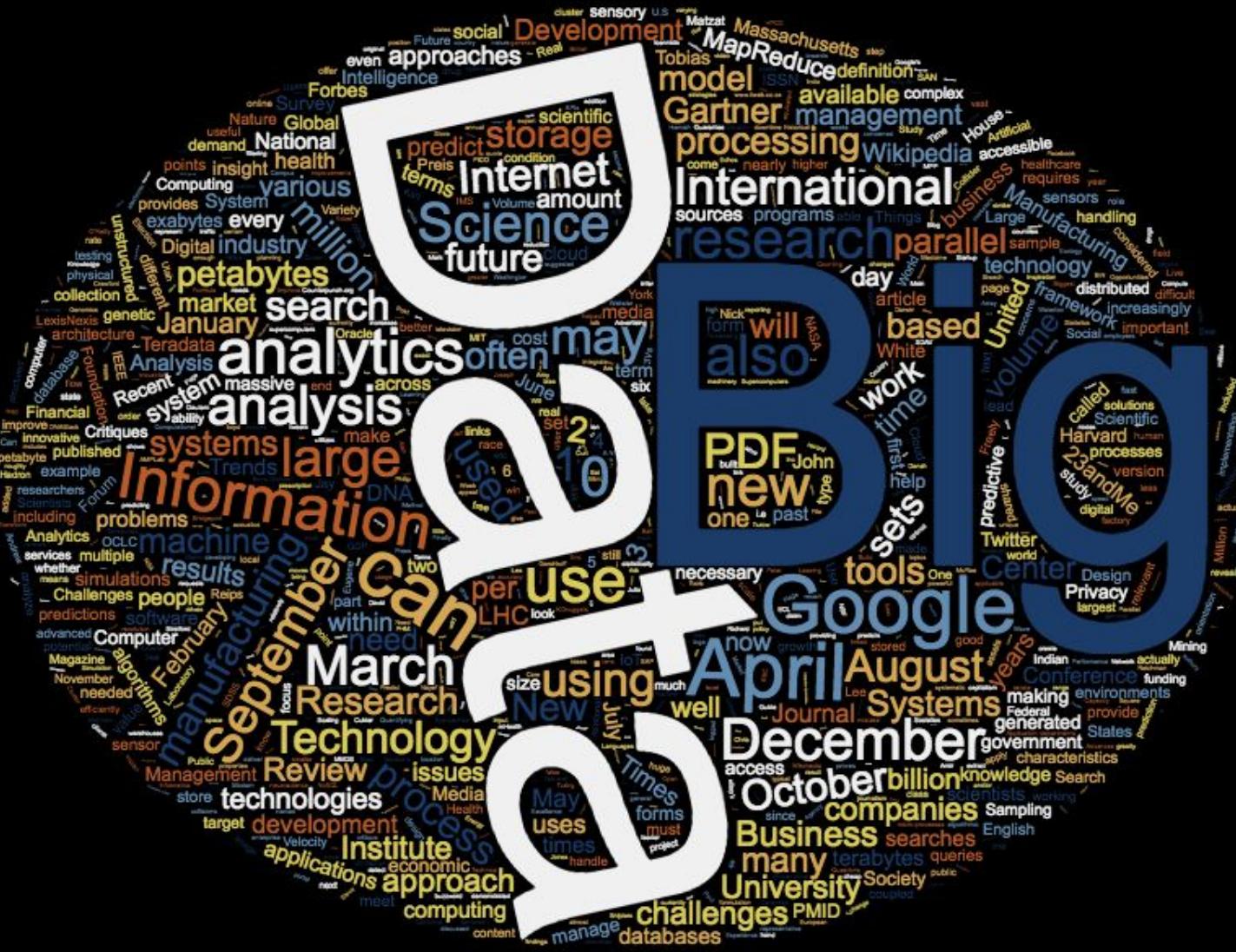


# Week 1: From Relational Data to Big Data

Big Data Management – Fall 2017  
Björn Þór Jónsson



# Relational Databases

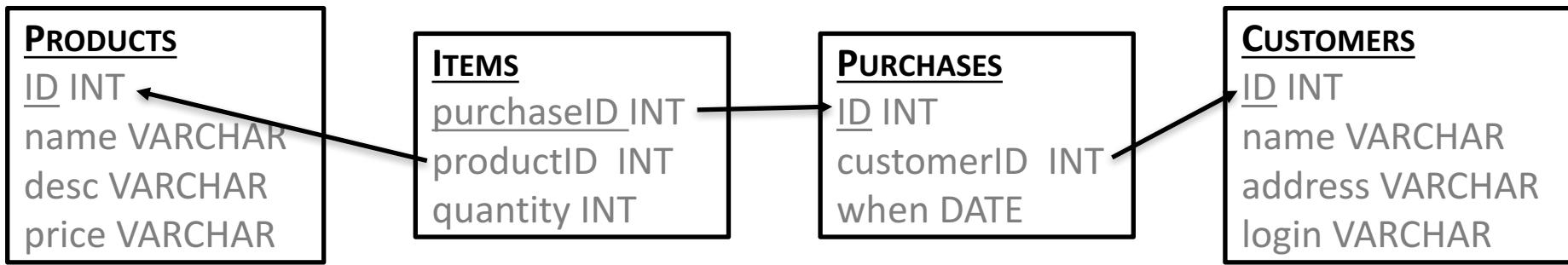
- What are the main techniques and concepts?

# Relational Database Design

- How are databases traditionally designed?

*For an online store we need to store information about customers, such as address, e-mail, and login information. We need to store information about our products, such as name, description, and price. Then we need to track all customer purchases, when they took place, which products were bought and how many copies of each.*

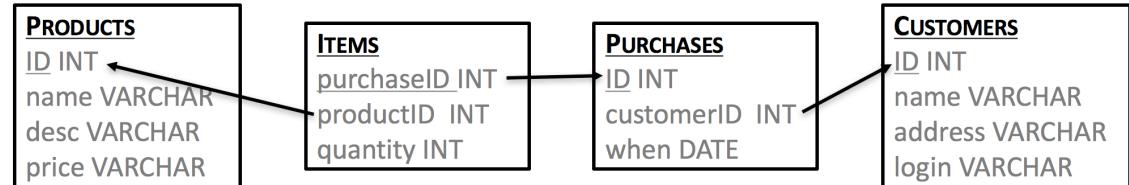
# Any problems with this design?



# Problem 1: History

- What if...

- Price changes?
  - Customer moves?



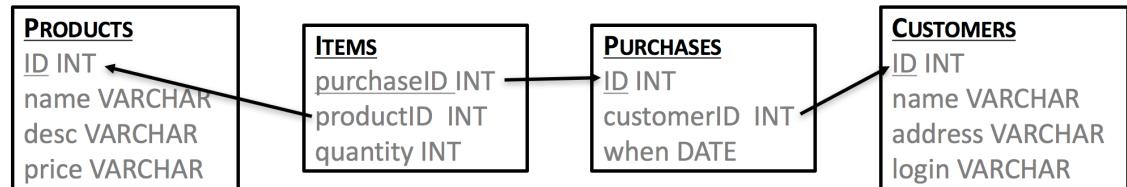
- Current state vs. Historical information
  - Traditional solution: Data warehouses
    - Snapshots of databases to keep track of history
  - Recent solution: Lambda architecture
    - Store everything and extract current state
  - Either way: Much more data!

# Problem 2: Coarseness

- What if...
  - Customer has problems with web-site?
  - Customer looks at items but does not purchase?
  - Do we want to know?
- Want much more detailed overview
  - Solution: Storing and analyzing logs
  - Logs are extensive: MUCH MORE data!

# Problem 3: Simple Structure

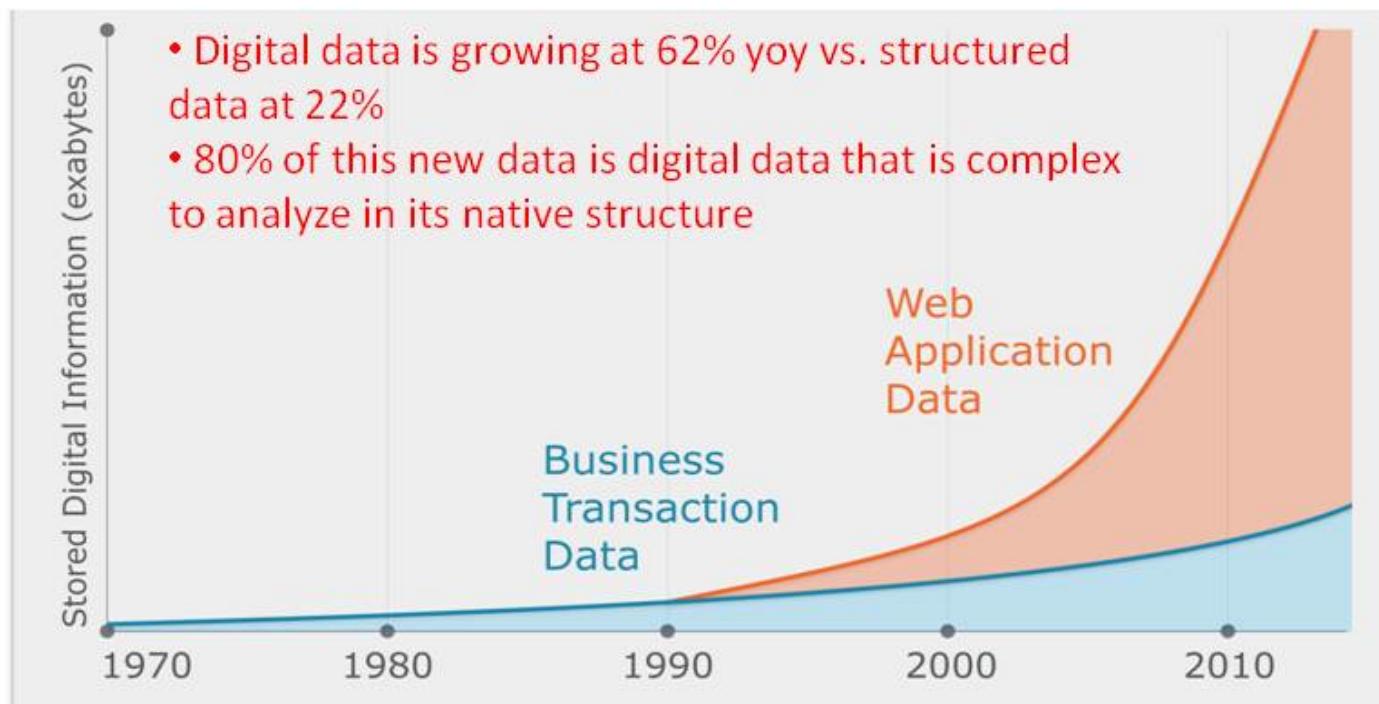
- What if...



- We want to allow (and keep track of) comments?
- Media (images/videos) is involved?
- Want to store any data type
  - Relational model forces you to **think inside the box!**
  - Unstructured and semi-structured data
  - Media is large: **MUCH MORE** data!

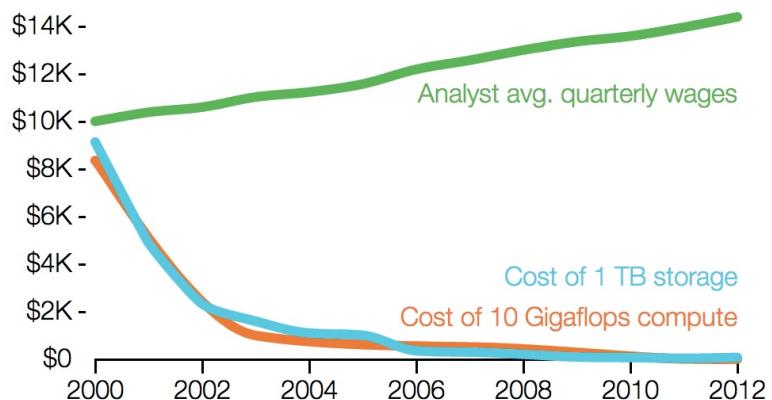
# Beyond Relational Data!

- Keep track of **all** the history
- Keep track of **all** interactions, also low-level
- Keep track of **all** data: media, text data, logs, ...



# Handling Data Growth?

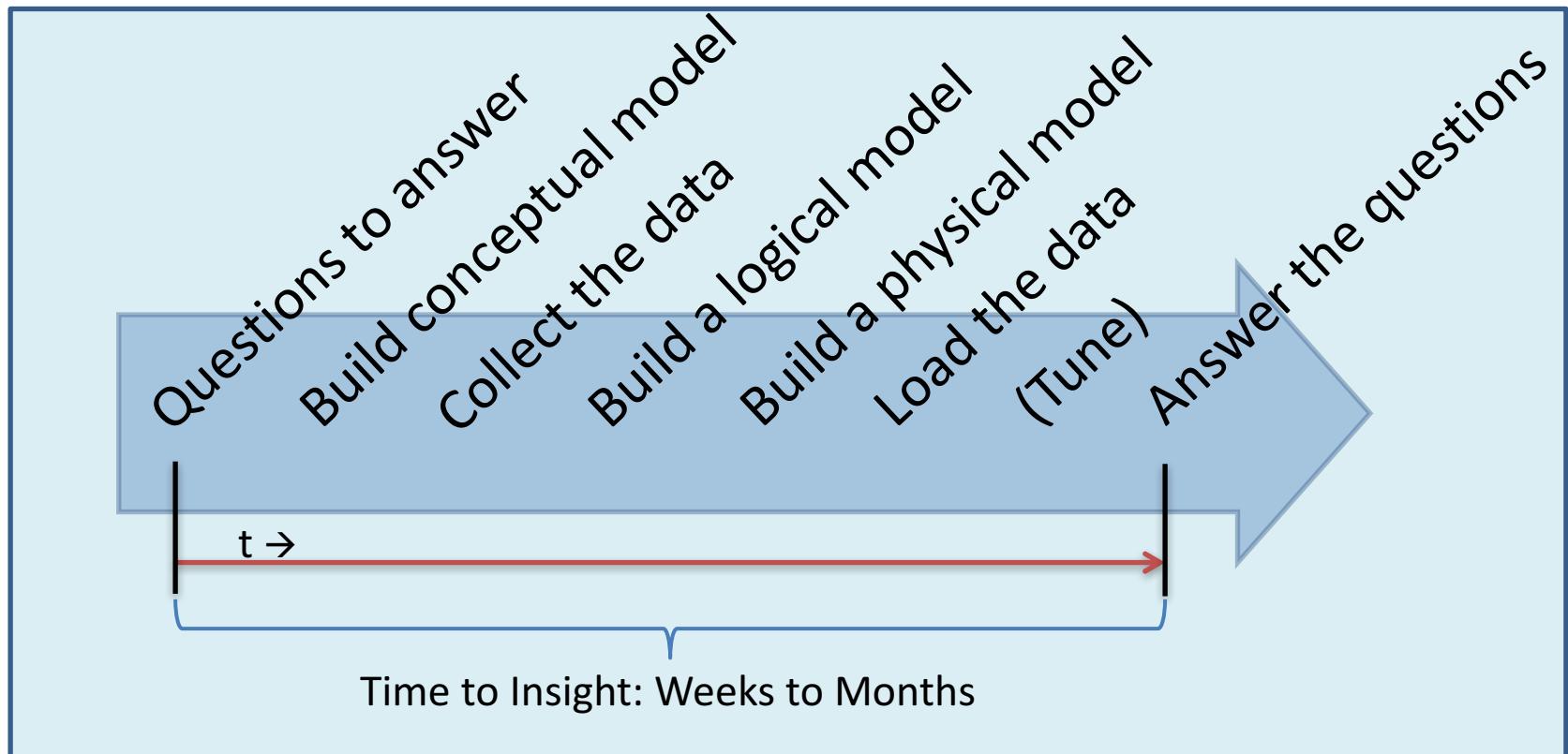
- Storage is cheap – people are not!



