

# Big Data Overview

## Presented to DAMA



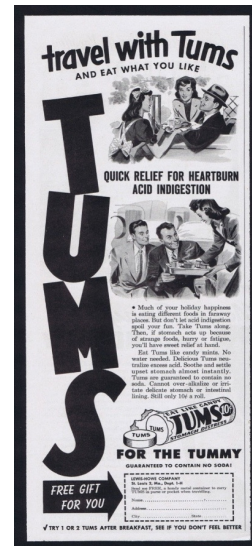
Lynn Hedegard  
 Lynn.Hedegard@us.ibm.com  
 Technical Sales Specialist  
 West Region

28<sup>th</sup> of May, 2015

© 2015 IBM Corporation

## When Was CRM Invented?

- As early as the 1930's some marketing companies were using coupons to track and measure response to printed advertising
- The coupon was usually located at the bottom of a full page ad.
- They said they were going to give you something for free...
- ... but they were actually measuring the effectiveness of the ad.
- Some of these ads are now collector's items



## Real-Time CRM in the Social World (Meet Lisa)



## How Did They Do That?

How Did  
They Do That?

## Technologies Used in Previous Example

- **They used Legacy Technology**
  - Online Transaction Processing Systems (OLTP)
  - Operational Data Stores (ODS)
  - Data Warehouse (DW)
- **They used “Newer” Technology**
  - Telecommunication
  - Mobile Computing
  - Wireless
  - Business Intelligence
- **They used Big Data Technologies**
  - Hadoop
  - NoSQL Database Systems
  - Data Science
  - Real-Time Analytic Processing

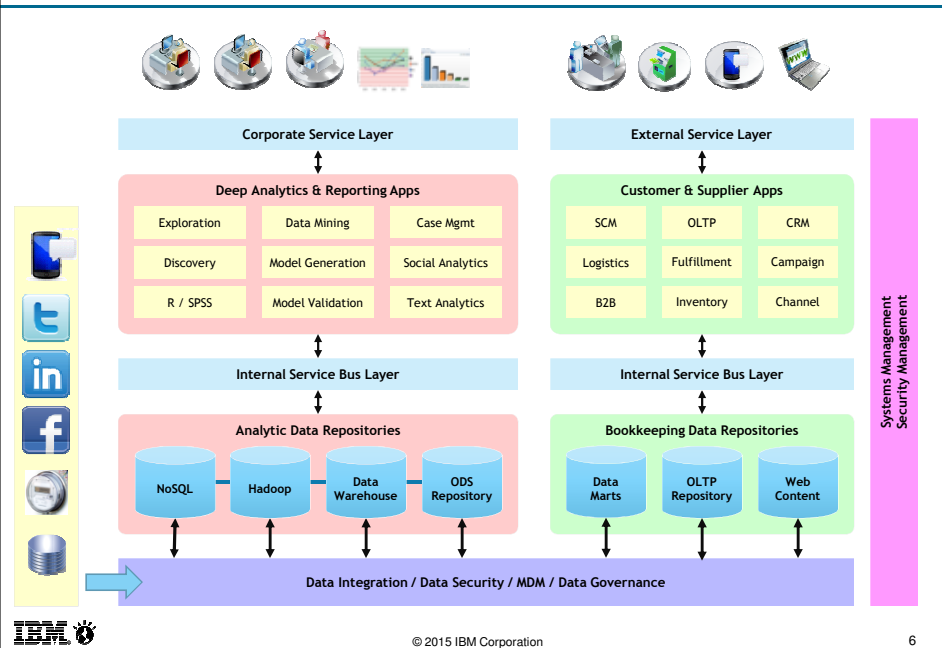
We will focus on this topic of Big Data



© 2015 IBM Corporation

5

## Reference Architecture for the On-Line Enterprise



6

## Problem Statement

# Problem Statement

(or why do we need Big Data)



© 2015 IBM Corporation

7

## Problem Statement – Complex Business Environment

- **The Local Environment is Complex:**
  - A single large retail store (1.5 million SKUs)
  - Large manufacturing floor (~6 million parts)
  - Vegas Casino (20 million card carrying customers)
- **The Global Environment is Complex:**
  - The number of variables affecting business performance is huge.
  - US citizens (*source: google population*)
    - 300+ Million total
    - (21M+ teenagers) + (40M+ in their 20's) (that's a lot of calls & text messages!)
  - The interrelationships between these variables is very complex (e.g.,  $N^2$  problem)
    - Multiple customer touch points
    - Multiple suppliers & distribution methods
    - Market forces (cost of raw goods & services, pricing dynamics, supply/demand)
- **Working Premise:** Few people in the enterprise can make “good” Operational Decisions – consistently & quickly
  - Few people can “see” all the necessary data.
  - Few people can “analyze” all the necessary data.
  - Few people understand all the inter-relationships between business variables.

Businesses can no longer tolerate inconsistent Business Processes



© 2015 IBM Corporation

8

But wait ...

... We're just getting started



© 2015 IBM Corporation

9

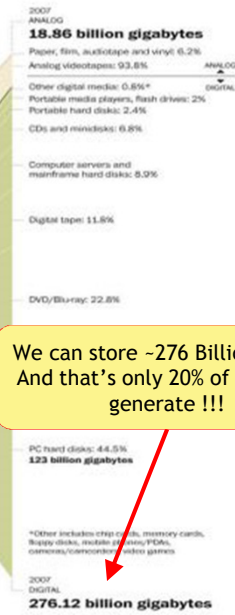
## Global Capacity to Store & Communicate Information

Martin Hilbert conducted detailed research into the amount of data that can be stored, analyzed, and communicated in recent history ...

