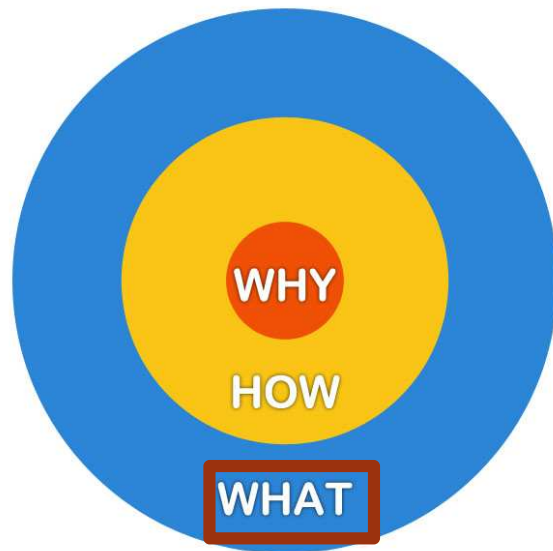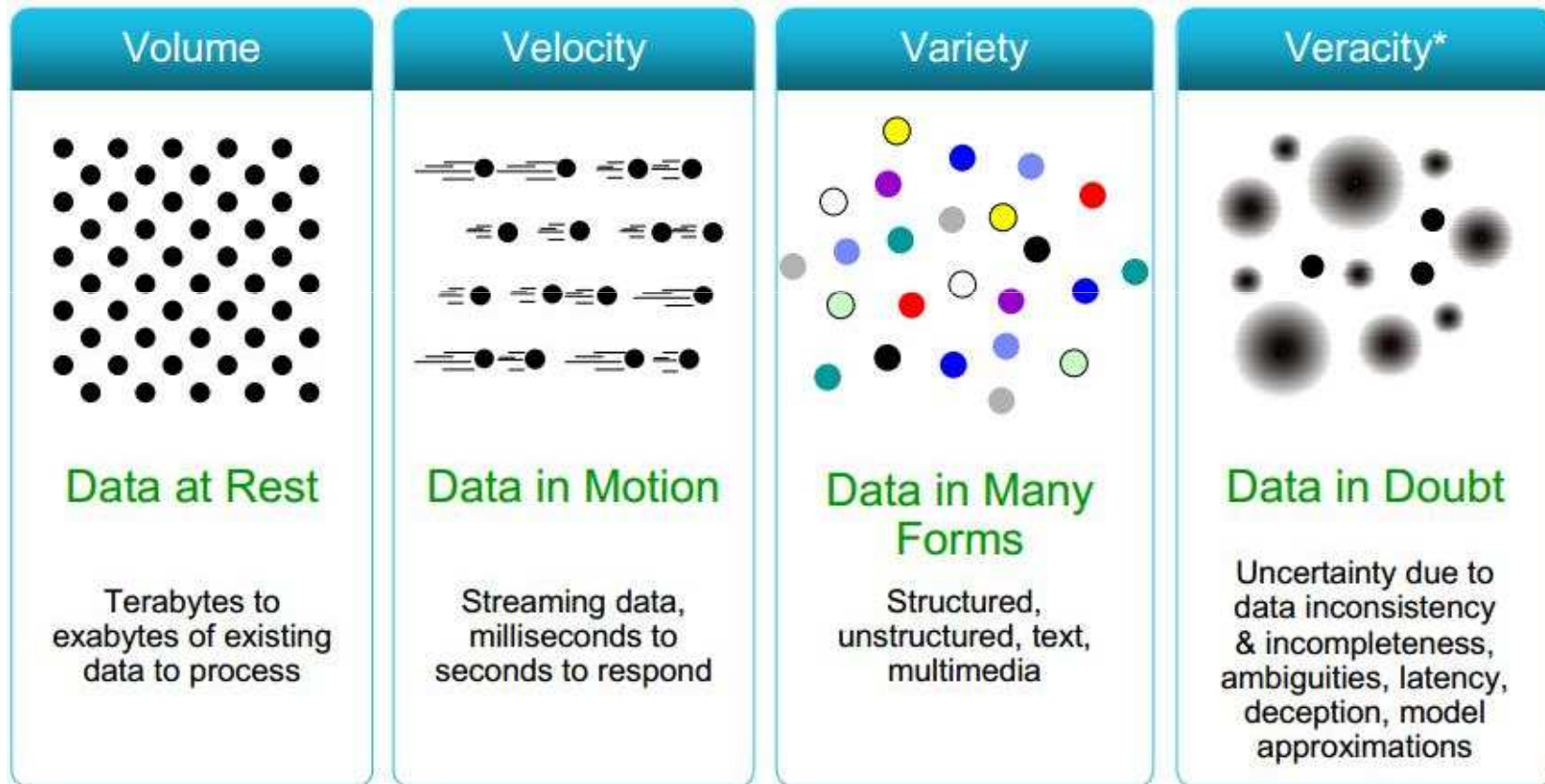# Introduction to Big Data
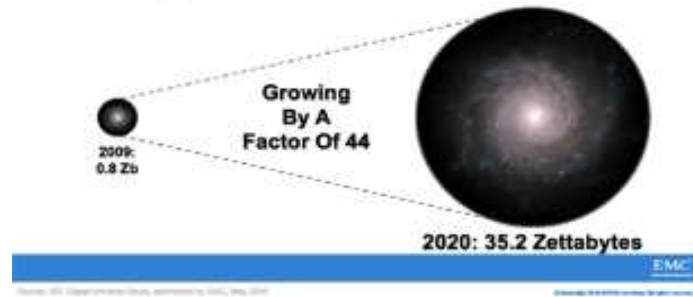