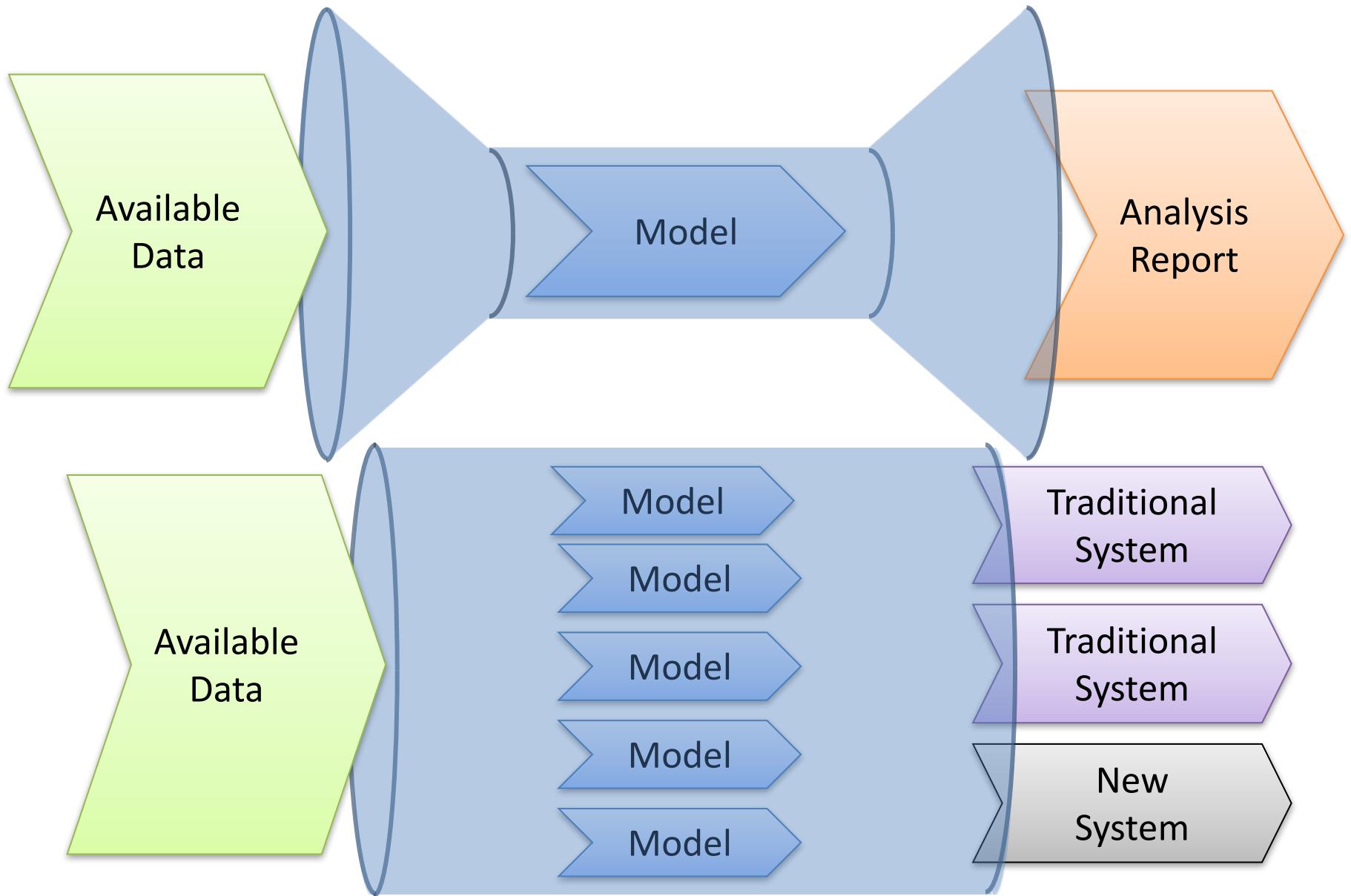
- Need better technology to handle the data

→ *Business Push*

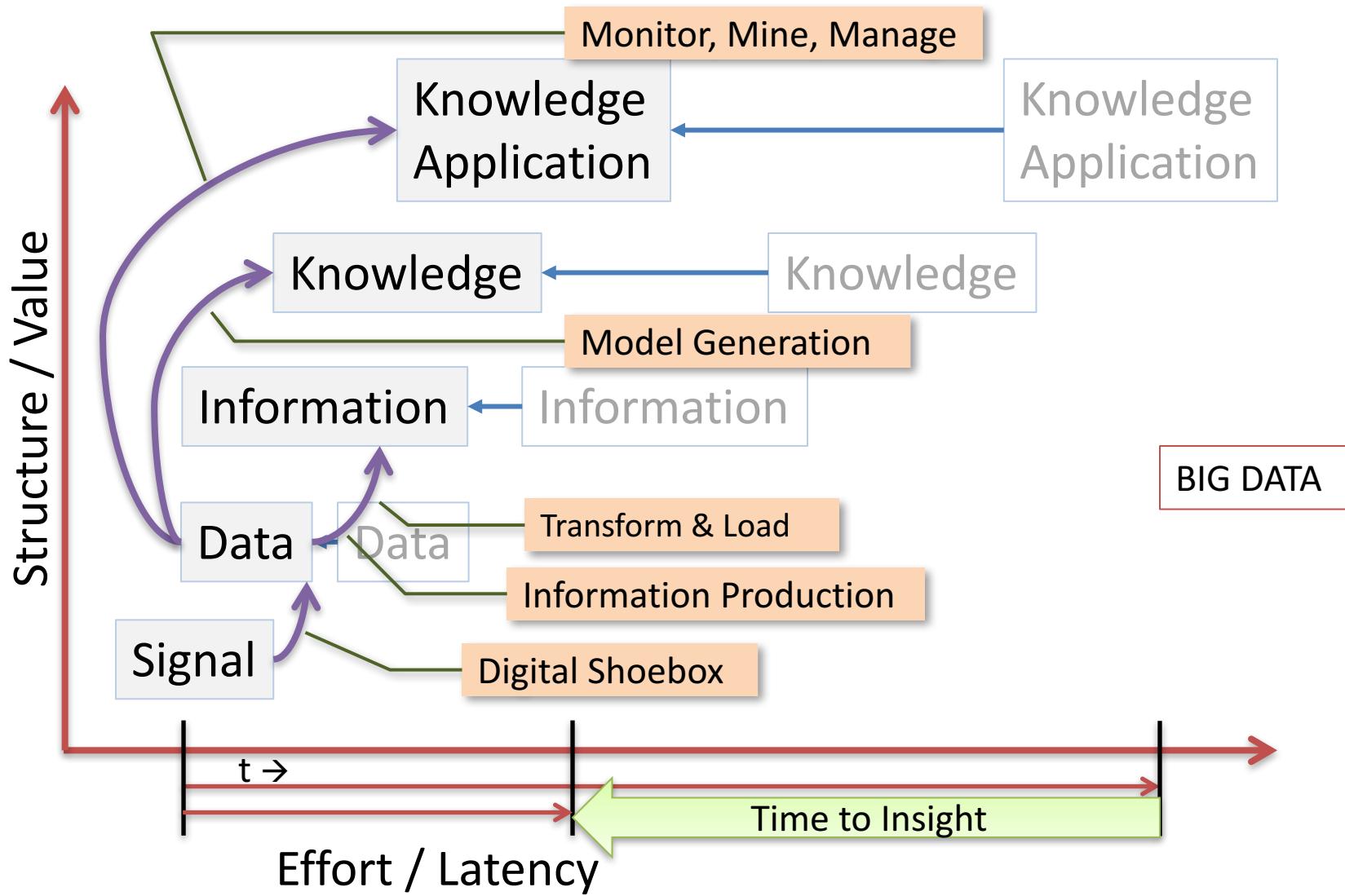
# Traditional “Sense Making”



OLD SCHOOL



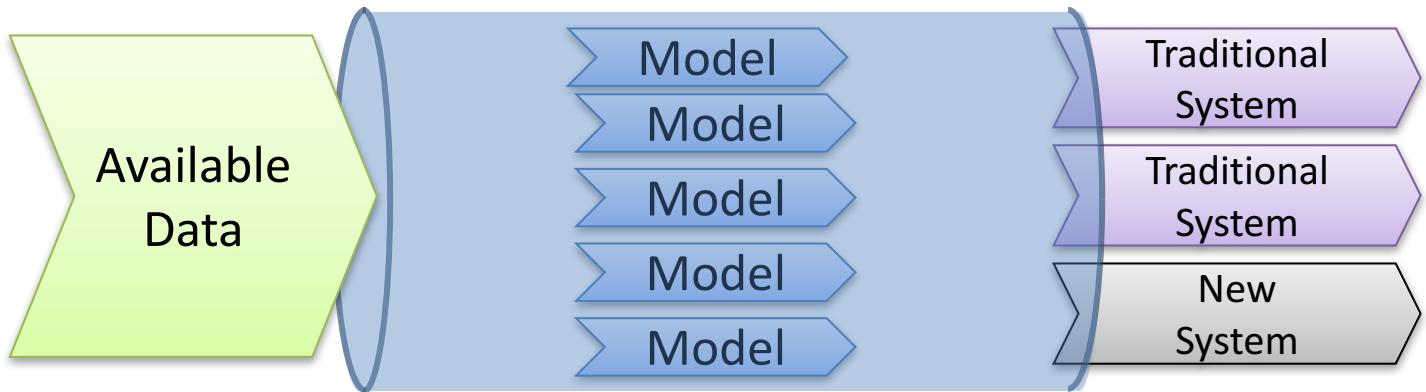
@ Dennis Shasha and Philippe Bonnet, 2013



@ Dennis Shasha and Philippe Bonnet, 2013

# Coping with a Faster World?

- Decisions must be made ever faster!

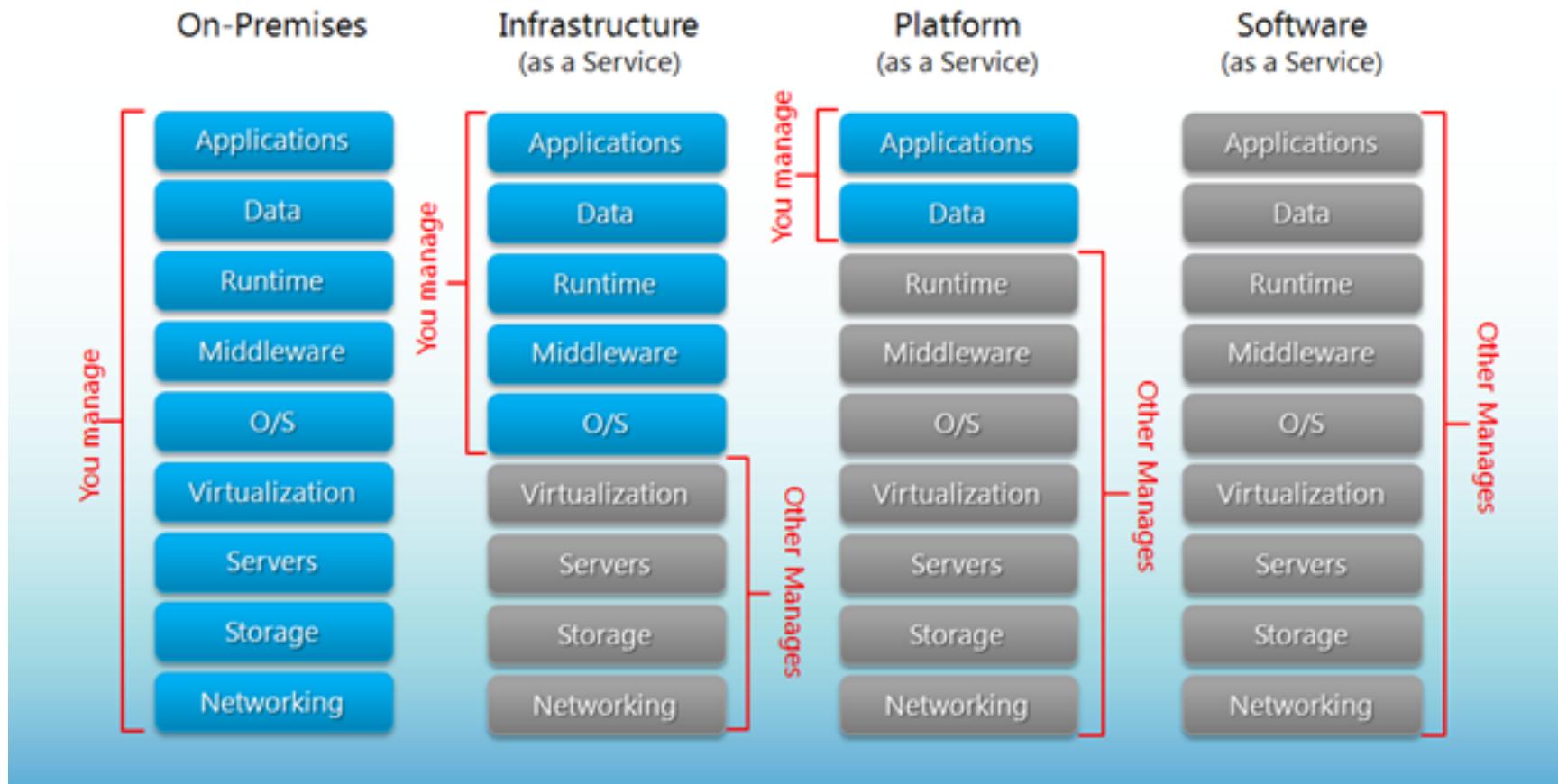


- Need better technology to use data **now!**

→ *Application Pull*

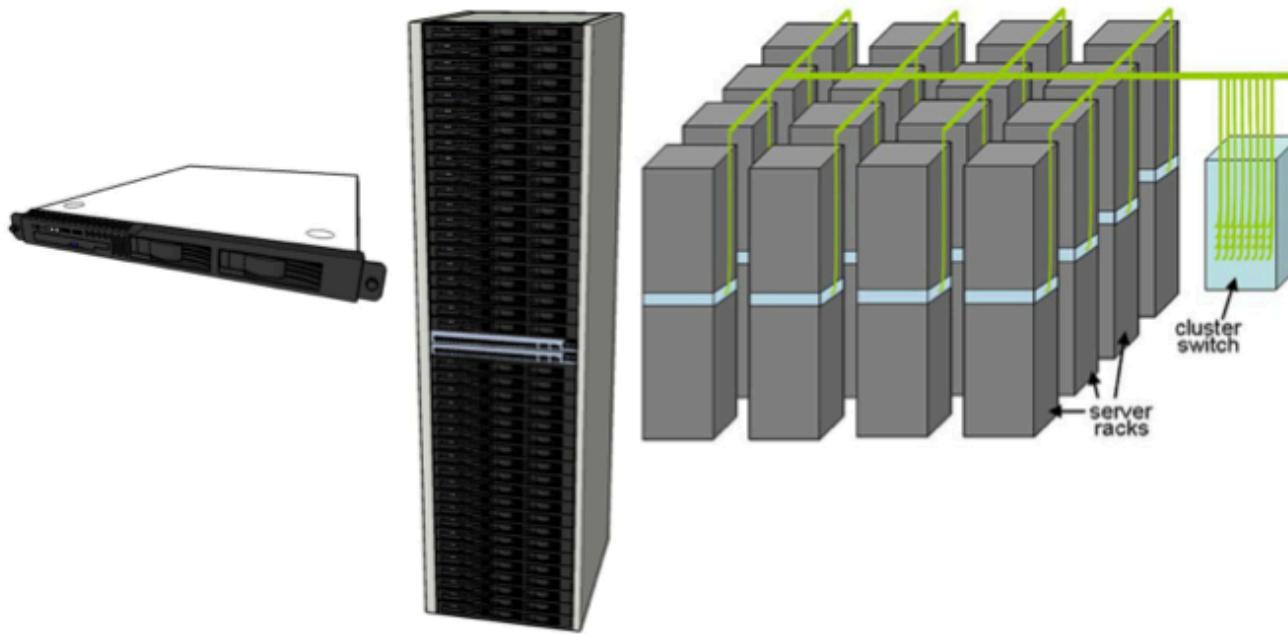
# Technology Trends

## Separation of Responsibilities



# Warehouse-Scale Computer

## THE DATACENTER AS A COMPUTER



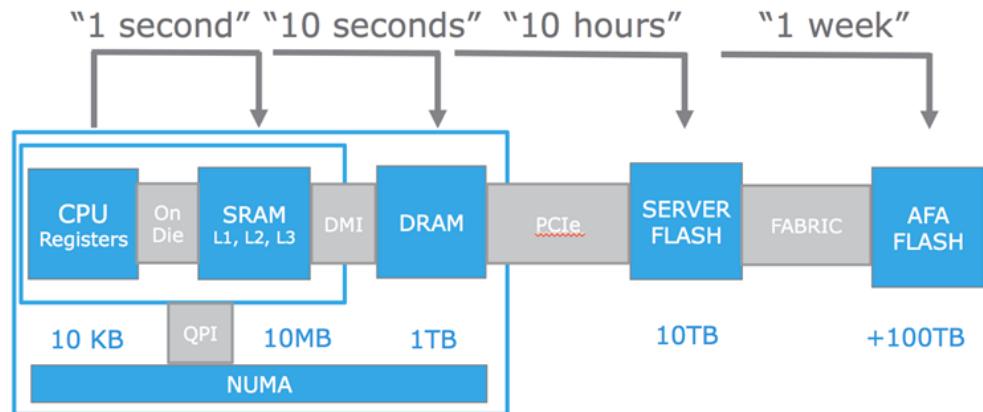
**FIGURE 1.1:** Typical elements in warehouse-scale systems: 1U server (left), 7' rack with Ethernet switch (middle), and diagram of a small cluster with a cluster-level Ethernet switch/router (right).

