

# Big Data Meets DNA

How Biological Data Science is improving our health, foods, and energy needs

Michael Schatz

June 18, 2014

CSHL Public Lecture Series



# DNA: The secret of life



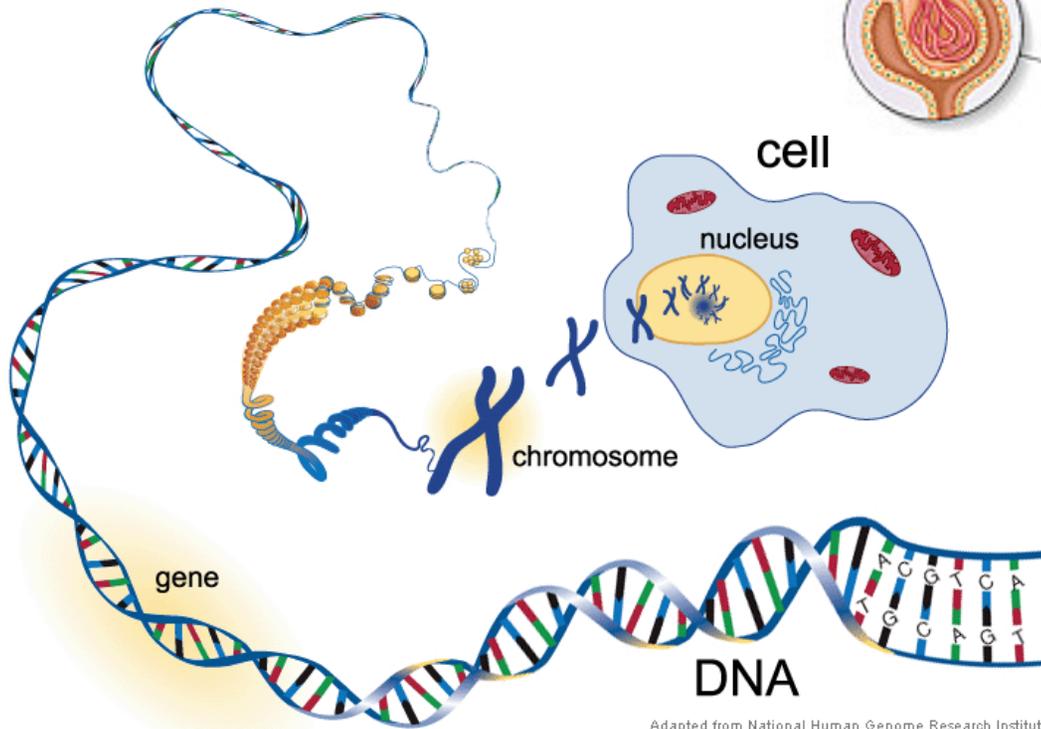
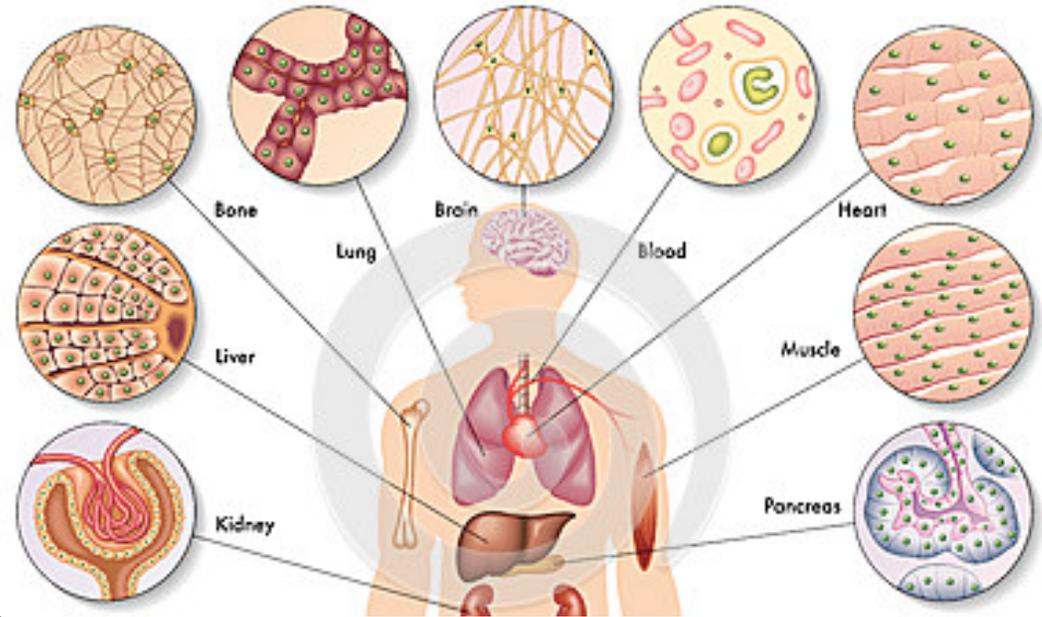
***Your DNA, along with your environment and experiences, shapes who you are***

- 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

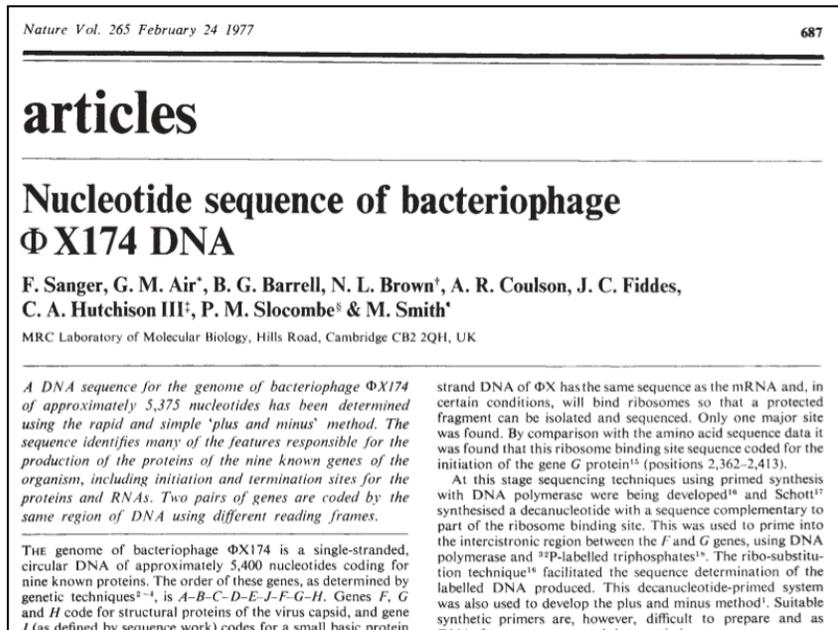
# Cells & DNA

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



Sanger et al. (1977) Nature  
1<sup>st</sup> Complete Organism  
Bacteriophage  $\phi$  X174; 5375 bp  
**Awarded Nobel Prize in 1980**



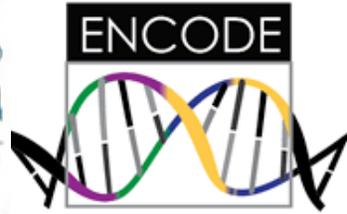
Radioactive Chain Termination  
5000bp / week / person  
<http://en.wikipedia.org/wiki/File:Sequencing.jpg>  
<http://www.answers.com/topic/automated-sequencer>

