

# Applied Comparative Genomics

Michael Schatz

January 31, 2017

Lecture 1: Course Overview



# Welcome!

**The primary goal of the course is for students to be grounded in theory and leave the course empowered to conduct independent genomic analyses.**

- We will study the leading computational and quantitative approaches for comparing and analyzing genomes starting from raw sequencing data.
- The course will focus on human genomics and human medical applications, but the techniques will be broadly applicable across the tree of life.
- The topics will include genome assembly & comparative genomics, variant identification & analysis, gene expression & regulation, personal genome analysis, and cancer genomics.

**Course Webpage:**

<https://github.com/schatzlab/appliedgenomics>

**Course Discussions:**

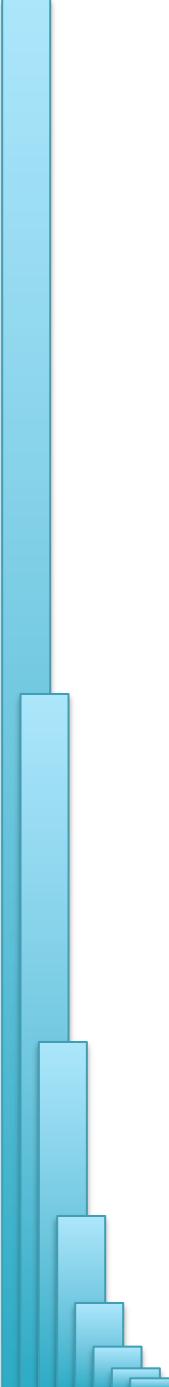
<http://piazza.com>

**Class Hours:**

Tues + Thurs @ 1:30p – 2:45p, Shaffer 304

**Office Hours:**

Tues + Thurs @ 3-4p and by appointment  
Please try Piazza first!



# Prerequisites and Resources

## **Prerequisites**

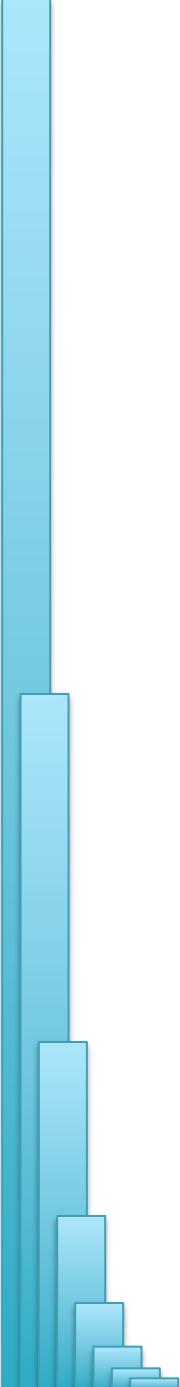
- No formal course requirements
- Access to an Apple or Linux Machine, or Install VirtualBox
- Familiarity with the Unix command line for exercises
  - bash, ls, grep, sed, + install published genomics tools
- Familiarity with a major programming language for project
  - C/C++, Java, R, Perl, Python

## **Primary Texts**

- None! We will be studying primary research papers

## **Other Resources:**

- Google, SEQanswers, Biostars, StackOverflow
- Applied Computational Genomics Course at UU: Spring 2017
- <https://github.com/quinlan-lab/applied-computational-genomics>
- Ben Langmead's teaching materials:
  - <http://www.langmead-lab.org/teaching-materials/>



# Grading Policies

## **Assessments:**

- ~5 Assignments: 30% Due at 11:59pm a week later  
***Practice using the tools we are discussing***
- 1 Exam: 30% In class (Tentatively 4/6)  
***Assess your performance, focusing on the methods***
- 1 Class Project: 40% Presented last week of class  
***Significant project developing a novel analysis/method***
- In-class Participation: Not graded, but there to help you!

## **Policies:**

- Scores assigned relative to the highest points awarded
- Automated testing and grading of assignments
- ***Late Days:***
  - Three (3) chances to extend the deadline for assignments by 24 hours without any penalty

# Course Webpage

Materials for JHU EN.600.649: Computational Genomics: Applied Comparative Genomics

10 commits 1 branch 0 releases 1 contributor MIT

Branch: master New pull request Create new file Upload files Find file Clone or download

File	Action	Time
mschatz committed on GitHub Update README.md		Latest commit abc@cf7 23 minutes ago
assignments	Add	5 hours ago
lectures	add	an hour ago
policies	Add	5 hours ago
LICENSE	Initial commit	5 hours ago
README.md	Update README.md	23 minutes ago

JHU EN.600.649: Computational Genomics: Applied Comparative Genomics

Michael Schatz (mschatz @ cs.jhu.edu)  
Class Hours: Tuesday + Thursday @ 1:30p - 2:45p in Shaffer 304  
Office Hours: Tuesday + Thursday @ 3-4p in Malone 323 and by appointment

The primary goal of the course is for students to be grounded in theory and leave the course empowered to

<https://github.com/schatzlab/appliedgenomics>

# Piazza

The screenshot shows a web browser window with the URL <https://piazza.com/class/lly77snf0k4zu?cid=6>. The browser title bar indicates the page is secure (https) and has a port number of 600.649. The address bar also shows the URL. The browser toolbar includes various icons for navigation, search, and bookmarks.

The main content area is the Piazza platform. At the top, there's a navigation bar with links for "Q & A", "Resources", "Statistics", and "Manage Class". On the right side of the header, there's a profile picture for "Michael Schatz" and a "Logout" button.

The left sidebar lists course materials: hw1, hw2, hw3, hw4, hw5. Below this, there are filters for "Unread", "Updated", "Unresolved", and "Following". A prominent blue button says "New Post". A search bar says "Search or add a post...".

The main content area starts with a note titled "Welcome to Piazza!" by Michael Schatz. The note text reads:

Welcome to Piazza! We'll be conducting all class-related discussion here this term. The quicker you begin asking questions on Piazza (rather than via emails), the quicker you'll benefit from the collective knowledge of your classmates and instructors. We encourage you to ask questions when you're struggling to understand a concept—you can even do so anonymously.

Below the note, it says "-Michael Schatz". There are "other" and "edit" buttons, and a note that it was "Updated 3 hours ago by Michael Schatz".

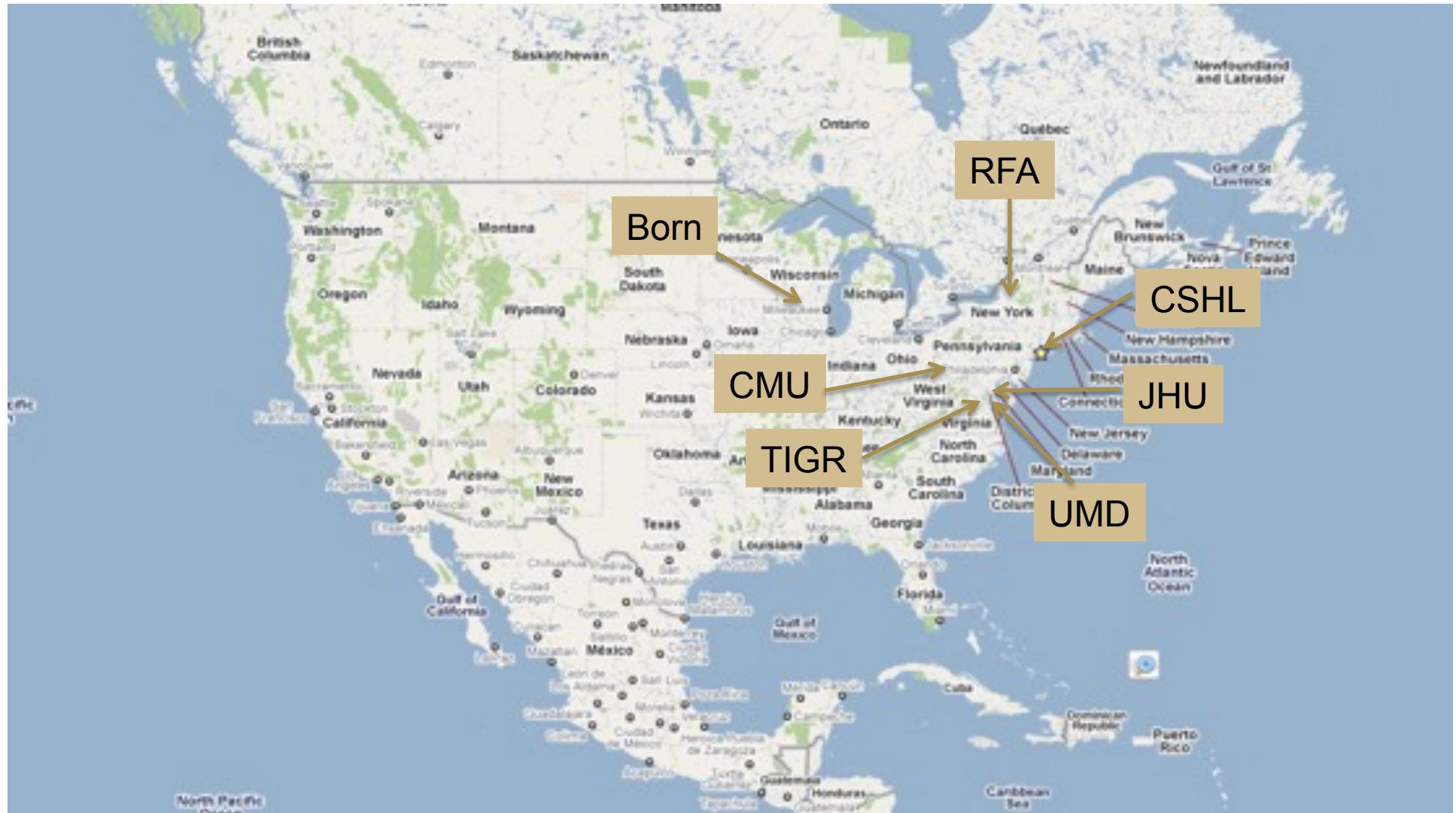
Below the note, there's a section titled "followup discussions" with a link "for lingering questions and comments". It includes a "Start a new followup discussion" button and a "Compose a new followup discussion" input field.

At the bottom of the main content area, there are statistics: "Average Response Time: N/A", "Special Mentions: There are no special mentions at this time.", "Online Now: 1", and "This Week: 1".

At the very bottom of the page, there's a footer with copyright information: "Copyright © 2013 Piazza Technologies, Inc. All Rights Reserved. Privacy Policy Copyright Policy Terms of Use Blog Report Bug".

<https://piazza.com/jhu/spring2017/600649/home>

# A Little About Me



# A Little About Me



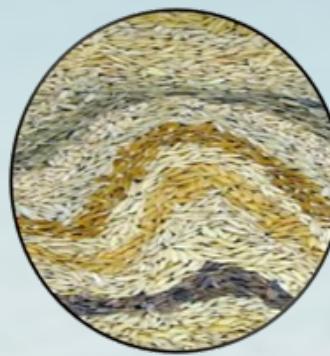
# Schatzlab Overview



## Human Genetics

Role of mutations in disease

Fang *et al.* (2016)  
Narzisi *et al.* (2015)



## Agricultural Genomics

Genomes & Transcriptomes

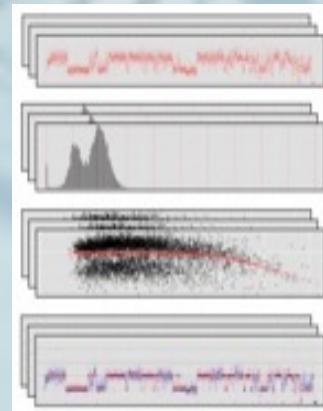
Lemmon *et al.* (2016)  
Ming *et al.* (2015)



## Algorithmics & Systems Research

Ultra-large scale biocomputing

Stevens *et al.* (2015)  
Marcus *et al.* (2014)



## Biotechnology Development

Single Cell + Single Molecule Sequencing

Chin *et al.* (2016)  
Garvin *et al.* (2015)

