# Intro to Big Data

## ITU ASP COE TRAINING ON
## "Developing the ICT ecosystem to harness IoTs"
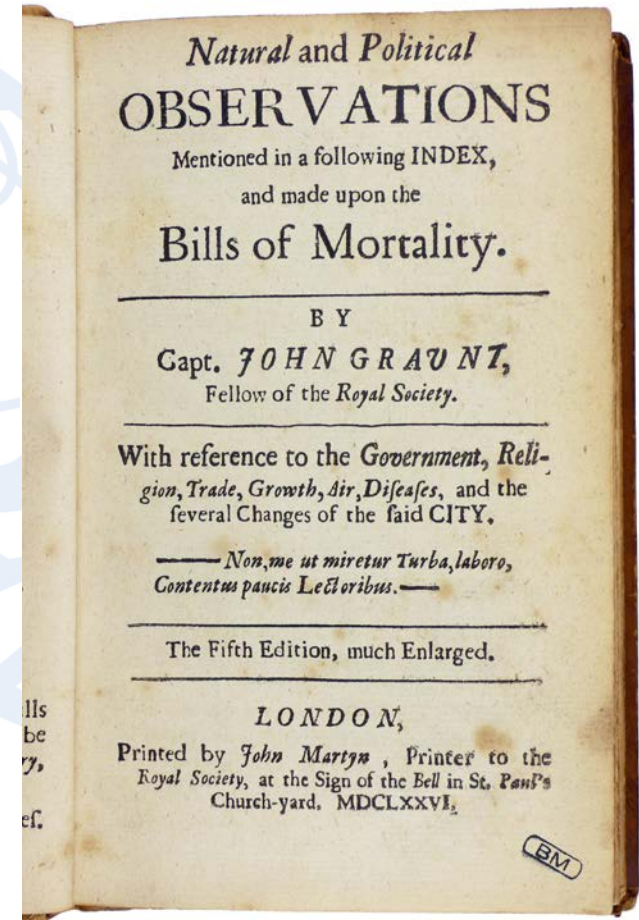
**Marco Zennaro, PhD**
**13-15 December 2016**
**Bangkok, Thailand**

# History of Big Data

- **300 BC – 48 AD** The Library of Alexandria is the world's largest data storage center – until it is destroyed by the Romans.

- 180,000 books, for a total of 20 Gb

# History of Big Data

- **1663** John Graunt conducts the first recorded statistical-analysis experiments in an attempt to curb the spread of the bubonic plague in Europe.



Natural and Political

OBSERVATIONS

Mentioned in a following INDEX,
and made upon the

Bills of Mortality.

BY

Capt. JOHN GRAUNT,

Fellow of the Royal Society.

With reference to the Government, Religion, Trade, Growth, Air, Diseases, and the several Changes of the said CITY.

—— Non, me ut miretur Turba, laboro,
Contentus paucis Lectoribus.——

The Fifth Edition, much Enlarged.

LONDON,

Printed by John Martyn, Printer to the Royal Society, at the Sign of the Bell in St. Paul's Church-yard. MDCLXXVI.

# History of Big Data

- **1881** Herman Hollerith creates the Hollerith Tabulating Machine which uses punch cards to vastly reduce the workload of the US Census. He is one of the founders of IBM.

- **1926** Nikola Tesla predicts that in the future, a man will be able to access and analyze vast amounts of data using a device small enough to fit in his pocket.

# History of Big Data

- **1965** The US Government plans the world's first data center to store 742 million tax returns and 175 million sets of fingerprints on magnetic tape.

- **1989** Early use of term Big Data in magazine article by fiction author Erik Larson – commenting on advertisers' use of data to target customers.
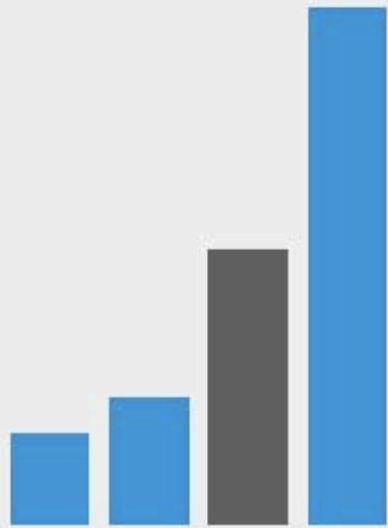
# History of Big Data

- **2010** Eric Schmidt, executive chairman of Google, tells a conference that as much data is now being created every two days, as was created from the beginning of human civilization to the year 2003.

