

Deciphering Big Data Stacks: An Overview of Big Data Tools

Tomislav Lipic¹, Karolj Skala¹, Enis Afgan^{1,2}

¹Ruđer Bošković Institute (RBI)

²Johns Hopkins University

Big Data Analytics: Challenges and Opportunities (BDAC-14)

New Orleans, LA

Nov 2014



-
1. Tool catalog
 2. Functional tool comparison
 3. Tool deployment dependency graph

Existing Technologies as Trends

Batch processing

MapReduce

- Hadoop
- Cloudera
- Hortonworks
- MapR
- AWS EMR

HDFS

- Replicated data store

Query processing

SQL-like interface

- Pig
- Hive

NoSQL databases

- Scalable storage
- Query interfaces
- Document, Key-value, Columnar, Graph

Low-latency processing

Near real-time query

- Impala
- Drill
- Dremel

Data store

- HDFS
- NoSQL

Continuous processing

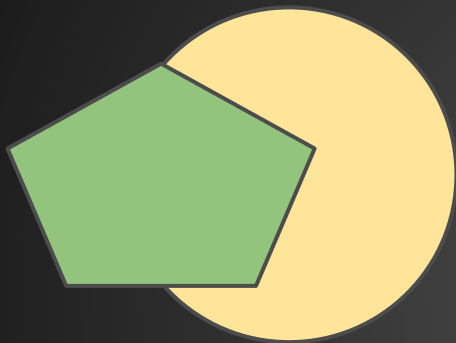
Stream processing

