



# Big Data and Advanced Analytics Technologies and Use Cases

*Colin White  
President, BI Research  
DAMA Portland  
February 2013*

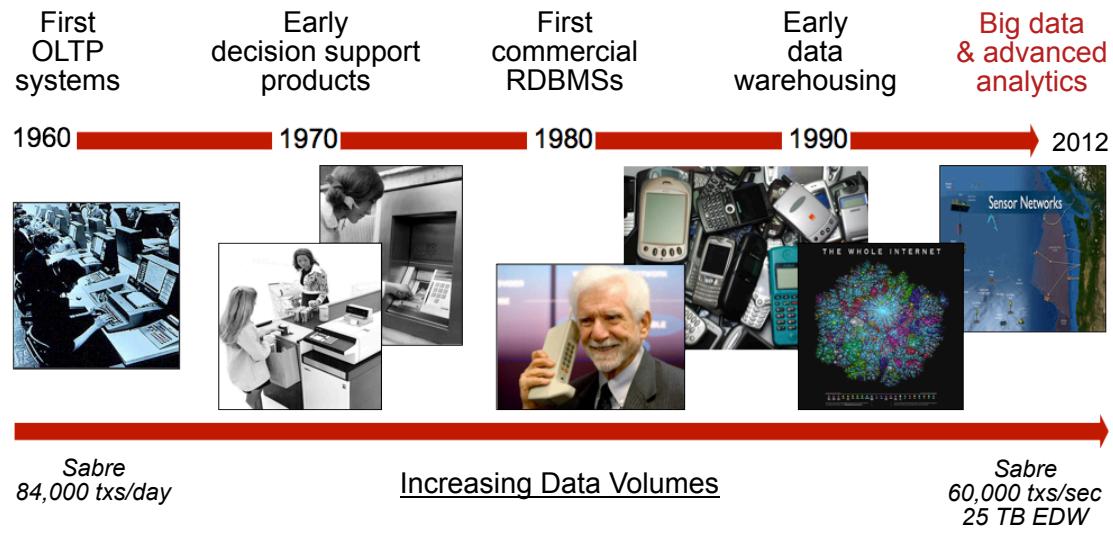


## Agenda

There is considerable interest at present on the topic of big data. Much of the discussion about this topic, however, is focused on the technology supporting big data, rather than on how analytics generated from big data can be leveraged for business benefit. One of the most exciting aspects of big data technology is that it allows organizations to support advanced analytic workloads and applications that were not previously possible for cost or performance reasons. New and evolving big data solutions provide significant business benefits because they help remove these cost and performance barriers. The objectives of this presentation are to discuss the benefits of big data and to present use cases and case studies that demonstrate the value of advanced analytics. It also explains how the existing data warehousing environment can be extended to support big data solutions. Topics that will be covered include:

- Review the history and evolution of big data and advanced analytics
- Explain the role of the data scientist in developing advanced analytics
- Look at the technologies that support big data
- Explain how the existing data warehousing environment can be extended to support big data and advanced analytics
- Discuss big data use cases and the benefits they bring to the business

# The Evolution of Digital Data



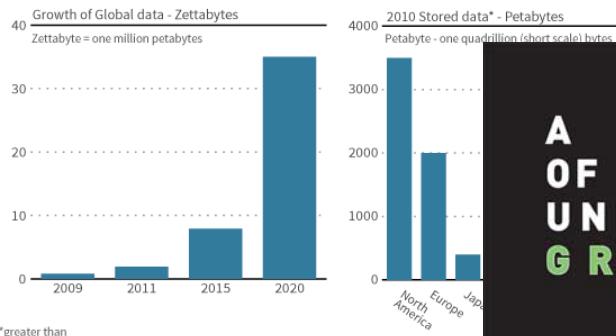
Copyright © BI Research, 2013

3

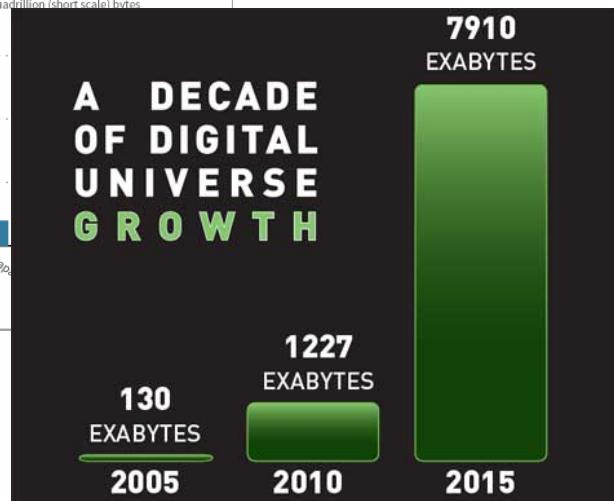
## Data Growth: Choose an Analyst!

### Big data growth

Big data market is estimated to grow 45% annually to reach \$25 billion by 2015



Reuters graphic/Catherine Trevethan 05/10/12



Copyright © BI Research, 2013

4

## Data Growth: Multi-Structured Data

Definition: data that has unknown, ill-defined or overlapping schemas

- Machine generated data, e.g., sensor data, system logs
- Internal/external web content including social computing data
- Text, document and XML data
- Graph, map and multi-media data

Volume increasing faster than structured data

Usually not integrated into a data warehouse

Increasing number of analytical techniques to extract useful information from this data

This information can be used to extend traditional predictive models and analytics

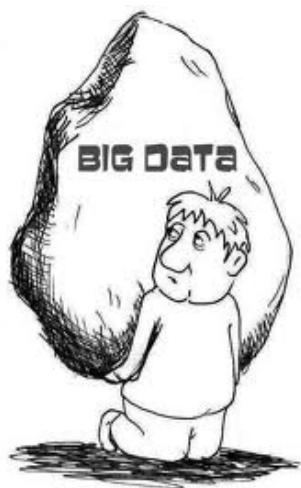
```
96.255.99.50 - - [01/Jun/2010:05:28:07 +0000] "GET /origin-log.enquisit...com/d.js?id=a1a3af-ly6l645&referrer=http://www.google.com/search/?hl=en&q=budget+planner&aq=5&aqi=g10&aqd=sog&budget+&gs_rfai=&location=https://money.strands.com/content/simple-and-free-monthly-budget-planner&ua=Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.0; SLC1; .NET CLR 2.0.50727;.NET CLR 3.0.30618;.NET CLR 3.5.30729; InfoPath.2)kpc=pgyv63w0xgn102in8m37wka8quxe74e&sc=cr1kot0wmxqik1wlr9p9weh6yy8q8sa&r=0.07550191624904945 HTTP/1.1" 200 380 "-" "Mozilla/4.0 (compatible; MSIE 7.0; Windows NT 6.0; SLC1; .NET CLR 2.0.50727;.NET CLR 3.0.30618;.NET CLR 3.5.30729; InfoPath.2)" "ac=bd76aad174480000679a040e005b130000"
```



Copyright © BI Research, 2013

5

## Data Growth: Big Data



Big data technologies apply to all types of digital data not just multi-structured data

“Big” is a relative term and is different for each organization and application

What you do with big data and how you use it for business benefit should be the main consideration – analytics play a key role here

Copyright © BI Research, 2013

6

# The Value of Data: IBM 2012 Study

IBM Center for Applied Insights

Outperforming in a data-rich, hyper-connected world



We live in the era of “big data” and connectivity. Collectively, our planet generates fifteen petabytes of new information every day<sup>1</sup> – about eight times the information housed in all the academic libraries in the United States.<sup>2</sup> People are connecting and communicating more than at any time in human history. In 2011 alone, an estimated seven trillion text messages were sent.<sup>3</sup> There are almost six billion mobile phone subscriptions.<sup>4</sup> Today, two billion people worldwide are plugged into the Internet,<sup>5</sup> interacting and sharing information on sites like Facebook and Twitter.

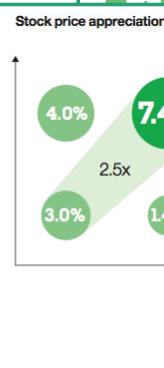
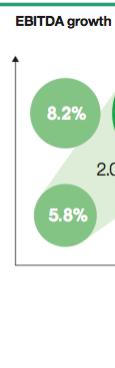
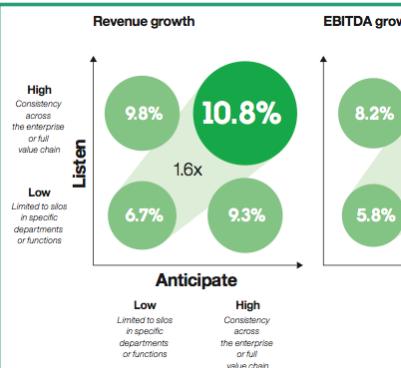
## About the Study

The IBM Center for Applied Insights, in cooperation with the Economist Intelligence Unit and the IBM Institute of Business Value, recently surveyed 1,168 executives (two thirds of which were CxOs) across nine industries in 64 countries. We investigated their ability to capture data, share insights and take actions based on what they learn, and used a binary logistic regression analysis to understand the statistically significant correlations with financial performance.