# Definition

"Big data is a broad term for data sets so large or complex that traditional data processing applications are inadequate. Challenges include analysis, capture, curation, search, sharing, storage, transfer, visualization, and information privacy. The term often refers simply to the use of predictive analytics or other certain advanced methods to extract value from data, and seldom to a particular size of data set." – Wikipedia
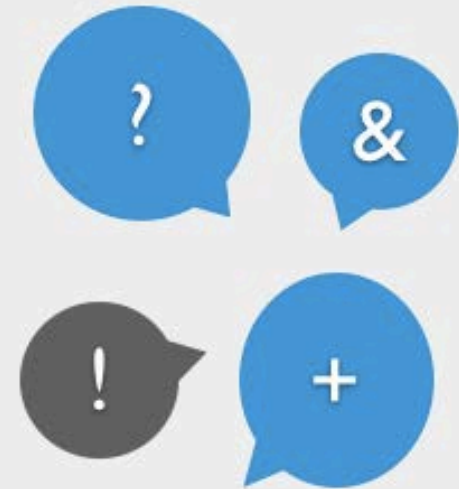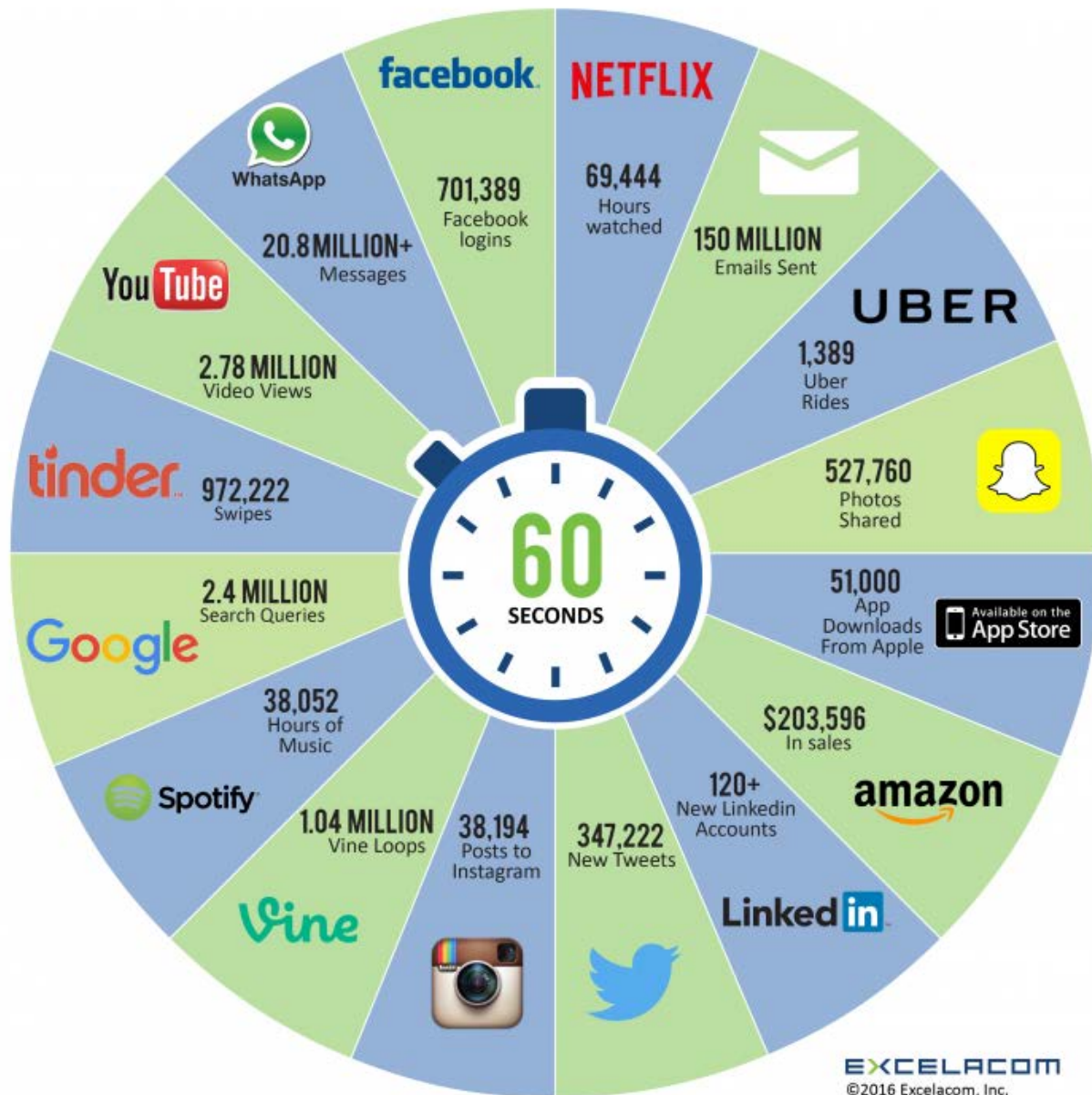
# Three Characteristics of Big Data
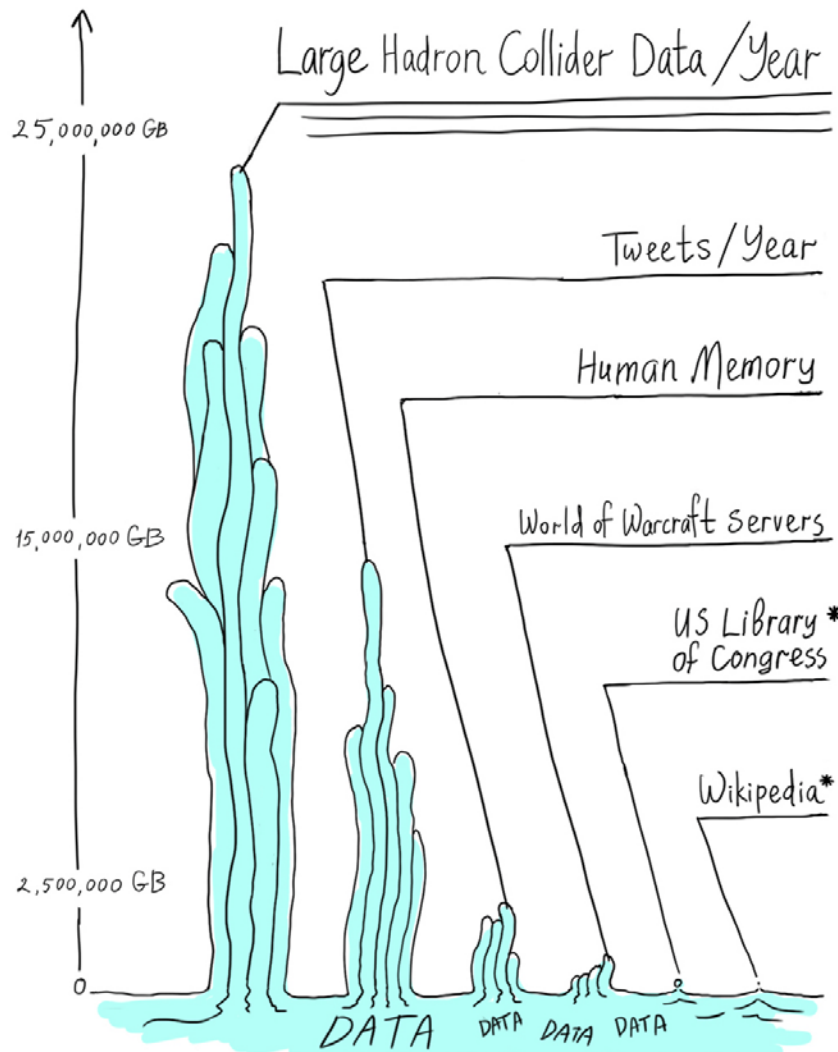


Volume    Velocity    Variety

# Volume

- The sheer size of data in terms of storage and access.
- For example: unstructured data from social media in form of posts, video, audio with relational data such as comments, discussions, likes, etc.

