

Einstieg Big Data für DBAs

Patrick Kramer
DOAG Regio 10.09.2019



- Was verstehen wir unter Big Data?
- Welche Probleme sollen gelöst werden?
- Überblick über Hadoop
- HDFS
- Hive

Wie "big" ist Big Data?

- Big Data Technologien sind entstanden, um riesige Datenmengen
 - zu speichern
 - zu analysieren
- Gerechnet wird in Terabyte (TB) oder Petabyte (PB)!
- Die Technologien können auch für kleinere Datenmengen genutzt werden
- Einige Beispiele (Stand 2017)
 - NYSE: Hadoop Cluster mit 20+ PB und 30 TB neuen Daten je Tag
 - CERN: Der Teilchenbeschleuniger erzeugt 80+ TB Daten je Tag
 - Facebook: Speichert mehr als 15+ TB je Tag in seinem Hadoop Warehouse

Volume (Volumen)

- Gigabyte
- Terabyte
- Petabyte

Velocity (Geschwindigkeit)

- Batch
- Near-time
- Realtime

Variety (Vielfalt)

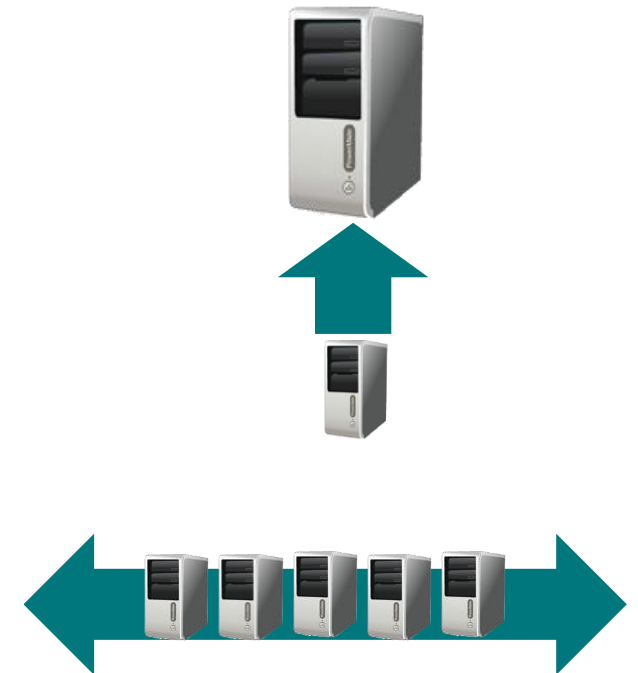
- Strukturiert
- Semistrukturiert
- Unstrukturiert

Agenda

- Was verstehen wir unter Big Data?
- Welche Probleme sollen gelöst werden?
- Überblick über Hadoop
- HDFS
- Hive

Welche technischen Probleme sollen gelöst werden?

- Klassische IT Systeme skalieren...
 - Vertikal sehr gut
 - Schnellere CPU
 - Mehr Speicher
 - Horizontal oft nur sehr begrenzt
 - Oracle RAC Cluster (3 Knoten sind üblich, 12 sind viel)
 - JBoss Cluster (4 Knoten sind üblich, 64 sind viel)
- Big Data Systeme skalieren...
 - Vertikal sehr gut
 - Aufrüstung existierender Knoten
 - Ersatz von Knoten durch neue und schnellere Hardware
 - Horizontal sehr gut
 - Yahoo hat einen Hadoop Cluster mit 4.500 Knoten
 - Netflix nutzt 80 verschiedene Cassandra Cluster mit insgesamt 2.500 Knoten



Welche fachlichen Probleme sollen gelöst werden?

