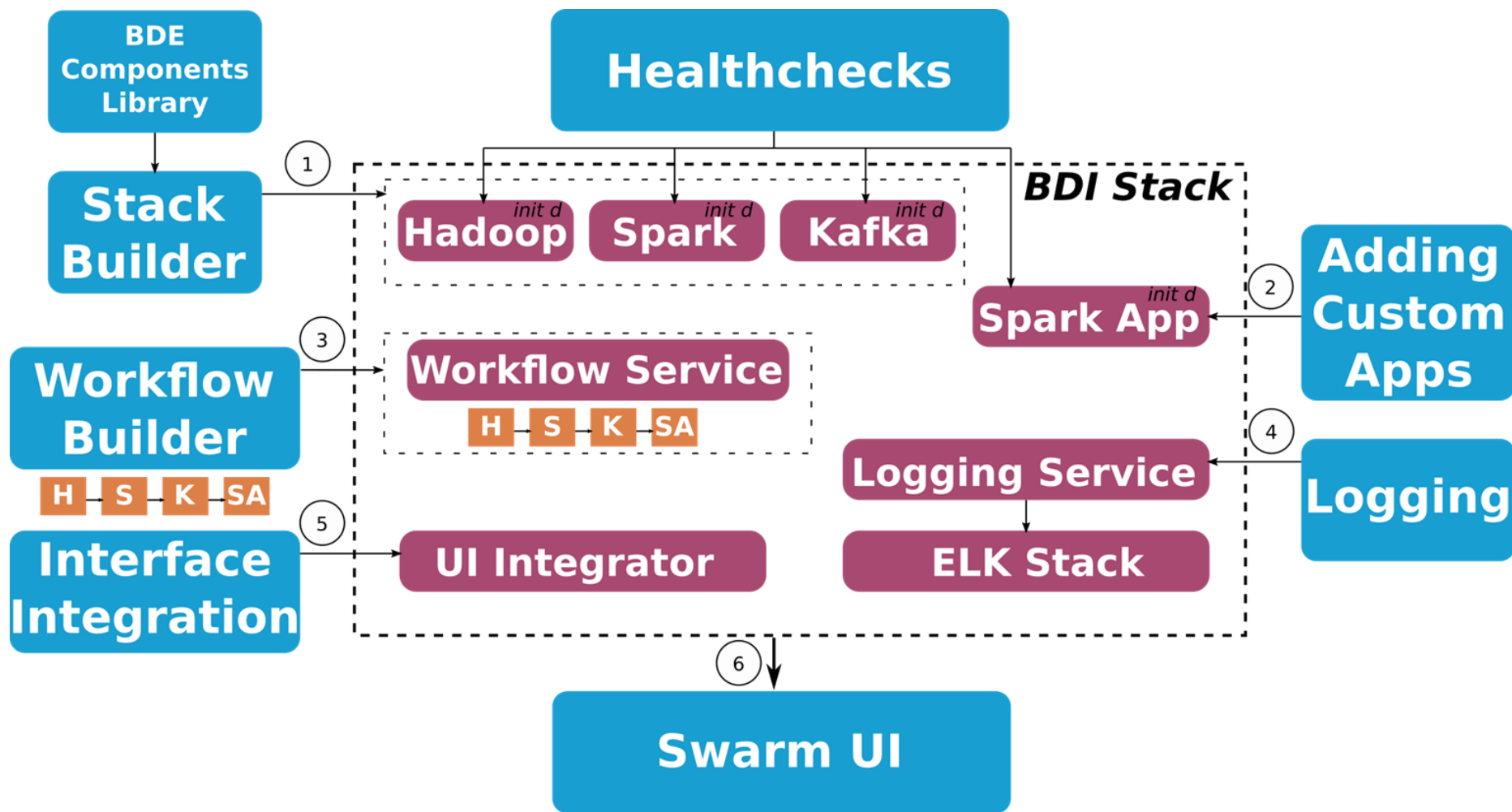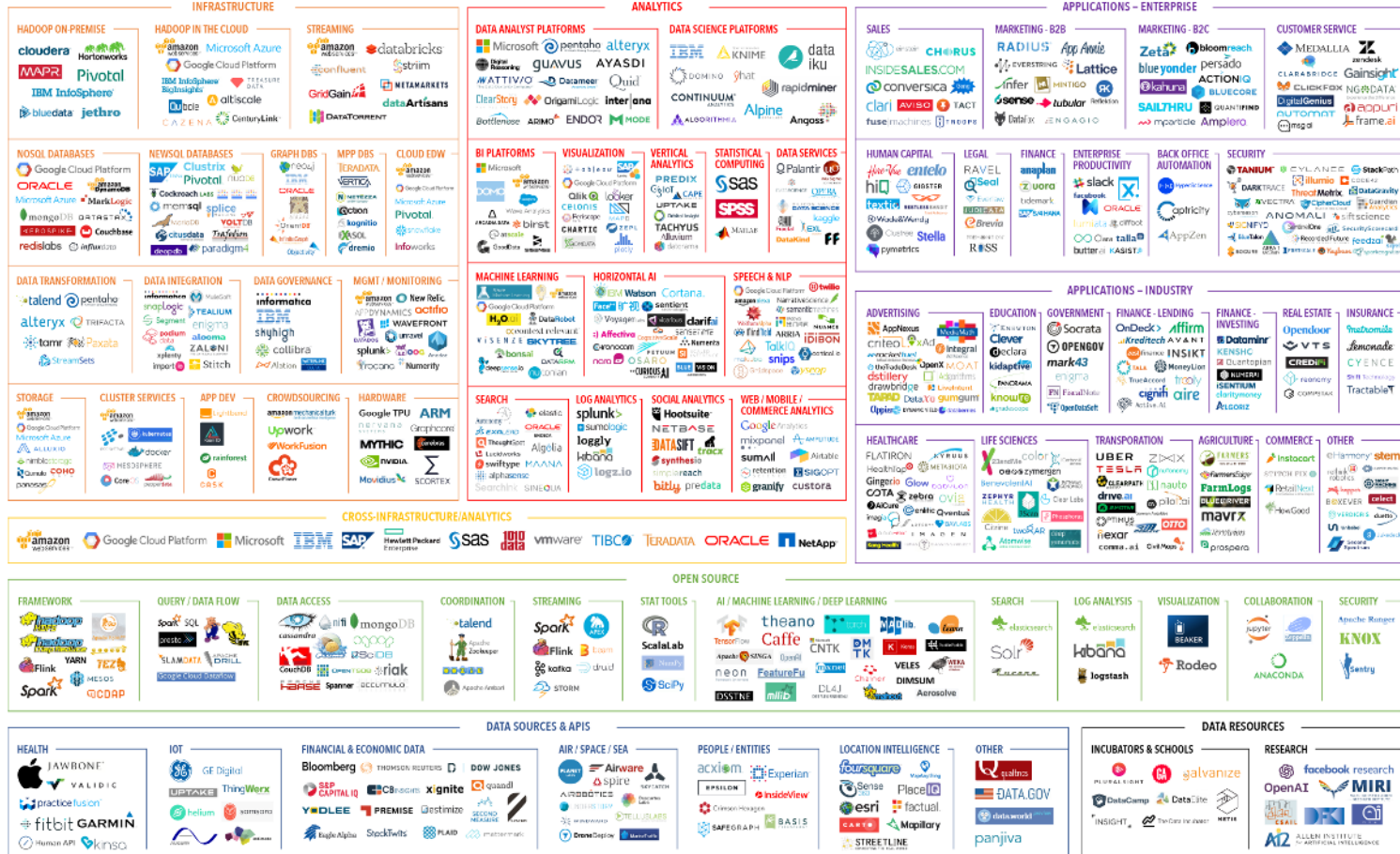
# BDI Stack Assembly