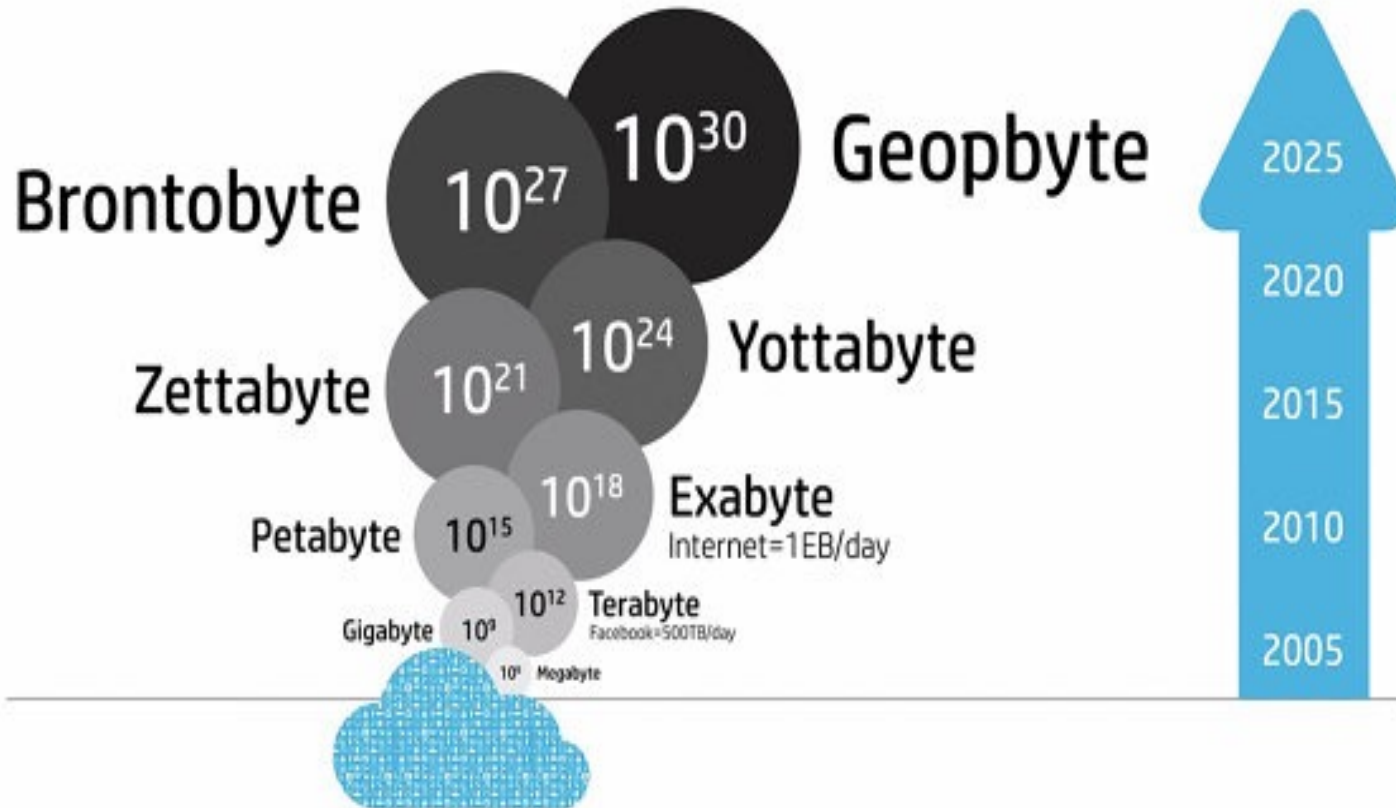# 2016 What happens in an INTERNET MINUTE?

**facebook** — 701,389 Facebook logins

**WhatsApp** — 20.8 MILLION+ Messages

**You Tube** — 2.78 MILLION Video Views

**tinder** — 972,222 Swipes

**Google** — 2.4 MILLION Search Queries

**Spotify** — 38,052 Hours of Music

**Vine** — 1.04 MILLION Vine Loops

**Instagram** — 38,194 Posts to Instagram

**Twitter** — 347,222 New Tweets

**Linked in** — 120+ New Linkedin Accounts

**amazon** — $203,596 In sales

**App Store** — 51,000 App Downloads From Apple

**Snapchat** — 527,760 Photos Shared

**UBER** — 1,389 Uber Rides

**Email** — 150 MILLION Emails Sent

**NETFLIX** — 69,444 Hours watched

## 60 SECONDS

EXCELACOM
©2016 Excelacom, Inc.

- The Large Hadron Collider (LHC) will generate 60 terabytes of data per day, 25 petabytes annually

- Wallmart generates 2.5 petabytes per hour

# Volume

# Sensor data from a cross-country flight

20 TB ✖ 2 ✖ 6 ✖ 28,537 ✖ 365

20 terabytes of information per engine every hour

twin-engine Boeing 737

six-hour, cross-country flight from New York to Los Angeles

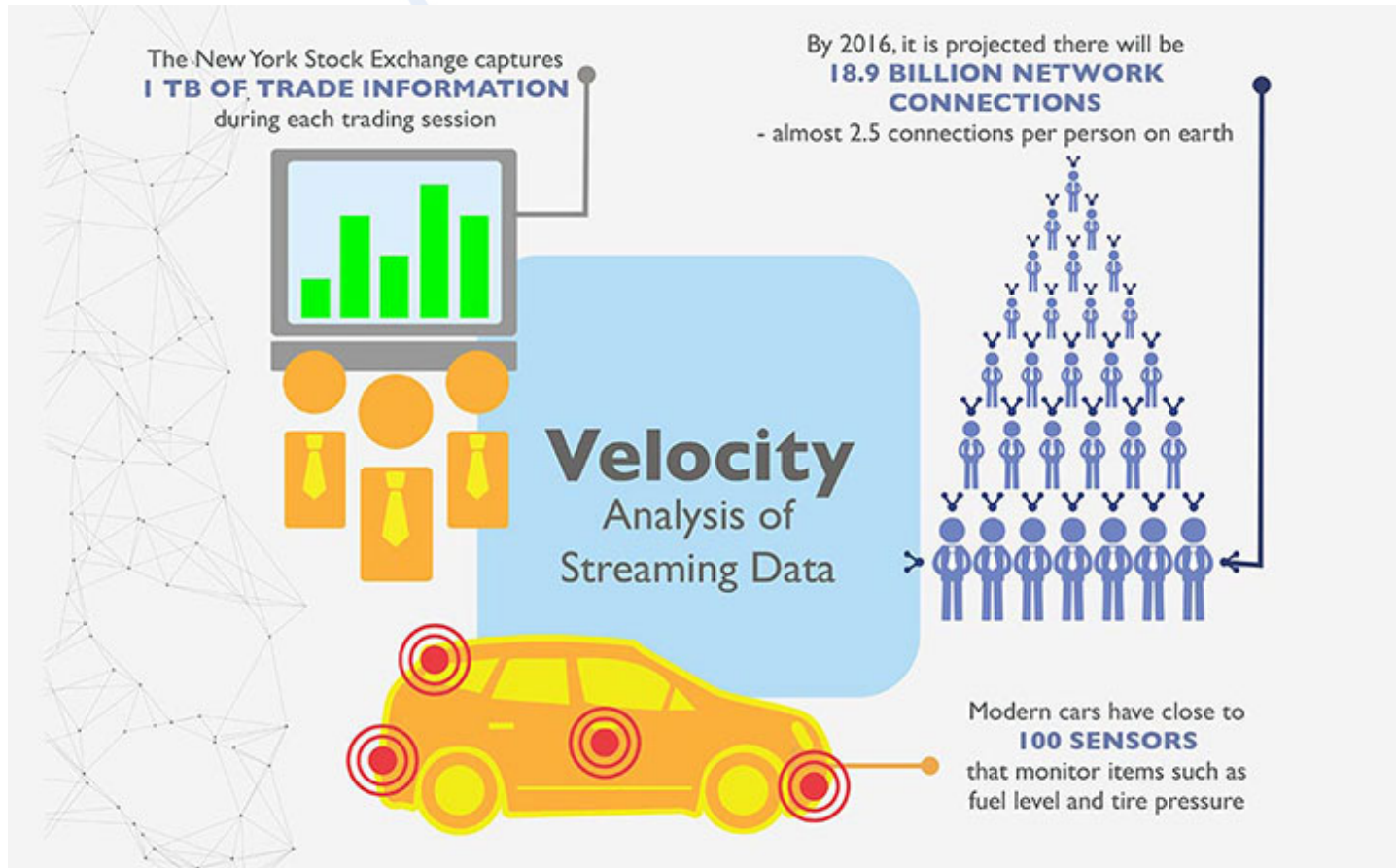# of commercial flights in the sky in the United States on any given day.

days in a year

= **2,499,841,200 TB**

# Velocity

- The speed of incoming data and the time it takes to process it.

