

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



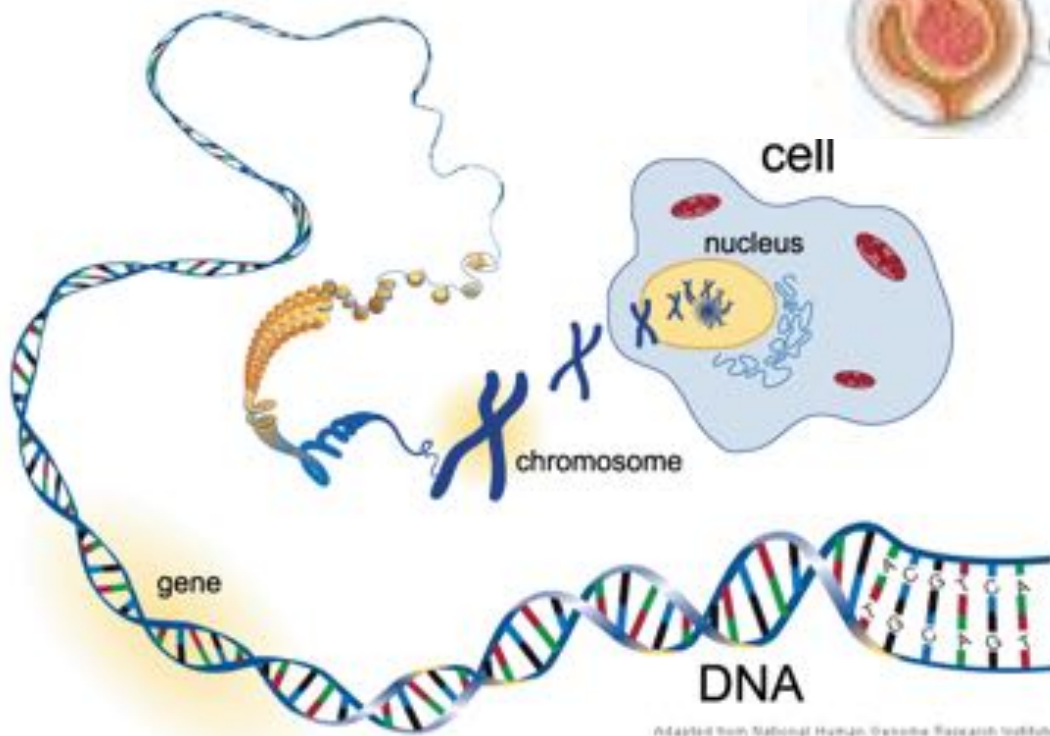
Your DNA, along with your environment and experiences, shapes who you are

- 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

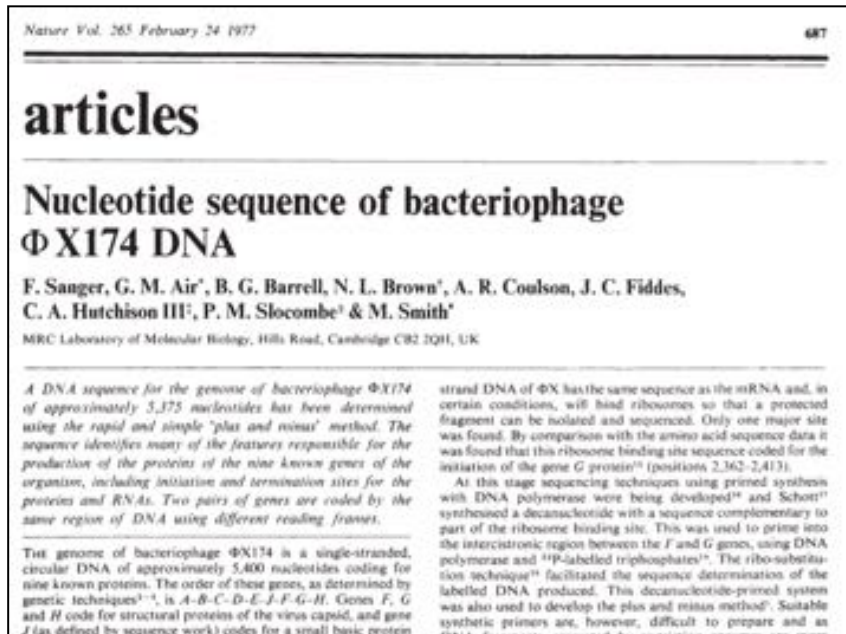
Cells & DNA

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
1st Complete Organism
Bacteriophage ϕ 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.