BIG DATA LANDSCAPE 2017

Last updated 4/5/2017 · © Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark (@firstmarkcap) · mattturck.com/bigdata2017

20

# Simplified Workflow

- ◎ Outline your application requirements
- ◎ Pick up components from BDE github repo
  - ○ If it is not there search Docker Hub
  - ○ Else dockerize it yourself
- ◎ Create docker-compose.yml
- ◎ Test with simple application
- ◎ Develop your application on top of it
- ◎ Proceed to enhancement step (if necessary)

# Application Requirements

⊚ Create user stories

- ○ User wants to see the most trending recent hashtags
- ○ User wants to see the most recent visualization of the hashtags
- ○ User wants to see past visualizations as well

# Core Functionality

◎Fetch tweets

- Spark streaming can do that

◎Store tweets somewhere (big data)

- HDFS

◎Store trends (not too big)

- MongoDB

◎Visualize the trends

- Kibana or custom application

# Components from BDE github

◎HDFS
- Store data

◎Spark
- Streaming
- Transformation
- Save/load data

# Searching Docker Hub

Search - Docker Hub

hub.docker.com/search/?isAutomated=0&isOfficial=0&page=1&pullCount=0&q=kafka&starCount=0

Docker Store is the new place to discover public Docker content. Check it out →

kafka

Explore    Help    Sign up    Sign in

## Repositories (3097)

All

| | STARS | PULLS | |
|---|---|---|---|
| **wurstmeister/kafka**<br>public \| automated build | 417<br>STARS | 10M+<br>PULLS | DETAILS |
| **sheepkiller/kafka-manager**<br>public \| automated build | 80<br>STARS | 1M+<br>PULLS | DETAILS |
| **ches/kafka**<br>public \| automated build | 82<br>STARS | 1M+<br>PULLS | DETAILS |
| **confluentinc/cp-kafka-connect**<br>public | 12<br>STARS | 1M+<br>PULLS | DETAILS |
| **spotify/kafka**<br>public \| automated build | 225<br>STARS | 500K+<br>PULLS | DETAILS |
| **cgswong/confluent-kafka**<br>public \| automated build | 7<br>STARS | 1M+<br>PULLS | DETAILS |
| **solsson/kafka** | | 500K+ | |

# Searching Docker Hub

◎ Create a table (example for HBase)

  ○ Name

  ○ Java Version

  ○ Docs

  ○ Configurable

  ○ Standalone

  ○ Pseudodistributed

  ○ Distributed

# Searching Docker Hub

◉ Select the best docker image from the table

○ Pay attention to the docker image license!

