

# Introduction

## **BCB 5200 Introduction to Bioinformatics I**

Fall 2017

**Tae-Hyuk (Ted) Ahn**

Department of Computer Science  
Saint Louis University



**SAINT LOUIS  
UNIVERSITY™**

— EST. 1818 —

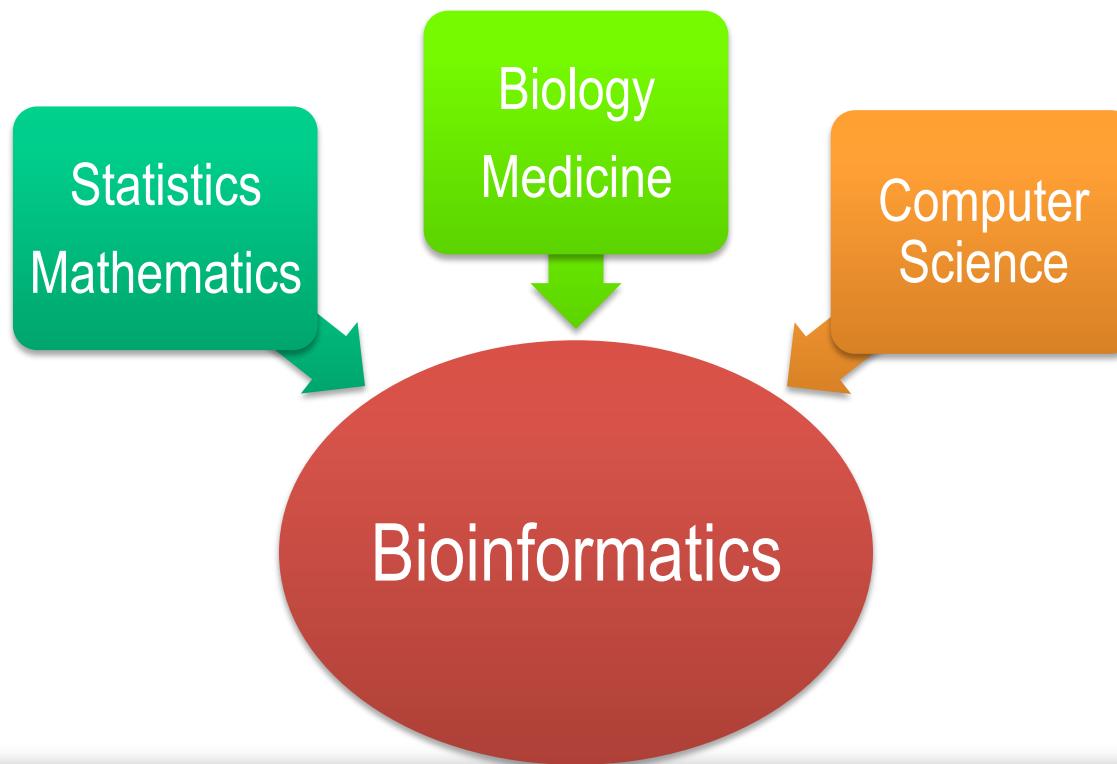
# Intro

- Please read the syllabus carefully!
- We will use Blackboard for posting lectures/assignments and grading.
- Please keep practicing Linux and editor (VI, emacs, nano)
  - Contact us if you have any problem!!!!

# Bioinformatics?

## What is Bioinformatics?

Application of techniques from computer science to problems from biology



# Bioinformatics vs Computational Biology

## Bioinformatics

An interdisciplinary scientific field that develops methods and software tools for storing, retrieving, organizing and analyzing biological data.

## Computational Biology

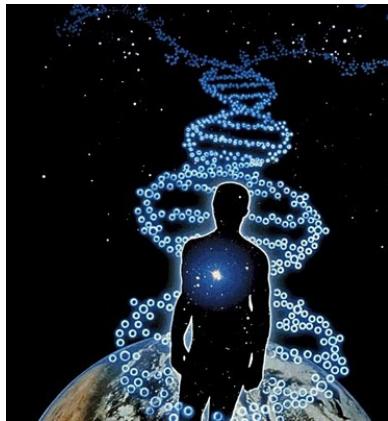
The development of theoretical methods, mathematical modeling, and computational simulation techniques to study of biological systems.

# Computational Goals of Bioinformatics

- Learn & Generalize: Discover conserved patterns (models) of sequences, structures, interactions, metabolism & chemistries from well-studied examples.
- Prediction: Infer function or structure of newly sequenced genes, genomes, proteins or proteomes from these generalizations.
- Organize & Integrate: Develop a systematic and genomic approach to molecular interactions, metabolism, cell signaling, gene expression...
- Simulate: Model gene expression, gene regulation, protein folding, protein-protein interaction, protein-ligand binding, catalytic function, metabolism...
- Engineer: Construct novel organisms or novel functions or novel regulation of genes and proteins.
- Gene Therapy: Target specific genes, or mutations, RNAi to change a disease phenotype.

[http://cmgm.stanford.edu/  
biochem218/01Genomics&Bioinformatics.pdf](http://cmgm.stanford.edu/biochem218/01Genomics&Bioinformatics.pdf)

# Myriad Applications of Bioinformatics



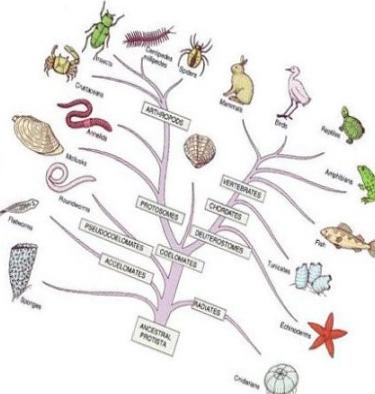
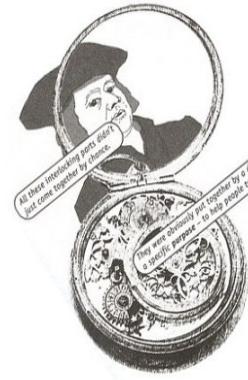
Human



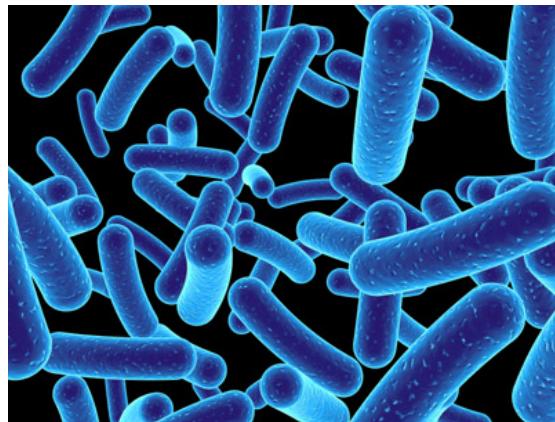
Agriculture



Livestock



Ancestry

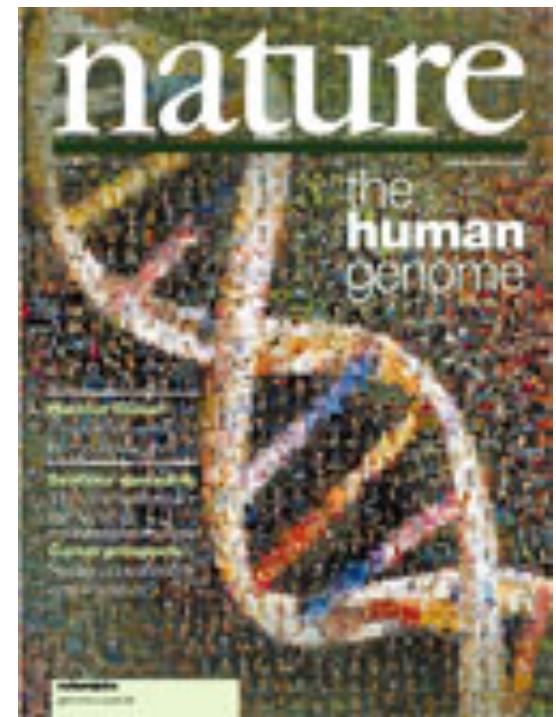
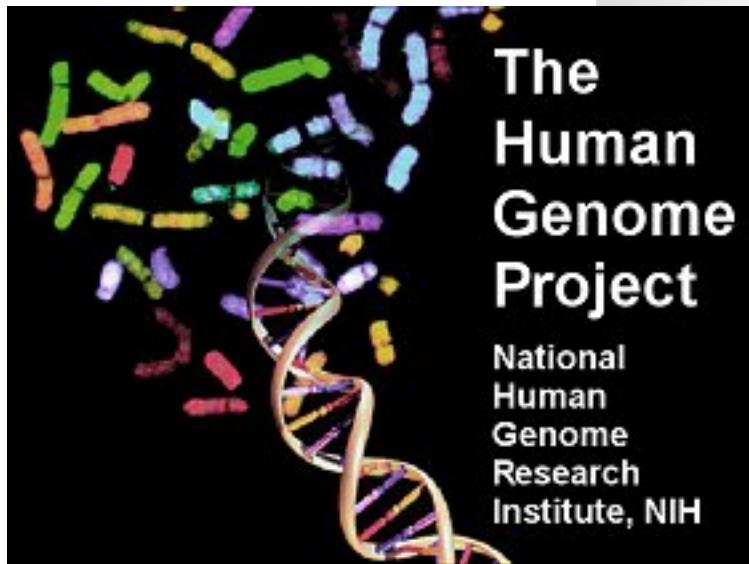