- Storm
- S4
- Samza
- AWS Kinesis

Data delivery

- Kafka
- Flume

Crossing the Trends

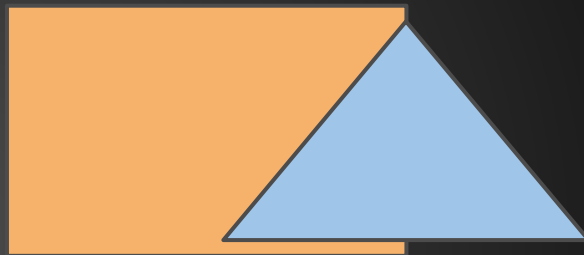


Spark

- Low-latency queries
- Data streams

Stratosphere Nephel

- Additional processing operators: join, union



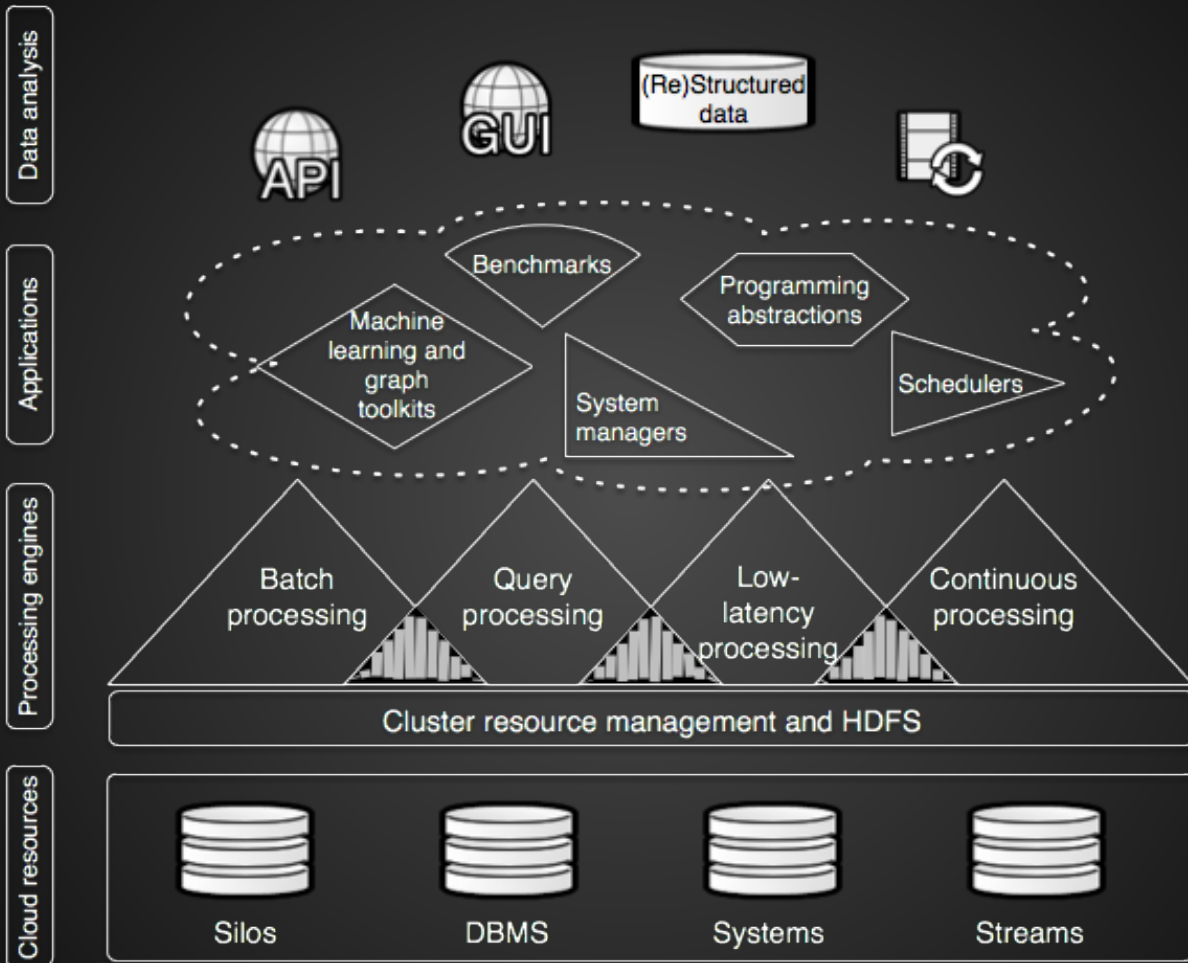
Meta-resource managers

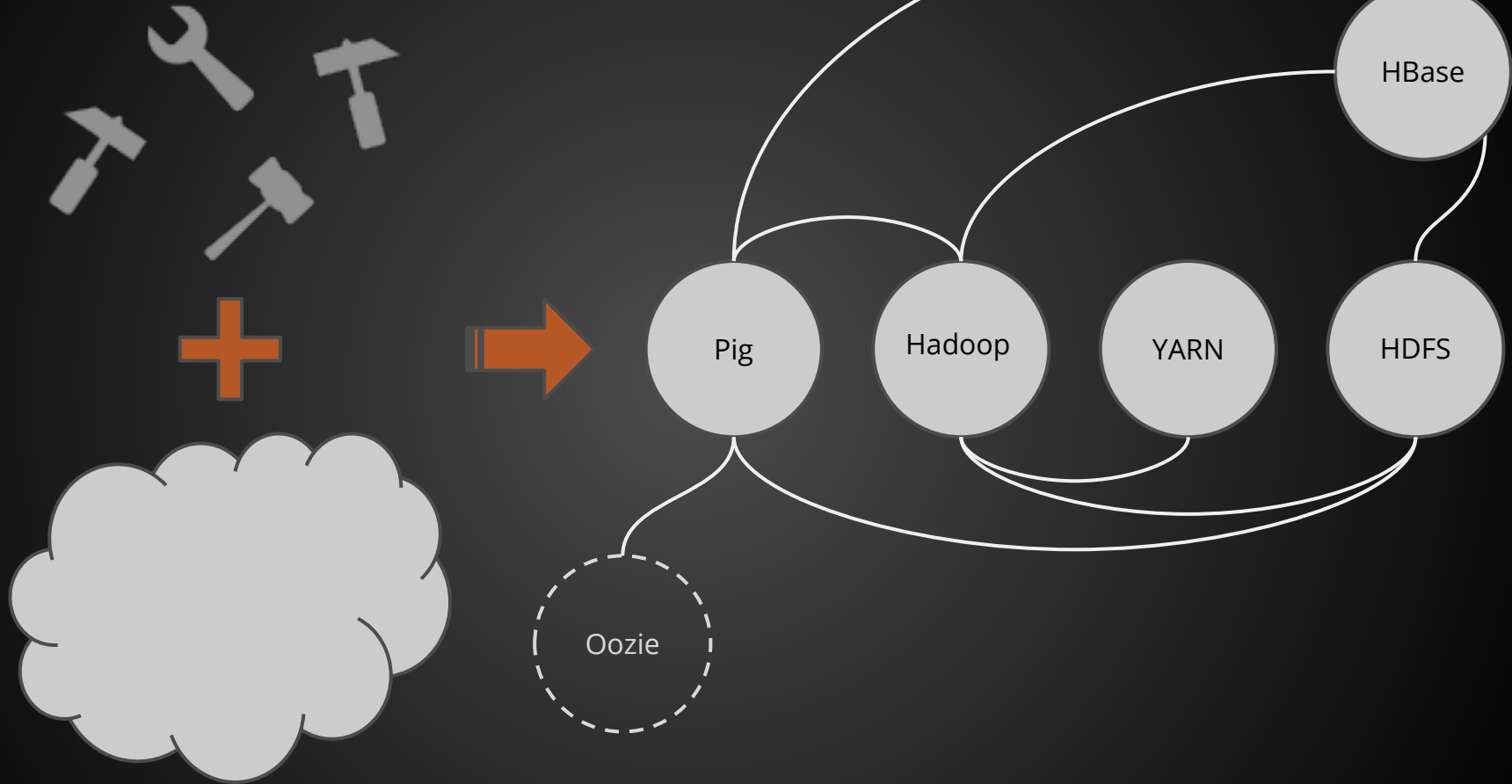
- YARN
- Mesos

One cluster per datacenter

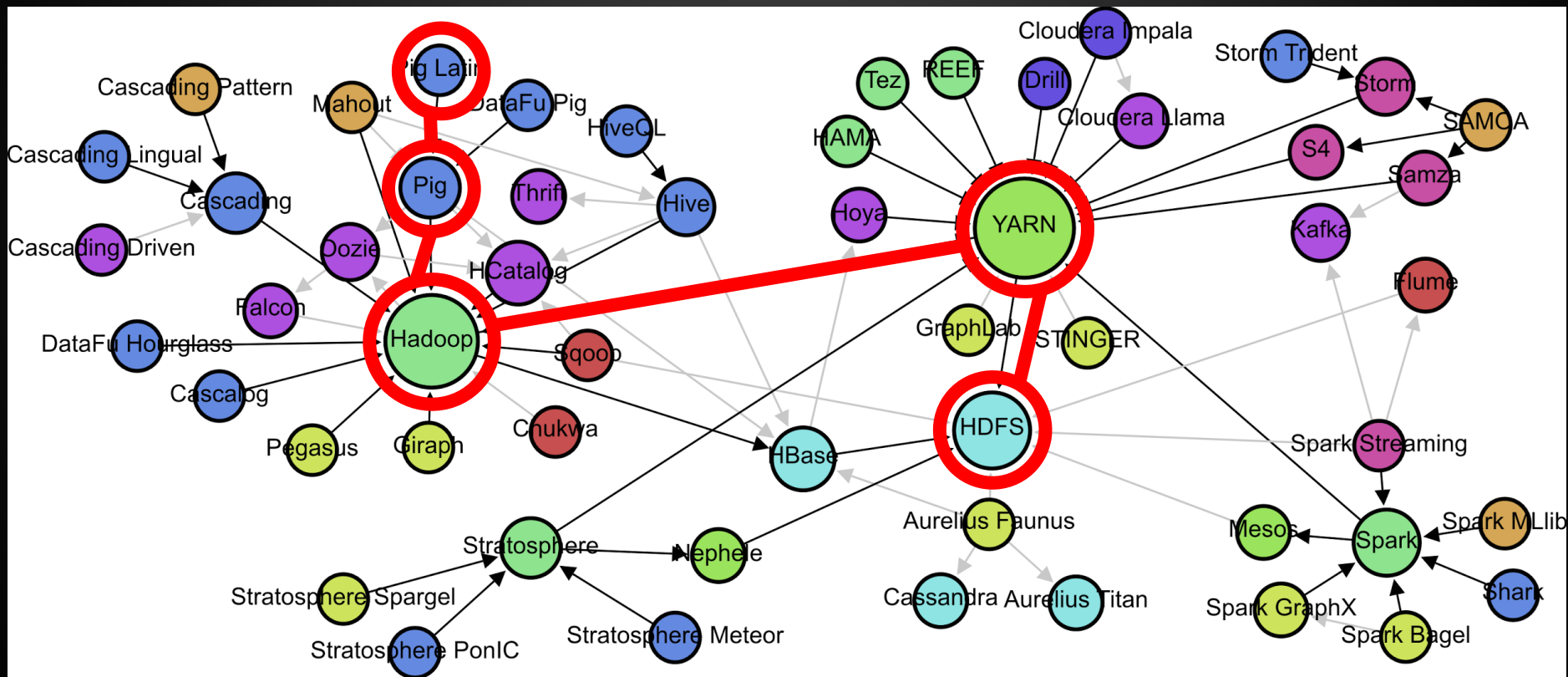
- On-demand resource provisioning

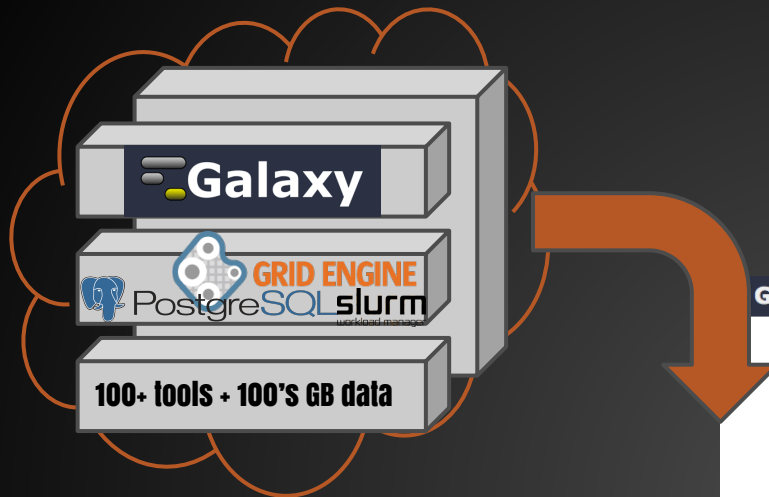
Functional classification	Existing Big Data Technologies
<i>High-level programming abstractions</i>	Apache Pig, Apache DataFu, Cascading Lingual, Cascalog
	Shark (for Spark), Trident (for Storm), Meteor, Sopremo and PonIC (for Stratosphere)