... a tipping point occurred around 2001

The World's Technological Capacity to Store, Communicate, and Compute Information; Martin Hilbert; Science Magazine, April 2011; Vol. 332 no. 6025 pp. 60-65

beginning of the digital age



## How Did the Big Data Wave Get Started?

Prior to 2000, not many people were using the internet.

- Modems were still slow.
- Video resolution was poor.
- Computer interfaces were clumsy.
- The average person might use e-mail, but not much more
- Xerox created the “*window interface*” and the “*mouse*” at the **Palo Alto Research Center (PARC Place)**, but not many people even know about it.
- I had a large black & white X-Terminal on my desk
- I was in heaven!!!



© 2015 IBM Corporation

11

## Society Experiences a “Tipping Point” in 2001

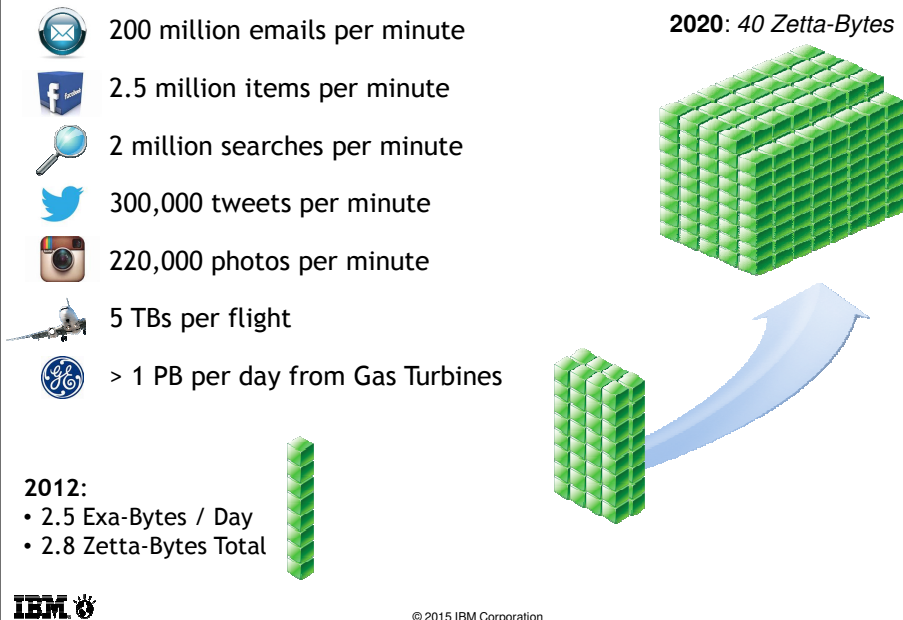
- Significant advances in technology cause a tipping point
  - Telecommunication capabilities advanced to the point where an “*average person*” could tolerate using modems.
  - We move from command-line interfaces to window/mouse based interface (created by Xerox).
  - The average person starts using the internet for ...
    - E-Mail, Instant Messaging, Shopping,
    - and of course ... Web Surfing.
  - Most adults (and many teenagers) have a cell-phone
  - Internet companies like, AOL, Yahoo, MicroSoft, & Google get overwhelmed with large quantities of DATA
  - **Conventional database system were ill-suited of storing this data**
    - Data was semi-structured
    - Data was not clean
    - And there was a LOT of this data.



