



Big Data zum Anfassen

Spark, Hive, Kafka & Co.

Wolfgang Pleus – **PLEUS** Consulting
www.pleus.net



IT Solution Architect , Entwickler, Agile Coach

> 20 Jahre Projekterfahrung

> 10 Jahre agile Projekte

DAX - Startup

Wolfgang Pleus – **PLEUS** Consulting

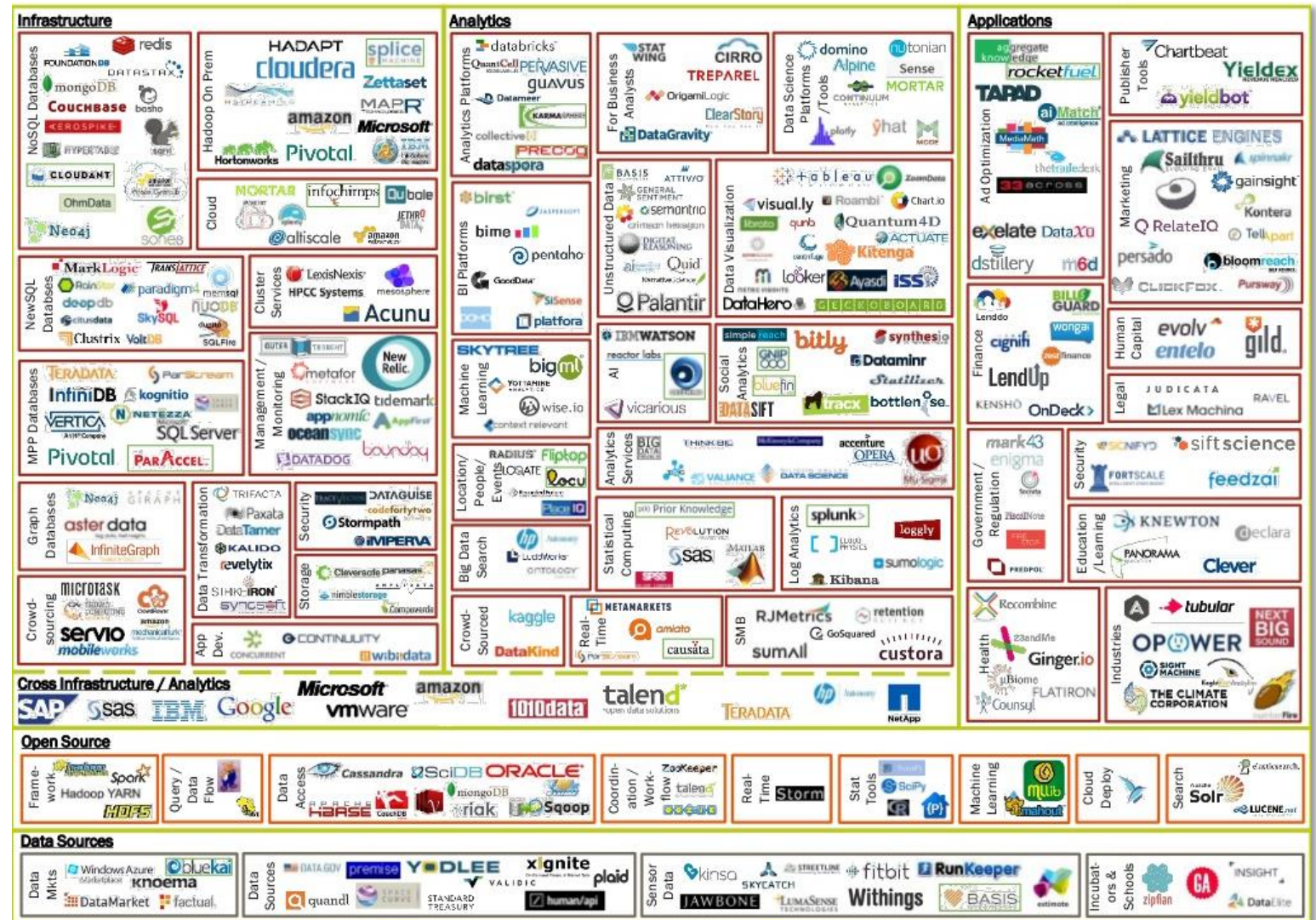
www.pleus.net

Agenda

- Big Data Stack
- Small Data
- Lambda Architektur
- Lauffähiger Prototyp
- Grüne Wiese on Premise



