



NOSQL im BigData-Land: Von Analog-Istan nach Dig-Italien

Frank Pientka, Bonn, 14.06.2012



Vorstellung des Referenten: Frank Pientka



Dipl.-Informatiker (TH Karlsruhe)

Senior Architekt in Dortmund

heise.de/developer/Federlesen-Kolumne

Über 20 Jahre Erfahrung mit Datenbanken

Veröffentlichungen und Vorträge dazu



MATERNA
Information & Kommunikation

MATERNA-Unternehmensgruppe

Dr. Winfried Materna Helmut an de Meulen

155 Mio. € Umsatz (2011, vorläufig)

1.300 Mitarbeiter

Gegründet 1980

IT'S VALUE

© MATERNA GmbH 2012 www.materna.de 3

MATERNA
Information & Kommunikation

Die Reiseroute

- Analog-Istan: Ein Blick zurück. Wie wichtig ist Physik?
- Die Herausforderungen: Schwimmen lernen im Datenstrom
- Dig-Italien: Besuch in einem unbekanntem Land?
 - Leute: Neue Daten(-Banken) benötigt?
 - Sprache: Was ist NOSQL, MR?
 - Land: Welche Einsatzgebiete gibt es?
- Bewertung
- Ein Blick voraus

IT'S VALUE

© MATERNA GmbH 2012 www.materna.de 4

Analog-Istan: Warum IBM keine Platten mehr herstellt

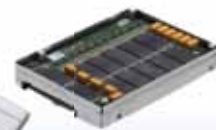
4. September 1956, die erste
Festplatte IBM 350 (5 MB)



2009
Western Digital WD
Caviar Green 2TB



2012
Hitachi Ultrastar
SSD400S.B 400GB



Datenexplosion

1,3 RFID tags
in 2005
30 Billion today

2 Billion
internet users
by 2011

4,6 Billion
mobile phones
world wide

Capital market
data volumes
grew 1750%
(2003-06)

Twitter
processes 7
Terabytes of
data per day

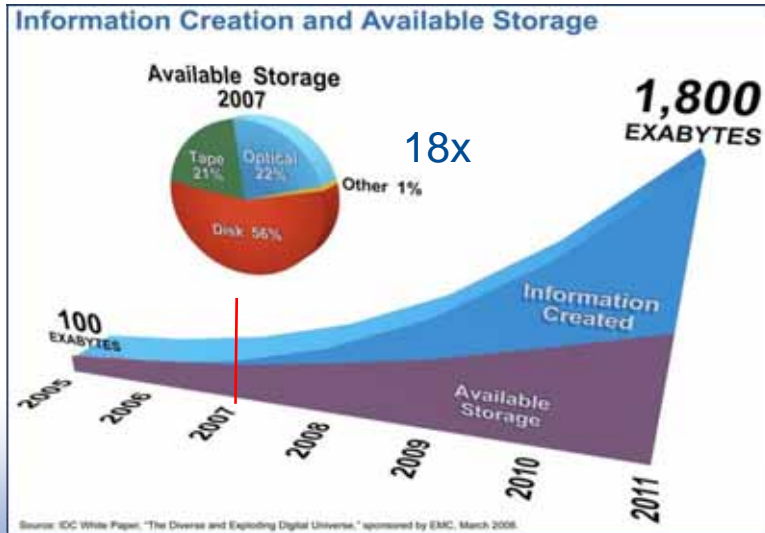
World Data Centre for
Climate
220 terabytes of
Webdata, 9 Petabytes of
additional data

