# Essential Tools For
# Your Big Data Arsenal
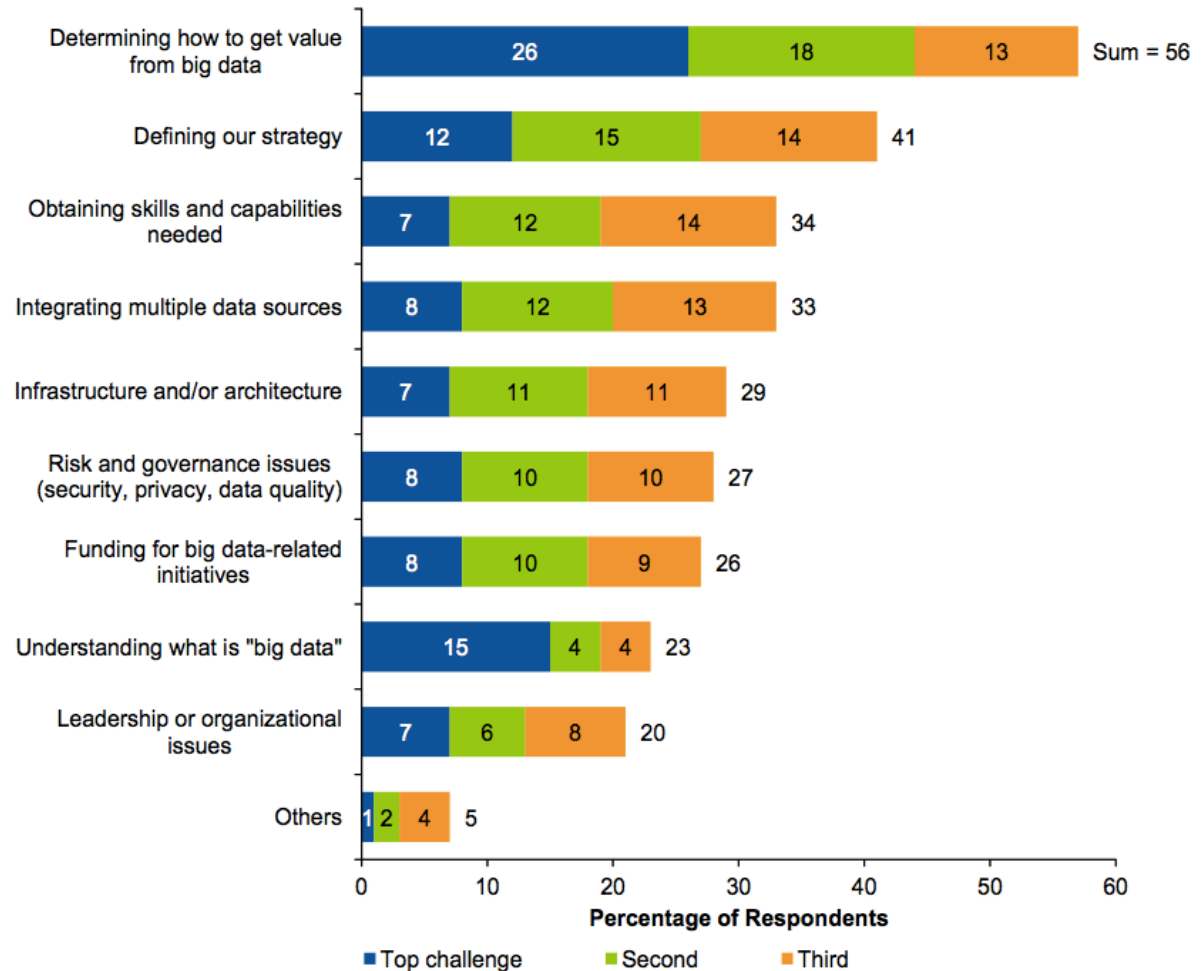
Matt Asay (@mjasay)
VP, Business Development & Strategy, MongoDB