- With the advent of IoT, streaming data is driving the need to process and analyze data in near-real time.

# Velocity

# Variety

- The type of files and format of data as well as sources.

- Data can be **structured,** such as a traditional database (pre-formatted data collected over time), or **unstructured** (unrelated data from unstructured sources such as social media, email, etc)

# Data Definition Framework

## Data Format

|  | Structured | Unstructured |
|---|---|---|
| **Internal** | **Human-Generated**<br>• Survey ratings<br>• Aptitude testing<br>**Machine-Generated**<br>• Web metrics from Web logs<br>• Product purchase from sales Records<br>• Process control measures | **Human-Generated**<br>• Emails, letters, text messages<br>• Audio transcripts<br>• Customer comments<br>• Voicemails<br>• Corporate video/communications<br>• Pictures, illustrations<br>• Employee reviews |
| **External** | **Human-Generated**<br>• Number of Retweets, Facebook likes, Google Plus +1s<br>• Ratings on Yelp<br>• Patient ratings<br>**Machine-Generated**<br>• GPS for tweets<br>• Time of tweet/updates/postings | **Human-Generated**<br>• Content of social media updates<br>• Comments in online forums<br>• Comments on Yelp<br>• Video reviews<br>• Pinterest images<br>• Surveillance video |

**Data Source**

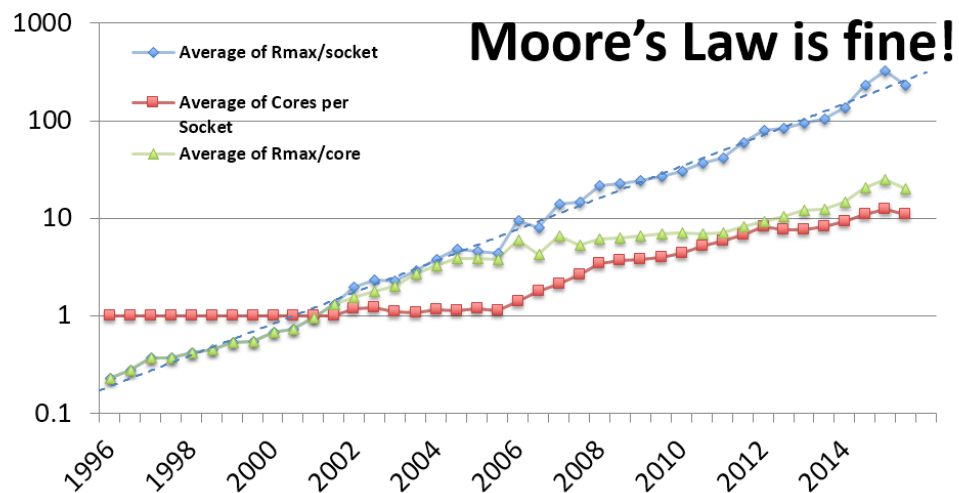# Big Data – Supporting Trends

- **Moore's Law**: an observation that the number of transistors on integrated circuits doubles every two years.
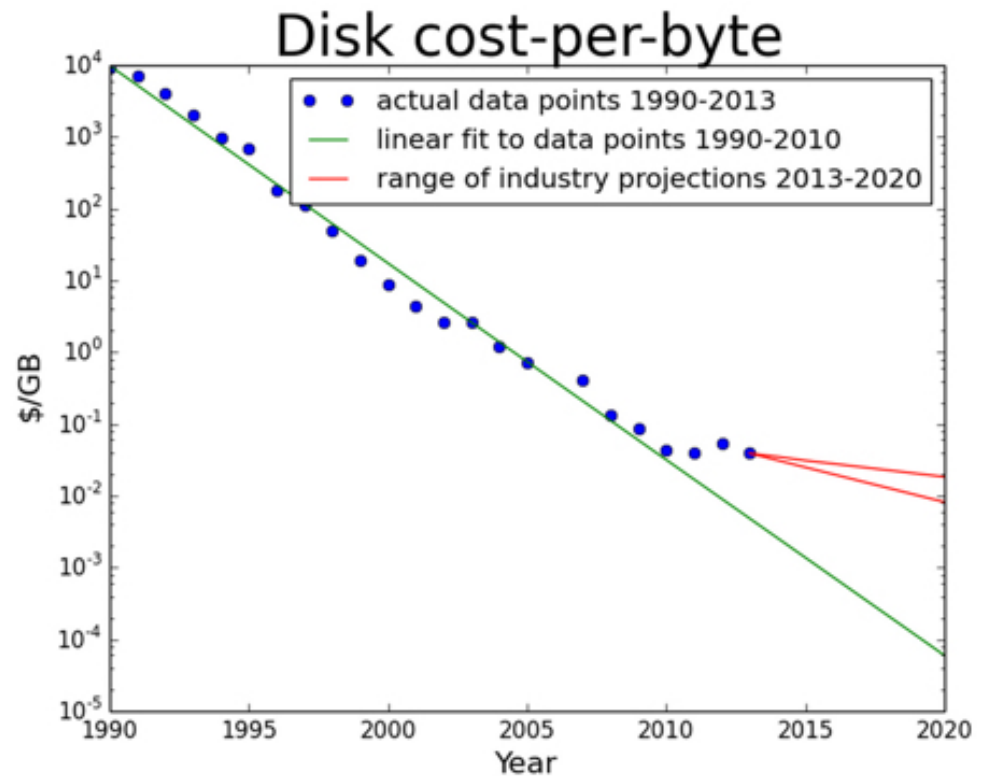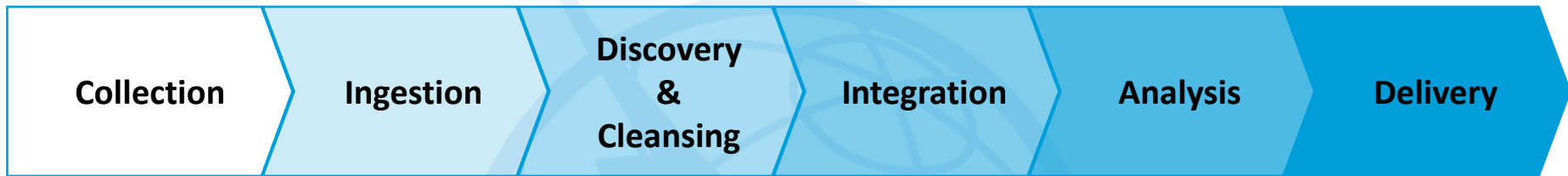
# Big Data – Supporting Trends