Copyright © BI Research, 2013

7

# The Value of Data: IBM 2012 Study



Listen ➤



Integrate data using common data definitions

Integrate data using end-to-end processes

Capture real-time data using automated devices

Anticipate ➤

Trigger alerts based on potential impact

Predict impact of future outcomes

Model best outcome based on unknown conditions

Copyright © BI Research, 2013

8

# The Changing World of BI Analytics

## Advanced Analytics

- Improved analytic tools and techniques for statistical and predictive analytics
- New tools for exploring and visualizing new varieties of data
- Operational intelligence with embedded BI services and BI automation

## Data Management

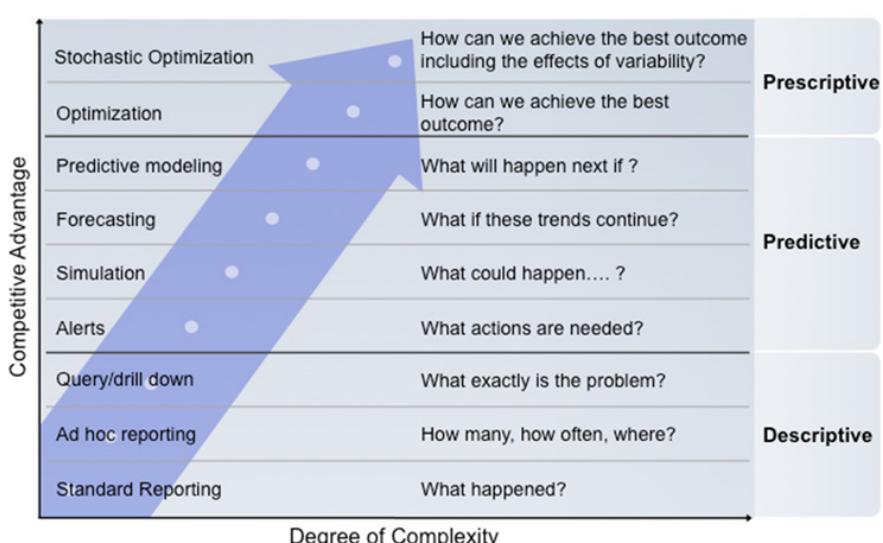
- Analytic relational database systems that offer improved price/performance and libraries of analytic functions
- In-memory computing for high performance
- Non-relational systems such as Hadoop for handling new types of data
- Stream processing/CEP systems for analyzing in-motion data



Copyright © BI Research, 2013

9

## Advanced Analytics Example: SC Digest 2012

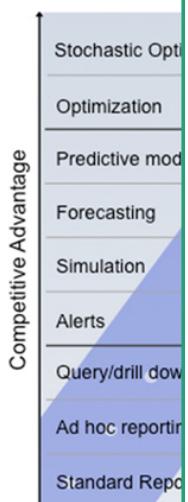


Based on: Competing on Analytics, Davenport and Harris, 2007

Copyright © BI Research, 2013

10

## Advanced Analytics Example: SC Digest 2012



1. **Descriptive analytics** – using historical data to describe the business. This is usually associated with Business Intelligence (BI) or visibility systems. In supply chain, you use descriptive analytics to better understand your historical demand patterns, to understand how product flows through your supply chain, and to understand when a shipment might be late.
2. **Predictive analytics** – using data to predict trends and patterns. This is commonly associated with statistics. In the supply chain, you use predictive analytics to forecast future demand or to forecast the price of fuel.
3. **Prescriptive analytics** – using data to suggest the optimal solution. This is commonly associated with optimization. In the supply chain, you use prescriptive analytics to set your inventory levels, schedule your plants, or route your trucks.

Advanced Analytics in Supply Chain, Dr. Michael Watson, Supply Chain Digest, November 2012

Copyright © BI Research, 2013

11

## The Role of the Data Scientist: CITO Interviews

“Data scientists turn big data into big value, delivering products that delight users, and insight that informs business decisions.”

Strong analytical skills are a given: above all, a data scientist needs to be able to derive robust conclusions from data. But a data scientist also needs to possess creativity and strong communication skills.”

*Daniel Tunkelang, Principal Data Scientist, LinkedIn*

“A data scientist is someone who can obtain, scrub, explore, model and interpret data, blending hacking, statistics and machine learning. Data scientists not only are adept at working with data, but appreciate data itself as a first-class product.”

*Hilary Mason, Chief Scientist at bitly*

“... someone who has the both the engineering skills to acquire and manage large data sets, and also has the statistician’s skills to extract value from the large data sets and present that data to a large audience.”

*John Rauser, Principal Engineer, Amazon.com*

[Source: citoresearch.com/content/growing-your-own-data-scientists](http://citoresearch.com/content/growing-your-own-data-scientists)

Copyright © BI Research, 2013

12

# Data Science Skills Requirements

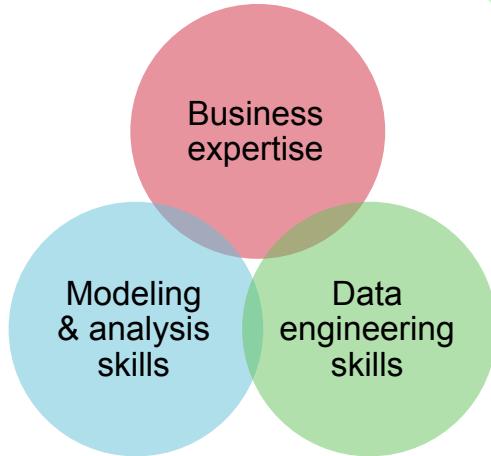
Business domain subject matter expert with strong analytical skills

Creativity and a good communications

Knowledgeable in statistics, machine learning and data visualization

Able to develop data analysis solutions using modeling/analysis methods and languages, such as MapReduce, R, SAS, etc.

Adept at data engineering, including discovering and mashing/blending large amounts of data

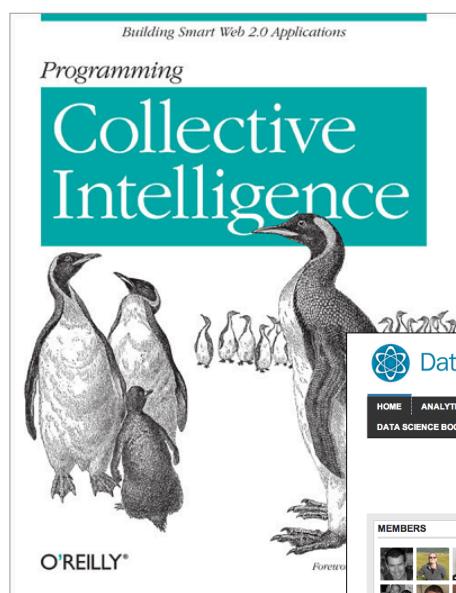


*Is this one person or a team of specialists?*

Copyright © BI Research, 2013

13

## Data Science: Further Reading



A screenshot of the Data Science Central website. At the top, there is a navigation bar with links for HOME, ANALYTICS, BIG DATA NEWS, VISUALIZATION, INDUSTRY BUZZ, DATA INTEGRATION, JOBS, FORUMS, BLOGS, and MY PROFILE. Below the navigation is a banner for 'PREDICTIVE ANALYTICS THE NEXT HORIZON'. On the left, there is a 'MEMBERS' section showing a grid of user profiles. On the right, there is a welcome message for new users and social media sharing options.

Copyright © BI Research, 2013

14

# What Then is Big Data?

Represents **analytic and data management solutions** that could not previously be supported because of:

- Technology limitations – poor performance, inadequate analytic capabilities, etc.
- High hardware and software costs
- Incomplete or limited data for generating the required solutions

Set of overlapping technologies that enable customers to deploy analytic systems optimized to suite specific business needs and workloads

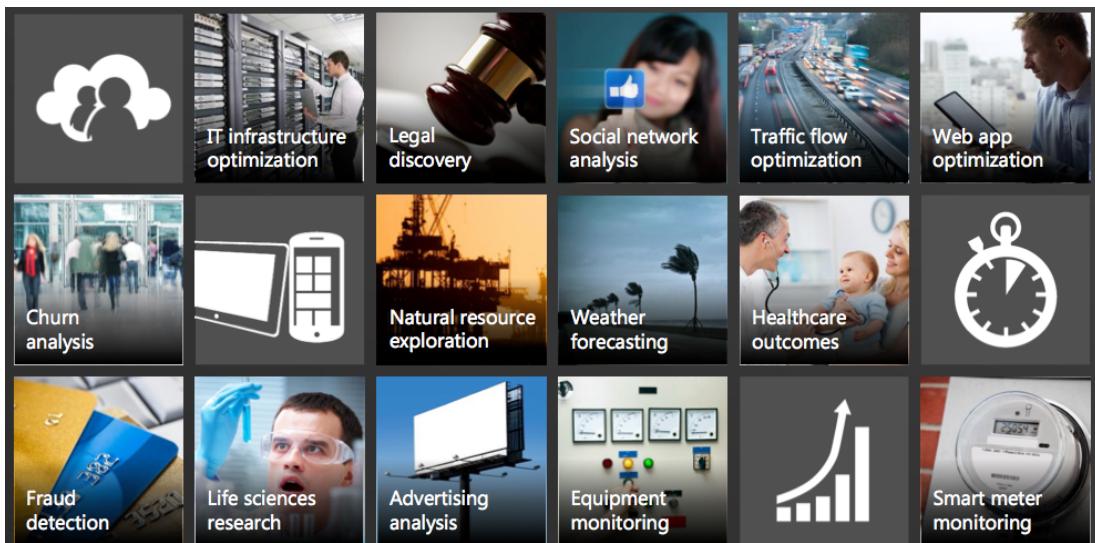


Optimization may involve improving performance, reducing costs, enabling new types of data to be analyzed, etc.