mongoDB

# The Big Data Unknown

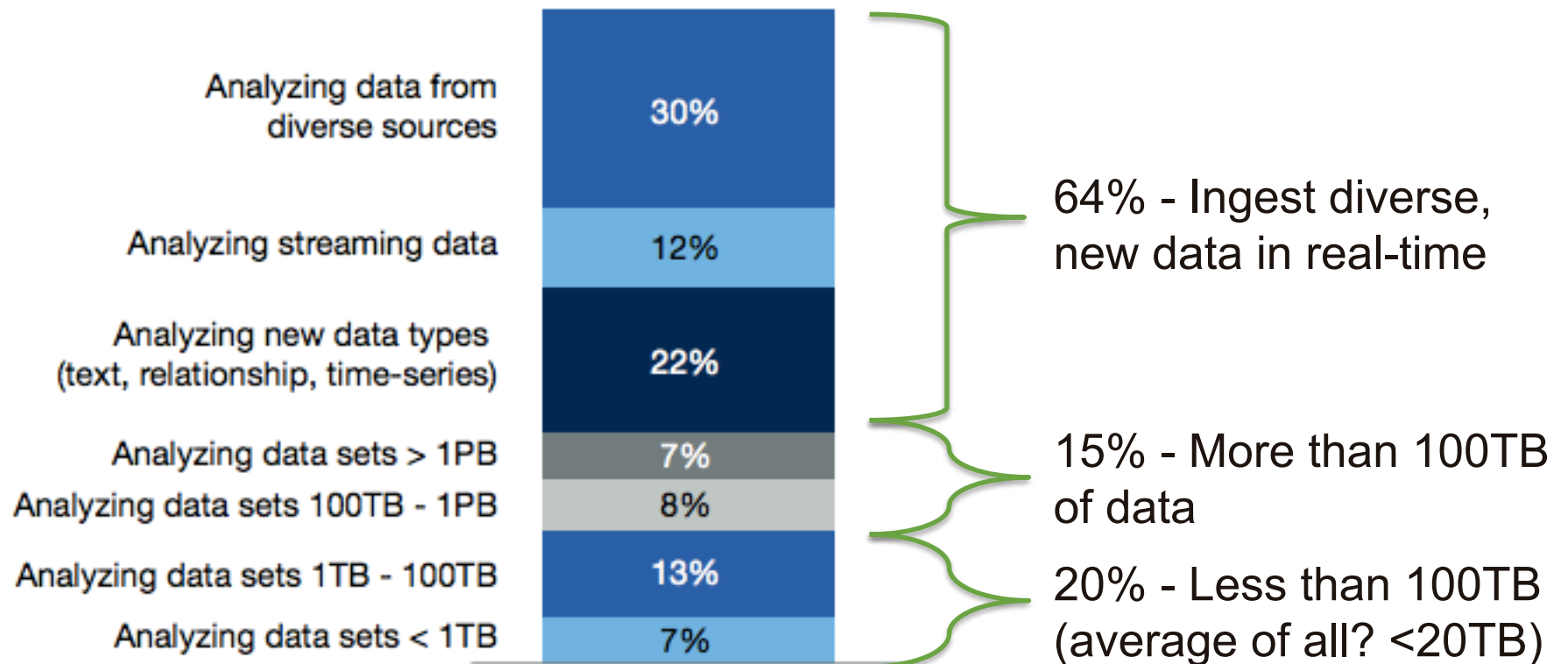# Top Big Data Challenges?

Translation?
Most struggle
to know what
Big Data is,
how to manage
it and who can
manage it



| Challenge | Top challenge | Second | Third | Sum |
|---|---|---|---|---|
| Determining how to get value from big data | 26 | 18 | 13 | Sum = 56 |
| Defining our strategy | 12 | 15 | 14 | 41 |
| Obtaining skills and capabilities needed | 7 | 12 | 14 | 34 |
| Integrating multiple data sources | 8 | 12 | 13 | 33 |
| Infrastructure and/or architecture | 7 | 11 | 11 | 29 |
| Risk and governance issues (security, privacy, data quality) | 8 | 10 | 10 | 27 |
| Funding for big data-related initiatives | 8 | 10 | 9 | 26 |
| Understanding what is "big data" | 15 | 4 | 4 | 23 |
| Leadership or organizational issues | 7 | 6 | 8 | 20 |
| Others | 1 | 2 | 4 | 5 |

Percentage of Respondents

■ Top challenge  ■ Second  ■ Third

N = 687 (excludes "don't know" responses)

*Source: Gartner*

3

mongoDB

# Understanding Big Data – It's Not Very "Big"

Analyzing data from diverse sources — 30%

Analyzing streaming data — 12%

Analyzing new data types (text, relationship, time-series) — 22%

Analyzing data sets > 1PB — 7%

Analyzing data sets 100TB - 1PB — 8%

Analyzing data sets 1TB - 100TB — 13%

Analyzing data sets < 1TB — 7%

64% - Ingest diverse, new data in real-time

15% - More than 100TB of data

20% - Less than 100TB (average of all? <20TB)

*from Big Data Executive Summary – 50+ top executives from Government and F500 firms*

mongoDB

# Innovation As Iteration

"I have not failed. I've just found 10,000 ways that won't work."
— Thomas A. Edison

# The New American Car.

This is the American Motors Gremlin. It is the kind of car this country has needed for a long, long time.