Rosaline Macharia

# What is Big Data?



"**Big Data**" is data whose **scale**, **diversity**, and **complexity** require new architecture, techniques, algorithms, and analytics to manage it and extract **value** and **hidden knowledge** from it
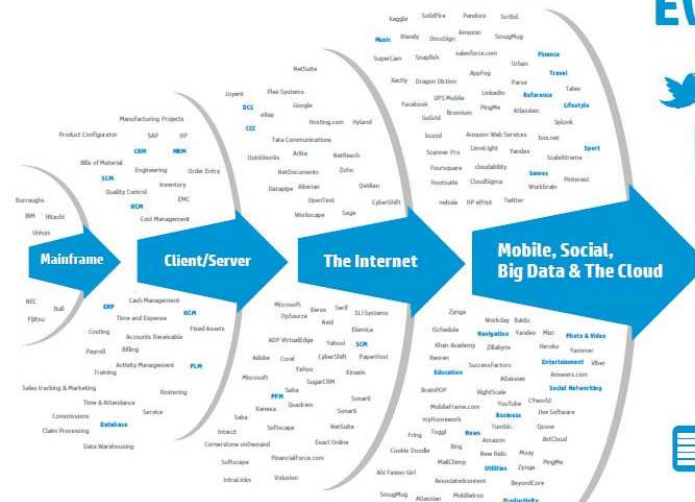
# Characteristics of Big Data



| Volume | Velocity | Variety | Veracity* |
|--------|----------|---------|-----------|
| **Data at Rest** | **Data in Motion** | **Data in Many Forms** | **Data in Doubt** |
| Terabytes to exabytes of existing data to process | Streaming data, milliseconds to seconds to respond | Structured, unstructured, text, multimedia | Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations |

# 1-Volume



The Digital Universe 2009-2020

Growing By A Factor Of 44

2009: 0.8 Zb

2020: 35.2 Zettabytes

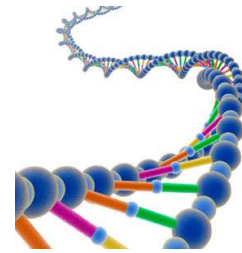44x increase from 2009 - 2020

A new style of IT emerging

Every 60 seconds

98,000+ tweets

695,000 status updates

11 million instant messages

698,445 Google searches

168 million+ emails sent

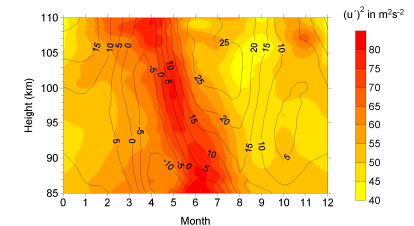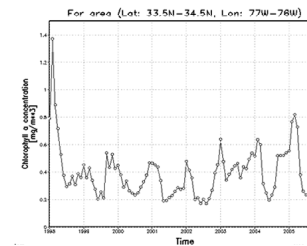1,820TB of data created