- **Kryder's Law:** the density of storage is increasing and the cost decreasing at a rate faster than Moore's Law



Disk cost-per-byte

- actual data points 1990-2013
- linear fit to data points 1990-2010
- range of industry projections 2013-2020

# Big Data Value Chain

| Collection | Ingestion | Discovery & Cleansing | Integration | Analysis | Delivery |
|---|---|---|---|---|---|

- **Collection** – Structured, unstructured and semi-structured data from multiple sources

- **Ingestion** – loading vast amounts of data onto a single data store

- **Discovery & Cleansing** – understanding format and content; clean up and formatting

- **Integration** – linking, entity extraction, entity resolution, indexing and data fusion

- **Analysis** – Intelligence, statistics, predictive and text analytics, machine learning

- **Delivery** – querying, visualization, real time delivery on enterprise-class availability

Source O'Reilly Strata 2012

# Big Data – Tools

- **Hadoop** is often used at the server level to organise the cluster along with a NoSQL database for data storage.

- **NoSQL** are databases that use looser consistency models than relational databases. Performance gains via simplification using key value stores.

# Examples of data generated by IoT

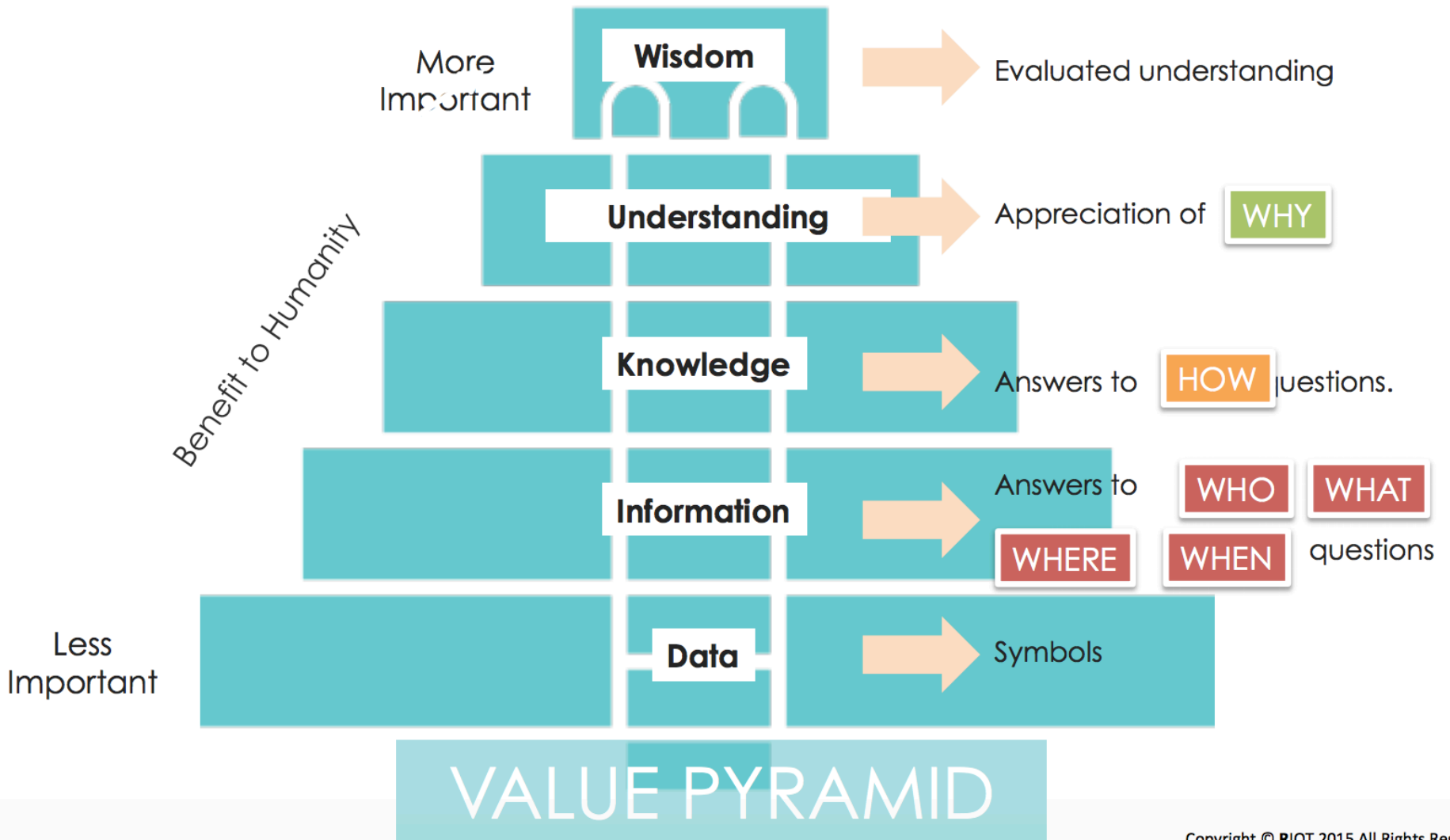| | Individual | Community | Society |
|---|---|---|---|
| **Level** | | | |
| **IoT** | Smart phones<br>Wearables | Connected Cars<br>Health devices<br>Smart homes | Smart Cities<br>Smart Grids |
| **Examples** | GPS, Fitbits<br>Visa PayWave<br>Mastercard Paypass<br>Employee passes | Intelligent Transport Systems<br>Event Data Recorders (EDRs)<br>Blood pressure monitors;<br>remote burglar/heating<br>systems | Smart metering;<br>Smart water meters<br>Traffic monitoring |
| **Data** | Mobile money<br>Fitness data, GPS<br>location-based data | Speed, distance, airbag,<br>crash locations/alerts;<br>Heart rate, blood pressure,<br>Diet, remote heating data | Electricity/water<br>consumption & billing;<br>Traffic flow data |
| **Intended Audience** | Individual person<br>Immediate friends/ family;<br>banks; employers | GP, health authorities;<br>health & car insurance;<br>police, social networks | Authorities/regulators<br>Utility companies;<br>Other citizens |