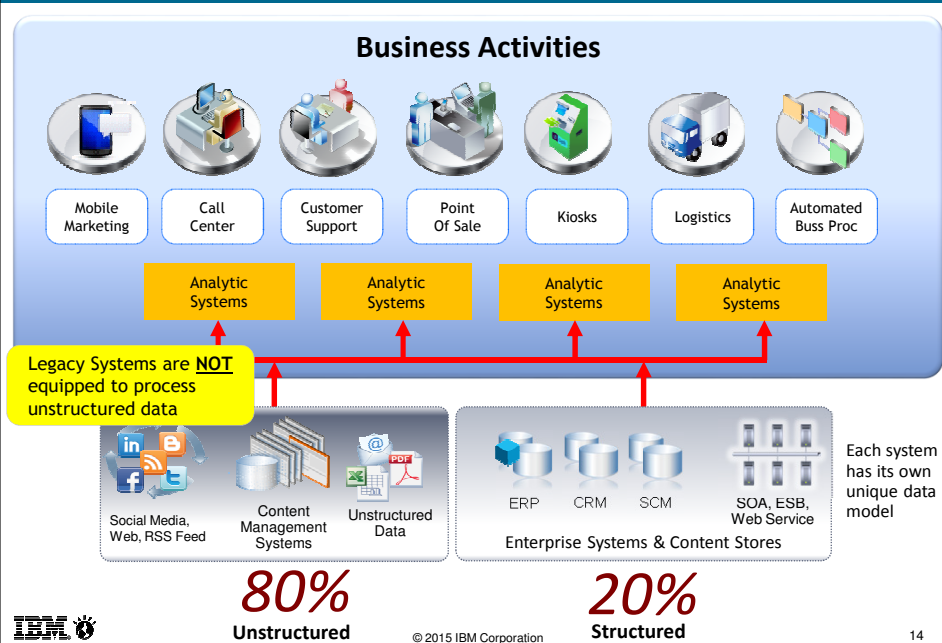
© 2015 IBM Corporation

12

## Growth of New Data

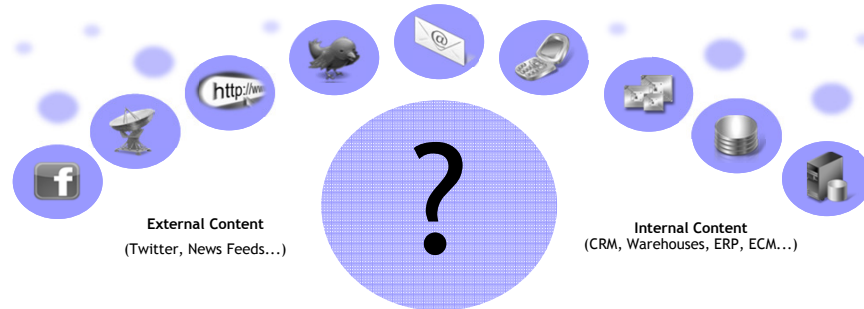


## 80% of Data Growth is Unstructured Data



## Data is Everywhere

Data is everywhere, but under-utilized at the point of impact



"I can't find the right answers fast enough to support my customers."



"I need to know what complications to expect, if I use this new medicine before the surgery."



"I know there's value in my data. How do I convert all this information into economic value?"



"I don't know what I don't know!!!  
Where is my business exposed?"



"It was a cold day in the city. The rain fell like tears on the pavement. I had a bad feeling about this."



© 2015 IBM Corporation

15

## Big Data Discussion

Big Data  
Discussion



© 2015 IBM Corporation

16



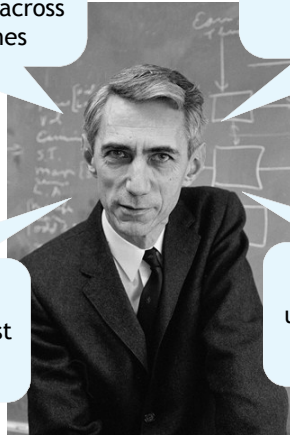
## Lets Make a Business Case for Big Data

Back in 1948 ... I showed you how to digitize voice & music (e.g., MP3, Compact Disk) **AND** how to send all this stuff across telecommunication lines

... I also showed you the **BUSINESS VALUE** of Data & Analytics

Next, I collaborated with Ed Thorp to create the first wearable computer ...

... We went to Vegas and used the wearable computer to win at Blackjack



Claude Shannon (Father of Information Theory)  
1948: First Paper on Information Theory