Big Data Ökosystem



Distribution

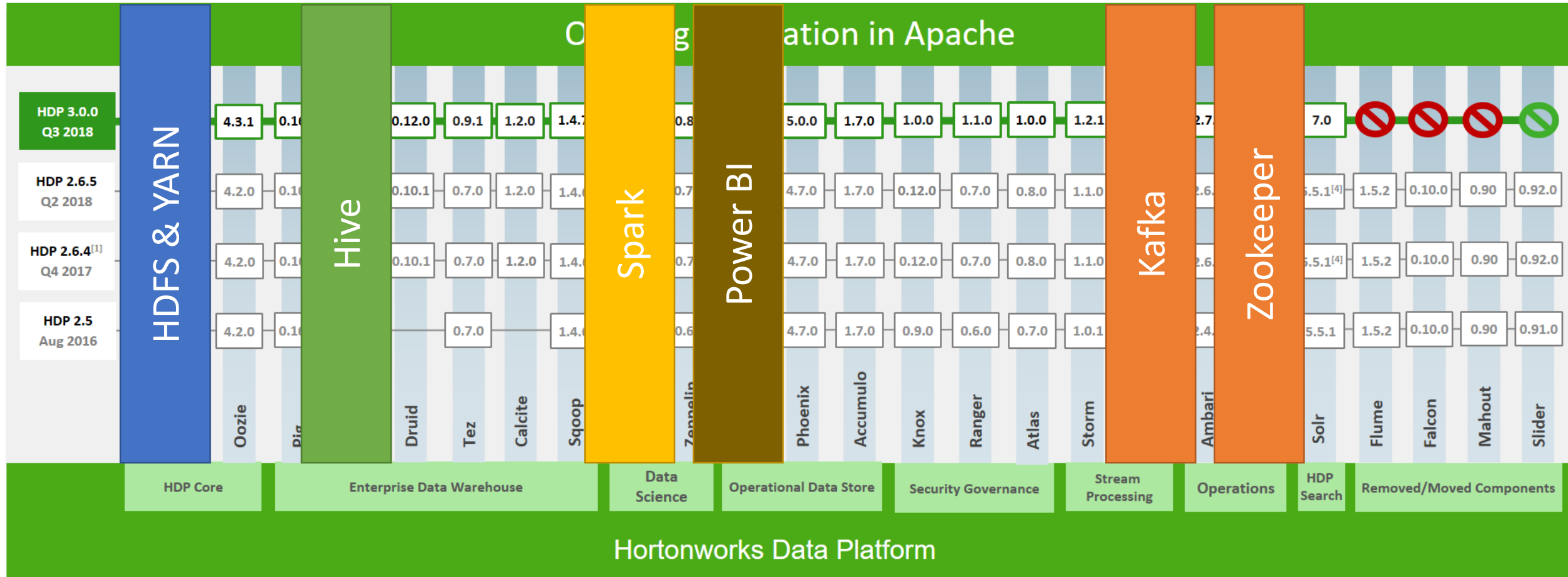
Ongoing Innovation in Apache

	Hadoop & YARN	Oozie	Pig	Hive	Druid	Tez	Calcite	Sqoop	Spark	Zeppelin	HBase	Phoenix	Accumulo	Knox	Ranger	Atlas	Storm	Kafka	Ambari	Zookeeper	Solr	Flume	Falcon	Mahout	Slider	
HDP 3.0.0 Q3 2018	3.1.0	4.3.1	0.16.0	3.0.0	0.12.0	0.9.1	1.2.0	1.4.7	2.3	0.8.0	2.0.0	5.0.0	1.7.0	1.0.0	1.1.0	1.0.0	1.2.1	1.0.1	2.7.0	3.4.6	7.0					
HDP 2.6.5 Q2 2018	2.7.3	4.2.0	0.16.0	1.2.1+ 2.1 ^[3]	0.10.1	0.7.0	1.2.0	1.4.6	1.6.3+ 2.3	0.7.3	1.1.2	4.7.0	1.7.0	0.12.0	0.7.0	0.8.0	1.1.0	1.0.0	2.6.2	3.4.6	5.5.1 ^[4]	1.5.2	0.10.0	0.90	0.92.0	
HDP 2.6.4 ^[1] Q4 2017	2.7.3	4.2.0	0.16.0	1.2.1+ 2.1 ^[3]	0.10.1	0.7.0	1.2.0	1.4.6	1.6.3+ 2.2 ^[5]	0.7.3	1.1.2	4.7.0	1.7.0	0.12.0	0.7.0	0.8.0	1.1.0	0.10.1	2.6.1	3.4.6	5.5.1 ^[4]	1.5.2	0.10.0	0.90	0.92.0	
HDP 2.5 Aug 2016	2.7.3	4.2.0	0.16.0	1.2.1+ 2.1 ^[3]		0.7.0		1.4.6	1.6.2+ 2.0 ^[2]	0.6.0	1.1.2	4.7.0	1.7.0	0.9.0	0.6.0	0.7.0	1.0.1	0.10.0	2.4.0	3.4.6	5.5.1	1.5.2	0.10.0	0.90	0.91.0	
	HDP Core	Enterprise Data Warehouse				Data Science		Operational Data Store		Security Governance		Stream Processing		Operations		HDP Search	Removed/Moved Components									
Hortonworks Data Platform																										