Web

- Click Streams
- Soziale Netzwerke

Finanzwirtschaft

- Betrugsfallerkennung
- Risikomanagement

Internet-of-Things

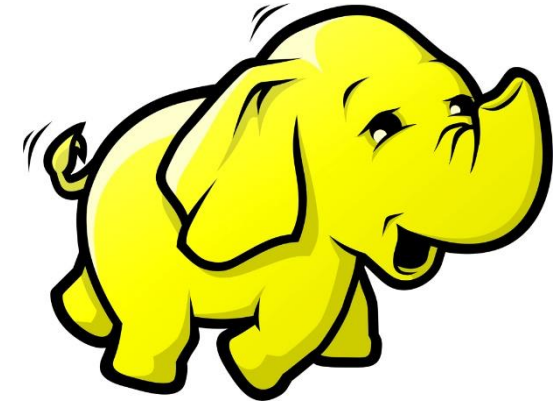
- Heimautomation
- Predictive Maintenance

...

Agenda

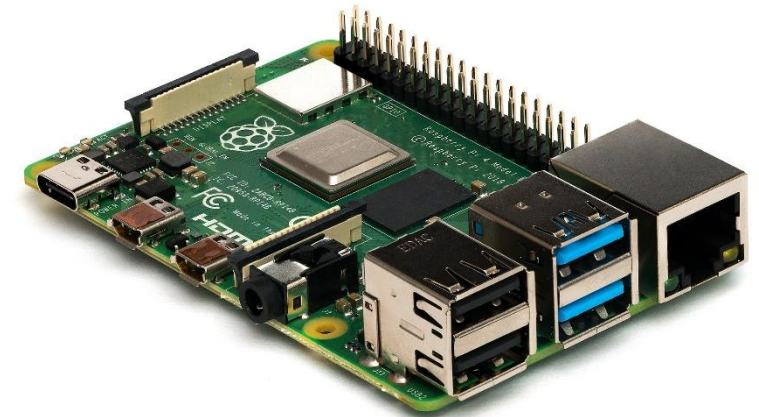
- Was verstehen wir unter Big Data?
- Welche Probleme sollen gelöst werden?
- Überblick über Hadoop
- HDFS
- Hive

- Verteilte Speicherung der Daten im Cluster
 - Erhöht die Geschwindigkeit beim Lesen und Schreiben
 - Striping (vergleichbar mit RAID0)
- Replikation der Daten im Cluster
 - Erhöht die Verfügbarkeit
 - 3-fach Replikation (vergleichbar mit RAID1)
 - Erasure Coding ab Hadoop 3 (vergleichbar mit RAID 5)
- Dateisystem angelehnt an Linux
 - `hadoop fs -ls`
 - Weitere Befehle: `-rm`, `-cat`, `-tail`
 - POSIX Rechte (rwx) und ACLs
- (Batch) Verarbeitung der Daten