# 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.

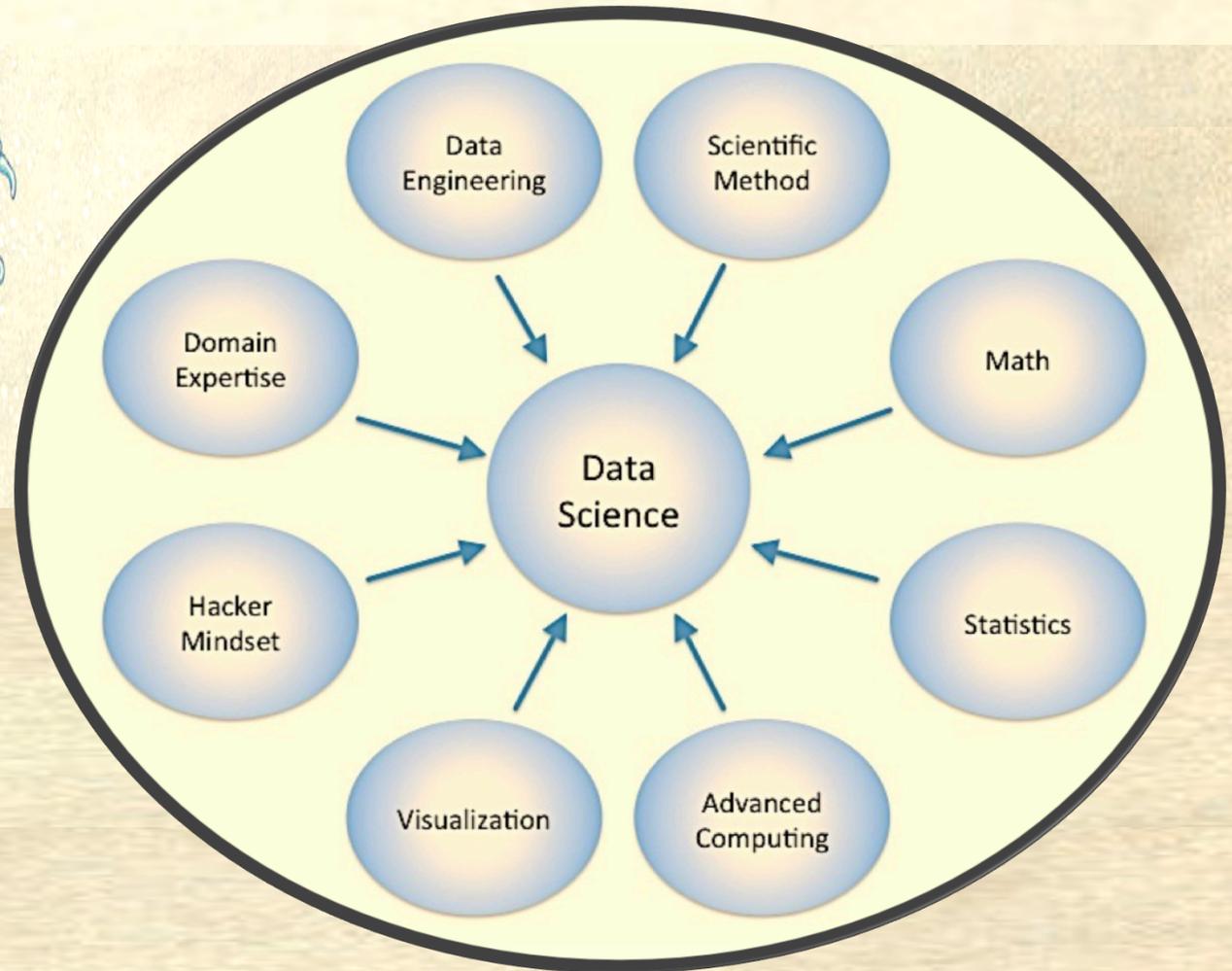
***What software and systems will?***

***And who will create them?***

- ***Plus hundreds and hundreds more***

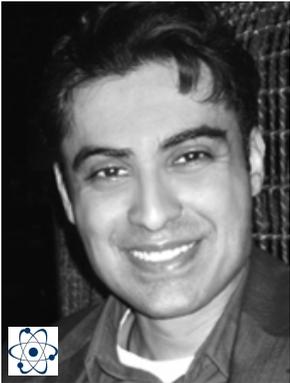


# Who is a Data Scientist?



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

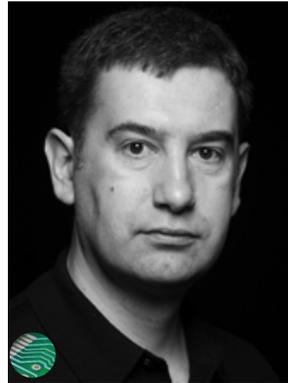
# CSHL Quantitative Biology



**Mickey Atwal**  
Population Genetics  
Cancer, Fertility



**Molly Hammel**  
Gene regulatory  
Networks, RNA Biology



**Ivan Iossifov**  
Human Genetics  
Molecular Networks



**Justin Kinney**  
Biophysics  
Machine learning



**Alexei Koulakov**  
Neurobiology  
Cortical design, Memory



**Alex Krasnitz**  
Genomics of Cancer  
Machine Learning



**Dan Levy**  
Human Genetics  
Phylogenetics, CNVs



**Partha Mitra**  
Neuroscience  
Neural Imaging & Disease

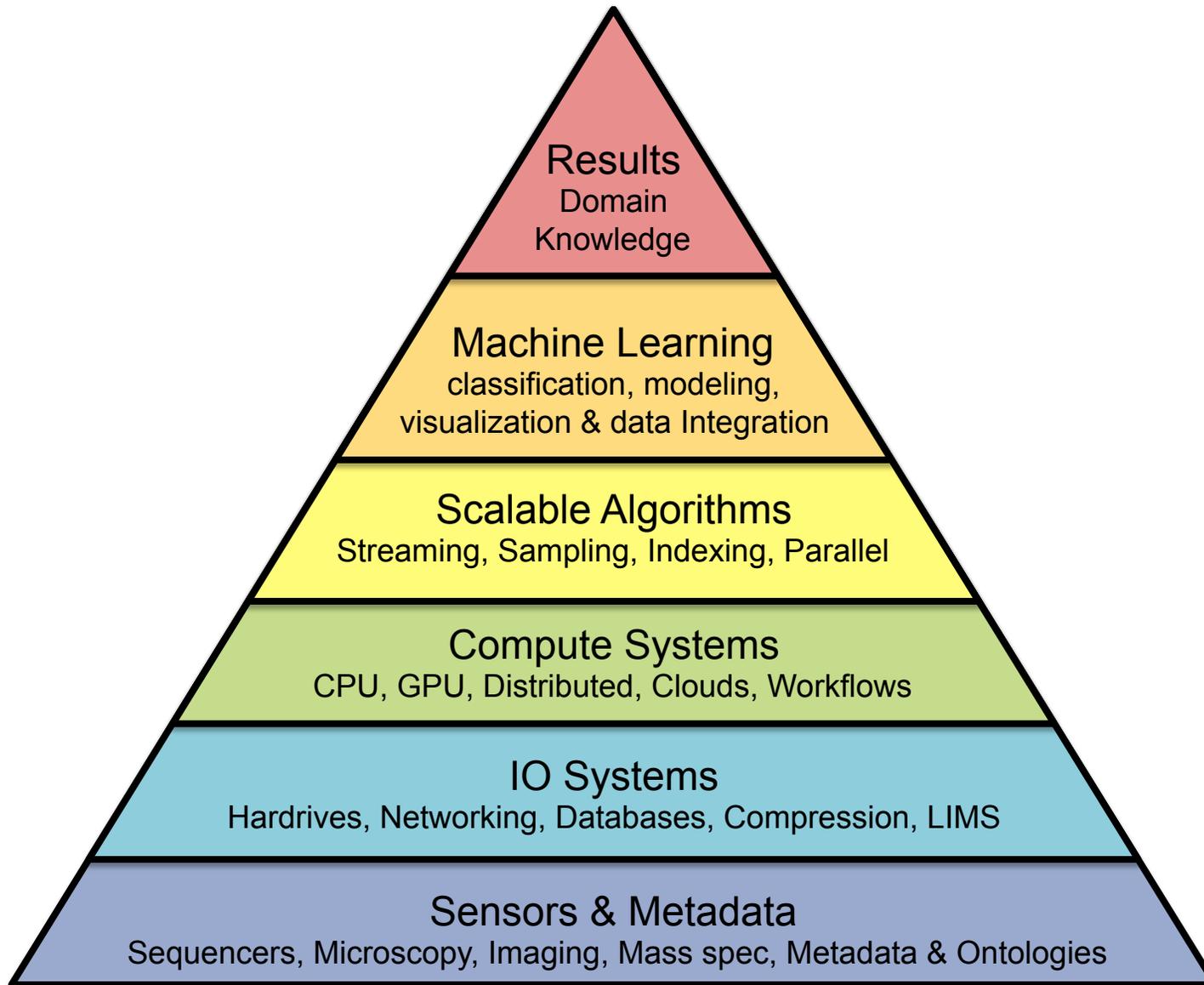


**Adam Siepel**  
Evolution  
Functional Annotation