It is designed to give the American motorist a car that is easy to buy, easy to handle, easy to take care of, and, at the same time, fun to drive.

The Gremlin is the smallest production car made in America.

It is 161 inches long, just 2½ inches longer than the Volkswagen.

Yet its turning circle, at 32 feet, 8 inches, is about 3 feet less than VW's.

Which makes the Gremlin about the easiest car in the world to park and handle.

The Gremlin gets the best gas mileage of any car made in America. It goes about 500 miles without stopping for gas.

This is great gas mileage, when you consider that the Gremlin has a bigger standard engine than any car near its size and price. 128 hp to VW's 57.

This engine gets from 0 to 60 in 15.3 seconds, the pickup you need on expressways.

And nobody's going to push you around in a Gremlin. It is 10 inches wider, 7 inches lower and 765 pounds heavier than a VW.

Which gives you about the smoothest, most stable ride possible in a car this size.

The Gremlin is remarkably easy to service and maintain.

Its normal oil change interval is 6 months or 6,000 miles; lubrication is normally needed only every 24,000 miles.

There are two basic Gremlin models.

A two-passenger, with storage area in the rear.

A four-passenger with fold-down rear seats for extra storage and flip-up rear window for easy access.

Both models cost about what you'd pay for an imported economy car.

The four-passenger lists for $1,959. The lowest list price of any car made in America.

Except for the two-passenger Gremlin. It lists for $1,879.

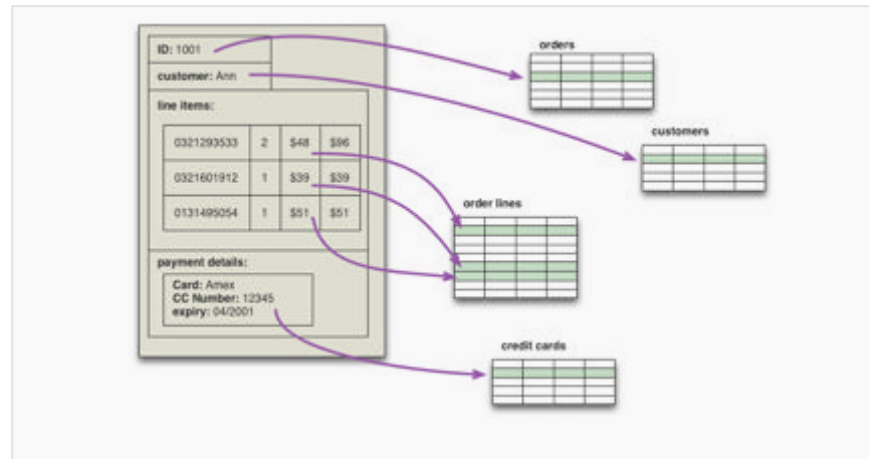Which is quite a bargain, when you consider what you get for your money.

The new American car.

**American Motors Gremlin**

$1,879  $1,959
2-Passenger  4-Passenger

# So Were Computers!

# Lots of Great Innovations Since 1970

# Including the Relational Database

# RDBMS Makes Development Hard



Code

XML Config

DB Schema

Application

Object Relational Mapping

Relational Database

mongoDB

# And Even Harder To Iterate



New Table

New Table

New Column

New Column

| Name | Pet | Phone | Email |
|------|-----|-------|-------|
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |
|      |     |       |       |

3 months later…

mongoDB

# From Complexity to Simplicity

**RDBMS**



**MongoDB**

```
{

    _id : ObjectId("4c4ba5e5e8aabf3"),
    employee_name: "Dunham, Justin",
    department : "Marketing",
    title : "Product Manager, Web",
    report_up: "Neray, Graham",
    pay_band: "C",
    benefits : [
            {  type :   "Health",
               plan : "PPO Plus" },
            {  type :    "Dental",
               plan : "Standard" }
                ]

}
```