© 2015 IBM Corporation

17

## We Need To Think Beyond Traditional Sources Of Data

### Transactional & Application Data



Volume  
Structured  
Throughput

### Machine Data



Velocity  
Semi-structured  
Ingestion

### Social Data



Variety  
Highly unstructured  
Veracity

### Enterprise Content



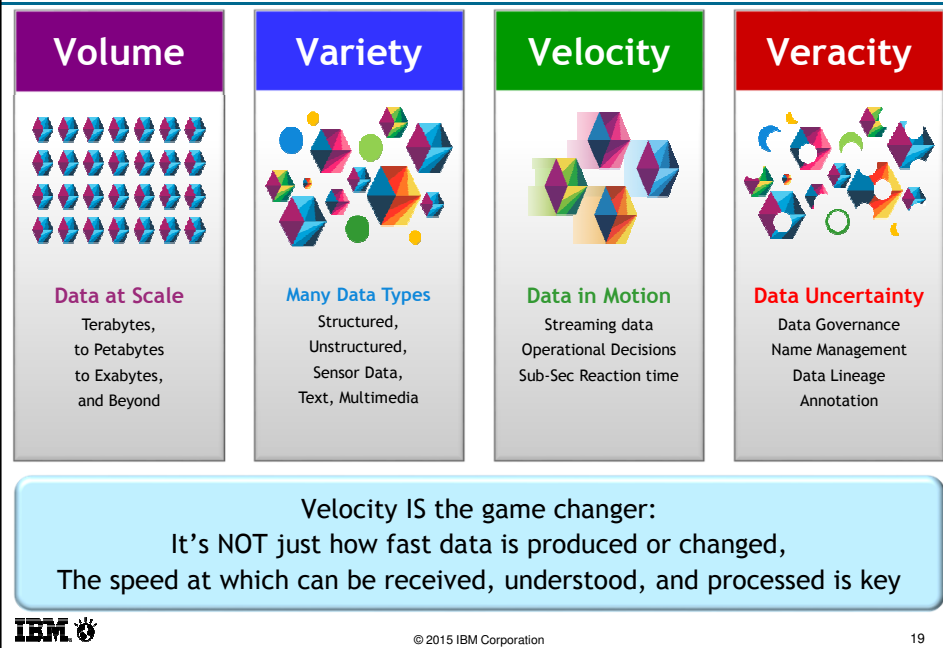
Variety  
Highly unstructured  
Volume



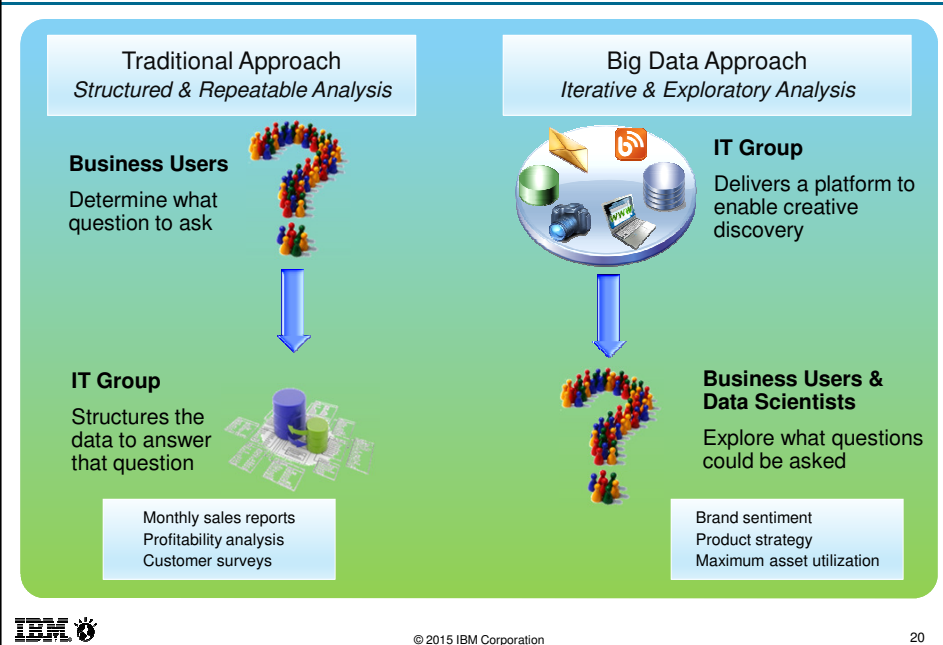
© 2015 IBM Corporation

18

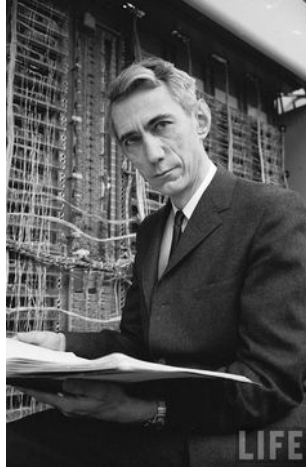
## We Need To Think About the Attributes of Data



## Merging the Traditional and Big Data Approaches



## Information & Uncertainty



*"Information is the resolution of uncertainty"*  
 Claude Shannon (Father of Information Theory)  
 1948: First Paper on Information Theory



© 2015 IBM Corporation

21

1 in 3

Business Leaders frequently make decisions based on information they don't trust, or don't have

1 in 2

Business Leaders say they don't have access to the information they need to do their jobs

83%

of CIOs cited "Business intelligence and analytics" as part of their visionary plans to enhance competitiveness

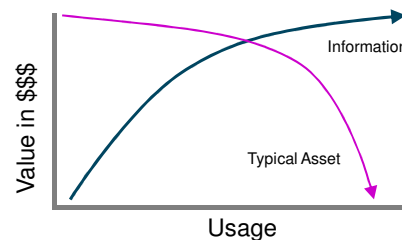
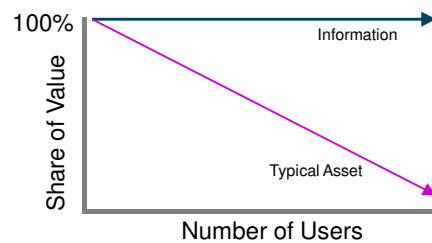
60%

of CEOs need to do a better job capturing and understanding information rapidly in order to make swift business decisions

**Analysis of Data, and Communication of insights reduces uncertainty.**

## Measuring the Value of Data – (The 7 Laws)\*

- **Law One:** Information is Infinitely Shareable.
- **Law Two:** The value of data increases with its use.
- Most other assets "wear out" with extended use.



\* Measuring The Value Of Information: An Asset Valuation Approach; Moody, Daniel & Walsh, Peter  
 European Conference on Information Systems (ECIS'99)

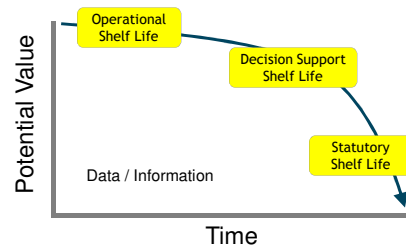


© 2015 IBM Corporation

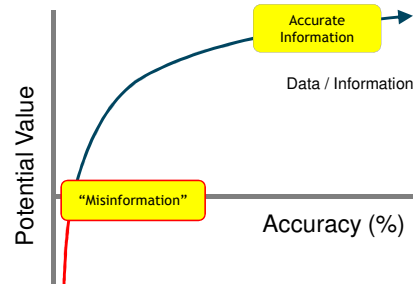
22

## Measuring the Value of Data (cont.)

- **Law Three:** The value of data tends to depreciate over time.
- “Old News”



- **Law Four:** The value of data Increases with Accuracy.
- In fact, **inaccurate** data has a **NEGATIVE** value.

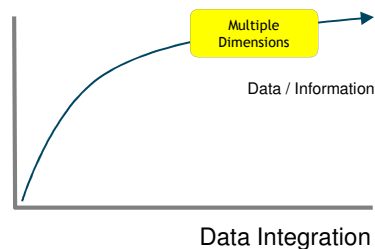


© 2015 IBM Corporation

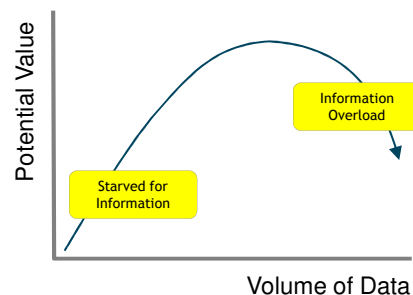
23

## Measuring the Value of Data (cont.)

- **Law Five:** The value of data increases when combined with other data.
- Combined value may exceed the sum of the parts



- **Law Six:** More data is **NOT** always a good thing
- “Information Overload”



- **Law Seven:** Information does not Deplete when used.