217 new mobile web users

# 2- Variety

- Various formats, types, and structures

  Text, numerical, images, audio, video, **sequences,** time series, social media data, multi-dim arrays, etc…
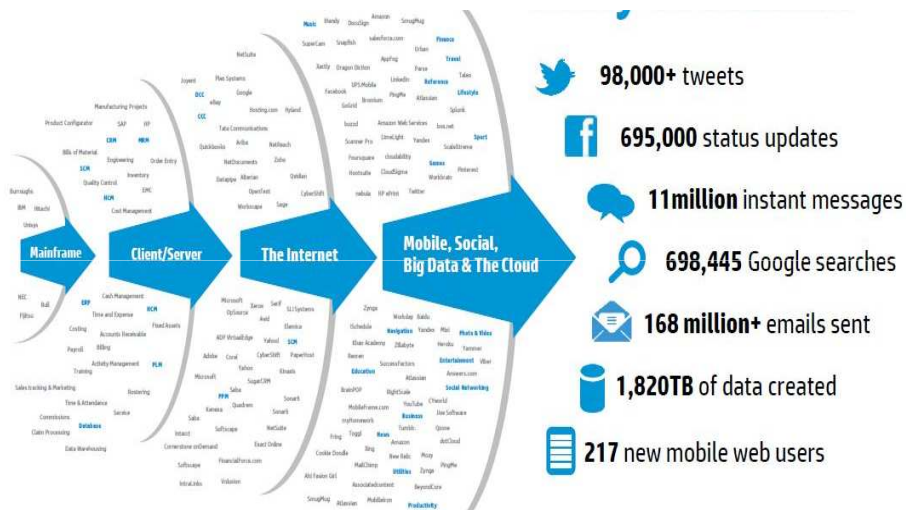
- Static data vs. streaming data



```
@SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
CAACGAGTTCACACCTTGGCCGACAGGCCCGGGTAA
+SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
BA@7>B=>:>>7@7@>>9=BAA?;>52;>:9=8.=A
@SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
CCAATGATTTTTTCCGTGTTTCAGAATACGGTTAA
+SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
BCCBA@BB@BBBBBAB@B9B@=BABA@A:@693:@B=
@SRR038845.53 HWI-EAS038:6:1:1:360 length=36
GTTCAAAAAGAACTAAATTGTGTCAATAGAAAACTC
+SRR038845.53 HWI-EAS038:6:1:1:360 length=36
BBCBBBBBB@@BAB?BBBBCBC>BBBAA8>BBBAA@
```
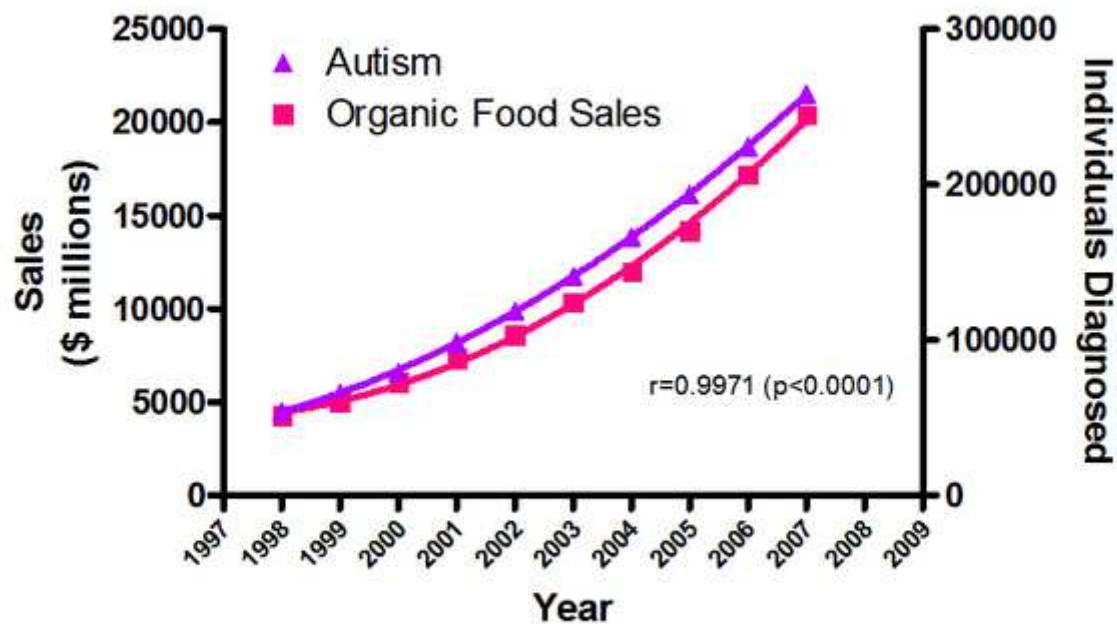
# 3-Velocity