Copyright © BI Research, 2013

15

## Big Data Application Examples



Source: Microsoft

Copyright © BI Research, 2013

16

# Big Data and Data Life Cycle Management

## Data Management and Analytic Performance

- Capacity planning
- Managing data warehouse growth
- Analytic performance management and optimization
- Service level agreements

## Data Governance

- Security: user access, encryption, masking, etc.
- Quality: governed/ungoverned data
- Backup and recovery
- Archiving and retention: historical analysis, compliance



Copyright © BI Research, 2013

17

# The Impact of Big Data on the Data Life Cycle

Need fast time to value to quickly gain business benefits from big data

- Impractical to use traditional EDW approach for all analytic solutions
- Extend existing data warehousing environment to support big data and accommodate data growth

Need high performance solutions for supporting big data analytic workloads

- One-size fits all data management is no longer viable
- Match technologies and costs to business needs and analytic workloads

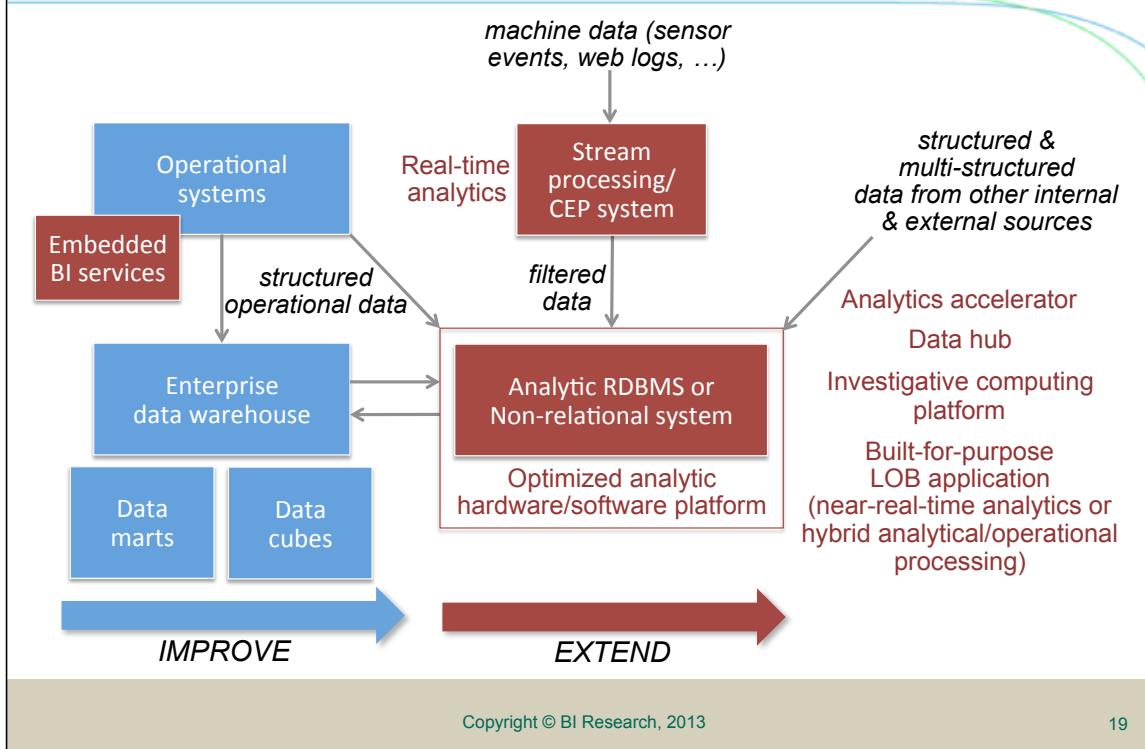
Need improved data governance to handle big data

- No longer practical to rigidly control and govern all forms of data – implement different levels of governance based on security, compliance and quality needs
- Determine data archiving and policies based on the possible future need to analyze historical data and data compliance requirements

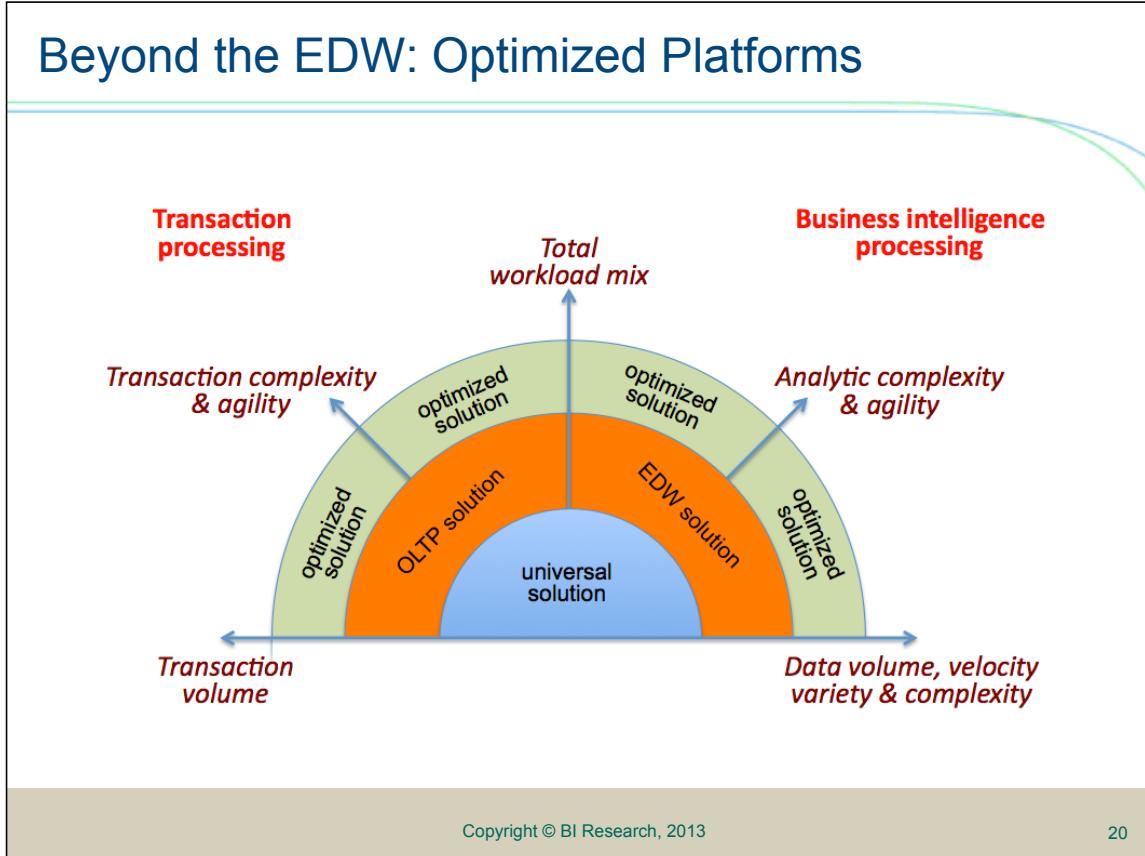
Copyright © BI Research, 2013

18

## The Extended (or Logical) Data Warehouse



## Beyond the EDW: Optimized Platforms



## Optimized Analytic Platforms: Variables

### Analytics Required to Meet Business Needs

- Complexity – reporting, OLAP or advanced analytics
- Agility – latency of data, analytics, decisions, recommendations and actions
- Workload mix – complexity of overall analytic workload; concurrent data modification

### Data Required to Meet Business Needs

- Volume – amount of data to be managed
- Velocity – rate of data generation or change
- Variety – types of data to be managed
- Complexity – number of data sources and relationships; quality and structure of data

Copyright © BI Research, 2013

21

## Optimized Platforms: Analytic RDBMSs

Extend traditional RDBMSs with features designed specifically for analytic processing and new analytic techniques



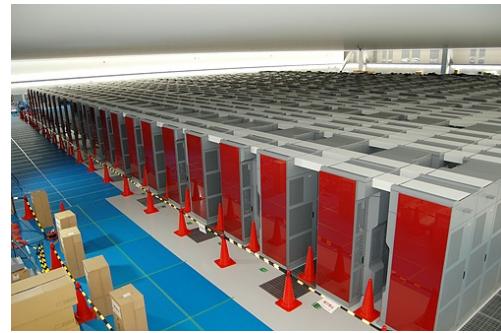
- Hardware exploitation
- Parallel computing
- New data types
- New storage structures
- Data compression
- Support for hybrid storage
- Intelligent workload management
- In-memory data
- In-memory analytics
- In-database aggregation & analytics

Copyright © BI Research, 2013

22

## Analytic RDBMSs: Hardware Exploitation

- Faster processors
- Multi-core processors
- Intelligent hardware
- 64-bit memory spaces
- Large-capacity disk drives
- Fast hard-disk and solid-state drives
- Hybrid storage configurations
- Scale-up/out parallel processing configurations
- Lower-cost hardware (blades, clusters)
- Reduced power and cooling requirements
- Packaged hardware/software appliances