# DNA: The secret of life



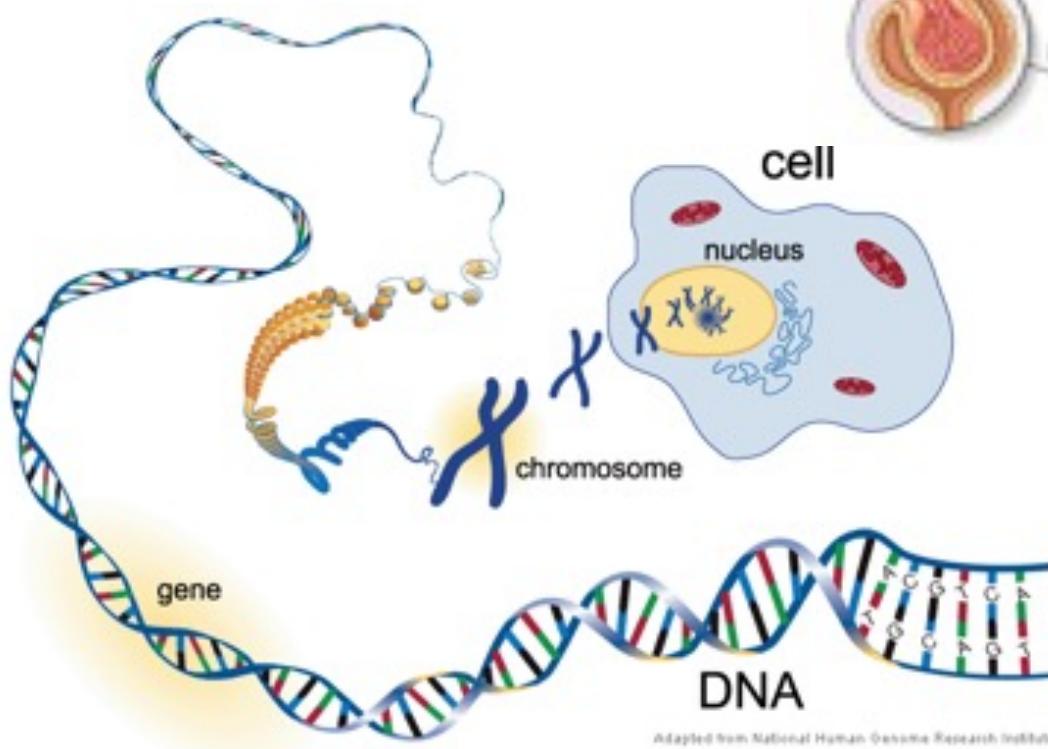
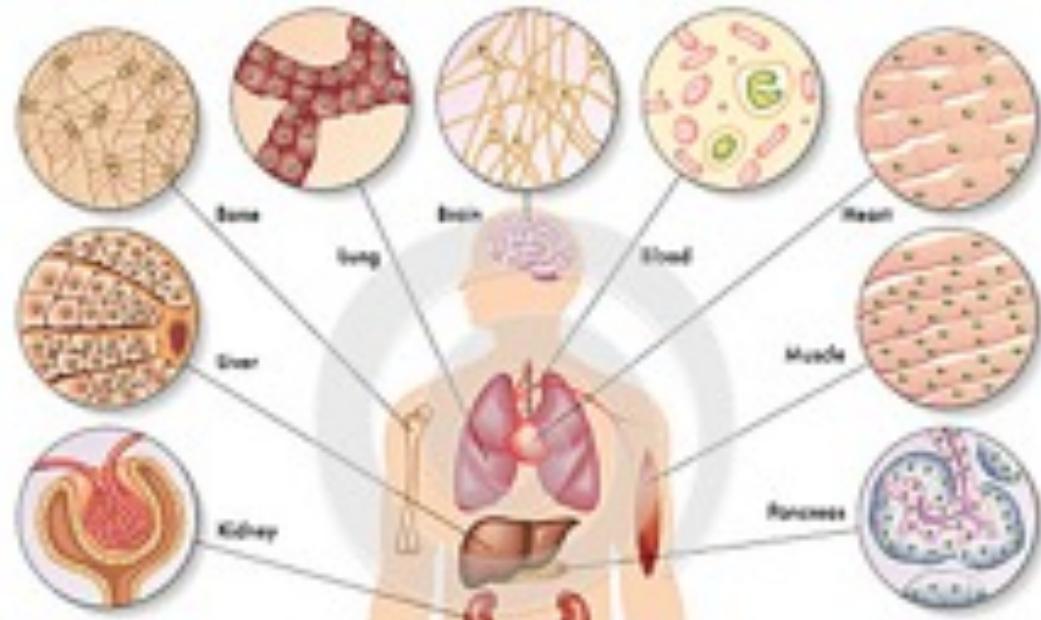
***Your DNA, along with your environment and experiences, shapes who you are***

- Height
- Hair, eye, skin color
- Broad/narrow, small/large features
- Susceptibility to disease
- Response to drug treatments
- Longevity and cognition

Physical traits tend to be strongly genetic, social characteristics tend to be strongly environmental, and everything else is a combination

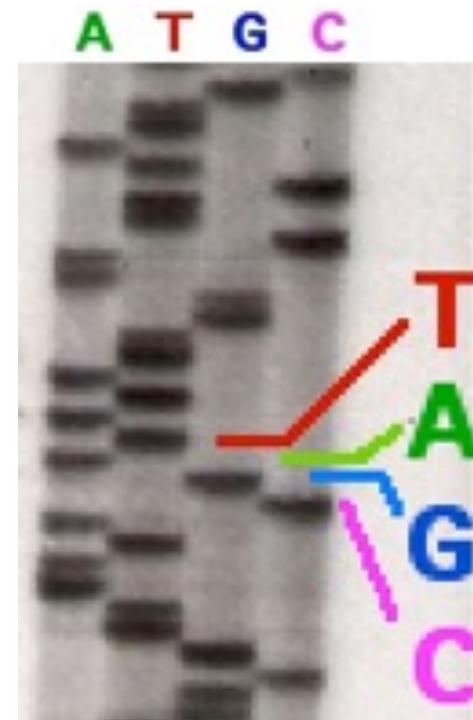
# Cells & DNA

Each cell of your body contains an exact copy of your 3 billion base pair genome.



Your specific nucleotide sequence encodes the genetic program for your cells and ultimately your traits

# The Origins of DNA Sequencing



Radioactive Chain Termination  
5000bp / week / person

**Nucleotide sequence of bacteriophage φX174 DNA**  
Sanger, F. et al. (1977) *Nature*. 265: 687 - 695

<http://en.wikipedia.org/wiki/File:Sequencing.jpg>  
<http://www.answers.com/topic/automated-sequencer>

# Milestones in DNA Sequencing



(TIGR/Celera, 1995-2001)

# Genomics across the tree of life



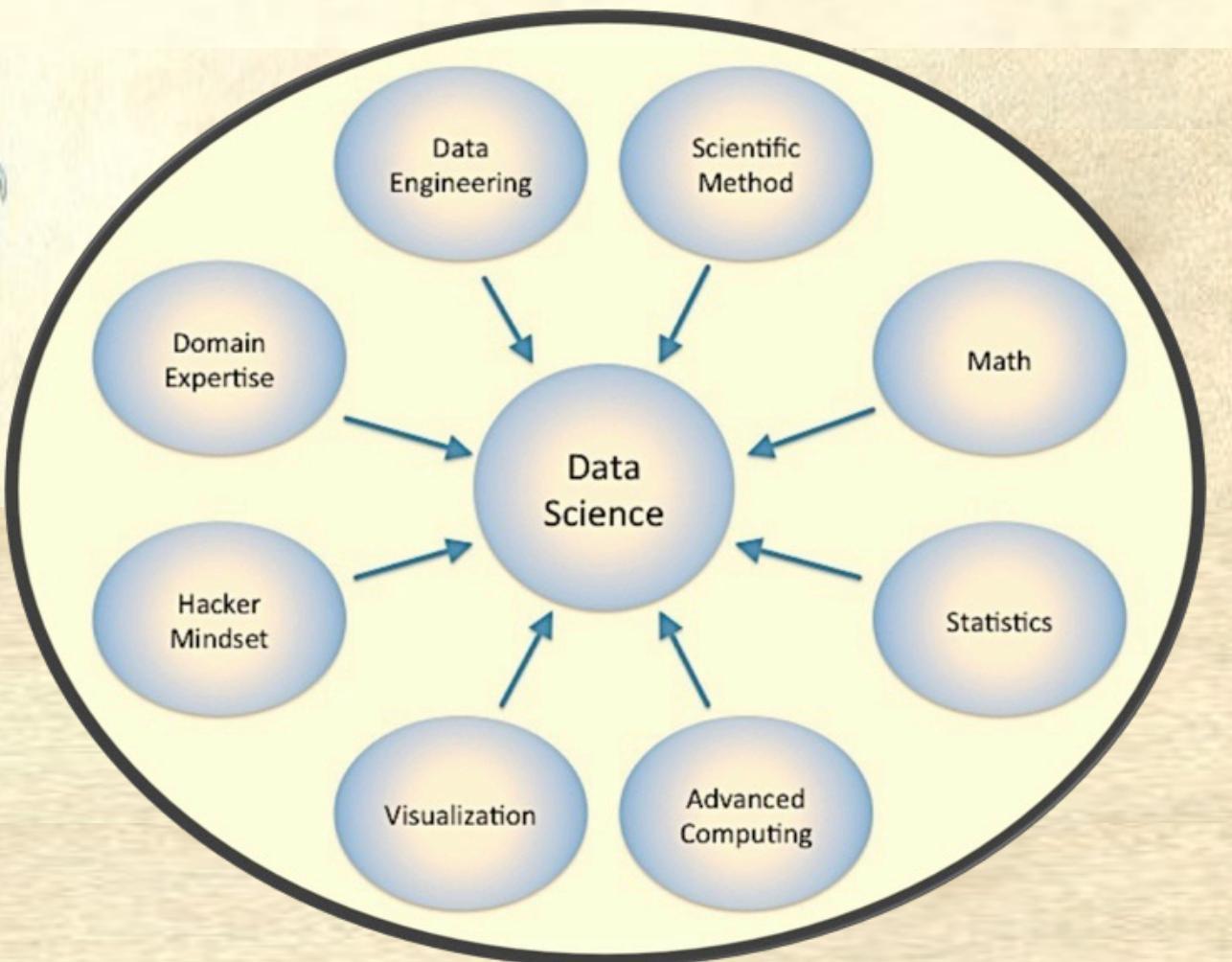
# Unsolved Questions in Biology

- What is your genome sequence?
- 
- 
- The instruments provide the data, but none of the answers to any of these questions.
- 
- 
- 
- 
- 
- 
- ***What software and systems will?***
- 
- 
- ***And who will create them?***
- 
- 
- ***Plus thousands and thousands more***