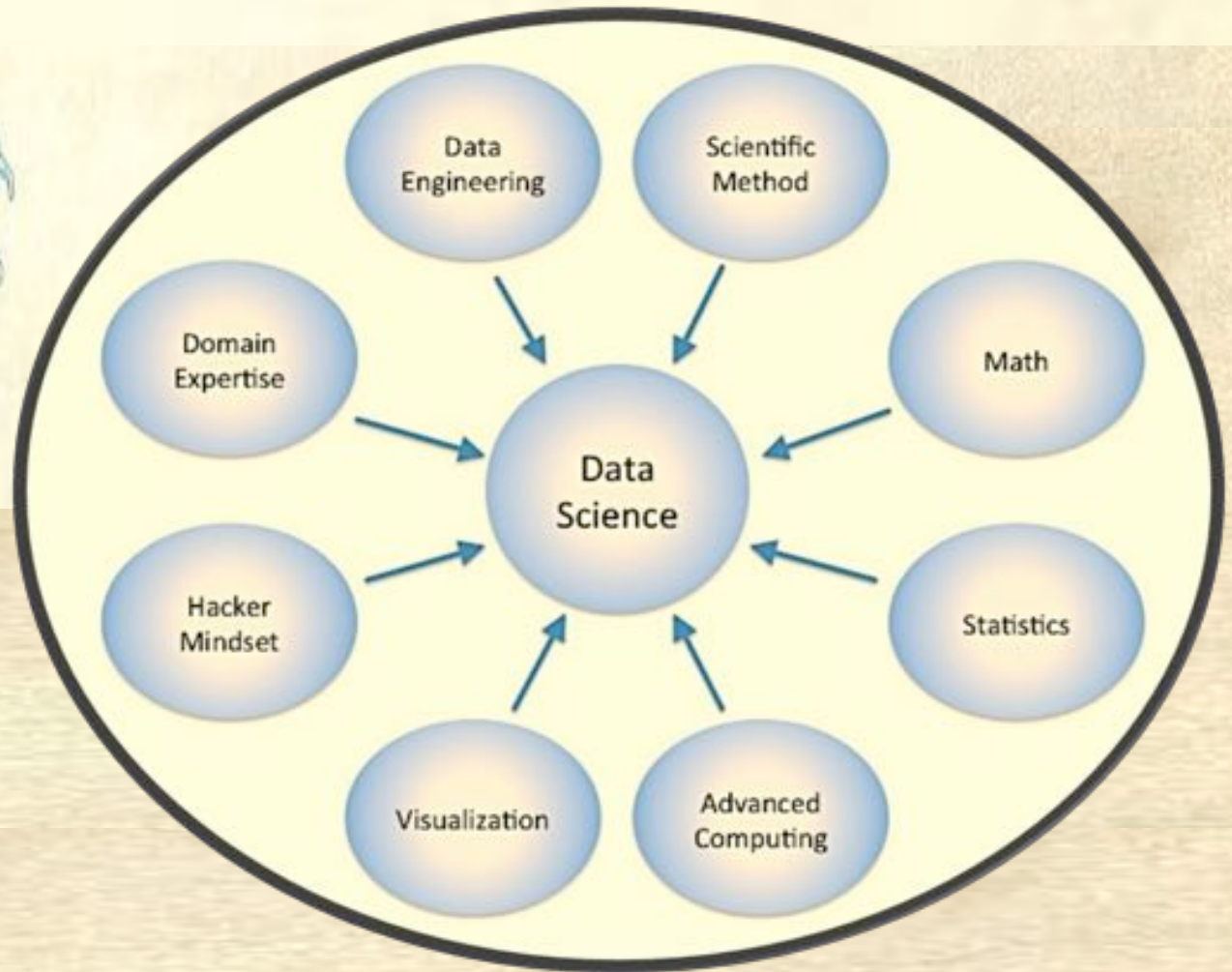
What software and systems will?

And who will create them?

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



Who is a Data Scientist?