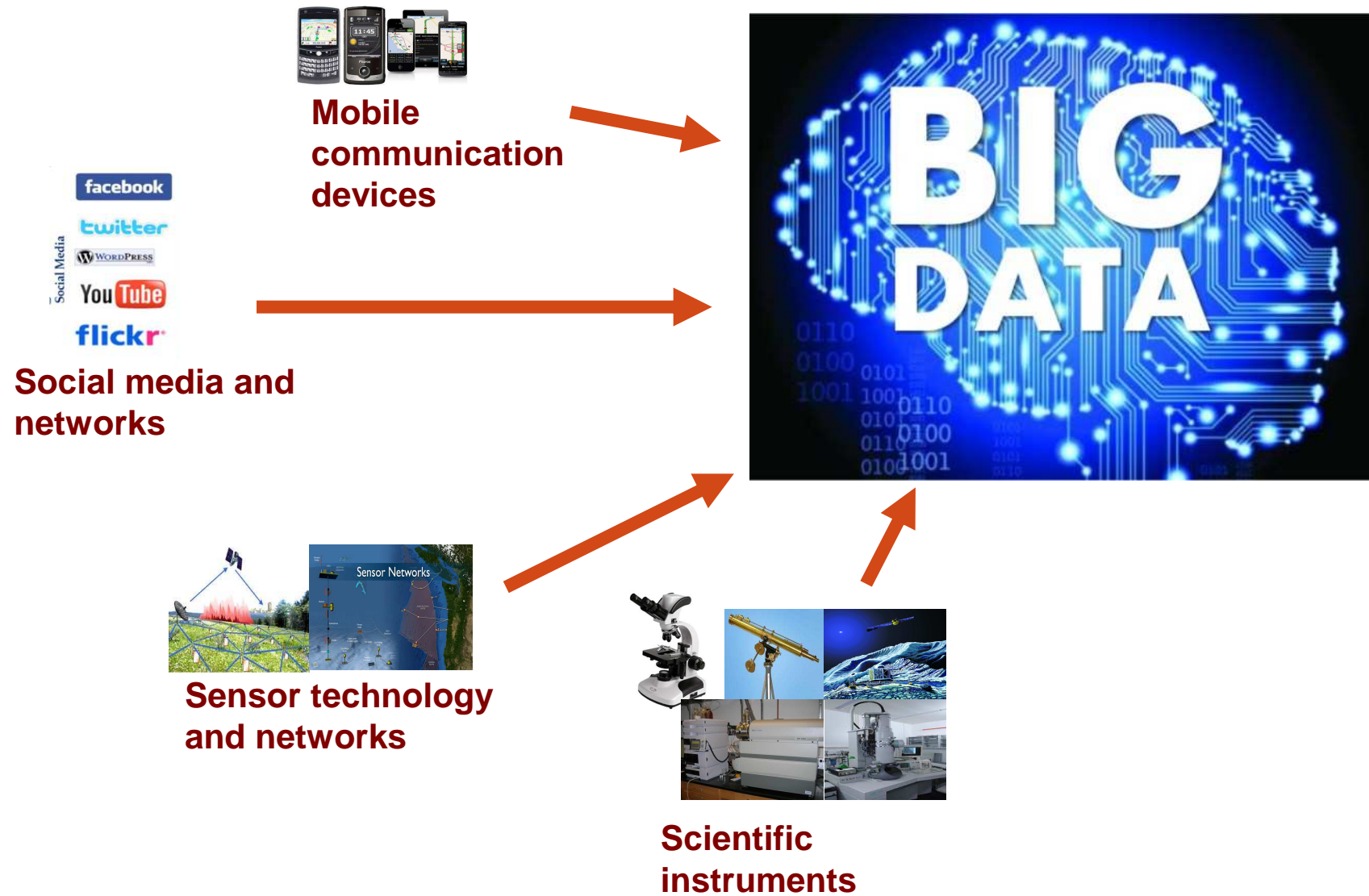Data is being generated fast and needs to be processed fast

# 4- Veracity

Does organic food consumption contribute to autism?