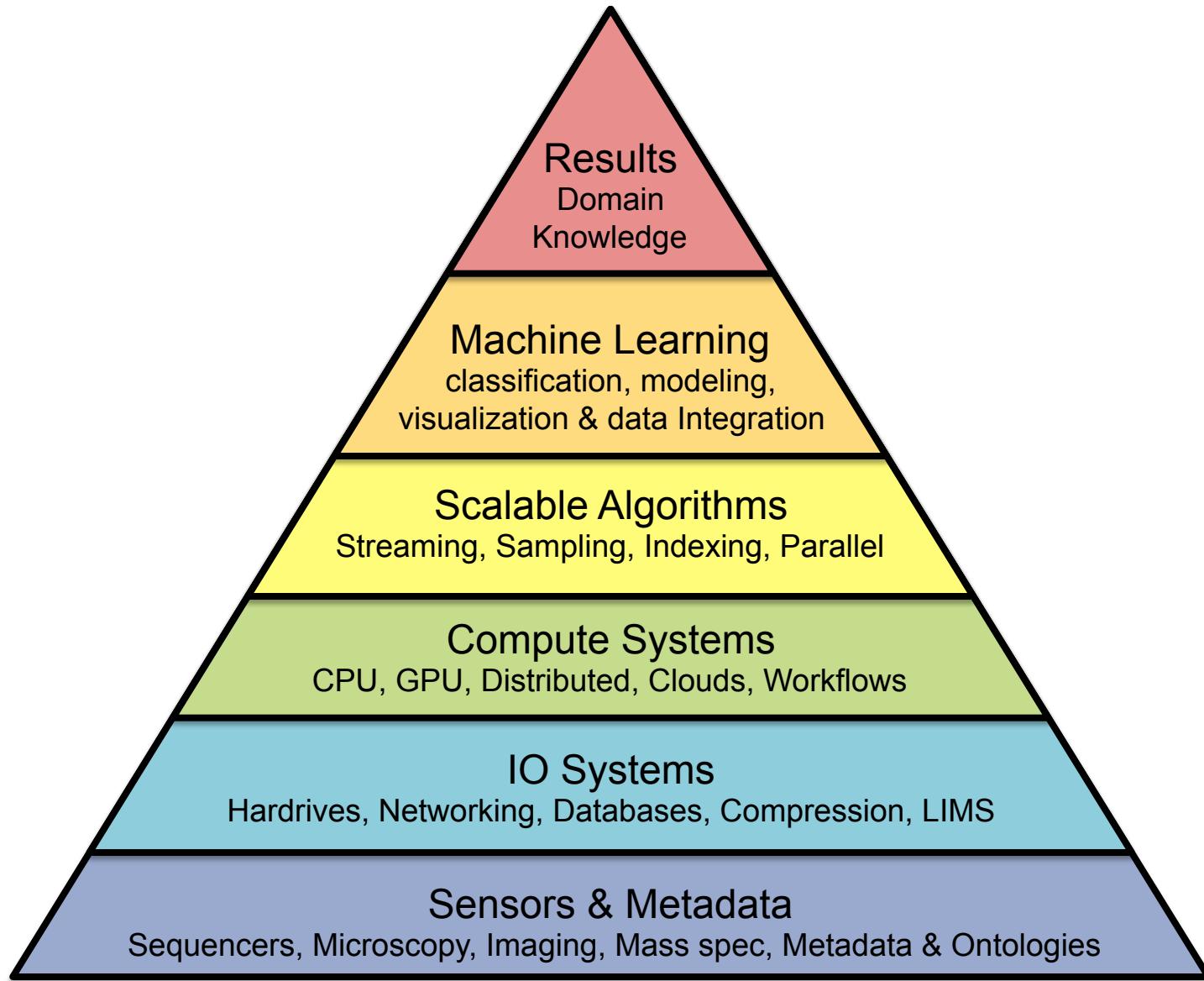


# Who is a Data Scientist?

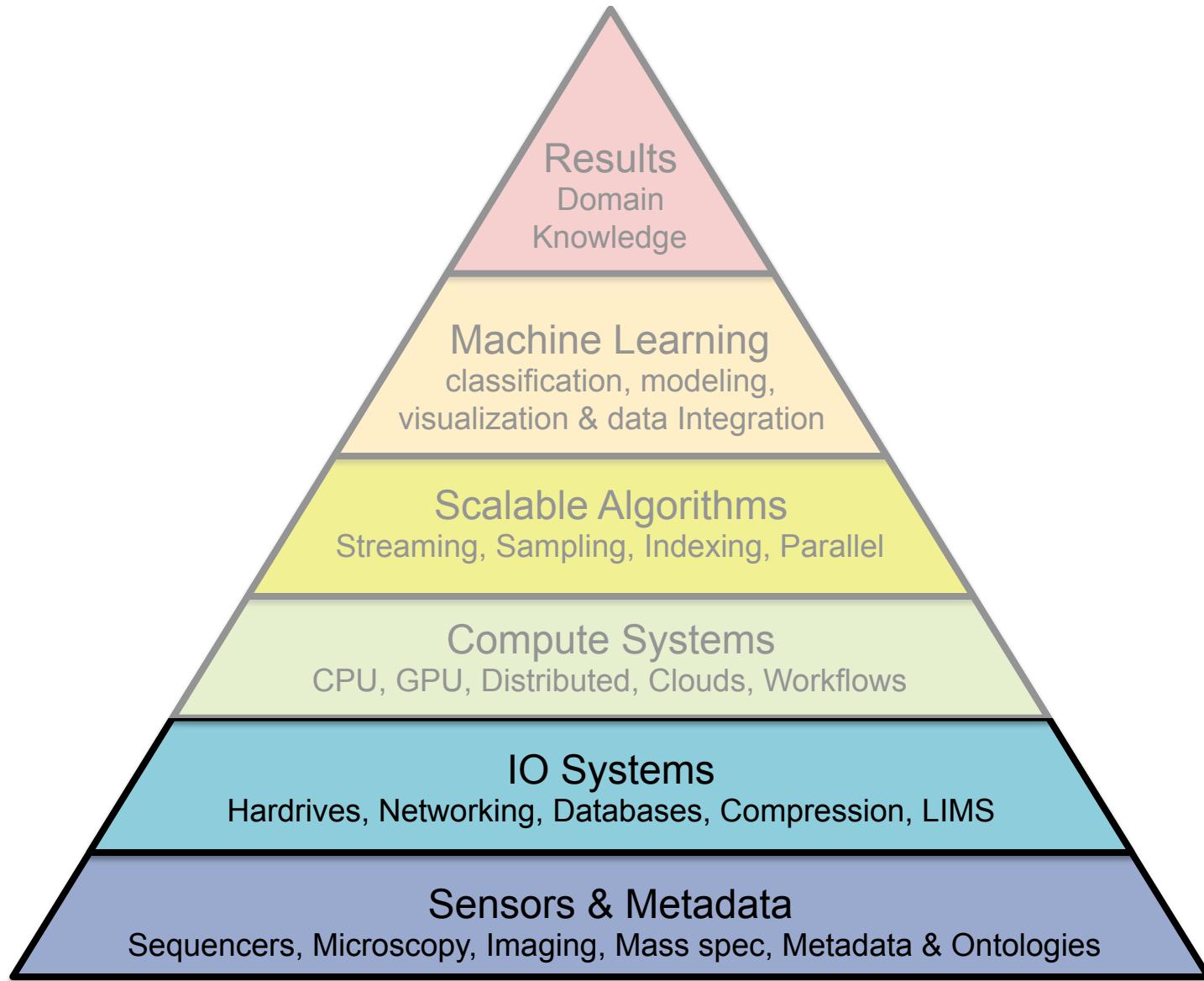


[http://en.wikipedia.org/wiki/Data\\_science](http://en.wikipedia.org/wiki/Data_science)

# Comparative Genomics Technologies



# Comparative Genomics Technologies

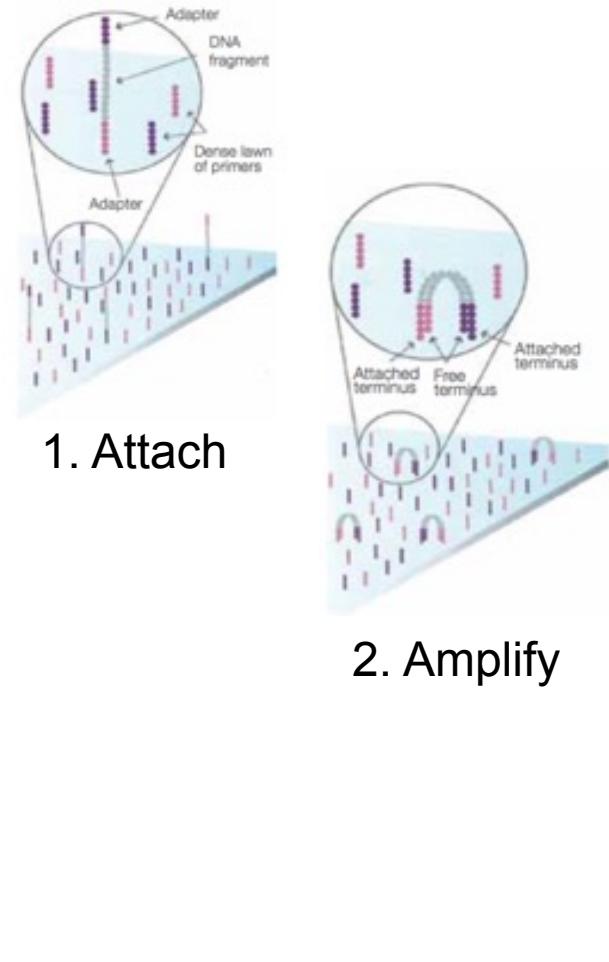


# Massively Parallel Sequencing



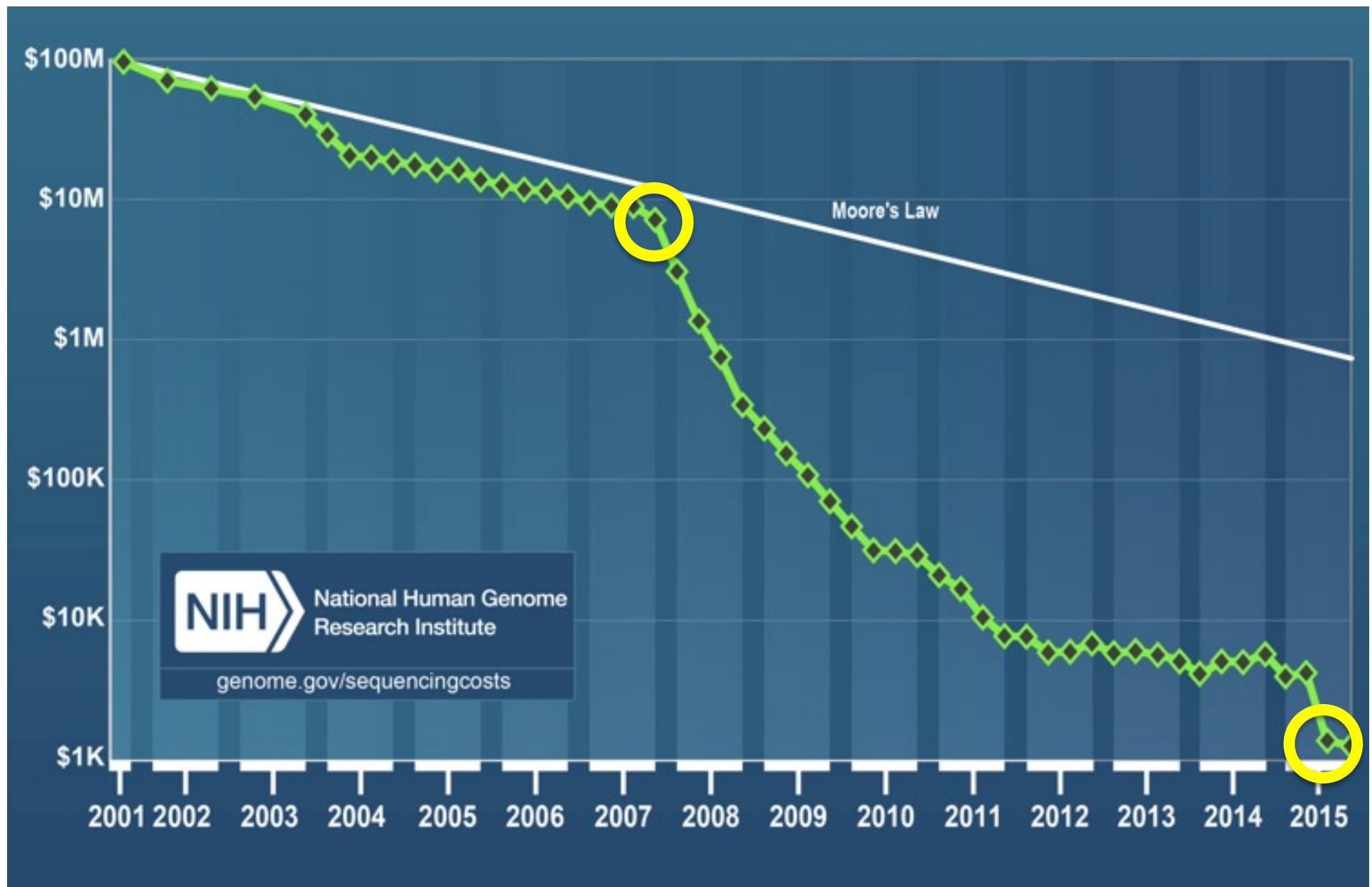
**Illumina HiSeq 2000**  
*Sequencing by Synthesis*

>60Gbp / day



Metzker (2010) Nature Reviews Genetics 11:31-46  
<http://www.youtube.com/watch?v=l99aKKHcxC4>

# Cost per Genome



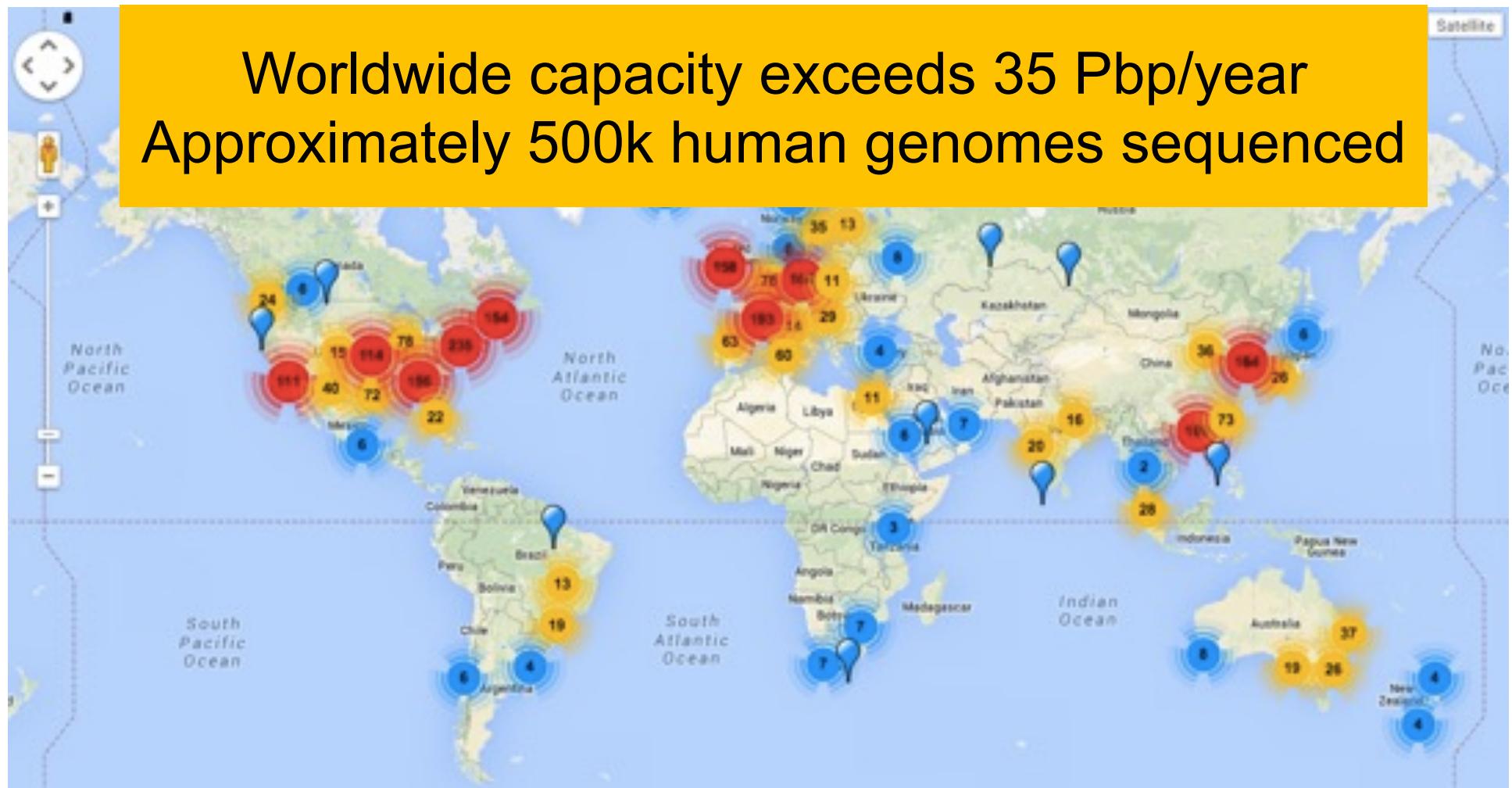
# HiSeq X Ten



320 genomes per week / 18,000 genomes per year  
\$1000 per genome / ~\$10 M per instrument

# Sequencing Centers

Worldwide capacity exceeds 35 Pbp/year  
Approximately 500k human genomes sequenced



**Next Generation Genomics: World Map of High-throughput Sequencers**  
<http://omicsmaps.com>

# How much is a petabyte?

Unit	Size
Byte	1
Kilobyte	1,000
Megabyte	1,000,000
Gigabyte	1,000,000,000
Terabyte	1,000,000,000,000
Petabyte	1,000,000,000,000,000