**LOOK UP:** [Werner Voegels on virtualization.](#)

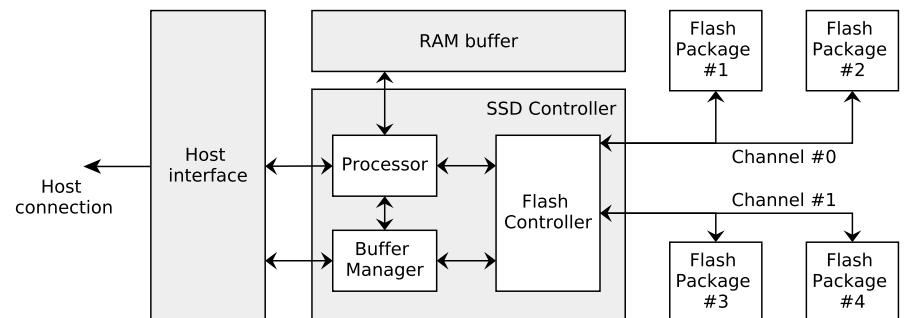
# Storage Trends

| Device                    | Read        | Write       |
|---------------------------|-------------|-------------|
| <b>Millisecond Scale</b>  |             |             |
| 10G Intercontinental RPC  | 100 ms      | 100 ms      |
| 10G Intracontinental RPC  | 20 ms       | 20 ms       |
| Hard Disk                 | 10 ms       | 10 ms       |
| 10G Interregional RPC     | 1 ms        | 1 ms        |
| <b>Microsecond Scale</b>  |             |             |
| 10G Intraregional RPC     | 300 $\mu$ s | 300 $\mu$ s |
| SATA NAND SSD             | 200 $\mu$ s | 50 $\mu$ s  |
| PCIe/NVMe NAND SSD        | 60 $\mu$ s  | 15 $\mu$ s  |
| 10Ge Inter-Datacenter RPC | 10 $\mu$ s  | 10 $\mu$ s  |
| 40Ge Inter-Datacenter RPC | 5 $\mu$ s   | 5 $\mu$ s   |
| PCM SSD                   | 5 $\mu$ s   | 5 $\mu$ s   |
| <b>Nanosecond Scale</b>   |             |             |
| 40 Gb Intra-Rack RPC      | 100 ns      | 100 ns      |
| DRAM                      | 10 ns       | 10 ns       |
| STT-RAM                   | <10 ns      | <10 ns      |

M.Wei et al. I/O speculation in the microsecond era. Usenix ATC'14.

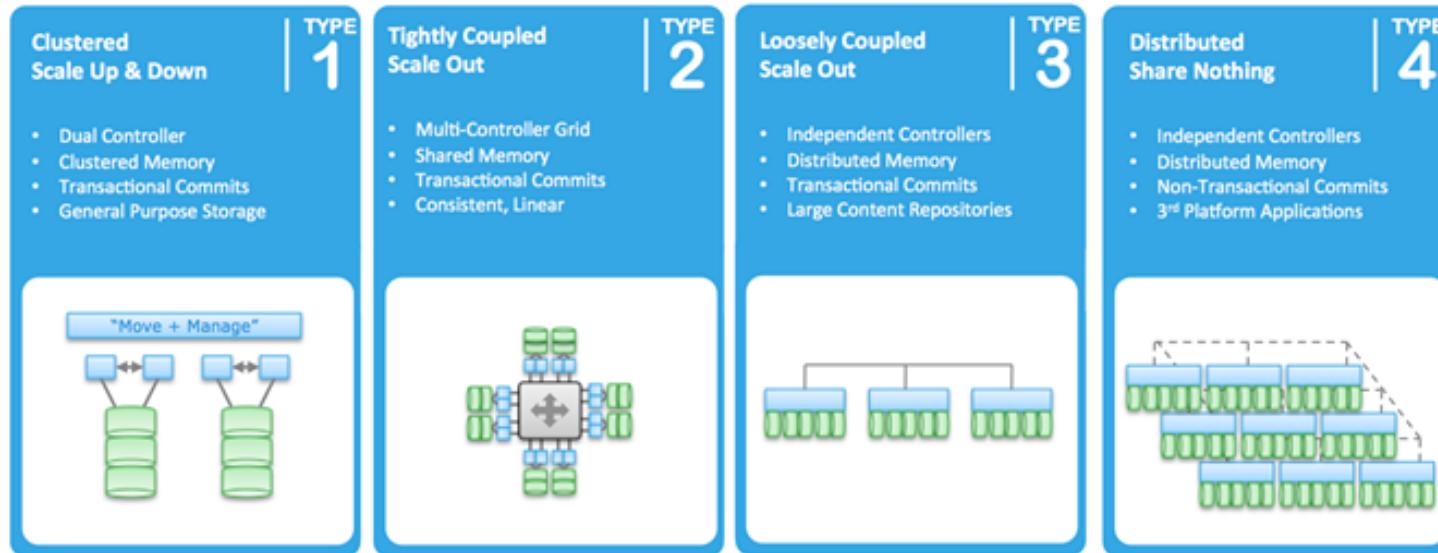


source: [Virtual Geek's take on storage tree of life](#)

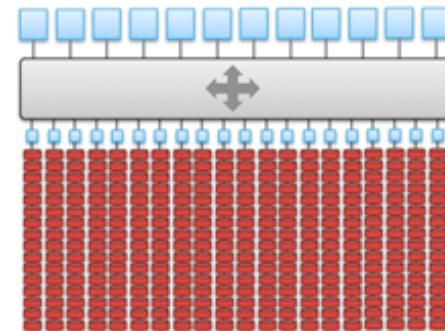


SSD Architecture

# Storage Architectures

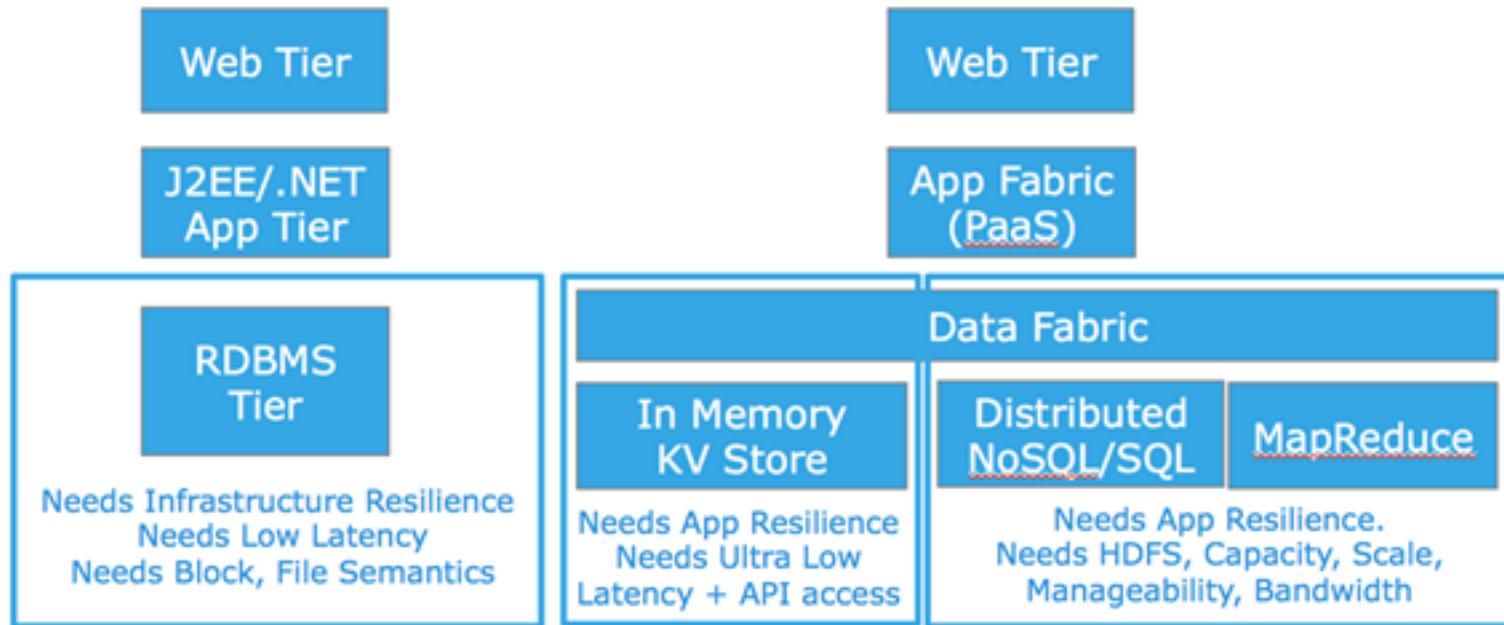


source: [Virtual Geek's take on storage tree of life – A MUST READ!!](#)



Storage    [light green square]    [dark red square]  
RAM        [light blue square]  
Interconnect [light blue square]

# Data-Intensive Applications: Server-side Architectures

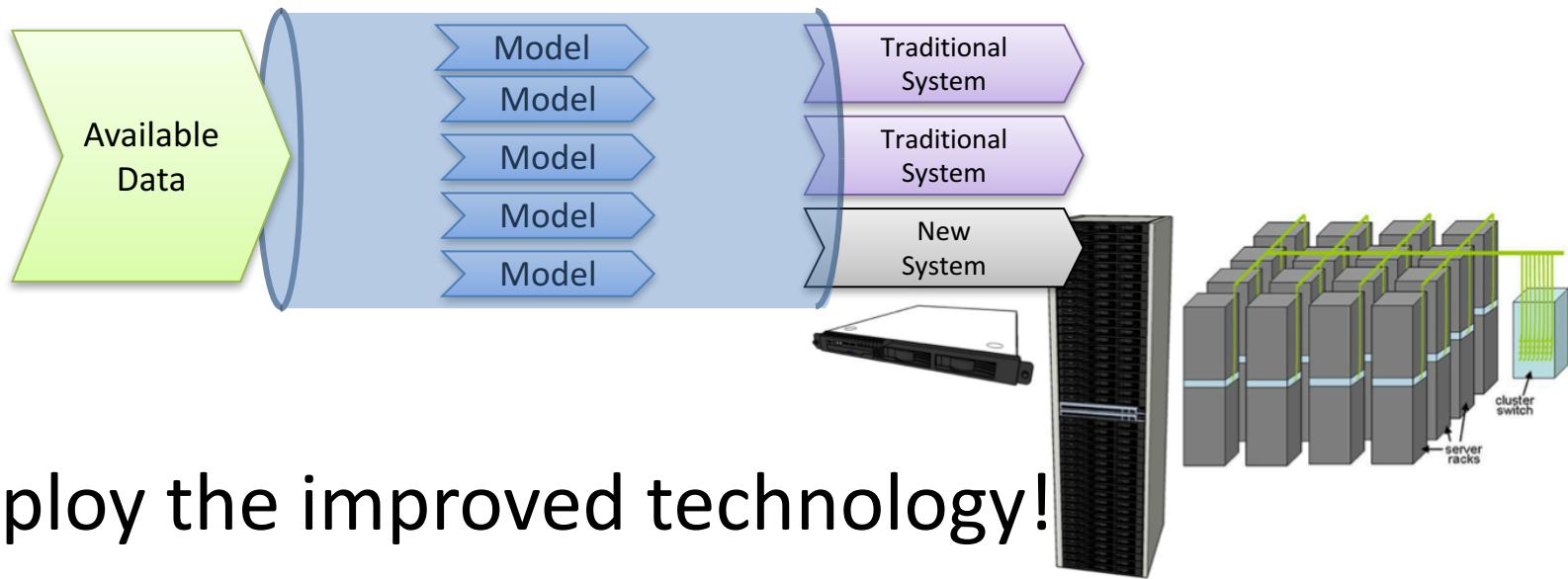


Look up [Fabric Computing](#) on Wikipedia.

source: [Virtual Geek's take on storage tree of life – A MUST READ!!](#)

# Utilizing a Faster World ☺

- Hardware architectures getting ever better!



- Employ the improved technology!

→ *Technology push!*

# Outline

- Trends Underlying Big Data
  - Business Push
  - Application Pull
  - Technology Push
- **What is Big Data?**
  - Relational vs. Big Data Management
  - Examples
  - Definitions (attempts?)
- Data models
- Course Outline

# Relational Database Management

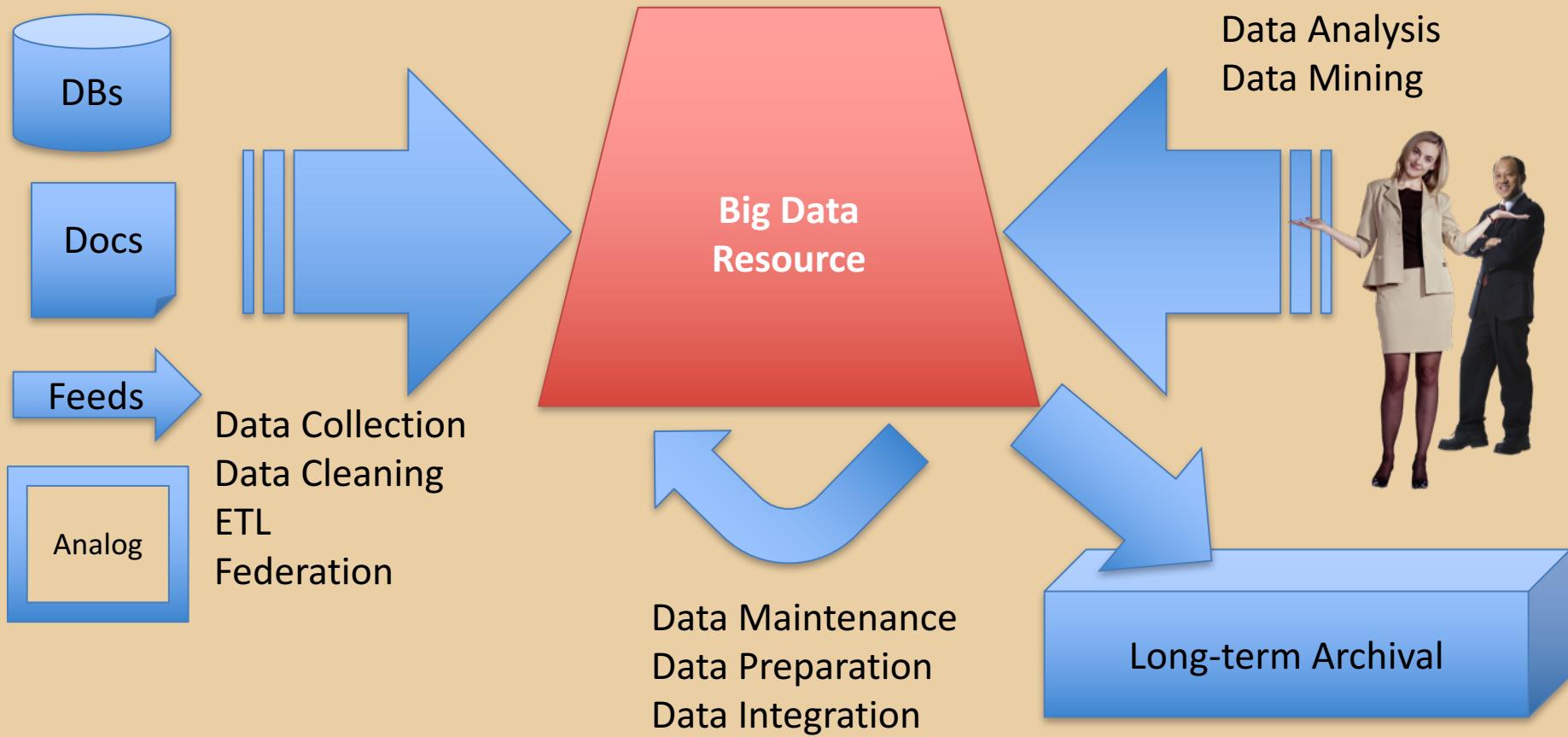
- **Narrow scope**
  - A database is created to serve a well defined purpose
- **Structured data**
  - Conceptual/Logical/Physical schema
  - Relational model dominates since 80s
  - Entity Relationship defines conceptual schema
- **Closed-world assumption**
  - Data as an instance of the schema
  - The data which is not part of an instance *does not exist*
    - *Any query on the database returns a value based on the current instance*
- **Data at rest**
  - Data is loaded and stored in the database, on disk.
- **Fully interactive architecture**
  - 3 tier architecture: Web server; App Server; Database Server.

