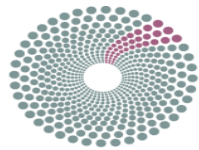# BIG DATA EUROPE

Empowering Communities
with Data Technologies

# Collaboration Opportunities in Big Data Platform, Data Spaces & Semantic Interoperability

Sören Auer, auer@cs.uni-bonn.de

PICASSO EU-US Collaboration, Minneapolis, June 19, 2017

# PICASSO ICT MEETING

Three possible streams for collaboration:

◎ A Big Data Platform for societal good

◎ Establishing data sharing and data value chains with the Industrial Data Space

◎ Semantic Domain Models (vocabularies, ontologies) for establishing a common understanding of the data

# Big Data Europe Platform

Empowering Communities with Data Technologies
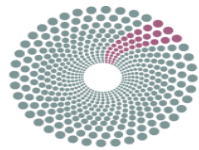Platform release

**BIG DATA EUROPE**

Empowering Communities
with Data Technologies

Smart Data Analytics

June 19 2017

# Platform Goals

- ◎ Opensource
- ◎ Ease of Use
- ◎ Support a variety of use cases
- ◎ Embrace emerging Big Data Technologies
- ◎ Simple integration with custom components

# Infrastructure

## Hadoop On-Premise
cloudera, Hortonworks, MAPR, Pivotal, IBM InfoSphere, splice, bluedata, jethro

## Hadoop in the Cloud
amazon web services, Microsoft Azure, Google Cloud Platform, IBM InfoSphere, CAZENA, altiscale, Qubole, xplenty

## Spark
databricks, GridGain, TACHYON NEXUS

## Cluster Services
amazon web services, docker, MESOSPHERE, Core OS, StackIQ

## NoSQL Databases
amazon DynamoDB, Google Cloud Platform, ORACLE, Microsoft Azure, MarkLogic, DATASTAX, mongoDB, AEROSPIKE, Couchbase, SequoiaDB, redislabs, influxdata

## NewSQL Databases
SAP, Clustrix, Pivotal, paradigm4, memsql, nuoDB, MariaDB, VOLTDB, citusdata, deepdb, Trafodion, Cockroach Labs