# So…Use Open Source
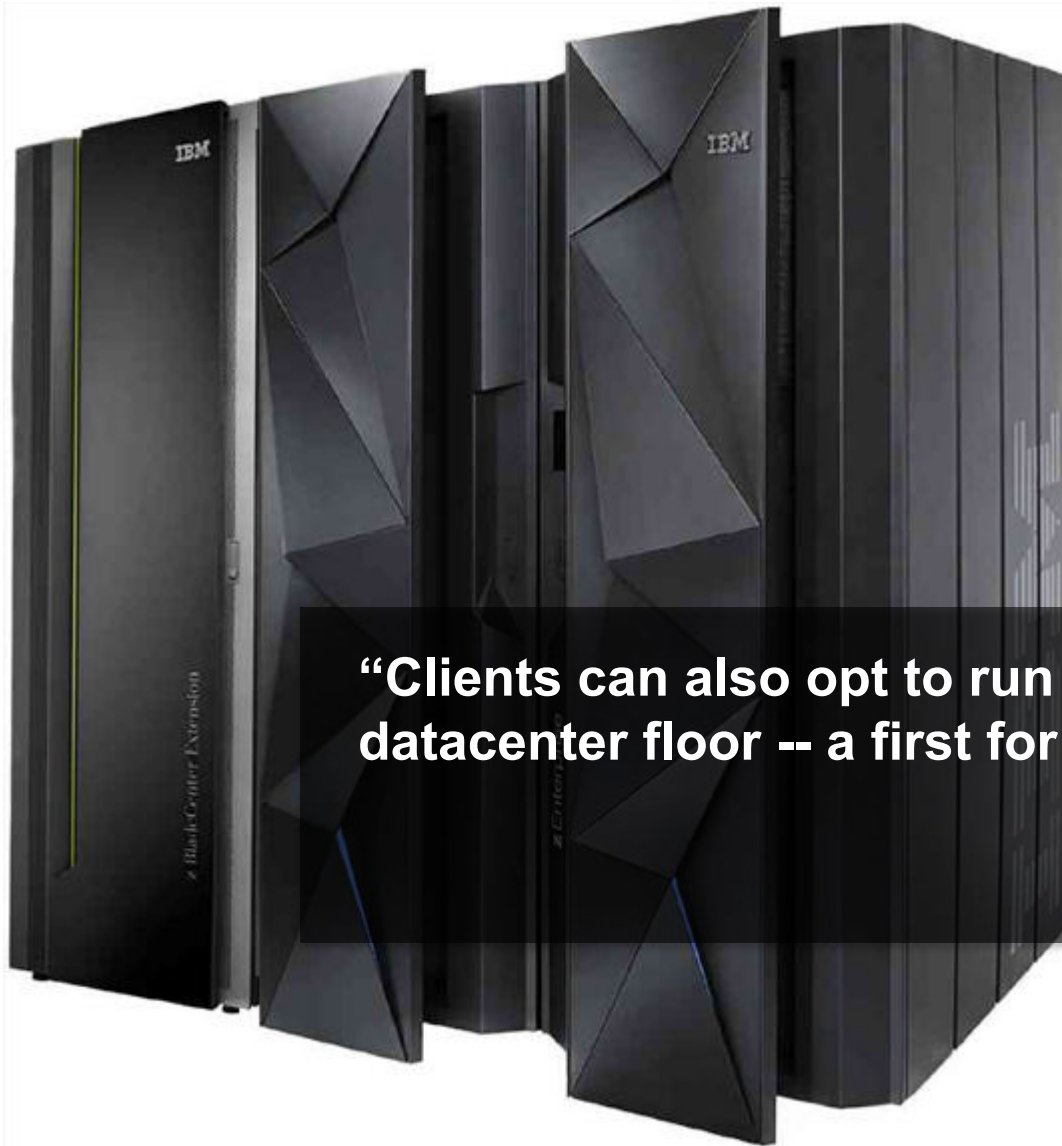
Modernizing IT
## ERIC KNORR

SEPTEMBER 18, 2012

## Open source in 2012: Bigger and better than ever

This year's Best of Open Source Software awards includes a whopping 125 products in 7 categories. The real story is the technology leadership so many of these products display