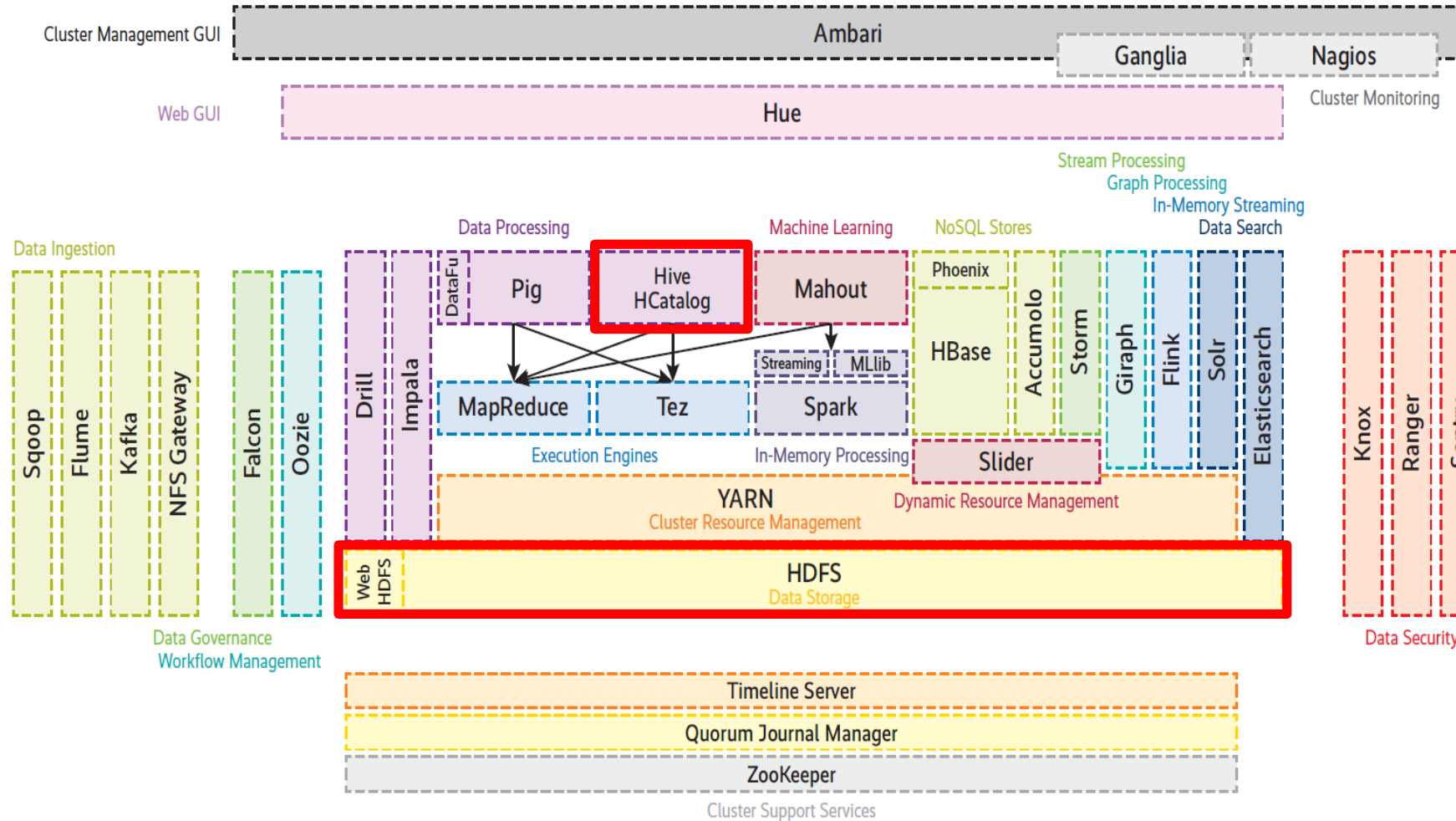


Apache Hadoop Design Prinzipien

- Scale Out Linear
 - Skaliert horizontal (und vertikal) nahezu linear
- Schema on Read
 - Erst beim Lesen wird Struktur in die Daten gebracht
- Write once, read many (WORM)
 - Optimiert für "Full Table Scans"
- Commodity Hardware
 - Nutzung günstiger und leicht verfügbarer Hardware
- Shared Nothing
 - Speicher und Platten werden nicht geteilt
 - Synchronisation der Prozesse über das Netzwerk
- Data Locality
 - Verarbeitung der Daten wo sie gespeichert sind
 - Verliert durch schnelle Netze und In-Memory Verarbeitung an Bedeutung