## Graph Databases
neo4j, VERTICA, OrientDB, kognitio, dremio

## MPP Databases
TERADATA, VERTICA, NETEZZA

## Cloud EDW
amazon web services, Google Cloud Platform, Microsoft Azure, Pivotal, Snowflake

## Data Transformation
alteryx, TRIFACTA, tamr, Paxata, Infoworks

## Data Integration
informatica, MuleSoft, snapLogic, BedrockData

## Management / Monitoring
New Relic, APPDYNAMICS, amazon web services, actifio, Numerify, splunk

## Security
TANIUM, illumio, CODE42, DataGravity, CipherCloud, VECTRA, BlueTalon

## Storage
amazon web services, Google Cloud Platform, Microsoft Azure, panasas, Qumulo

## App Dev
apigee, CASK, Typesafe, CONCURRENT

## Crowd-sourcing
CrowdFlower, WorkFusion

# Analytics

## Analyst Platforms
Palantir, AYASDI, Quid, enigma, Digital Reasoning

## Analytics Platforms
Microsoft, guavus, Datameer, inter/ana

## Data Science Platforms
CONTINUUM, DataRobot, Alpine, MODE, plotly, dataiku, DOMINO, sense, yhat

## Visualization
tableau, Google, Roambi, Qlik, CHARTIO

## BI Platforms
Power BI, amazon web services, DOMO, birst, GoodData, plotera, looker

## Statistical Computing
SAS, SPSS, MATLAB

## Log Analytics
splunk, sumologic, kibana, CLOUD PHYSICS, loggly

## Social Analytics
NETBASE, DATASIFT, trackx, bitly, synthesio, bottlenose, simplereach

## Real-Time
amazon web services, METAMARKETS, confluent, data Artisans

## Machine Learning
H2O.ai, DataRobot, SKYTREE, PredictionIO

## Speech & NLP
NarrativeScience, aplai, NUANCE, semantria, MindMeld, iDIBON

## Horizontal AI
IBM Watson, Cortana, VIV, nara, MetaMind, clarifai

## Search
ORACLE, HP, EXA, Lucidworks, elastic, MAANA, swiftype, Algolia

## Data Services
OPERA, DATA SCIENCE, kaggle, DataKind

## For Business Analysts
OrigamiLogic, ClearStory, CIRRO, import io

## SMB / Commerce
Google Analytics, RJMetrics, BLUECORE, sumall, granify, Airtable, retention, custora

# Applications

## Sales & Marketing
RADIUS, Gainsight, bloomreach, Zeta, livefyre, blue yonder, kahuna, Lattice, SAILTHRU, persado, infer, sense, AVISO, ACTIONIQ, QUANTIFIND, ENGAGIO

## Customer Service
MEDALLIA, ATTENSITY, CLARABRIDGE, STELLAService, Preact, NGDATA, Wise.io, epinion, Digital Genius, fuse machines

## Human Capital
gild, Connectifier, textio, entelo, hiQ

## Legal
RAVEL, JUDICATA, Everlaw, eBrevia

## Ad Optimization
MediaMath, Integral, OpenX, rocket fuel, theTradeDesk, LiveIntent, dstillery, Data.xu, Dapper, TAPAD

## Security
CYLANCE, CounterTack, cyberreason, ThreatMetrix, AREA 1, Guardian Analytics, Recorded Future, FORTSCALE, sift science, feedzai, SIGNIFYD

## Vertical AI Applications
facebook, X.ai, Clara, KASISTO, lumata

## Publisher Tools
outbrain, mixpanel, Chartbeat, yieldbot, Yieldmo

## Govt / Regulation
Socrata, OPENGOV, FiscalNote, enigma, mark43, OpenDataSoft

## Finance
Affirm, LendingClub, OnDeck, Kreditech, ZestFinance, LendUp, Kabbage, tidemark, INSIKT, Zuora, Dataminr, Lenddo, KENSHO, AIDYIA, iSENTIUM, Quantopian, sentient

## Education / Learning
Knewton, Clever, declara, PANORAMA, knowre

## Life Sciences
23andMe, Counsyl, Recombine, KYRUUS, FLATIRON, zymergen, HealthTap, METABIOTA, ZEPHYR HEALTH, Ginger.io, AiCure, Glow

## Industries
OPOWER, eHarmony, RetailNext, STITCH FIX, duetto, WorkFusion, BLUERIVER, TACHYUS, Seeq, FarmLogs, HowGood, celect, stat

# Cross-Infrastructure/Analytics
amazon web services, Google, Microsoft, IBM, SAP, SAS, HP, vmware, talend, TIBCO, TERADATA, ORACLE, NetApp

# Open Source

## Framework
Hadoop, YARN, Spark, MESOS, TEZ, Flink, CDAP

## Query / Data Flow
HIVE, PIG, SLAMDATA, DRILL, Google Cloud Dataflow

## Data Access
cassandra, HBASE, mongoDB, accumulo, SciDB, kafka, CouchDB, riak, nifi

## Coordination
talend, Zookeeper, Apache

## Real-Time
STORM, Spark, APEX, Flink, TACHYON

## Stat Tools
R, Scala, SciPy

## Machine Learning
mllib, Apache SINGA, MADlib, Aerosolve, Caffe, CNTK, FeatureFu, VELES, DIMSUM, DL4J

## Search
elasticsearch, Solr

## Security
Apache Ranger, Visualization

# Data Sources & APIs

## Health
Apple, JAWBONE, GARMIN, practice fusion, fitbit, Withings, VALIDIC, netatmo, kinsa, Human API

## IOT
UPTAKE, ThingWorx, helium, samsara

## Financial & Economic Data
Bloomberg, DOW JONES, S&P CAPITAL IQ, YODLEE, PREMISE, quandl, xignite, CB INSIGHTS, estimize, PLAID