\*Technically a kilobyte is  $2^{10}$  and a petabyte is  $2^{50}$

# How much is a petabyte?



100 GB / Genome  
4.7GB / DVD  
~20 DVDs / Genome

X

10,000 Genomes

=

1PB Data  
200,000 DVDs



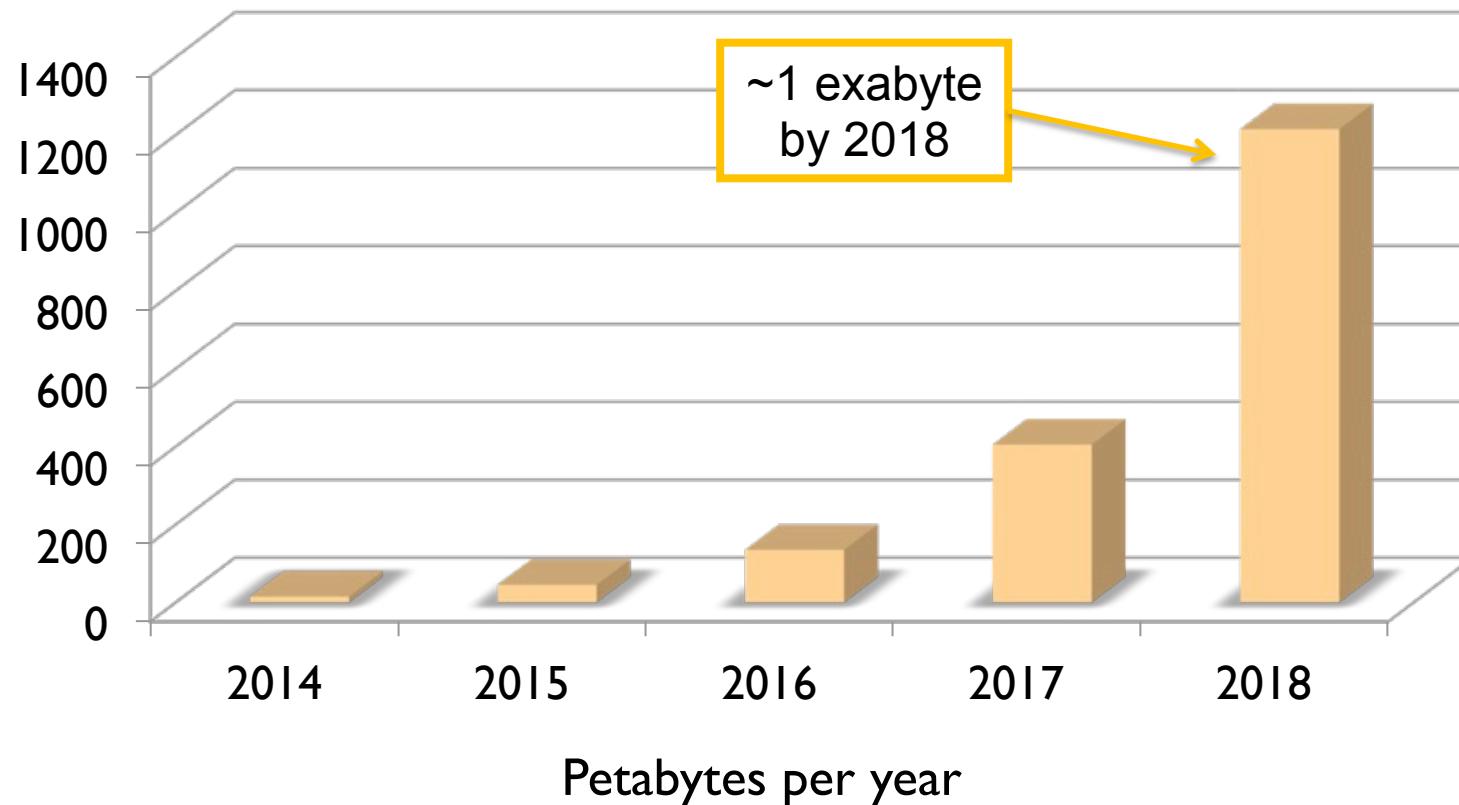
787 feet of DVDs  
~1/6 of a mile tall



500 2 TB drives  
\$500k

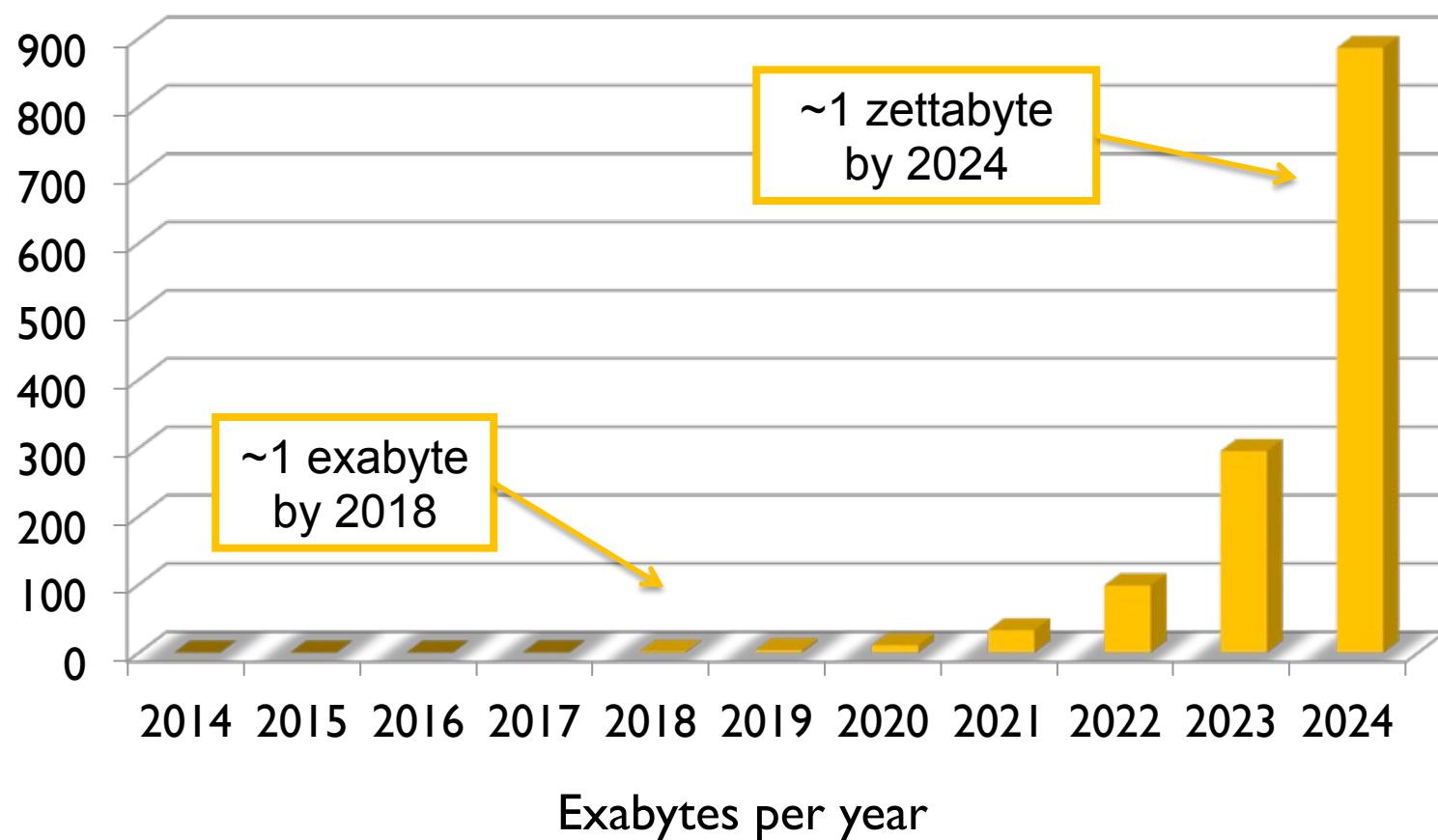
# DNA Data Tsunami

*Current world-wide sequencing capacity is growing at ~3x per year!*



# DNA Data Tsunami

*Current world-wide sequencing capacity is growing at ~3x per year!*



# How much is a zettabyte?

Unit	Size
Byte	1
Kilobyte	1,000
Megabyte	1,000,000
Gigabyte	1,000,000,000
Terabyte	1,000,000,000,000
Petabyte	1,000,000,000,000,000
Exabyte	1,000,000,000,000,000,000
Zettabyte	1,000,000,000,000,000,000,000

# How much is a zettabyte?



100 GB / Genome  
4.7GB / DVD  
~20 DVDs / Genome

X

10,000,000,000 Genomes

=

1ZB Data  
200,000,000,000 DVDs



150,000 miles of DVDs  
~ ½ distance to moon

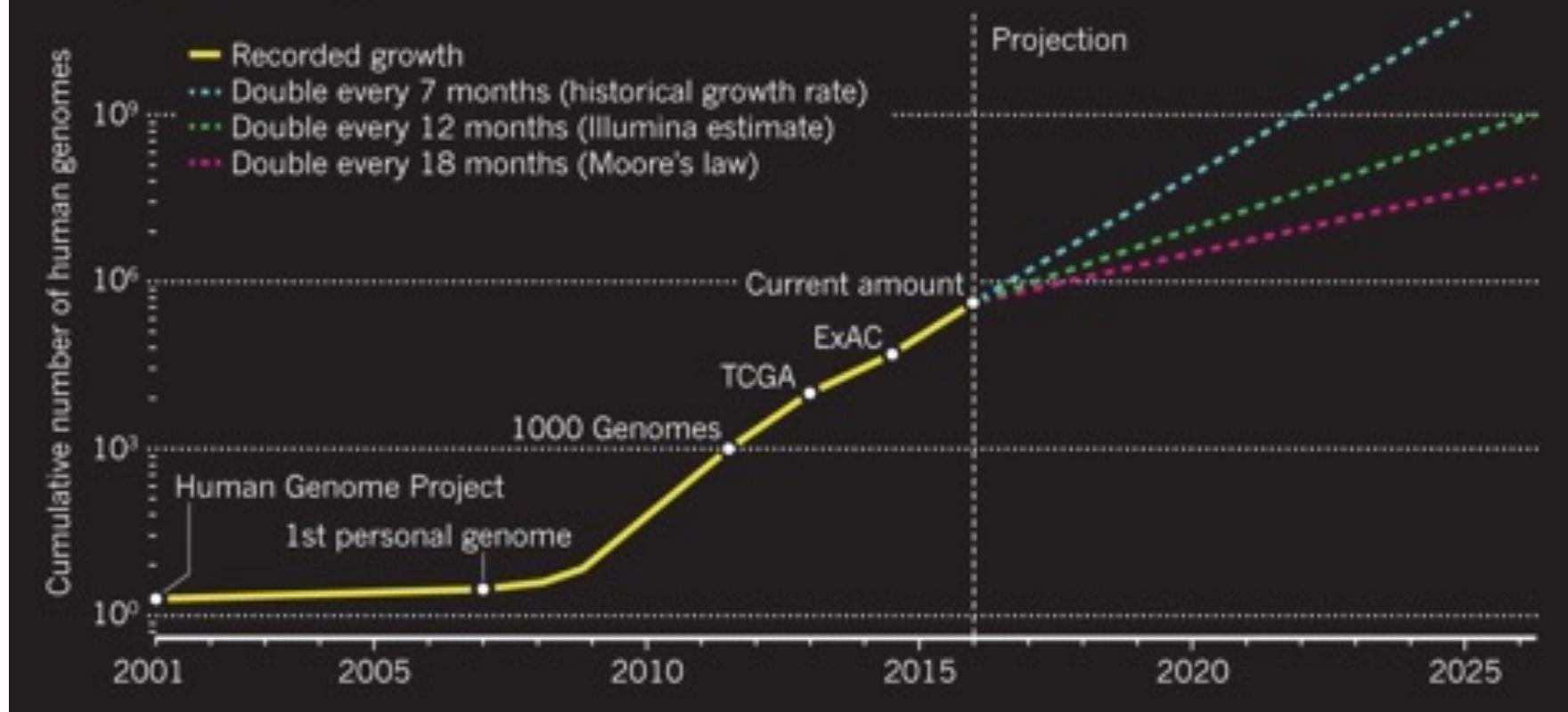


Both currently ~100Pb  
And growing exponentially

# Sequencing Capacity

## DNA SEQUENCING SOARS

Human genomes are being sequenced at an ever-increasing rate. The 1000 Genomes Project has aggregated hundreds of genomes; The Cancer Genome Atlas (TCGA) has gathered several thousand; and the Exome Aggregation Consortium (ExAC) has sequenced more than 60,000 exomes. Dotted lines show three possible future growth curves.



## Big Data: Astronomical or Genomical?

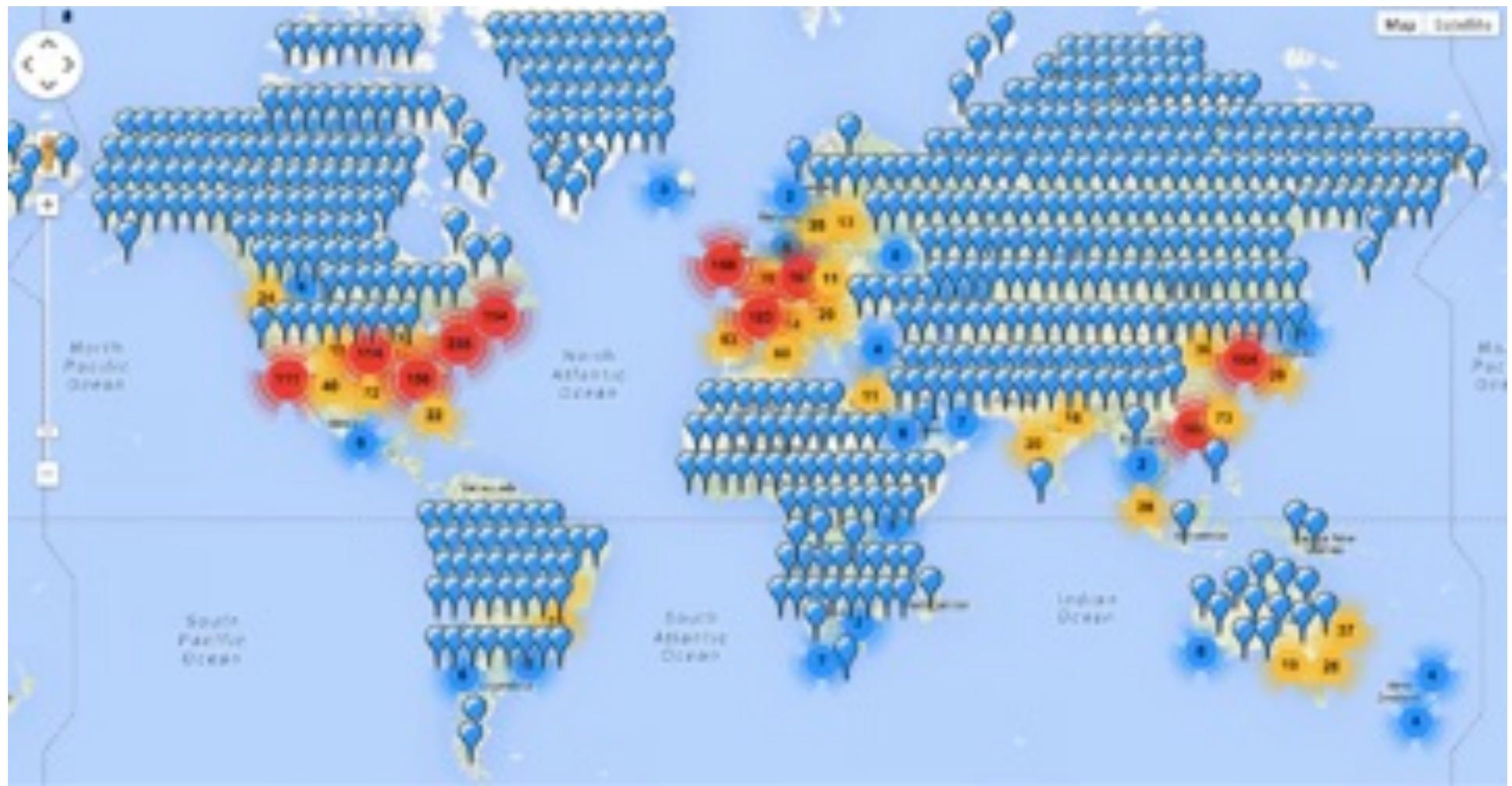
Stephens, Z, et al. (2015) PLOS Biology DOI: [10.1371/journal.pbio.1002195](https://doi.org/10.1371/journal.pbio.1002195)