http://en.wikipedia.org/wiki/Data_science

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

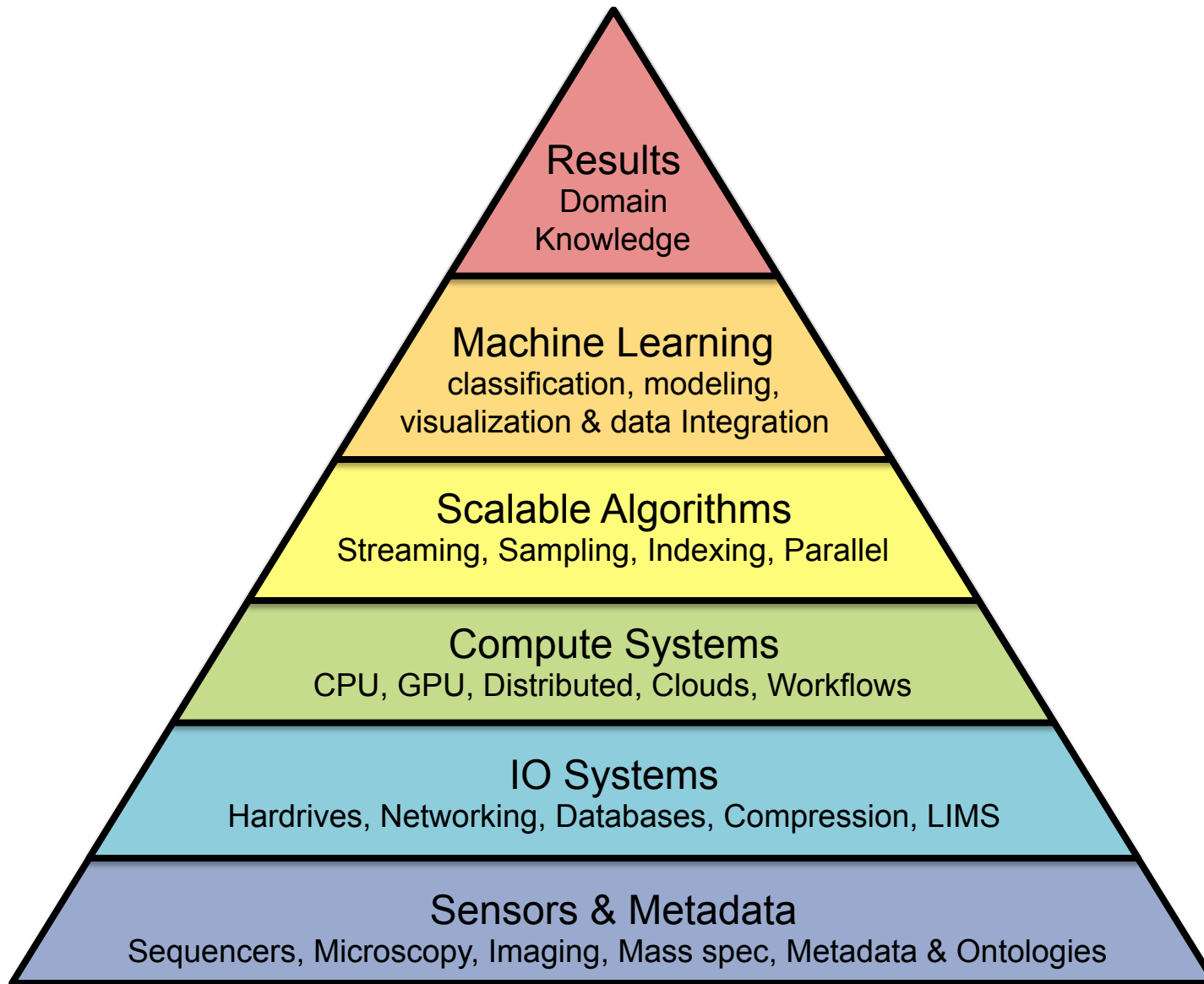


Adam Siepel
Evolution
Functional Annotation

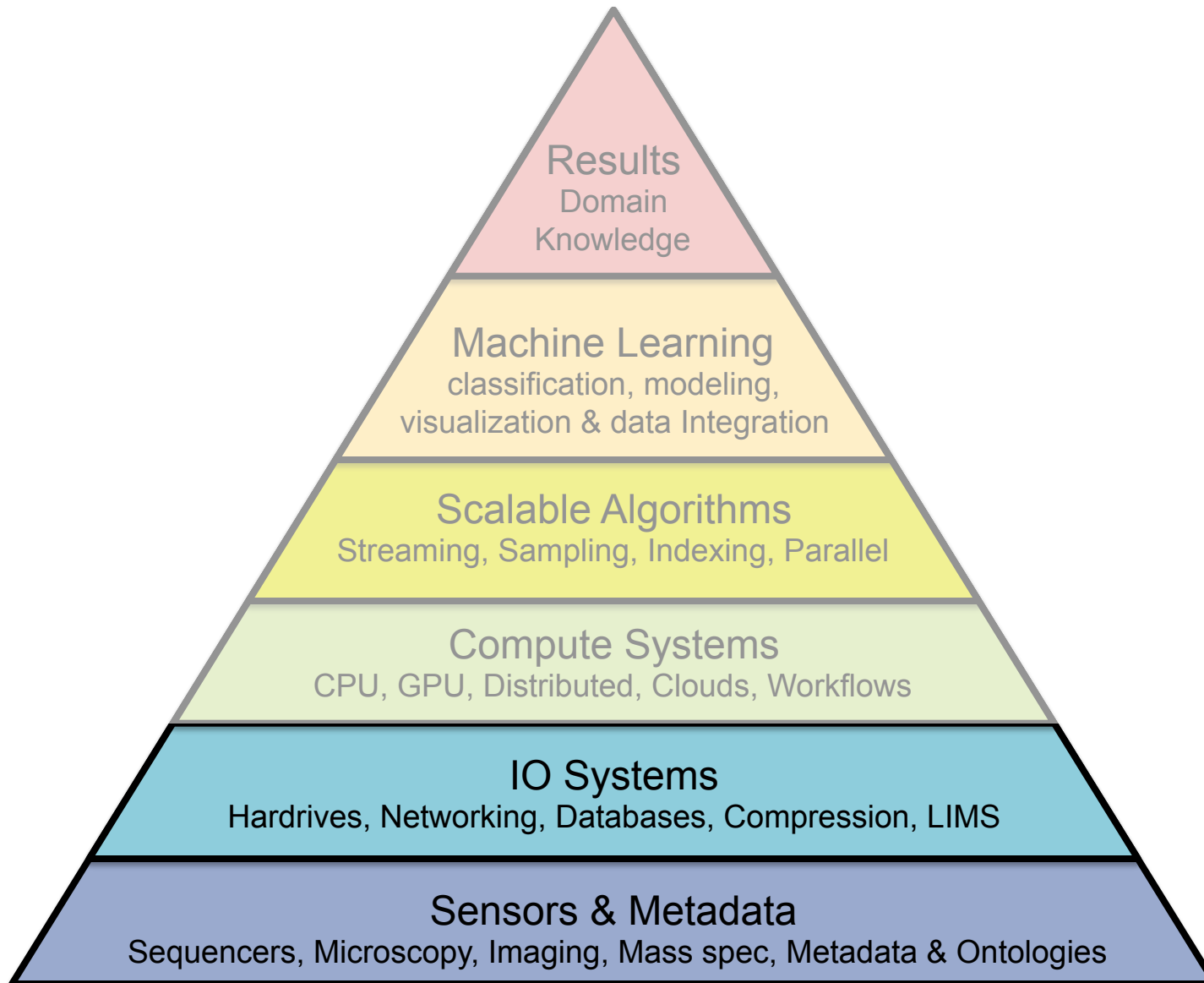


Michael Wigler
Genetic Disorders
Cancer, Autism