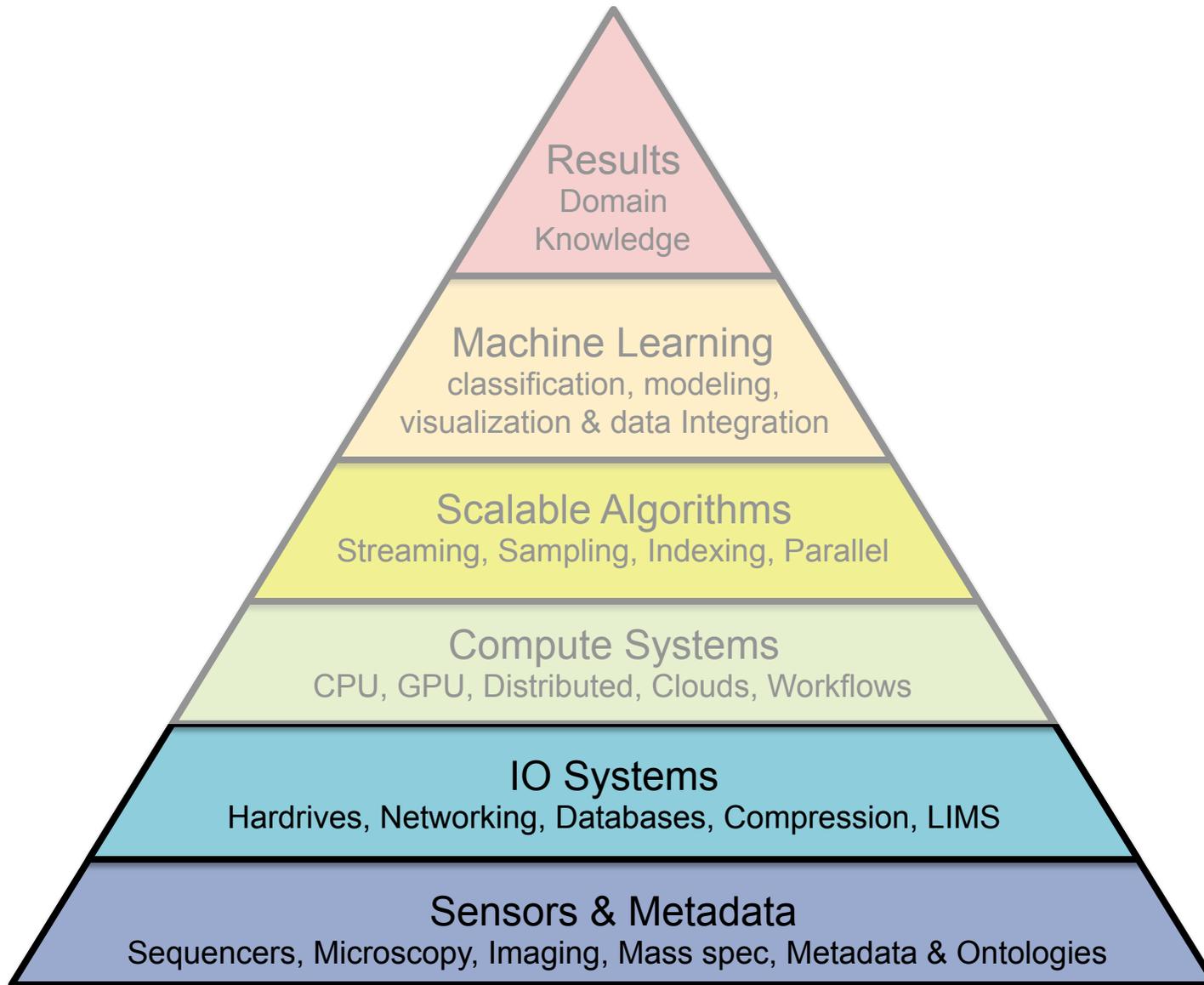


**Michael Wigler**  
Genetic Disorders  
Cancer, Autism

# Quantitative Biology Technologies



# Quantitative Biology Technologies

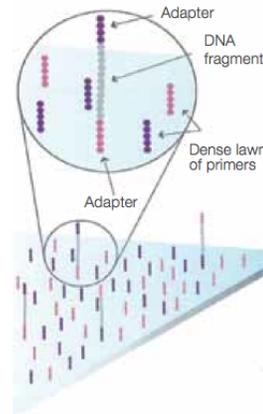


# Massively Parallel Sequencing

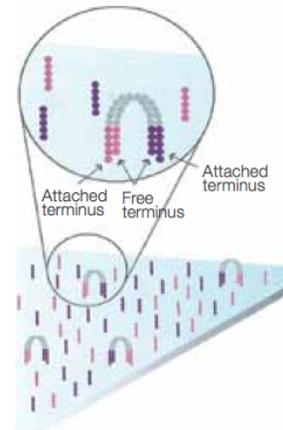


**Illumina HiSeq 2000**  
*Sequencing by Synthesis*

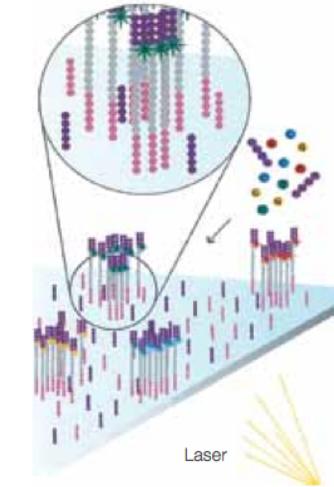
>60Gbp / day



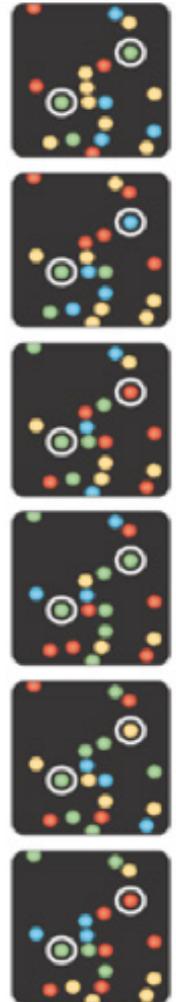
1. Attach



2. Amplify

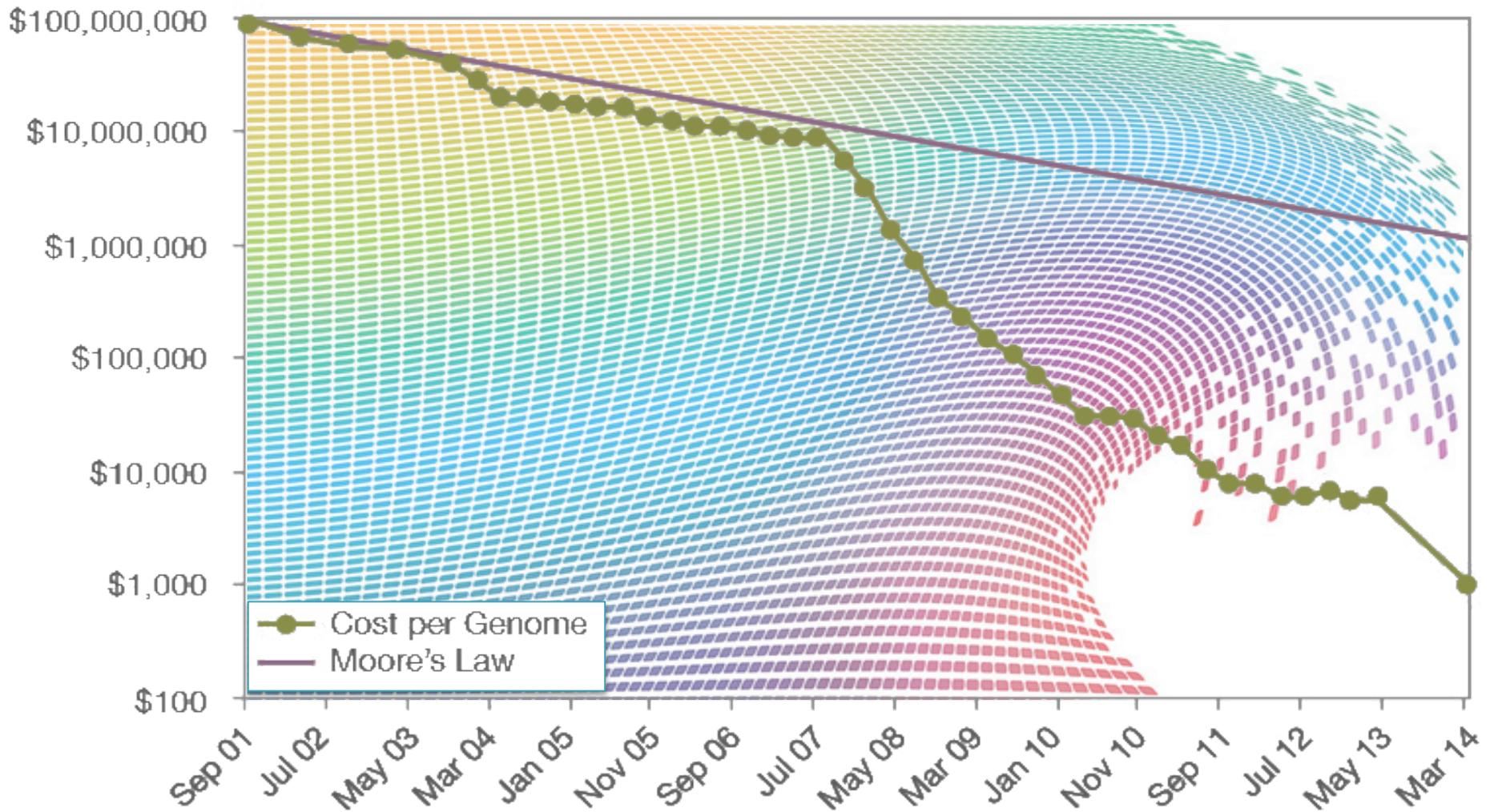


3. Image



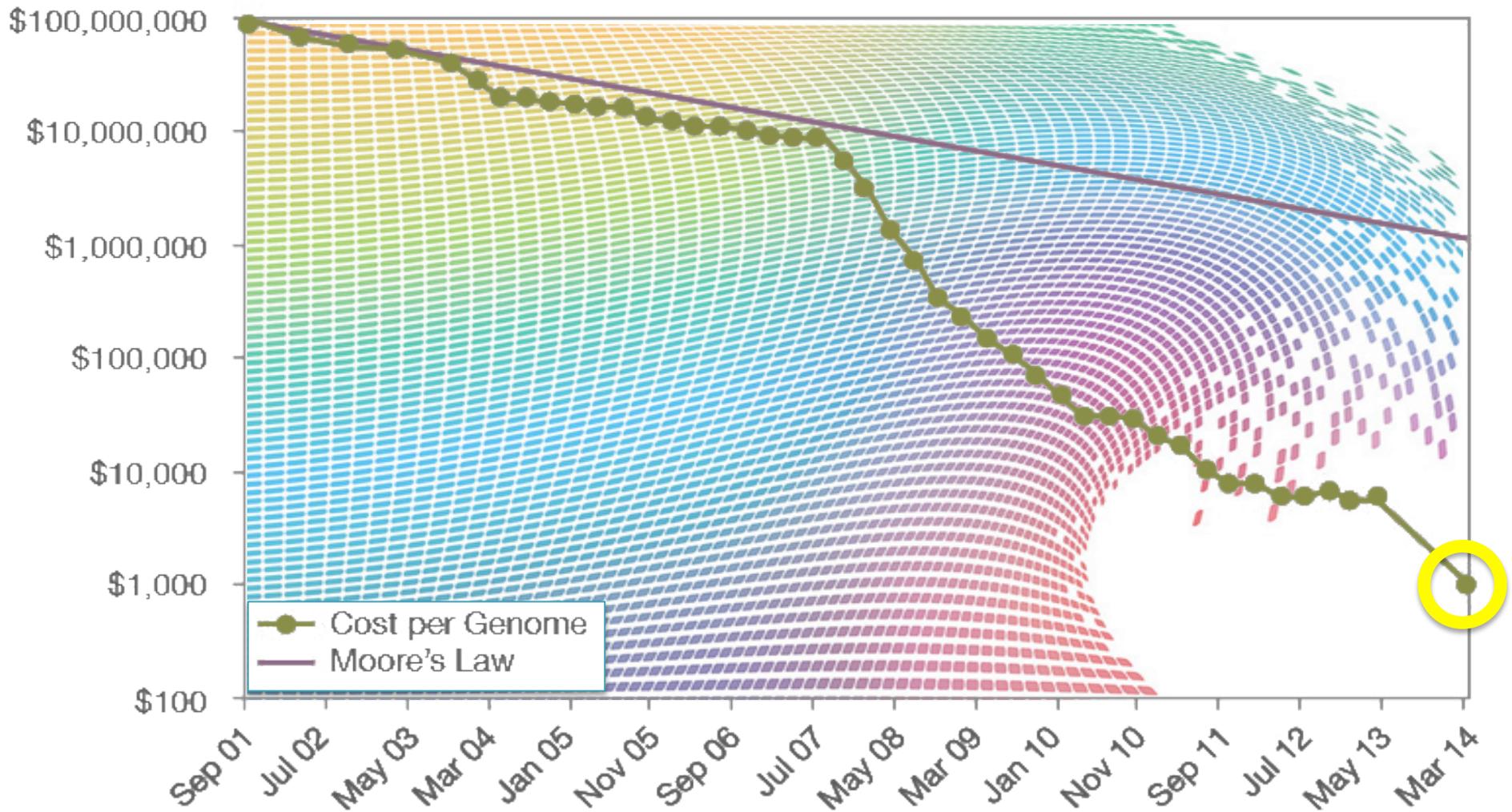
Metzker (2010) Nature Reviews Genetics 11:31-46  
<http://www.youtube.com/watch?v=I99aKKHcxC4>

# Cost per Genome



<http://www.genome.gov/sequencingcosts/>

# Cost per Genome



<http://www.genome.gov/sequencingcosts/>

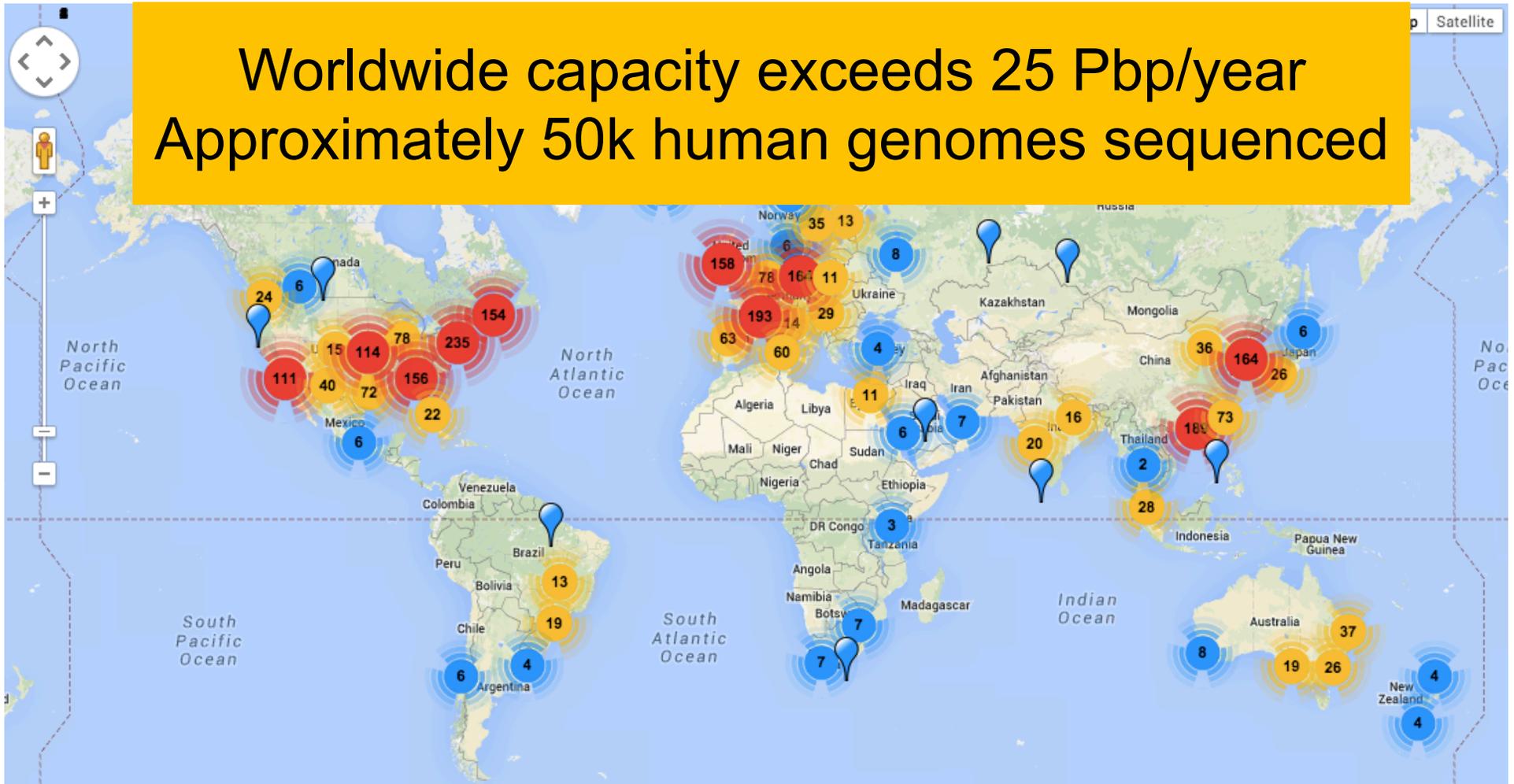
# HiSeq X Ten



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

# Sequencing Centers

Worldwide capacity exceeds 25 Pbp/year  
Approximately 50k human genomes sequenced