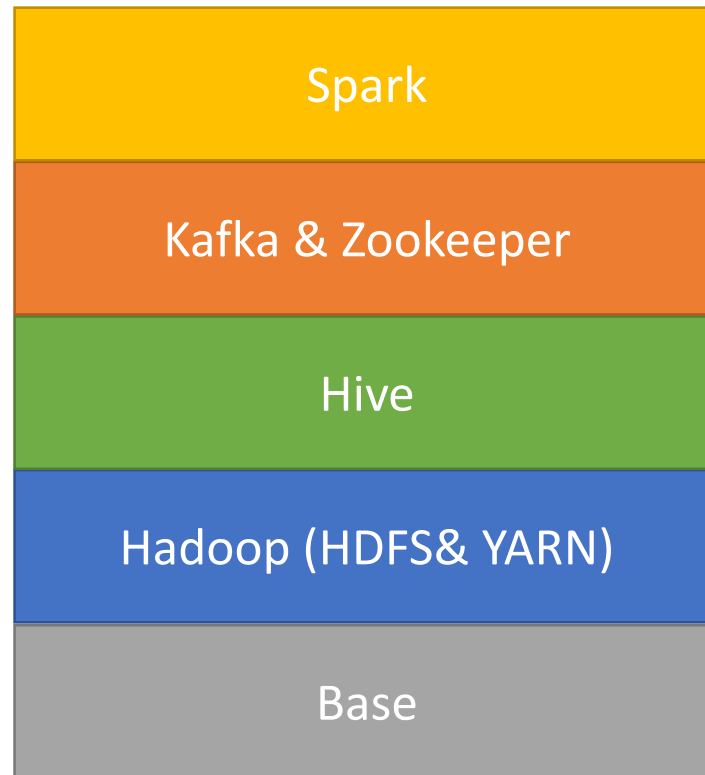
Hortonworks Data Platform



Minimales Setting

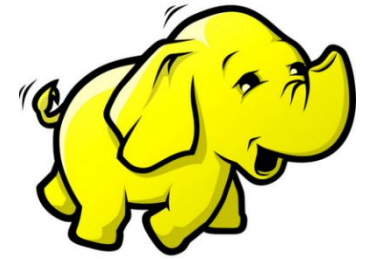


Docker Images



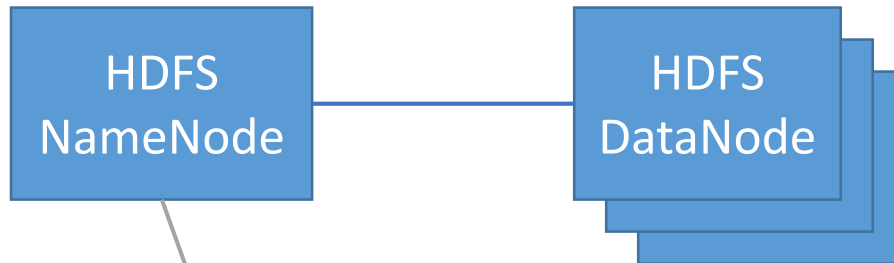
Demo: Basis

Hadoop: HDFS & YARN



- Skalierbarer Speicher, CPU und Memory

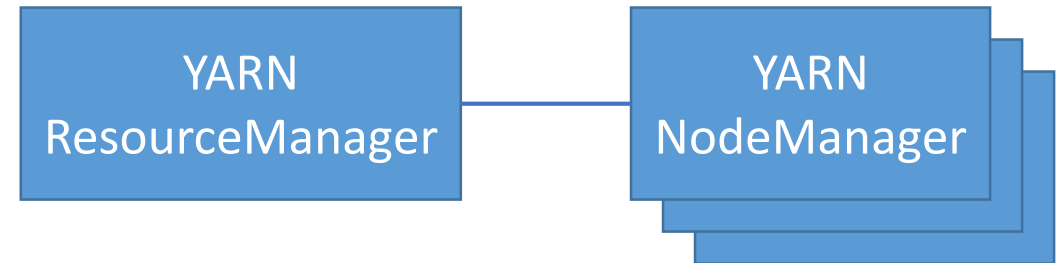
HDFS



Daten

```
hdfs dfs -put data.csv
```

YARN



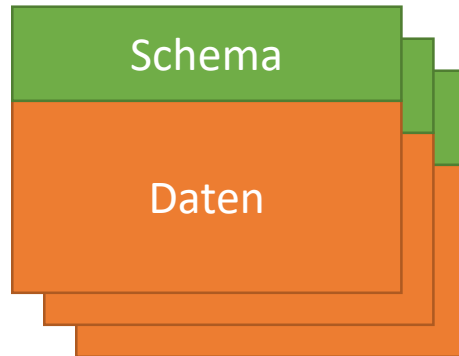
CPU + RAM

Demo: Hadoop

Avro



- Schema für Datasets
- Serialisierung
- Schemaevolution



jobs.avro

```
{ "namespace": "net.pleus.bddemo",  
  "name": "jobs",  
  "version": "1.0",  
  "doc" : "Jobs",  
  "type": "record",  
  "fields": [  
    { "name": "created", "type": "string", "default": "unknown", "doc": "Record creation date" },  
    { "name": "published", "type": ["string", "null"], "default": "unknown", "doc": "Published date" },  
    { "name": "source", "type": ["string", "null"], "default": "unknown", "doc": "Source of the job" },  
    { "name": "title", "type": ["string", "null"], "default": "unknown", "doc": "Name of the job" },  
    { "name": "type", "type": ["string", "null"], "default": "unknown", "doc": "Offer type" },  
    { "name": "link", "type": ["string", "null"], "default": "unknown", "doc": "Link to the job" },  
    { "name": "region", "type": ["string", "null"], "default": "unknown", "doc": "Work region" },  
    { "name": "geolocation", "type": ["string", "null"], "default": "unknown", "doc": "Geolocation" },  
    { "name": "longitude", "type": ["string", "null"], "default": "unknown", "doc": "Longitude" },  
    { "name": "latitude", "type": ["string", "null"], "default": "unknown", "doc": "Latitude" }  
  ]  
}
```