# What is a Big Data Resource?

- A big data resource is a **collection of data which is made available for analysis**
  - Analysis of data is a *process of inspecting, cleaning, transforming, and modeling* data with the goal of discovering useful information, suggesting conclusions, and supporting decision making.
- For example:
  - The databases that underly the learnit blog or FM's Building Management System are not big data resources
  - The data made available by Eurostat on unemployment in Europe can be considered a big data resource.
    - More examples of this kind at [Google Public Data Explorer](#)

# What is Big Data Analysis?

- Need for insight based on data not currently available for analysis
  - How can we set goals for energy management at the IT University?
    - We do not know how electricity is consumed.
    - We do not know where electricity is consumed.
    - We do not know how to link electricity consumption and people's practices.
  - How long does it take you to answer your mails?



Big data is not a product, but a collection of **processes**

# Big Data Management

- **Wide scope**
  - Data is made available for yet-to-be-defined analysis
- **Data Variety**
  - Time series are highly structured; Text is not
- **Open world assumption**
  - Data sources might be added or removed
  - So any analysis is only valid based on the current state of the big data resource
- **Data in movement, and at rest**
  - Data streams complements stored data
  - Some data streams are stored, others are not
- **Lambda Architecture (or variant thereof)**
  - Batch layer; serving layer; speed layer; Analytics

# SMALL DATA

# BIG DATA

Specific questions

***GOAL***

Broad concerns

One location

***LOCATION***

Many locations

Structured

***STRUCTURE***

Varied, unstructured

Single user

***SOURCE***

Many providers

Transient

***LONGEVITY***

Durable

Focused

***MEASUREMENTS***

Broad

Can be recreated

***REPRODUCIBILITY***

Gone if not captured

Small risk

***STAKES***

Big risk

Simple

***INTROSPECTION***

Metadata is vital

Complete

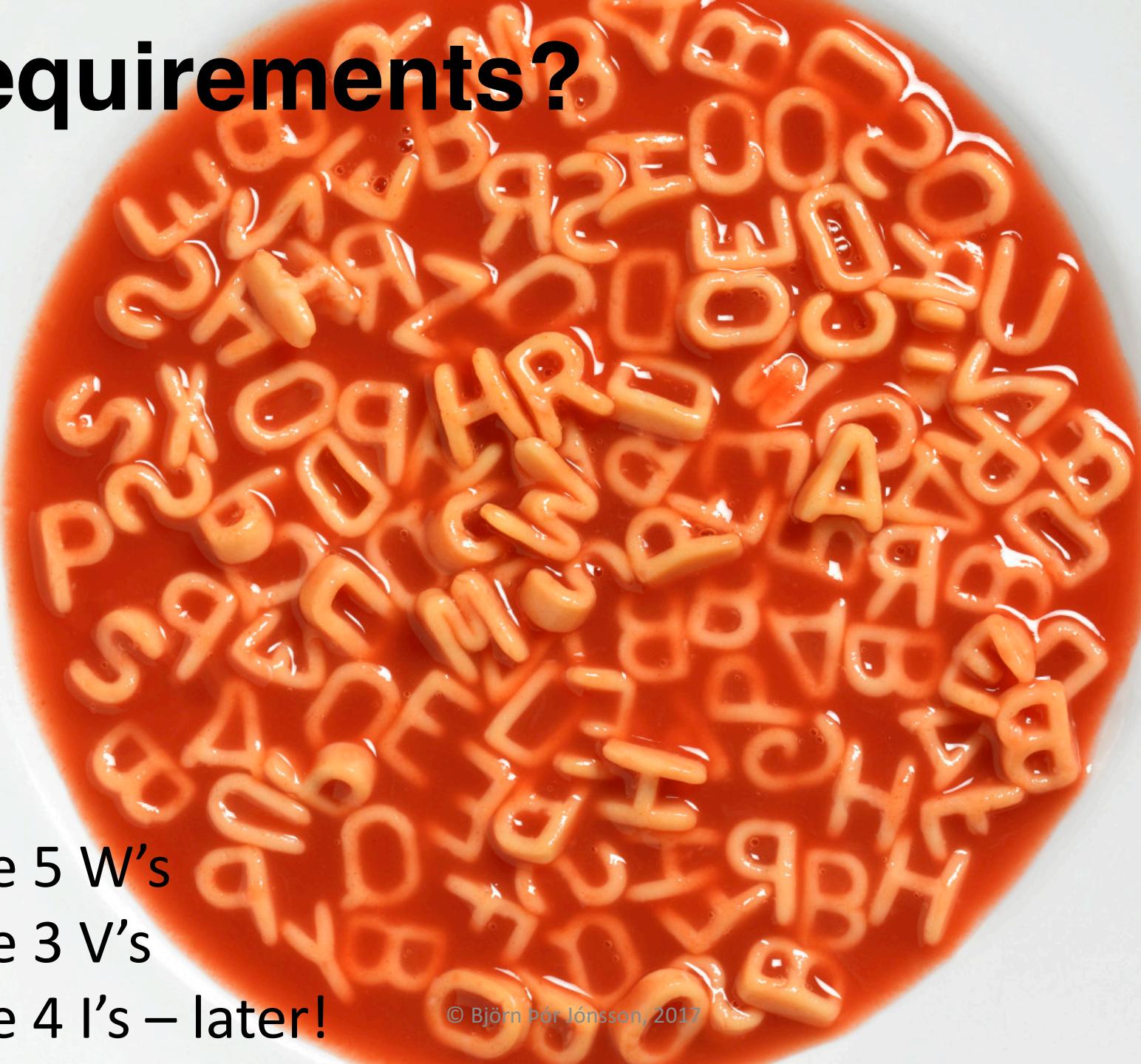
***ANALYSIS***

Incremental

# Why Big Data

- “THE MOST COMMON PURPOSE OF BIG DATA IS TO PRODUCE SMALL DATA”
  - Berman, page xxiv

# Requirements?



- The 5 W's
- The 3 V's
- The 4 I's – later!

# The Three “V”s



Volume



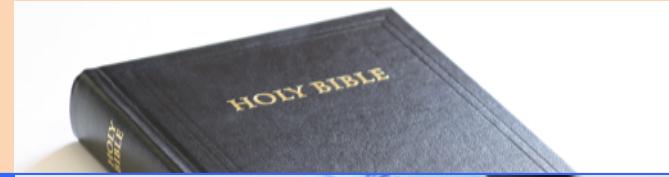
Velocity



Variety



Veracity



Validity



Viability



Value



# The 3 Vs

- Volume, Velocity, Variety: What does it mean?
- At best, the Vs are dimensions to structure:
  - Non-functional requirements
  - Capacity sizing
  - Performance evaluation



# Big Data Analytics: Making Government Data Work

"Big data" comes with many promises, but the data alone is not a silver bullet. True, it holds the potential for extracting business or mission intelligence and improving decision-making, but without the application of expert domain knowledge to give data contextual meaning, big data is nothing but a whole lot of dark figures.

**csc**

A Briefing from GIC  
Industry Insights  
November 2012

**Making Big Data Work for Government**

"Big data" comes with many promises. But the least of these is access to genuinely meaningful information. "Dark data," and large volumes of unstructured information, may estimate place the amount of unstructured data worldwide at 90 percent of all data, and as the average federal agency stores 1.6 petabytes of data,<sup>1</sup> this represents a potentially useful source of extremely unstructured information for analysis.<sup>2</sup> But big data alone is not a silver bullet. True, it holds the potential for extracting business or mission intelligence and improving decision-making, but without the application of expert domain knowledge to give data contextual meaning, big data will be nothing but a whole lot of dark figures.<sup>3</sup>

Currently, federal agencies cannot make use of all their data because they do not (or cannot afford) employ enough data scientists—that is, experts who possess domain knowledge and can use big data analytic technologies to ask the right questions and extract business or mission intelligence from vast pools of data. Making use of big data under these circumstances presents a unique challenge.

**Subduing Big Data**

The technological challenges of capturing, securing and managing the worldwide explosion of data are not insignificant. The amount of worldwide data currently measures about 1.7 petabytes and is projected to double every two years.<sup>4</sup> Yet, controlling big data is as much an organizational challenge as a technological one. The explosion of data has transformed business intelligence.<sup>5</sup> Companies have been investing much time and money in the former,<sup>6</sup> but estimates that big data initiatives in 2010 will total \$34 billion,<sup>7</sup> and big data will drive \$102 billion in spending over the next five years.<sup>8</sup>

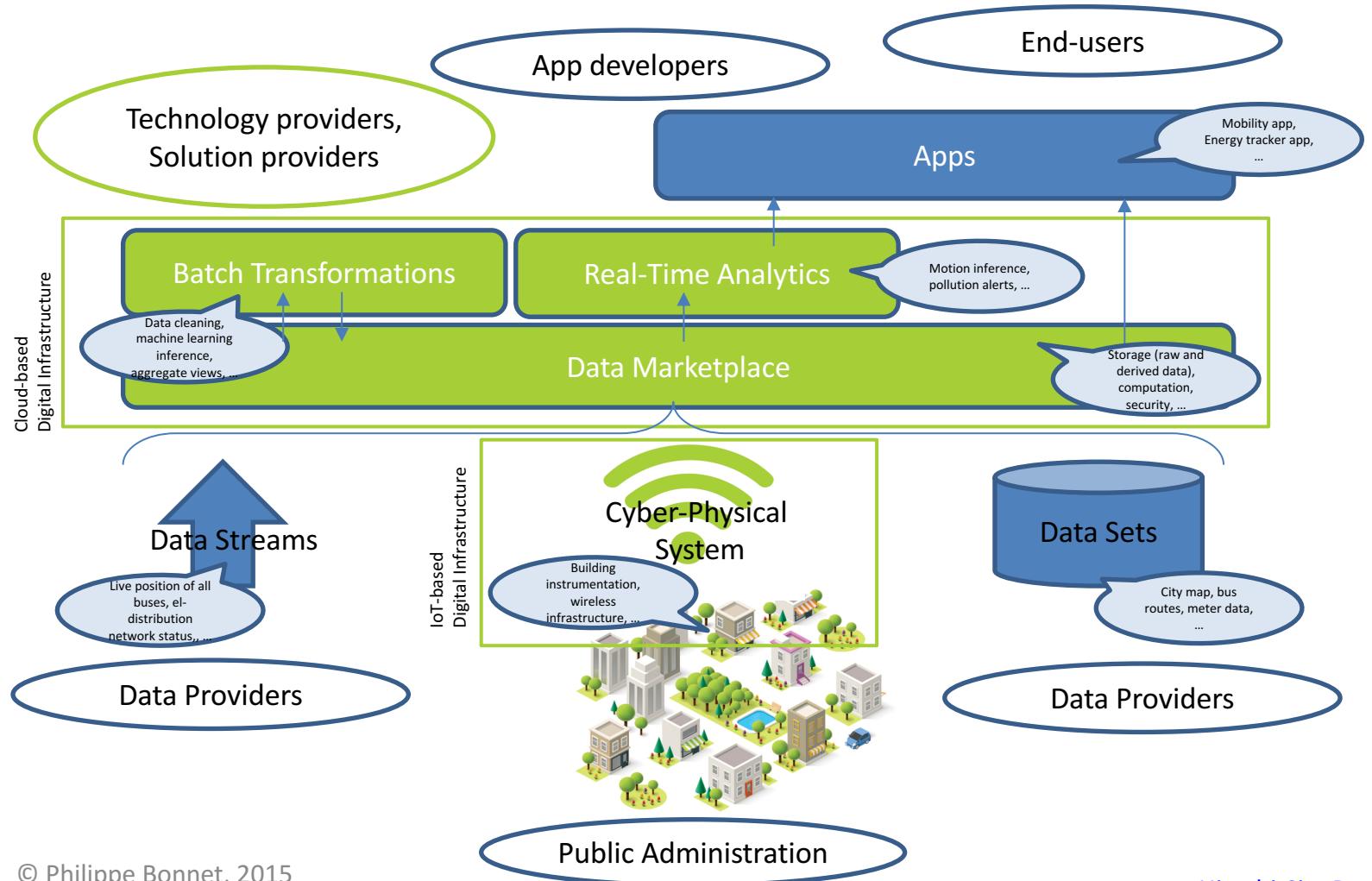
**csc**

Information & Communications Technology Solutions  
Delivery Insights

## Government Big Data

Currently, federal agencies cannot make use of all their government data because they do not (or cannot afford to) employ enough data scientists—that is, experts who possess domain knowledge and can use government big data analytic technologies to ask the right questions and extract business or mission intelligence from vast pools of data. Making use of big data under these circumstances presents a unique challenge.

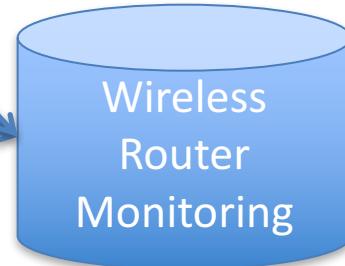
# Big Data at City Scale



# Big Data at Building Scale



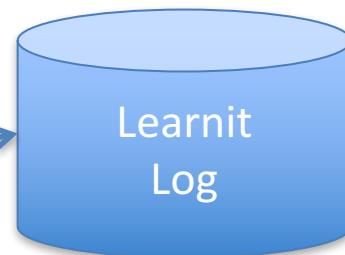
Light on/off events not logged



IT Dept



Facility Management

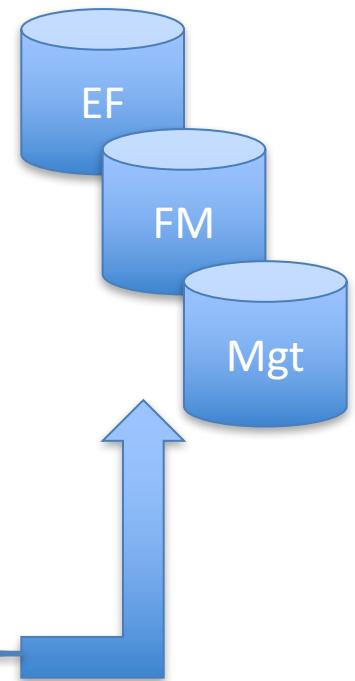
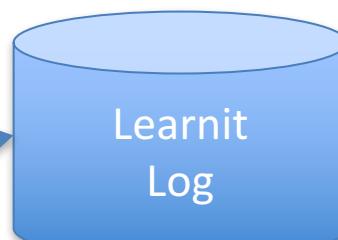
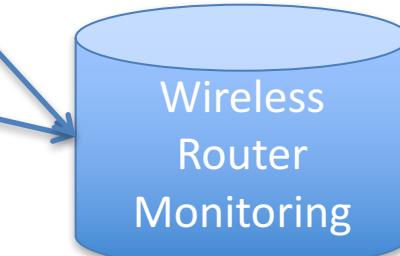


Admin  
Teachers  
Students?