## Air / Space / Sea
spire, PLANET LABS, WINDWARD, Airware, DroneDeploy

## Location/People/Entities
GARMIN, foursquare, InsideView, esri, STREETLINE, factual, Place IQ, placemeter, BASIS

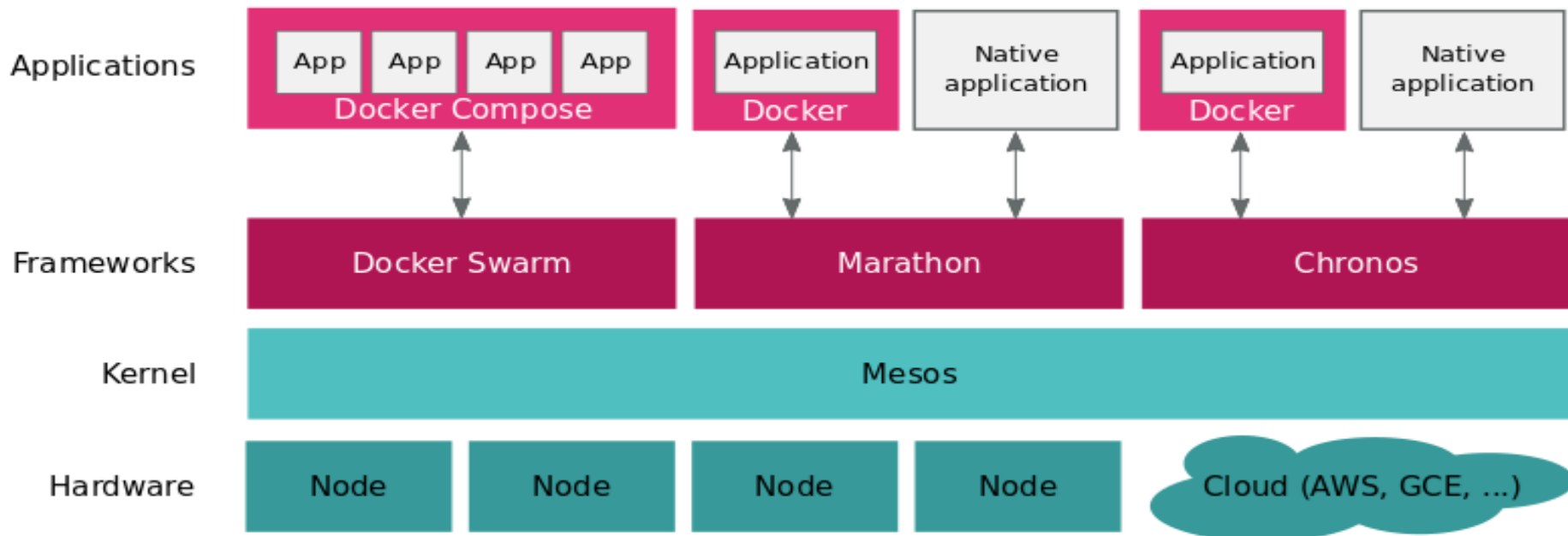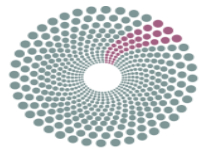## Other
qualtrics, panjiva, DATA.GOV