Microbes



Bioenergy

# Human Genome Project (1990 – 2003)



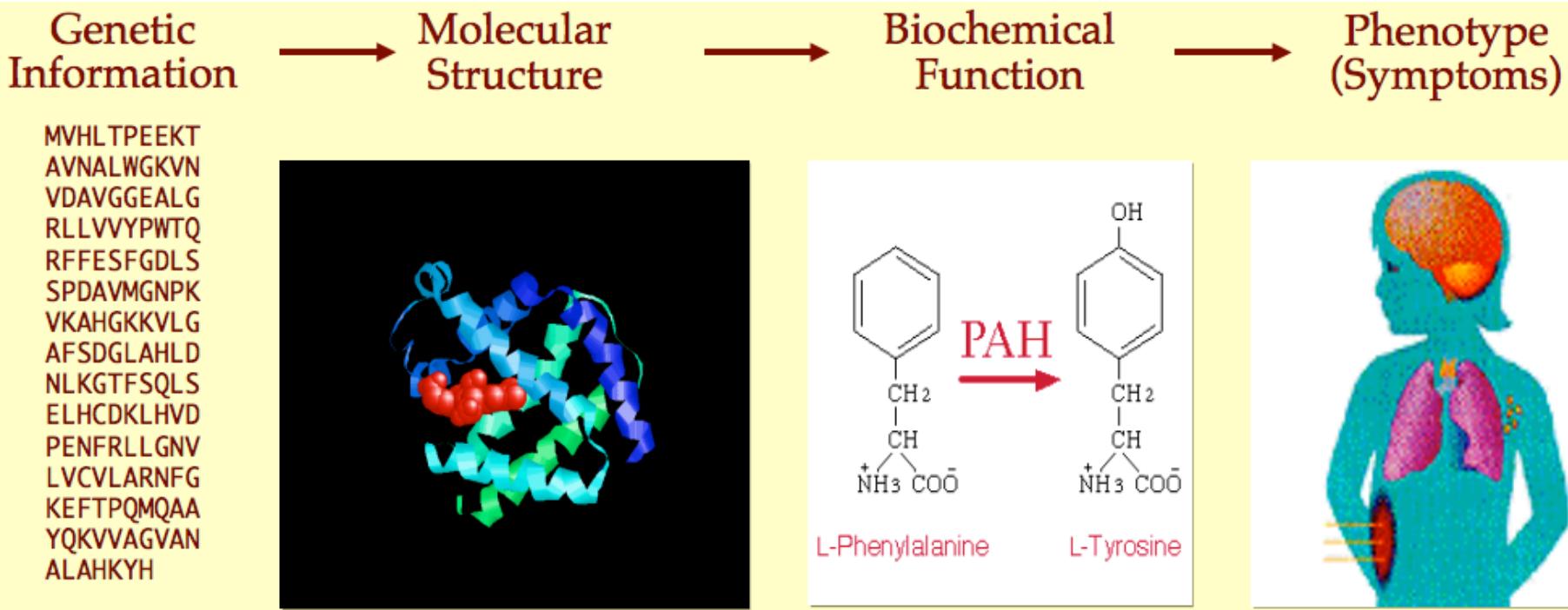
# How to Sequence the Human Genome (Mac)



# Health, Diseases, and Medicine



# Central Paradigm of Bioinformatics



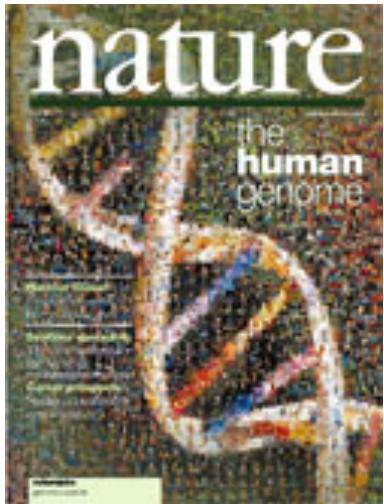
<http://cmgm.stanford.edu/biochem218/01Genomics&Bioinformatics.pdf>

# Genomic Medicine

- An emerging medical discipline that involves using genomic information about an individual as part of their clinical care (e.g., for diagnostic or therapeutic decision-making) and the other implications of that clinical use.



# The Path to Genomic Medicine



## Human Genome Project

## Personal Genome

# NIH (National Institutes of Health)

- The National Institutes of Health (NIH) is a biomedical research facility primarily located in Bethesda, Maryland.
- An agency of the United States Department of Health and Human Services, it is the primary agency of the United States government responsible for biomedical and health-related research.



- Yearly Budget: \$ 30 Billion
- More than 5000 researchers are working

## BUDGET AND APPROPRIATIONS

[Building on Opportunities in Cancer Research](#)

[The NCI Director's Message](#)

[The Changing Cancer Landscape](#) +

[Building on the National Cancer Program](#)

[New Approaches to Funding Researchers](#)

[NCI-Designated Cancer Centers](#)

[NCI's National Clinical Trials Enterprise](#)

[Overcoming Cancer Health Disparities](#)

[NCI's Intramural Research Program](#)

**[Bioinformatics to Accelerate Research](#)**

[Frederick National Laboratory for Cancer Research](#)

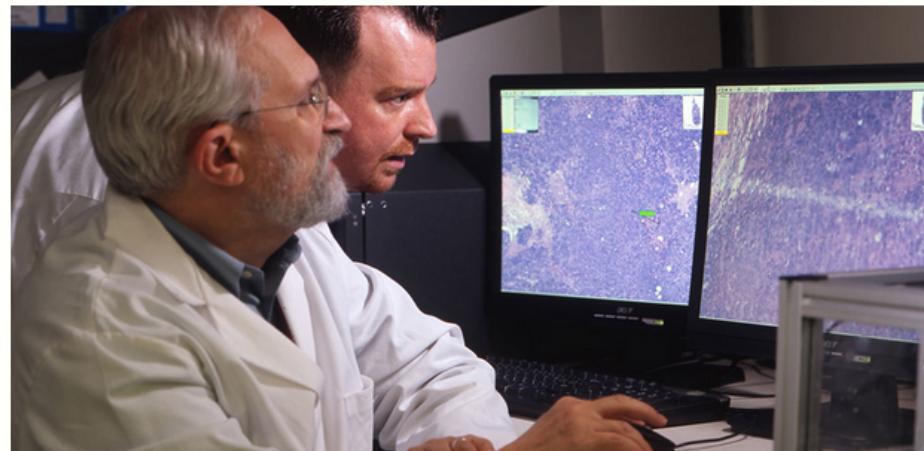
[Opportunities in Cancer](#)

## Bioinformatics to Accelerate Research

Bioinformatics, which enables the management and use of very large sets of molecular and clinical data, has become a core component of the NCI's research enterprise.

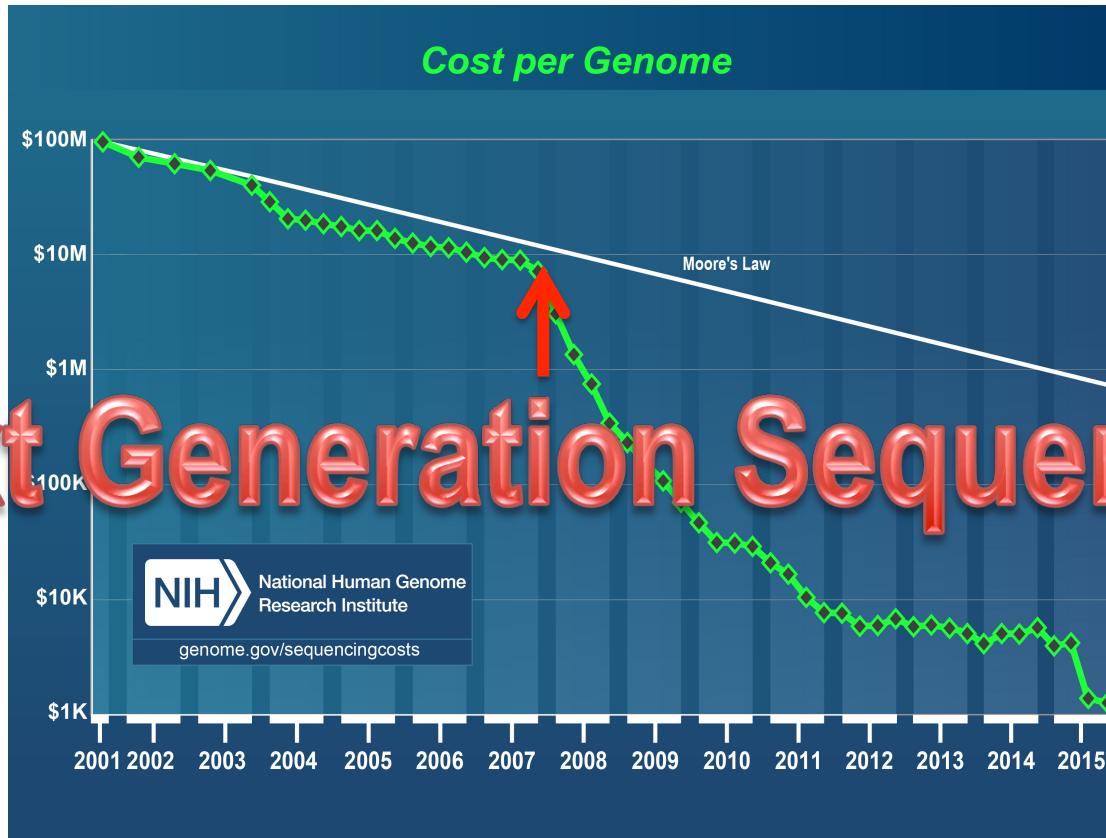
The [National Cancer Informatics Program](#) (NCIP), part of the NCI Center for Biomedical Informatics and Information Technology (CBIIT), is the institute's main bioinformatics initiative.