Apache Hadoop Ökosystem



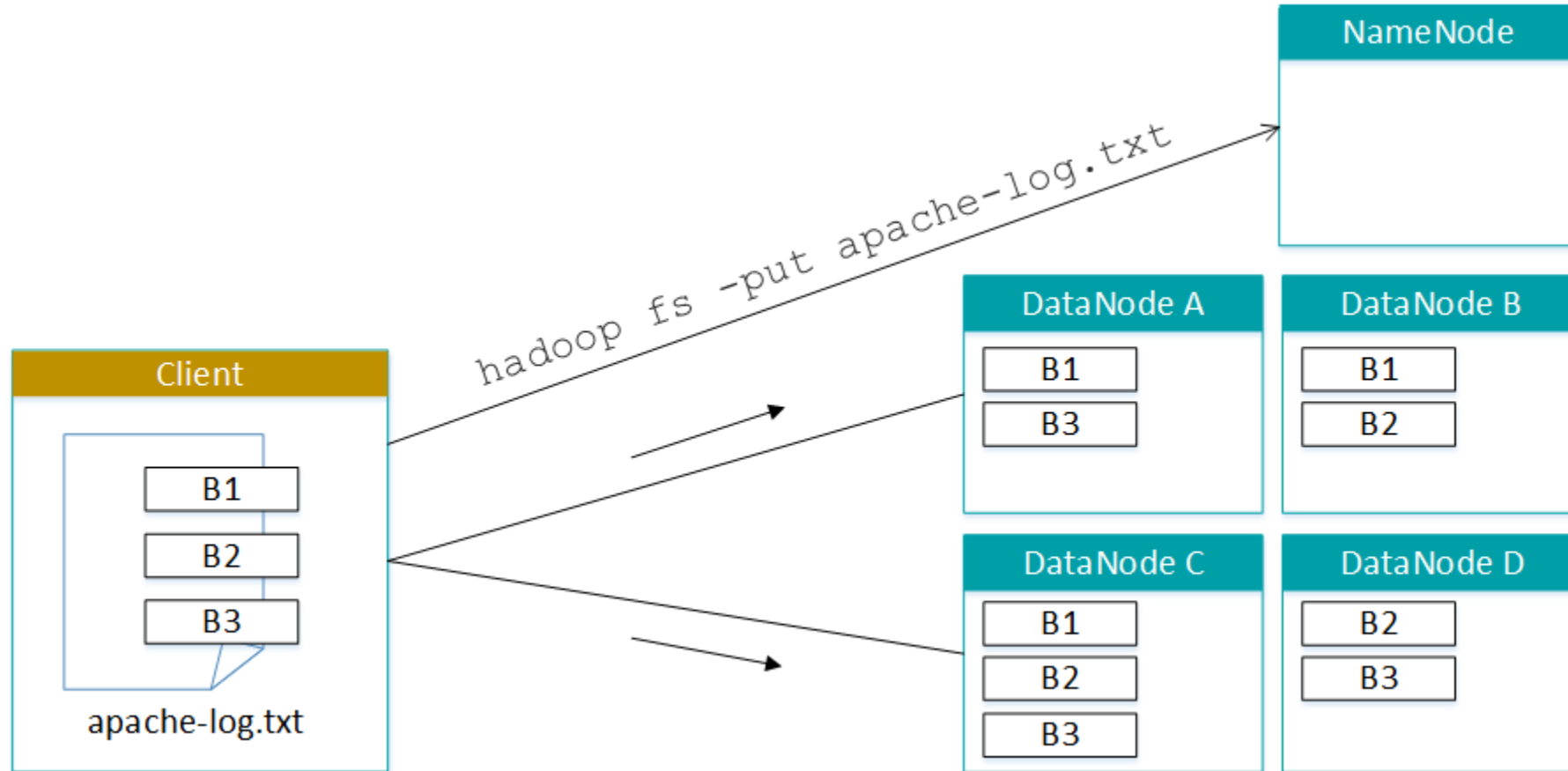
Quelle: iX Developer 2015 – Big Data

Agenda

- Was verstehen wir unter Big Data?
- Welche Probleme sollen gelöst werden?
- Überblick über Hadoop
- HDFS
- Hive