Quantitative Biology Technologies



Quantitative Biology Technologies

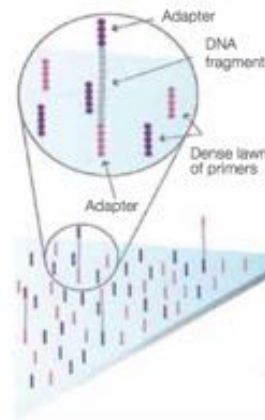


Massively Parallel Sequencing

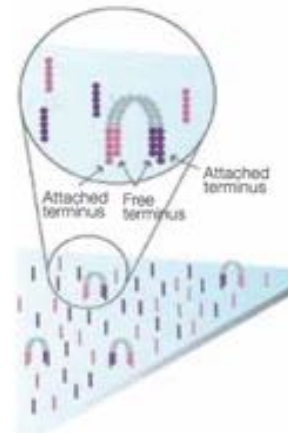


Illumina HiSeq 2000
Sequencing by Synthesis

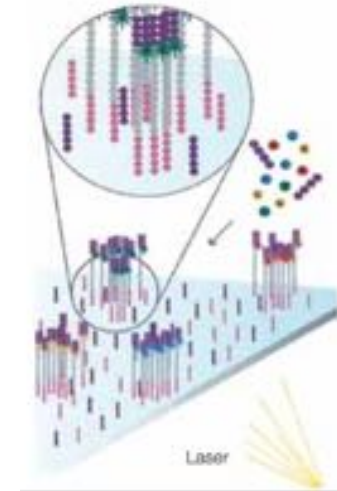