Facebook
processes 10
Terabytes of
data per day

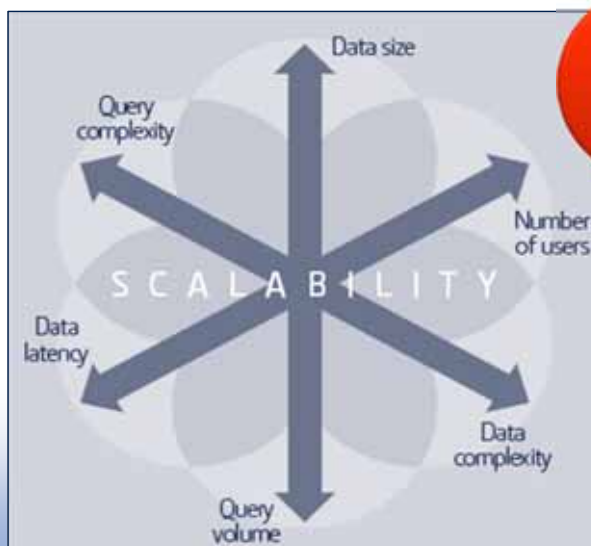
Datenexplosion - das digitale Universum

Information
Overkill
Overload

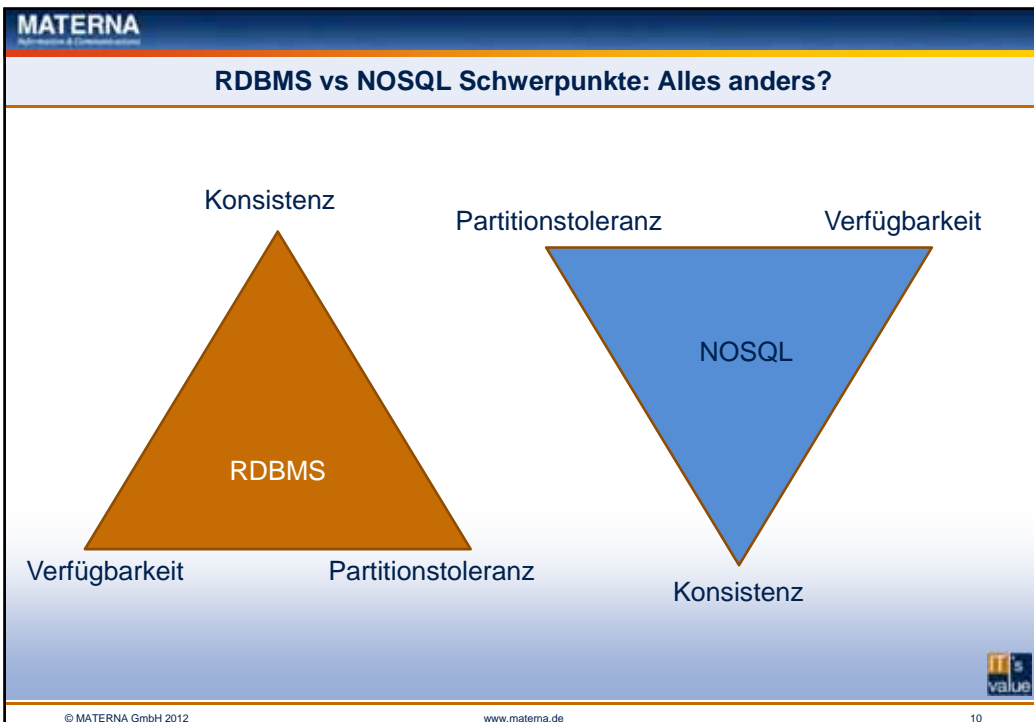
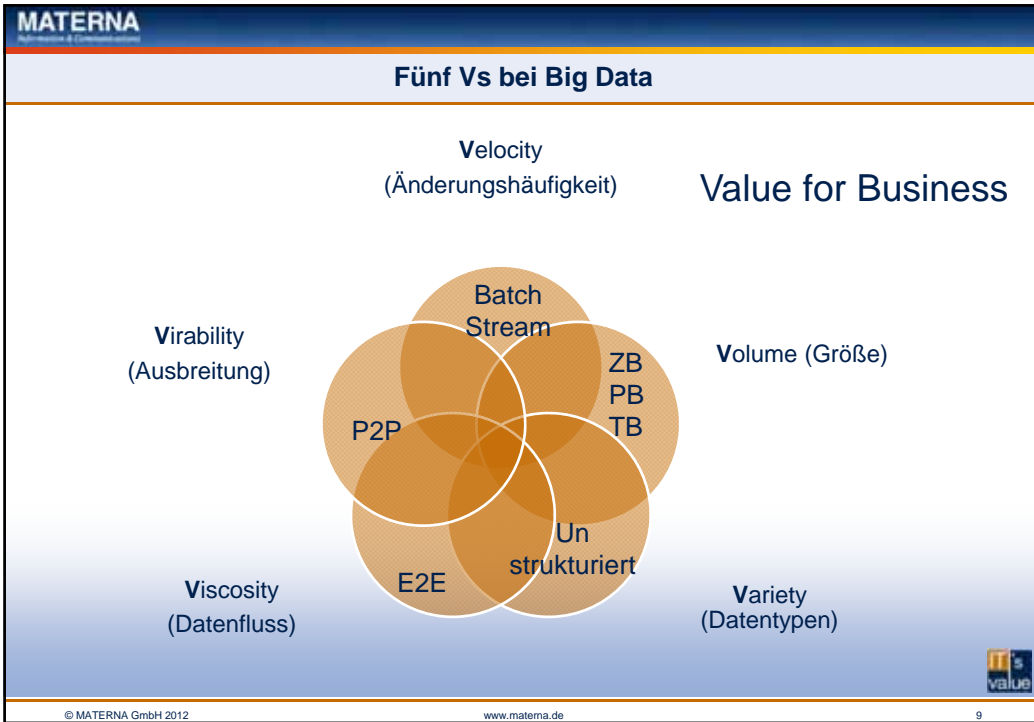
Die erzeugte Datenmenge
(gespeichert, gesendet)
übersteigt die aktuellen
Speicherkapazitäten



Datenwachstum ist ein mehrdimensionales Problem



Wachsende
Daten
wirken
sich auf alle
Bereiche
aus!!!

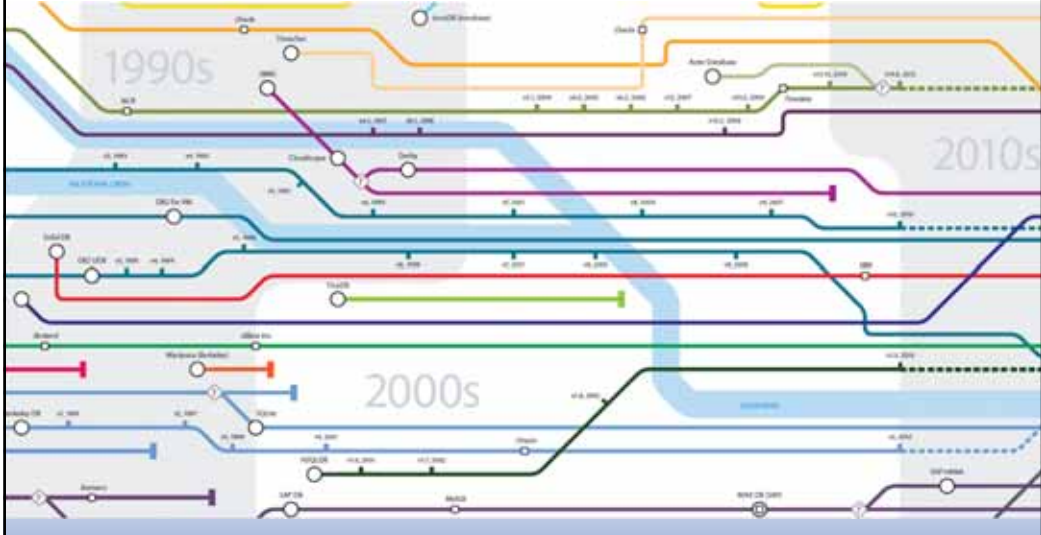


Die Reiseroute

- Analog-Istan: Ein Blick zurück. Wie wichtig ist Physik?
- Die Herausforderungen: Schwimmen lernen im Datenstrom
- Dig-Italien: Besuch in einem unbekanntem Land?
 - Leute: Neue Daten(-Banken) benötigt?
 - Sprache: Was ist NOSQL, MR?
 - Land: Welche Einsatzgebiete gibt es?
- Bewertung
- Ein Blick voraus



An relationalen Flüssen und Codd-Seen....



www.hpi.uni-potsdam.de/naumann/projekte/rdbms_genealogy



Eine Ära geht zu Ende...