# 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



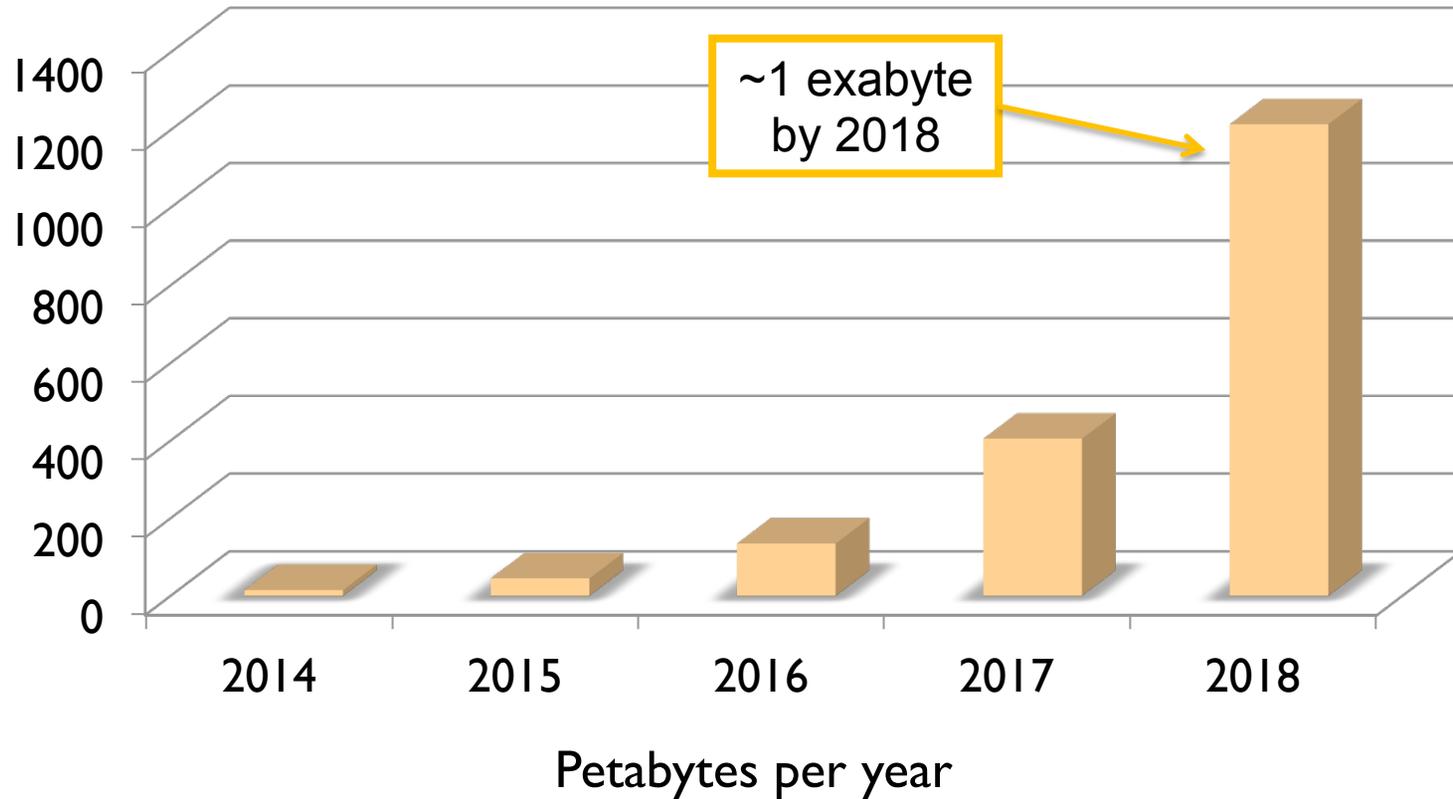
787 feet of DVDs  
~1/6 of a mile tall



500 2 TB drives  
\$500k

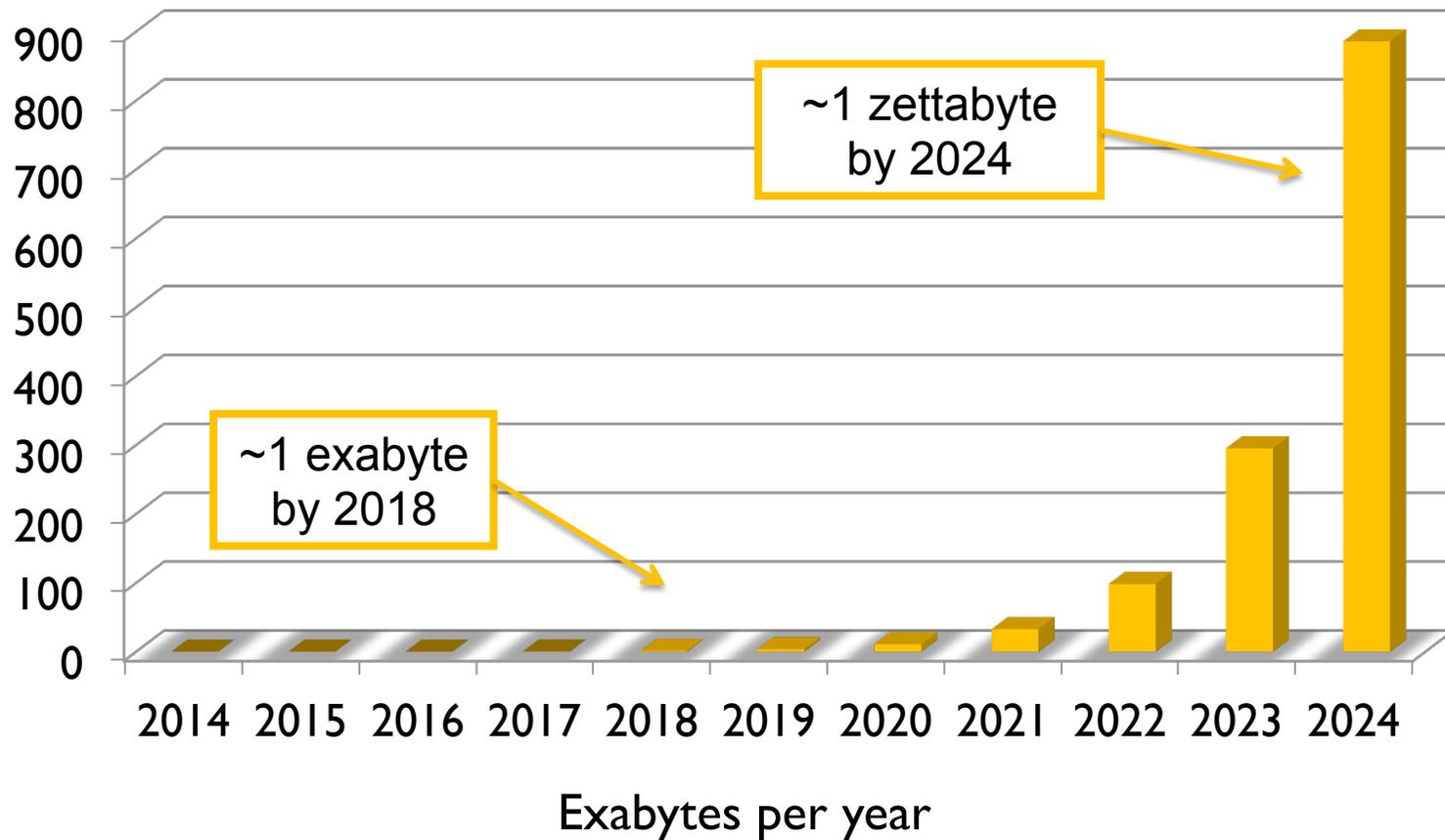
# DNA Data Tsunami

*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  
~ 1/2 distance to moon



Both currently ~100Pb  
And growing exponentially





# 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

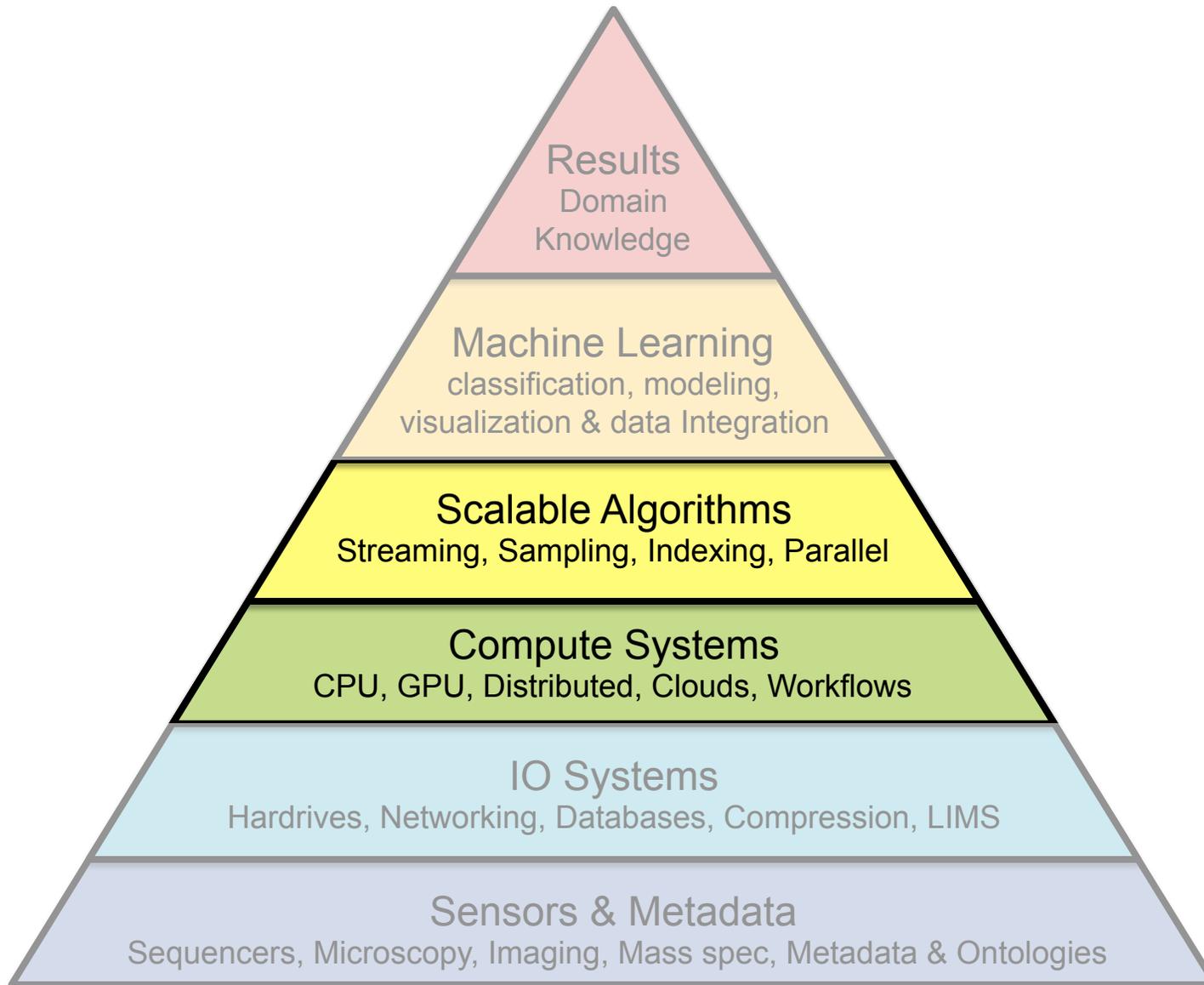
- Major efforts in human health and disease, agriculture, bioenergy
- Genomic information coupled with medical records and other medical data