>60Gbp / day



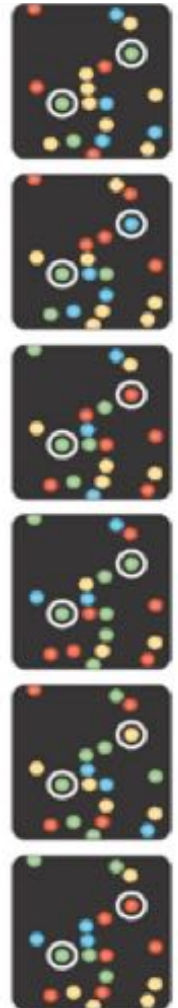
1. Attach



2. Amplify

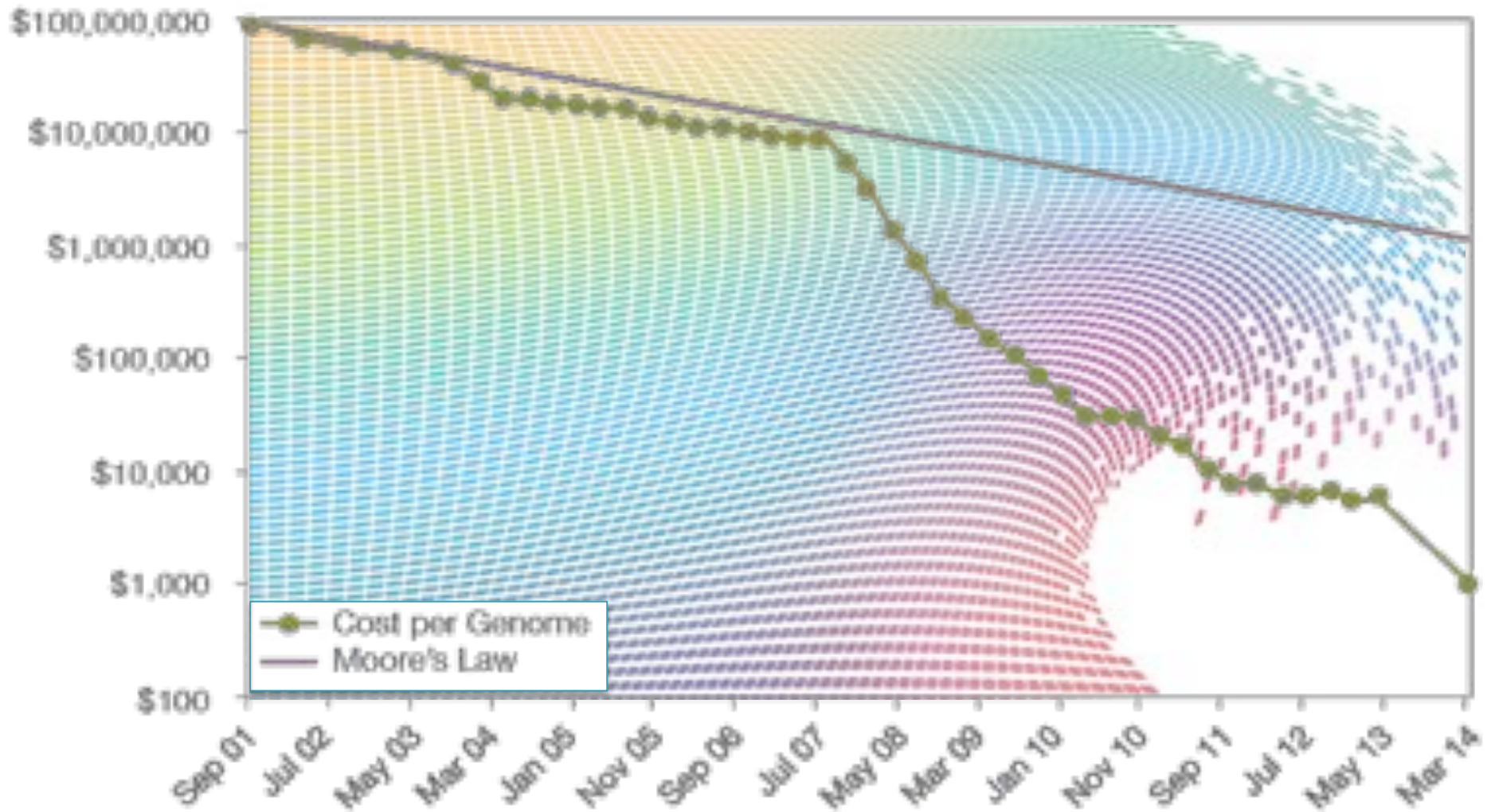


3. Image



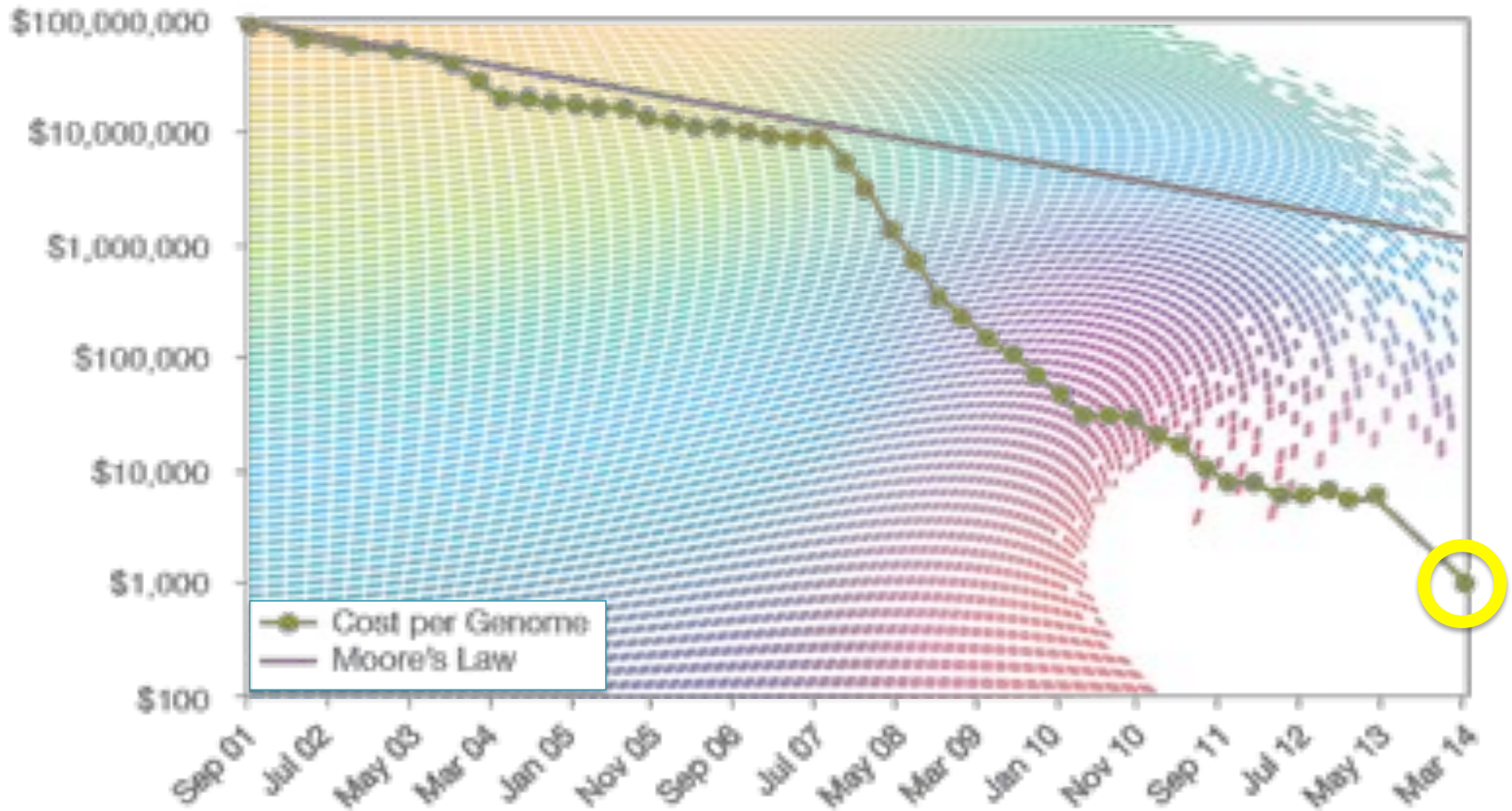
Metzker (2010) Nature Reviews Genetics 11:31-46
<http://www.youtube.com/watch?v=I99aKKHcxC4>