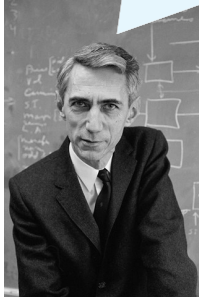


© 2015 IBM Corporation

24

## Data & Information is a Valuable Asset

$$\text{Growth} = E_e[\log^d W] - H(E|\vec{G}) - D_{KL}(\vec{P}(e|g) \| P(e|m)) - I(E; \vec{G})$$



Huh !?!?!

- Growth of the Enterprise is proportional to the amount of **DATA acquired & analyzed**, and how well you **communicate** that insight throughout the Enterpris (i.e., A reduction in Entropy)
- “*The more you know ... the more you grow*”
- Just look at the market value of companies that capture extreme volumes of data. (QED)
- **Key Point:** Data, and the Insights gleaned from analysis of that data is one of the most valuable assets in you Enterprise
- If you’re **NOT** gathering data about your business processes and other related events, you’re throwing money away!!!



© 2015 IBM Corporation

25

## A Profile for Big Data

- **Information Flow**
  - **Human generated data**
    - Internet Search (Google, Yahoo, Bing)
    - Social Media Posts, Tweets, Blogs, etc.
  - **Machine Generated Data**
    - Machine Generated Logs (Linux, Network, Manufacturing, etc.)
    - Internet of Things (IOT) (Telematics, Appliances, etc.)
    - Infrastructure (e.g. Wireless, Toll Roads, Cameras, etc.)
- **Information Retention**
  - Amount of data retained doubles every 3 years
  - 2.5 Exabytes in 1986 (~1% digitized)
  - 300 Exabytes in 2007 (~94% digitized)
  - **About 80% of data that has been captured will be discarded !!!**
- **Information Analysis**
  - We can only analyze a small portion of the data retained
  - The majority of data is Semi-structured -or- non-structured



© 2015 IBM Corporation

26

## What is Hadoop

# What Is Hadoop

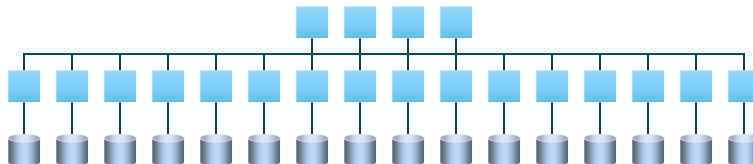


© 2015 IBM Corporation

27

## Two Key Elements of Hadoop

- **Parallel Execution Framework (MapReduce)**
  - How Hadoop understands and assigns work to the nodes (machines)
- **Hadoop Distributed File System (HDFS)**
  - Where Hadoop stores data
  - The file system spans all nodes in a Hadoop cluster
  - Links together the file systems on many local nodes to make them into one big file system



© 2015 IBM Corporation

28

## Understanding MPP

- Imagine I give you the phone book for Los Angeles
- I ask you to make a list of all the entries where...
  - First Name = "John"
  - Area Code = "323"
- You will need to look at ALL of the pages sequentially
- Now imagine we have 10 clones ready to do the same job
- We divide the phone book into 10 equal size sub-sections
- The "team" operates in parallel.
- The "team" will be able to make the list 10 times faster
- Note:
  - Data Warehousing, Hadoop, & NoSQL benefit from MPP
  - OLTP does NOT benefit from MPP



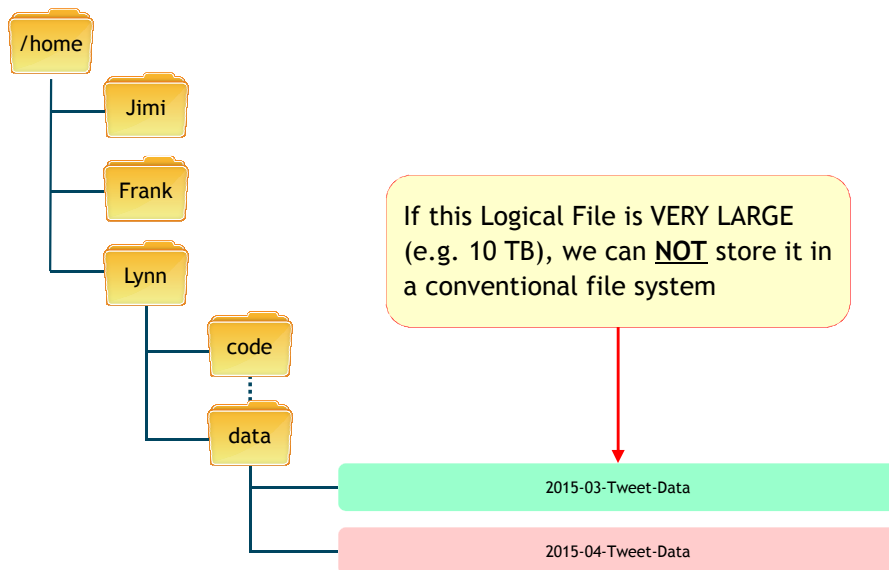
If I had 1000 clones,  
the work would finish  
1000 times faster ...  
... *Linear Scalability*