Copyright © BI Research, 2013

23

## Analytic RDBMSs: New Storage Structures

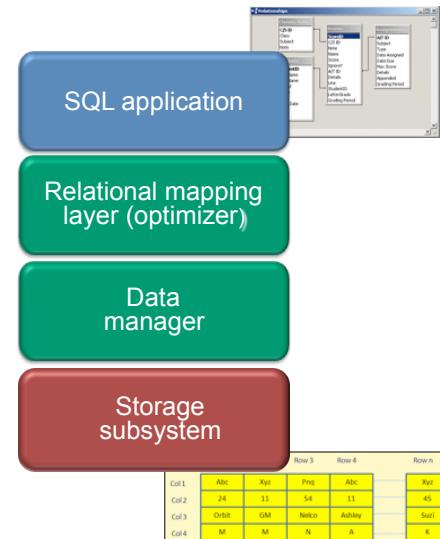
DBMS vendors are enabling new *physical storage structures* to improve performance, reduce storage requirements and support new types of analyses

Examples: compressed columnar, XML, time-series, multi-media

These enhancements and their implementation vary by vendor

It is important to recognize that *physical* storage structures should be independent of the *logical* data model and the data manipulation language (DML)

- True for the relational model and SQL
- Often not true for non-relational systems



Copyright © BI Research, 2013

24

## Analytic RDBMSs: Data Storage Options



### Large-capacity hard-disk drives (HDD)

- More economical, less reliable, slower performance, e.g., SATA drives in white-box H/W
- Often “short-stroked” to improve performance

### High-performance hard-disk drives (HDD)

- More expensive, more reliable, better performance, e.g., enterprise SAS drives

### Solid-state drives (SSDs)

- High and consistent performance
- Better reliability and more energy efficient
- Distinguish between commodity SSDs and enterprise SSDs

### Dynamic RAM (DRAM)

- Best performance - eliminates I/O overheads
- Use by in-memory computing systems

Copyright © BI Research, 2013

25

## Memory versus Storage

### Memory

- Data is directly addressable by a CPU via a memory bus
- Eliminates I/O overhead and provides fast access to data
- Types of solid-state memory:
  - Processor cache(s) – very fast, volatile data
  - Dynamic RAM – fast (nanoseconds), volatile data

### Storage

- Data is addressable via a device interconnect or network protocol
- Several types of storage for persisting data:
  - Commodity HDD: high capacity, less reliable, low cost (e.g., SATA)
  - Enterprise HDD: more reliable, higher cost (e.g., SAS)
  - NAND flash memory devices: fast, very reliable, high cost (e.g., PC SSD, enterprise SSD, flash storage array, hybrid SSD/HDD)

Copyright © BI Research, 2013

26

## What is In-Memory Computing?

A workload where all the data being processed is stored in a computer memory that is directly addressable via the CPU's memory bus

Provides high-speed performance for OLTP and BI workloads by eliminating I/O to storage devices

Especially beneficial for *interactive* and *iterative* BI analytic workloads

Several types of BI-related in-memory computing



## The Changing World of BI Analytics

### Advanced Analytics

- Improved analytic tools and techniques for statistical and predictive analytics
- New tools for exploring and visualizing new varieties of data
- Operational intelligence with embedded BI and BI automation

In-memory analytics

### Big Data Management

- Analytic relational database systems that offer improved price/performance and libraries of analytic functions
- Non-relational systems such as Hadoop for handling new types of data
- Stream processing/CEP systems for analyzing in-motion data

In-memory data

# Why In-Memory Computing for BI Analytics?

## Benefits

- Technology answer: Improved speed and performance, e.g., quickly run complex analyses on the fly
- Business answer: What if you could do ...? e.g., real-time fraud detection

## Considerations

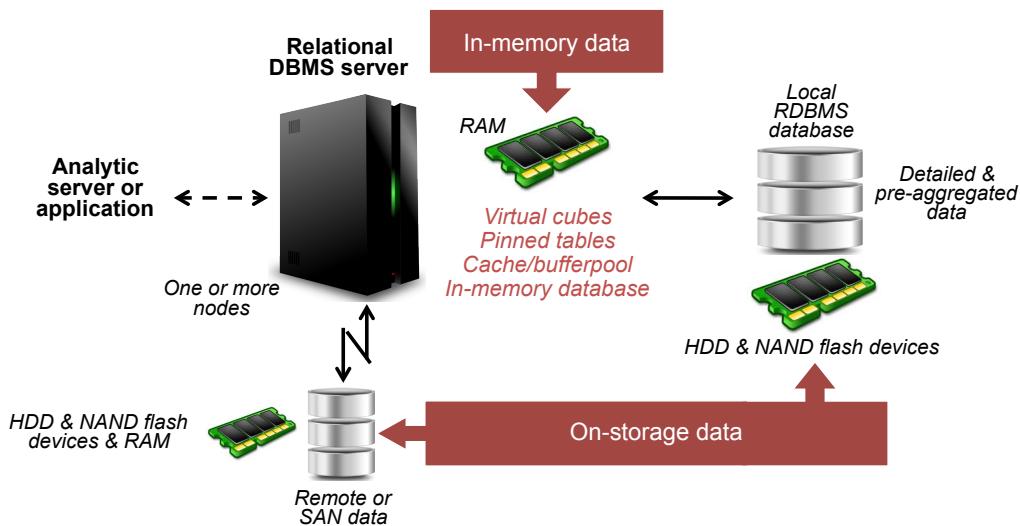
- Types of *in-memory data* and *in-memory analytics* and their benefits
- Relationship to in-database processing, e.g., in-database aggregation and in-database analytics



Copyright © BI Research, 2013

29

# In-Memory Data



Copyright © BI Research, 2013

30

# In-Memory Database Systems: Vendor Examples

## Relational DBMSs

EXASOL EXASolution

IBM solidDB and Informix Warehouse Accelerator

Kognitio Analytical Platform

Microsoft xVelocity and “Hekaton”

Oracle TimesTen

SAP HANA and Sybase ASE

VoltDB

Many vendors also support “pinned” tables

## Non-Relational DBMSs

Memcached

## Multi-Dimensional DBMSs

IBM Cognos TM1

## In-Memory Data Grids

Cloud platforms from Amazon, Google, IBM, Microsoft, VMware, etc.

Copyright © BI Research, 2013

31

# In-Memory Data: Important to Note

InformationWeek  
reports

## IT Pro Impact: In-Memory Analytics and Databases

It's taken six years or so of effort, but vendors have finally aligned the technology and the economics to bring capabilities once limited to telcos and trading floors to companies looking to maximize the value of their big data. The falling prices of DRAM make the business case that much stronger and mean that DBMS architects can, finally, dare to think beyond disk. We'll discuss the state of the market and connect platform choices with real-world deployment scenarios.

By Sreedhar Kajepeta

Report ID:55330712

July 2012 \$99

Merely converting a disk-based RDBMS to work with a RAM disk doesn't automatically make it an in-memory database.

Copyright © BI Research, 2013

32

# In-Database Technologies

## In-Database Aggregation

- Some RDBMSs support pre-aggregation to enhance BI performance
- Should be transparent to user – optimizer decides when to use aggregate
- Various names – materialized views, materialized query tables, etc.

## In-Database Analytics

- Brings the processing to the data rather than the data to the processing
- Consists primarily of predefined analytic functions – created by RDBMS vendor, third-party vendor, open source community, user developed

# In-Database Analytic Functions

Analytical functions stored in an RDBMS offer several benefits

- Users (e.g., data scientists) only need to understand what a function does and how to use it - they do not need to know how to develop the functions
- Functions can exploit the parallel processing capabilities of an RDBMS – moves the processing to the data rather than the data to the processing
- Important to understand the level of parallel processing and how a function is run, e.g., external to the RDBMS, in RDBMS protected memory, etc.

Several approaches to using in-database functions

- RDBMS built-in functions (arithmetic, string, date, statistical functions)
- Functions provided by a 3rd-party vendor, e.g., FuzzyLogix
- Open source functions, e.g., Apache Mahout, R
- Development options for creating user-defined functions (scripting language, Java, C++, SQL MapReduce, etc.)

# Analytic RDBMSs: Vendor Examples

## Traditional RDBMS Products

IBM DB2: PureData for Operational Analytics  
IBM Informix: Ultimate Warehouse Edition (with Warehouse Accelerator)  
Microsoft SQL Server: Parallel Data Warehouse  
Oracle Database: Exadata  
SAP Sybase ASE and IQ  
Teradata Database: Active EDW, DW Appliance, Extreme Data Appliance, etc.

## Other Solutions

EMC Greenplum Database and Distributed Computing Appliance  
HP Vertica Analytics Platform  
IBM Netezza: PureData for Analytics, DB2 Analytics Accelerator  
Kognitio Analytical Platform  
Oracle Exalytics (Oracle TimesTen & Oracle Essbase)  
ParAccel Analytic Platform  
SAP HANA  
Teradata Aster Database: MapReduce Platform, Big Data Analytics Appliance  
InfoBright, MySQL, PostgreSQL, etc.