Cost per Genome



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

Cost per Genome



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

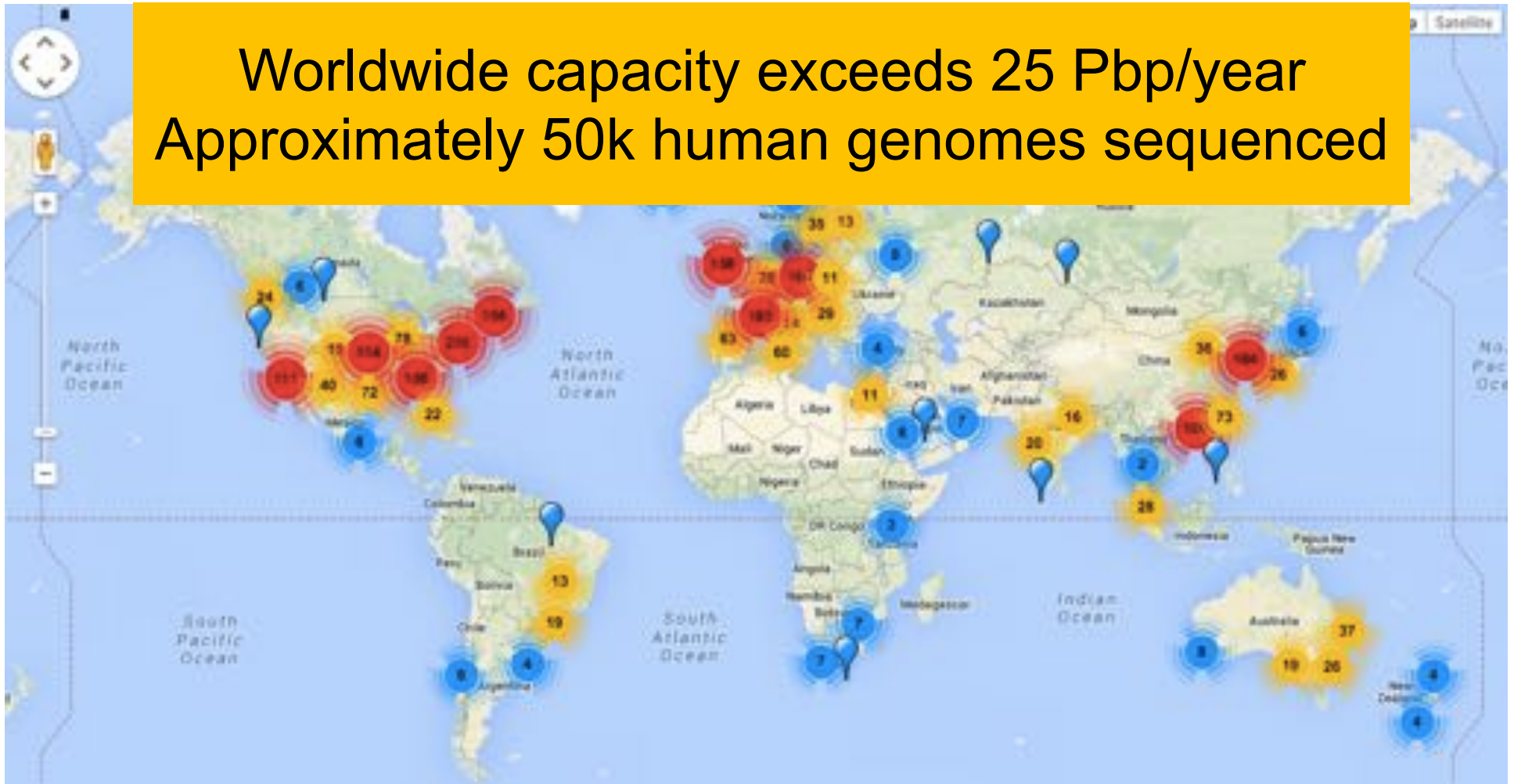
HiSeq X Ten



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



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



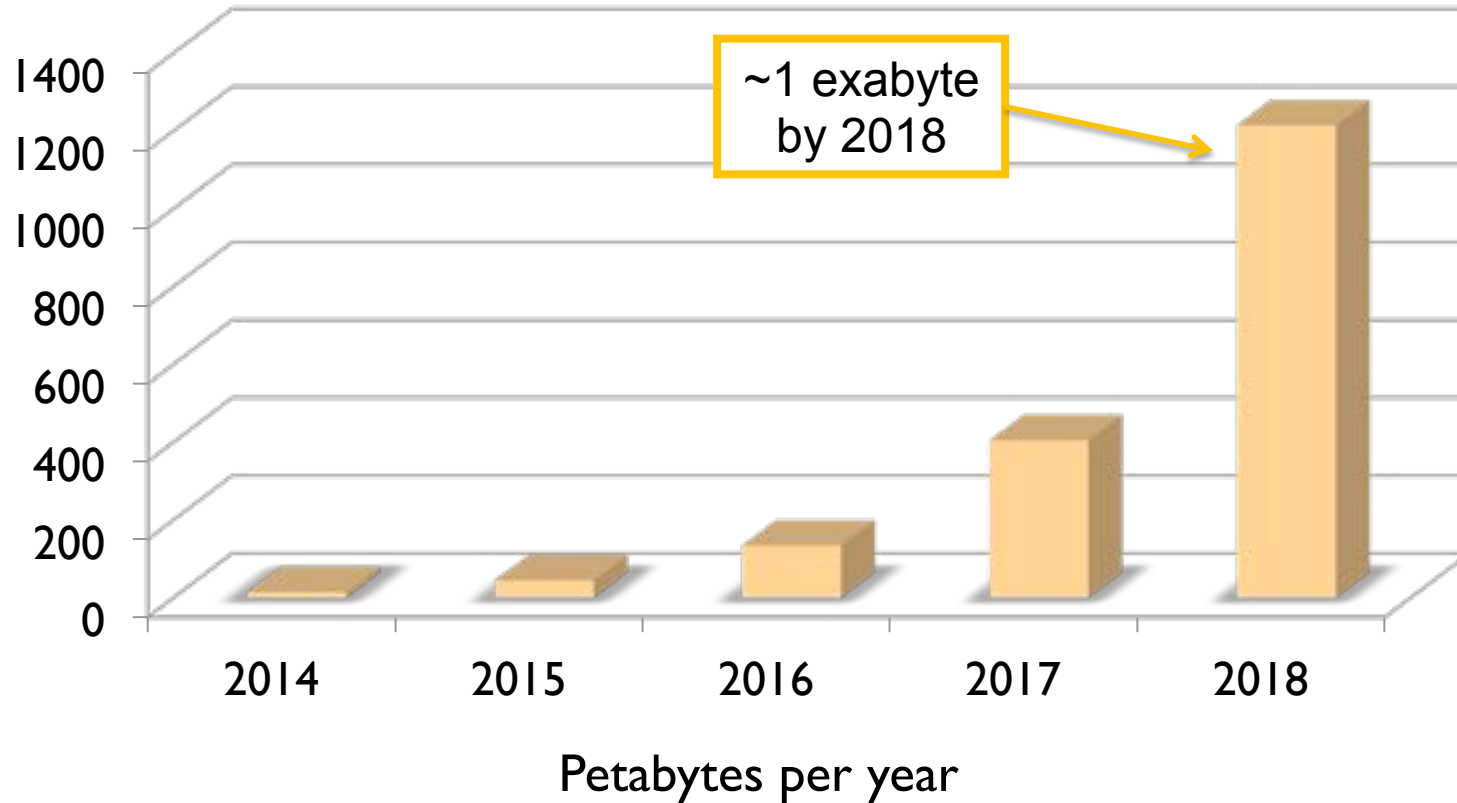
787 feet of DVDs
~1/6 of a mile tall



500 2 TB drives
\$500k

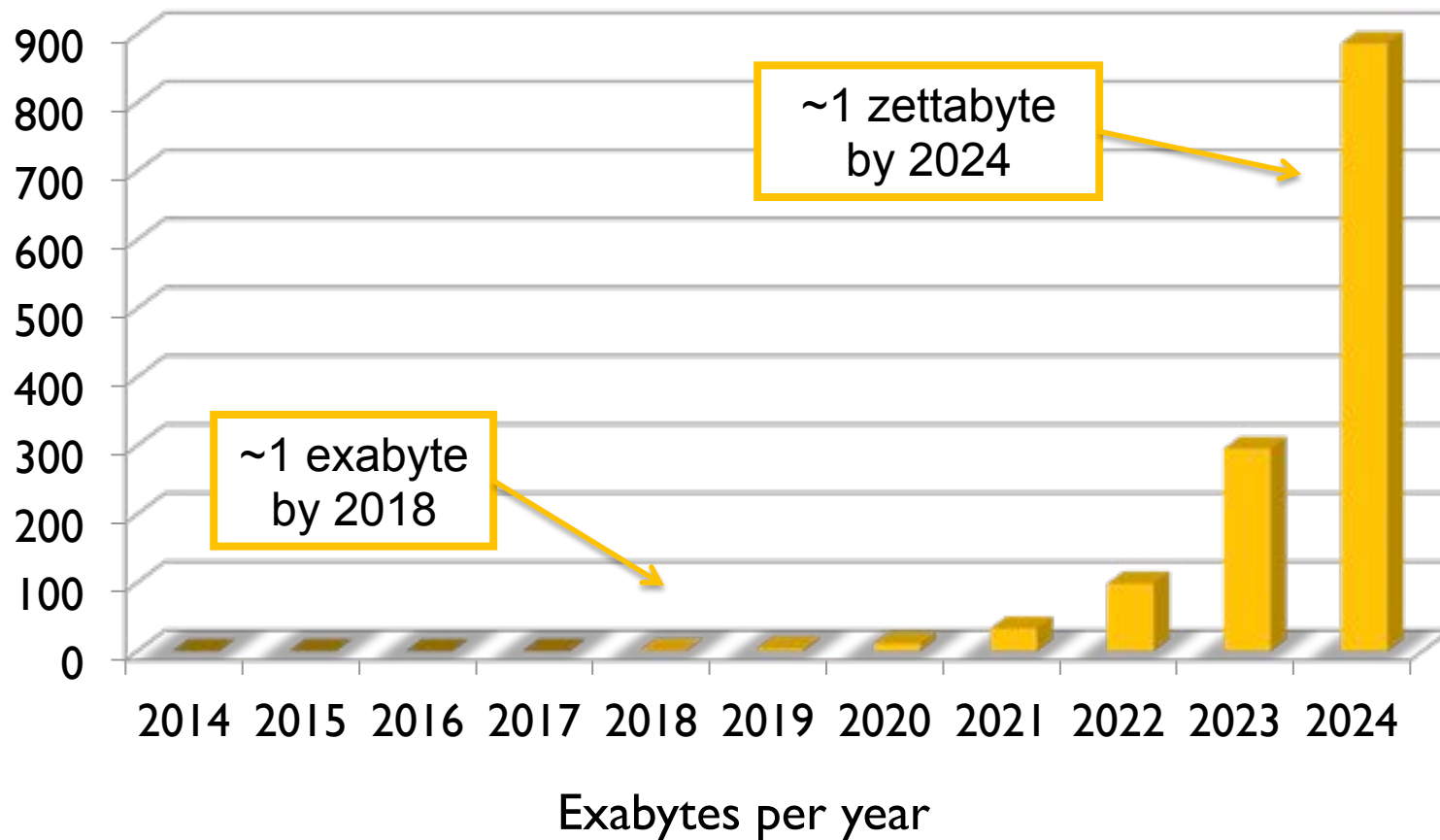
DNA Data Tsunami