# Sequencing Centers 2017



**Next Generation Genomics: World Map of High-throughput Sequencers**  
<http://omicsmaps.com>

# Sequencing Centers 2027



**Next Generation Genomics: World Map of High-throughput Sequencers**  
<http://omicsmaps.com>

# Biological Sensor Network



Oxford Nanopore



DC Metro via the LA Times

***The rise of a digital immune system***  
Schatz, MC, Phillippy, AM (2012) GigaScience 1:4

# Biological Sensor Network



@JasonWilliamsNY



Aspyn @ CSH High School

***The rise of a digital immune system***  
Schatz, MC, Phillippy, AM (2012) GigaScience 1:4

# Data Production & Collection

**Expect massive growth to sequencing and other biological sensor data over the next 10 years**

- Exascale biology is certain, zettascale on the horizon
- Compression helps, but need to aggressively throw out data
- Requires careful consideration of the “preciousness” of the sample

**Major data producers concentrated in hospitals, universities, agricultural companies, research institutes**

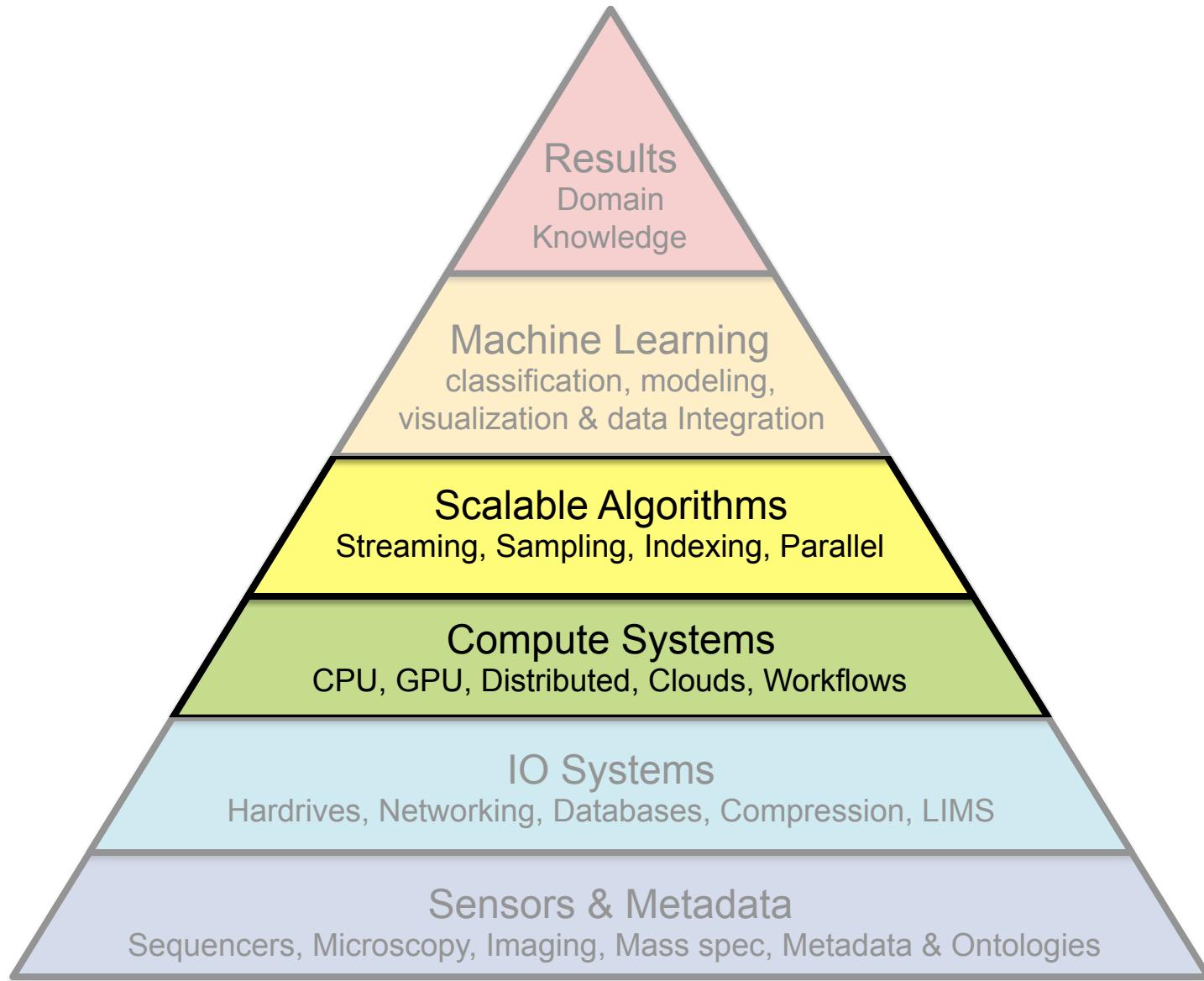
- Major efforts in human health and disease, agriculture, bioenergy
- Genomic information coupled with medical records and other medical data

**But also widely distributed mobile sensors**

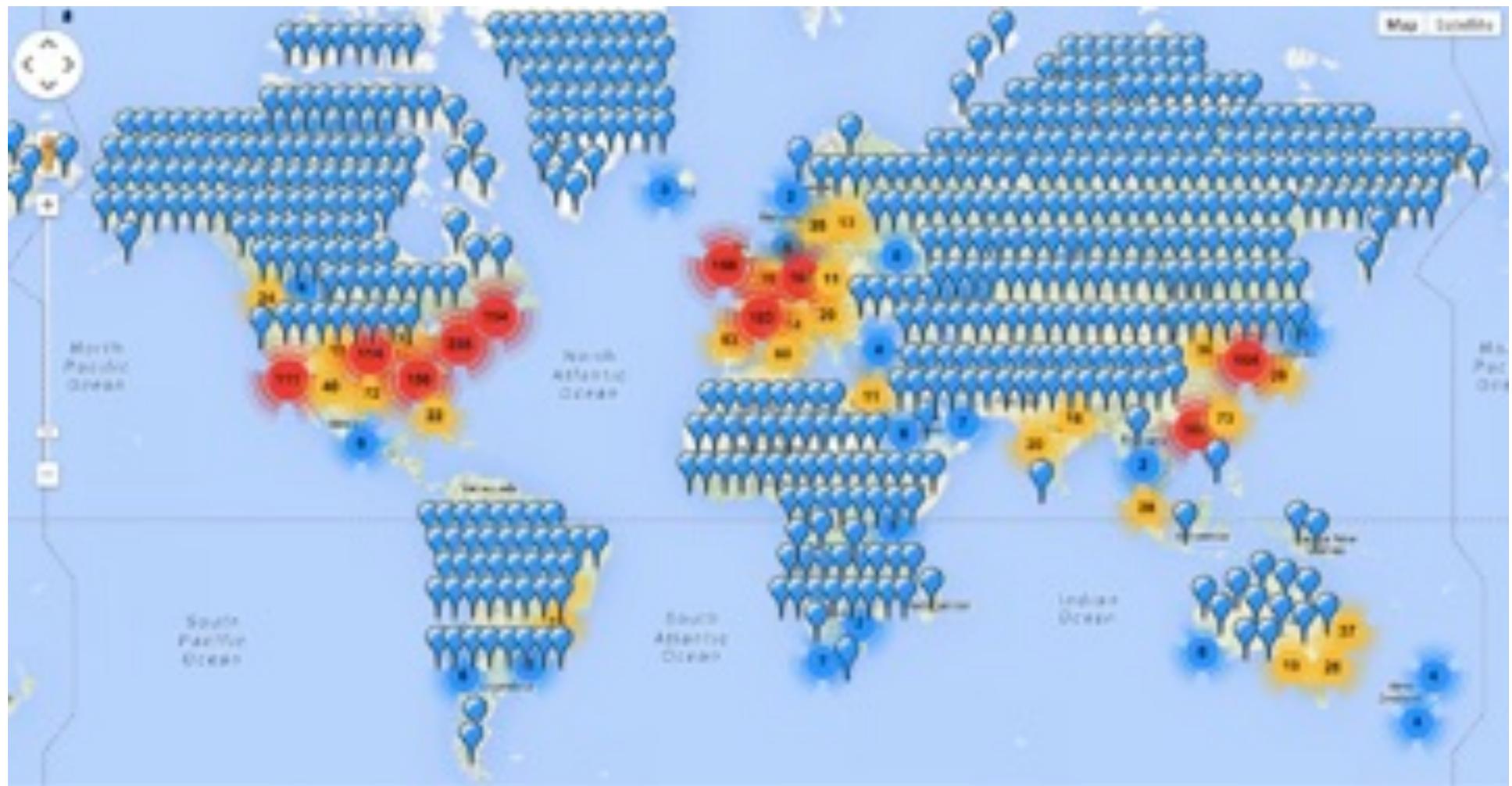
- Schools, offices, sports arenas, transportation centers, farms & food distribution centers
- Monitoring and surveillance, as ubiquitous as weather stations
- The rise of a digital immune system?