The collection, analysis, storage, retrieval, and distribution of "big data" are essential for many aspects of cancer research—especially for cancer genomics, in which millions of data points are frequently collected on each patient—and the monitoring of clinical trials.



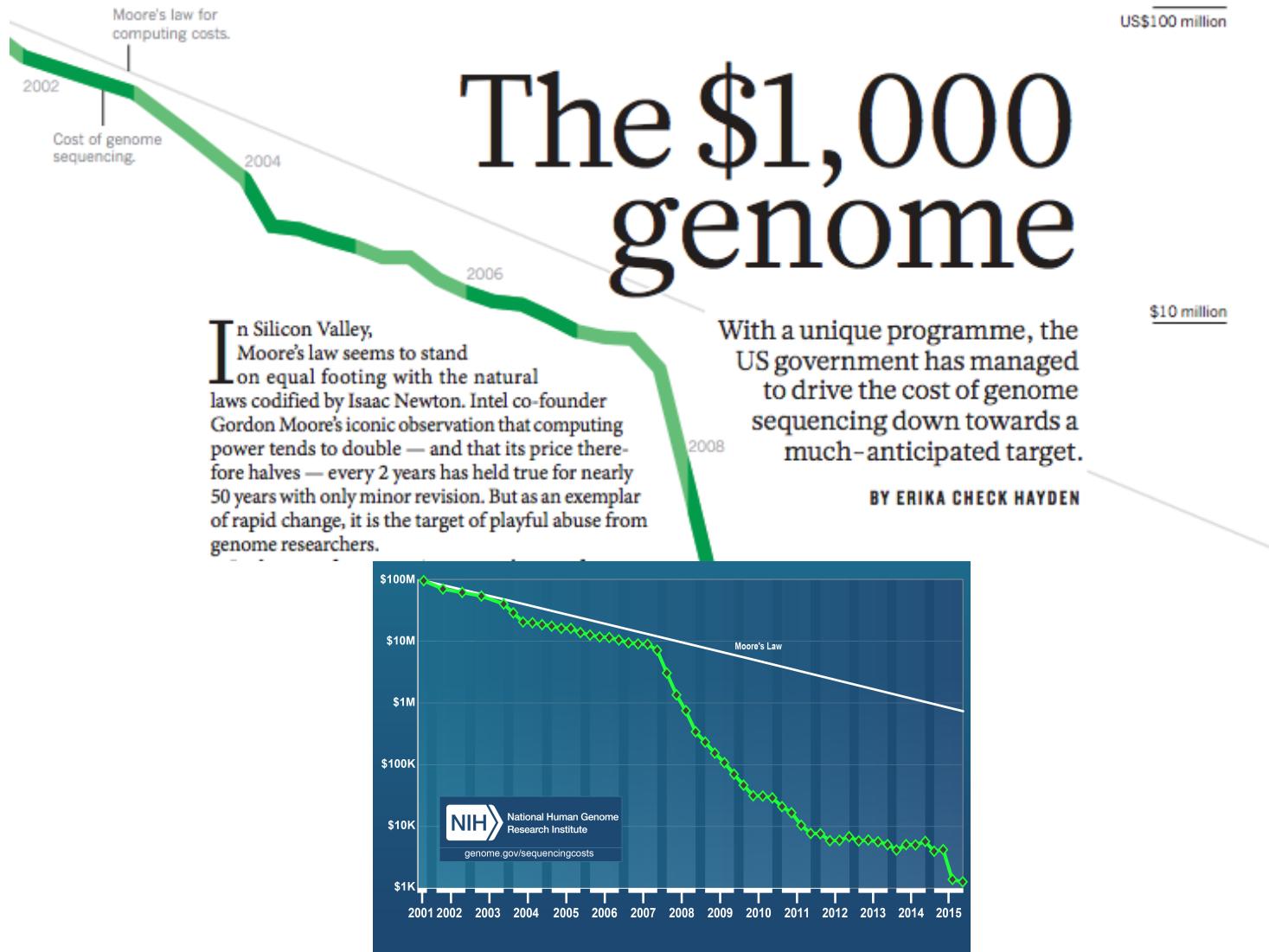
# DNA Sequencing Cost

- Cost of sequencing a human-sized genome  
\$1,000,000,000 US dollars (2001) → \$1,000 (2015)



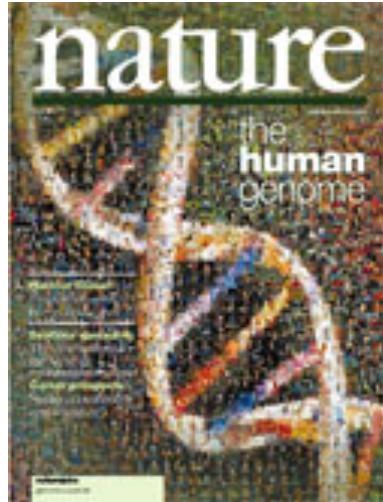
## Next Generation Sequencing

# The \$1000 Genome



# Sequencing a Human Genome

Human Genome Project  
2003



Today  
2017



~\$1B  
~ 6-8 years

~\$2-4K  
~ 1-3 days

# The Current Bottleneck of Bioinformatics



# Data Explosion!

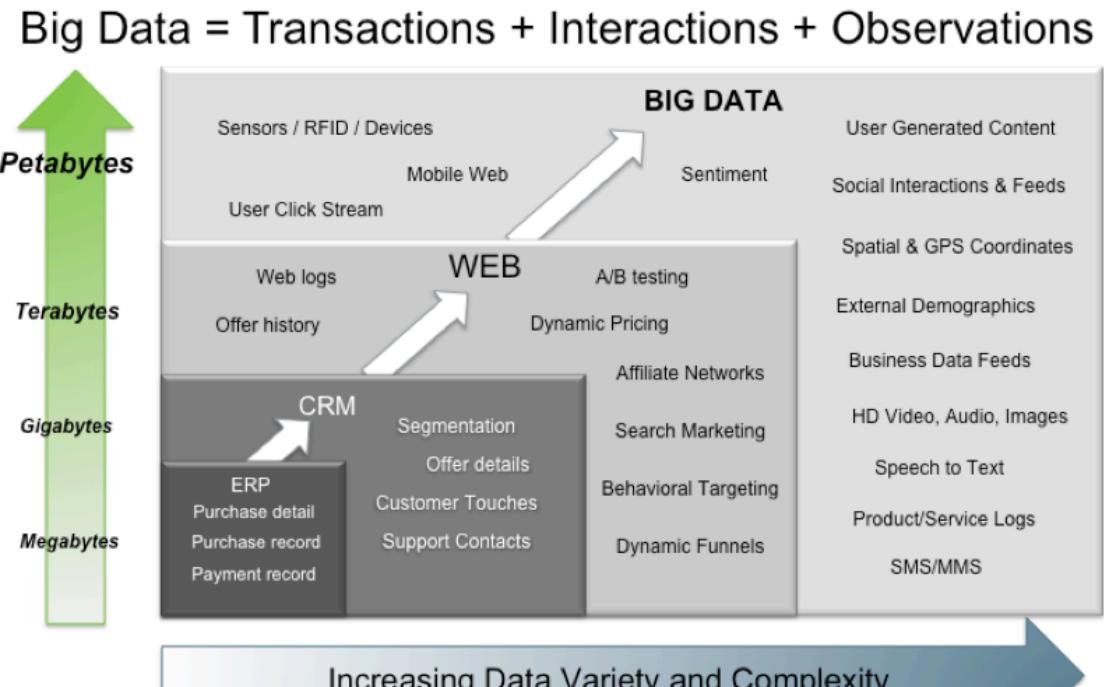
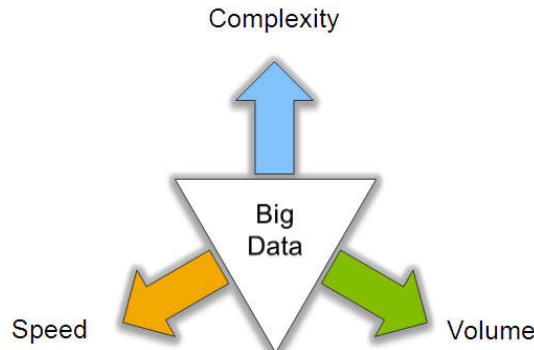
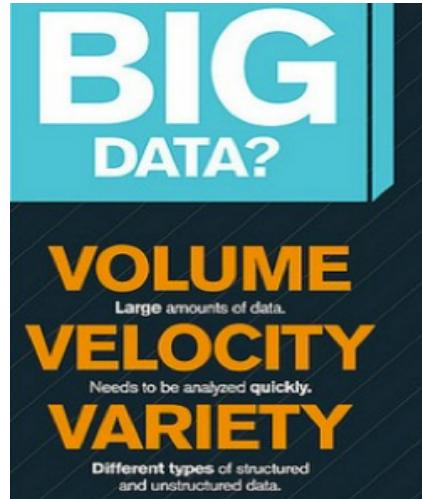
## Every minute:

- Google receives over 4 million search queries
- Facebook users share nearly 2.5 million pieces of content.
- Twitter users tweet nearly 300,000 times.
- Instagram users post nearly 220,000 new photos.
- YouTube users upload 72 hours of new video content.
- Apple users download nearly 50,000 apps.
- Email users send over 200 million messages.
- Amazon generates over \$80,000 in online sales.