# Big Data at Building Scale

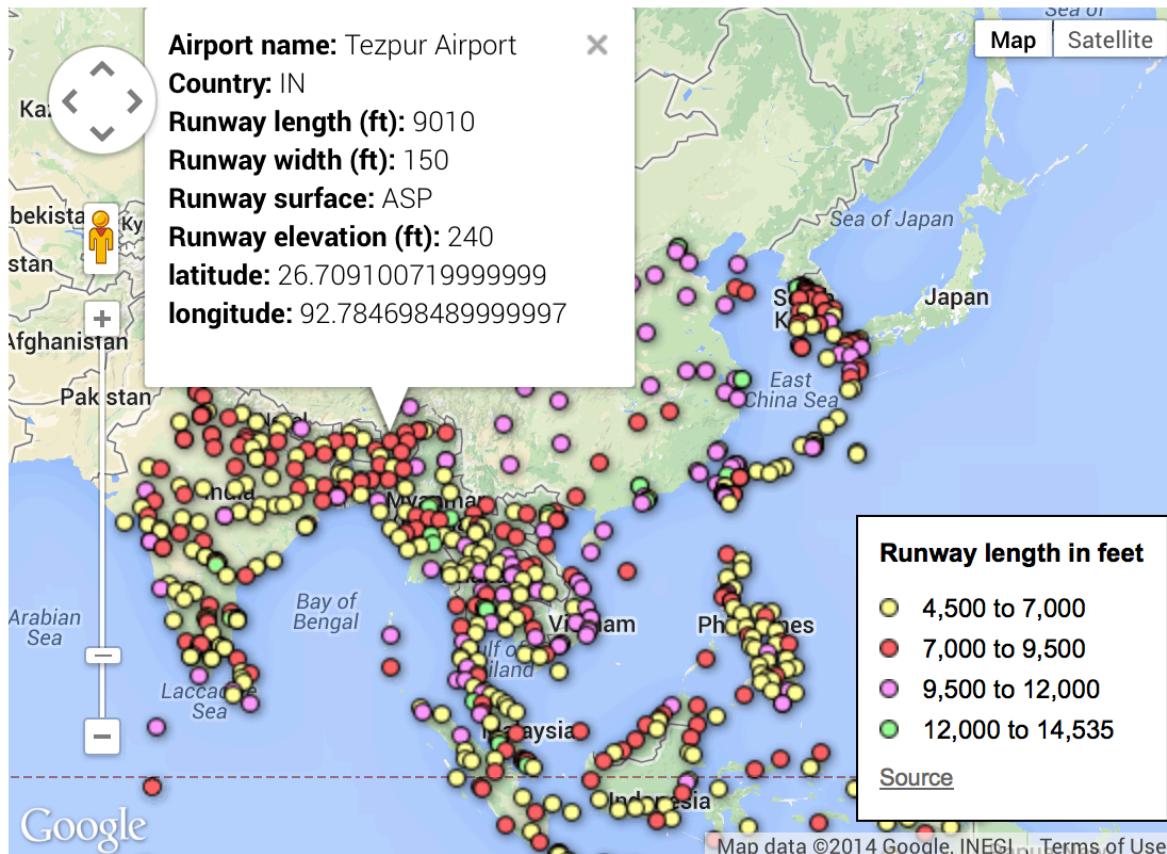


Light on/off events



ITU Big Data Resources available for Analytics within and outside the IT University

# Search and Rescue Example

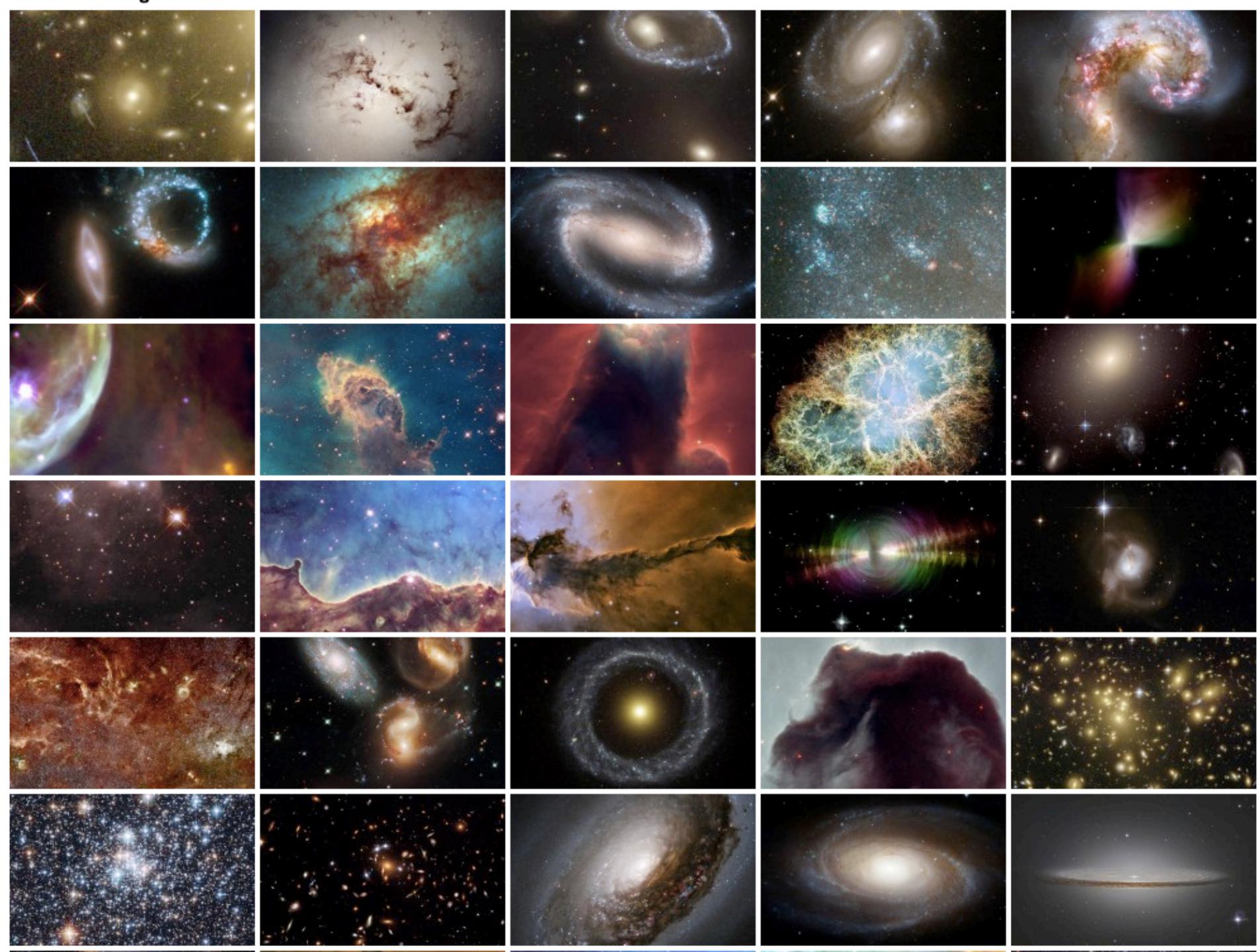


Map by David Strip based on Google maps and OurAirports.com

See [Jame Fallows' article](#) at the Atlantic.



© Björn Þór Jónsson, 2017





© Björn Þór Jónsson, 2017

# The Five “W”s



What?

Who?

Where?

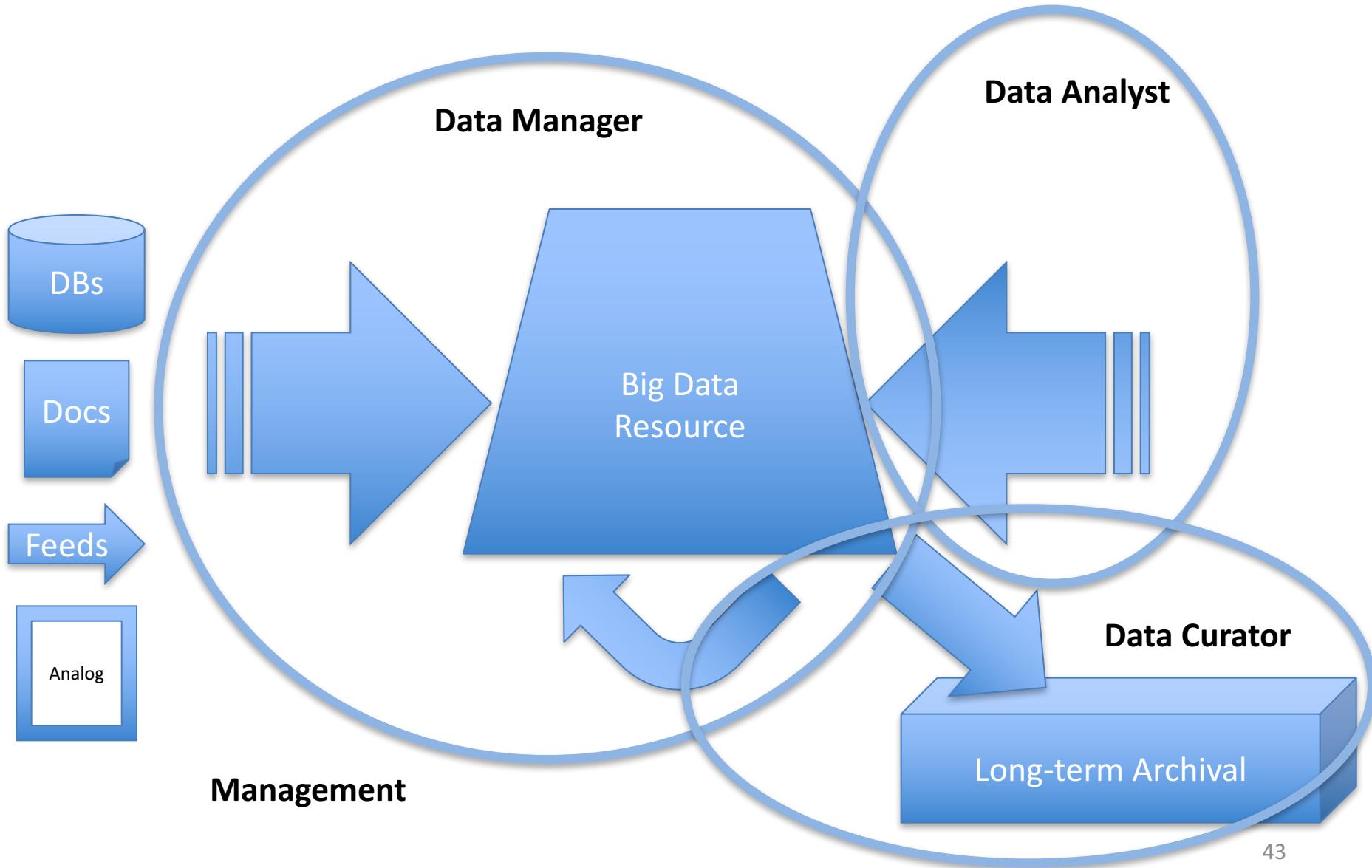
When?

Why?

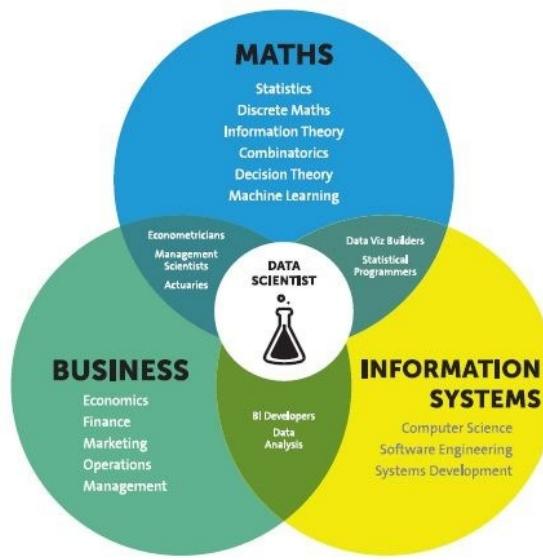
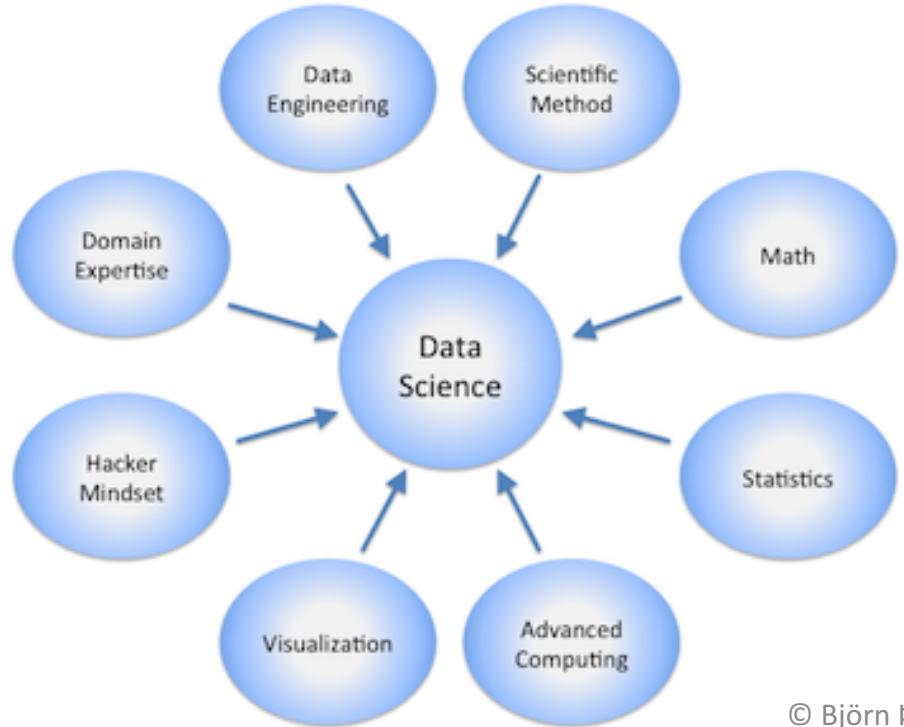
# When does Big Data Make Sense?

- Need for insight based on data not currently available for analysis

# Who is involved?



# Data Science



# Outline

- Trends Underlying Big Data
- What is Big Data?
- **Data Models**
  - Concepts
  - Naming/Identification/Ontologies
  - Structured/Semi-structured/Unstructured
- Course Outline

# Conceptual Model

- Entities
  - An entity is a “thing” or object
  - An entity set is a collection of similar entities
  - An attribute is a property of an entity
- Relationships
  - A relationship connects two or more entity sets
  - Relationship sets are collection of similar relationships
  - An attribute is a property of a relationship

# Logical Data Model

## 1. Mathematical representation of data.

- Examples: set/multiset (relations); trees/graphs (linked data).

## 2. Operations on data.

- Example: Select all items whose acceleration, defined as  $(\sqrt{x^2+y^2+z^2})$  is greater than 1,5

## 3. Constraints

- Example: All valid items should have an acceleration, defined as  $(\sqrt{x^2+y^2+z^2})$  greater than 1,5

# Physical Data Model

- Relational databases
  - SQL tables; indices; partitionning
  - Row or column store mapped onto files
- NoSQL systems
  - Columnar representation of nested data structures
    - Parquet
  - Serialisation/deserialisation (SerDe) of linked data structures onto files
    - Avro, thrift, protocol buffers
    - Includes compression

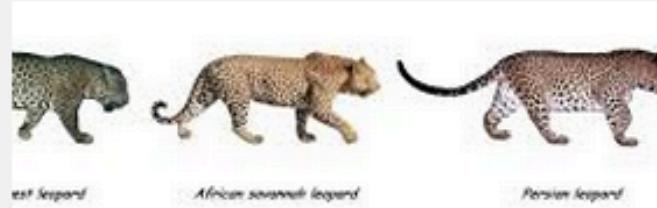
# What is Hard?

- Scoping
  - What is an entity, what is an attribute?
- Naming
  - How to name entities and properties
  - Identification
- Dealing with change
  - Integration across data sources: ontologies
  - In time: Immutability

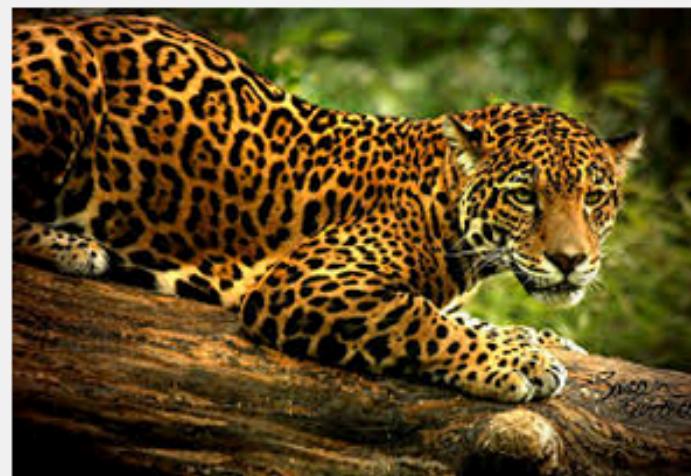
# What's in a name?



Car



Vs Leopard



[https://www.google.com/search?site=&tbo=isch&source=hp&biw=1440&bih=802&q=jaguar&oq=jaguar&gs\\_l=img.3..0l10.1339.2325.0.2503.6.4.0.2.2.0.108.298.1j2.3.0.msedr...0...1ac.1.61.img..1.5.307.a-kqsMXizB4](https://www.google.com/search?site=&tbo=isch&source=hp&biw=1440&bih=802&q=jaguar&oq=jaguar&gs_l=img.3..0l10.1339.2325.0.2503.6.4.0.2.2.0.108.298.1j2.3.0.msedr...0...1ac.1.61.img..1.5.307.a-kqsMXizB4)

# What's in a name?

in® ≡ Search for people, jobs, companies, and more... Ad

Home Profile Connections Jobs Interests

Global Leadership Program - For the executive who's ready to lead beyond b



**Björn Pór Jónsson**  
Software Engineer at Siminn  
Iceland | Computer Software

Current      Siminn  
Previous     HugurAx  
Education    Háskóli Íslands

126

[9] B. T. Jónsson, G. Tómasson, H. Sigurthórsson, Á. Eiríksdóttir, L. Amsaleg, and M. K. Lárusdóttir. A multi-dimensional data model for personal photo browsing. In *Proc. MMM*, Sydney, Australia, 2015.