# Comparative Genomics Technologies



# Sequencing Centers 2027



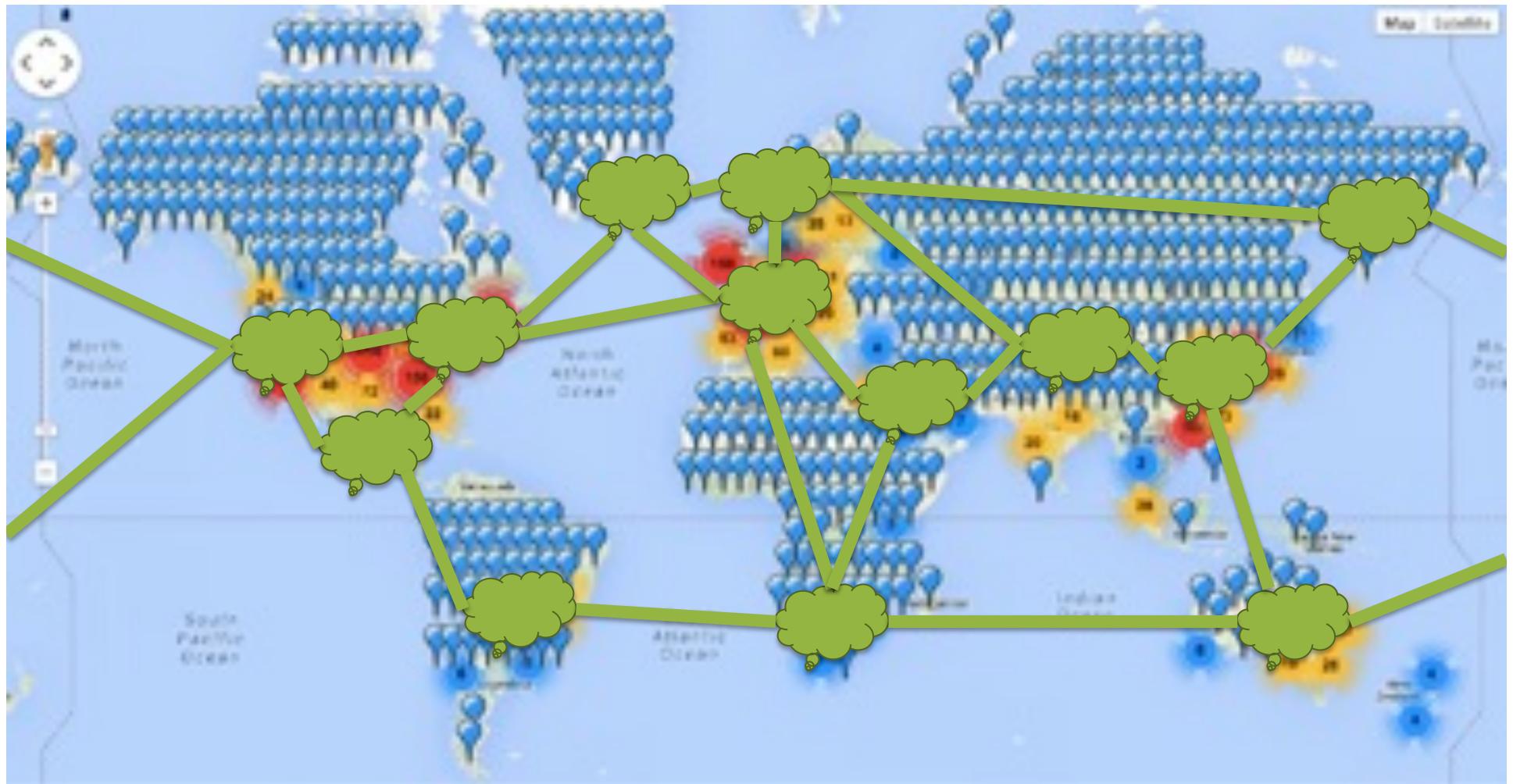
# Informatics Centers 2027



***The DNA Data Deluge***

Schatz, MC and Langmead, B (2013) *IEEE Spectrum*. July, 2013

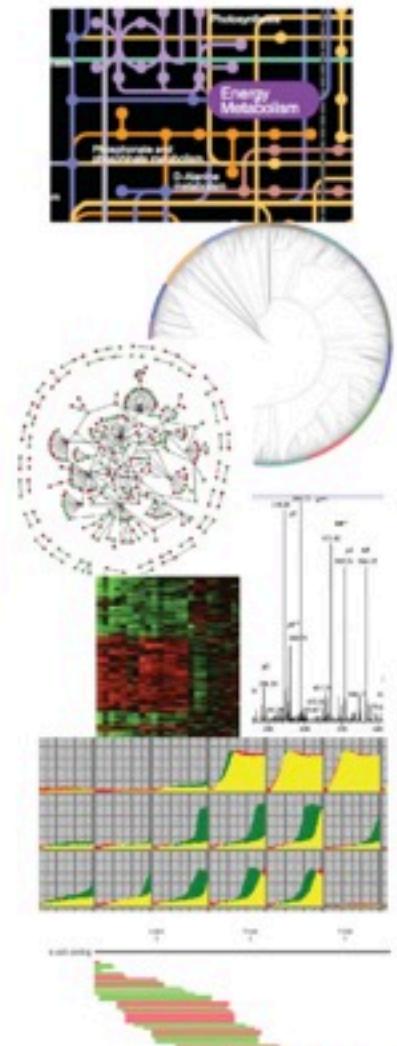
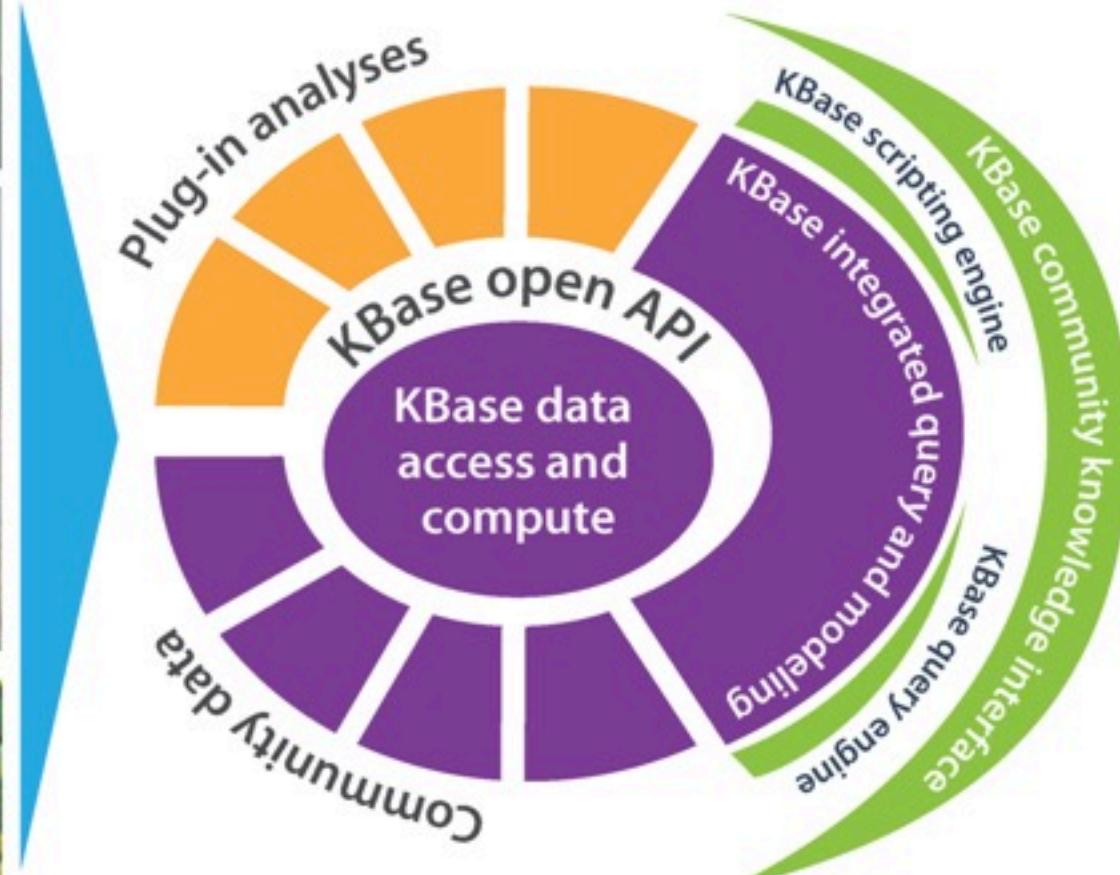
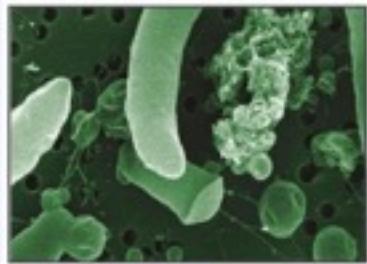
# Informatics Centers 2017



## ***The DNA Data Deluge***

Schatz, MC and Langmead, B (2013) *IEEE Spectrum*. July, 2013

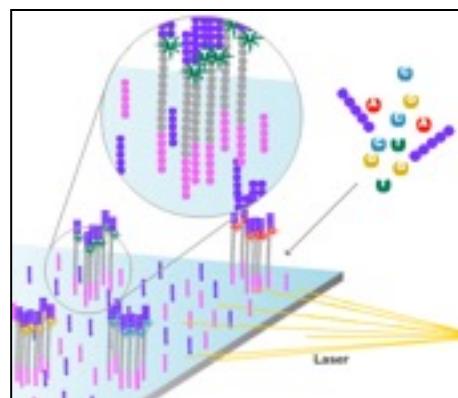
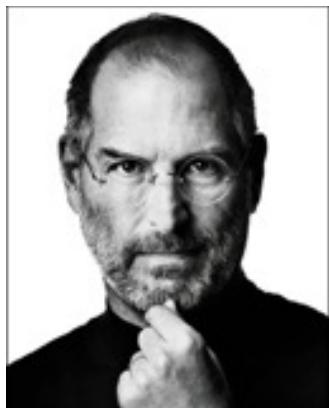
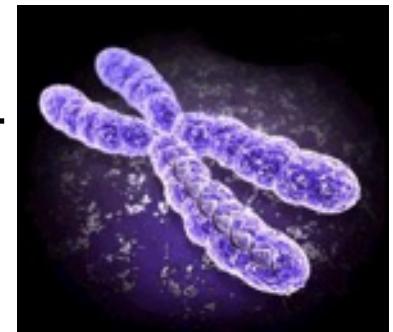
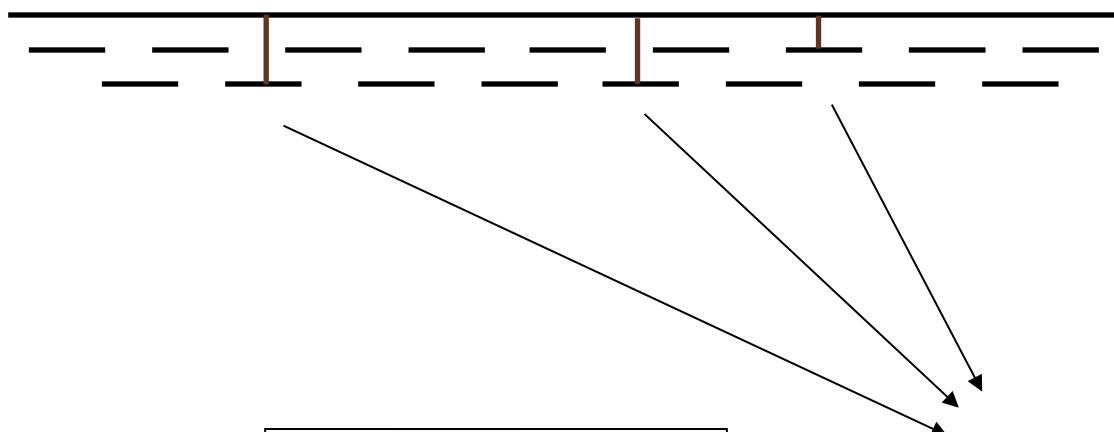
# DOE Systems Biology Knowledgebase



<http://kbase.us>: Predictive Biology in Microbes, Plants, and Meta-communities

# Personal Genomics

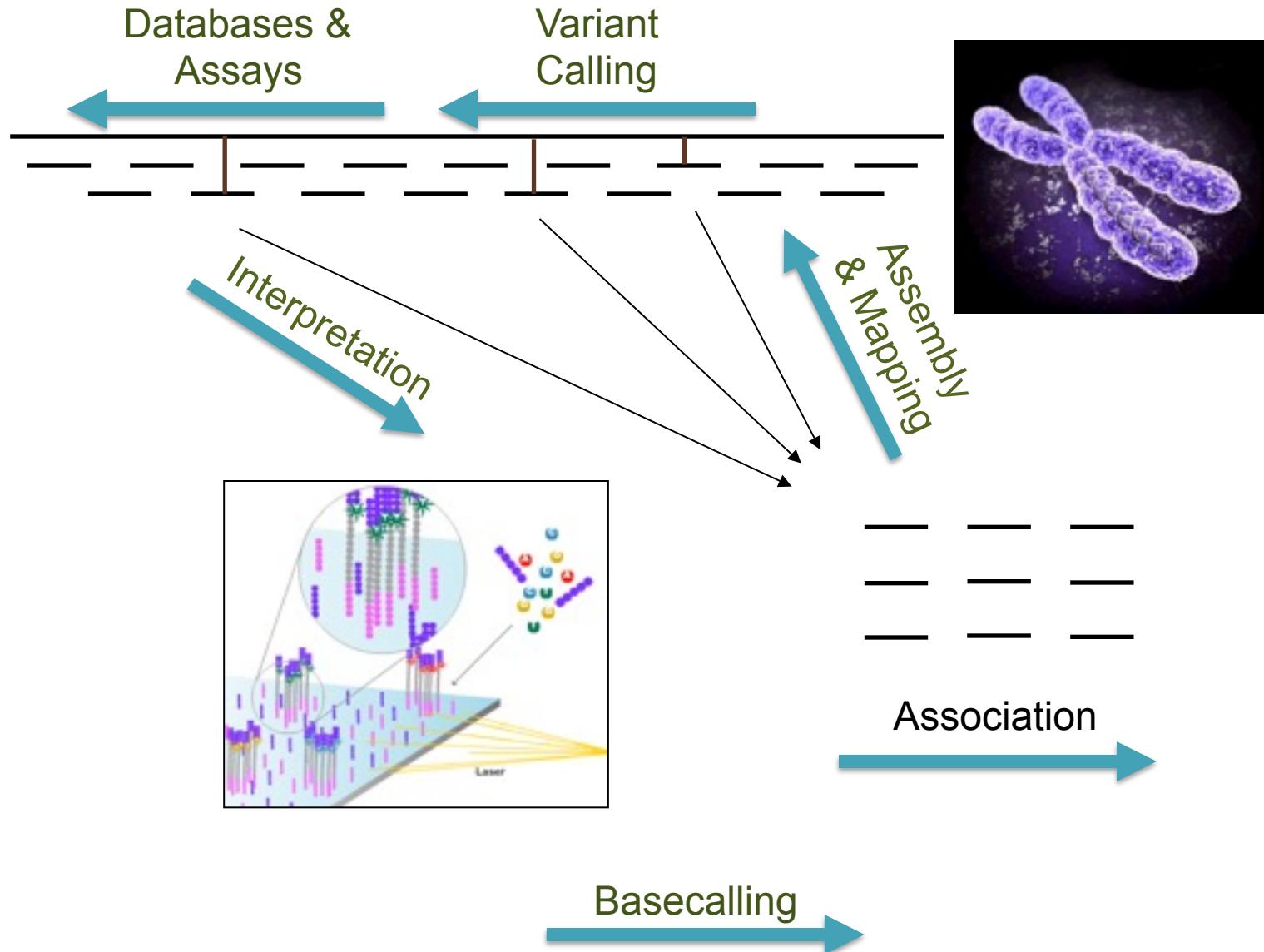
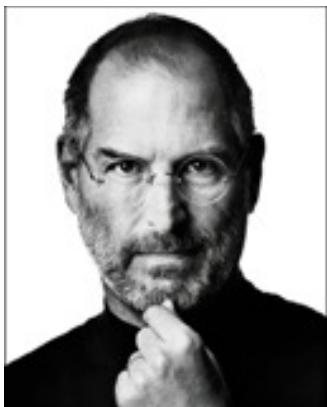
How does your genome compare to the reference?