*“The relational model of the 70s is not necessarily the answer”
(Michael Stonebraker, 2007)*



Von Analog-Istan ins wilde Dig-Italien...



Verbindung von un-/strukturierten Inhalten



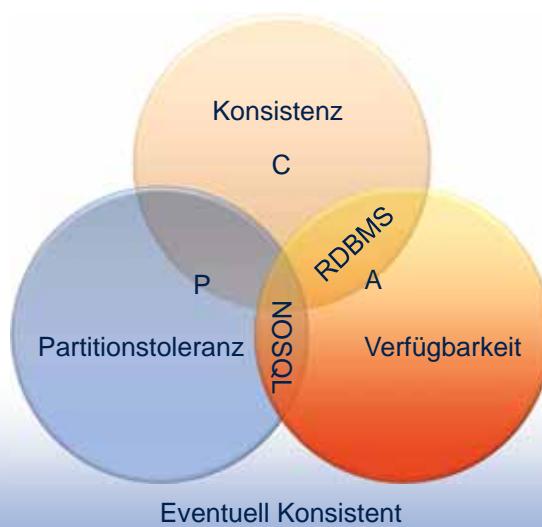
Die Reiseroute

- Analog-Istan: Ein Blick zurück. Wie wichtig ist Physik?
- Die Herausforderungen: Schwimmen lernen im Datenstrom
- Dig-Italien: Besuch in einem unbekanntem Land?
 - Leute: Neue Daten(-Banken) benötigt?
 - Sprache: Was ist NOSQL, MR?
 - Land: Welche Einsatzgebiete gibt es?
- Bewertung
- Ein Blick voraus

Kennzeichen von NOSQL-Systemen

- 1998 Carlo Strozzi No SQL, 2009 Johan Oskarsso Not only SQL
- Andere Datenstrukturen: einfach, geschachtelt, hierarchisch, vernetzt
- Andere Verteilung: CAP-Theorem
- Andere Verarbeitung: MR
- Andere Protokolle, Austauschformate: REST, JSON
- Partielle Konsistenz, statt verteilte Transaktionen: BASE vs ACID
- Mehr Flexibilität, weniger Kontrolle, Einschränkungen: Schemafrei, kein RI
- Weniger Verwaltungs-, mehr Programmieraufwand: Keine Sperren, Standard-QL
- Horizontale, lineare Skalierbarkeit: Standardhardware

CAP-Theorem: nur 2-aus3 (Brewer 2000)



ACID versus BASE Transaktionen

ACID:
Atomic
Consistent
Isolated Transactions
cannot interfere with
each other.
Durable

BASE:
Basic Availability
Soft-state
Eventual consistency

Verfügbarkeit wichtiger als
Konsistenz für schwächere
Transaktionsanforderungen

*"Your Coffee Shop doesn't use
Two-Phase-Commit"
(Gregor Hohpe, 2005)*



NOSQL-Systeme



Kommen nicht
aus den
Forschungslabors

Aus der Praxis des
Webs für das Web

Amazon, Google,
Facebook als
Ideengeber

OpenSource
verfügbar und
erweiterbar
(Ökosystem,
Partnerschaften)

Developed-by-
Community vs
Designed-by-
Company/
Committee



MATERNA
Information & Kommunikation

NOSQL-Produkte

Key Value Stores	Document Databases	BigTable Clones	Graph Database	Data Grid Caching

© MATERNA GmbH 2012 www.materna.de 21

MATERNA
Information & Kommunikation

NotOnlySQL-Datenbanken

N★SQL

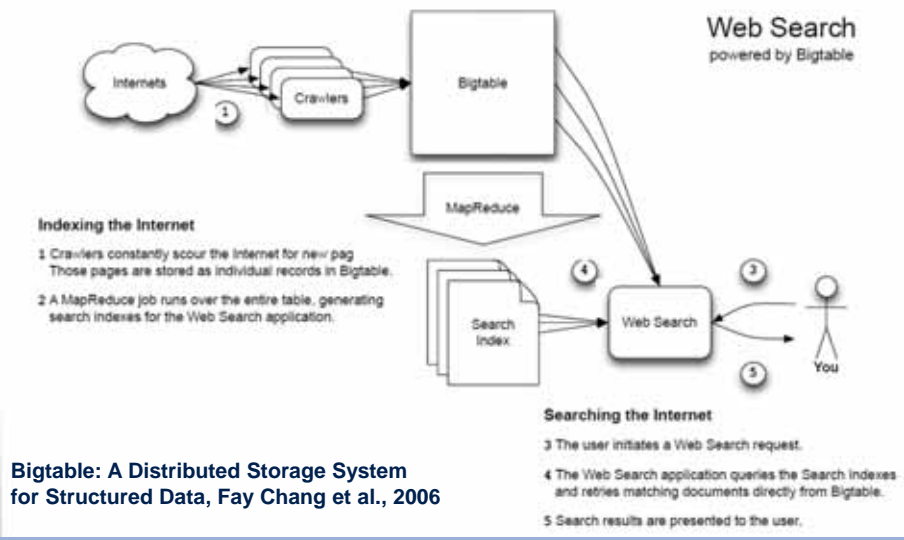
Gab es das nicht schon mal?