mongoDB

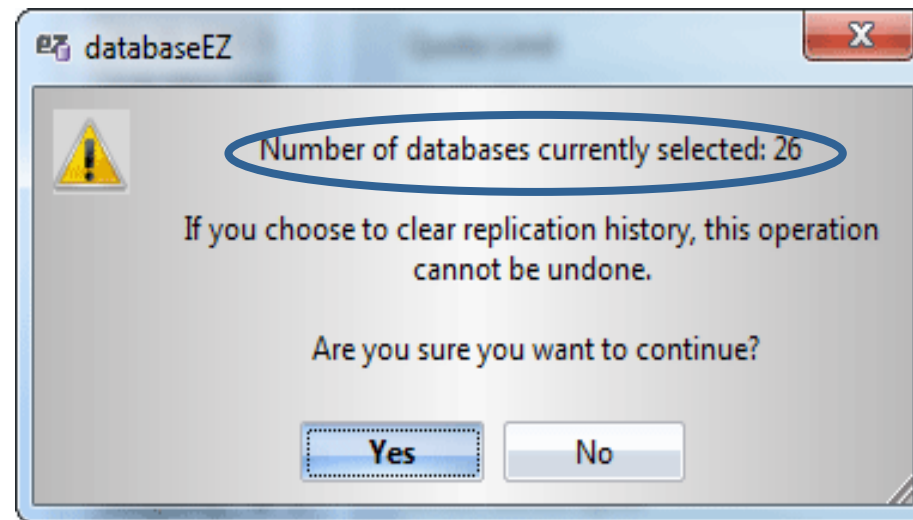# Big Data != Big Upfront Payment

# RDBMS Is Expensive To Scale

"Clients can also opt to run zEC12 without a raised datacenter floor -- a first for high-end IBM mainframes."

*IBM Press Release 28 Aug, 2012*

# Spoiled for choice

### DB-Engines.com Database Ranking

| | | | | |
|---|---|---|---|---|
| 1 | Oracle | Relational DBMS | 1583.84 | 54.23 |
| 2 | MySQL | Relational DBMS | 1331.34 | 25.58 |
| 3 | Microsoft SQL Server | Relational DBMS | 1207 | -106.78 |
| 4 | PostgreSQL | Relational DBMS | 177.01 | -5.22 |
| 5 | DB2 | Relational DBMS | 175.83 | 3.58 |
| 6 | MongoDB | NoSQL Document Store | 149.48 | -2.71 |
| 7 | Microsoft Access | Relational DBMS | 142.49 | -4.21 |
| 8 | SQLite | Relational DBMS | 77.88 | -4.9 |
| 9 | Sybase | Relational DBMS | 73.66 | -1.68 |
| 10 | Teradata | Relational DBMS | 54.41 | 3.32 |

databaseEZ

Number of databases currently selected: 26

If you choose to clear replication history, this operation cannot be undone.

Are you sure you want to continue?