# Usecase: Traffic

- **Collect** traffic data and transportation data from sensors

- **Build** a model of traffic patterns

- **Predict** the traffic and congestions

- **Act**: divert traffic, adjust troll, adjust traffic lights
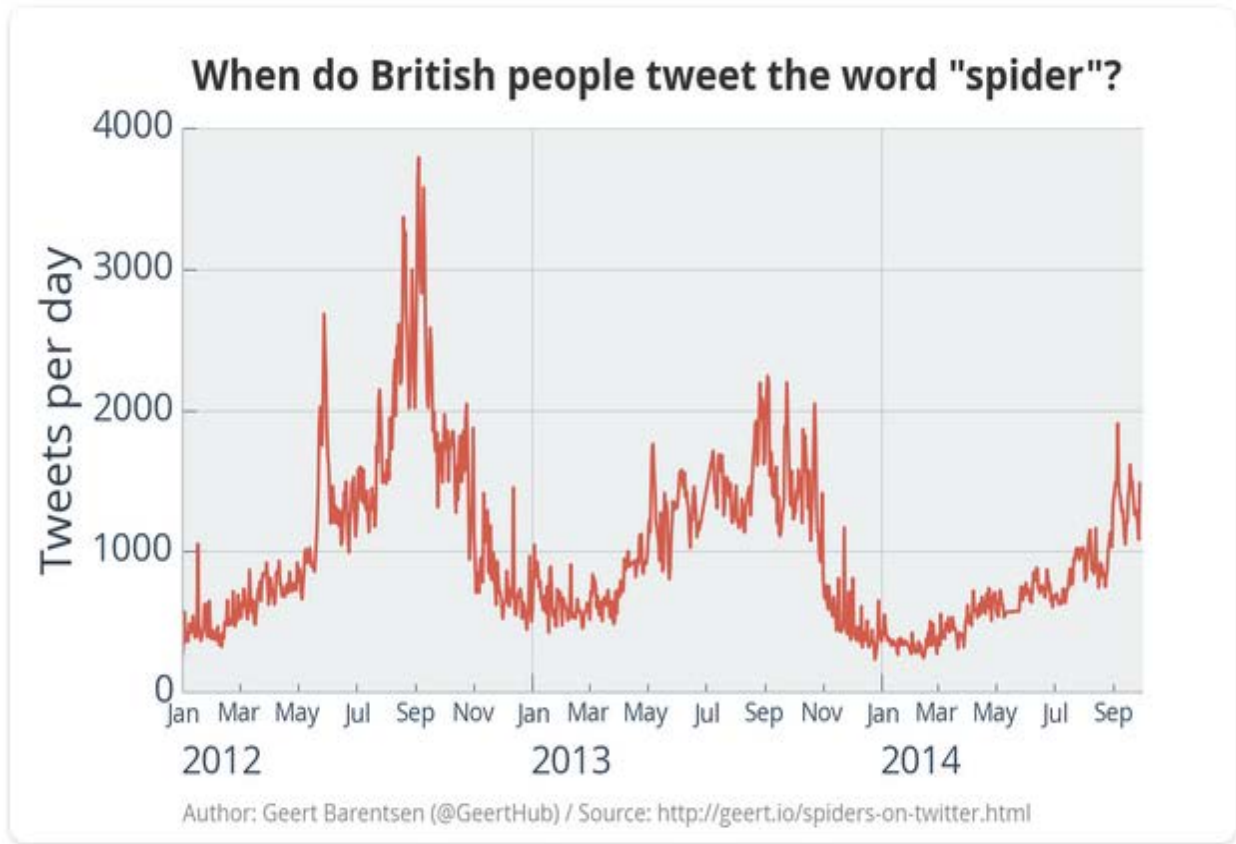
Source: http://inrix.com/

More Important

Less Important

Benefit to Humanity

Wisdom → Evaluated understanding

Understanding → Appreciation of **WHY**

Knowledge → Answers to **HOW** questions.

Information → Answers to **WHO** **WHAT** **WHERE** **WHEN** questions

Data → Symbols

VALUE PYRAMID

Source: REDtone IOT

**When do British people tweet the word "spider"?**

Number of tweets per day in Britain that contained the word "spider". Retweets, replies, and tweets about Spider-Man have been excluded.

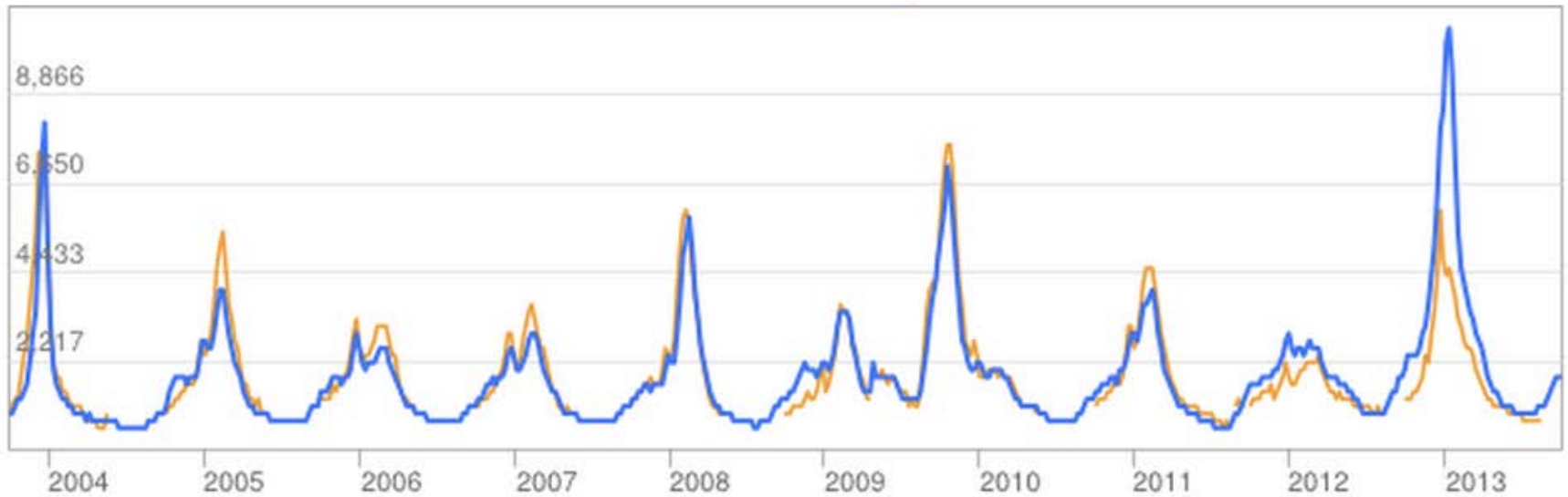# The number of tweets about spiders in Britain is correlated with the local mean temperature:



Author: Geert Barentsen (@GeertHub) / Source: http://geert.io/spiders-on-twitter.html

*Number of British spider tweets per day (top), shown against the mean temperature in Central England obtained from the UK Met Office (bottom).*