- 60er hierarchisch: IMS, LDAP, XML (Tamino, Xindice)
- 70er netzwerkartig: CODASYL
- 80er Hash-DB: xDBM Ken Thompson, **BerkleyDB**
- 80er dokumentenorientiert: Lotus Notes, Jackrabbit
- 90er objektorientiert: db4o, Versant
- 90er objekt-relational: ORACLE, DB2, mySQL

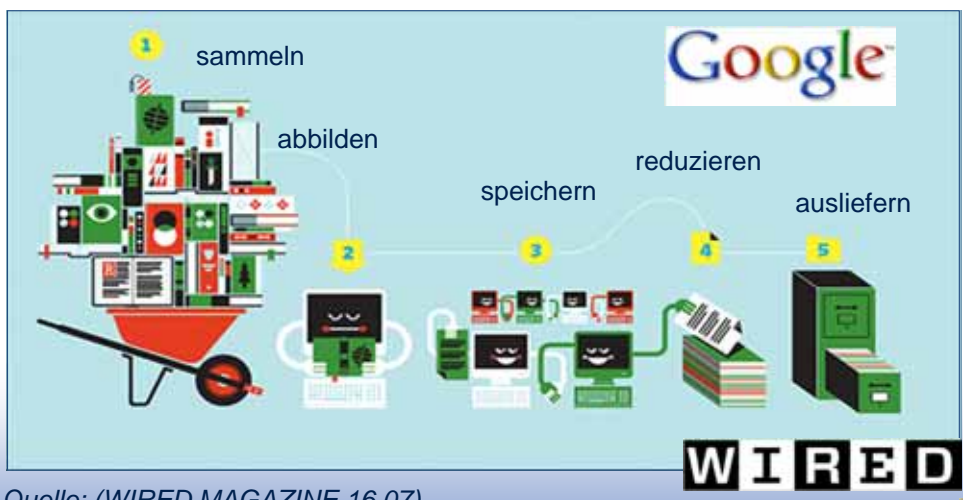
→ Hybride, temporale Datenbanken SQL:2011
Was ist die meist-installierte SQL-Datenbank?