# Big Data



# Big Data 3V's



*Source:* Contents of above graphic created in partnership with Teradata, Inc.

# Data Analysis Bottleneck



# Big Data to Knowledge (BD2K)

Thursday, October 9, 2014

NIH invests almost \$32 million to increase utility of biomedical research data

## Big Data to Knowledge (BD2K)



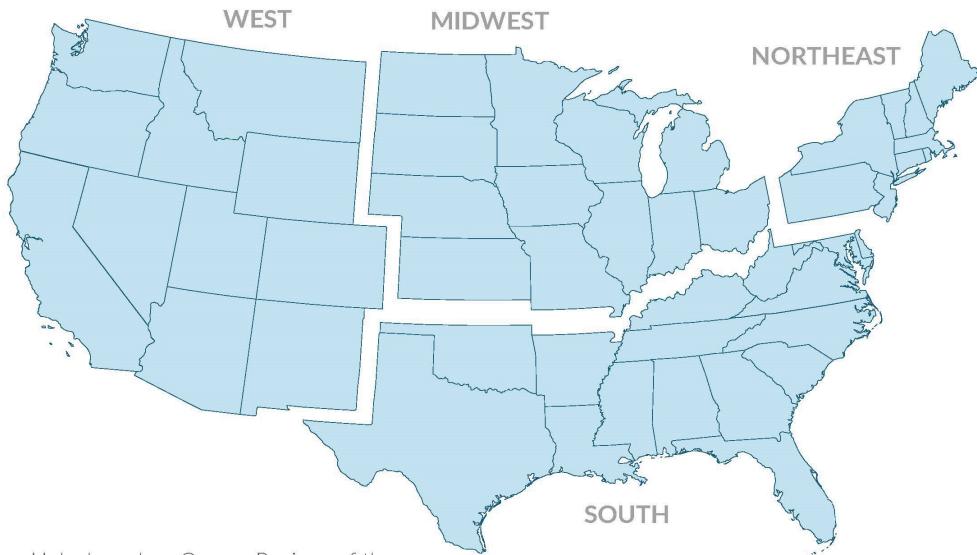
1. Facilitating Broad Use
2. Developing and Disseminating Analysis Methods and Software
3. Enhancing Training
4. Establishing Centers of Excellence

<http://bd2k.nih.gov>



the  
constitute an  
ected to

# NSF Big Data Research Initiative



Hubs based on Census Regions of the  
United States

Alaska & Hawaii are part of the West Region  
US Territories can participate in any region



18

National Science Foundation



## 4 Regional Big Data Hub

# Challenges or Opportunities?

- ✓ Large amount of data.
- ✓ The need for the relevant models.
- ✓ The computational complexity of these models.
- ✓ Various simulation methodologies for these models.

Biological, mathematical and computational challenges related to algorithm development and efficient computation are essential.



# Not Just Nucleotide Sequence Data...

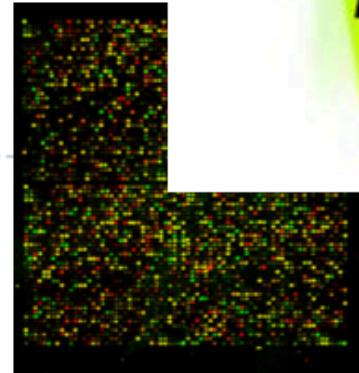


MKLTLKNLSMAIMMSIVMGS  
SAMAADSNEKTEVQGLNHS...GAS  
GYLPEHTLF...  
ADYLEQD  
LHDHYLDI  
DRARKDG  
DEIKSLKF  
QTYPGRFPMGK  
HTFEEEIEFVQGLLNHSTGK  
NIGIYPEIKA...PWFHQEGKDI  
AAKTLLEV...LKKYGYTGKDDKV