# United States Flu Activity

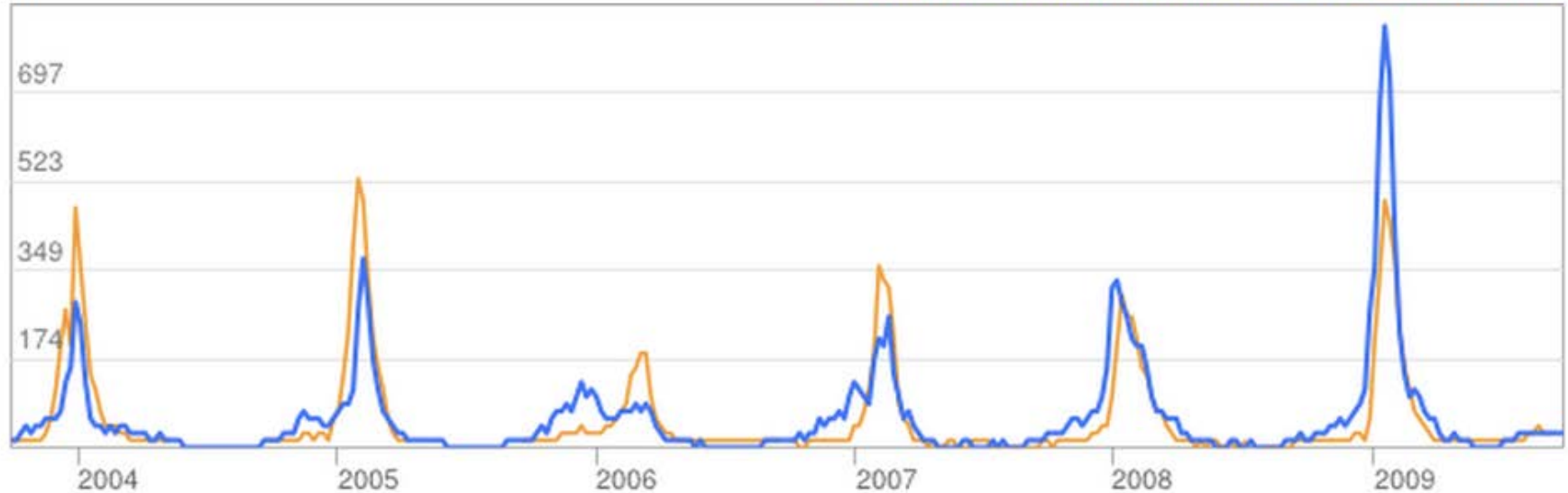Influenza estimate    ● Google Flu Trends estimate ● United States data

United States: Influenza-like illness (ILI) data provided publicly by the U.S. Centers for Disease Control.
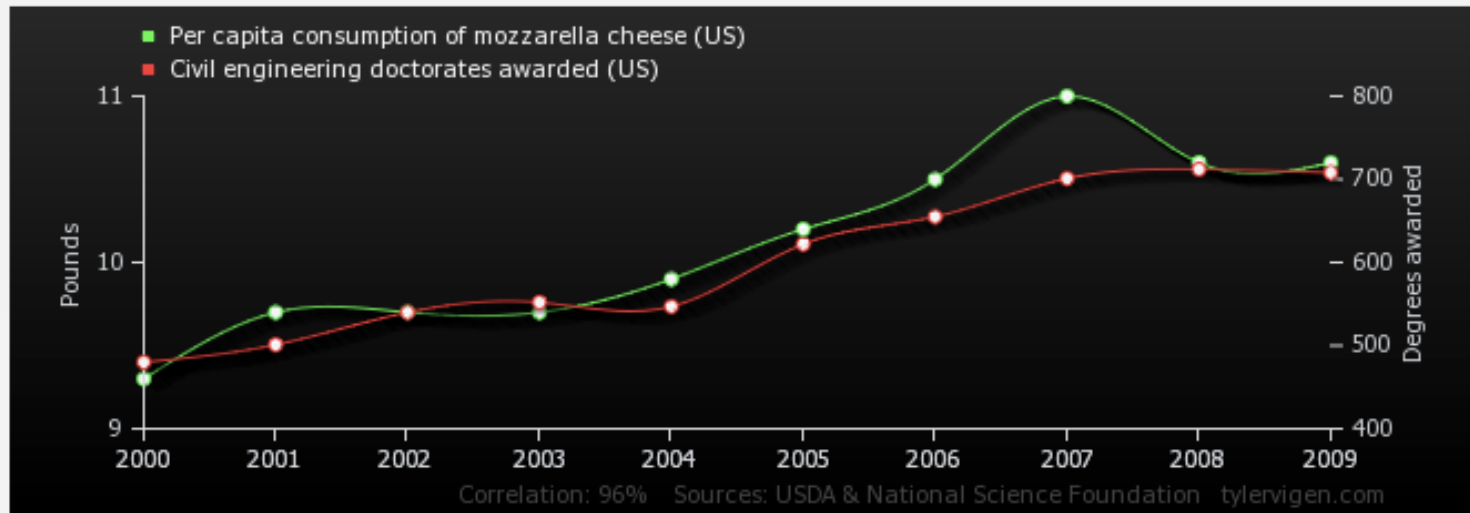
# Switzerland Flu Activity

Influenza estimate                           ● Google Flu Trends estimate  ● Switzerland data



Switzerland: Influenza-like illness (ILI) data provided publicly by the European Influenza Surveillance Network of the European Centre for Disease Prevention and Control.

# Per capita consumption of mozzarella cheese (US)
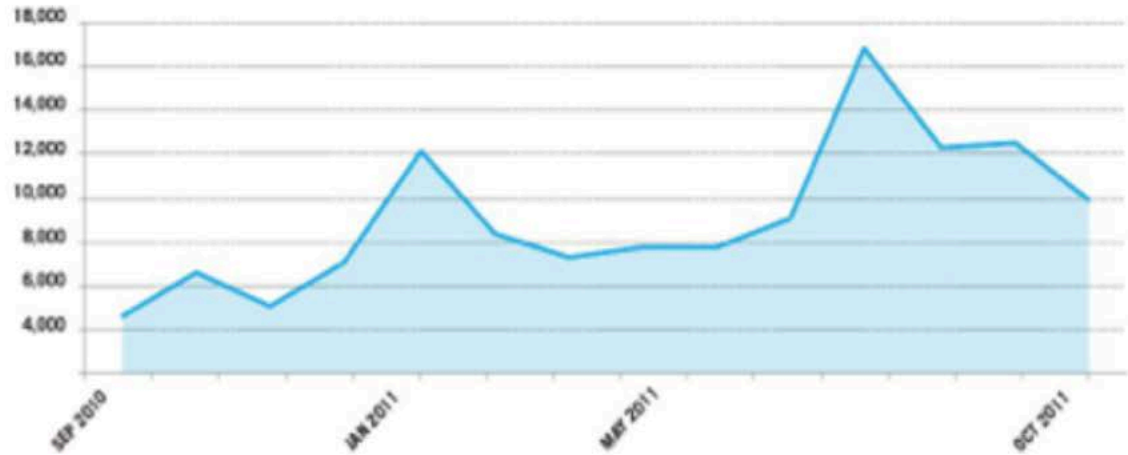## correlates with
# Civil engineering doctorates awarded (US)