## But also widely distributed mobile sensors

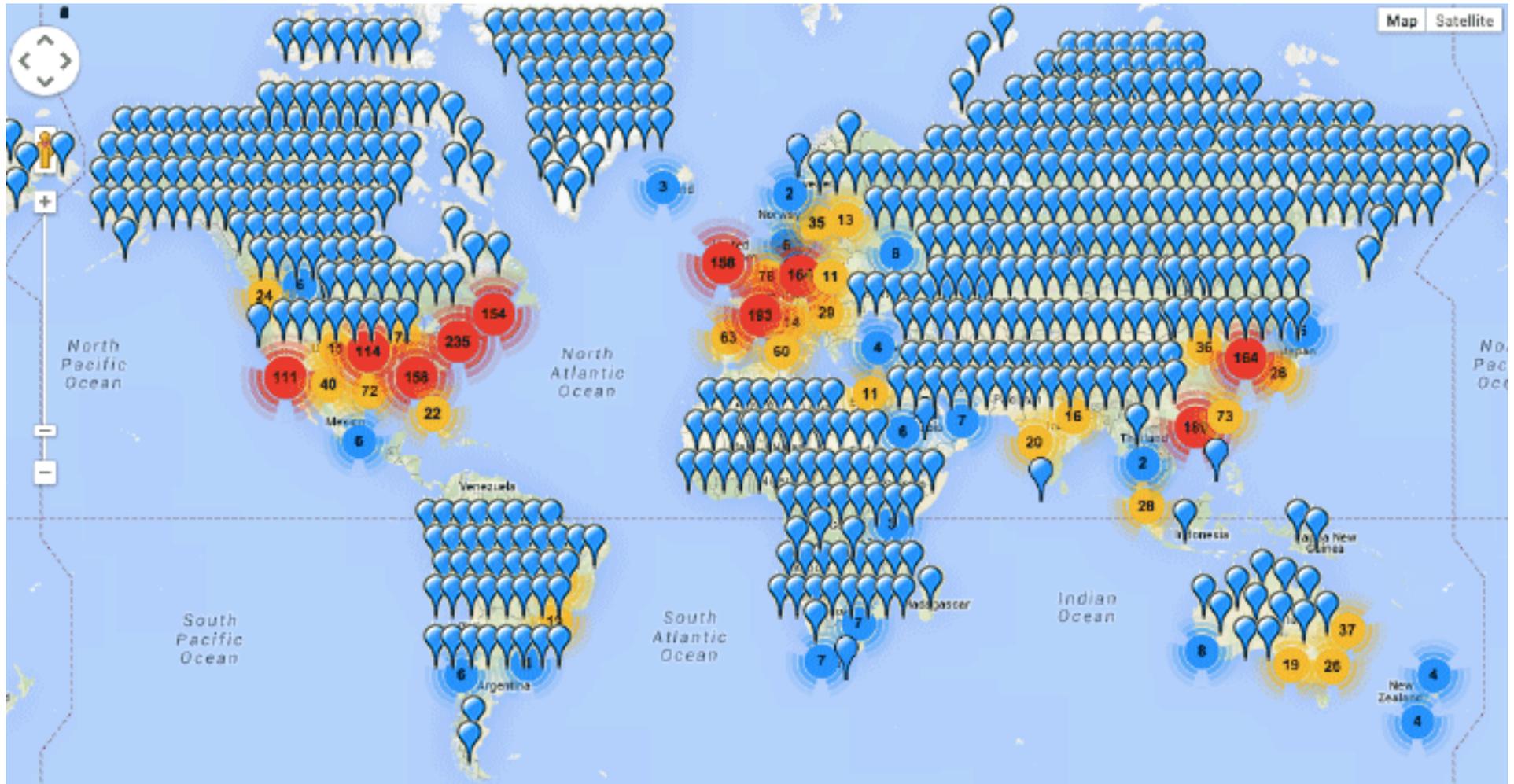
- Schools, offices, sports arenas, transportations centers, farms & food distribution centers
- Monitoring and surveillance, as ubiquitous as weather stations
- The rise of a digital immune system?