© 2015 IBM Corporation

29

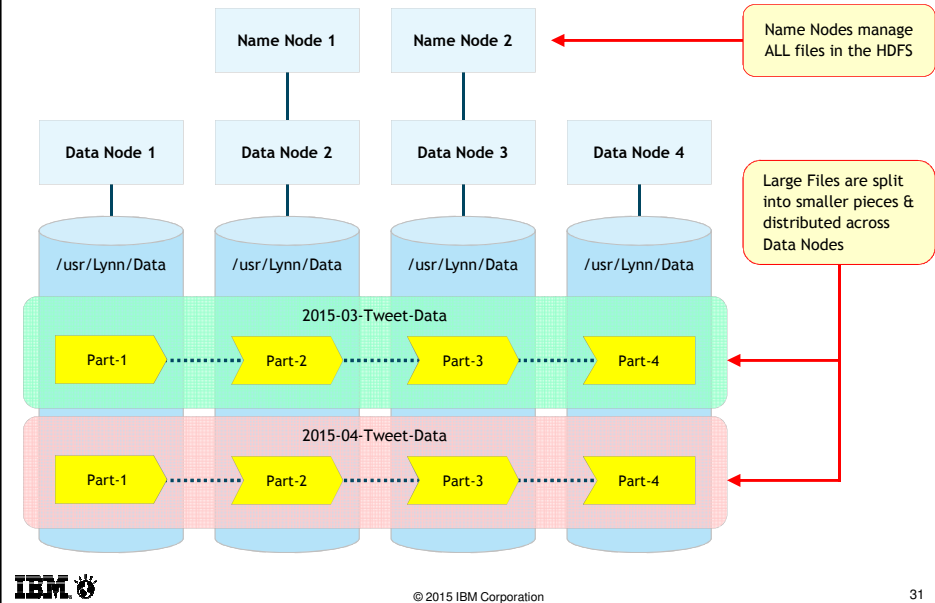
## Large Files



© 2015 IBM Corporation

30

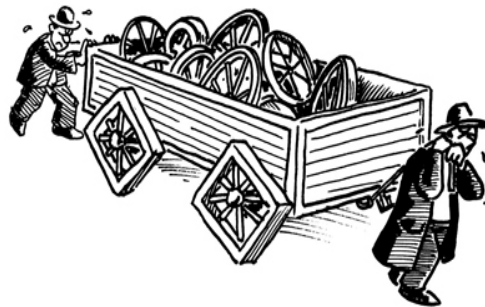
## Hadoop File System



## The Right Tool for the Job

*"People shouldn't get ideas about turning off their relational systems and replacing them with Hadoop..."*

*... As we start thinking about big data from the perspective of business needs, we're realizing that Hadoop isn't always the best tool for everything we need to do, and that using the wrong tool can sometimes be painful."*



**Ken Rudin**  
Head of Analytics at Facebook



© 2015 IBM Corporation

32



## Big Data Methodologies

# Big Data Methodologies

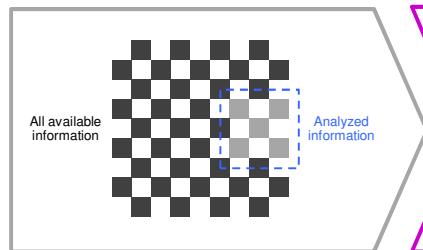


© 2015 IBM Corporation

33

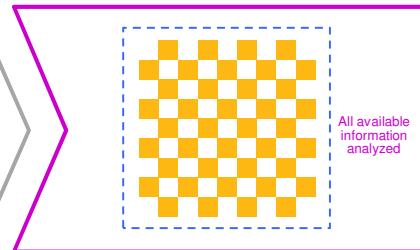
## Leverage More of the Data Being Captured

### TRADITIONAL APPROACH



Analyze small subsets of information

### BIG DATA APPROACH



Analyze **ALL** information

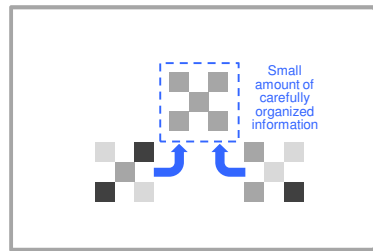


© 2015 IBM Corporation

34

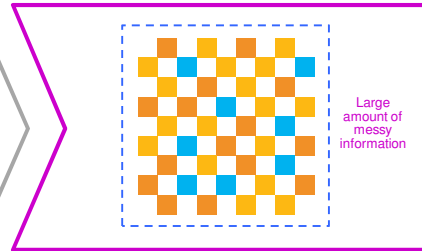
## Reduce Effort Required to Leverage Data

### TRADITIONAL APPROACH



Carefully cleanse information **before** any analysis

### BIG DATA APPROACH



Analyze information as is, cleanse as needed

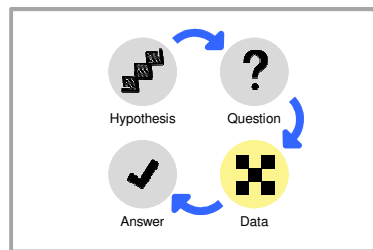


© 2015 IBM Corporation

35

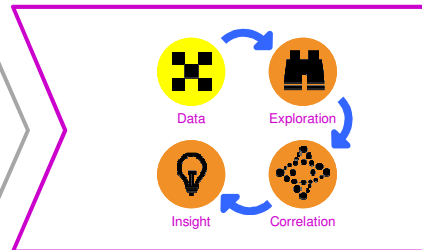
## Data Leads the Way

### TRADITIONAL APPROACH



Start with hypothesis and test against selected data

### BIG DATA APPROACH



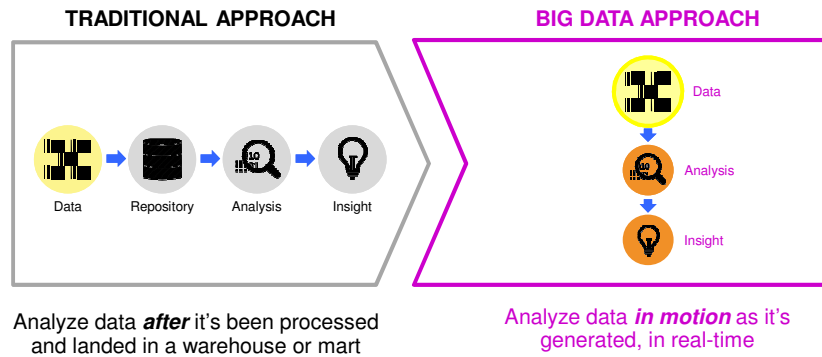
Explore **all** data and identify correlations