© MATERNA GmbH 2012 www.materna.de 22

Der Einsatz von BigTable bei Google

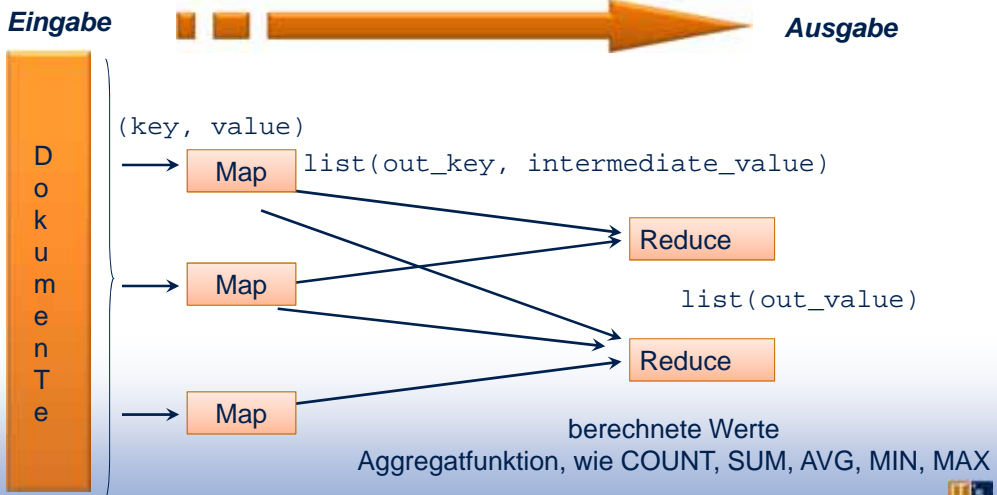


Sorting the World with MR

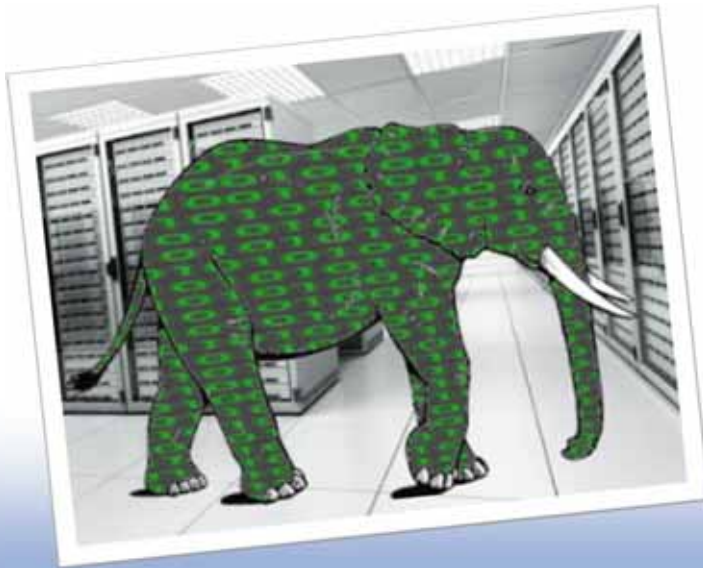


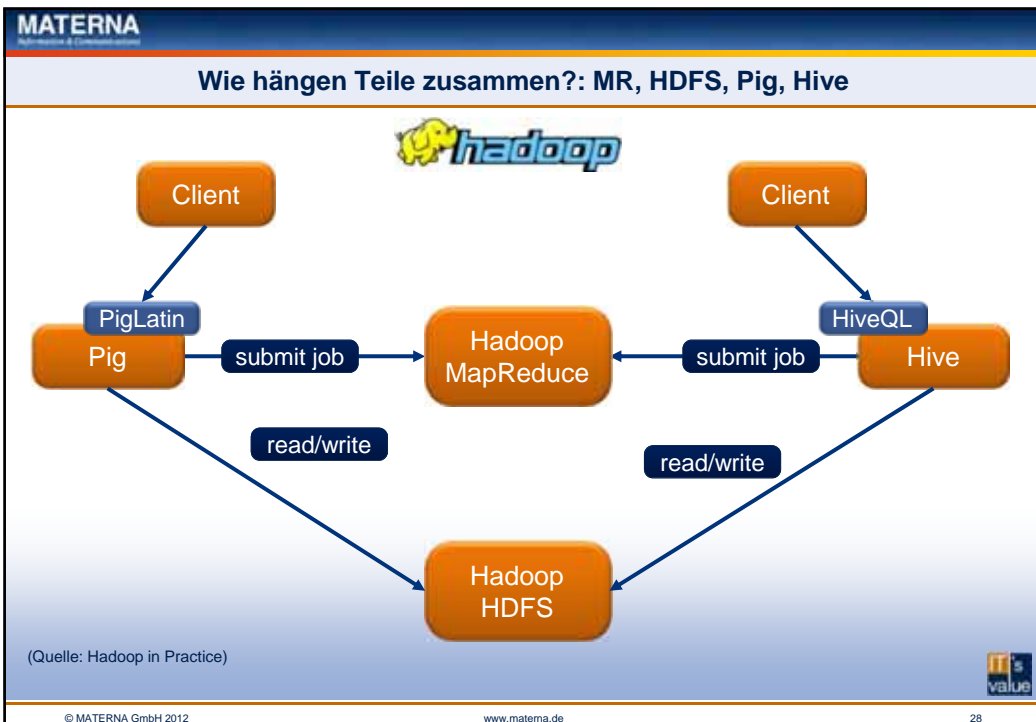
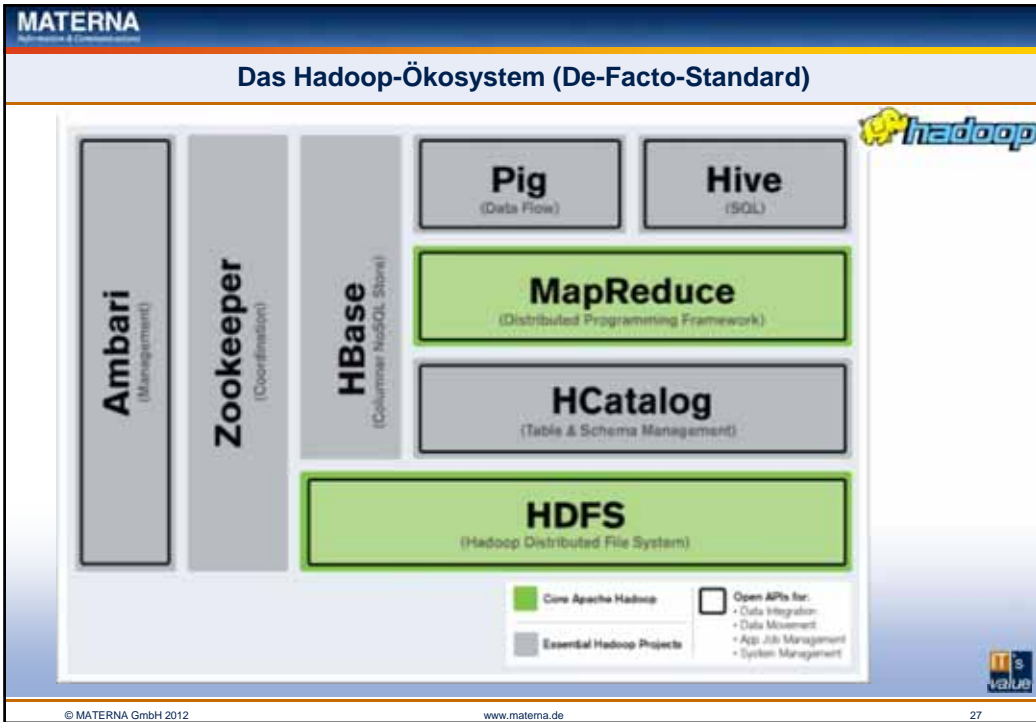
MapReduce-Verfahren

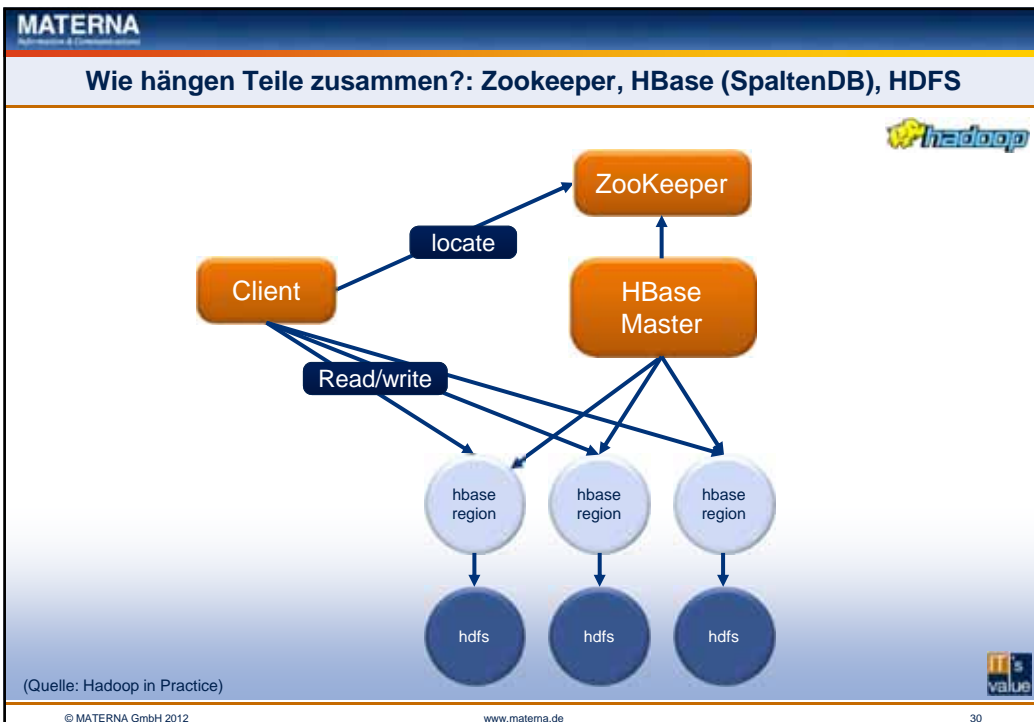
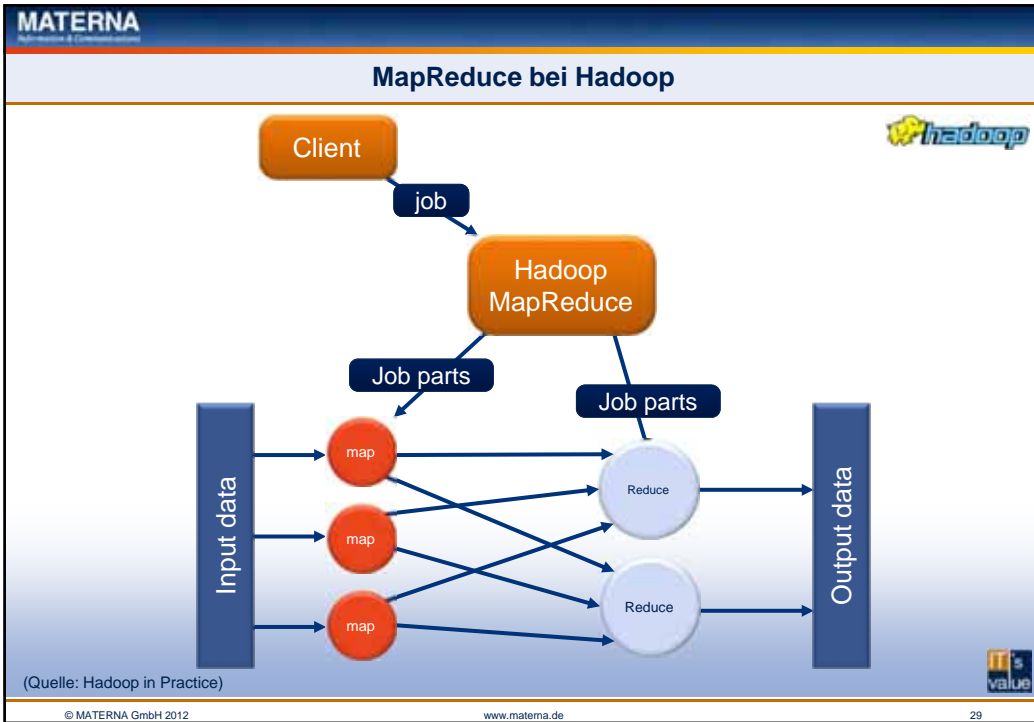
2 Phasen parallel ausgeführt, wie Pipes&Filter (UNIX)