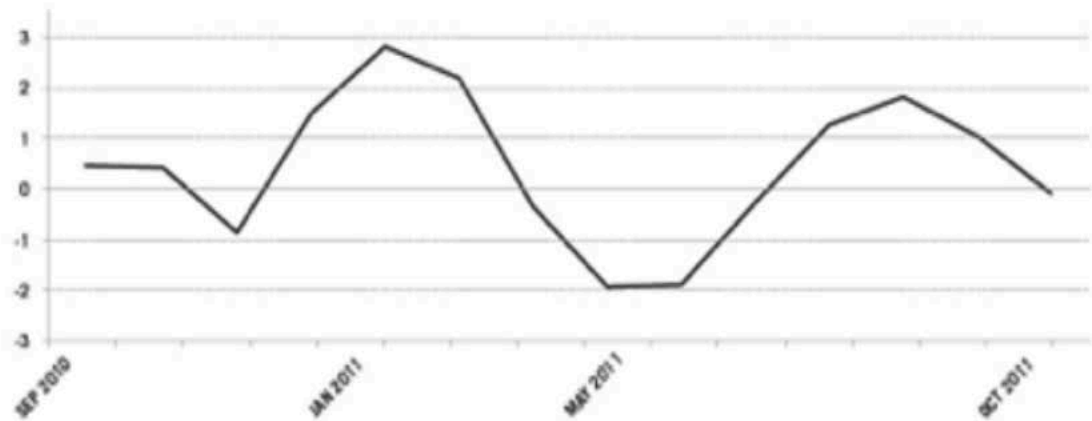| | 2000 | 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 |
|---|---|---|---|---|---|---|---|---|---|---|
| Per capita consumption of mozzarella cheese (US) Pounds (USDA) | 9.3 | 9.7 | 9.7 | 9.7 | 9.9 | 10.2 | 10.5 | 11 | 10.6 | 10.6 |
| Civil engineering doctorates awarded (US) Degrees awarded (National Science Foundation) | 480 | 501 | 540 | 552 | 547 | 622 | 655 | 701 | 712 | 708 |

**Correlation: 0.958648**

Source: http://tylervigen.com/
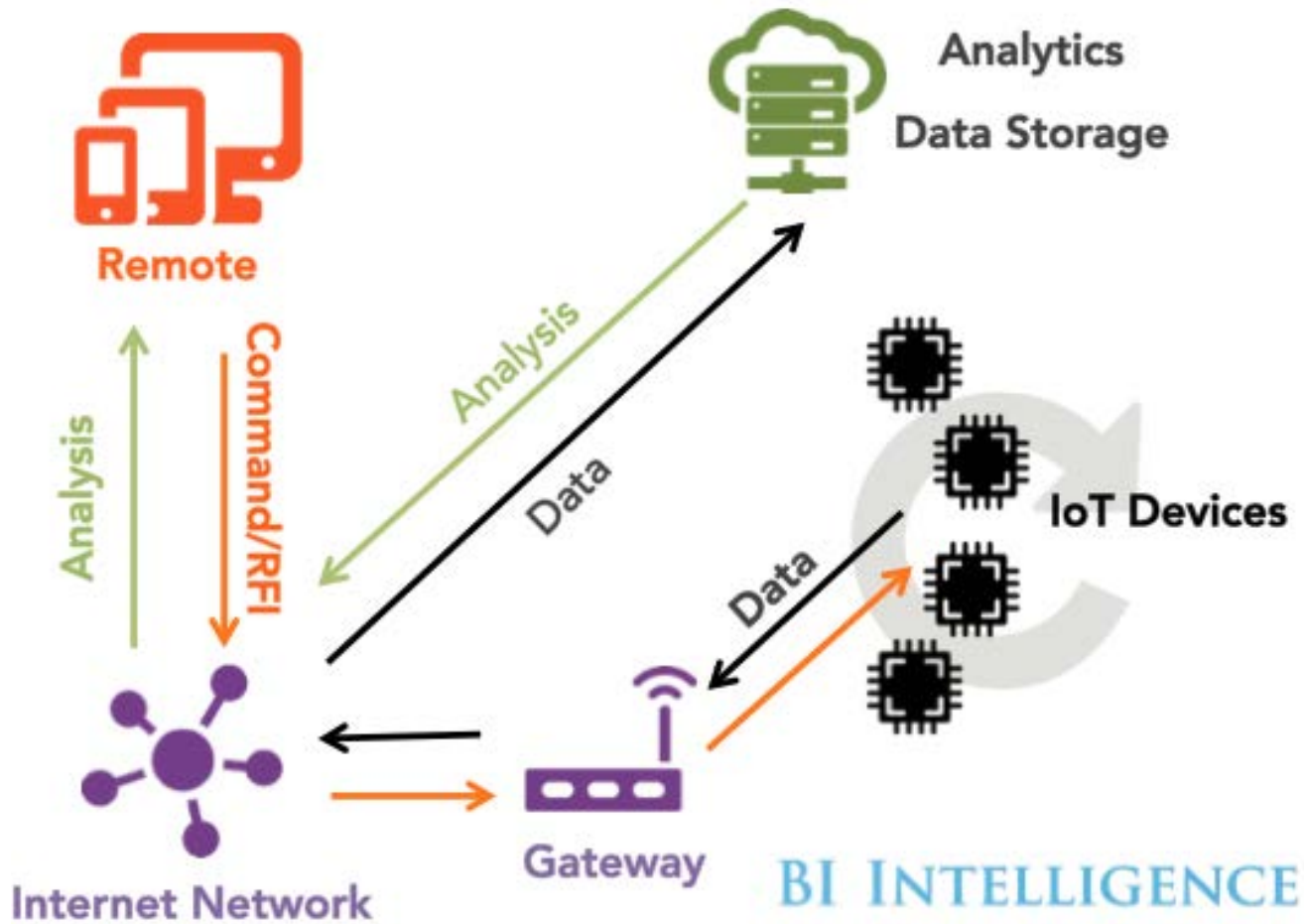
Tweets about the price of rice (per month)

Food Price Inflation

Source: UN Global Pulse

# The Internet of Things Ecosystem



Remote

Analytics
Data Storage

Analysis

Command/RFI

Analysis

Data

IoT Devices

Data

Internet Network

Gateway

BI INTELLIGENCE

# Demo

http://discover-iot.eu-gb.mybluemix.net/#/play/device/smartphone

# Thank You