HDFS - Hadoop Distributed File System

- Verteiltes Dateisystem optimiert für...
 - Sehr große Dateien
 - Häufiges Lesen und keine Updates
 - "Full Table Scans" bzw. "Streaming Data Access"
- Speichert Daten auf den DataNodes
 - In großen Blöcken (128 MiB default)
 - Redundant
 - 3-fach Replikation default
 - Erasure Coding (z.B. 6+3) [Hadoop 3]
- Zugriff über Shell oder API
- Befehle und Rechtemanagement an POSIX angelehnt
 - `hadoop fs -ls /user`
 - `hadoop fs -put local.txt hdfs.txt`
 - `hadoop fs -chmod go-r hdfs.txt`



- Verzeichnisinhalt auflisten

```
hadoop fs -ls /user/ordix
```

- Inhalt einer Datei anzeigen

```
hadoop fs -cat hdfs-file.txt  
hadoop fs -tail hdfs-file.txt
```

- Datei vom lokalen Dateisystem ins HDFS kopieren

```
hadoop fs -copyFromLocal local-dir/file.txt hdfs-dir/  
hadoop fs -put local-dir/file.txt hdfs-dir/
```

- Datei vom HDFS ins lokale Dateisystem kopieren

```
hadoop fs -copyToLocal hdfs-dir/file.txt local-dir/  
hadoop fs -get hdfs-dir/file.txt local-dir/
```

- Datei im HDFS kopieren

```
hadoop fs -cp file.txt file-copy.txt
```

Agenda

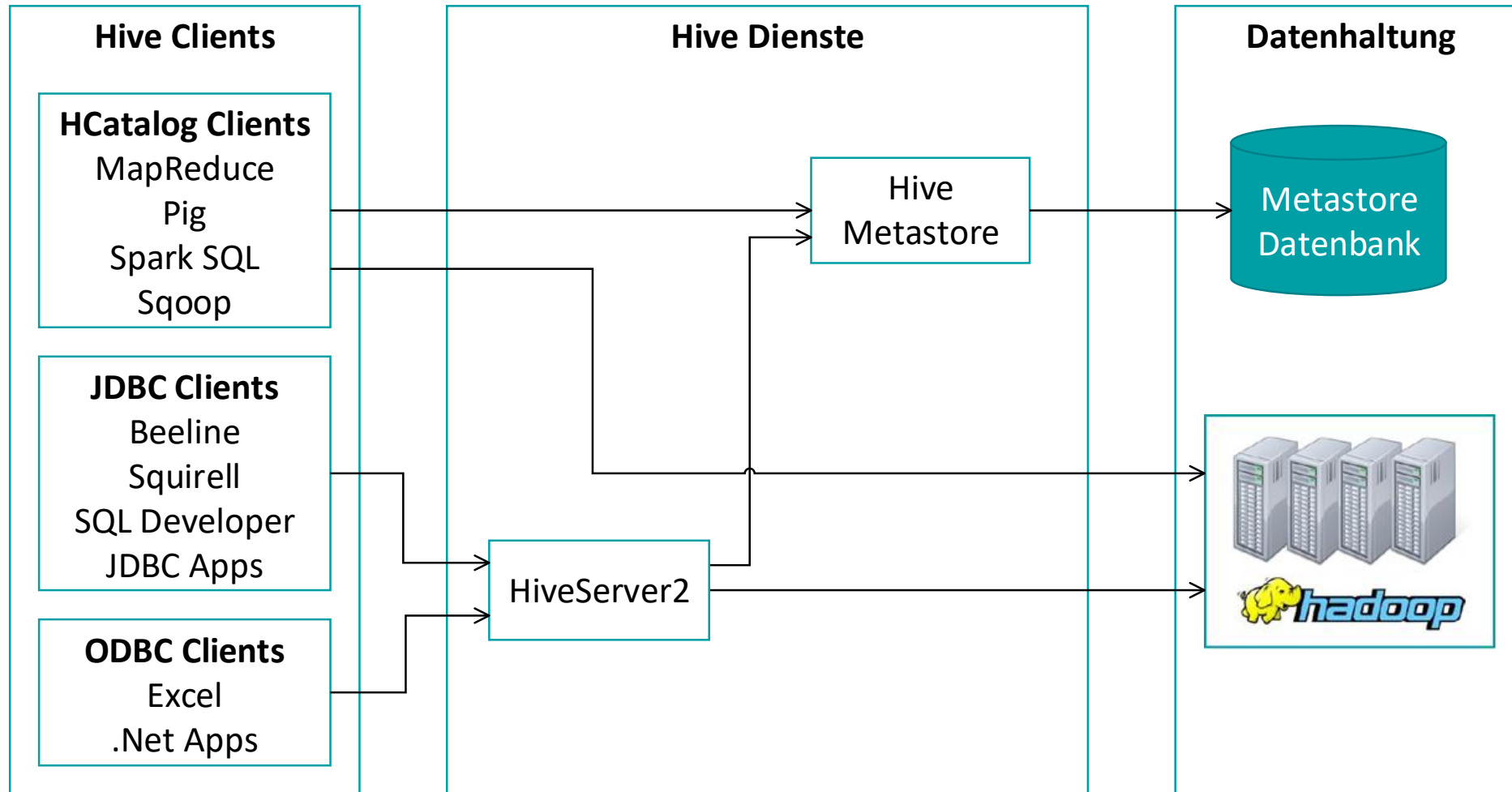
- Was verstehen wir unter Big Data?
- Welche Probleme sollen gelöst werden?
- Überblick über Hadoop
- HDFS
- Hive