# Quantitative Biology Technologies



# Sequencing Centers 2024



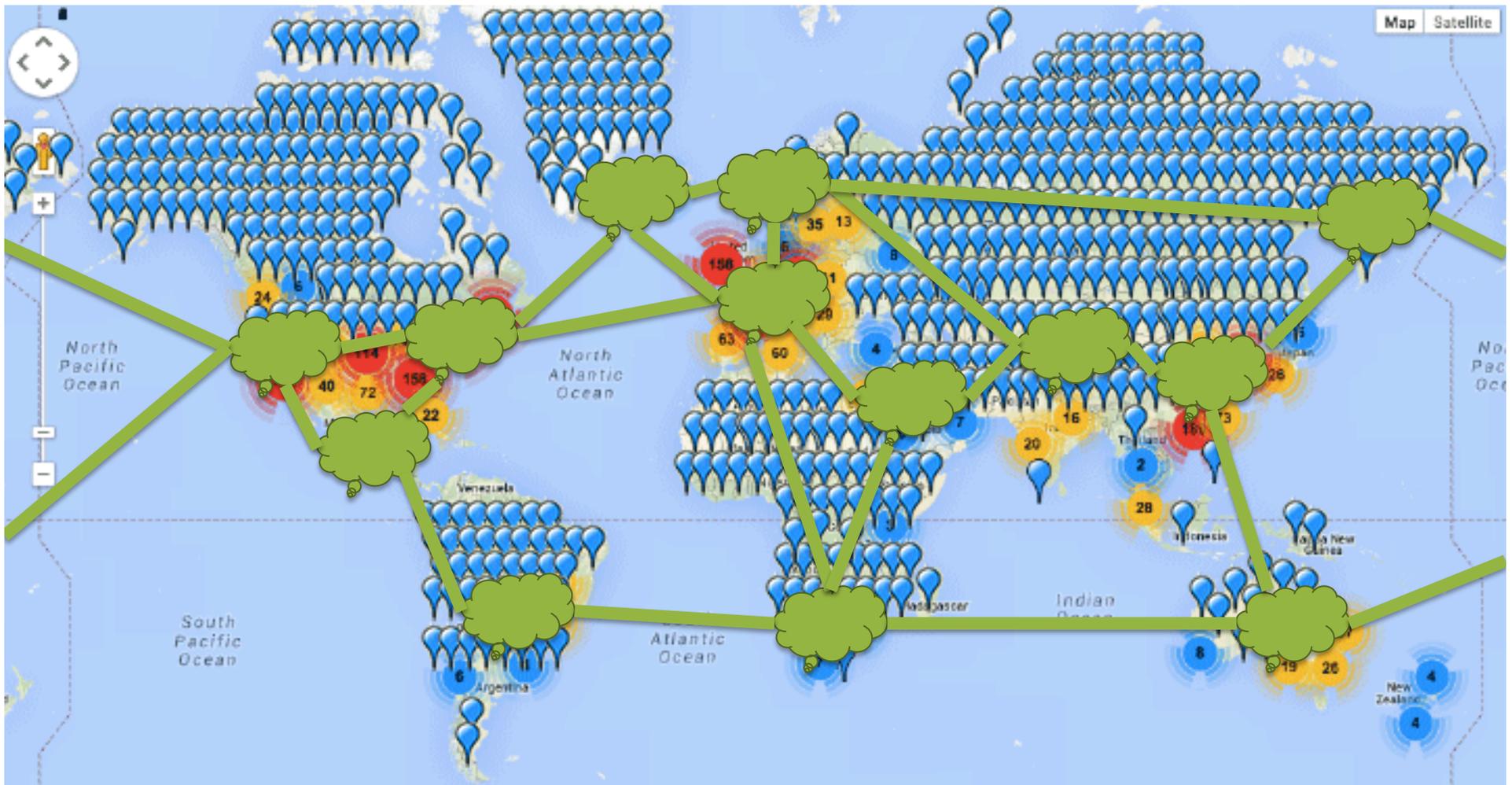
# Informatics Centers 2024



## **The DNA Data Deluge**

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

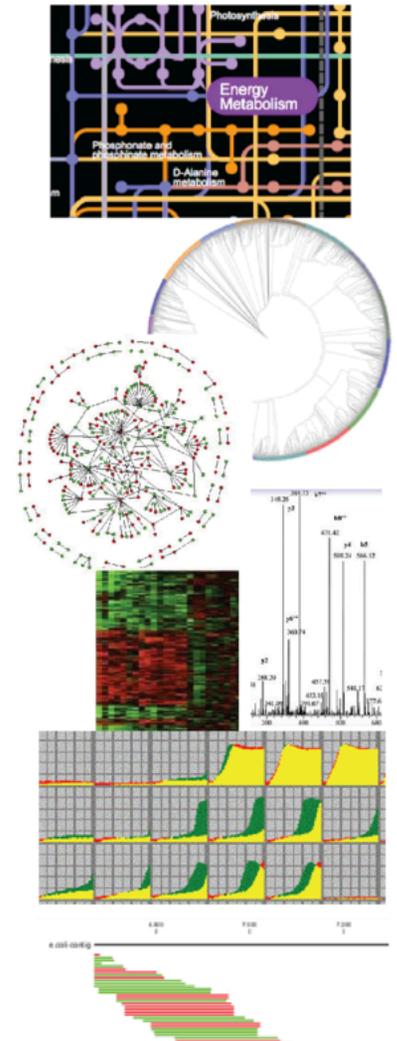
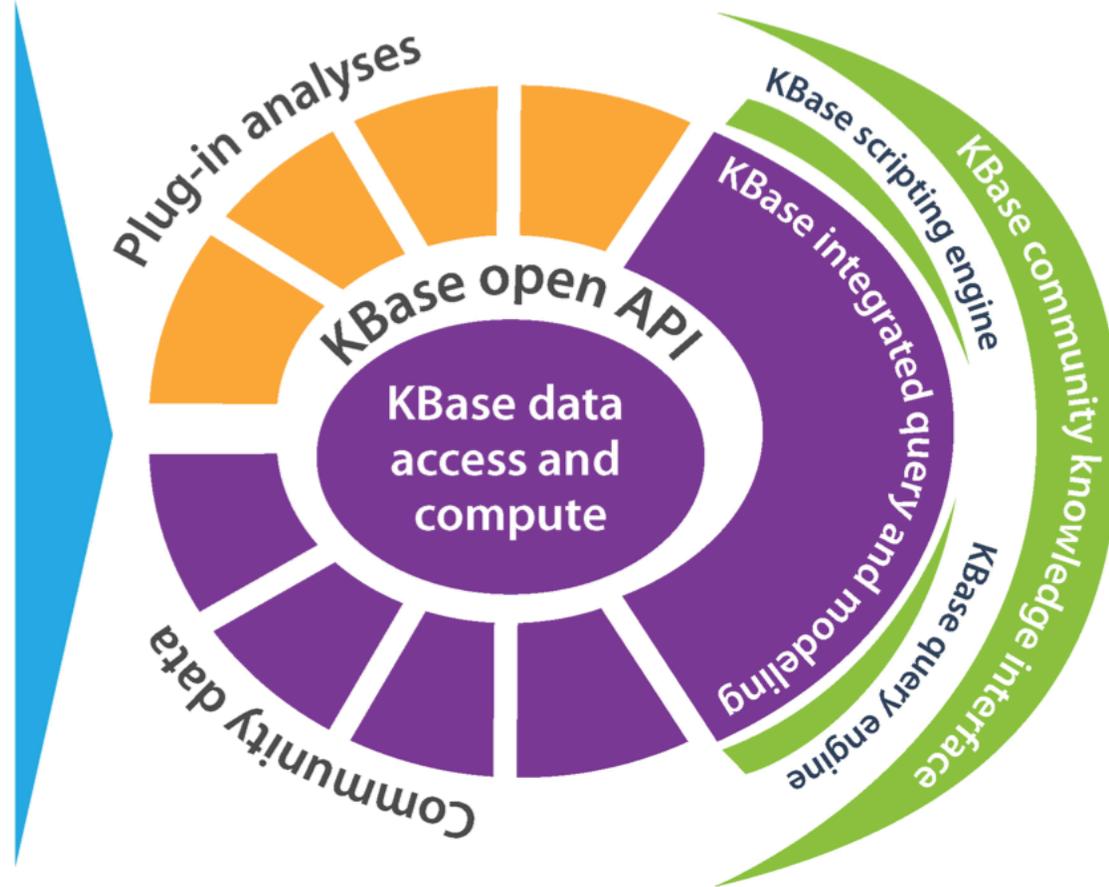
# Informatics Centers 2014



## **The DNA Data Deluge**

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

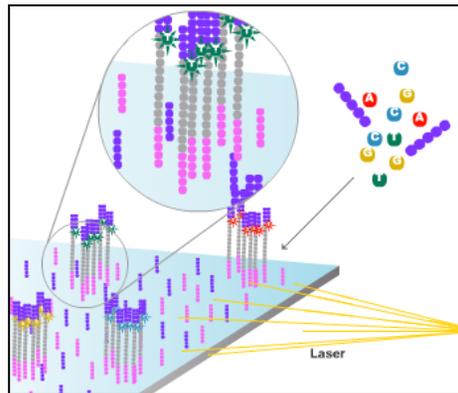
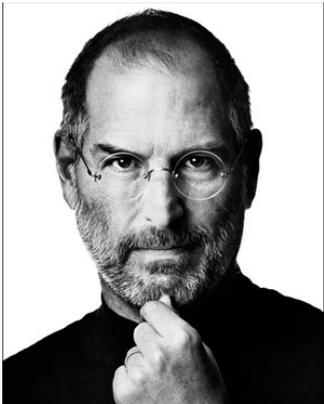
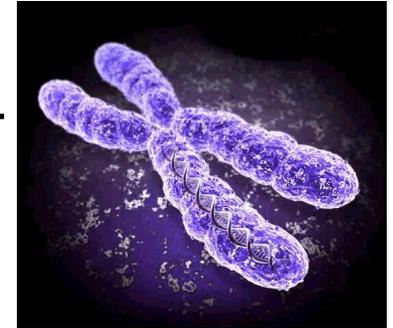
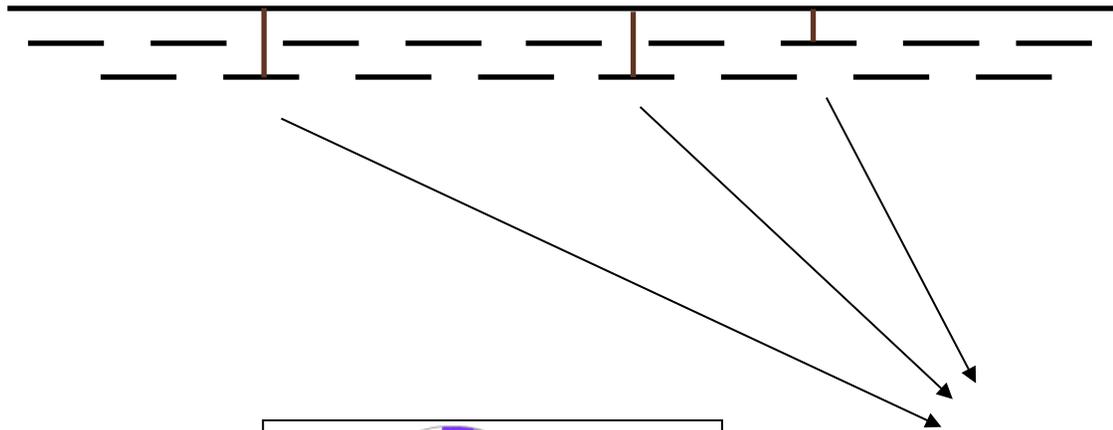
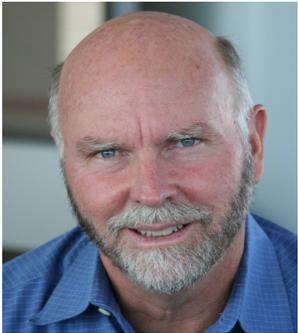
# DOE Systems Biology Knowledgebase



<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 \_\_\_\_\_  
\_\_\_\_\_



# Crossbow

<http://bowtie-bio.sourceforge.net/crossbow>

- Align billions of reads and find SNPs
  - Reuse software components: Hadoop Streaming
  - Mapping with Bowtie, SNP calling with SOAPsnp
- 4 hour end-to-end runtime including upload
  - Costs \$85; Today's costs <\$10

- Very compelling example of cloud computing in genomics
- Commercial vendors probably have better security than your institution
- Need more applications!

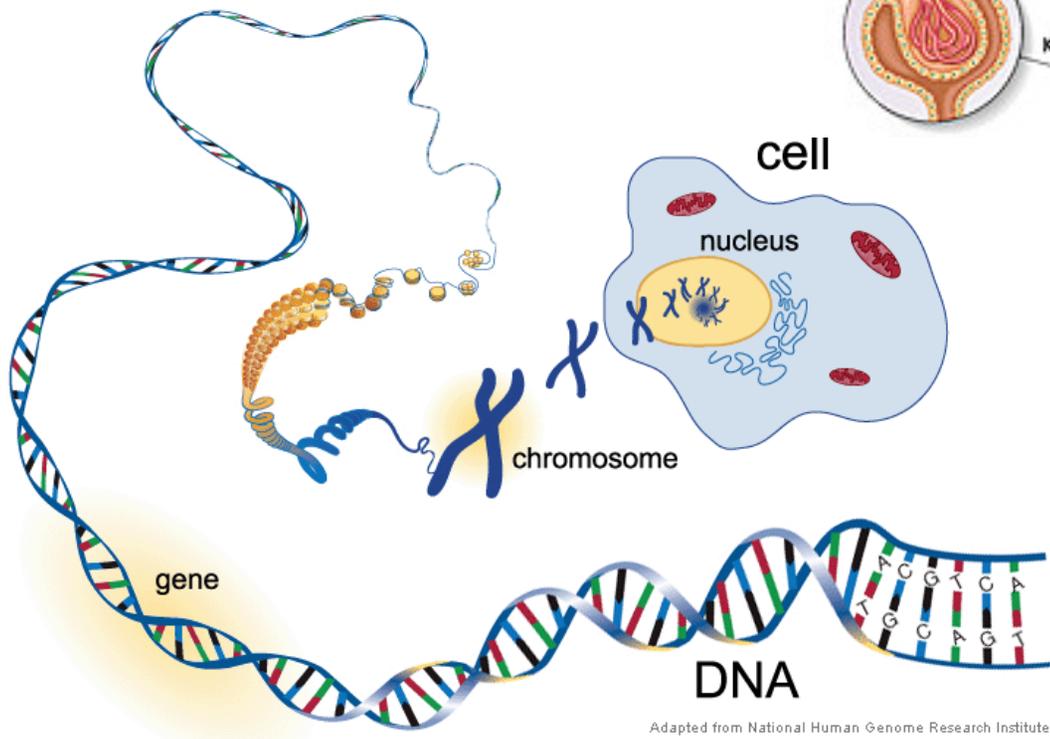
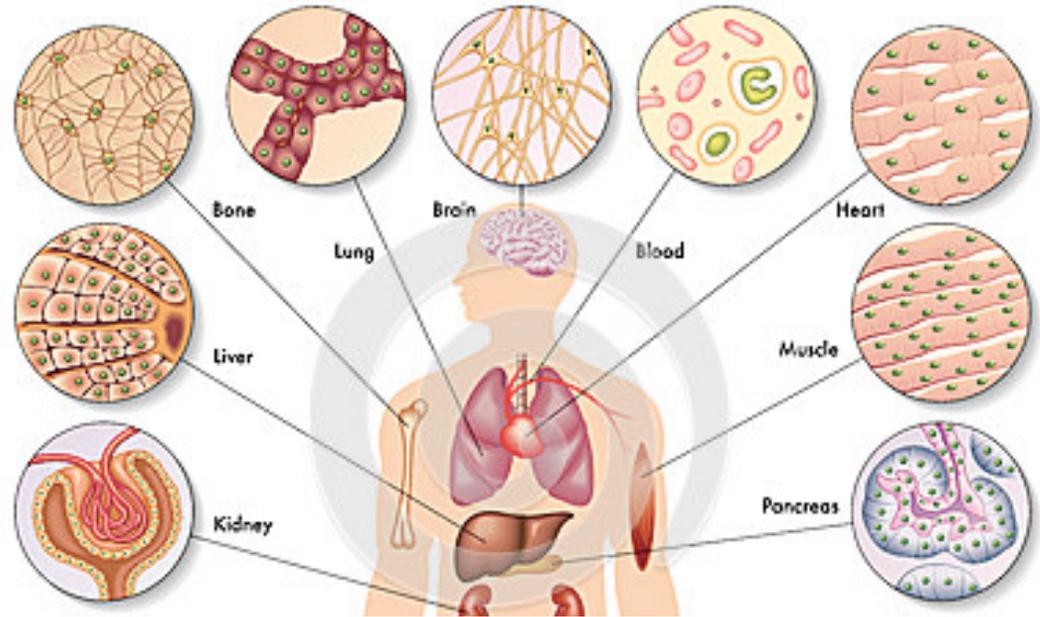