Connecting Research  
and Researchers

[FOR RESEARCHERS](#)[FOR ORGANIZATIONS](#)[ABOUT](#)[HELP](#)[SIGN IN](#)

## DISTINGUISH YOURSELF IN THREE EASY STEPS

ORCID provides a persistent digital identifier that distinguishes you from every other researcher and, through integration in key research workflows such as manuscript and grant submission, supports automated linkages between you and your professional activities ensuring that your work is recognized. [Find out more.](#)

**1****REGISTER**

Get your unique ORCID identifier [Register now!](#)  
Registration takes 30 seconds.

**2****ADD YOUR  
INFO**

Enhance your ORCID record with your professional information and link to your other identifiers (such as Scopus or ResearcherID or LinkedIn).

**3****USE YOUR  
ORCID ID**

Include your ORCID identifier on your Webpage, when you submit publications, apply for grants, and in any research workflow to ensure you get credit for your work.

### LATEST NEWS

Tue 01/13/2015  
New webinar: The metadata round trip

Mon 01/12/2015  
ORCID Partners with Hypothes.is and NIF on Helmsley Trust-Supported Open Annotation Project

Thu 01/08/2015



1930's - 1950's



1930's - 1950's, Private Vehicle



1930's - 1950's



1950 - 1958, Private Vehicle  
Jim Fox/Picture by Greg Gibson



1958 - 1967, Private Vehicle



1976 Series, Private Vehicle



2008 Series, Private Vehicle



Civil Defense (1)



1953 Series, Army



1953 Series, Airforce



2008, Unified Defense Force (2)



1953 Series, Navy



1976 Series, Army



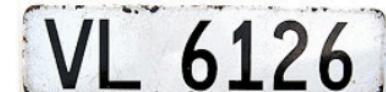
Airforce Home Guard (1)



Before 2012, N.A.T.O. Headquarters (77000 - 77999 Block)



1930's, Motorcycle Dealer



1950's. Dealer (1)



1950's, Commercial, Partially Tax Exempt (\*)



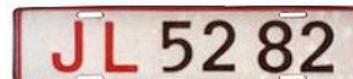
1950's, Temporary



1965, Temporary



2003, Temporary



Temporary Test Plate



Driver Education



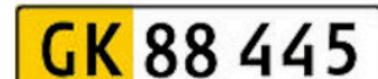
Mid-1990's, Airport Vehicle



1970's, Trailer



Diplomatic Corps (1)



Vehicle Paying Intermediate

Tax Rate



Emergency Services  
Ambulance

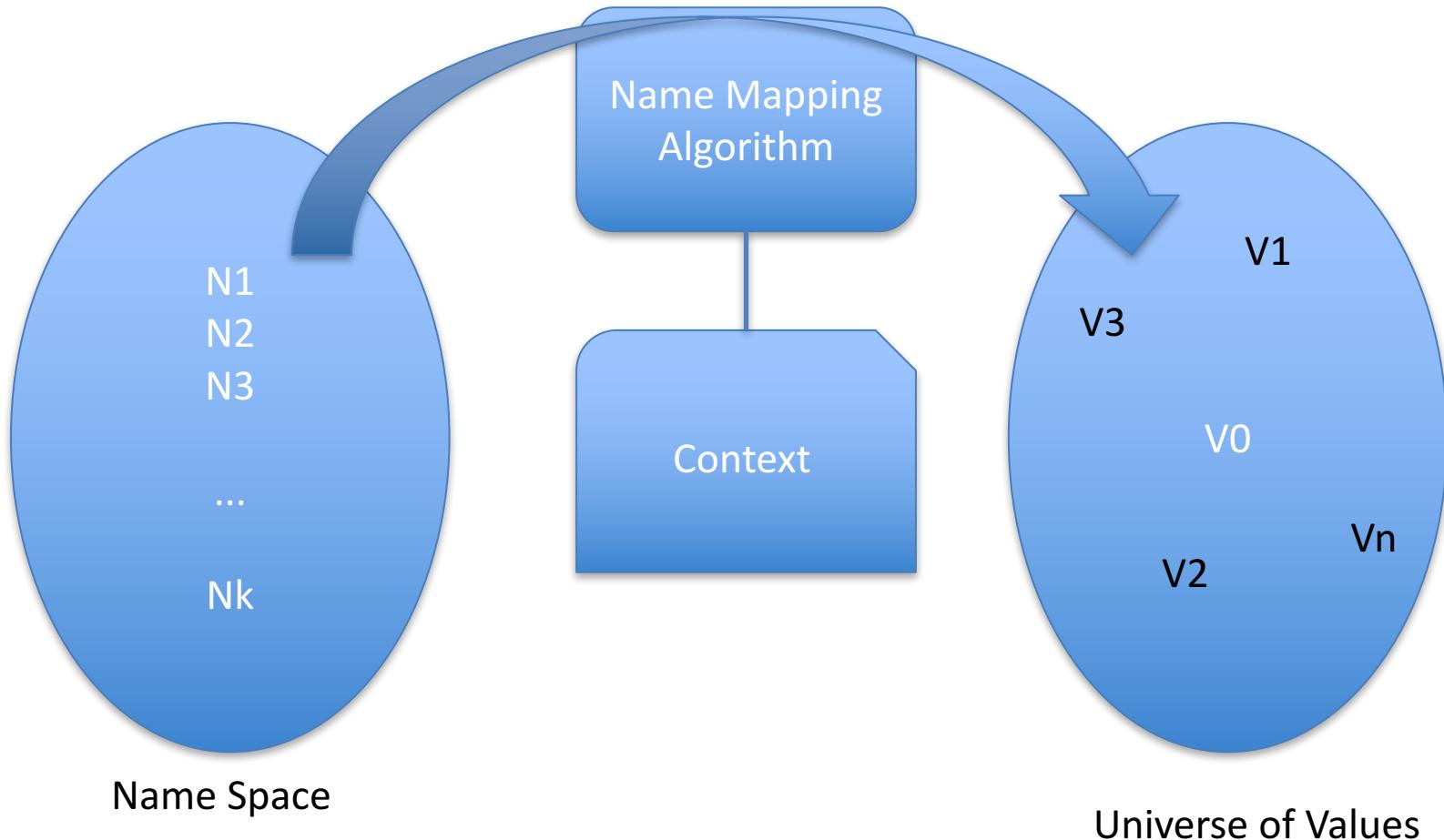


Emergency Services  
Copenhagen Firebrigade (1)



Emergency Services  
Copenhagen Firebrigade (1)

# Name Mapping



# Examples (Database)

- Key constraint
  - Name space: key attribute values
  - Universe of values: tuples in a table
  - Name mapping algorithm: Index
  - Context: A given table
- Sequence Id
  - Name space: sequence of integers
  - Universe of values: tuples in a table
  - Name mapping algorithm: Index
  - Context: A given table

# Example (File System)

- File System Name
  - Name space: Path name
  - Universe of values: File Identifiers
  - Name mapping algorithm: recursive mapping
  - Context: Root directory

# Example (Unix File System)

| Layer              | Names               | Values        | Context                 | Name-mapping algorithm      |
|--------------------|---------------------|---------------|-------------------------|-----------------------------|
| Symbolic link      | Path names          | Path names    | The directory hierarchy | PATHNAME_TO_GENERAL_PATH    |
| Absolute path name | Absolute path names | Inode numbers | The root directory      | GENERALPATH_TO_INODE_NUMBER |
| Path name          | Relative path names | Inode numbers | The working directory   | PATH_TO_INODE_NUMBER        |
| File name          | File names          | Inode numbers | A directory             | NAME_TO_INODE_NUMBER        |
| Inode number       | Inode numbers       | Inodes        | The inode table         | INODE_NUMBER_TO_INODE       |
| File               | Index numbers       | Block numbers | An inode                | INDEX_TO_BLOCK_NUMBER       |
| Block              | Block numbers       | Blocks        | The disk drive          | BLOCK_NUMBER_TO_BLOCK       |

The diagram illustrates the layers of a Unix file system, organized into seven horizontal layers. From top to bottom, the layers are:

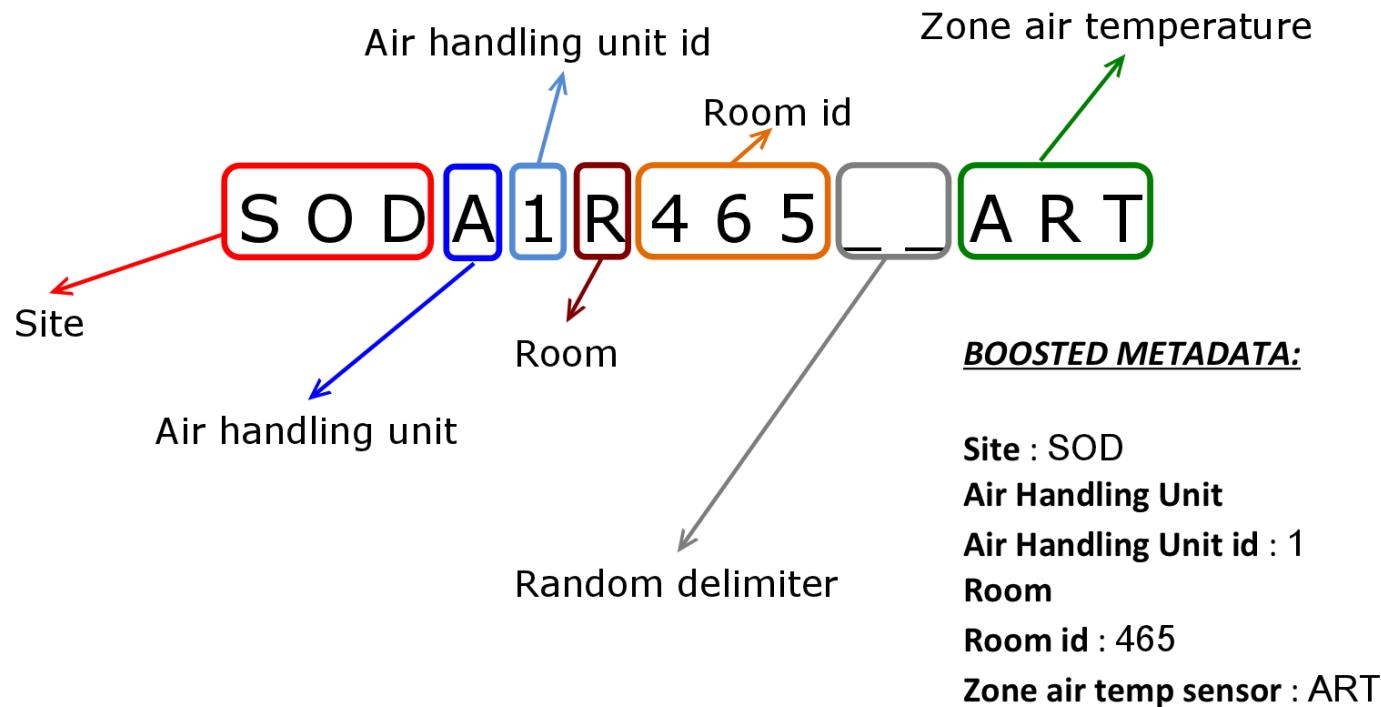
- Symbolic link**: Maps **Path names** to **Path names** within the **directory hierarchy**, using the **PATHNAME\_TO\_GENERAL\_PATH** algorithm.
- Absolute path name**: Maps **Absolute path names** to **Inode numbers** in the **root directory**, using the **GENERALPATH\_TO\_INODE\_NUMBER** algorithm.
- Path name**: Maps **Relative path names** to **Inode numbers** in the **working directory**, using the **PATH\_TO\_INODE\_NUMBER** algorithm.
- File name**: Maps **File names** to **Inode numbers** in a **directory**, using the **NAME\_TO\_INODE\_NUMBER** algorithm.
- Inode number**: Maps **Inode numbers** to **Inodes** in the **inode table**, using the **INODE\_NUMBER\_TO\_INODE** algorithm.
- File**: Maps **Index numbers** to **Block numbers** in an **inode**, using the **INDEX\_TO\_BLOCK\_NUMBER** algorithm.
- Block**: Maps **Block numbers** to **Blocks** in the **disk drive**, using the **BLOCK\_NUMBER\_TO\_BLOCK** algorithm.

Vertical arrows on the right side of the table indicate the flow of data:

- An upward arrow labeled **user-oriented names** points from the **Symbolic link** layer to the **Path name** layer.
- A downward arrow labeled **machine-user interface** points from the **Path name** layer to the **File name** layer.
- An upward arrow labeled **machine-oriented names** points from the **File** layer to the **Block** layer.
- A downward arrow labeled **machine-oriented names** points from the **Block** layer back up to the **File** layer.

Figure from Saltzer and Kaashoek

# Logical Identifiers



# Identifiers

## Properties

- Size
  - Number of possible values
  - Type
- Uniqueness
  - Deterministic
- Randomness
  - Non-deterministic
  - Unique
  - Deterministic

## Generation

- 1 place of generation

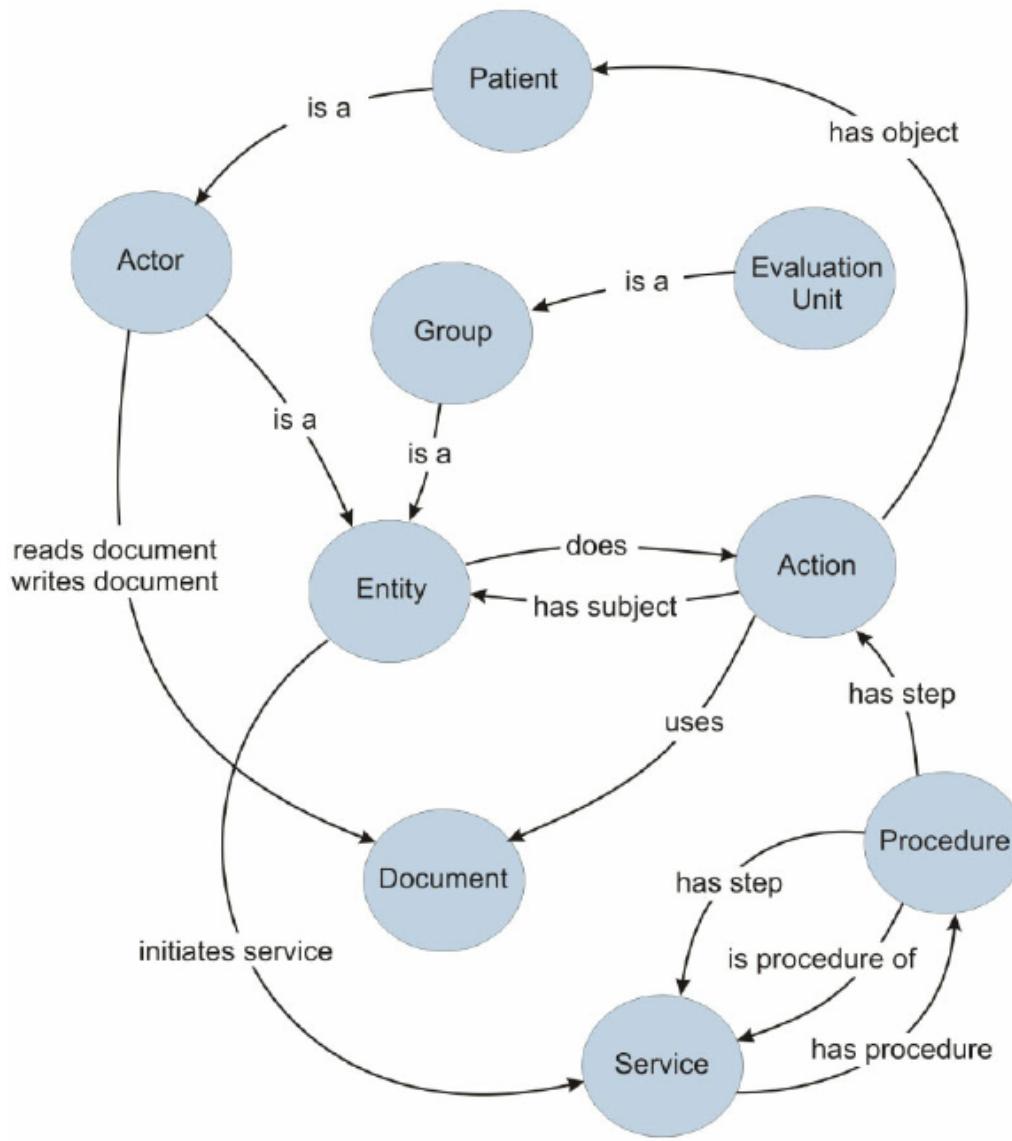


<http://becuo.com/despicable-me-minions-whaaaat>

# Ontology

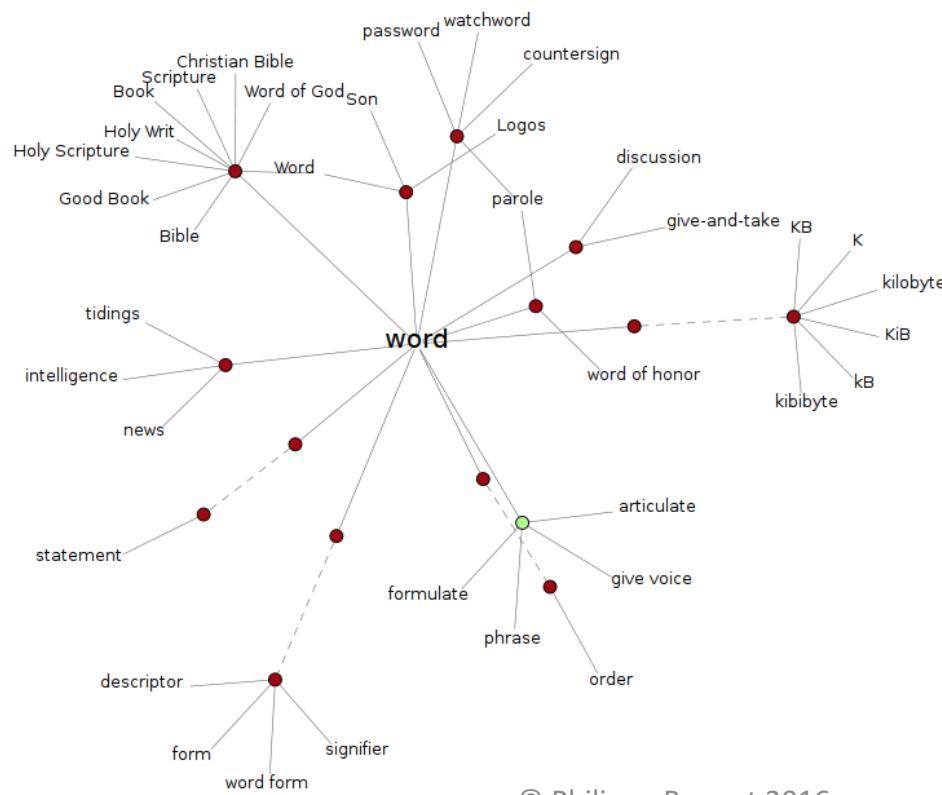
- Facts:
  - Concepts (entity sets) and individuals (entities)
    - attributes or properties
      - changed by some event
- Axioms:
  - Named relationships between concepts
    - Hierarchy
    - Direct Acyclic Graph
    - Arbitrary Graph
- Rules:
  - if-then rules to derive new facts