<i>Big Data-aware machine-learning toolkits</i>	Apache Mahout (for Hadoop), MLlib (for Spark), Cascading Pattern, GraphLab framework, Yahoo SAMOA
<i>Graph processing systems</i>	Apache Giraph (for Hadoop), Bagel (for Spark), Stratosphere Spargel, GraphX (for Spark), Pegasus, Aurelius Faunus, GraphLab PowerGraph
	Stinger, Neo4j, Aurelius Titan
<i>Data ingestion and scheduling systems</i>	Apache Sqoop, Apache Chukwa, Apache Flume
	Apache Falcon, Apache Oozie
<i>Systems management solutions</i>	Apache Hue
	Apache Ambari, Apache Helix, Apache Whirr, Cask Coopr
<i>Benchmarking and testing applications</i>	Berkeley Big Data benchmark, BigBench, BigDataBench,, Big Data Top 100, Apache Bigtop





Tool Interdependency Graph





Galaxy Info: [report bugs](#) | [wiki](#) | [screenshots](#)

Galaxy CloudMan Console

Welcome to the Galaxy Cloud Manager. This application will allow you to manage this cloud and the services provided within. If this is your first time running this cluster, you will need to select an initial data volume size. Once the data store is configured, default services will start and you will be add and remove additional services as well as 'worker' nodes on which jobs are run.

[Terminate cluster](#) [Add instances ▼](#) [Remove instances ▼](#) [Access Galaxy](#)

Status

Cluster name: cm_dev1		i-b6d427db
Disk status: 49M / 1014M (5%)		State: Ready
Worker status: Idle: 12 Available: 12 Requested: 12		Alive: 23m 8s
Service status: Applications Data		Filesystems Permissions Scheduler

[Cluster status log](#)



CloudMan

Questions?

Slides at slideshare.net/afgane

Paper at bit.ly/BDTools

Support from Ruđer Bošković Institute and Johns Hopkins University