## Searching for SNPs with Cloud Computing.

Langmead B, Schatz MC, Lin J, Pop M, Salzberg SL (2009) *Genome Biology*. **10**:R134

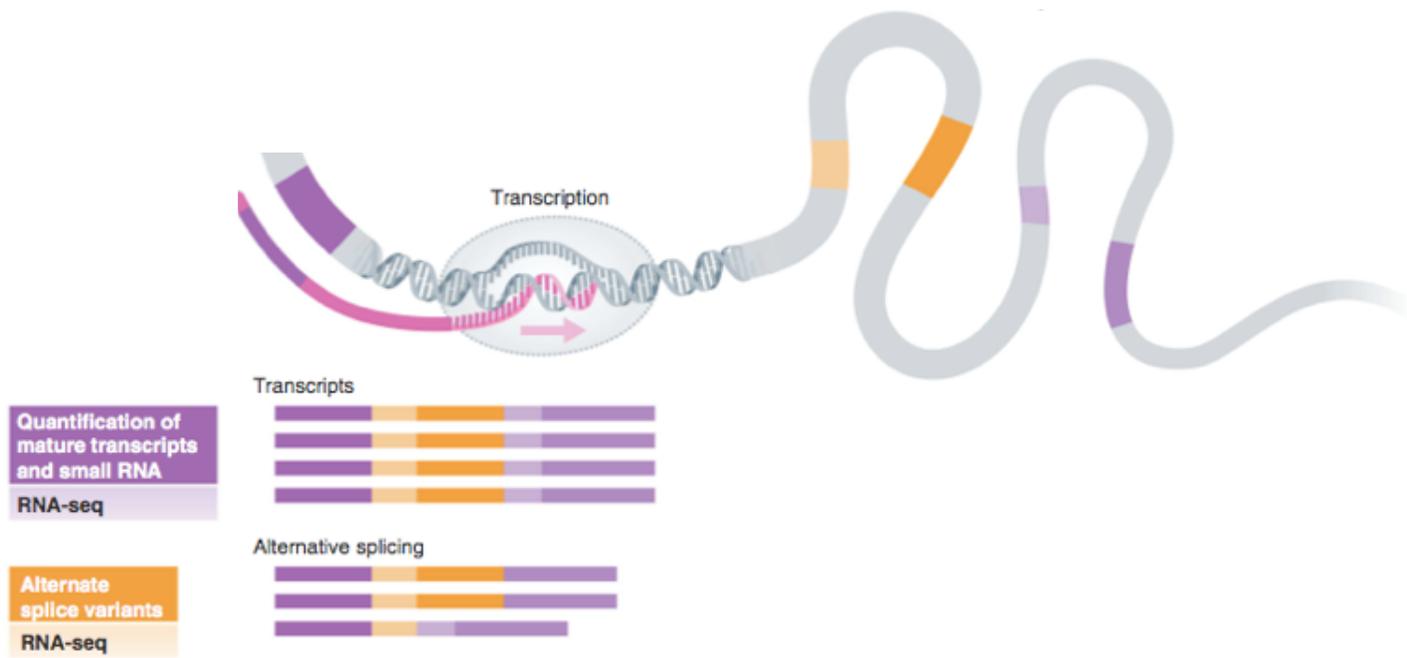
# Cells & DNA

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

Adapted from National Human Genome Research Institute



# 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**

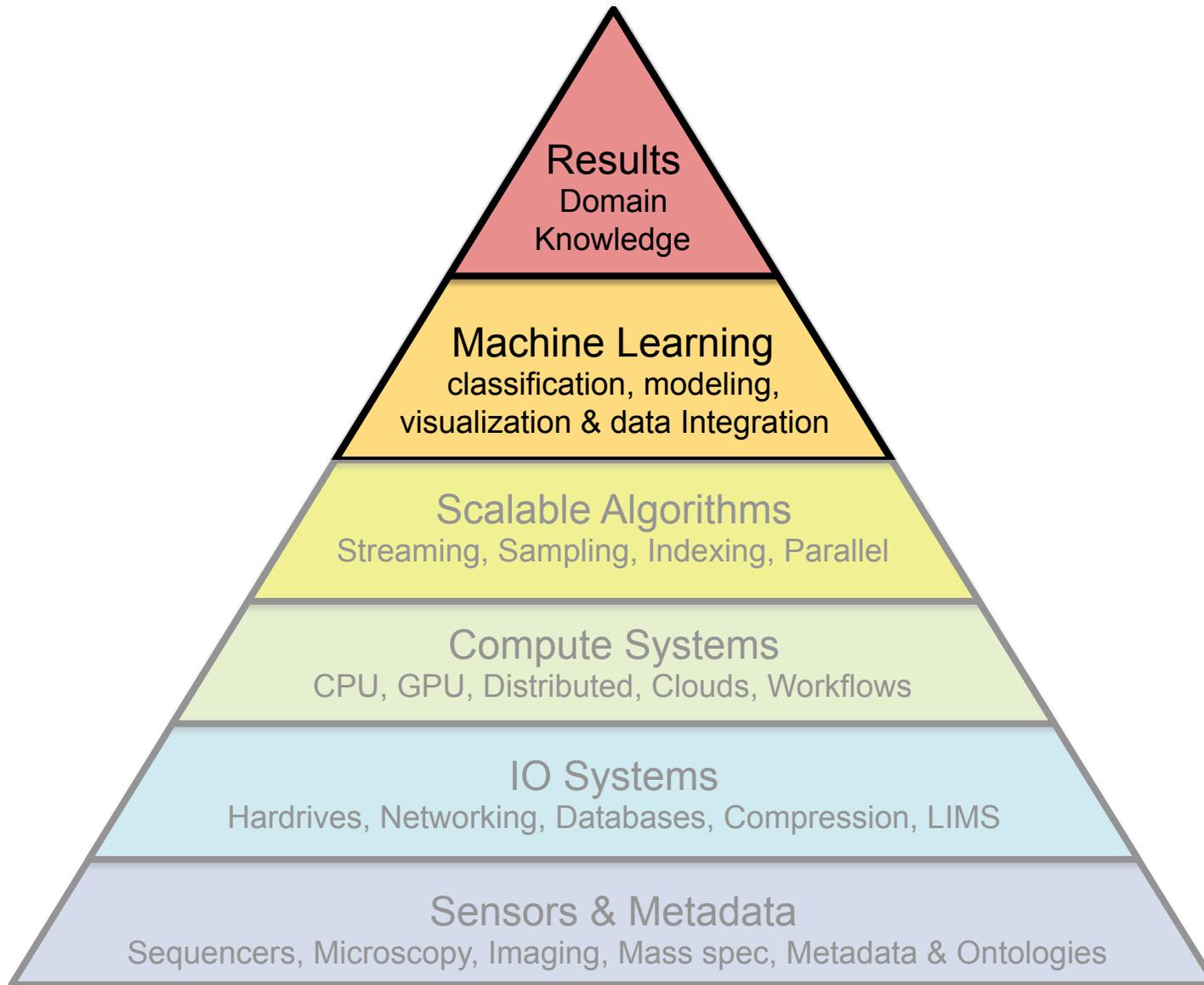
- Expect to see >1000 cores in a single computer
- Compute & IO needs to be considered together
- Rewriting efficient parallel software is complex and expensive

**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



# Quantitative Biology 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?**

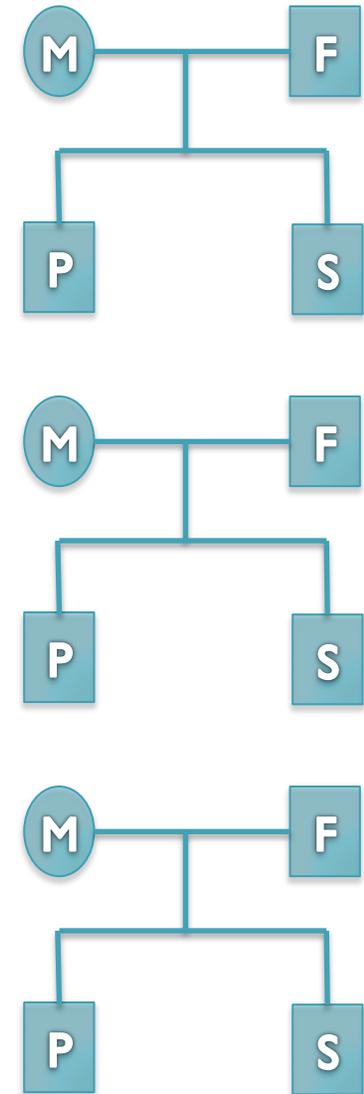
<http://www.autismspeaks.org/what-autism>