Sources: Organic Trade Association, 2011 Organic Industry Survey; U.S. Department of Education, Office of Special Education Programs, Data Analysis System (DANS), OMB# 1820-0043: "Children with Disabilities Receiving Special Education Under Part B of the Individuals with Disabilities Education Act

# Where is it from?

**Mobile communication devices**

**Social media and networks**

**Sensor technology and networks**

**Scientific instruments**
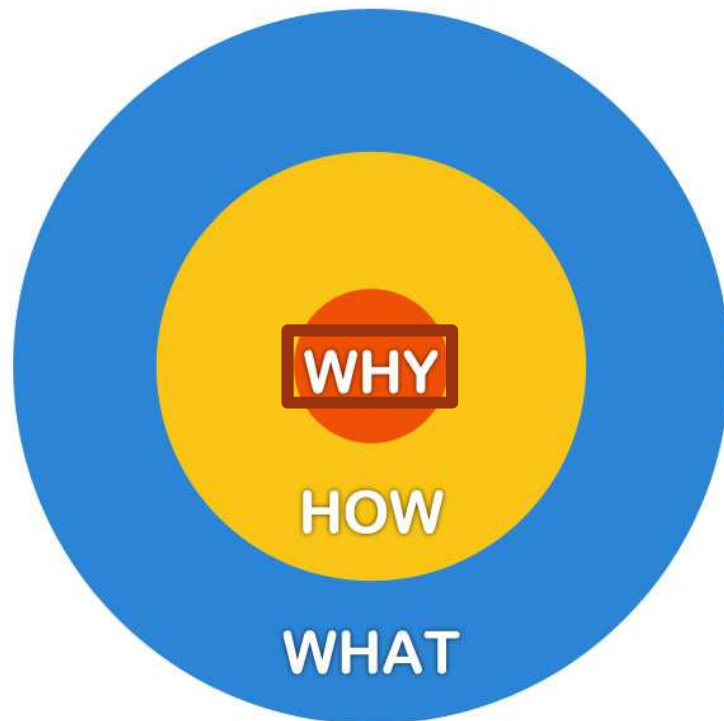
BIG DATA

# Data Generation & Consumption

**Old Model:** Few companies were generating data, while all others consume it



**New Model:** All of us are generating data, and all of us are consuming data

# Why Big data



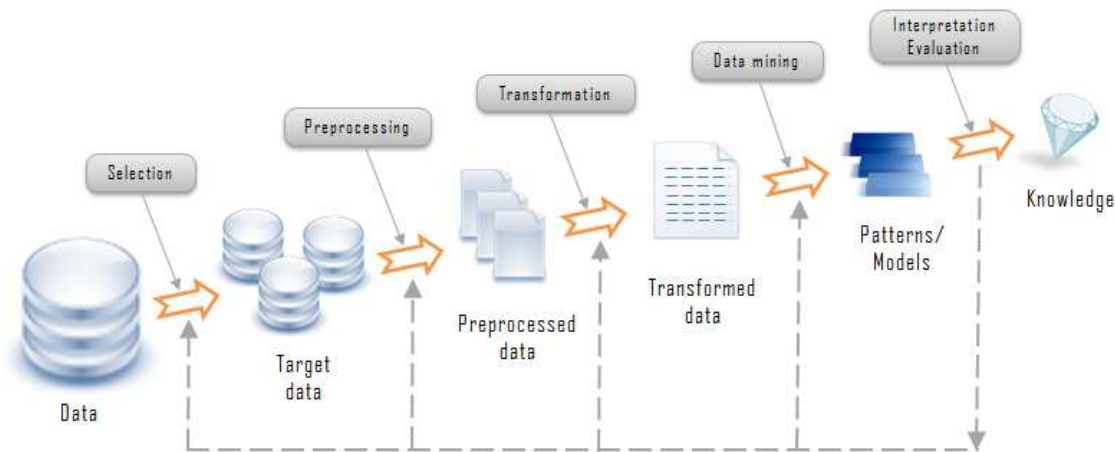Progress and innovation is no longer hindered by the ability to collect data…