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
~ 1/2 distance to moon



Both currently ~100Pb
And growing exponentially

Sequencing Centers 2014



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

Sequencing Centers 2024



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

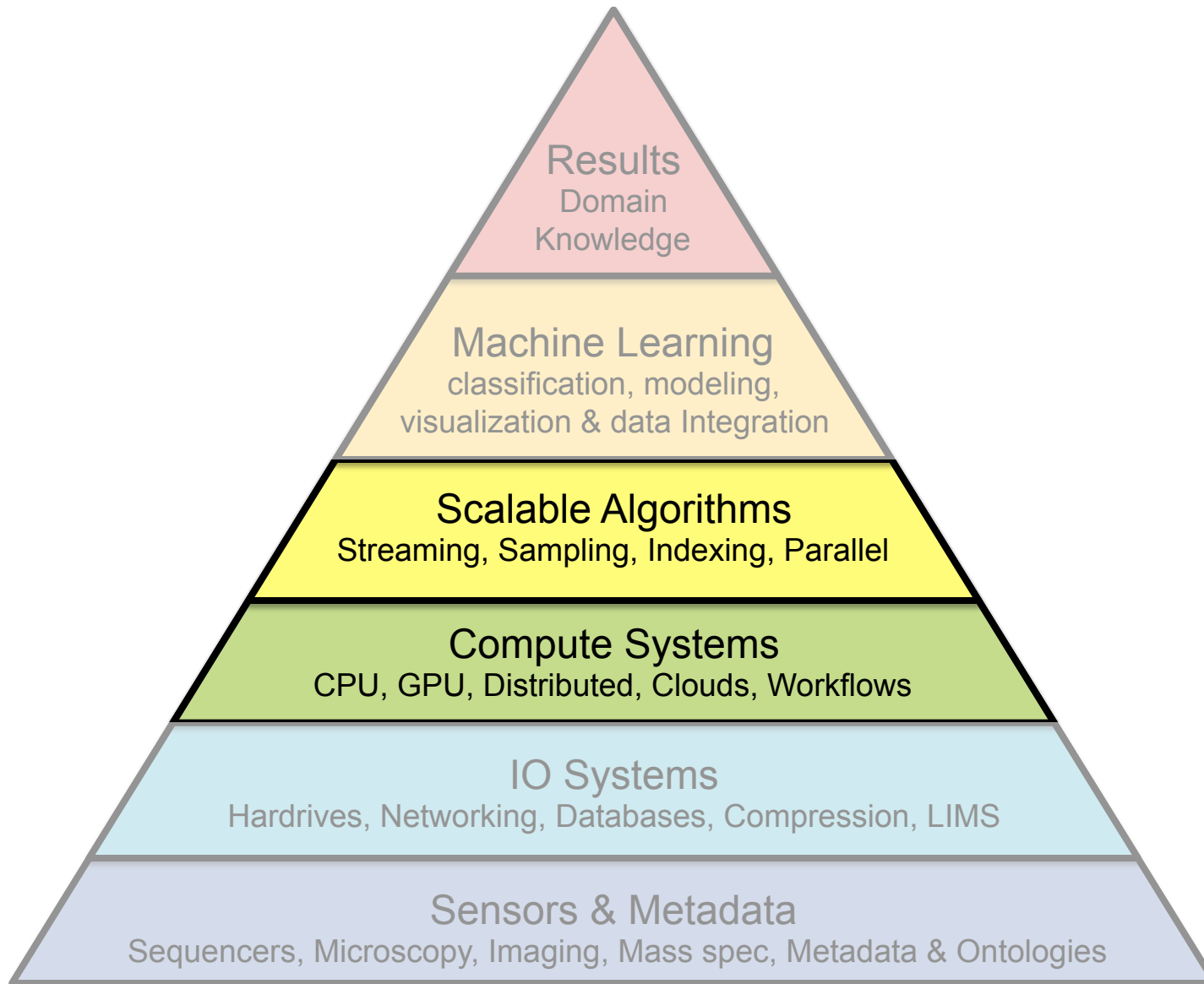
- Major efforts in human health and disease, agriculture, bioenergy
- Genomic information coupled with medical records and other medical data