Heart Disease  
Cancer  
Creates magical technology

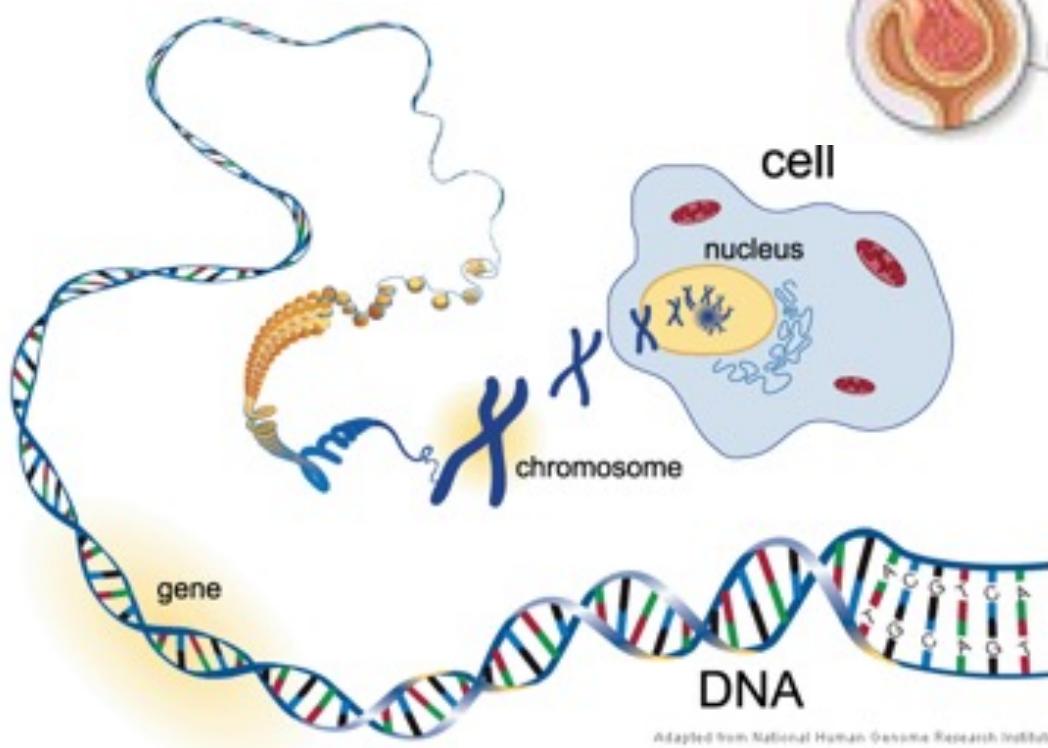
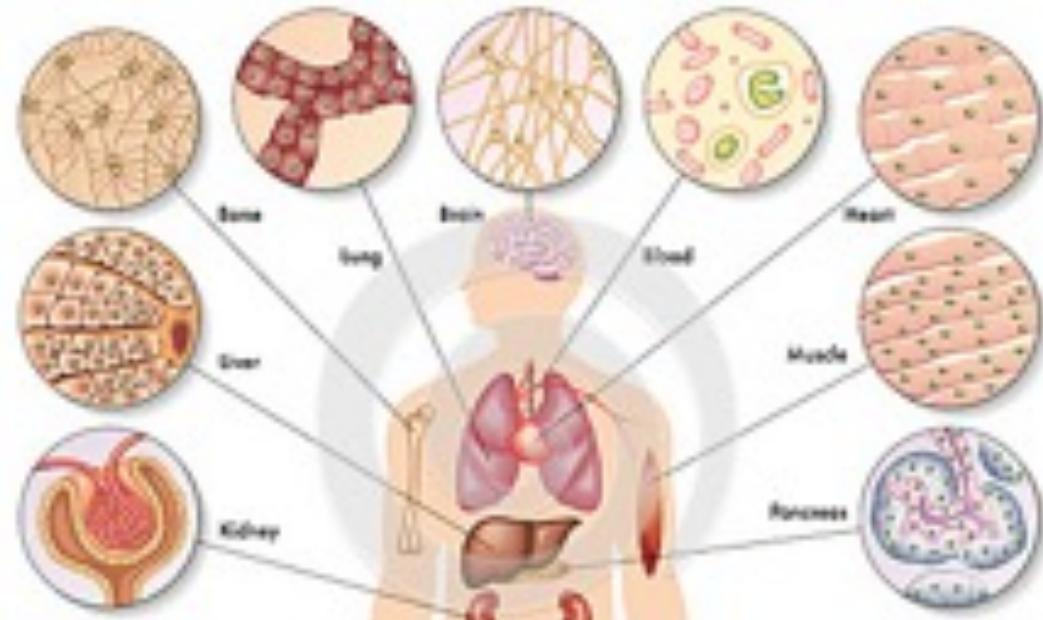
# Personal Genomics

How does your genome compare to the reference?

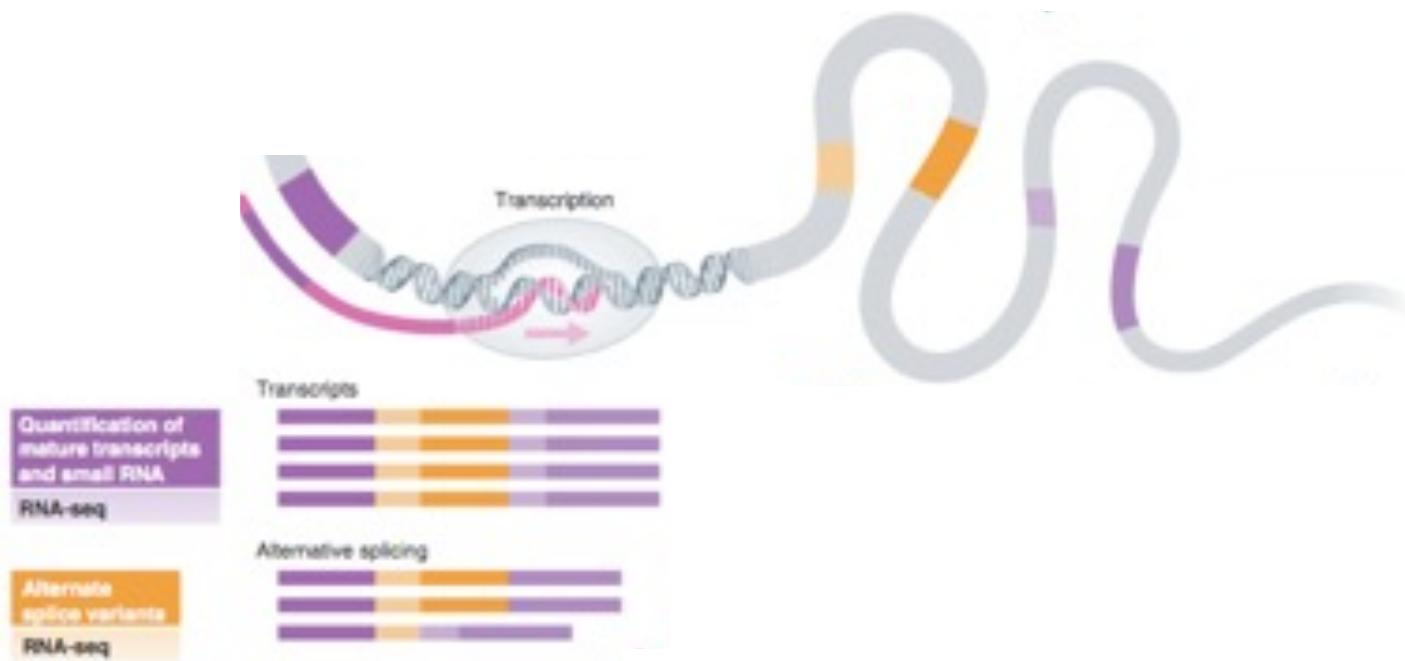


# Cells & DNA

Each cell of your body contains an exact copy of your 3 billion base pair genome.



Your specific nucleotide sequence encodes the genetic program for your cells and ultimately your traits

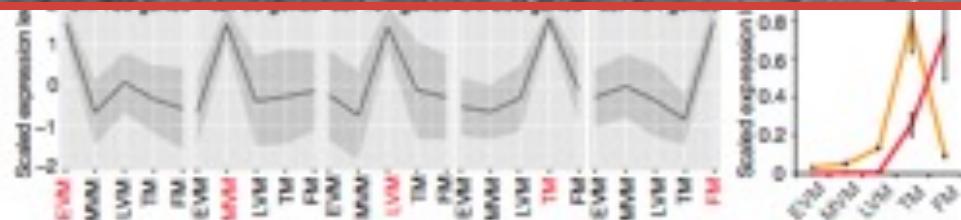


# Rate of meristem maturation determines inflorescence architecture in tomato

Soon Ju Park<sup>1</sup>, Ke Jiang<sup>1</sup>, Michael C. Schatz, and Zachary B. Lippman<sup>2</sup>

Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724

Edited by Maarten Koornneef, Wageningen University and Research Centre, Cologne, Germany, and approved November 28, 2011 (received for review September 12, 2011)



- When those genes are delayed or interrupted, tomato mutants take on very different branching patterns.

# Compute & Algorithmic Challenges

**Expect to see many dozens of major informatics centers that consolidate regional / topical information**

- Clouds for Cancer, Autism, Heart Disease, etc
- Plus many smaller warehouses down to individuals
- Move the code to the data

**Parallel hardware and algorithms are required**

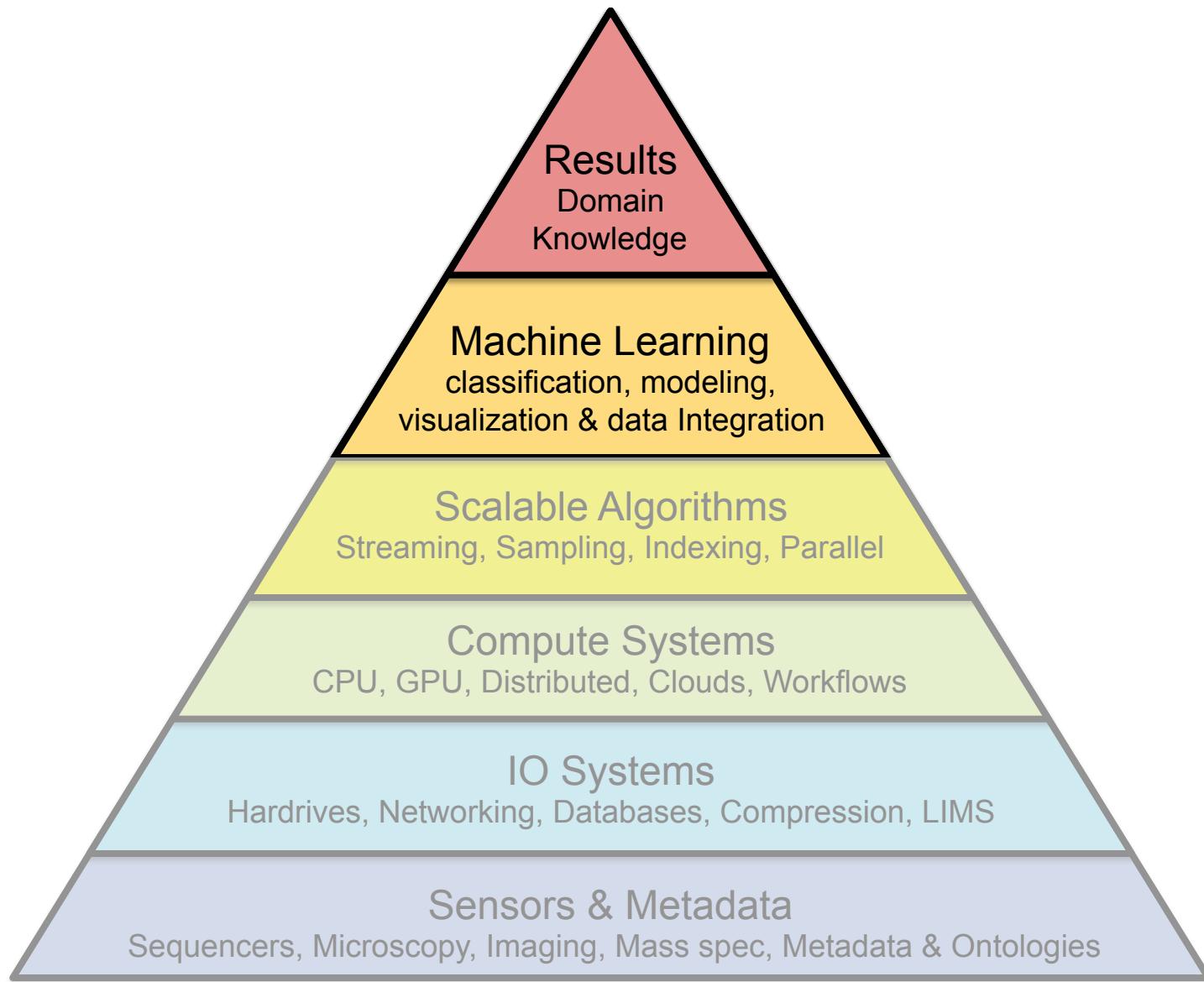
- Expect to see >1000 cores in a single computer
- Compute & IO needs to be considered together
- Rewriting efficient parallel software is complex and expensive

**Applications will shift from individuals to populations**

- Read mapping & assembly fade out
- Population analysis and time series analysis fade in
- Need for network analysis, probabilistic techniques



# Comparative Genomics Technologies



# Genetic Basis of Autism Spectrum Disorders



## ***Complex disorders of brain development***

- Characterized by difficulties in social interaction, verbal and nonverbal communication and repetitive behaviors.
- Have their roots in very early brain development, and the most obvious signs of autism and symptoms of autism tend to emerge between 2 and 3 years of age.

## ***U.S. CDC identify around 1 in 68 American children as on the autism spectrum***

- Ten-fold increase in prevalence in 40 years, only partly explained by improved diagnosis and awareness.
- Studies also show that autism is four to five times more common among boys than girls.
- Specific causes remain elusive

## **What is Autism?**

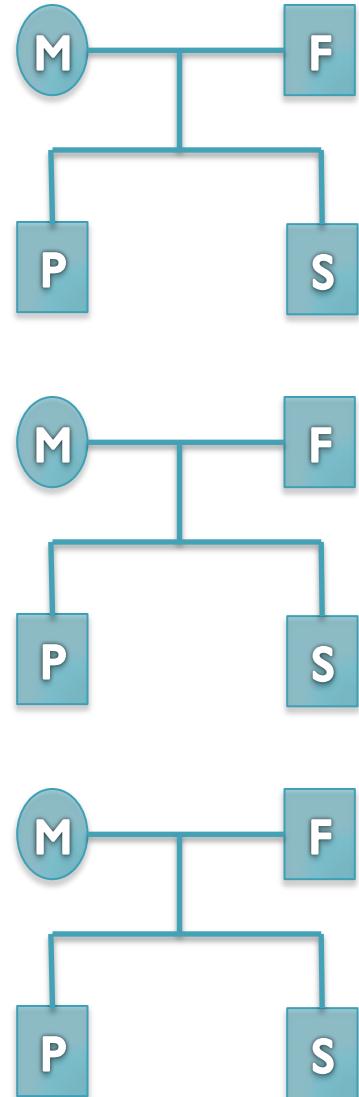
<http://www.autismspeaks.org/what-autism>

# Searching for the genetic risk factors

## Search Strategy

- Thousands of families identified from a dozen hospitals around the United States
- Large scale genome sequencing of “simplex” families: mother, father, affected child, unaffected sibling
- Unaffected siblings provide a natural control for environmental factors

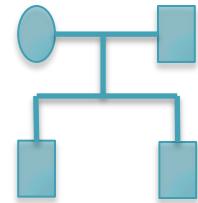
***Are there any genetic variants present in affected children, that are not in their parents or unaffected siblings?***



# De novo mutation discovery and validation

## De novo mutations:

Sequences not inherited from your parents.