Yes     No

mongoDB

# Remember the Long Tail?

# It Didn't Work Out So Well

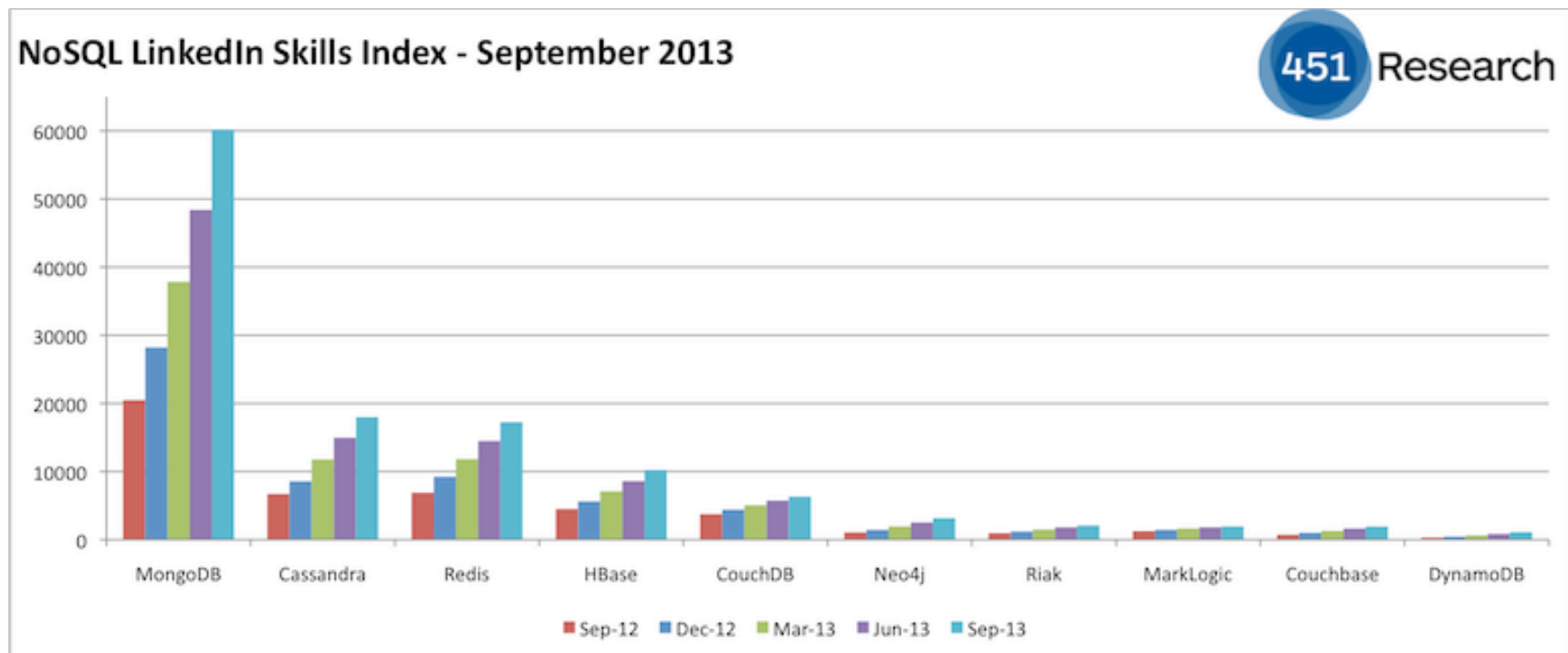# Use Popular, Well-Known Technologies



Source: Silicon Angle, 2012

# Ask the Right Questions…

"Organizations already have people who know their own data better than mystical data scientists….Learning Hadoop [or MongoDB] is easier than learning the company's business."
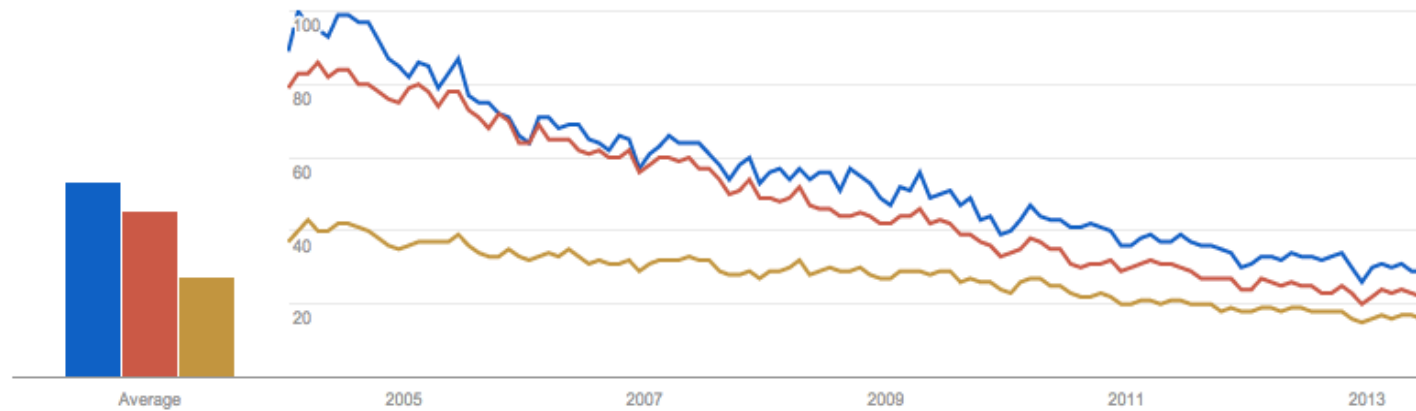
(Gartner, 2012)

# Leverage Existing Skills



NoSQL LinkedIn Skills Index - September 2013

451 Research

Legend: Sep-12, Dec-12, Mar-13, Jun-13, Sep-13

Categories: MongoDB, Cassandra, Redis, HBase, CouchDB, Neo4j, Riak, MarkLogic, Couchbase, DynamoDB