jobs.avsc

Sqoop

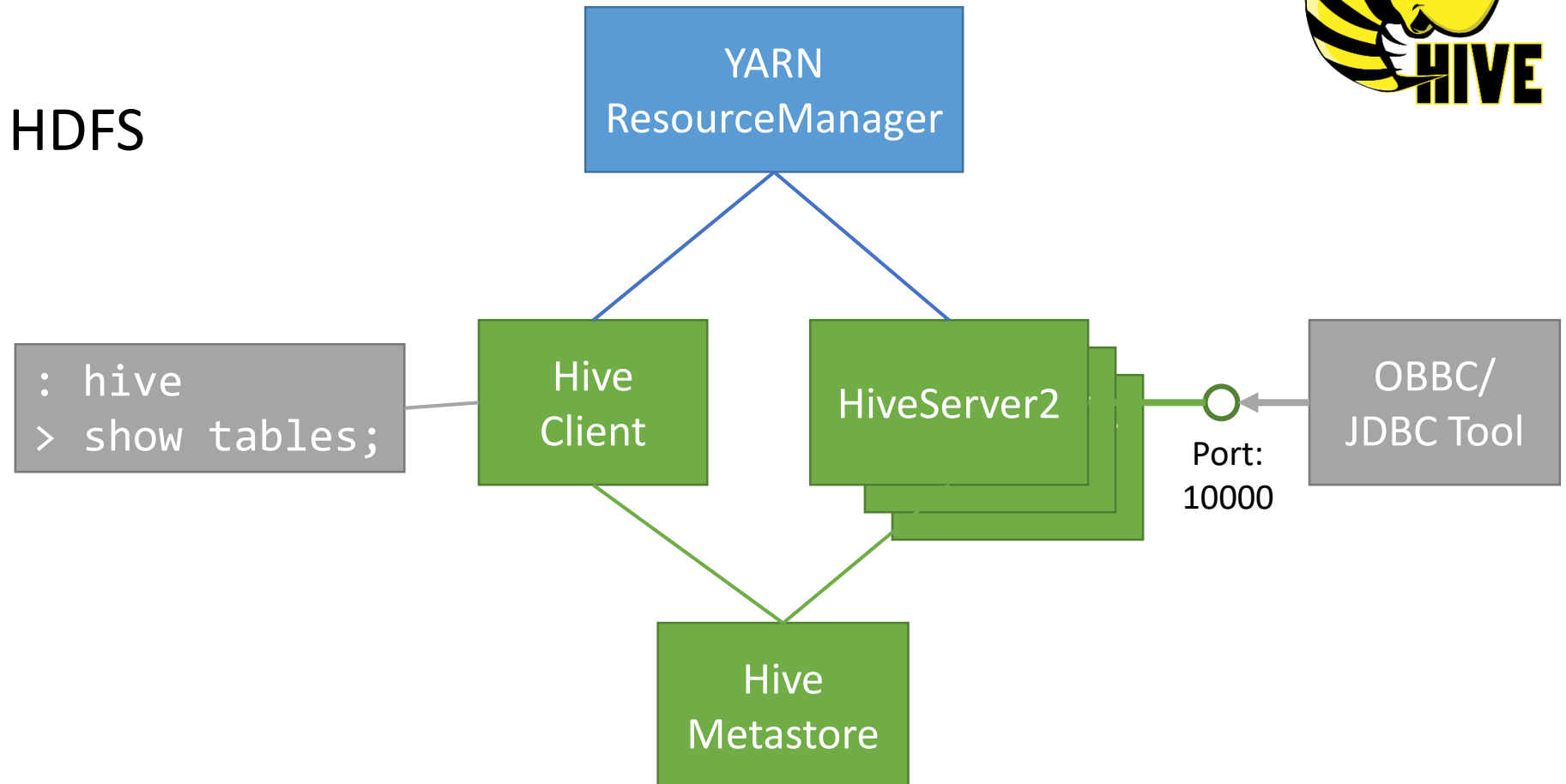


- SQL + Hadoop
- Import und Export relationaler Daten
- Paralleles Laden mit Map Reduce

```
sqoop import
--connect jdbc:oracle:thin:@server:1521:demo
--username wpl
--password=***
--create-hive-table
--table ORA.JOBS
-m 1
--hive-import
--hive-overwrite
--hive-table jobs
```