## Incubators & Schools
GA, DataCamp, INSIGHT, DataElite, The Data Incubator

© Matt Turck (@mattturck), Jim Hao (@jimrhao), & FirstMark Capital (@firstmarkcap)
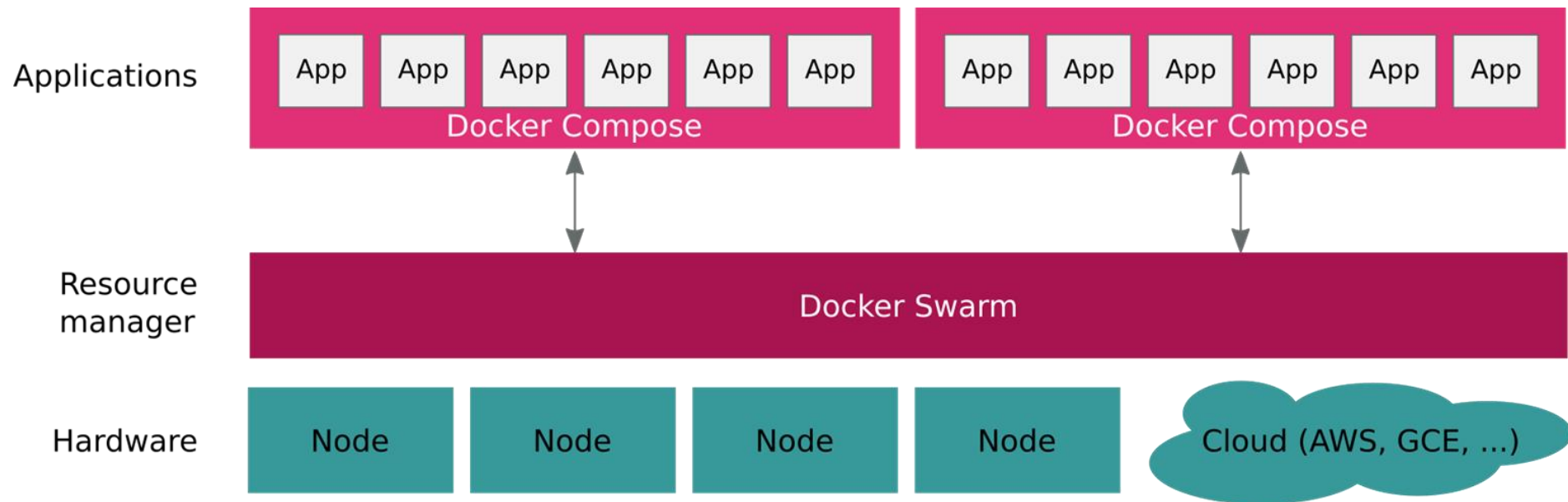
FIRSTMARK

# Platform Architecture Evolution

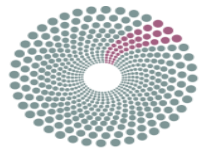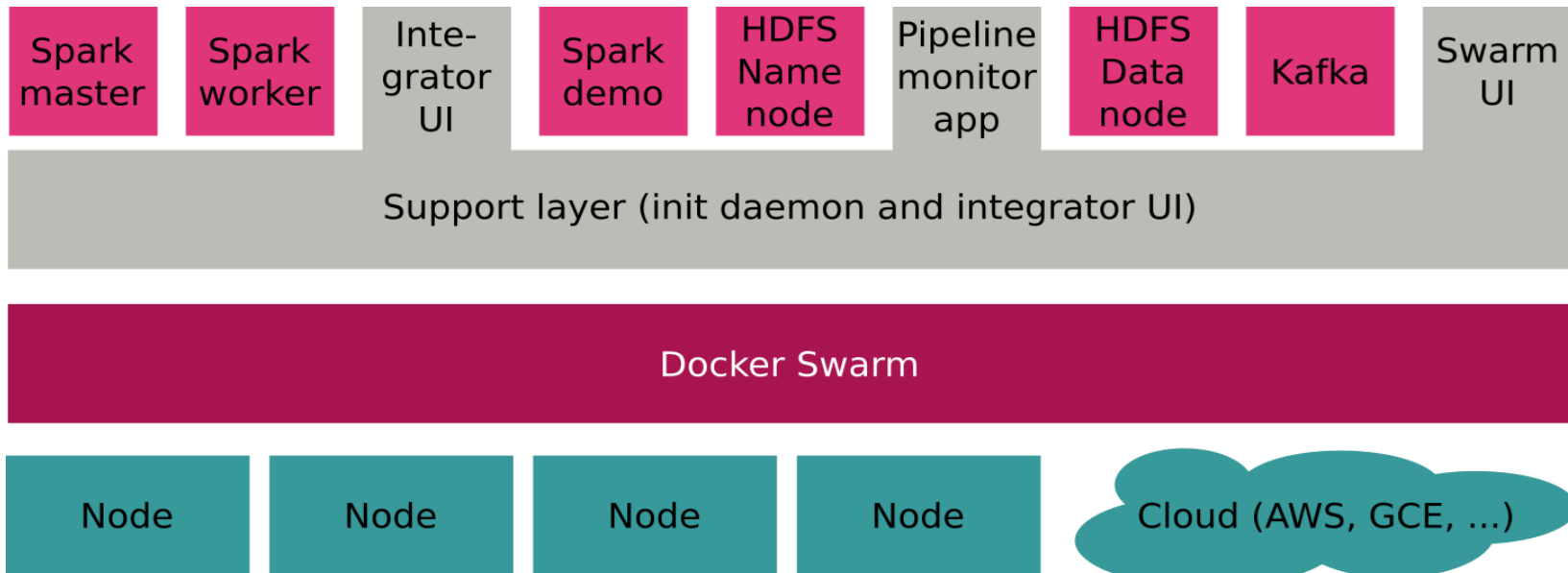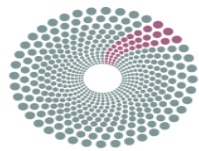# Platform Architecture Evolution