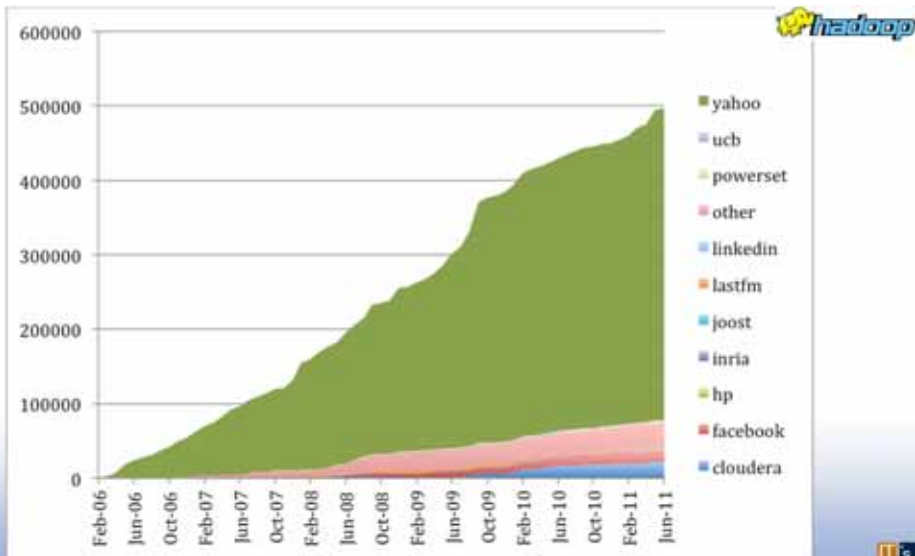
Der Elefant hat den Laden verlassen ...







Wer entwickelt Hadoop?



(Quelle: Hadoop in Practice)

Der BigData-Hadoop-Zoo



Pig

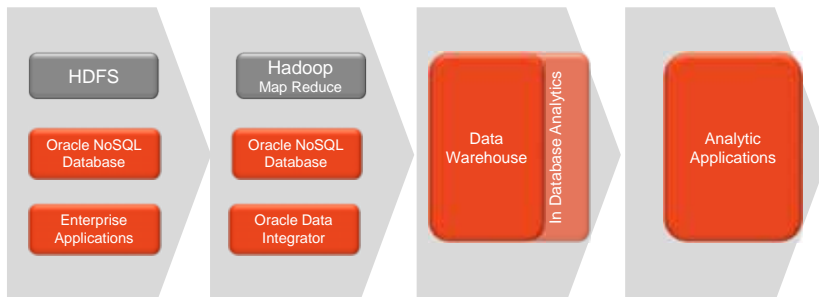


cloudera

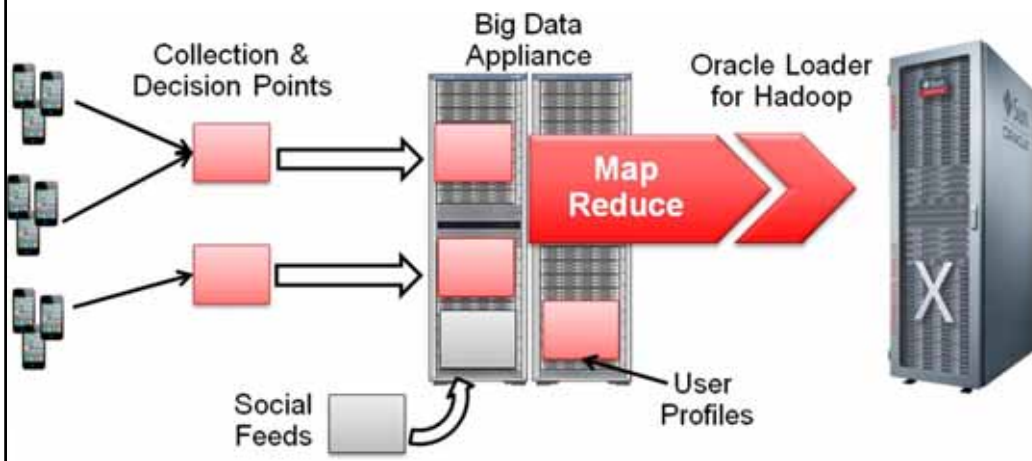
The Forrester Wave: Enterprise Hadoop Solutions, Q1 2012