◉ Extend if necessary

◉ <u>Example</u>

# Assembling docker-compose

- BDE provides docker-compose.yml snippets
- Docker images which follow the best practices does the same

# Assembling docker-compose

◎Copy/paste the snippets and adjust

```
version: '2'
services:
  namenode:
    image: bde2020/hadoop-namenode:1.1.0-hadoop2.8-
java8
    container_name: namenode
    volumes:
      - ./data/namenode:/hadoop/dfs/name
    environment:
      - CLUSTER_NAME=TwitterTrendsCluster
…
```

# Testing your BDI Stack

◎Manual testing

  ○All containers running?

  ○No errors in the initialization logs?

◎Automatic testing

  ○Deploy a simple application using the BDI stack along it

  ○Produce correct results?

# Ready to develop your app!

# Adding Application to the Stack

◎ Create Dockerfile

◎ Expose external interfaces

  ○ REST

  ○ SQL

  ○ SPARQL

◎ Upload to docker hub

  ○ Or your enterprise repository (e.g. gitlab)

# Demo (15 mins)

◎ BDI Stack for

- ○ Hadoop
- ○ Spark
- ○ TwitterTrends
- ○ VisualizationApp

# Packing existing application

# Example BDI Stack: Halyard

◎ Which BD components does Halyard use?
   ○ HDFS
   ○ YARN (for MapReduce jobs)
   ○ HBase
◎ Which interfaces are supported?
   ○ Shell scripts (bulkload, export etc)
   ○ RDF4J console
   ○ RDF4J REST Server + Workbench

# Halyard: BDI Stack

**Hadoop**

**DFS**
- Namenode
- Datanode
- Resource Manager

**YARN**
- Node Manager
- History Server

**HBase**
- Master
- Region Server

**Zookeeper**
- Zookeeper

◎ Hadoop
  - ○ DFS
  - ○ YARN
◎ HBase
◎ Zookeeper

```yaml
namenode:

  image: bde2020/hadoop-namenode:1.2.0-hadoop2.8-java8

  container_name: namenode

  networks:

    - hbase

  volumes:

    - ./data/hadoop/namenode:/hadoop/dfs/name

  environment:

    - CLUSTER_NAME=test

  ports:

    - "50070:50070"

  env_file:

    - ./hadoop.env
```

Simply execute the command:

docker-compose up -d

```
$ docker exec -it hbase /bin/bash

$ hbase shell

> list

> create 't1', 'f1'
```

```
FROM bde2020/hadoop-base:1.2.0-hadoop2.8-java8 as hadoop-base

FROM bde2020/hbase-base:1.0.0-hbase1.2.6 as hbase-base

FROM openjdk:8

MAINTAINER Ivan Ermilov <ivan.s.ermilov@gmail.com>

ENV HADOOP_VERSION=2.8.0

COPY --from=hadoop-base /opt/hadoop-$HADOOP_VERSION /opt/hadoop-$HADOOP_VERSION

RUN ln -s /opt/hadoop-$HADOOP_VERSION/etc/hadoop /etc/hadoop

ENV PATH /opt/hadoop-$HADOOP_VERSION/bin:$PATH

ENV HBASE_VERSION=1.2.6

COPY --from=hbase-base /opt/hbase-$HBASE_VERSION /opt/hbase-$HBASE_VERSION

RUN ln -s /opt/hbase-$HBASE_VERSION/conf /etc/hbase

ENV PATH /opt/hbase-$HBASE_VERSION/bin:$PATH

ENV HALYARD_VERSION 1.2

…
```

# Running Halyard SDK

```
$ docker run -it --rm --network hbase --
  env-file ./hbase.env bde2020/halyard-
  sdk:1.0.0-halyard1.2 /bin/bash
$ ./console
```

# Running Halyard SDK

```
> create hbase
> open halyard
> load http://danbri.org/foaf.rdf
```

# Running Halyard SDK

```
> sparql
select ?s ?p ?o {where ?s ?p ?o} .
```

# Halyard: BDI Stack (complete)

**Hadoop**
- **DFS**
  - Namenode
  - Datanode
  - Resource Manager
- **YARN**
  - Node Manager
  - History Server

**HBase**
- Master
- Region Server

**Zookeeper**
- Zookeeper

**Halyard**
- sdk
- rdf4j-server
- workbench

**BDE**
- UI Integrator
- Workflow
- Logging

# Thank you

## Questions?

Github: https://github.com/earthquakesan
@AKSW: http://aksw.org/IvanErmilov.html
Email: iermilov@informatik.uni-leipzig.de
Twitter: @earthquakesan
LinkedIn: https://www.linkedin.com/in/iermilov/