Hive

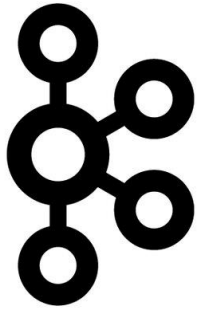
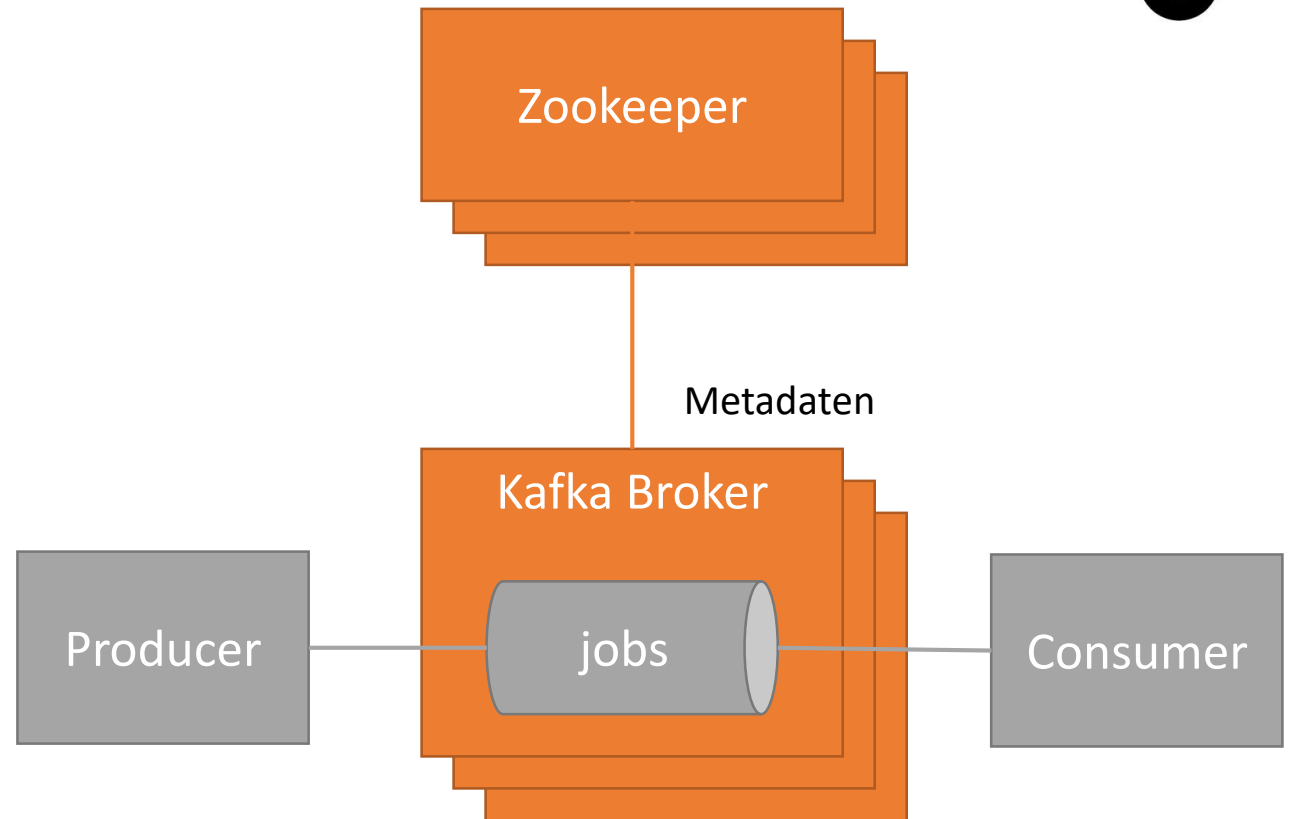
- SQL Layer für HDFS