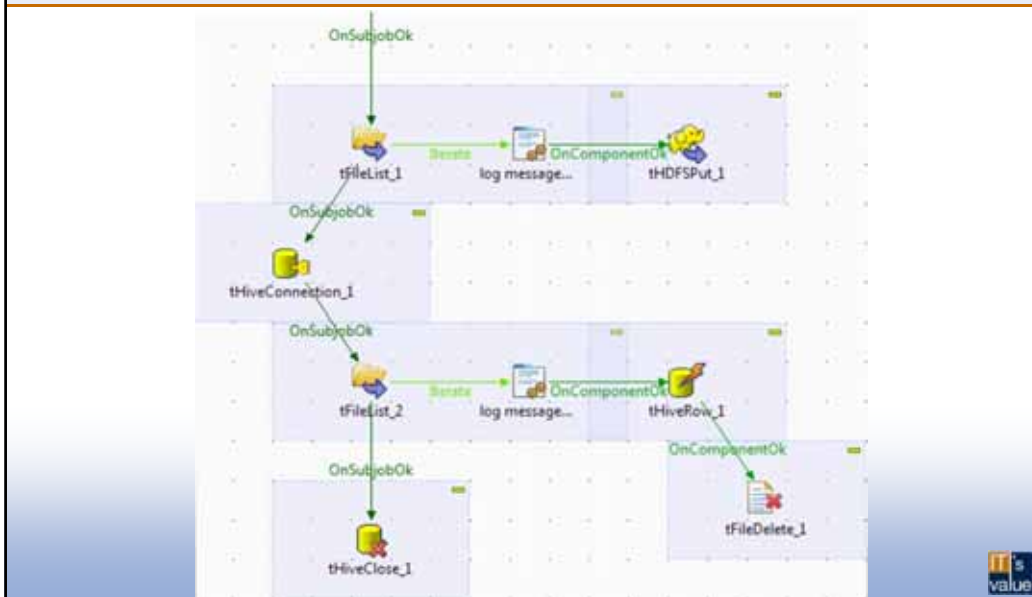
Oracle Big Data Appliance: Connectoren, Lader, NOSQL



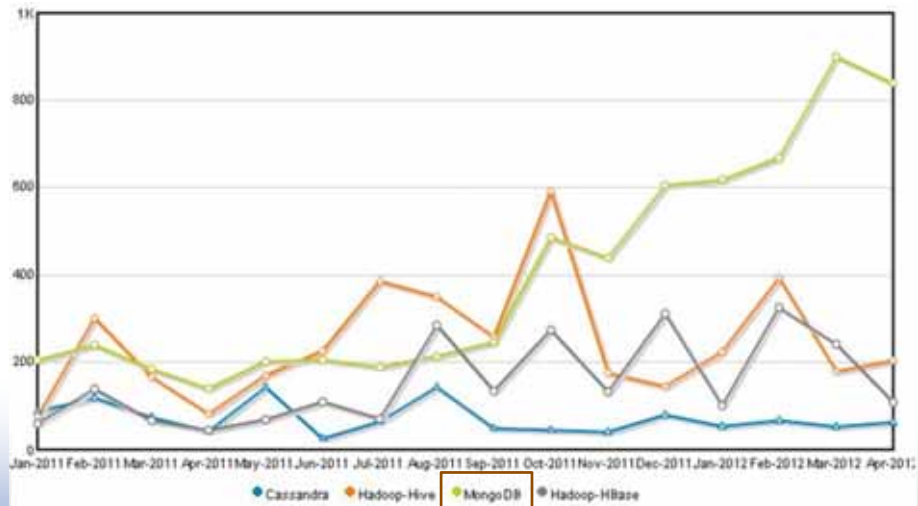
Einsatzmöglichkeit Oracle Big Data Appliance



Datenladen in HDFS und Hive mit Talend Open Studio



Monthly Downloads of Top Four Big Data Connectors (Jaspersoft)



Jaspersoft Big Data Index (JBDI)



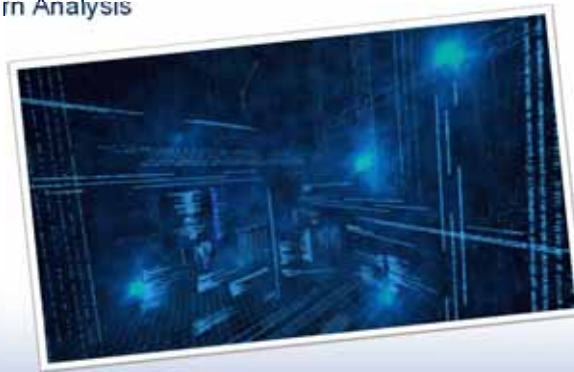
Die Reiseroute

- Analog-Istan: Ein Blick zurück. Wie wichtig ist Physik?
- Die Herausforderungen: Schwimmen lernen im Datenstrom
- Dig-Italien: Besuch in einem unbekanntem Land?
 - Leute: Neue Daten(-Banken) benötigt?
 - Sprache: Was ist NOSQL, MR?
 - Land: Welche Einsatzgebiete gibt es?
- Bewertung
- Ein Blick voraus



Big Data Use Cases

- Recommendation Engine
- Marketing Campaign Analysis
- Customer Retention and Churn Analysis
- Social Graph Analysis
- Capital Markets Analysis
- Predictive Analytics
- Risk Management
- Rogue Trading
- Fraud Detection
- Retail Banking
- Network Monitoring
- Research And Development
- Archiving