# Searching for the genetic risk factors

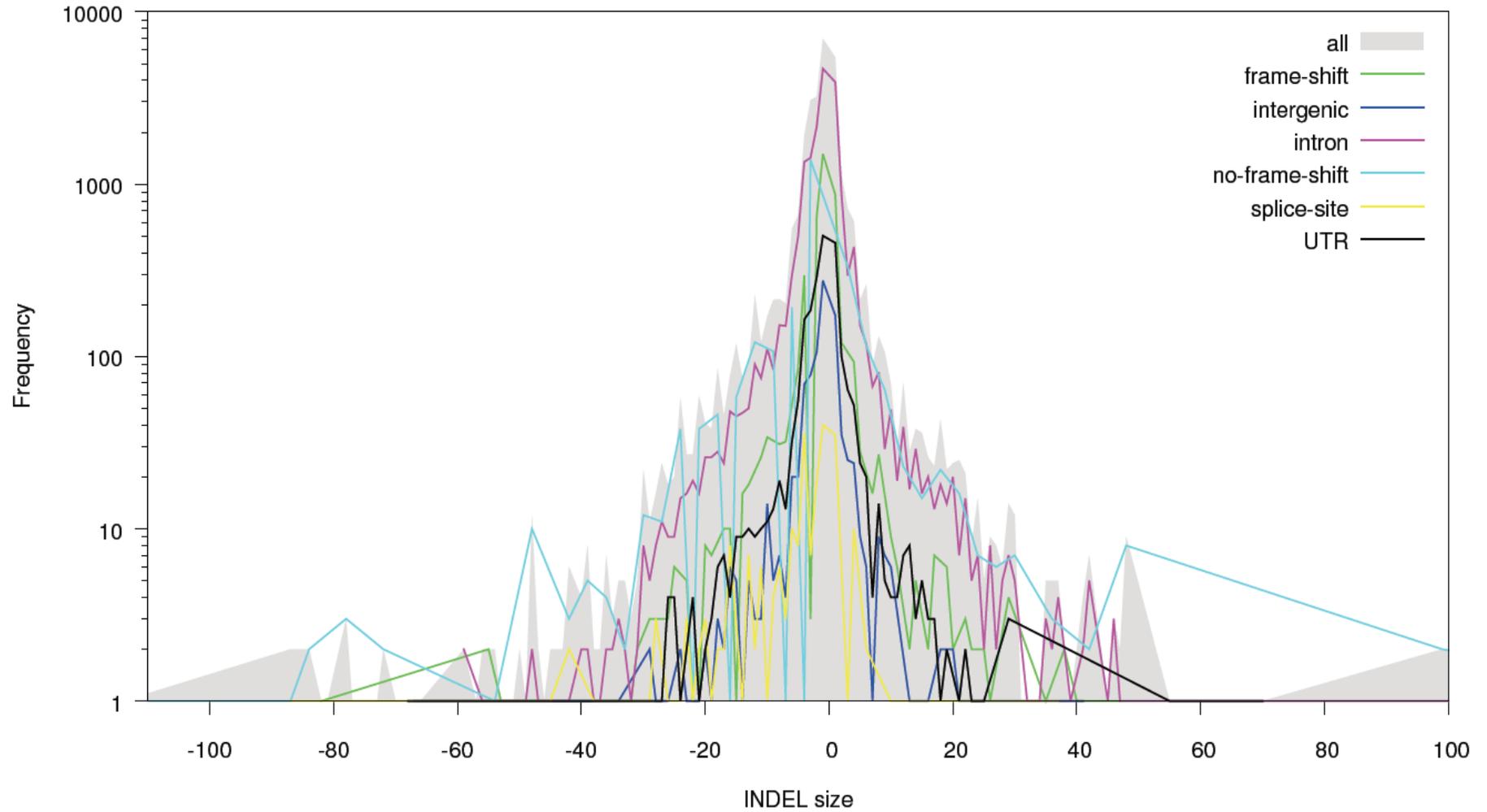
## Search Strategy

- 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?***



# Population Analysis of the SSC

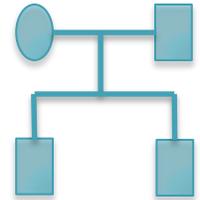


Constructed database of >IM transmitted and de novo genetic mutations

# De novo mutation discovery and validation

## De novo mutations:

Sequences not inherited from your parents.



**Reference:** . . . **TCAAATCCTTTTAAATAAAGAAGAGCTGACA** . . .

Father(1): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Father(2): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Mother(1): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Mother(2): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Sibling(1): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Sibling(2): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

Proband(1): . . . TCAAATCCTTTTAAATAAAGAAGAGCTGACA . . .

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

4bp heterozygous deletion at chr15:9352406 | CHD2

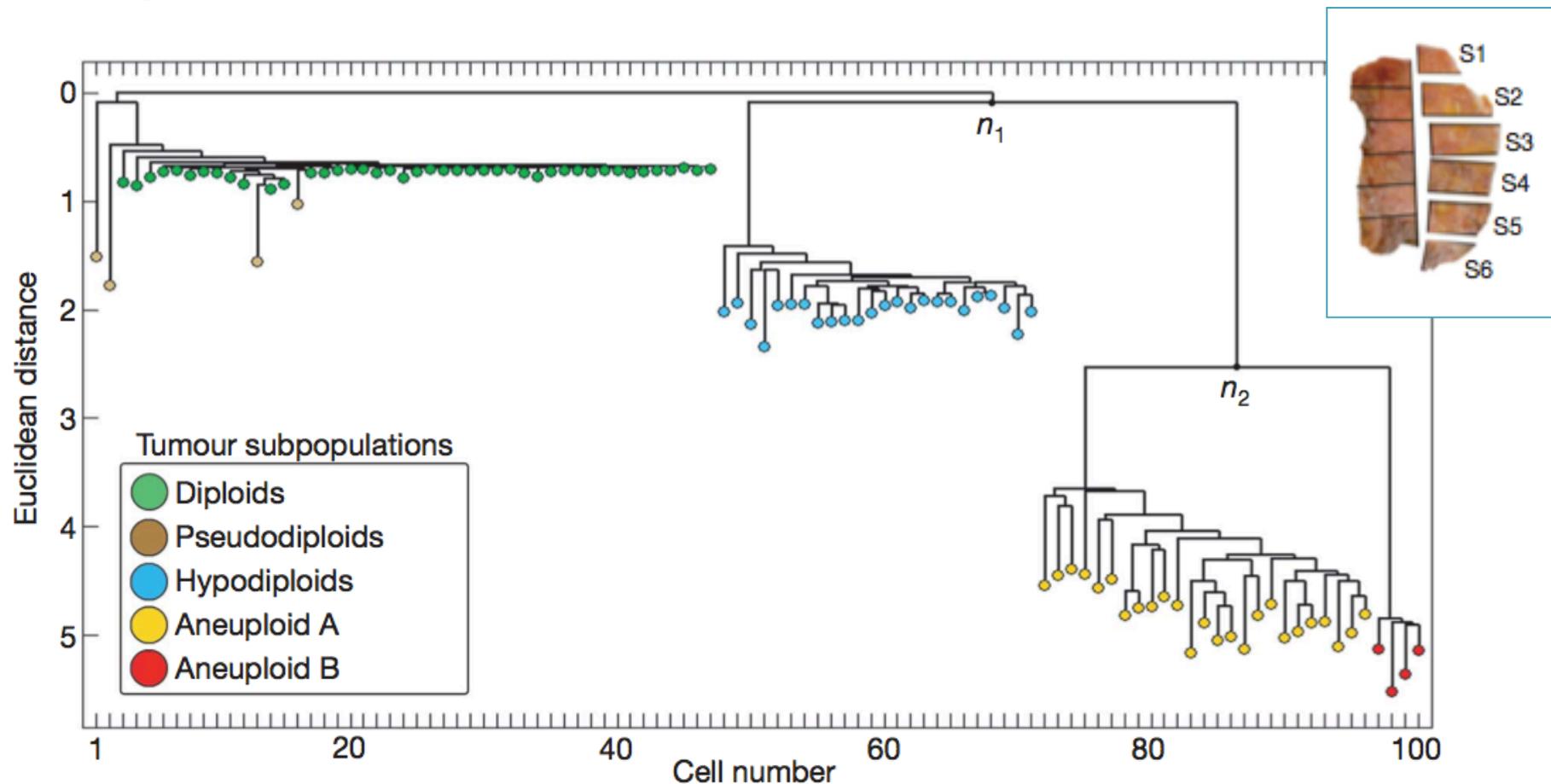
# De novo Genetics of Autism

- 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 FMR1
  - Related to neuron development and synaptic plasticity
  - Also strong overlap with chromatin remodelers

**Accurate detection of de novo and transmitted INDELs within exome-capture data using micro-assembly**  
Narzisi, G, O'Rawe, J, Iossifov, I, Lee, Y, Wang, Z, Wu, Y, Lyon, G, Wigler, M, Schatz, MC (2014) *In press*.

# Tumour evolution inferred by single-cell sequencing

Nicholas Navin<sup>1,2</sup>, Jude Kendall<sup>1</sup>, Jennifer Troge<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?

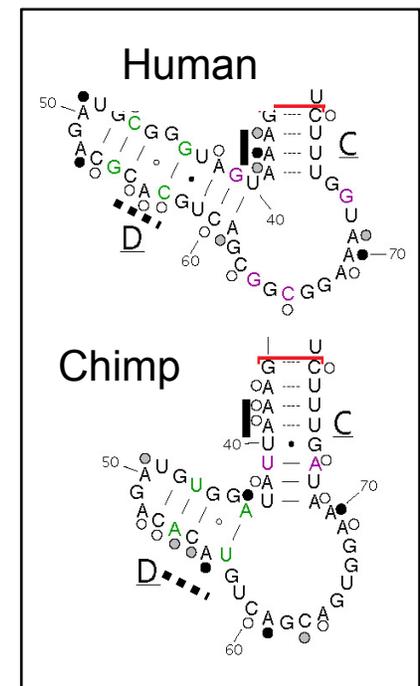
## “Human Accelerated Regions”



human	TTGATGGCTGTAGACCCAGTGCAGCGCGGAAATGGTTTCTATCAAAATCAAAGTGT	TTTAGAGATTTTCCTCAAATTTCAAATTA
chimp	TTTATGGCTGTAGACACATGTCAGCAGTGGAAATAGTTTCTATCAAAATCAAAGT	TTTAGAGATTTTCCTCAAATTTCAAATTA
dog	TTTATGGCTGTAGACACATGTCAGCGGTGCAAAACAGTTTCTATCAAAATCAAAGT	TTTAGAGATTTTCCTCAAATTTCAAATTA
mouse	TTTATGGCTGTAGACACATGTCAGCGCGGAAATGGTTTCTATCAAAATCAAAGT	TTTAGAGATTTTCCTCAAATTTCAAATTA
rat	TTTATGGCTGTAGACACATGTCAGCAGTGGAAATGGTTTCTATCAAAATCAAAGT	TTTAGAGATTTTCCTCAAATTTCAAATTA
chicken	TTTATGGCTGTAGACACATGTCAGCAGTGAAGAACAGTTTCTATCAAAATCAAAGT	TTTAGAGATTTTCCTCAAATTTCAAATTA

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 Golub,<sup>6</sup> Fran Halmage,<sup>7,8,9</sup> Yaniv Erlich<sup>1\*</sup>

Sharing sequencing data sets without identifying individuals is a challenge. Here, we report that surnames can be recovered from Y-chromosome STR repeats on the Y chromosome (Y-STRs) and that a combination of a surname and Y-STR data can be used to triangulate the identity of the father. This relies on free, publicly accessible Internet resources for identification for U.S. males. We further demonstrate that with high probability the identities of multiple

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

<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 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 Biology 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: yaniv@wi.mit.edu

www.s

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. Gitschier (11)

## Predicting Social Security numbers from public data

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

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.

have already left the barn: We demonstrate that it is possible to 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



# How can you participate?



## **Students**

- Learn python!
- Study math & statistics & computer science
- Visit the DNA Learning Center

## **Individuals**

- Personal Genome Project  
Harvard Medical School  
<http://www.personalgenomes.org>
- 23andMe  
Genetic testing and ancestry  
<http://www.23andme.com>
- CSHL Public Lectures & Events  
<http://www.cshl.edu>

# Acknowledgements

## Schatz Lab

Giuseppe Narzisi  
Shoshana Marcus  
James Gurtowski  
Srividya  
Ramakrishnan  
Hayan Lee  
Rob Aboukhalil  
Mitch Bekritsky  
Charles Underwood  
Tyler Gavin  
Maria Nattestad  
Alejandro Wences  
Greg Vulture  
Eric Biggers  
Aspyn Palatnick

## CSHL

Hannon Lab  
Gingeras Lab  
Jackson Lab  
Iossifov Lab  
Levy Lab  
Lippman Lab  
Lyon Lab  
Martienssen Lab  
McCombie Lab  
Tuveson Lab  
Ware Lab  
Wigler Lab  
  
IT Department



# *CSHL Public Lecture*

June 24, 2014 @ 7-9pm

Understanding Autism Spectrum Disorders: Focus on the Facts

Michael Ronemus, Ph.D. & Rebecca Sachs, Ph.D.



# Thank you!

<http://schatzlab.cshl.edu>

@mike\_schatz