Demo: Hive

Kafka

- Message Broker
- Streamverarbeitung
- Skalierbar

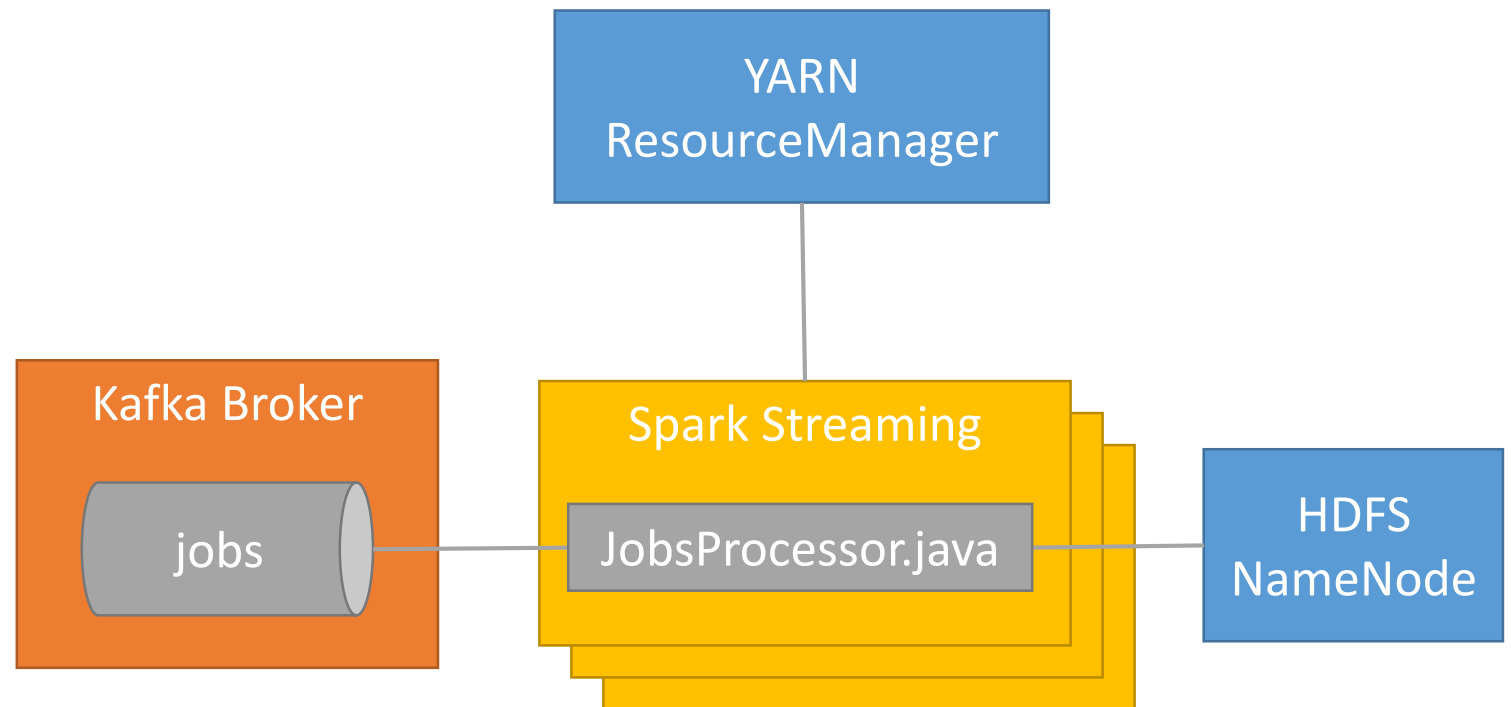


Demo: Kafka

Spark



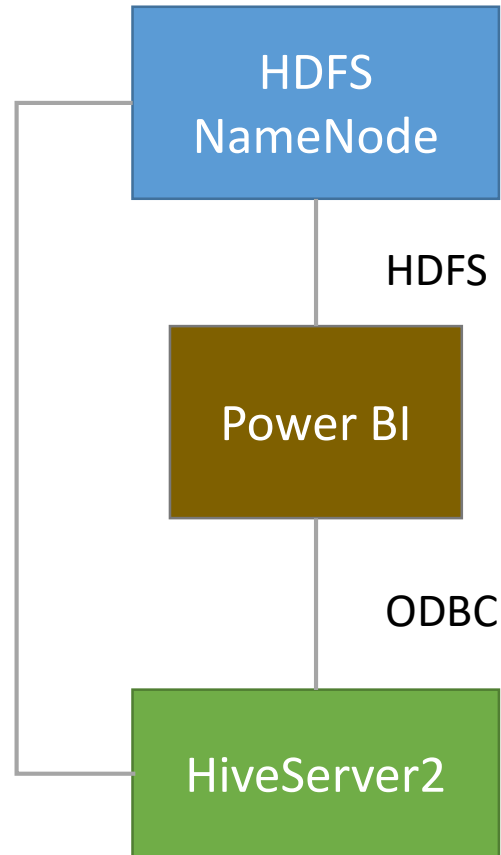
- Streaming
- SQL
- Transformationen
- Clusterfähig