Copyright © BI Research, 2013

35

# Data Warehouse DBMSs: Gartner 2013 MQ



Copyright © BI Research, 2013

36

# Optimized Platforms: Non-Relational Systems - 1

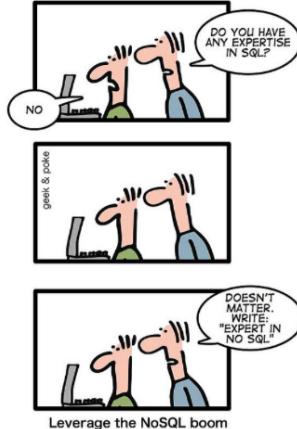
Several Internet companies developed their own non-relational (NoSQL or NewSQL) systems to support extreme data volumes

- Google example: Google file system, MapReduce, BigTable, BigQuery
- Main goal was the processing of large volumes of multi-structured data
- Several of these developments found their way into the open source community

Non-relational systems are not new, but modern versions are often open source

- Deployed on low-cost white-box hardware in a large-scale distributed computing environment
- Several types of systems & data stores
- Key industry focus area is Hadoop

## HOW TO WRITE A CV



Copyright © BI Research, 2013

37

# Optimized Platforms: Non-Relational Systems - 2

Many types of products, APIs and languages

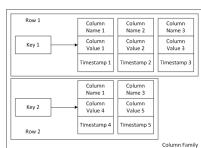
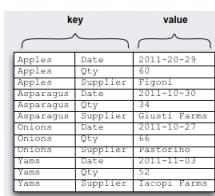
Volume ← → Complexity

Key/Value Pair

Column Family

Document

Graph



```
{
  "_id": "biking",
  "_rev": "A19EB7654",
  "title": "Biking",
  "body": "My biggest hobby is mountainbiking. The other day...",
  "date": "2009/01/30 18:04:11"
}
```



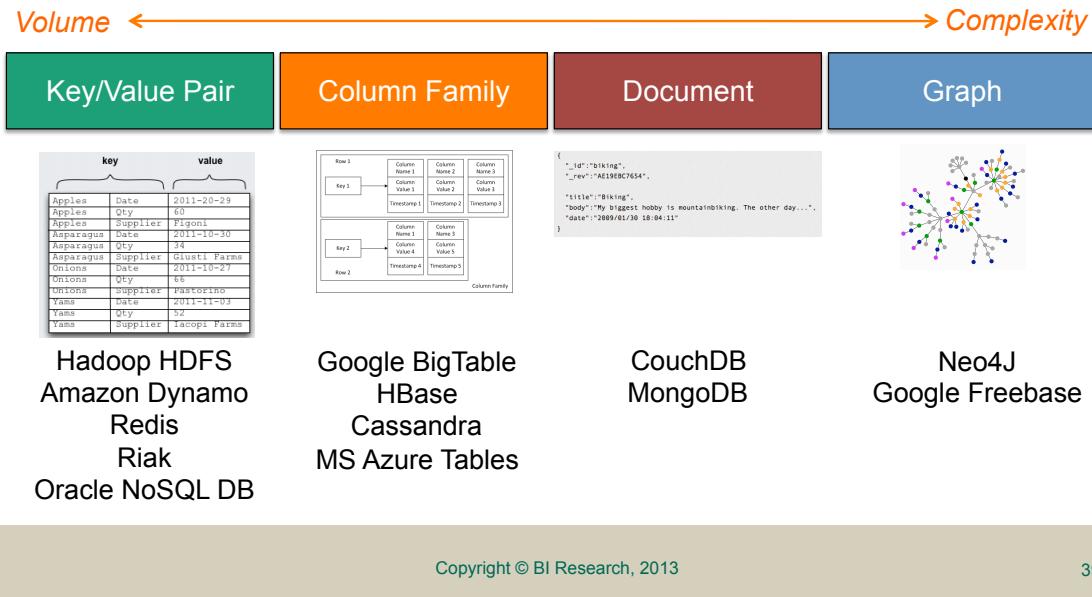
Can handle varieties of data and processing that are difficult to support using a traditional RDBMS

Copyright © BI Research, 2013

38

## Optimized Platforms: Non-Relational Systems - 2

Many types of products, APIs and languages



39



"A framework for running applications on a large hardware cluster built of commodity hardware."  
[wiki.apache.org/hadoop/](http://wiki.apache.org/hadoop/)

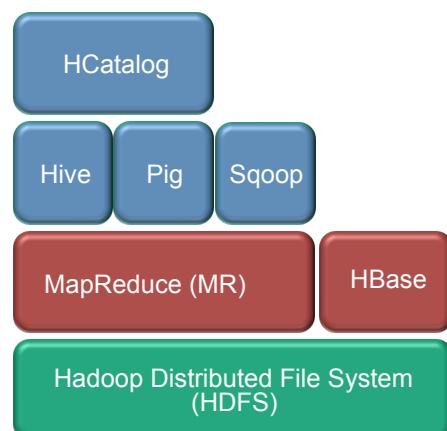
Provides a distributed file system (**HDFS**) that stores data across the nodes of the cluster to provide high performance

Includes a programming model called **MapReduce (MR)** where the processing is divided into small fragments of work that can be executed on any node in the cluster

**Hive** and **Pig** are high-level languages for MR development

Related components include HBase, Sqoop, HCatalog, Flume, Storm, Mahout, Impala, etc.

Major distributions come from Apache, Cloudera, Hortonworks and MapR

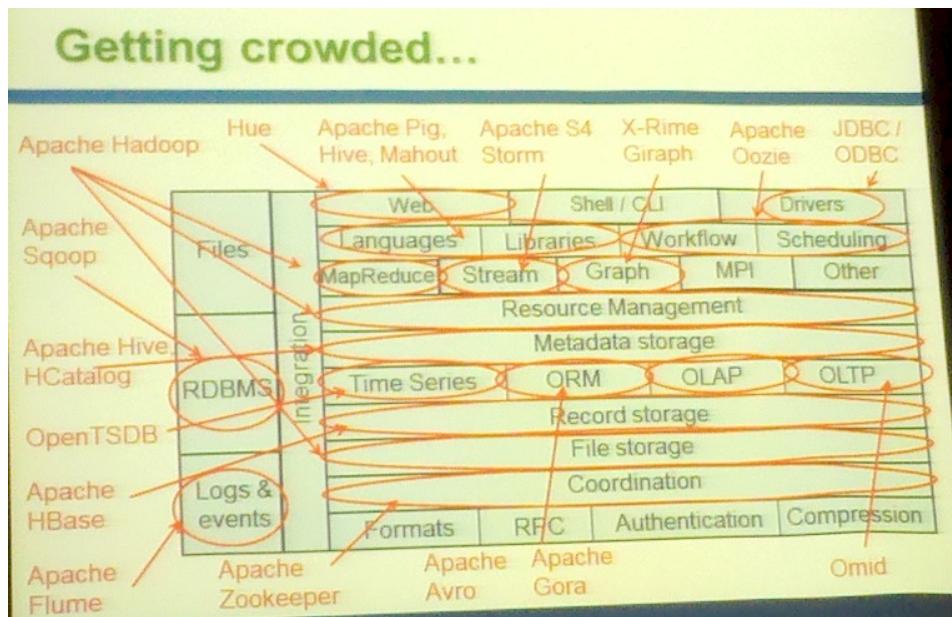


Copyright © BI Research, 2013

40



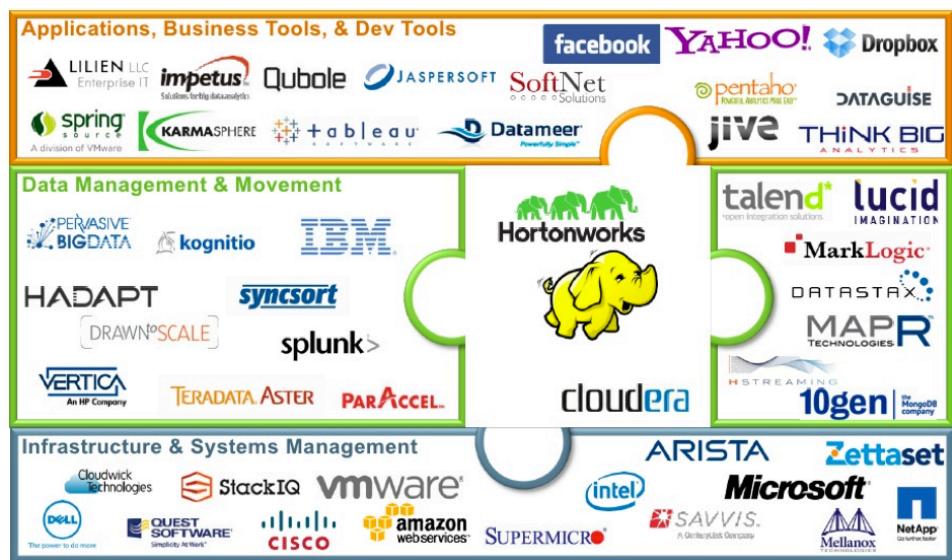
## Hadoop Components



Copyright © BI Research, 2013

41

## The Hadoop Ecosystem



Source: Hortonworks

Copyright © BI Research, 2013

42

## Data Management: Hadoop Option

### Distributions

Apache: Hadoop, HBase, Hive, Pig, Sqoop, Cassandra, Mahout

Cloudera: Enterprise Free (CDH), Enterprise Core, RTD for HBase, RTQ for Impala

HortonWorks Data Platform (includes Talend)

MapR: M3 (free), M5, M7

### Other solutions

EMC Greenplum: HD, Isilon NAS for HD, HD Distributed Computing Appliance (DCA)\*

Hadapt Adaptive Analytical Platform\*

HP AppSystem for Apache Hadoop

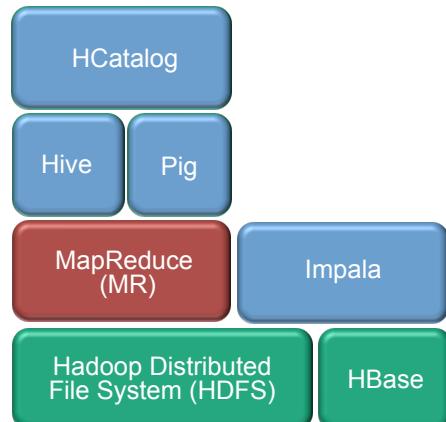
IBM InfoSphere BigInsights Basic and Enterprise Editions