Protein sequence data

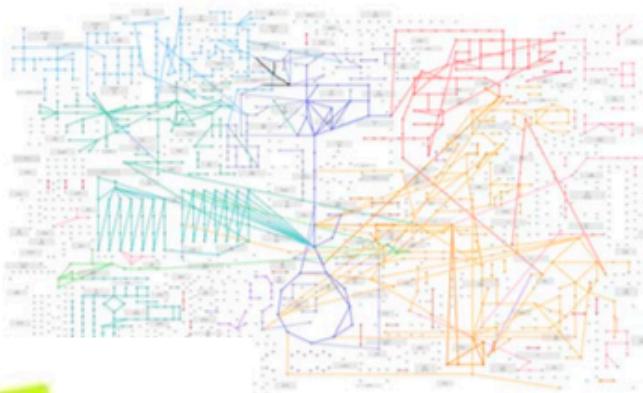


Protein str

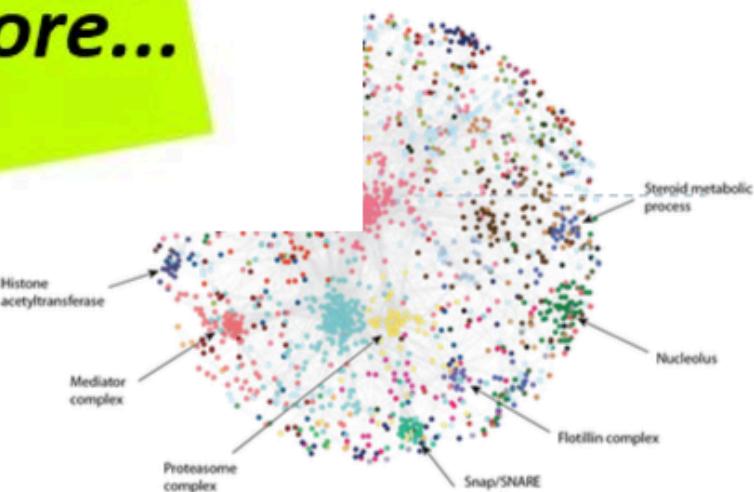


Gene expression data

*and many*  
**More...**

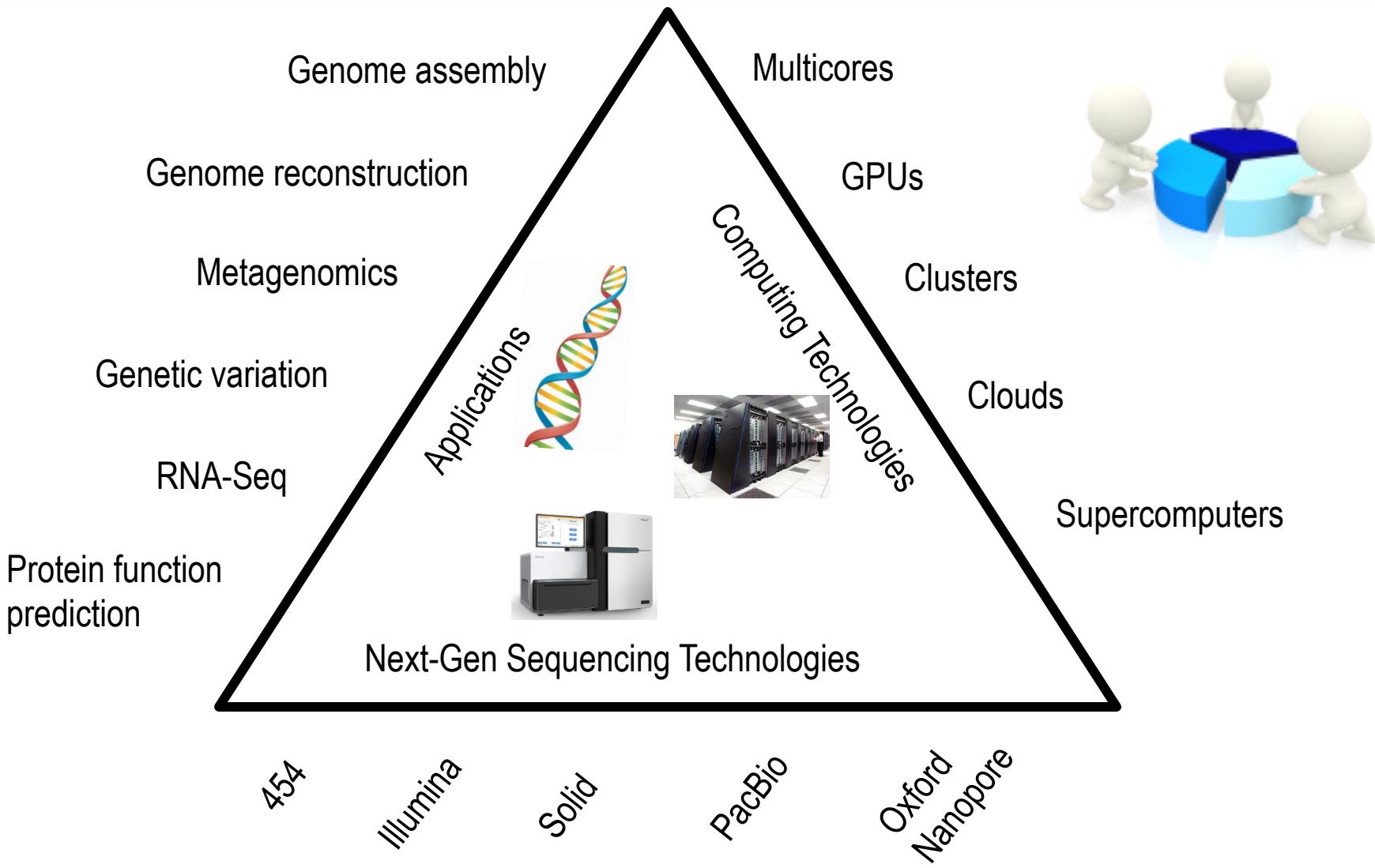


biological pathway data



Protein interaction network data

# Research Initiative



# Genome Size

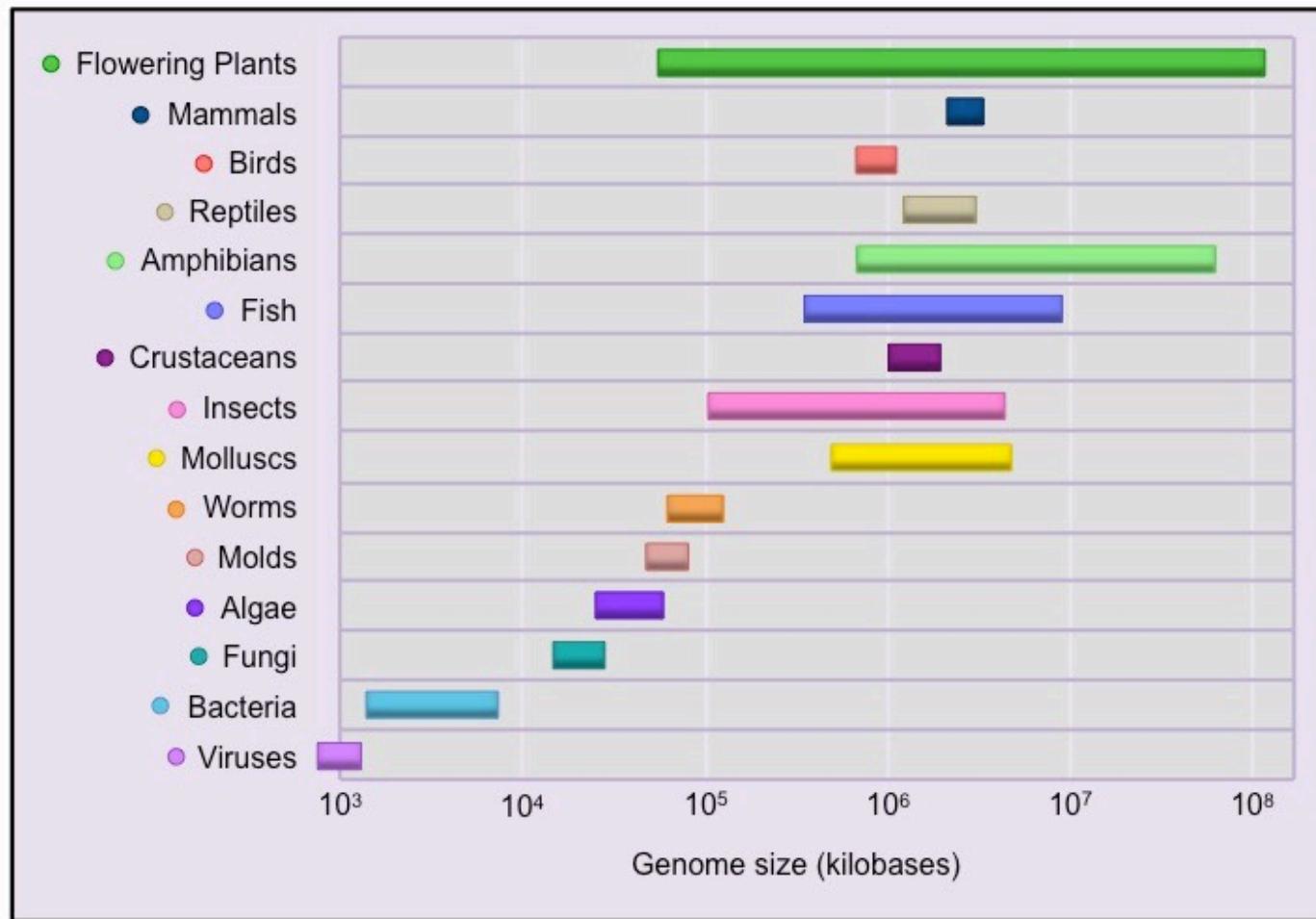
## Comparison of Genome Size in Different Organisms

Species	<i>T2 phage</i>	<i>Escherichia coli</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Paris japonica</i>
Genome Size					
Common Name	 Virus	 Bacteria	 Fruit fly	 Human	 Canopy Plant

<http://ib.bioninja.com.au/standard-level/topic-3-genetics/32-chromosomes/genome-size.html>

# Genome Size

## Variation in Genome Sizes For Different Types of Organisms



<http://ib.bioninja.com.au/standard-level/topic-3-genetics/32-chromosomes/genome-size.html>

# Assembling a Genome

1. Shear & Sequence DNA



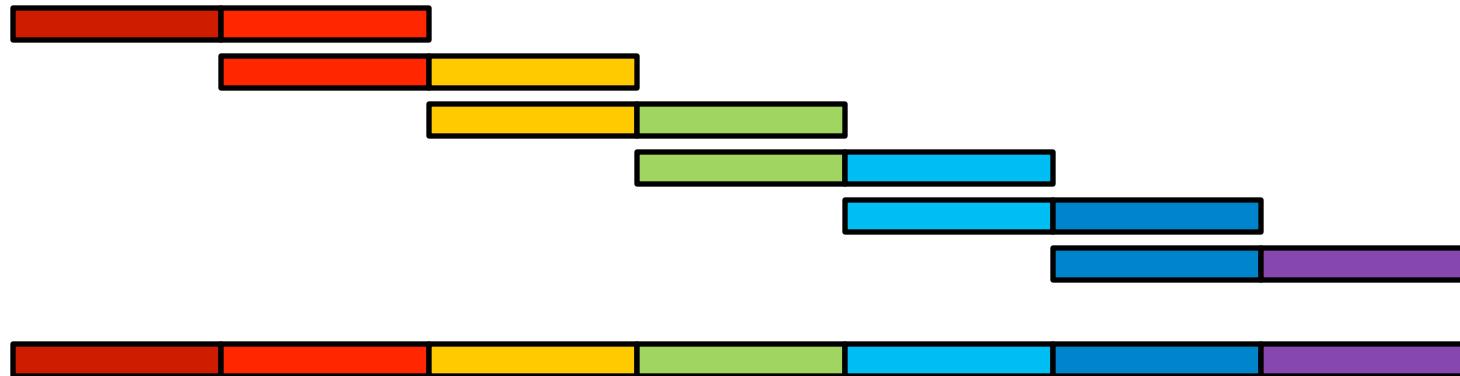
2. Construct assembly graph from overlapping reads

...AGCCTAGGA**TGCGCGACACGT**

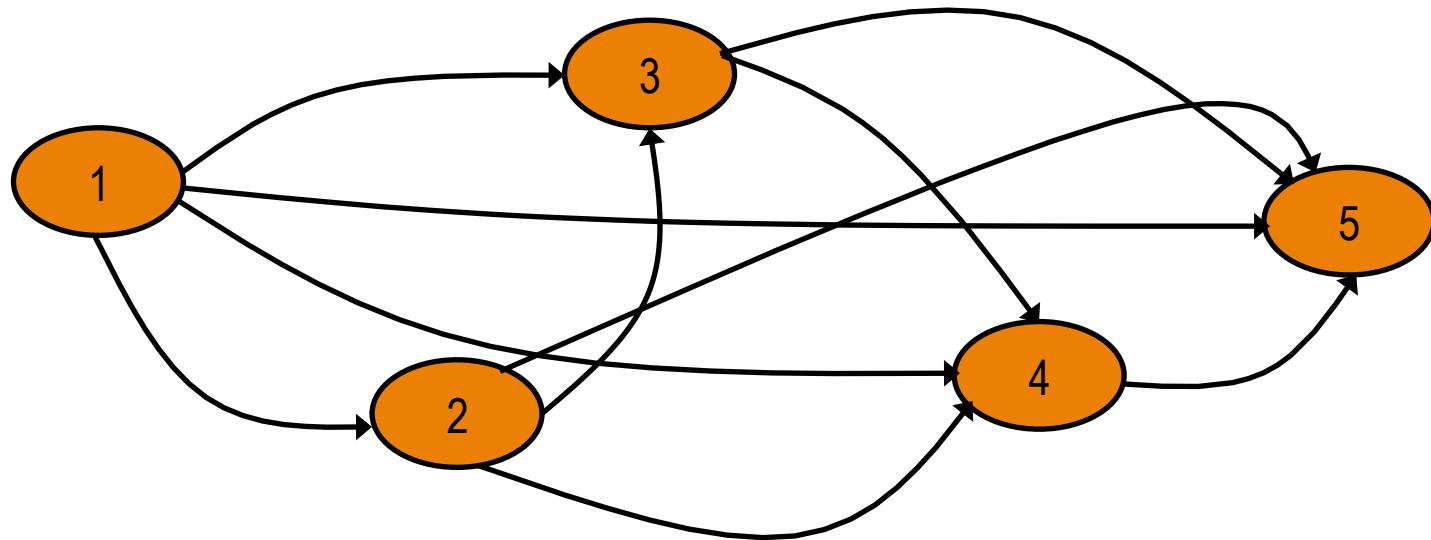
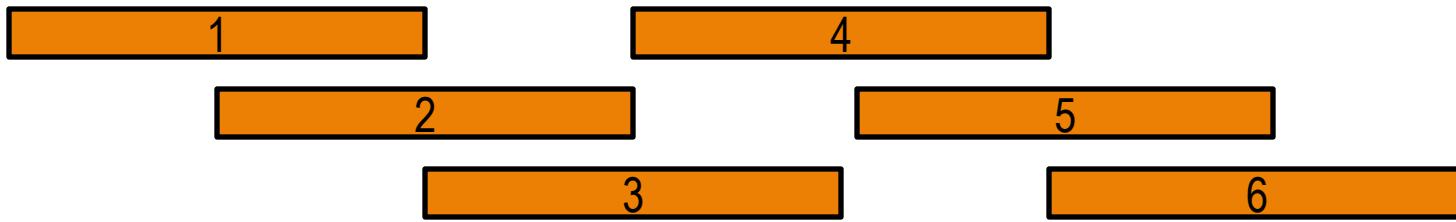
**TGCGCGACACGT**CGCATA**TCCGGTTGAT**