# Search as a Sign?

# When To Use Hadoop, NoSQL

# Enterprise Big Data Stack

**Management & Monitoring**

**Security & Auditing**

## Applications
CRM, ERP, Collaboration, Mobile, BI

## Data Management

### Online Data
mongoDB          RDBMS

### Offline Data
Hadoop          EDW

## Infrastructure
OS & Virtualization, Compute, Storage, Network

mongoDB

# Consideration – Online vs. Offline

**Online**          *vs.*          **Offline**

- Real-time
- Low-latency
- High availability

- Long-running
- High-Latency
- Availability is lower priority

mongoDB

# Consideration – Online vs. Offline

**Online**　　　*vs.*　　　**Offline**



mongoDB　　　　hadoop

mongoDB

# Hadoop Is Good for…

Risk Modeling

Churn Analysis

Recommendation Engine

Ad Targeting

Transaction Analysis

Trade Surveillance

Network Failure Prediction

Search Quality

Data Lake

mongoDB

# MongoDB/NoSQL Is Good for...

360° View of the Customer

Mobile & Social Apps

Fraud Detection

User Data Management
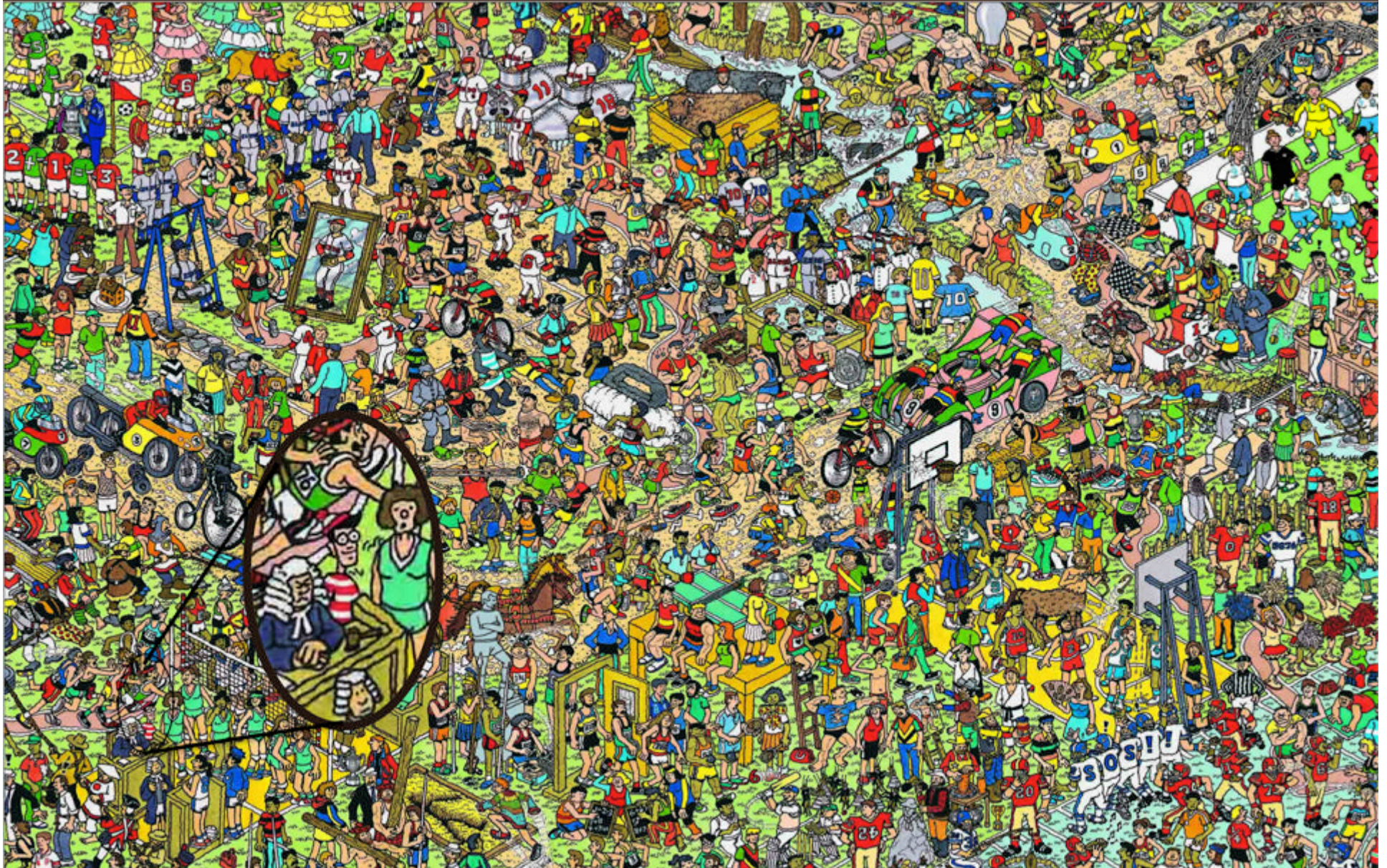
Content Management & Delivery