# Healthcare Ontology: K4Care



# Wordnet

- Synset (synonym ring)
  - Group of entities that are semantically equivalent



# Structured Data

- Instance and schema are clearly separated
  - relational, object-oriented

*Data part:*

```
(#123, [ {[“John”, s111111111, [123,”Main St”]],  
          [“Joe”, s222222222, [321, “Pine St”]] }  
    ] )
```

*Schema part:*

```
PersonList[ LastName: String,  
            Contents: [ Name: String,  
                         Id: String,  
                         Address: [Number: Integer, Street: String] ]  
        ]
```

# Semi-Structured Data

- *Self-describing:*
  - Attribute names embedded in the data itself, *but are distinguished from values*
  - Doesn't need schema to figure out what is what (but schema might be useful nonetheless)

```
(#12345,  
 [ListName: "Students",  
  Contents: { [ Name: "John Doe",  
              Id: "s111111111",  
              Address: [Number: 123, Street: "Main St."] ] ,  
              [Name: "Joe Public",  
               Id: "s222222222",  
               Address: [Number: 321, Street: "Pine St."] ] }  
 ] )
```

# XML (<https://www.w3.org/XML/>)

- Must have a *root element*
- Every *opening tag* must have matching *closing tag*
- Elements must be *properly nested*  
`<foo><bar></foo></bar>` is a no-no
- An *attribute* name can occur *at most once* in an opening tag. If it occurs,
  - It *must have an explicitly specified value* (Boolean attrs, like in HTML, are not allowed)
  - The value *must be quoted* (with " or ')
- *XML processors are not supposed to try and fix ill-formed documents (unlike HTML browsers)*

# JSON (<http://json.org>)

- A JSON object is an unordered collection of name-value pairs.
  - Each name is followed by :
  - A value is a string (in between quotes), a number, true/false, null, an object or an array; values can be nested
  - Name-value pairs are separated by ,
  - An object starts with { and ends with }
- A JSON array is an ordered collection of values
  - An array starts with [ and ends with ]

```
{  
  "/sensor0": {  
    "Metadata": {  
      "SourceName": "Test Source",  
      "Location": { "City": "Berkeley" }  
    },  
    "Properties": {  
      "Timezone": "America/Los_Angeles",  
      "UnitofMeasure": "Watt",  
      "ReadingType": "double"  
    },  
    "Readings": [[1351043674000, 0], [1351043675000, 1]],  
    "uuid": "d24325e6-1d7d-11e2-ad69-a7c2fa8dba61"  
  }  
}
```

# Representing Documents

- BLOB: unstructured
  - Need to extract structure for further processing
  - Term extraction, NLP
- Entity with text attributes
  - Structured data
  - Semi-structured data

# Representing Media

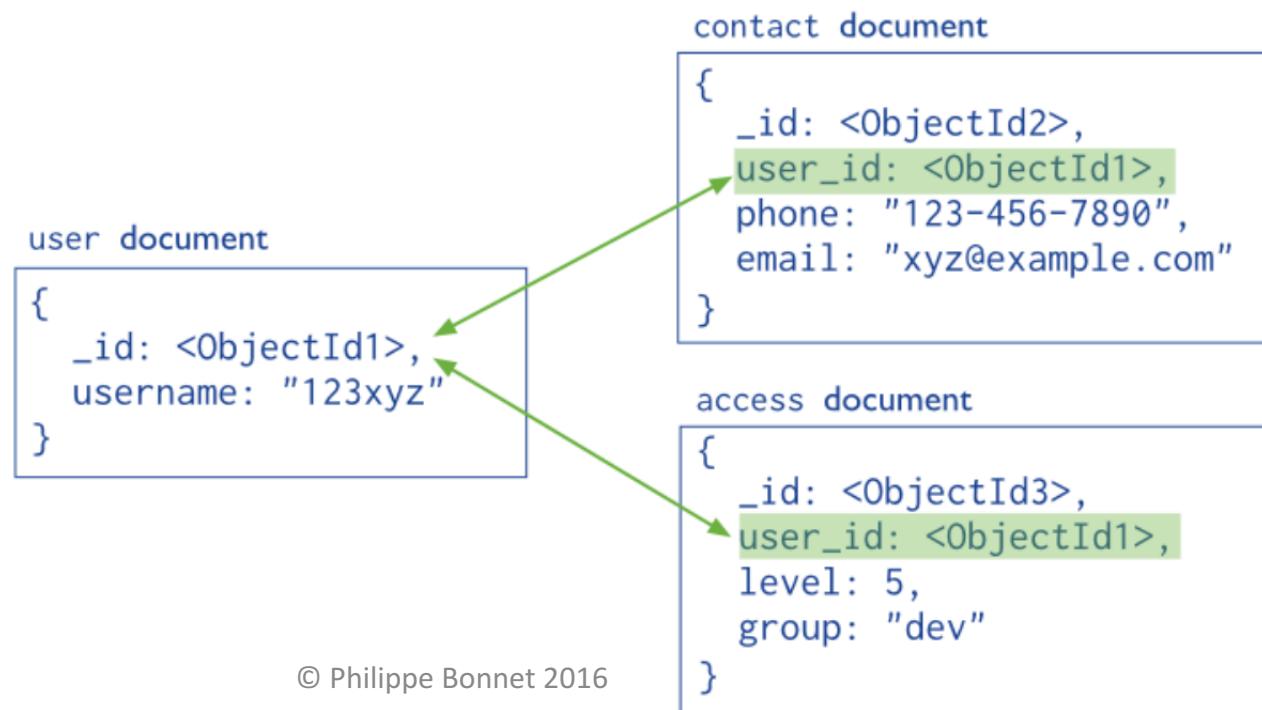
- BLOB: unstructured
  - Need to extract structure for further processing
  - Concept extraction, deep learning
- Entity with text attributes
  - Structured data
  - Semi-structured data
  - Always need the BLOB!?

# Linked Data (Tree, Graph) in SQL

- Adjacency List  $R(N1, N2)$
- Denormalization  $R(N1, ListOfPaths)$
- Preorder traversal  $R(N, Lft, Rgt)$

# Linked Data in MongoDB

- Entities represented as Json objects, called documents
  - Identification through `_id` attribute
- Relationships represented as references



# Linked Data in Key-Value Stores

- Entities and relationships represented as associative arrays (or maps)
  - (Ordered) collection of (key,value) pairs
  - Each key occurs only once in the collection
  - Interface
    - add/remove (key,value) pairs
    - modify value for a given key (or a range of keys)
    - Look up a value with a given key
- Systems:
  - [Cassandra](#), [LevelDB](#), RocksDB, Aerospike ...

# Linked Data in Big Table

- A Bigtable is a sparse, distributed, persistent multidimensional sorted map.
  - The map is indexed by a dimension key, and a timestamp; each value in the map is an uninterpreted array of bytes.
- Well-suited for
  - Relationships that can be denormalized onto a multidimensional space

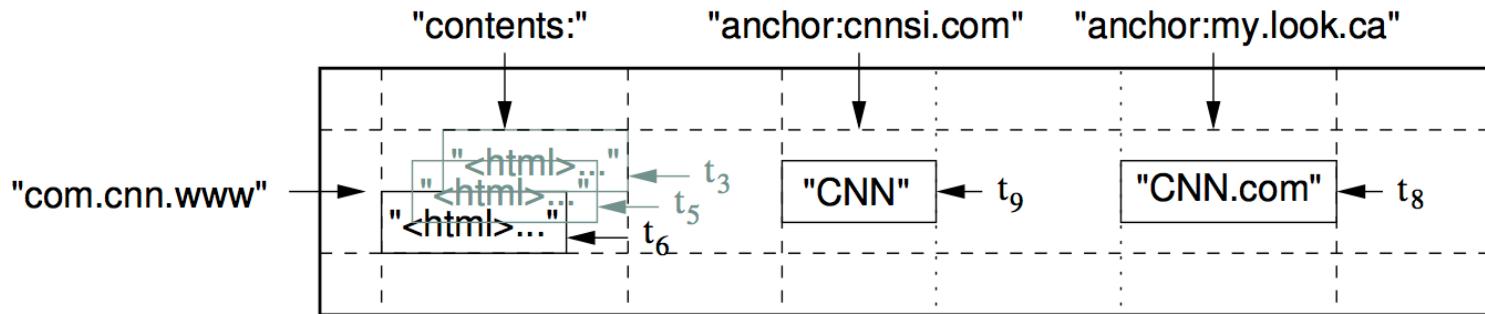
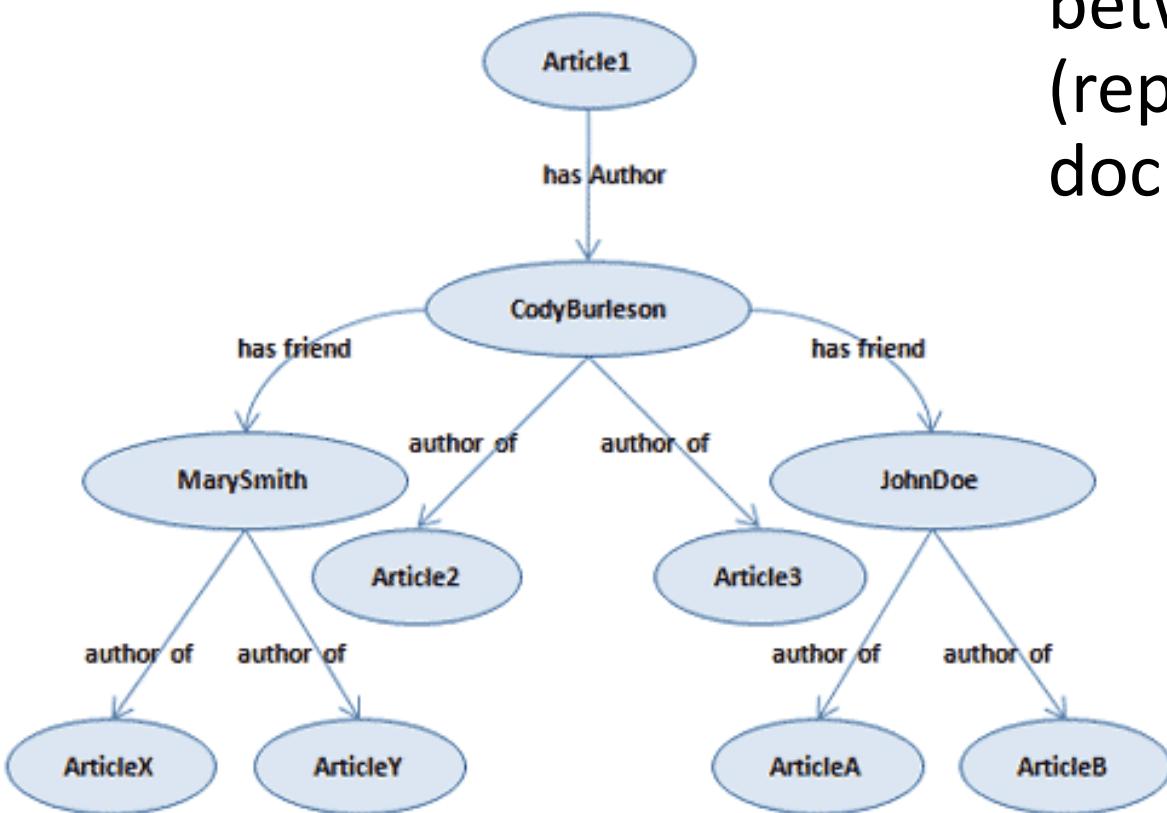


Figure 1: A slice of an example table that stores Web pages. The row name is a reversed URL. The **contents** column family contains the page contents, and the **anchor** column family contains the text of any anchors that reference the page. CNN's home page is referenced by both the Sports Illustrated and the MY-look home pages, so the row contains columns named **anchor:cnnsi.com** and **anchor:my.look.ca**. Each anchor cell has one version; the contents column has three versions, at timestamps  $t_3$ ,  $t_5$ , and  $t_6$ .

# Linked Data in RDF

Directed, labeled graph, where the edges represent relationships between two entities (represented as XML documents).



# Linked Data as Graphs

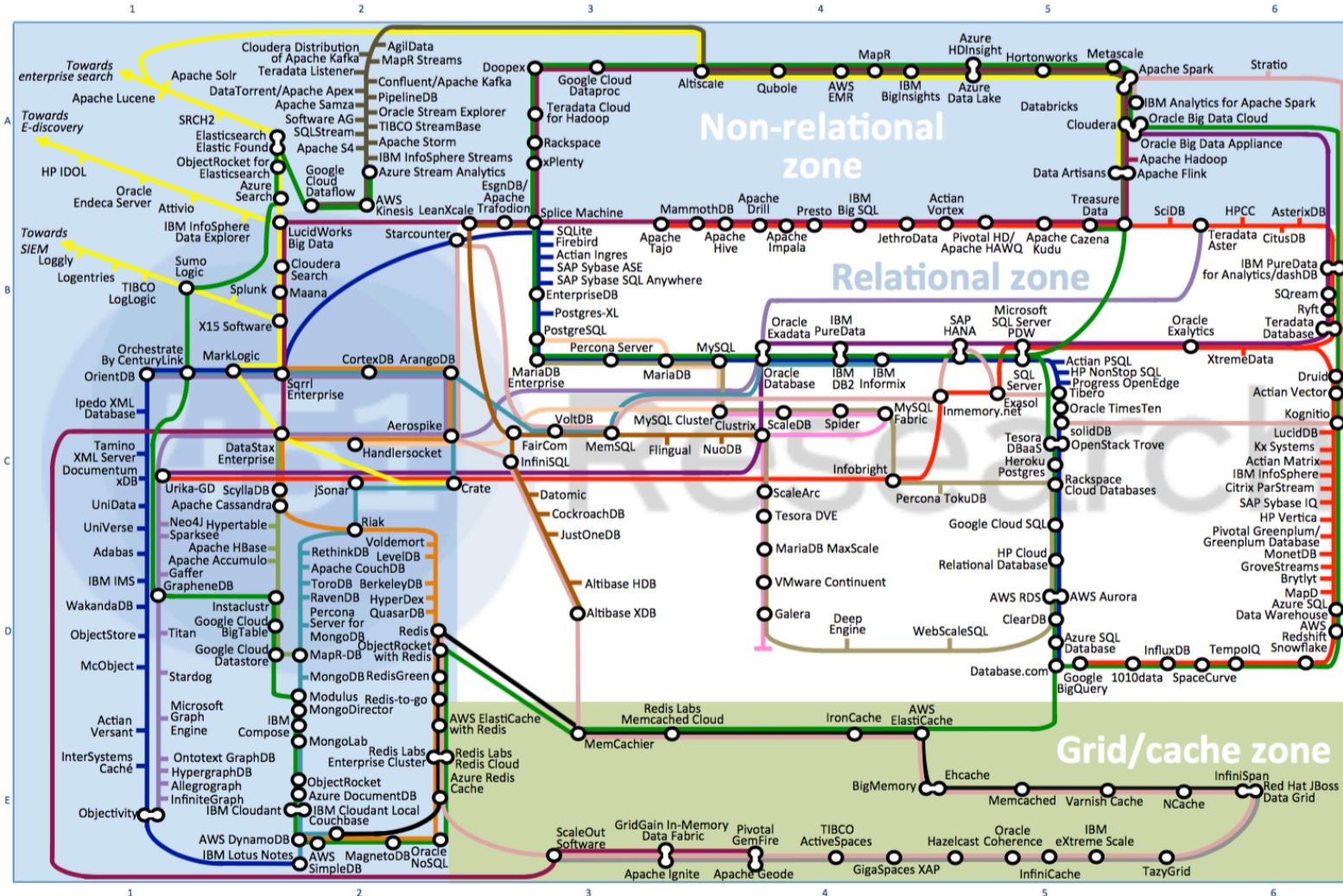
- Graphs can be represented with adjacency lists
  - Indexed adjacency list (relations)
    - Stored as a map or a relation
  - Index-free adjacency list (graphs)
    - Pointers from one node to the other
    - Data structure underlying graph databases
- Systems: [Neo4J](#), [GraphX \(Spark\)](#), [GraphLab](#), ...