**Reference:** ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Father(1): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Father(2): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Mother(1): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Mother(2): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Sibling(1): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Sibling(2): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Proband(1): ...TCAAATCCTTTAATAAAAGAAGAGCTGACA...

Proband(2): ...TCAAATCCTTTAAT\*\*\*\*AAGAGCTGACA...

4bp heterozygous deletion at chr15:9352406 | CHD2

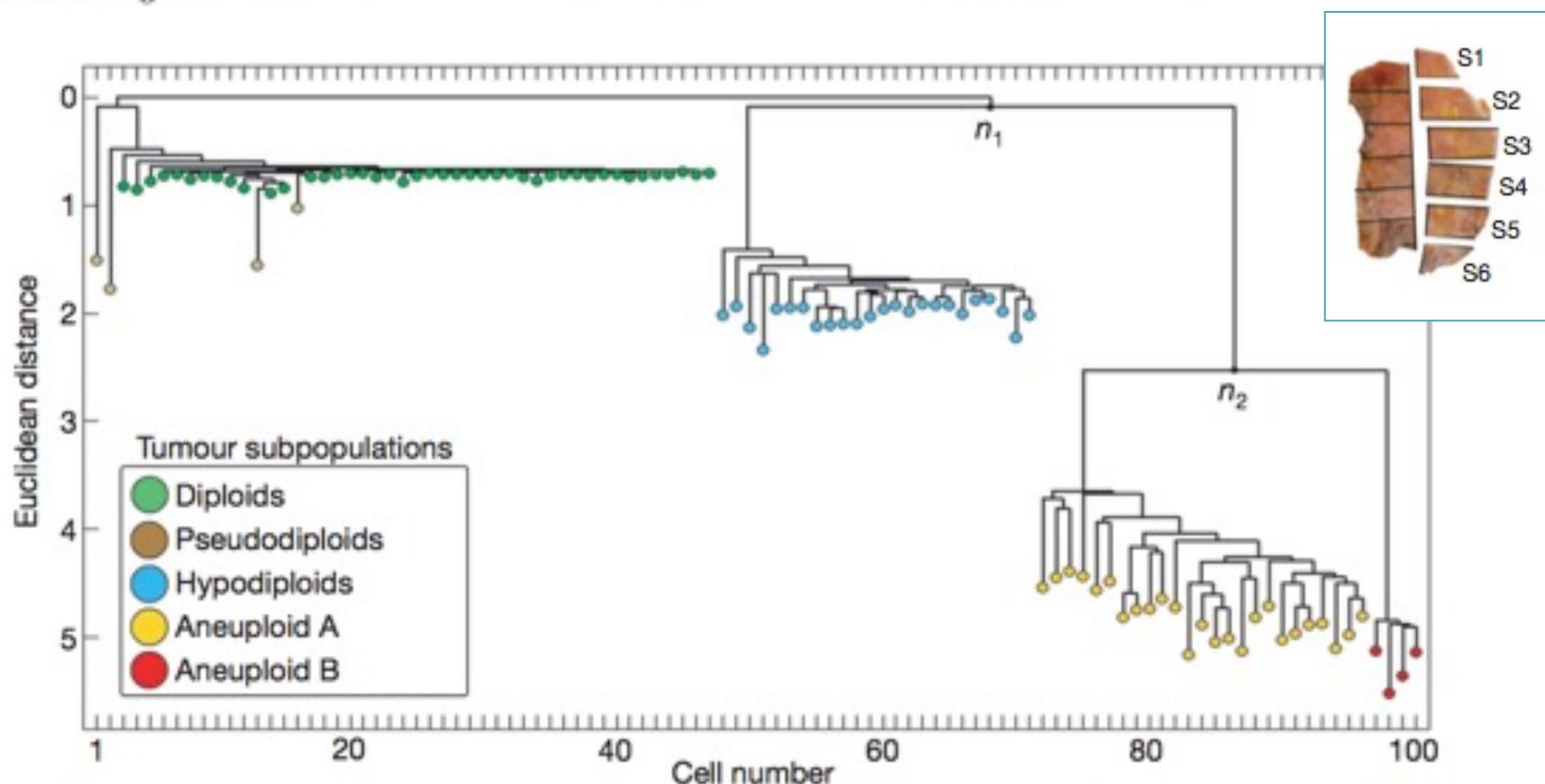
# De novo Genetics of Autism

- In 593 family quads so far, we see significant enrichment in de novo ***likely gene killers*** in the autistic kids
  - Overall rate basically 1:1
  - 2:1 enrichment in nonsense mutations
  - 2:1 enrichment in frameshift indels
  - 4:1 enrichment in splice-site mutations
  - Most de novo originate in the paternal line in an age-dependent manner (56:18 of the mutations that we could determine)
- Observe strong overlap with the 842 genes known to be associated with fragile X protein FMRP
  - Related to neuron development and synaptic plasticity
  - Also strong overlap with chromatin remodelers

Accurate de novo and transmitted indel detection in exome-capture data using microassembly.  
Narzisi et al (2014) Nature Methods doi:10.1038/nmeth.3069

# Tumour evolution inferred by single-cell sequencing

Nicholas Navin<sup>1,2</sup>, Jude Kendall<sup>1</sup>, Jennifer Trogan<sup>1</sup>, Peter Andrews<sup>1</sup>, Linda Rodgers<sup>1</sup>, Jeanne McIndoo<sup>1</sup>, Kerry Cook<sup>1</sup>, Asya Stepansky<sup>1</sup>, Dan Levy<sup>1</sup>, Diane Esposito<sup>1</sup>, Lakshmi Muthuswamy<sup>3</sup>, Alex Krasnitz<sup>1</sup>, W. Richard McCombie<sup>1</sup>, James Hicks<sup>1</sup> & Michael Wigler<sup>1</sup>



# What makes us human?

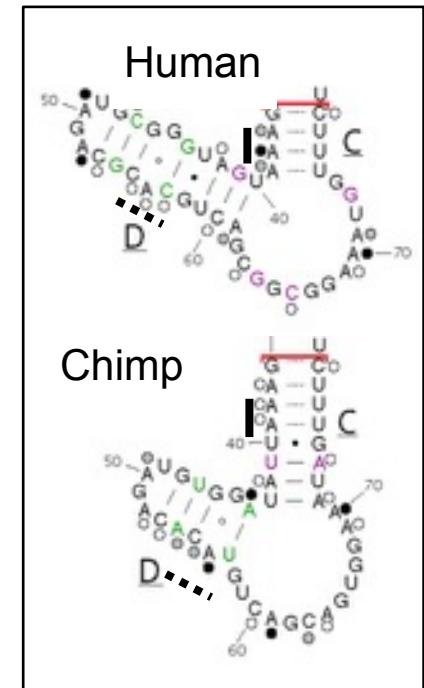
## “Human Accelerated Regions”



human TCGATGGCGTAGACCCACGTCAGCGGGAAATGGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAGTTCAAAATGA  
chimp TTATAGCTGTAGACCATGTCAGCAGGGAAATAGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAAATTTCAAAATTA  
dog TTATAGCTGTAGACCATGTCAGCGCGCAACAGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAAATTTCAAAATTA  
mouse TTATAGCTGTAGACCATGTCAGCGCTGGAAATGGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAAATTTCAAAATTA  
rat TTATAGCTGTAGACCATGTCAGCAGGGAAATGGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAAATTTCAAAATTA  
chicken TTATAGCTGTAGACCATGTCAGCAGTAGAACAGTTCTATCAAATCAAAGTCTTAGAGATTTCCCTCAAAATTTCAAAATTA

Systematic scan of recent human evolution identified the gene *HAR1F* as the most dramatic “human accelerated region”.

Follow up analysis found it was specifically expressed in Cajal-Retzius neurons in the human brain from 6 to 19 gestational weeks.



(Pollard et al., *Nature*, 2006)

# Genetic Privacy

## Identifying Personal Genomes by Surname Inference

Melissa Gymrek,<sup>1,2,3,4</sup> Amy L. McGuire,<sup>5</sup> David Golan,<sup>6</sup> Evan Mandel,<sup>7,8,9</sup> and Ryan Kirchner<sup>10</sup>

By combining other pieces of demographic information, such as date and place of birth, they fully exposed the identity of their biological fathers. Lunshof *et al.* (*10*) were the first to speculate that this technique could expose the full identity of participants in sequencing projects. Giachieri (*11*)

Sharing sequencing data sets without identity can be recovered. Here, we report that surnames can be recovered from short tandem repeats on the Y chromosome (Y-STRs) and mitochondrial DNA (mtDNA). We show that a combination of a surname and mtDNA can be used to triangulate the identity of the individual. This relies on free, publicly accessible Internet resources for genome identification for U.S. males. We further demonstrate that it is possible to predict with high probability the identities of multiple individuals.

Surnames are paternally inherited in most human societies, resulting in their rapid segregation with Y-chromosome haplotypes (*1–5*). Based on this observation, multiple genealogical genealogy companies offer services to reunite distant patrilineal relatives by genotyping a few dozen Y-STRs.

<sup>1</sup>Whitehead Institute for Biomedical Research, 9 Cambridge Center, Cambridge, MA 02142, USA. <sup>2</sup>Harvard–Massachusetts Institute of Technology (MIT) Division of Health Sciences and Technology, MIT, Cambridge, MA 02139, USA. <sup>3</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA. <sup>4</sup>Department of Molecular Medicine and Diabetes Unit, Massachusetts General Hospital, Boston, MA 02114, USA. <sup>5</sup>Center for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, TX 77030, USA. <sup>6</sup>Department of Statistics and Operations Research, Tel Aviv University, Tel Aviv 69978, Israel. <sup>7</sup>School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel. <sup>8</sup>Department of Molecular Microbiology and Biotechnology, Tel-Aviv University, Tel Aviv 69978, Israel. <sup>9</sup>The International Computer Science Institute, Berkeley, CA 94704, USA.

\*To whom correspondence should be addressed. E-mail: yanki@mit.edu

PNAS

www.pnas.org

## Predicting Social Security numbers from public data

Alessandro Acquisti<sup>1</sup> and Ralph Gross<sup>2</sup>

Carnegie Mellon University, Pittsburgh, PA 15213

Communicated by Stephen E. Fienberg, Carnegie Mellon University, Pittsburgh, PA, May 5, 2009 (received for review January 18, 2009)

Information about an individual's place and date of birth can be exploited to predict his or her Social Security number (SSN). Using only publicly available information, we observed a correlation between individuals' SSNs and their birth data and found that for younger cohorts the correlation allows statistical inference of private SSNs. The inferences are made possible by the public availability of the Social Security Administration's Death Master File and the widespread accessibility of personal information from multiple sources, such as data brokers or profiles on social networking sites. Our results highlight the unexpected privacy consequences of the complex interactions among multiple data sources in modern information economies and quantify privacy

number (SN). The SSA openly provides information about the process through which ANs, GNs, and SNs are issued (*1*). ANs are currently assigned based on the zipcode of the mailing address provided in the SSN application form [RM00201.030] (*1*). Low-population states and certain U.S. possessions are allocated 1 AN each, whereas other states are allocated sets of ANs (for instance, an individual applying from a zipcode within New York state may be assigned any of 85 possible first 3 SSN digits). Within each SSA area, GNs are assigned in a precise but nonconsecutive order between 01 and 99 [RM00201.030] (*1*). Both the sets of ANs assigned to different states and the sequence of GNs are publicly available (see [www.socialsecurity.gov/employer/](http://www.socialsecurity.gov/employer/)).

**EXTRAPOLATING TO THE U.S. LIVING POPULATION, THIS WOULD IMPLY THE POTENTIAL IDENTIFICATION OF MILLIONS OF SSNS FOR INDIVIDUALS WHOSE BIRTH DATA WERE AVAILABLE. SUCH FINDINGS HIGHLIGHT THE HIDDEN PRIVACY COSTS OF WIDESPREAD INFORMATION DISSEMINATION AND THE COMPLEX INTERACTIONS AMONG MULTIPLE DATA SOURCES IN MODERN INFORMATION ECONOMIES (11), UNDERSCORING THE ROLE OF PUBLIC RECORDS AS BREEDER DOCUMENTS (12) OF MORE SENSITIVE DATA.**

### Methods

have already left the barn: We demonstrate that it is possible to predict and day of application. Empirical observation of SSA's policies—

SEE COMMENTARY

# Learning and Translation

## Tremendous power from data aggregation

- Observe the dynamics of biological systems
- Breakthroughs in medicine and biology of profound significance

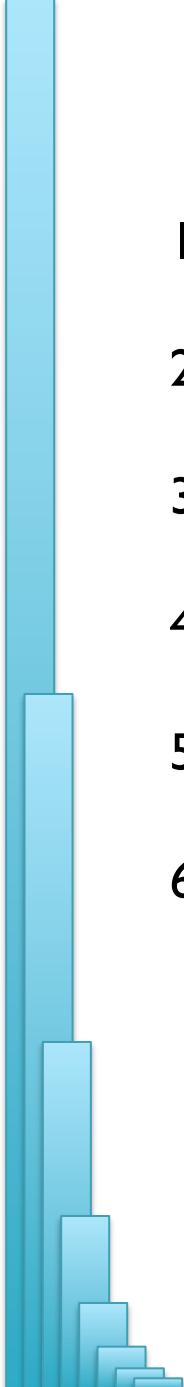
## Be mindful of the risks

- The potential for over-fitting grows with the complexity of the data, statistical significance is a statement about the sample size
- Reproducible workflows, APIs are a must
- Caution is prudent for personal data

## The foundations of biology will continue to be observation, experimentation, and interpretation

- Technology will continue to push the frontier
- Feedback loop from the results of one project into experimental design for the next





# Next Steps

1. Reflect on the magic and power of DNA 😊
2. Check out the course webpage
3. Register on Piazza
4. Set up Dropbox for yourself!
5. Set up Linux, set up Virtual Machine
6. Get comfortable on the command line



**Welcome to Applied Comparative Genomics**

<https://github.com/schatzlab/appliedgenomics>

**Questions?**