…but, by the ability to manage, analyze, summarize, visualize, and discover knowledge from the collected data in a timely manner and in a scalable fashion
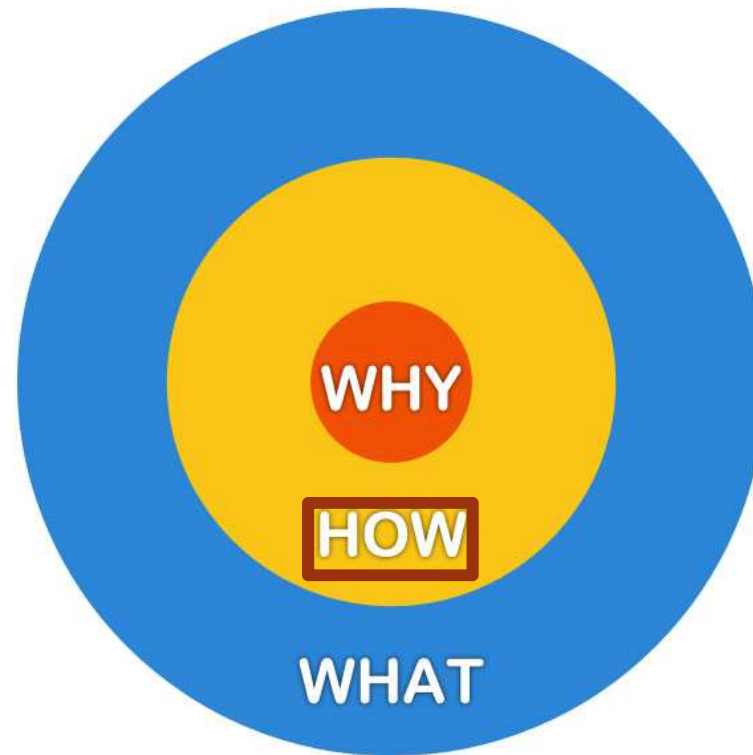
# Data to Knowledge

What drives big data?

**Technology –** "Big data: Only for the privileged"

**Value** **–** "Data is the new oil/ currency"

# Big data in Biology



How is it generated?