**TCCGGTTGAT**CACGAATA...CG...

3. Simplify assembly graph



# Assembly Complexity

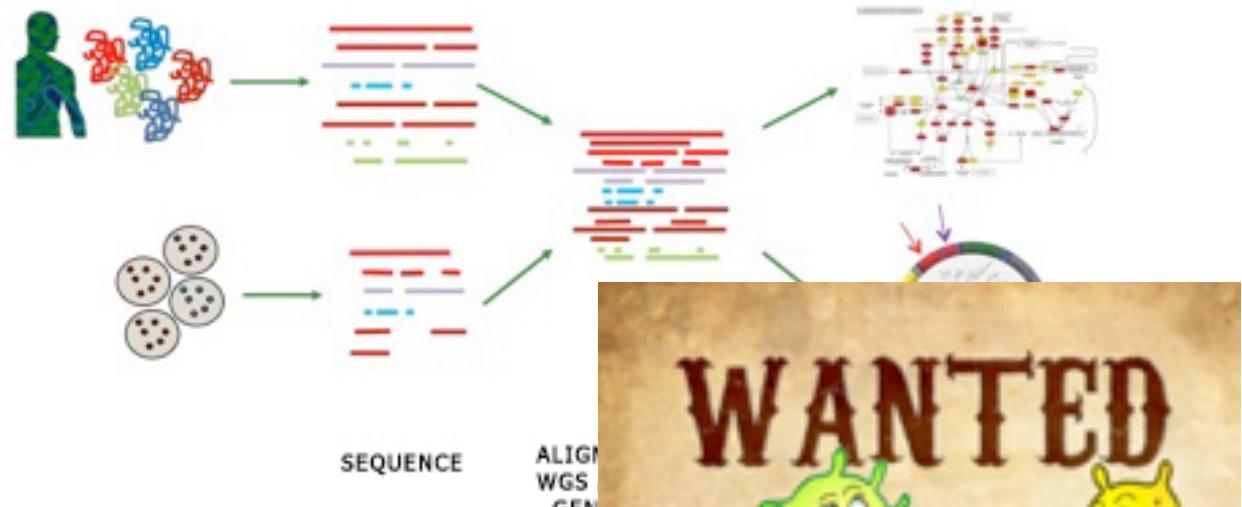
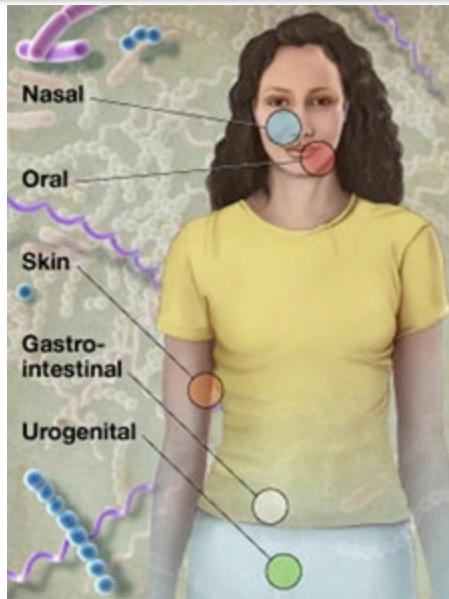


# What is Metagenomics?

**Culture-independent** genomic analysis of a community of microorganisms from **environmental sample**.



# Human Microbiome Project



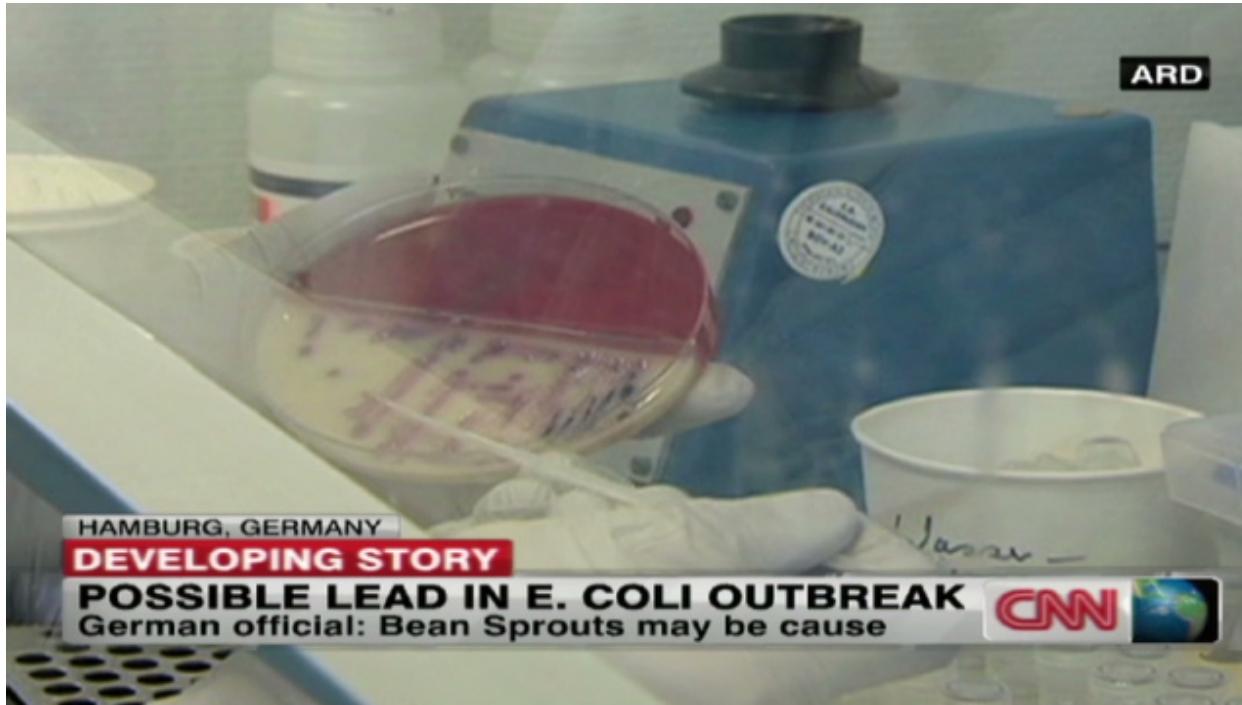
- From 242 healthy U.S. adults
- More than 5,000 samples
- 3000 reference genomes  
reference genomes are available



image courtesy of the NIH Common Fund

# 2011 German *E. coli* Outbreak.

In May through June 2011, a novel strain of *Escherichia coli* O104:H4 bacteria caused a serious outbreak of foodborne illness.