Was ist Apache Hive?

- Hive ist ein Datawarehouse System
 - Zur Verarbeitung strukturierter Daten
 - Für verteilte Systeme (Hadoop, HBase, Cassandra, ...)
- HiveQL
 - SQL Dialekt von Hive
 - Stark beeinflusst von MySQL
- Beeline Shell
- Wichtige Features
 - Viele Dateiformate
 - Komplexe Datentypen (Array, Map, Struct)
 - Partitionierung
 - User-Defined Functions (UDFs)



Hive Architektur



- **Tabelle anlegen**

```
create table t1 (key int, value string);
```

- **Metadaten anzeigen**

```
show tables;  
describe t1;  
show create table t1;  
describe formatted t1;
```

- **Schema kopieren**

```
create table t2 like t1;
```

- **Tabelle löschen**

```
drop table t2;
```

Livedemo / Managed Tables und External Tables

- Table Type: MANAGED_TABLE
 - Ist der Standard
 - Hive verwaltet die Tabelle
 - DROP TABLE löscht Metadaten und Dateien

```
create table t1 (key int, value string);
```

- Table Type: EXTERNAL_TABLE
 - Hive verwaltet nur die Metadaten
 - DROP TABLE löscht nur Metadaten
 - Wenn andere Systeme direkt auf die Dateien zugreifen
z.B. Pig oder MapReduce

```
create external table t2 (key int, value string)  
location '/user/ordix/hivedwh/t2';
```

**Vielen Dank für
Ihre Aufmerksamkeit**

ORDIX[®] best practice
einfach. gut. beraten.

ORDIX AG
Aktiengesellschaft für
Softwareentwicklung, Schulung,
Beratung
und Systemintegration

Zentrale Paderborn
Karl-Schurz-Straße 19a
33100 Paderborn
Tel.: 05251 1063-0
Fax: 0180 1 67349 0

Seminarzentrum Wiesbaden
Kreuzberger Ring 13
65205 Wiesbaden
Tel.: 0611 77840-00

info@ordix.de
www.ordix.de

Bildquellen

- Seite 9: svn.apache.org/repos/asf/hadoop/logos/out_rgb/elephant_rgb.jpg
- Seite 10: By Miiichiaieil Heinizilieir / Wikimedia Commons, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=80140656>
- Seite 11: iX Developer 2015 – Big Data
- Seite 17: By Davod - Own work, using File:Apache Hive logo.jpg as base., Apache License 2.0, <https://commons.wikimedia.org/w/index.php?curid=44338923>