# A historical perspective of Bioinformatics data

- The 1960s: the birth of bioinformatics

  - computer languages

  - Academic access to computers

  - First protein database



**IBM 7090 computer**

# Data Processing Developments

**1960**

**2015**

**IBM 7090 computer**

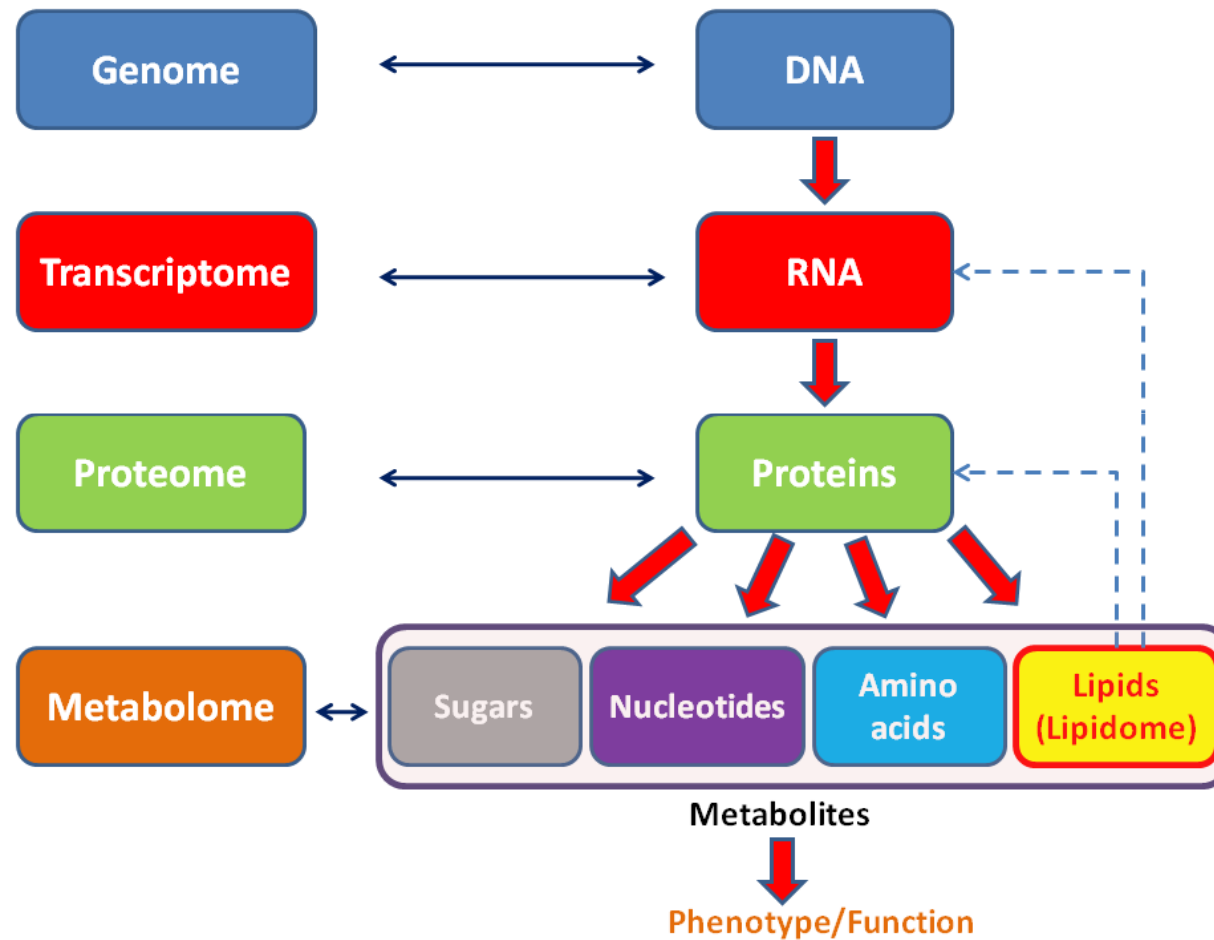32 Kbytes RAM

2.18 µHz

$2,900,000 in 1960
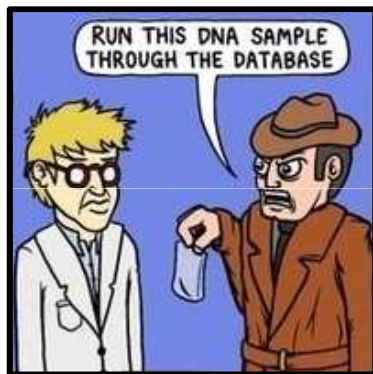
**21.5" Apple iMac**

8 GB RAM

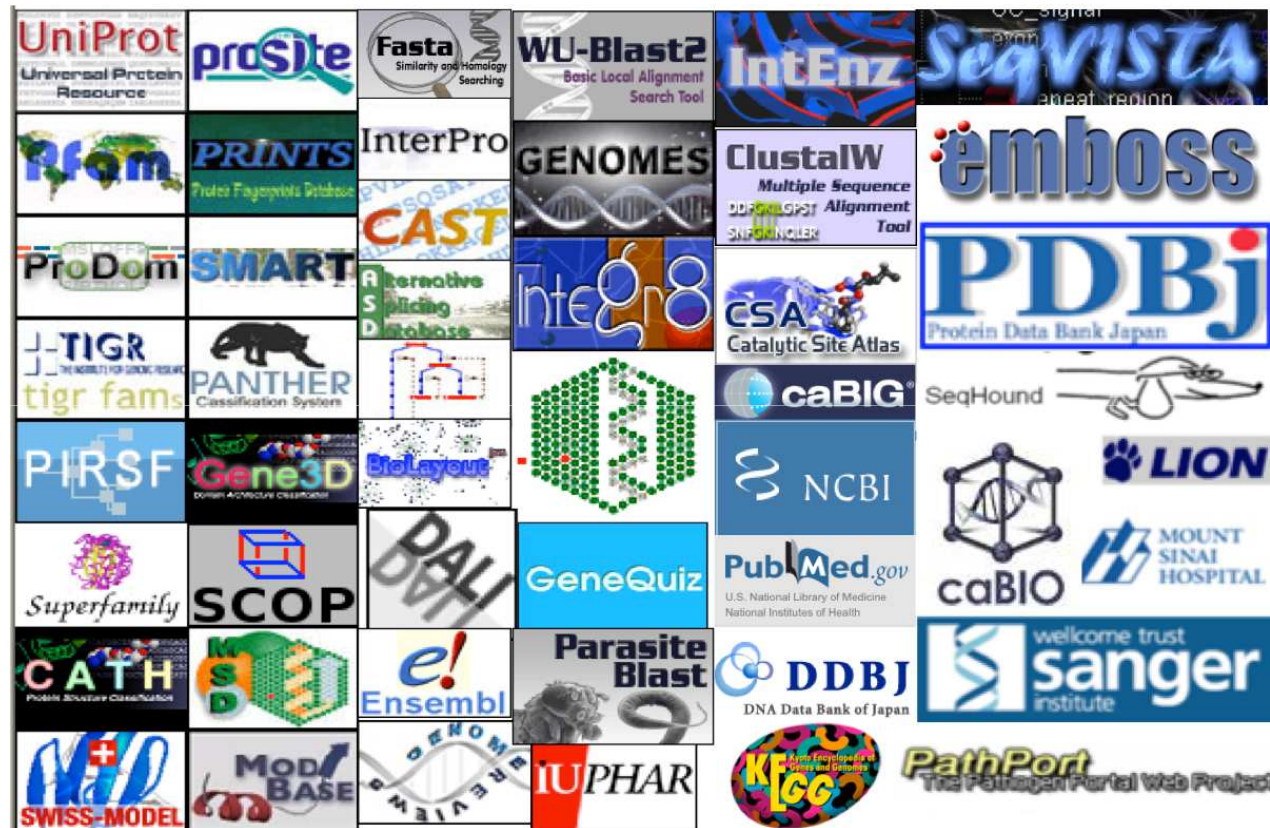2.7 GHz

$1,149 in 2015

# BIG "Omics" DATA

# Generation of OMICS DATA



**NGS technologies**

# Biological Databases & Tools

# Challenges in Handling Big Data

- **The Bottleneck in technology**