**In all, 3,950 people were affected and 53 died.**

# 2014 Ebola Virus Disease



**BREAKING NEWS**  
**EBOLA PATIENT ARRIVES AT EMORY  
UNIVERSITY HOSPITAL IN ATLANTA**

**FROM SPREADING CNN EBOLA HAS INFECTED MORE**

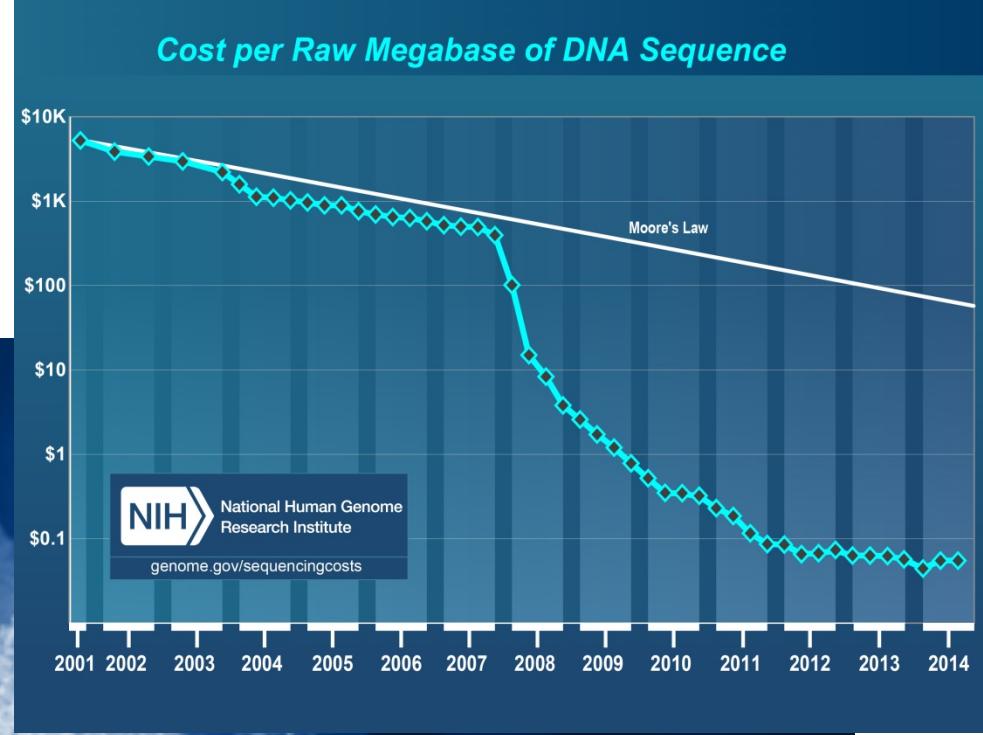
**LIVE  
CNN**

**12:29 PM ET**

© CNN/WSB

# Metagenomics and Biosurveillance

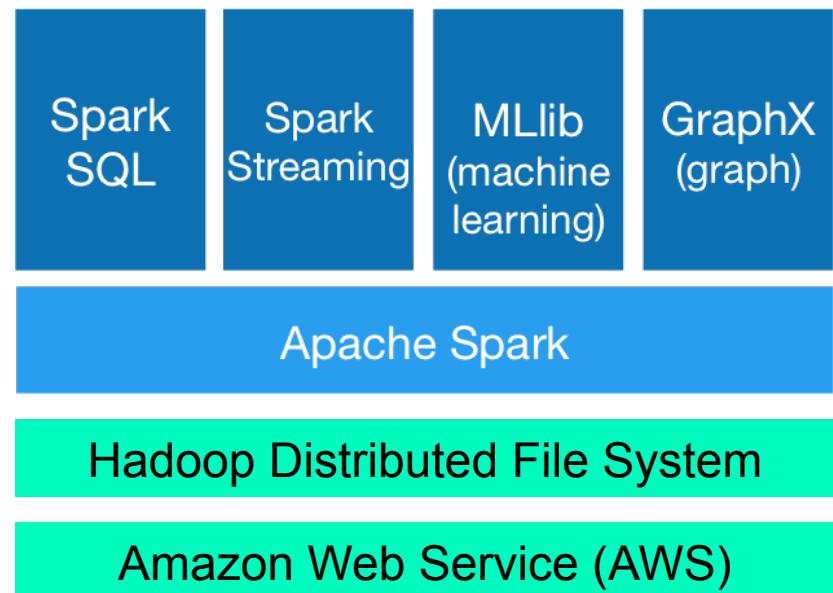
- Conventional methods to identify pathogens: laboratory culture, immunoassays, and genotyping



- NGS metagenomic analysis can be a promising new method for biosurveillance.

# New Big Data Technology - Spark

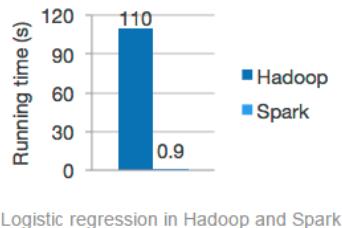
- An emerging **in-memory** processing framework
  - Iterative machine learning
  - Interactive data analytics
  - Scala based implementation
  - Master-Slave, HDFS, Fault Tolerance



## Speed

Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

Spark has an advanced DAG execution engine that supports cyclic data flow and in-memory computing.

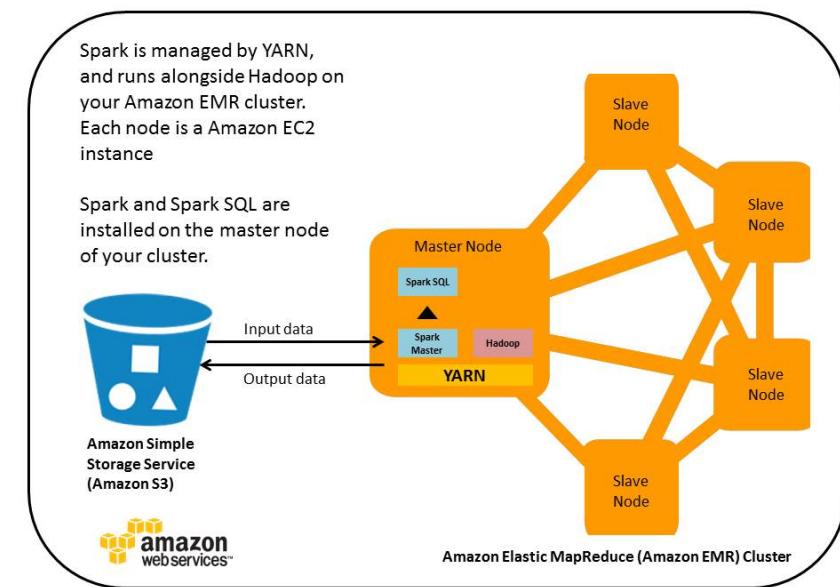
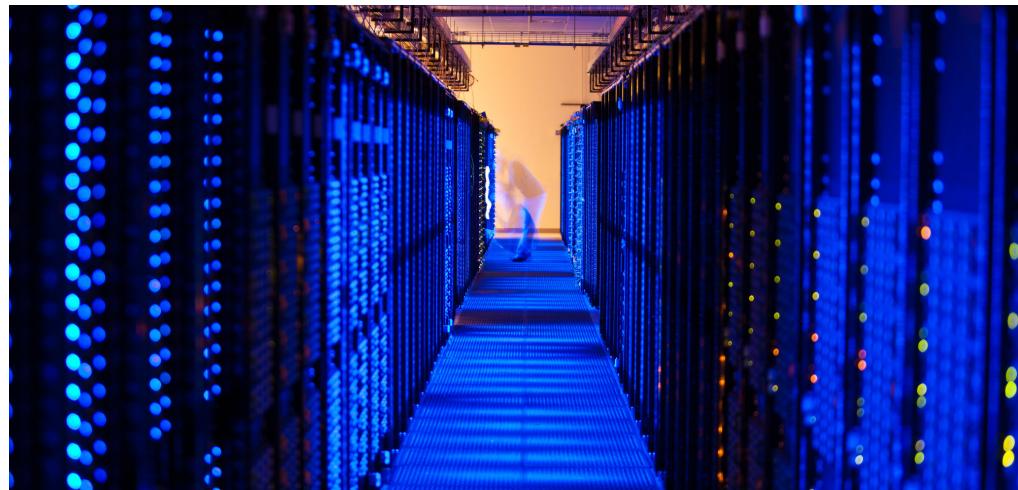


## Ease of Use

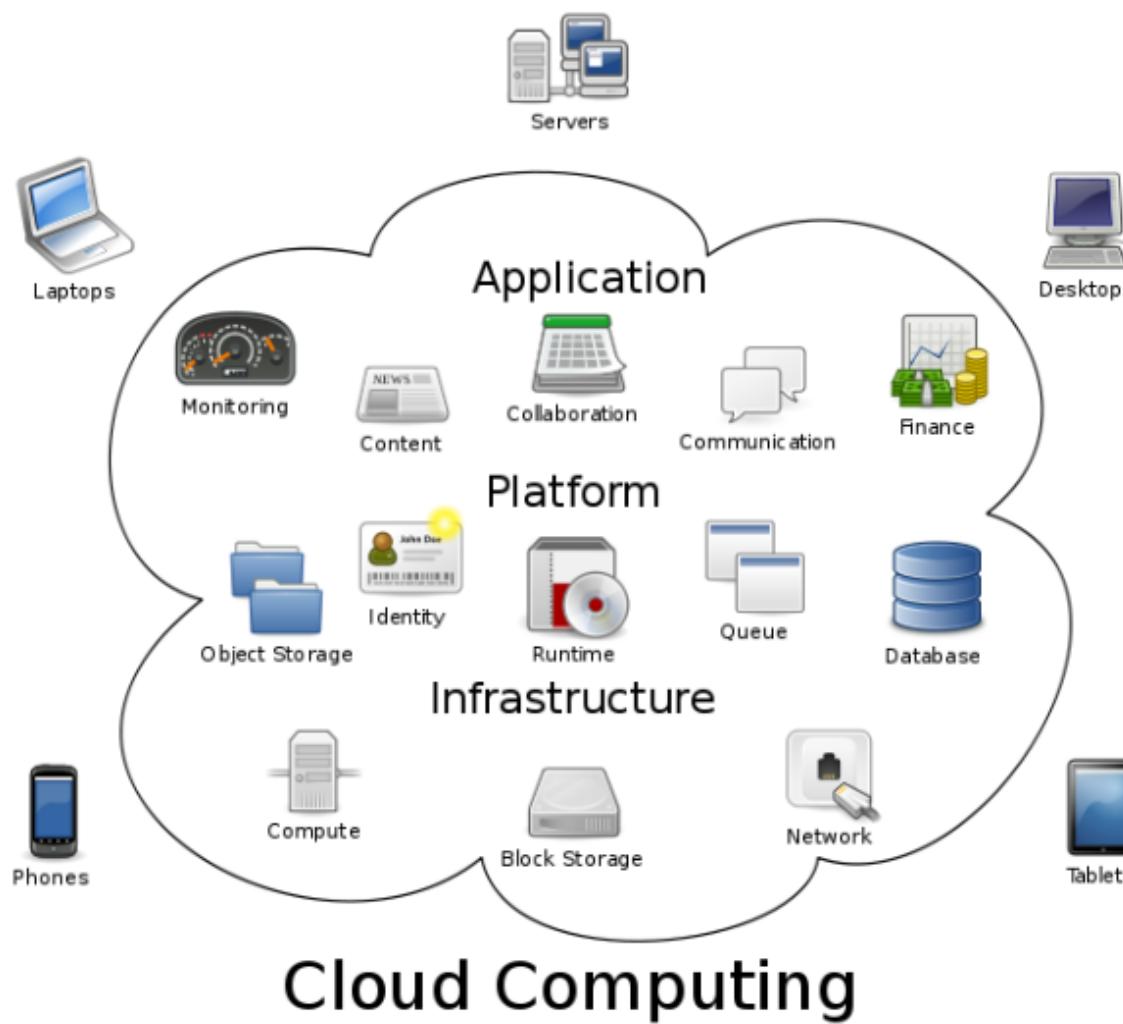
Write applications quickly in Java, Scala or Python.

```
file = spark.textFile("hdfs://...")  
  
file.flatMap(lambda line: line.split())  
    .map(lambda word: (word, 1))  
    .reduceByKey(lambda a, b: a+b)
```