Wirtschaftlicher Nutzen von BigData-Technologien

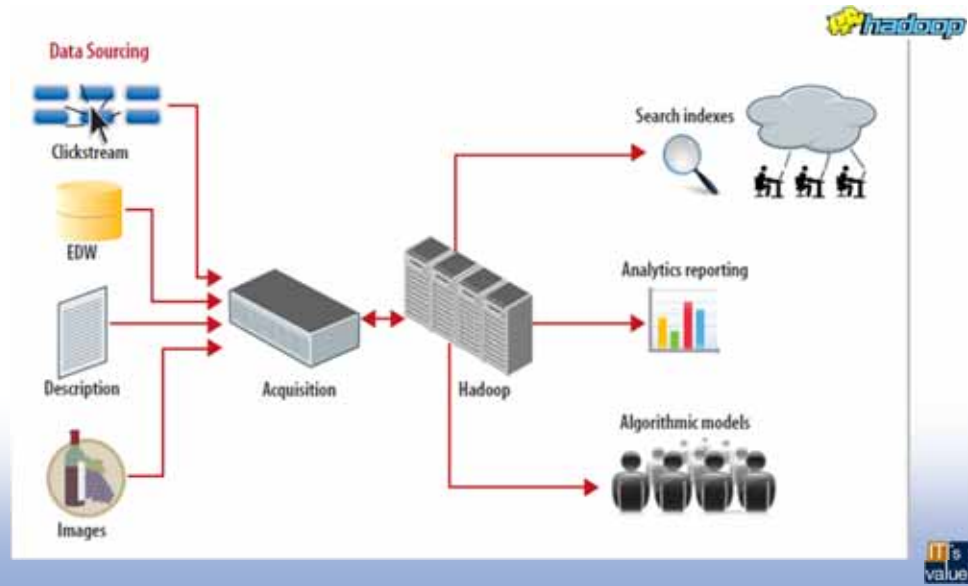
McKinsey&Company



Wie mit Big Data arbeiten? (Latenzzeiten)

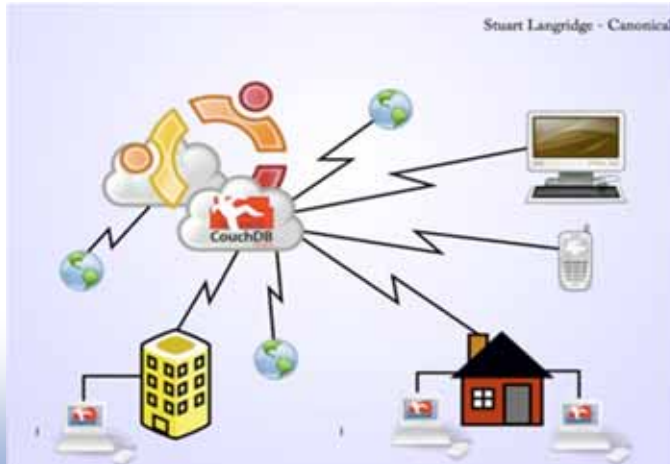


Clickstream-Analyse mit Hadoop bei eBay



CouchDB Ubuntu One: Cloud, PC, Smartphone

CouchDB auf über millionen PCs und in der Cloud mit Ubuntu 9.10 10.2009
EOL mit Ubuntu One 12.04 11.2011



Die Reiseroute

- Analog-Istan: Ein Blick zurück. Wie wichtig ist Physik?
- Die Herausforderungen: Schwimmen lernen im Datenstrom
- Dig-Italien: Besuch in einem unbekanntem Land?
 - Leute: Neue Daten(-Banken) benötigt?
 - Sprache: Was ist NOSQL, MR?
 - Land: Welche Einsatzgebiete gibt es?
- Bewertung
- Ein Blick voraus



Zwischen Analog-Istan und Dig-Italien



Bewertung und Ausblick

- BigData: große Herausforderung, große Chance
- Die Hadoop-Plattform ist Defacto-Standard für Auswertung großer Datenmengen
- Standards (UnQL, Jaql, JSON Schema, HQL, BSON) fehlen
- MR, JSON-Datentyp werden in Datenbanken integriert
- Connectoren und Standardschnittstellen überbrücken Welten
- Professionalisierung der NOSQL-Produkte
- Marktkonsolidierung und Innovation birgt Risiken
- Fachkräftemangel und erste BestPractices-Erfahrungen
- "The Hadoop and MapReduce market will likely develop along the lines established by the development of the Linux ecosystem" (IDC 2012, Dan Vesset)



Fazit



- BigData geht uns alle an (technisch, rechtlich)
- Neben der Verarbeitung wird Auswertung unstrukturierte Daten beschleunigt
- Starke Verwandtschaft mit Suche, BI, ETL, DW, DataMining, Batch, Streaming
- OpenSource, Connectoren und Appliance senken Einstiegshürden
- Für langfristigen Nutzen müssen Konzepte verstanden und für eigene Anforderungen angepasst werden

Big Data: Technologies and Techniques for Large-Scale Data



“Ultimately, big data is more about attitude than tools; data-driven organizations look at big data as a solution, not a problem.”

Roger Magoulas und Ben Lorica, Februar 2009, O'Reilly

Datenwachstum ist nicht nur ein Datenbankproblem

*Das Problem des Datenwachstums
ist so alt wie die IT.
Kein All-Heilmittel, sondern nur
Methoden, Werkzeuge,
um Schmerzen lindern.*



Die Big Five

Schneller, Stärker,
Größer....



...haben keine
natürlichen
Feinde!



Ist BigData reif für die Praxis?



Es hängt davon ab, was man damit macht?

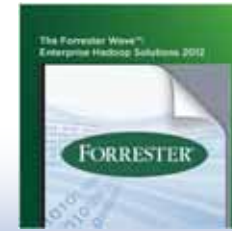


Weitere Informationen



<http://nosql-database.org>
<http://www.emc.com/leadership/programs/digital-universe.htm>

McKinsey&Company



Vielen Dank für Ihre/Eure Aufmerksamkeit! Besuchen Sie unseren Stand

MATERNA GmbH
Dipl. Inform. Frank Pientka
Senior Software Architect
Business Division Applications

Telefon: +49 231 5599-8854
Telefax: +49 231 5599-272
Frank.Pientka@materna.de
http://xing.to/frank_pientka