  New architecture, algorithms, techniques are needed
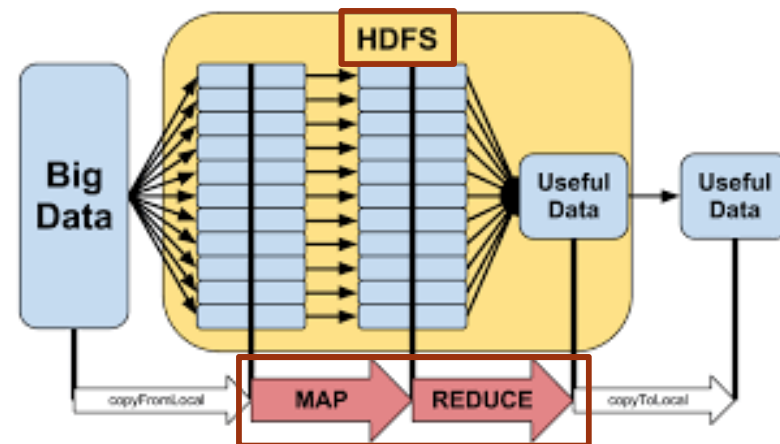
- **Technical skills**

  Experts in using the new technology and dealing with big data

# Big Data Analytics

**Challenges with existing DBMS tools**

- Fixed Schema
- Cost
- Saving/Accessing huge files
- Performing Analysis

Examining data to uncover hidden patterns

# Cloud computing

"computing in which dynamically scalable and virtualized resources are provided as a service over the internet"

| Pros | Cons |
|---|---|
| Scale | Security |
| Cost is relative on scale | Lack of control |
| Choice | Reliability |
| Access to NGS architecture | |

Bioinformatics Approaches for NGS Analysis Course