# Cloud (Amazon Web Service) + Spark



# Cloud Computing



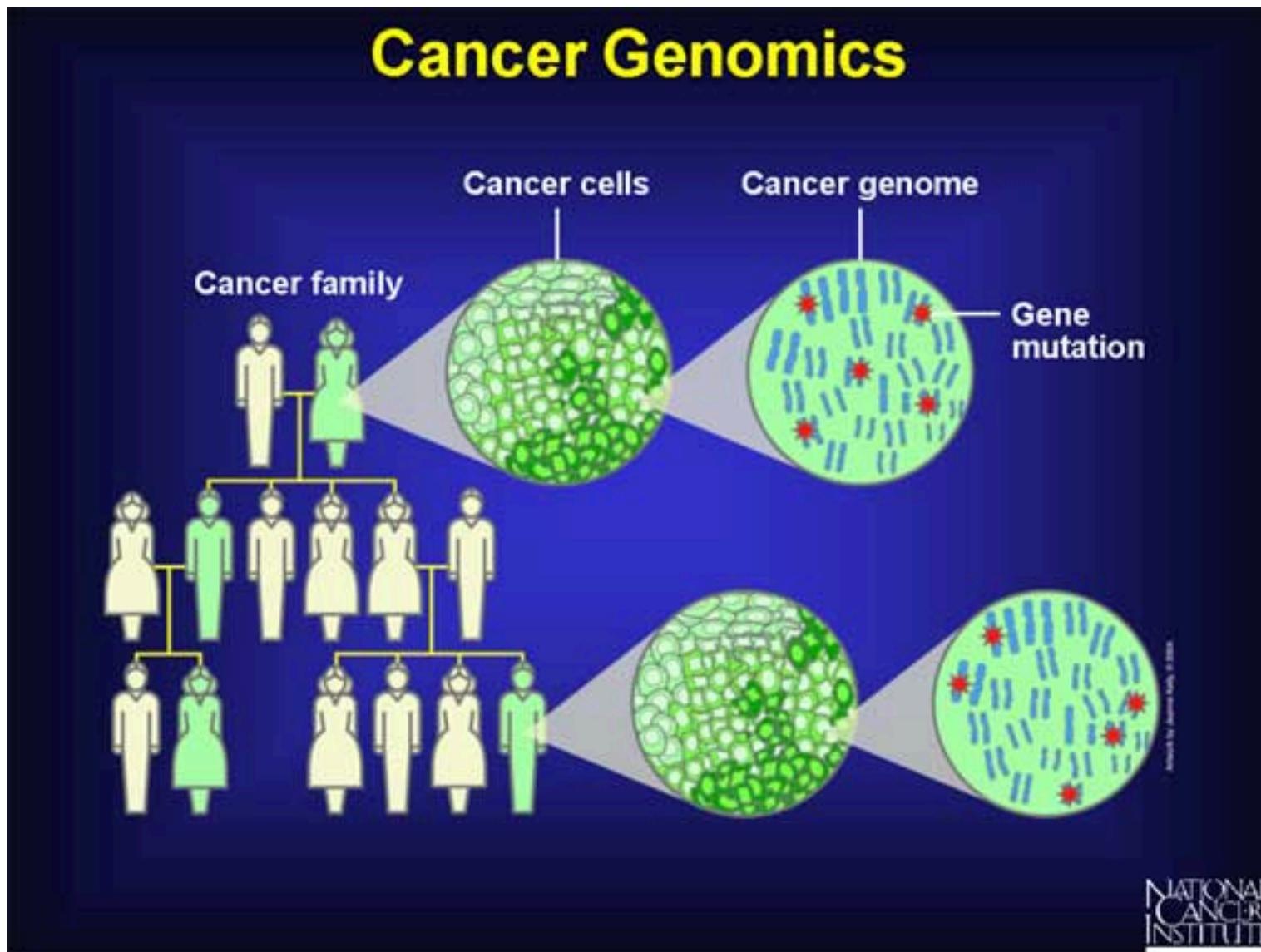
wikipedia: Cloud Computing

# Benefits

- Cost & management
  - Economies of scale, “out-sourced” resource management
- Reduced Time to deployment
  - Ease of assembly, works “out of the box”
- Scaling
  - On demand provisioning, co-locate data and compute
- Reliability
  - Massive, redundant, shared resources
- Sustainability
  - Hardware not owned

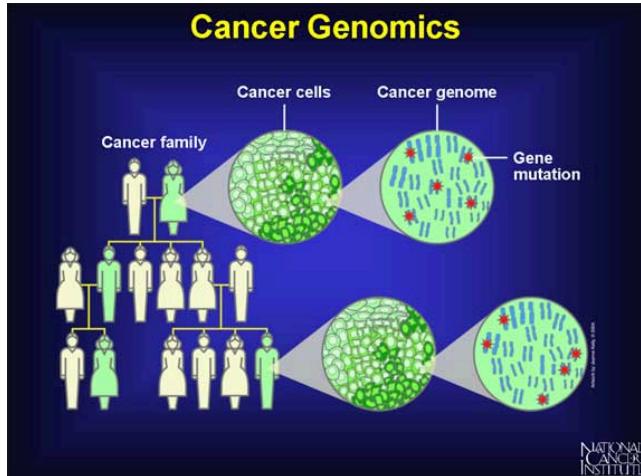
# Genomics and Cancer

## Cancer Genomics



NATIONAL  
CANCER  
INSTITUTE

# Hot Research Topics in Genomic Medicine



Cancer Genomics



Pharmacogenomics

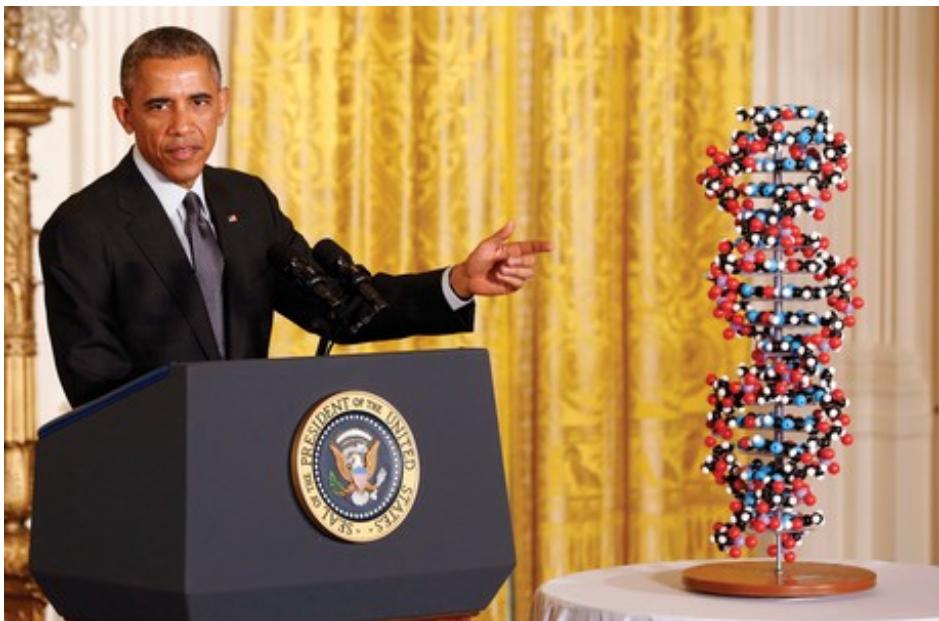
It is estimated that  
**350 MILLION**  
people worldwide suffer from one of over  
**7,000 rare diseases**

Ultra-Rare Disease Diagnosis



Prenatal & Newborn Genomic Analysis

# Precision Medicine



- A broader context for ‘individualizing’ medical care to advanced human health

# Precision Medicine

- Today: most medical care based on expected response of the average patient
- Tomorrow: medical care based on individual genomic, environmental, and lifestyle differences that enable more precise ways to prevent and treat diseases.

**nature** International weekly journal of science

Search  Advanced search

Home | News & Comment | Research | Careers & Jobs | Current Issue | Archive | Audio & Video | For Authors

News & Comment > News > 2016 > May > Article

NATURE | NEWS

Obama to seek \$215 million for precision-medicine plan

Details emerge as White House prepares to release budget request to Congress.

Sara Reardon

30 January 2015

E-alert RSS Facebook Twitter

Muzzled



Nine years of censorship

The screenshot shows a news article from the journal 'nature' about the US government's proposed budget for precision medicine. The article discusses the Obama administration's plan to seek \$215 million for a precision-medicine initiative. The URL of the article is visible in the browser's address bar. The page includes navigation links for the journal, such as Home, News & Comment, Research, and For Authors, along with social media sharing options. A sidebar features a cartoon illustration of a laboratory flask with a speech bubble, labeled 'Muzzled' and 'Nine years of censorship'.

# An integrative bioinformatics/biomedical knowledge of big data

- It is not easy to analyze the large-scale data to derive meaningful biological correlations by
  - absence of easy-to-use knowledge-based bioinformatics tools
  - difficulty to select appropriate software tools
  - lack of computing system resources



Bioinformatics + Bio Tech + Big Data + HPC = Gold Mine!!

# Question?

Thank you for listening and I will be happy to take questions!