Microsoft HDInsight Server & Service

Oracle Big Data Appliance

SAS High-Performance Analytics Server

Teradata Aster Big Data Analytics Appliance\*



\* Hybrid Hadoop/RDBMS system

Copyright © BI Research, 2013

43

## Data Management: Relational vs Non-Relational

Given the number of options and a fast changing marketplace comparisons are difficult

Focus is on analytic RDBMSs versus Hadoop HDFS – DBMS versus file system, which is an “apples to oranges” comparison

At a high-level, an analytic DBMS is suited to complex interactive workloads and Hadoop HDFS for batch processing of multi-structured

From an Hadoop perspective, HBase is becoming more important, but lack of SQL support is an inhibitor

- Hive support for HBase is in development but it still uses batch Map/Reduce (MR)
- Cloudera is developing Impala which supports Hive SQL syntax but eliminates MR – supports both HDFS and HBase

Workload suitability and performance are important, but development and administration effort, and tools support are also key considerations

Copyright © BI Research, 2013

44

# Data Management: Language Considerations

```

1. package org.myorg;
2.
3. import java.io.IOException;
4. import java.util.*;
5.
6. import org.apache.hadoop.fs.Path;
7. import org.apache.hadoop.conf.*;
8. import org.apache.hadoop.io.*;
9. import org.apache.hadoop.mapred.*;
10. import org.apache.hadoop.util.*;
11.
12. public class WordCount {
13.
14.     public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
15.         private final static IntWritable one = new IntWritable(1);
16.         private Text word = new Text();
17.
18.         public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, ReportProgress reporter) throws IOException {
19.             String line = value.toString();
20.             StringTokenizer tokenizer = new StringTokenizer(line);
21.             while (tokenizer.hasMoreTokens()) {
22.                 word.set(tokenizer.nextToken());
23.                 output.collect(word, one);
24.             }
25.         }
26.     }
27.
28.     public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
29.         public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, ReportProgress reporter) throws IOException {
30.             int sum = 0;
31.             while (values.hasNext()) {
32.                 sum += values.next().get();
33.             }
34.             output.collect(key, new IntWritable(sum));
35.         }
36.     }
37.
38.     public static void main(String[] args) throws Exception {
39.         JobConf conf = new JobConf(WordCount.class);
40.         conf.setJobName("WordCount");
41.
42.         conf.setOutputKeyClass(Text.class);
43.         conf.setOutputValueClass(IntWritable.class);
44.
45.         conf.setMapperClass(Map.class);
46.         conf.setCombinerClass(Reduce.class);
47.         conf.setReduceClass(Reduce.class);
48.
49.         conf.setInputFormat(TextInputFormat.class);
50.         conf.setOutputFormat(TextOutputFormat.class);
51.
52.         FileInputFormat.setInputPaths(conf, new Path(args[0]));
53.         FileOutputFormat.setOutputPath(conf, new Path(args[1]));
54.
55.         JobClient.runJob(conf);
56.     }
57. }
58.

```

```

##### define the explanatory variable with two levels:
##### 1-one or more parents smoke, 0-no parents smoke
parentsSmoke-as.factor(c(1,0))

##### NOTE: if we do parentsSmoke=c(1,0) R will treat this as
##### a numeric and not categorical variable

response<-cbind(yes=c(816,188),no=c(3203,1168))
response

##### fit the logistic regression model

smoke.logistic<-glm(response~parentsSmoke, family=binomial(link=logit))

##### OUTPUT

smoke.logistic
summary(smoke.logistic)
anova(smoke.logistic)

```

```

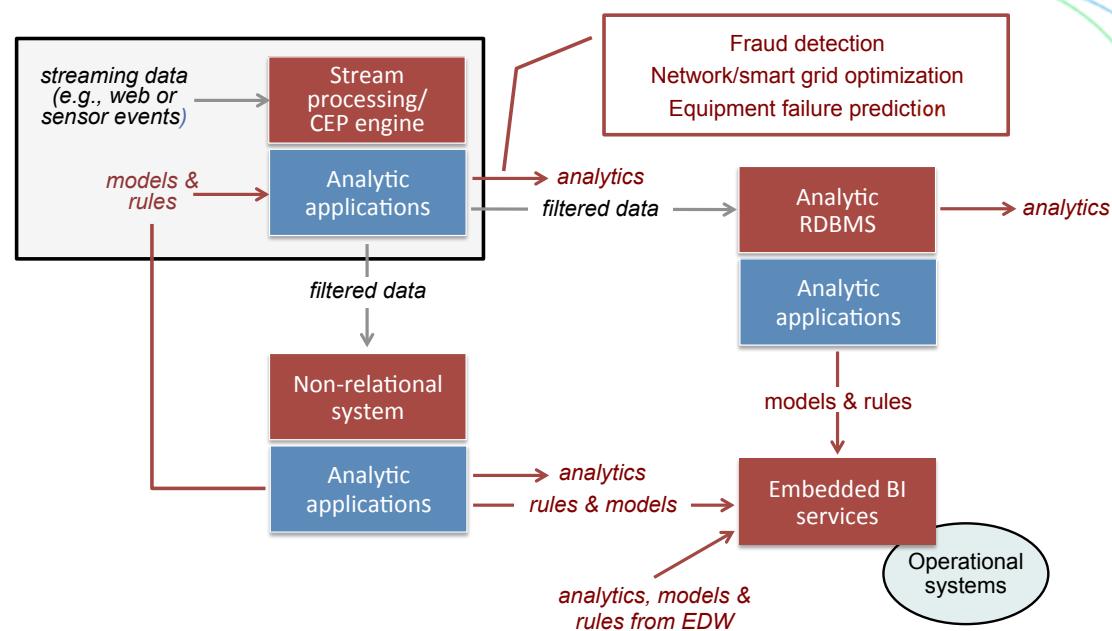
--selfjoin.pig
-- For each stock, find all dividends that increased between two dates
divs1    = load 'NYSE_dividends' as (exchange:chararray, symbol:chararray,
                                 date:chararray, dividends);
divs2    = load 'NYSE_dividends' as (exchange:chararray, symbol:chararray,
                                 date:chararray, dividends);
jnd     = join divs1 by symbol, divs2 by symbol;
increased = filter jnd by divs1::date < divs2::date and
           divs1::dividends < divs2::dividends;

```

Copyright © BI Research, 2013

45

# Optimized Platforms: Stream Processing/CEP



Copyright © BI Research, 2013

46

# Optimized Platforms: Cloud Option

Two approaches:

Virtual machine image – may be provided by the user, a DBMS vendor, or a cloud vendor (e.g., Amazon)

Database as a Service (DBaaS) offered by a public cloud vendor

Relational DBMS DBaaS

Amazon Redshift (ParAccel) and Relational Database Service (MySQL, Oracle, MS SQL Server)

Google Cloud SQL (MySQL) and BigQuery Service

HP Cloud Relational Database for MySQL

Kognitio Cloud

Microsoft Azure SQL Database

Oracle Cloud

Non-Relational DBMS DBaaS

Amazon DynamoDB, SimpleDB, ElastiCache

Microsoft Azure Tables and Blob storage

SalesForce database.com

Copyright © BI Research, 2013

47

## Data Management in the Cloud Example



Netflix's on-premises IT infrastructure was too fragile and the traditional operations model didn't respond fast enough to business needs

- Rapidly growing and highly variable data-center requirements
- Inability to automate data-center operations

As a result the company migrated from an on-premises environment to an Amazon Web Services infrastructure

It evaluated how the new environment would affect the IT infrastructure and redesigned applications as appropriate

Spreads its processing across many different Amazon data centers and regions to enhance reliability and availability

Different service environments are randomly taken offline to confirm that the environment can continue operating in the face of a resource failure

Conclusion: Netflix changed its approach because it recognized that the future of its business required a different way of doing things

Source: Netflix ([slideshare.net/adrianco/netflix-in-the-cloud-at-sv-forum](http://slideshare.net/adrianco/netflix-in-the-cloud-at-sv-forum)) and CIO Magazine

Copyright © BI Research, 2013

48

## Summary: Big Data Benefits

| Traditional Decision-Making Environment<br><i>(determine and analyze current business situation)</i> | Big Data Extensions<br><i>(provide more complete answers, predict future business situations, investigate new business opportunities)</i> |
|--|---|
| Integrated data sources  | Virtualized and blended data sources  |
| Structured data  | Multi-structured data   |
| Aggregated and detailed data (with limits)   | Large volumes of detailed data (no limits)  |
| Relational EDW with at rest data<br>Dimensional cubes/marts with at rest data                        | Non-relational stores with at rest data<br>Streaming/CEP systems with in motion data  |
| One-size fits all data management  | Flexible & optimized data management  |
| Reporting and OLAP   | Advanced analytic functions & predictive models   |
| Dashboards and scorecards  | Sophisticated visualization of large result sets  |
| Structured navigation (drill, slice/dice)  | Flexible exploration of large result sets   |
| Humans interpret results, patterns and trends  | Sophisticated trend and pattern analysis  |
| Manual analyses, decisions and actions   | Analytics & model-driven recommendations & actions  |