© 2015 IBM Corporation

36

## Leverage Data as it is Captured



© 2015 IBM Corporation

37

## NoSQL Databases

NoSQL  
Databases



© 2015 IBM Corporation

38

## NoSQL Databases

- NoSQL: A term describing a class of Big Data Databases
- Its strange to define something by what it is NOT ?!?!?
- The term; “NoSQL” was created by Johan Oskarsson as a hash-tag to advertise an upcoming Meetup about BigData
- NoSQL does NOT mean NO SQL ...  
... it means; “Not Only SQL.”
- Characteristics of NoSQL
  - Non-Relational
  - Cluster Friendly
  - Schema-less
  - Designed to enable modern Web



© 2015 IBM Corporation

39

## Different Categories of noSQL Databases

NoSQL Category	Use this when....	Application Examples	Vendors
<b>Document</b> 63% revenue share*	<ul style="list-style-type: none"> <li>• Schema is not well defined</li> <li>• Schema is very likely to change, need to maintain flexibility</li> <li>• Commonly described with JSON</li> </ul>	<ul style="list-style-type: none"> <li>• Frequently changing product catalogs</li> </ul>	<ul style="list-style-type: none"> <li>• Cloudant**</li> <li>• MongoDB</li> <li>• Couchbase</li> <li>• MarkLogic</li> </ul>
<b>Key-Value</b> 24% revenue share*	<ul style="list-style-type: none"> <li>• Your data is not highly related</li> <li>• All you need is basic Create, Read, Update, Delete (CRUD)</li> <li>• Rapid Scaling for simple data collections</li> </ul>	<ul style="list-style-type: none"> <li>• User Sessions</li> <li>• Shopping Cart</li> </ul>	<ul style="list-style-type: none"> <li>• Redis</li> <li>• Aerospike</li> <li>• AWS (DynamoDB)</li> <li>• Basho Technologies (Riak)</li> </ul>
<b>BigTable Columnar</b> 9% revenue share*	<ul style="list-style-type: none"> <li>• High volume, low latency write</li> <li>• Big Data, sparse data</li> <li>• Need compression or versioning</li> </ul>	<ul style="list-style-type: none"> <li>• Telco, heavy ingest, petabyte scale</li> <li>• User Activity logs</li> <li>• Sensor data</li> </ul>	<ul style="list-style-type: none"> <li>• HBase (Hadoop)**</li> <li>• BigTable</li> <li>• Cassandra</li> </ul>
<b>Graph DB</b> 4% revenue Share*	<ul style="list-style-type: none"> <li>• Your data looks like a graph</li> <li>• Have highly interconnected data, need to trace relationships</li> </ul>	<ul style="list-style-type: none"> <li>• Website Purchase Recommendations</li> <li>• Social Network Processing</li> </ul>	<ul style="list-style-type: none"> <li>• Titan**</li> <li>• Neo Technology (Neo4J)</li> </ul>

\* Source: IBM study 2013 estimated by splitting total noSQL revenue (\$288m) by ratio of top 10 vendors reported 2013 revenue.  
Total 2013 noSQL database revenue estimated \$343m

\*\* IBM Solutions of Choice.



© 2015 IBM Corporation

40

## Key – Value Database

- Based on Amazon's Dynamo technology
- Data Model:
  - Key; Usually a scalar value
  - Value; Usually a scalar value
- Good for VERY Fast lookup
  - User-ID : Password
  - Account : Max Charge
  - SKU : Cost
- Examples: Dynamite, Voldemort, Tokyo



© 2015 IBM Corporation 41

41

## Column Based DB

- Based on Google's BigTable technology
- Data Model
- Example: HBase, Hypertable, Cassandra

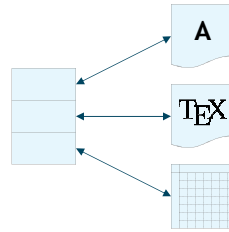


© 2015 IBM Corporation 42

42

## Document Database

- Inspired by Lotus Notes
- Data model:
  - Collections of K-V collections
- Examples:
  - CouchDB
  - MongoDB
  - Riak