# Platform Architecture Existing

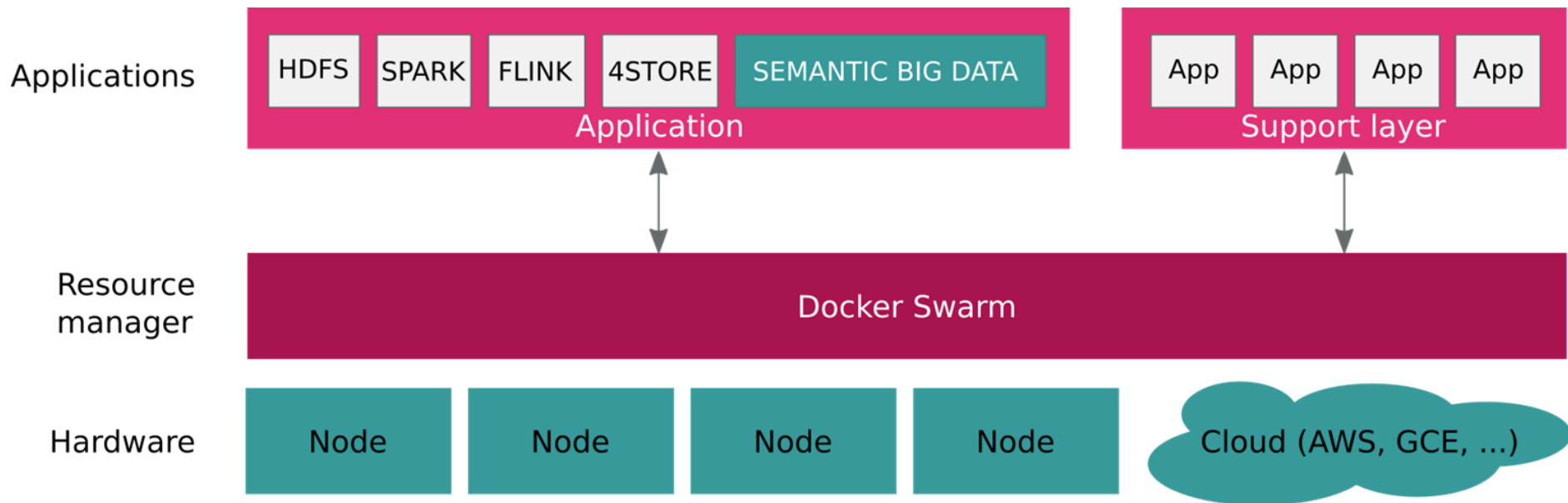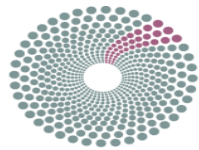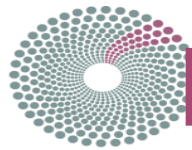| Spark master | Spark worker | Inte-grator UI | Spark demo | HDFS Name node | Pipeline monitor app | HDFS Data node | Kafka | Swarm UI |

**Support layer (init daemon and integrator UI)**

**Docker Swarm**

| Node | Node | Node | Node | Cloud (AWS, GCE, …) |

# Platform Architecture Alternate View

**Support Layer**

Init Daemon

GUIs

Monitor

**App Layer**

Traffic Forecast

Satellite Image Analysis

Real-time Stream Monitoring

...

**Platform Layer**

Spark

Flink

Kafka

...

**Semantic Layer**

Ontario

SANSA

Semagrow

**Data Layer**

Hadoop

RDF Store

NOSQL Store

Elasticsearch

Cassandra

...