Copyright © BI Research, 2013

49

## Choosing the Right Solution

Organizations will likely use multiple analytic solutions and data management systems – the challenge is deciding which to use when and how to interconnect the systems

Copyright © BI Research, 2013

50

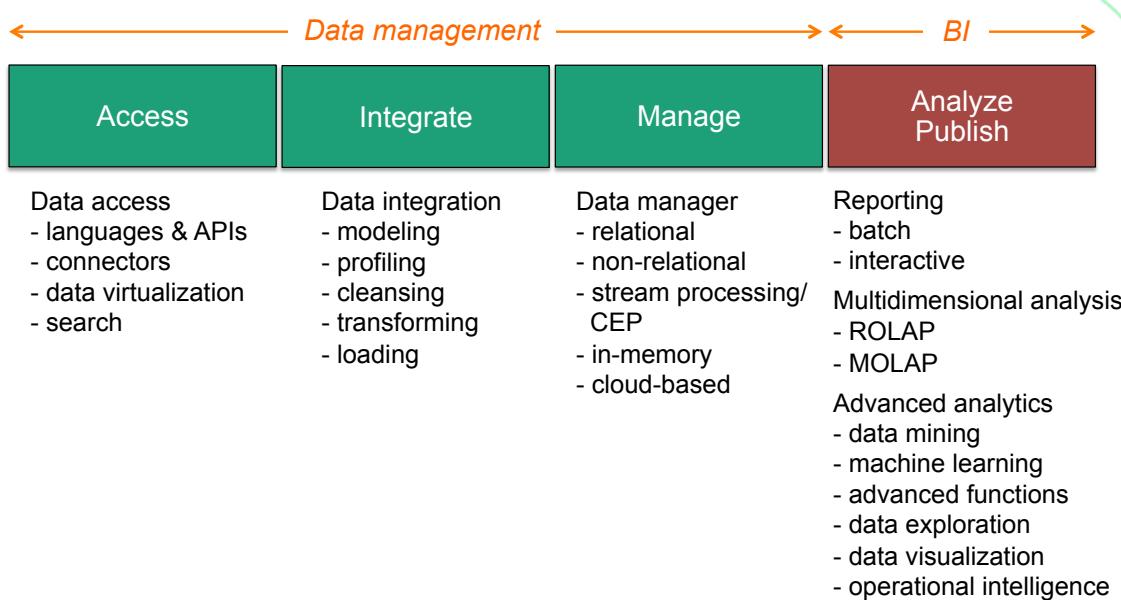
## Use Cases and Application Examples

| Use Case                         | Application Example  |
|----------------------------------|--|
| Real-Time Monitoring & Analytics | In-line fraud detection to reduce financial losses caused by stolen credit cards                                 |
| Near-Real-Time Analytics         | Next best customer offer to the channel to increase customer satisfaction & reduce churn                         |
| Data Integration Hub             | Collect and manage all sales-related detailed data (POS, web, supply chain) for down stream analysis             |
| Analytics Accelerator            | Offload & boost the performance of selected financial analyses to increase satisfaction/retention of key clients |
| New LOB Analytic Application     | Manage & monitor spot buying on web advertising exchanges  |
| Investigative Computing Platform | Evaluate the effectiveness of different social computing channels  |

Copyright © BI Research, 2013

51

## Software Selection: Some Key Options



Copyright © BI Research, 2013

52

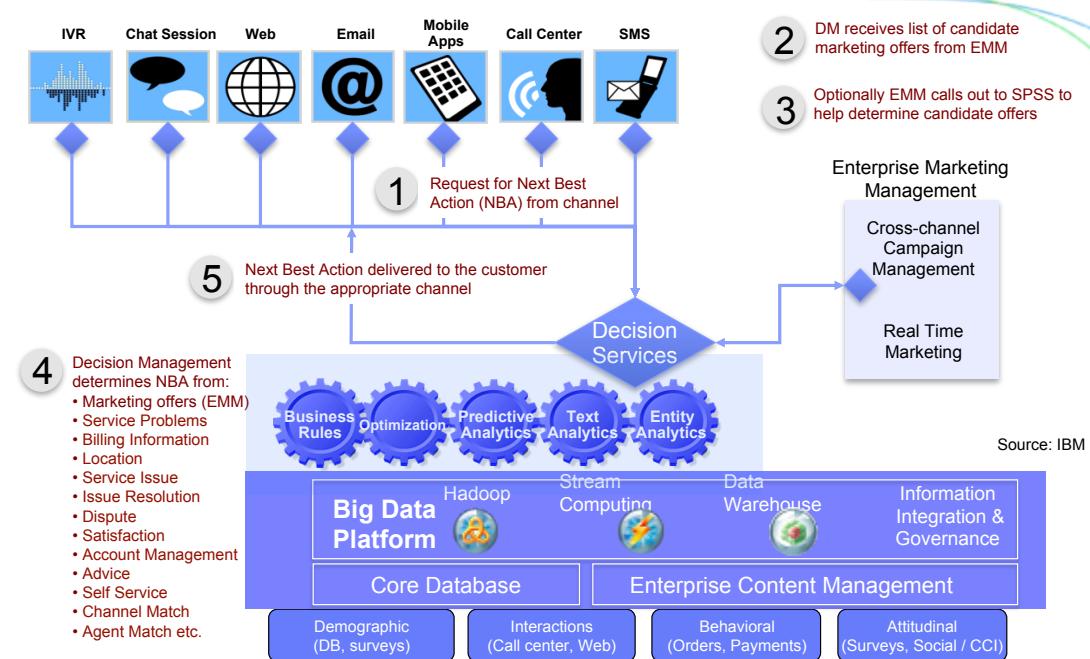
# Use Cases and Technologies

| Use Case                         | Stream Processing/<br>CEP system | Embedded BI<br>Services | Enterprise Data<br>Warehouse | Analytic Relational<br>DBMS | Hadoop System |
|----------------------------------|----------------------------------|-------------------------|------------------------------|-----------------------------|---------------|
| Real-Time Monitoring & Analytics | ✓                                | ✓                       |                              |                             |               |
| Near-Real-Time Analytics         |                                  |                         | ✓                            | ✓                           | ✓             |
| Data Integration Hub             |                                  |                         |                              |                             | ✓             |
| Analytics Accelerator            |                                  |                         |                              | ✓                           | ✓             |
| New LOB Analytic Application     |                                  |                         |                              | ✓                           | ✓             |
| Investigative Computing Platform |                                  |                         |                              | ✓                           | ✓             |

Copyright © BI Research, 2013

53

## Example Telco Provider: Real-Time Embedded BI



Copyright © BI Research, 2013

54

## Example

SEARS HOLDINGS



MetaScale

\$43 billion retail organization with over 4,000 stores (Sears and Kmart)

Numerous legacy systems with applications written in COBOL and Assembler (over 100 million lines)

Running out of capacity but at the current cost of \$3K-\$7K per MIP per year another solution had to be found

Requirements:

- Cost effectively manage increasing data volume
- Reduce the number of data warehouses and ETL jobs
- Reduce analytical processing times and provide intra-day analytics
- Capture and store all detailed transaction data (POS data, web clicks, supply chain events, etc.) for analysis
- Limit changes to existing user interfaces

*Primary source: Presentation by Dr. Phillip Shelley (CTO Sears Holdings and CEO MetaScale) at the Hadoop Summit, June 2012*

## Example

SEARS HOLDINGS



MetaScale

### Solution: Hadoop Data Hub and Analytics Accelerator

Enhanced pricing application

- Issue: only about 10% of the sales data is in the EDW; pricing models were taking 8 weeks to setup and run
- Hadoop MapReduce solution analyzes price elasticity based on 100% of the sales data
- Pricing models can now be run weekly (or daily if required)

Improved customized offers to loyal customers

- Issue: existing system was not scalable; only a small subset of the data could be analyzed
- Replaced 6,000 COBOL application with 400 lines of Pig and Java UDFs; implemented in 6 weeks
- Application can now be run multiple times per day per store per line item per customer – reduces impact of competitors such as Amazon

## Example

SEARS HOLDINGS



MetaScale

### **Solution: Hadoop Data Hub and Analytics Accelerator /cont.**

Reduce time to run batch BI applications

- Existing mainframe window of 3.5 hours was becoming insufficient to run 64 batch pricing jobs against 500 million rows of data
- Batch jobs were rewritten in Pig and run on data FTP'd to Hadoop from the mainframe and results FTP'd back to the mainframe
- Jobs run 100% faster (run in 8 minutes; FTP is the main overhead)

Reduce time to run batch and interactive BI applications

- Existing batch and interactive BI applications were taking too long to run and could handle only a subset of the data
- Data from over 50 sources now stored and analyzed on Hadoop
- Datameer is used for analytics
- Pig is used for ETL and for creating output for Excel

## Example

SEARS HOLDINGS



MetaScale

### **Conclusions**

Pleased with Hadoop's ability to run enterprise workloads – enables detailed data to be stored and analytics to be run that were not previously possible

Hadoop is only one component of the BI/DW ecosystem and strategy

Hadoop requires significant education and implementation effort and is lacking tools for enterprise integration

New Sears subsidiary (MetaScale) formed to help enterprises integrate existing systems with Hadoop because “75% of CEOs and CIOs don't even know what Hadoop is”

## Example: International Bank Trading Desk

### Solution: Analytic RDBMS as an Analytics Accelerator

This international bank offers a wide range of services to its over 40 million customers