Demo: Spark

Visualisierung

- Z.B. Power BI

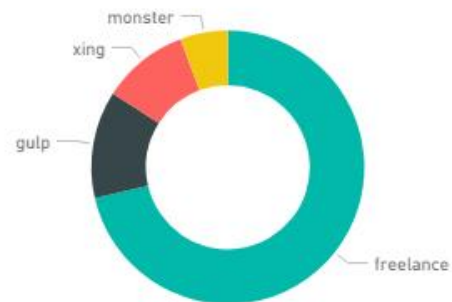


Demo: Visualisierung

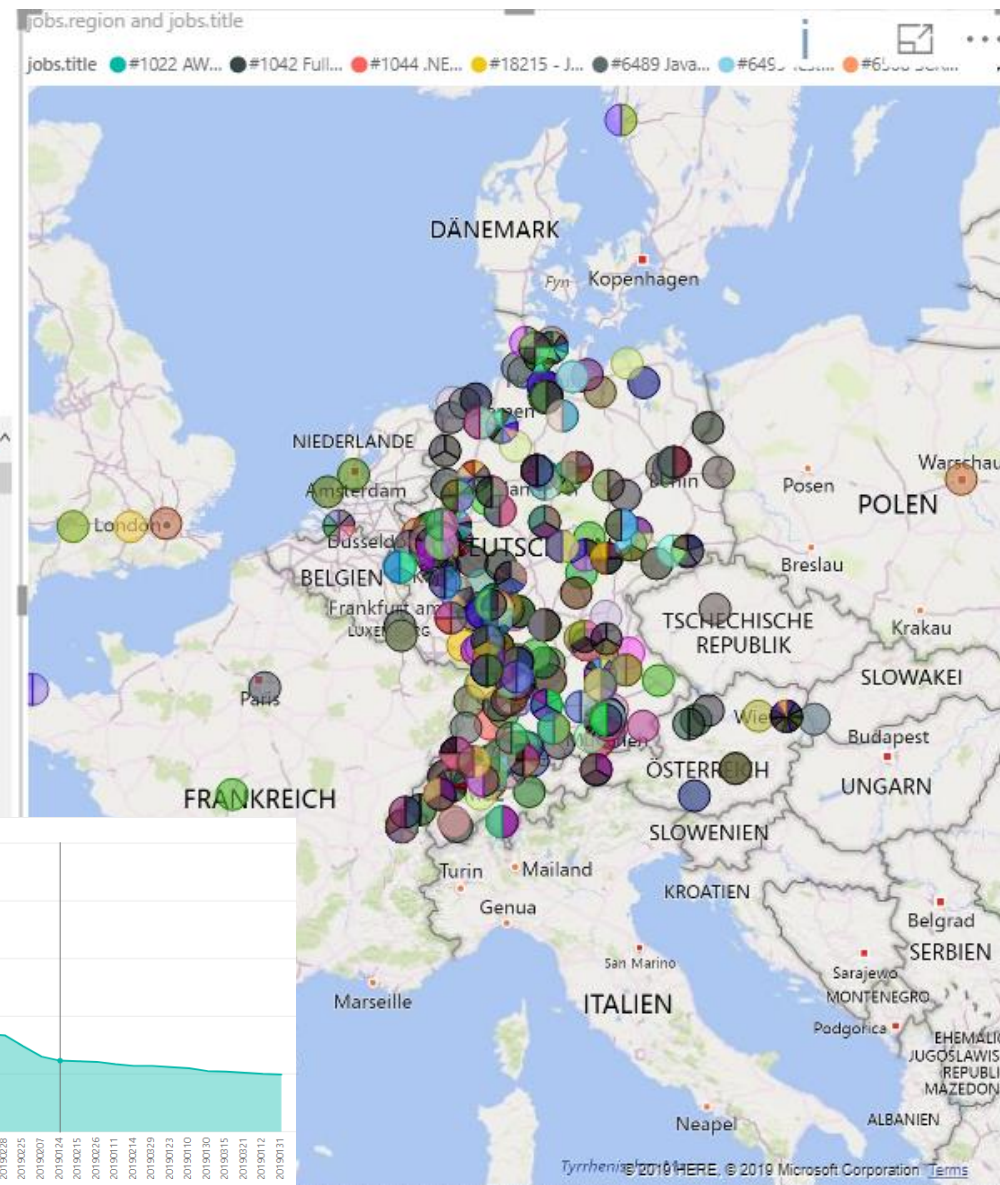
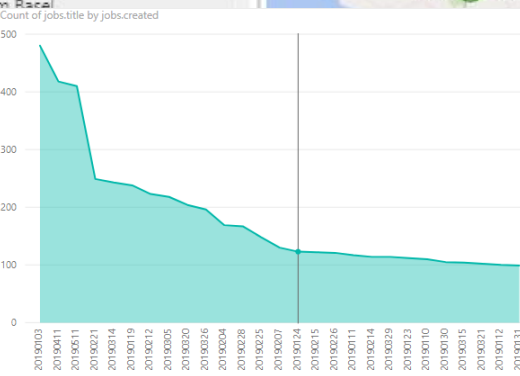
5005