**Resource Management Layer (Swarm)**

**Hardware Layer**

Premises

Cloud (AWS, GCE, MS Azure, …)

# BDE Supported Frameworks

| Search/indexing | Data processing |
|---|---|
| Apache Solr | Apache Spark |
| **Data acquisition** | Apache Flink |
| Apache Flume | **Semantic Components** |
| **Message passing** | Strabon |
| Apache Kafka | Sextant |
| **Data storage** | GeoTriples |
| Hue | Silk |
| Apache Cassandra | SEMAGROW |
| ScyllaDB | LIMES |
| Apache Hive | 4Store |
| Postgis | OpenLink Virtuoso |

# Platform features

◎ BDE Development Environment
- o Stack builder
- o Workflow builder
- o Instructions to add custom components to the BDE stack

◎ Administrator Interface
- o SwarmUI

◎ UI Integrator
- o Workflow monitor
- o Integrated web interface

# What BDE Provides ?

◎Platform Installation Instructions

◎Usage Instructions

- Creating a stack
- Creating a workflow
- Monitoring the Stack
- Integration of Custom Components

# Platform installation

Manual installation guide

Using Docker Machine

- On local machine (VirtualBox)
- In cloud (AWS, DigitalOcean, Azure)
- Bare metal

Screencasts

https://www.big-data-europe.eu/platform/

https://github.com/big-data-europe

# Deploying a Big Data Stack

◎ Stack Builder

◎ Stack
- Collection of communicating components to solve a specific problem

◎ Described in Docker Compose
- Component configuration
- Application topology

# Creation of WorkFlows

◎Pipeline Builder
- Allows creation of dependencies among different applications

◎WorkFlow Monitor
- Monitoring of pipeline-workflow using

# Integrating Custom Components

◎ Instructions

- o Orchestrator required for initialization process (init_daemon)

  - ❖ Components may depend on each other

  - ❖ Components may require manual intervention

- o User Interface Integration

  - ❖ Standard Interfaces from components

  - ❖ Combine and align the interfaces

# User Interfaces

◎ Target: Facilitate the use of the platform

- ○ User Interface Adaption

◎ Available interfaces

- ○ Workflow UIs
  - ❖ Workflow Builder
  - ❖ Workflow Monitor
- ○ Swarm UI
- ○ Integrator UI

# Platform Architecture



Applications

| HDFS | SPARK | FLINK | 4STORE | SEMANTIC BIG DATA |

Application

| App | App | App | App |

Support layer

Resource manager

Docker Swarm

Hardware

| Node | Node | Node | Node | Cloud (AWS, GCE, …) |

# Pilot Show Cases

Health — SC1

SC7

Food & Agriculture — SC2

Energy — SC3

Transport — SC4

Climate — SC5

Social Sciences — SC6

Security

SC1 - Open PHACTS discovery platform relating to biological/medical questions

SC2 - Discovery and Linking of Viticulture-relevant information

SC3 - System monitoring in energy production units

SC4 - Short-Term traffic flow forecasting.