One of the bank's trading desks uses the appliance for an analytic solution that handles the ad hoc analysis of billions of rows of detailed loan/bond data

Month-end loading was reduced from days to 2 hours

Key customer queries were reduced from 3-4 days to about 7 minutes

Appliance was treated as a black box by the IT group for compliance reasons



Copyright © BI Research, 2013

59

## Example



### Solution: Analytic RDBMS for New LOB Application

MediaMath is a leader in the billion dollar display advertising business

Provides a platform called TerminalOne that enables ad agencies and large-scale advertisers to identify, bid on, buy, and optimize ad impressions

Automatically matches each impression in real time with ads that are meaningful and relevant to users

Analyzes upwards of 15 billion ad impressions a day and calculates the fair value of more than 50,000 ads/sec



Copyright © BI Research, 2013

60

## Example @WalmartLabs

### Solution: Hadoop System for Investigative Computing

Phase 1 of the project was to use Hadoop and MapReduce to consolidate web site data (from 10 sites) for e-commerce analysis

Phase 2 involved providing users with Hive access to the web (and social) data for investigative purposes

- Demand for access and use of the system grew dramatically
- People forgot this was an experimental system!
- Requirements grew: larger cluster, resource management, SLAs, real-time data, metadata catalog
- Extended Hadoop to support a data hub containing 10 years of detailed data and reduce data stored in existing DW systems

*Primary source: Presentation by Stephen O'Sullivan and Jeremy King at the Hadoop Summit, June 2012*

## Barriers to Success

Educating IT and the business about the use cases and business benefits of big data

Lack of skills for enabling data science and investigative computing projects

Understanding and selecting the components that are required to build and support a big data analytics ecosystem

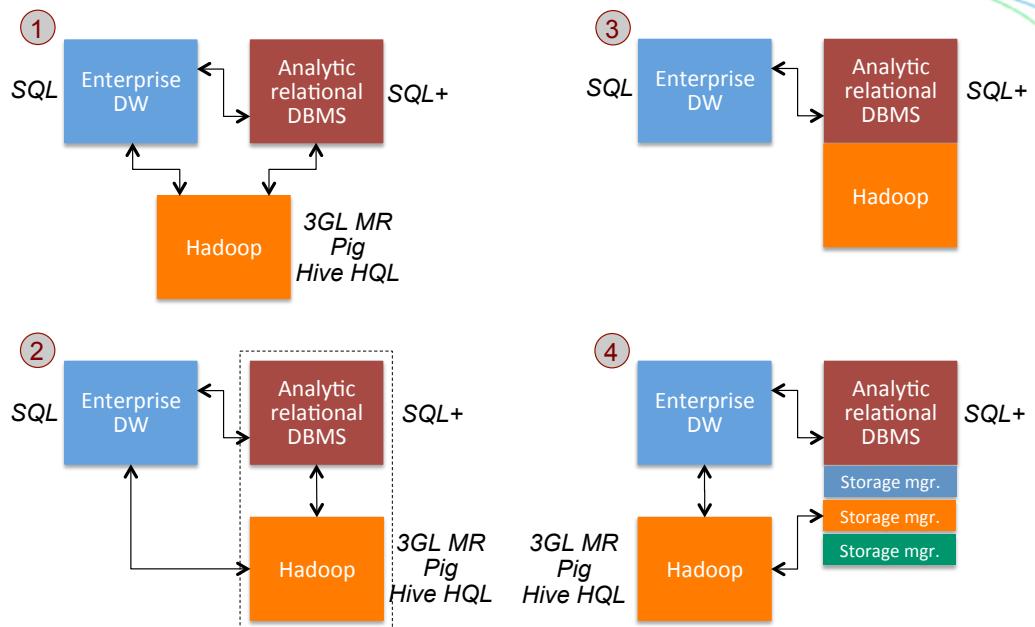
The immaturity of new non-relational systems and the level of IT development and administration resources and skills required for supporting them

The amount of data integration and level of data movement required in a big data environment

Developing data governance and data retention approaches to support the big data environment

Providing business users with a single and seamless user interface

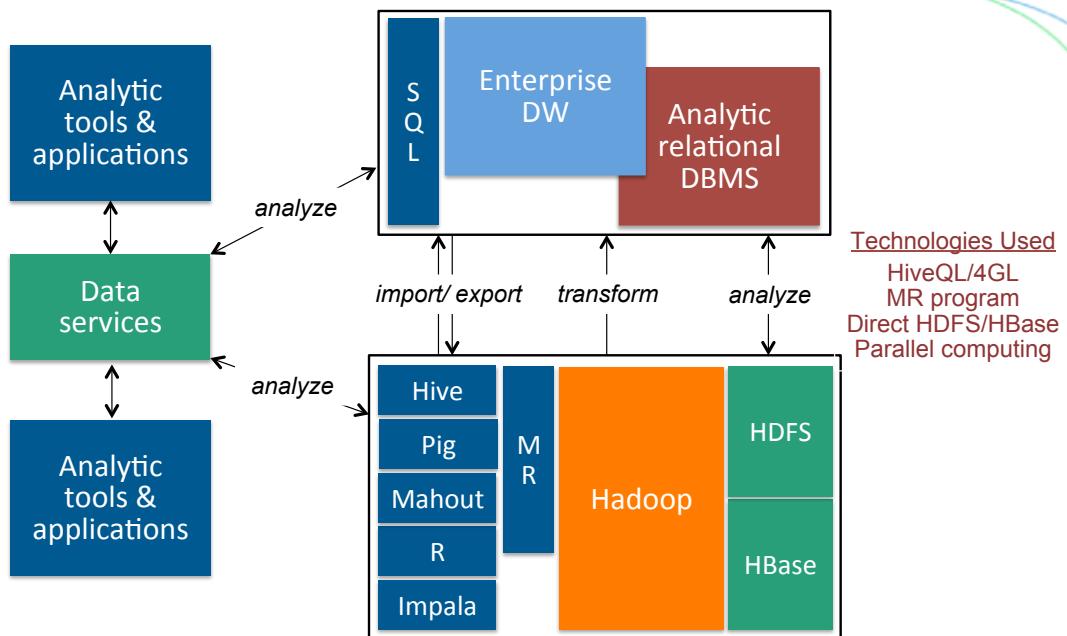
## Integration Example: RDBMS + Hadoop



Copyright © BI Research, 2013

63

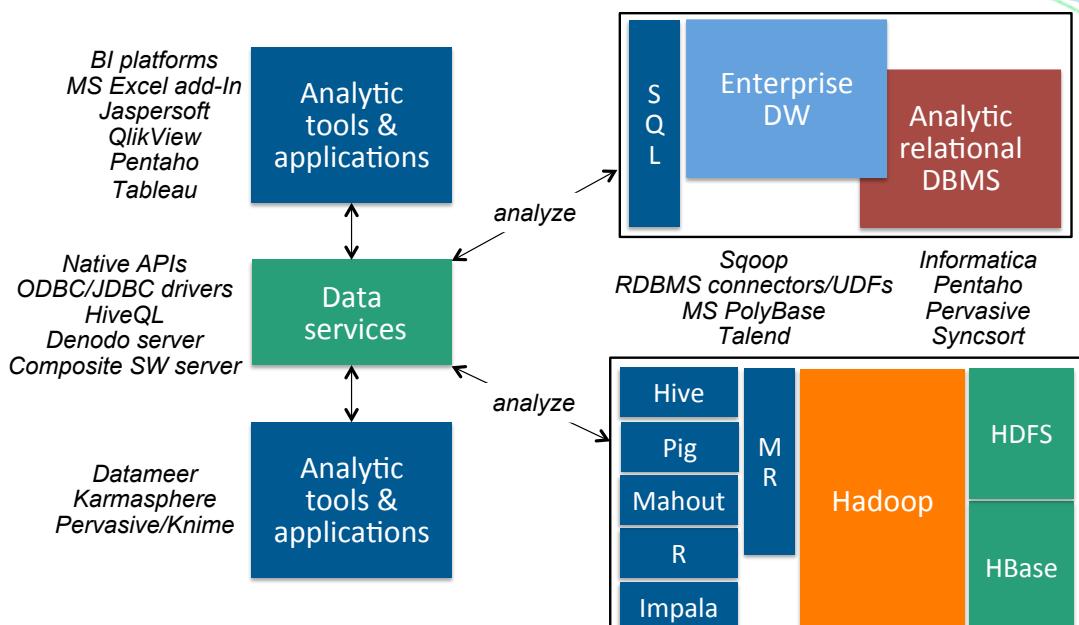
## Hybrid Option 1: Techniques



Copyright © BI Research, 2013

64

## Hybrid Option 1: Product Examples



Copyright © BI Research, 2013

65

## Conclusions

Big data represents a new and evolving analytics ecosystem – it is not one technology, but a set of overlapping technologies

Big data represents a new wave of analytic value for organizations

Big data initiatives at present are primarily LOB use-case driven

True value is gained from a hybrid of existing and new data systems

One size does not fit all and it is becoming impossible for any given vendor to satisfy all requirements

Integration with existing enterprise systems will be a key vendor differentiator

The vendors that win will be those that can best tackle cost, product and data integration, and/or advanced analytics needs

It is important for organizations to design an extended (or logical) data warehouse and advanced analytics ecosystem in order to be able to evolve and grow the BI environment

Copyright © BI Research, 2013

66