Count of jobs.title

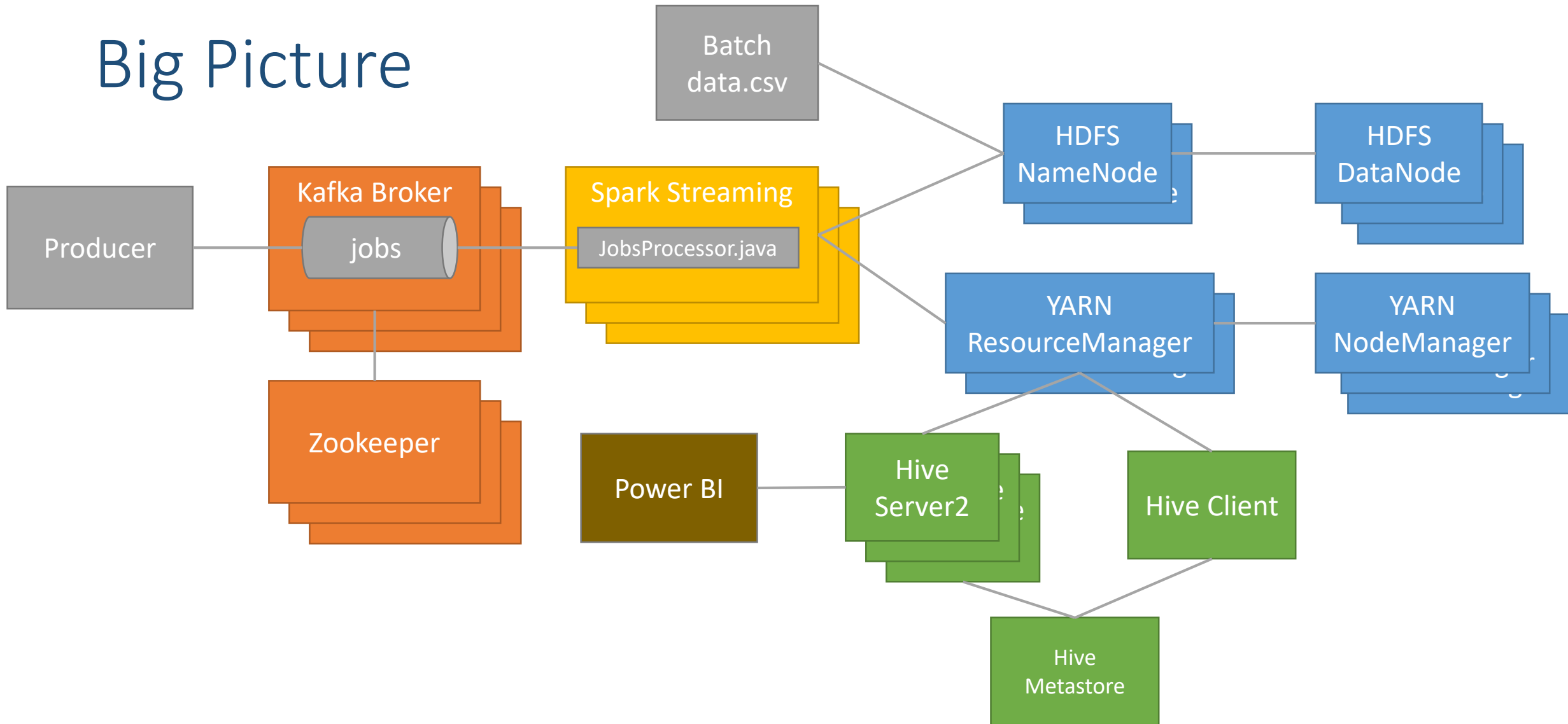
Count of jobs.source by jobs.source



jobs.title	jobs.created	jobs.region
ILC RPG - Entwickler (m/w/x)	20190103	1445 Luxembourg
Agile Full-Stack Developer (m/w)	20190103	A-Wien
Angular2+-Webentwickler (m/w) 393254/6	20190103	A-Wien
FREELANCER4VIENNA - Salesforce Vlocity Developer (f/m)	20190103	A-Wien
MS Dynamics Entwickler	20190103	A-Wien
Entwickler - webbasierte Anwendungen f��r eine Robo Advisor-Plattform (m/w)	20190103	Berlin
Java Entwicklung (w/m)	20190103	Bern
Project Manager in the Scientific and Manufacturing field	20190103	CH-Aarau
Agile .NET Developer 100% (m/w) 394208/19	20190103	CH-Basel
Python Developer / Business Analyst - Data Modelling/SQL Server	20190103	CH-Basel
Java/Angular Entwickler (m/w), Pos. 1756	20190103	CH-Bern
Network Engineer / Migration Expert (w/m)	20190103	CH-Bern
Professional Software Quality Specialist and Integration Developer (m/f) 393789/11	20190103	CH-Gro��raubach
Business Analyst (m/w)	20190103	CH-Gro��raubach
Senior Testautomatisierer (m/w), Pos. 1757	20190103	CH-Gro��raubach
Scrum Master / Scrum Coach (5866)	20190103	CH-Gro��raubach
DevOps Engineer (100%) - C#/Powershell/HyperV/SCVMM /Lucerne	20190103	CH-Gro��raubach
Applikations Architekten/Entwickler - Webservice (60%, 80% & 100%) (5961)	20190103	CH-Gro��raubach
Cobol Entwickler - zOS / DB2 (5946)	20190103	CH-Gro��raubach
IT Business Analyst - Agile (5935)	20190103	CH-Gro��raubach
IT Business Analyst / RE - Kollektivleben (5948)	20190103	CH-Gro��raubach
Java Backend Developer (5699)	20190103	CH-Gro��raubach
Java Entwickler (5712)	20190103	CH-Gro��raubach
Scrum Master (m/w)	20190103	CH-Gro��raubach



Big Picture



Fazit

- Der Einstieg ist einfach
- Wirtschaftlich und technisch skalierbar
- Ausprobieren lohnt sich



Vielen Dank



www.pleus.net



wolfgang.pleus@pleus.net

Download Folien und Code www.pleus.net/blog



Creative Software Workbench www.cswob.de

Lambda Architektur