SC5 - Supporting data-intensive climate research

SC6 - Citizens & Researchers Budget on Municipal Level

SC7 - Ingestion of remote sensing images and social sensing data to detect and verify changes on the Earth surface for security applications
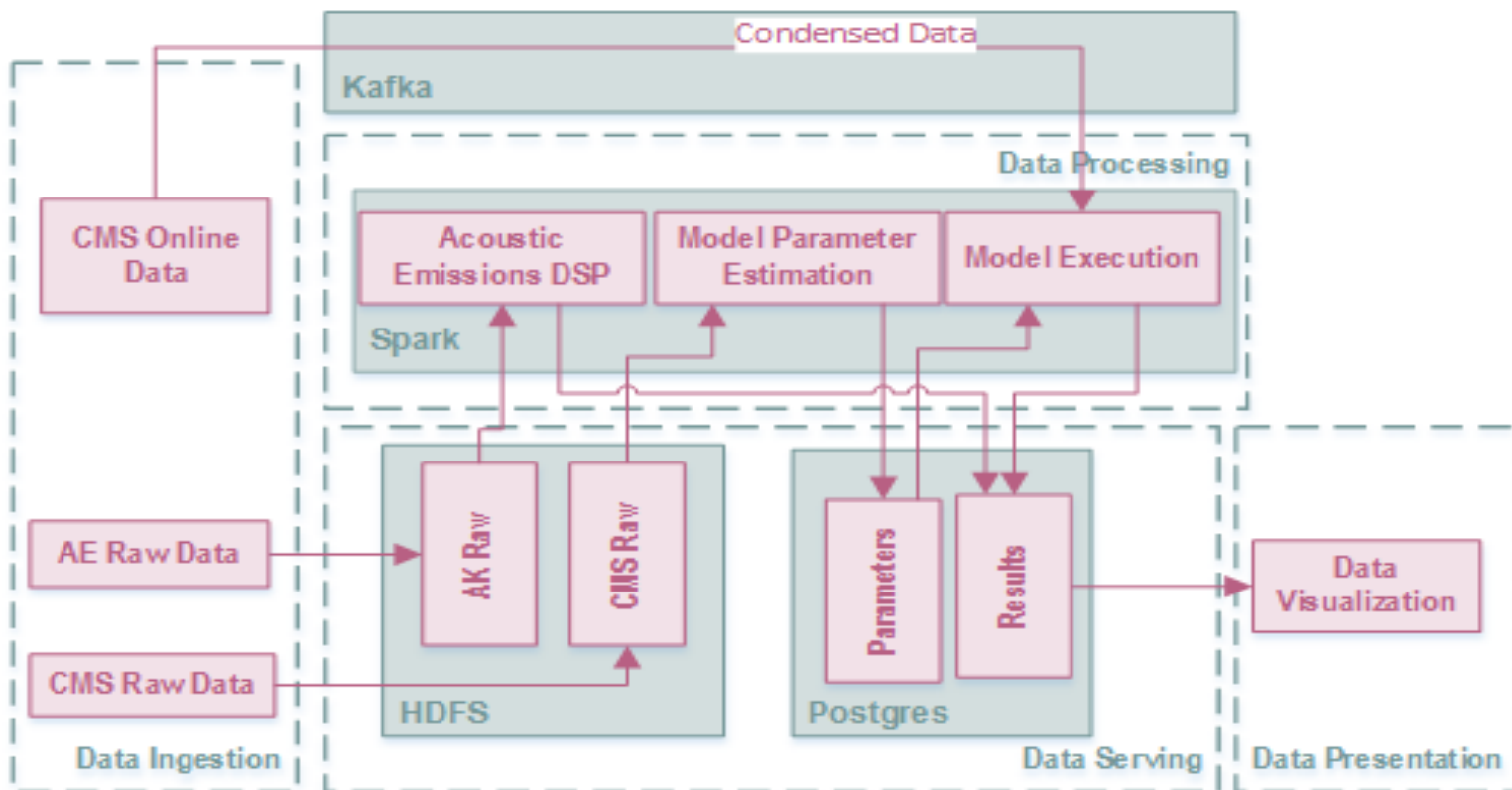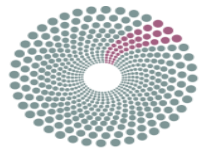
# SC1- Health