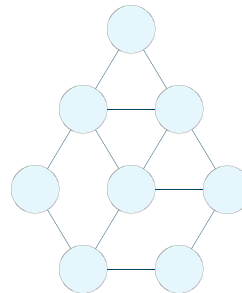


© 2015 IBM Corporation

43

## Graph Database

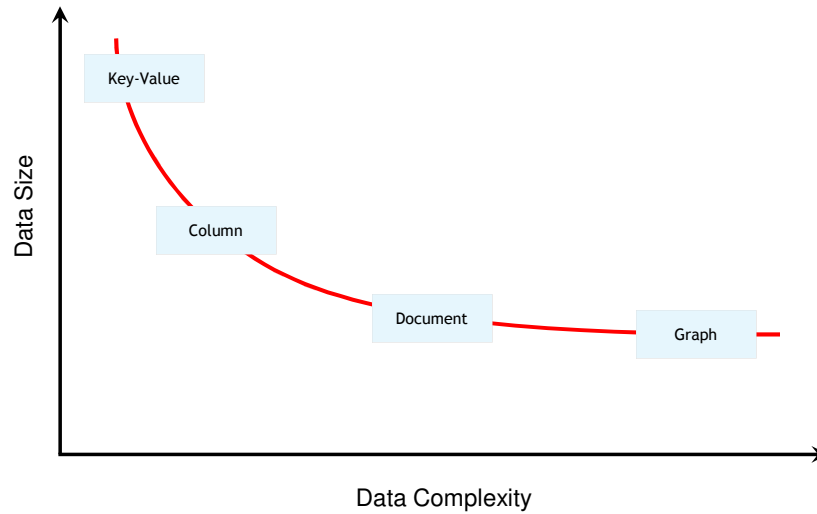
- Derived from Euler & CS Graph Theory
- Data model:
  - Nodes,
  - Relations
  - K-V on both
- Examples: AllegroGraph, Sones, Neo4j



© 2015 IBM Corporation

44

## NoSQL Data Models – Size ~ Complexity

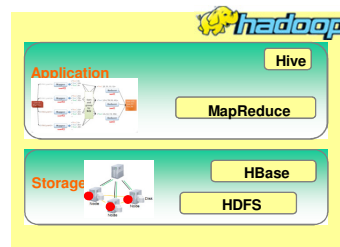


© 2015 IBM Corporation

45

## NoSQL and Hadoop 101

- NoSQL is a generic term for a very nongeneric landscape of data stores
  - Evolved over time to serve a very different world of requirements
  - Social, blogs, clickstream, Cloud apps, mobile apps, complex relationships...
- Hadoop is a *framework* for processing large and varied data sets with low cost at a high degree of fault tolerance. Components include (among others):
  - A file system (Hadoop File System - HDFS) -
  - A framework for processing data (Map-Reduce)
  - May include a NoSQL Database - HBase



© 2015 IBM Corporation

46

## The NoSQL Revolution

- Different requirements require different tools

- Document stores
- Key/value stores
- BigTable implementations (columnar)
- Graph databases

- Values (there are exceptions)

- Huge data volumes - easy scale-out
- Developers code integrity if it's needed
- Relaxed (eventual) consistency
- Semi-structured data
- Schema on read



Cloudant  
an IBM Company



Cassandra



© 2015 IBM Corporation

47

## Categories of NoSQL

Type	Examples
Document store	
Column store	
Key-value store	
Graph store	

Source: Akmal Chaudhri's NoSQL presentation.



© 2015 IBM Corporation

48



## Database Landscape Overview

	SQL	noSQL database	Hadoop
Description	<ul style="list-style-type: none"> <li>Relational SQL (RDBMS)</li> <li>Operational and Analytic</li> <li>E.g. DB2, Oracle, Microsoft, Teradata, etc.</li> </ul>	<ul style="list-style-type: none"> <li>noSQL database</li> <li>Mainly operational</li> <li>E.g. Cloudant, MongoDB, Redis, Riak, Aerospike, Amazon Dynamo DB, etc.</li> </ul>	<ul style="list-style-type: none"> <li>SQL on Hadoop (mainly analytic)</li> <li>HBase (evolving OLTP, ACID)</li> <li>E.g. BigInsights, Cloudera, Hortonworks, MapR, Pivotal</li> <li>HP Labs Trafodion</li> </ul>
Typical Infrastructure	<ul style="list-style-type: none"> <li>Proprietary database storage</li> <li>Unix, Linux, Windows</li> <li>SMP, MPP, SAN, Integrated Systems, Appliances</li> </ul>	<ul style="list-style-type: none"> <li>Proprietary database storage</li> <li>Linux</li> <li>Commodity clusters</li> <li>Local attach disks, NAS</li> <li>Cloud</li> <li>Mobile</li> </ul>	<ul style="list-style-type: none"> <li>HDFS files</li> <li>Linux</li> <li>Commodity clusters</li> <li>Local attach disks</li> </ul>
Primary Driver	<ul style="list-style-type: none"> <li>Traditional I/T</li> <li>ACID</li> </ul>	<ul style="list-style-type: none"> <li>Developer</li> <li>Agility, scalability, workload, cost</li> </ul>	<ul style="list-style-type: none"> <li>Lower Cost</li> <li>All types of data</li> </ul>



© 2015 IBM Corporation

49

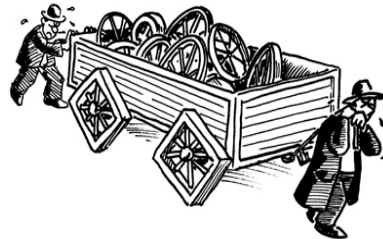
## The Right Tool for the Job

Big Data



≠

Hadoop



*"There's a belief that if you want big data, you need to go out and buy Hadoop and then you're pretty much set. People shouldn't get ideas about turning off their relational systems and replacing them with Hadoop..."*

*As we start thinking about big data from the perspective of business needs, we're realizing that Hadoop isn't always the best tool for everything we need to do, and that using the wrong tool can sometimes be painful."*



Ken Rudin  
Head of Analytics at Facebook



© 2015 IBM Corporation

50

# THANK YOU

“The art of progress is to preserve order amid change  
and to preserve change amid order”

Alfred North Whitehead