# Data Platforms Map

January 2016

**Key:**

- General purpose
- Specialist analytic
- as-a-Service
- BigTables
- Graph
- Document
- Key value stores
- Key value direct access
- Hadoop
- MySQL ecosystem
- Advanced clustering/sharding
- New SQL databases
- Data caching
- Data grid
- Search
- Appliances
- In-memory
- Stream processing



[https://451research.com/  
state-of-the-database-landscape](https://451research.com/state-of-the-database-landscape)

© 2016 by 451 Research LLC.  
All rights reserved

# Metadata

Metadata is 'data about data' and can be divided into four basic classifications:

- **Business Metadata** - the business meaning of data. It includes business definitions of the objects and metrics, hierarchies, business rules, and aggregation rules.
- **Operational Metadata** - Operational metadata stores information about who accessed what and when.
- **Technical Metadata** - Technical Metadata describes the data structures and formats such as table types, data types, indexes, and partitioning method.
- **Process Metadata** - Process Metadata describes the data input process.

## Typical Use Cases

- **Introspection** - discovering and harvesting of information into a metadata repository
- **Impact Analysis** - analyse dependencies across the architecture. Requires end to end view of interdependencies
- **Data Lineage** - description of the origins of a piece of data and the proves by which it arrived in the database. The 'provenance' or 'pedigree.'

# Metadata

- Only if every user has a common and exact understanding of the data can it be exchanged trouble-free.
  - ISO/IEC 11179 Metadata Registry Specification
- “Microsoft says nobody uses metadata!”
  - <https://social.technet.microsoft.com/Forums/en-US/1fb1adc5-b7ec-48e3-bb1a-ebeff70e4c0e/microsoft-says-nobody-used-metadata-so-we-don?forum=sharepointgeneralprevious>

# Why Metadata?

## Small Data

- Specific questions
- One location
- Structured
- Single user
- Transient
- Focused
- Can be recreated
- Small risk
- Simple
- Complete

*GOAL*

*LOCATION*

*STRUCTURE*

*SOURCE*

*LONGEVITY*

*MEASUREMENTS*

*REPRODUCIBILITY*

*STAKES*

*INTROSPECTION*

*ANALYSIS*

## Big Data

- Broad concerns
- Many locations
- Varied, unstructured
- Many providers
- Durable
- Broad
- Gone if not captured
- Big risk
- Metadata is vital
- Incremental

# Take Away Points

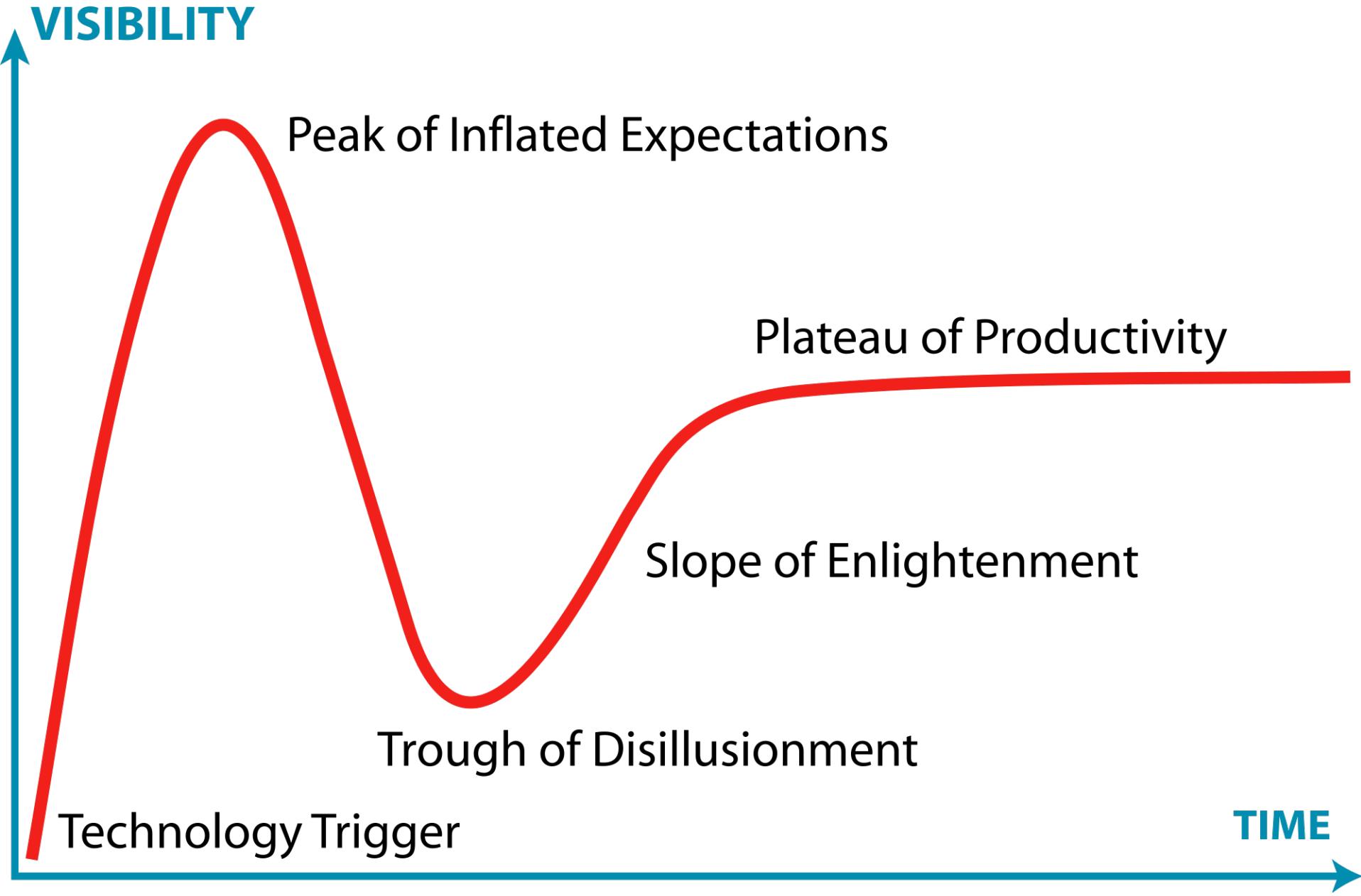
- Conceptual model based on entities and relationships
  - Focus on naming, identification, versioning (changes in time) and ontologies (integration)
- Structured data maps ER onto relations
- Semi-structured data maps ER onto linked data represented with MongoDB, KV store, BigTable, RDF, Graphs

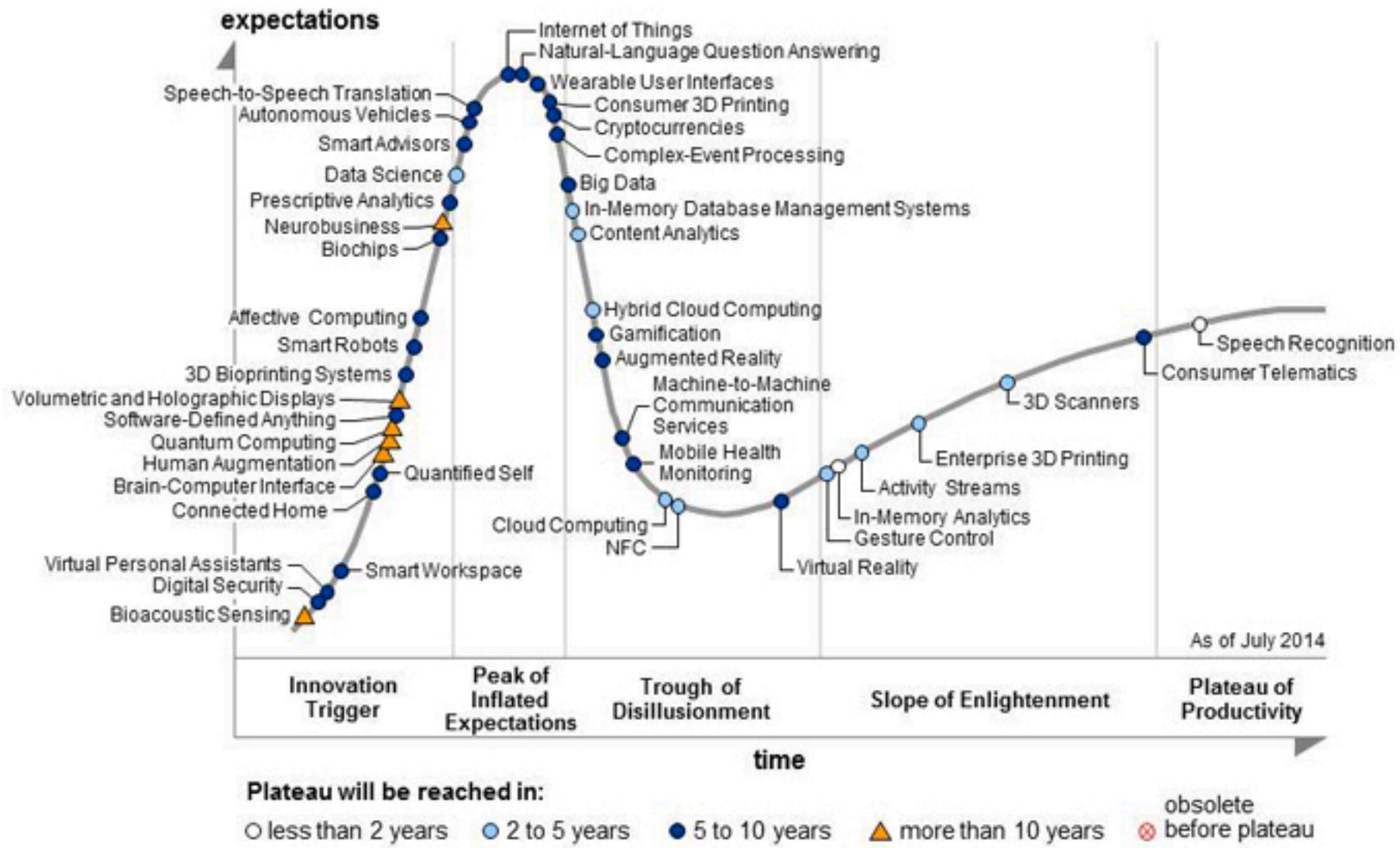
# Big Data: Who's doing it?

[Gartner](#): “64% of enterprises surveyed indicate that they're deploying or planning Big Data projects. Yet even more acknowledge that they still don't know what to do with Big Data.”

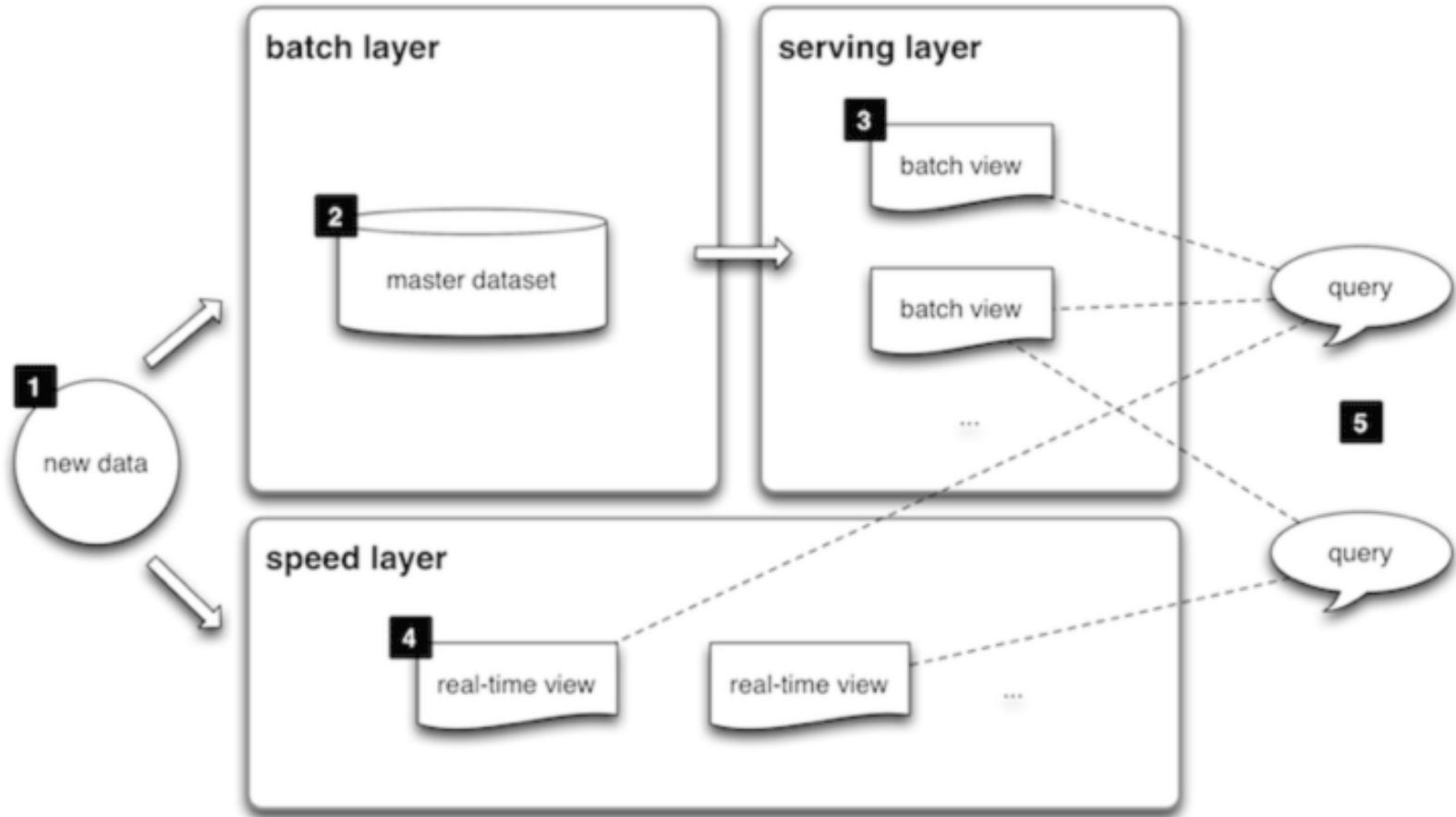
**“Big data is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it”**

Popular Saying





# Lambda Architecture



<http://www.drdobbs.com/database/applying-the-big-data-lambda-architectur/240162604>

# “teenage sex” analogy cont’d...

- From [LinkedIn](#):
  - After the initial excitement you wonder what all the fuss was about?
  - The positive side of this analogy, those who really know what they are doing, do fall in love and conceive some great ideas and make their data exploration a wonderful experience.

# Take Away Points

1. Big data is not a product but a collection of processes centered around big data resources
  - Collections of data made available for analysis
2. Lambda Architecture is a good way to organize Big Data management
3. Joint focus on data manager and data analyst
  - Not a data mining class, however!

# Outline

- Trends Underlying Big Data
- What is Big Data?
- Data Models
- **Course Outline**