# SC2 - Food

# SC3 - Energy

# SC4 - Transport



FCD: Floating Car Data
NRT: Near Real Time

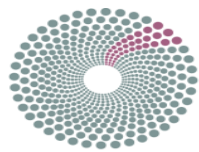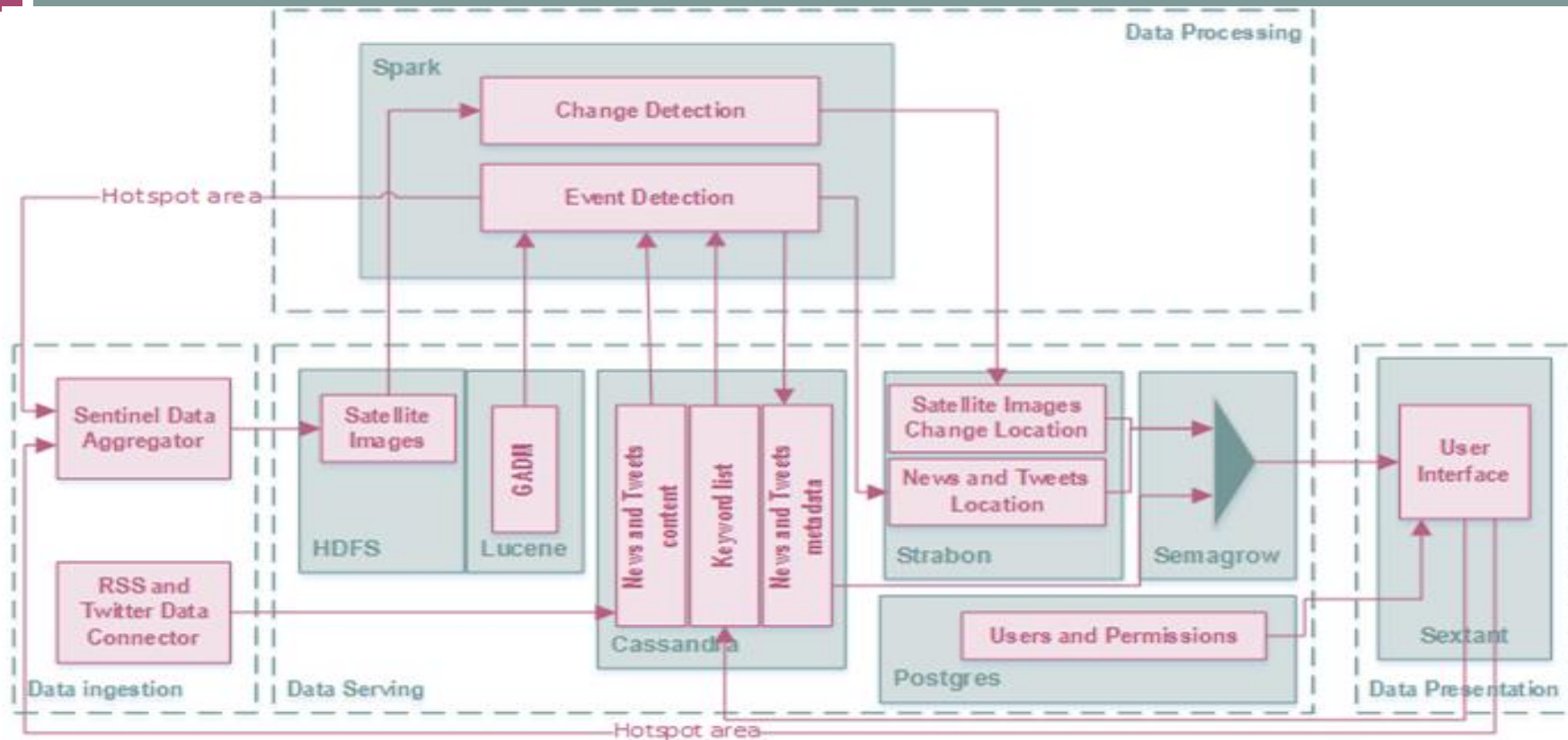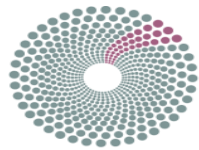# SC6 - Social Sciences

# SC7 - Security

# BDE vs Hadoop distributions

|  | Hortonworks | Cloudera | MapR | Bigtop | BDE |
|---|---|---|---|---|---|
| *File System* | HDFS | HDFS | NFS | HDFS | HDFS |
| *Installation* | Native | Native | Native | Native | lightweight virtualization |
| *Flexible Modular Architecture* | no | no | no | no | yes |
| *High Availability* | Single failure recovery (yarn) | Single failure recovery (yarn) | Self healing, mult. failure rec. | Single failure recovery (yarn) | 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 | Windows/Mac/Linux |
| *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

# Wrap Up, thanks for your Attention

Three possible streams for collaboration:

◎ A Big Data Platform for societal good

◎ Establishing data sharing and data value chains with the Industrial Data Space

◎ Semantic Domain Models (vocabularies, ontologies) for establishing a common understanding of the data

Please get in touch: Sören Auer (coordinator Big Data Europe), auer@cs.uni-bonn.de