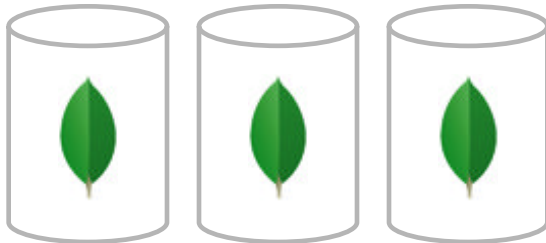Reference Data

Product Catalogs

Machine to Machine Apps
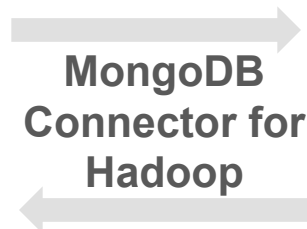
Data Hub

mongoDB

# How To Use The Two Together?

# Finding Waldo

# Customer example: Online Travel

**Travel**



**MongoDB Connector for Hadoop**

**Algorithms**



- **Flights, hotels and cars**
- **Real-time offers**
- **User profiles, reviews**
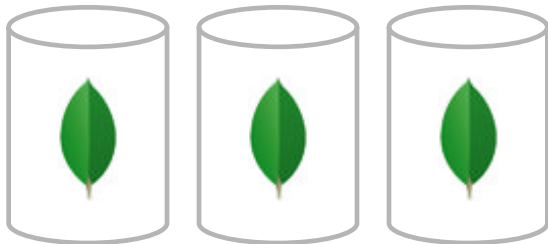- **User metadata (previous purchases, clicks, views)**

- **User segmentation**
- **Offer recommendation engine**
- **Ad serving engine**
- **Bundling engine**

# Predictive Analytics

### Government

### Algorithms

**MongoDB**

**+ Hadoop**

*MapReduce*

- **Predictive analytics system for crime, health issues**
- **Diverse, unstructured (incl. geospatial) data from 30+ agencies**
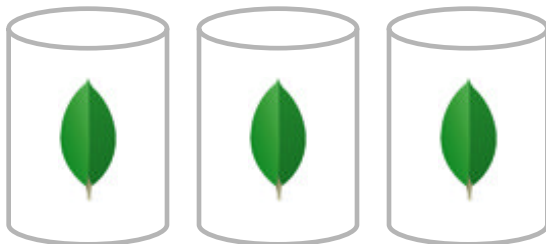- **Correlate data in real-time**

- **Long-form trend analysis**
- **MongoDB data dumped into Hadoop, analyzed, re-inserted into MongoDB for better real-time response**

# Data Hub

**Insurance**

**Churn Analysis**

MongoDB Connector for Hadoop

- **Insurance policies**
- **Demographic data**
- **Customer web data**
- **Call center data**
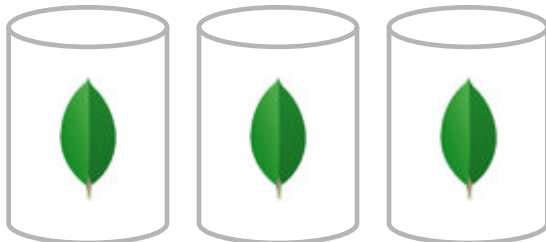- **Real-time churn detection**

- **Customer action analysis**
- **Churn prediction algorithms**

mongoDB

# Machine Learning



## Ad-Serving

- Catalogs and products
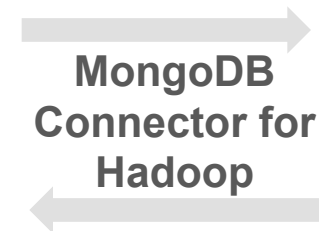- User profiles
- Clicks
- Views
- Transactions

## Algorithms

- User segmentation
- Recommendation engine
- Prediction engine

**MongoDB Connector for Hadoop**

# MongoDB + Hadoop Connector

- Makes MongoDB a Hadoop-enabled file system

- Read and write to live data, in-place

- Copy data between Hadoop and MongoDB

- Full support for data processing

  - Hive

  - MapReduce

  - Pig

  - Streaming

  - EMR

**MongoDB Connector for Hadoop**