But also widely distributed mobile sensors

- Schools, offices, sports arenas, 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



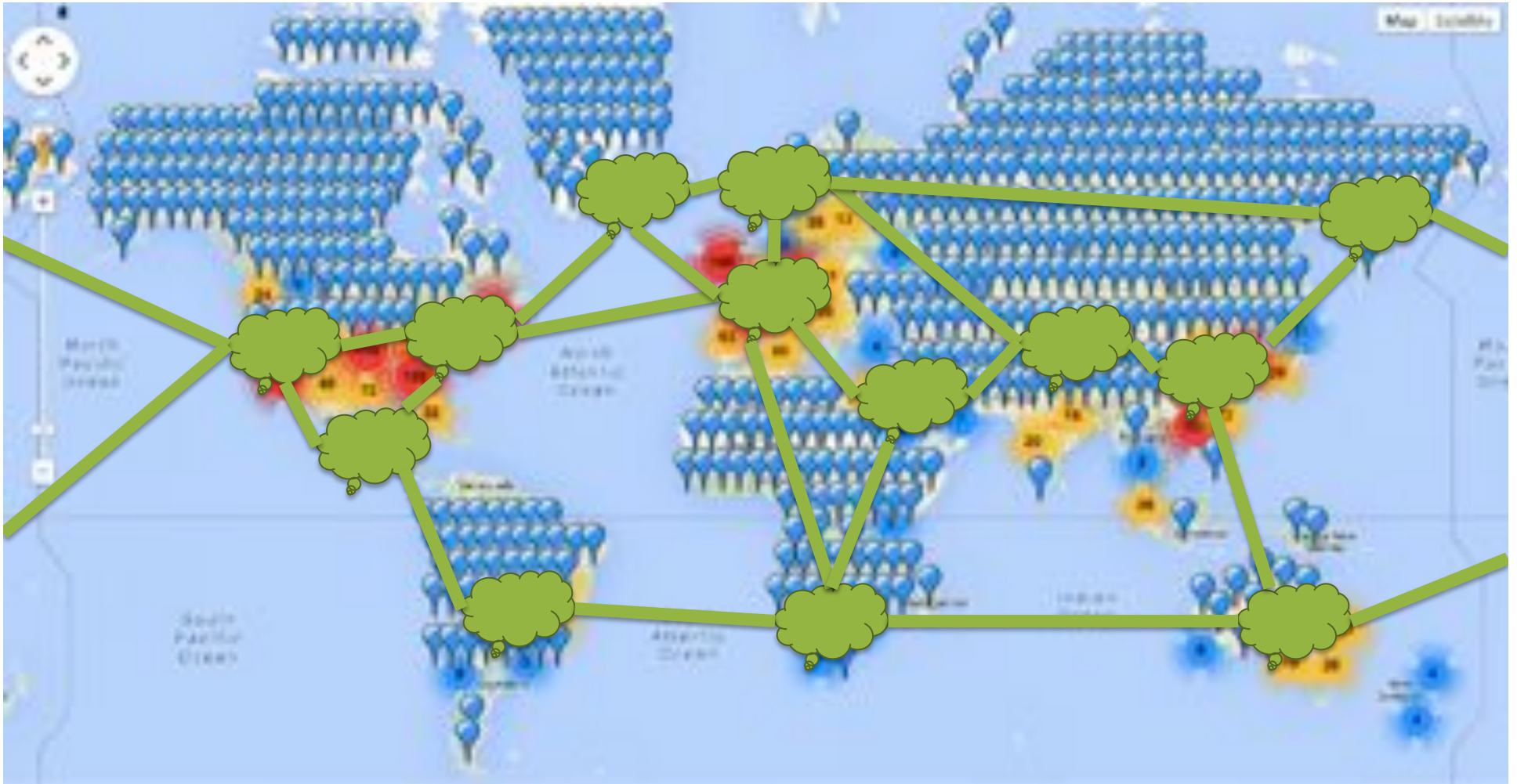
Informatics Centers 2024



The DNA Data Deluge

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

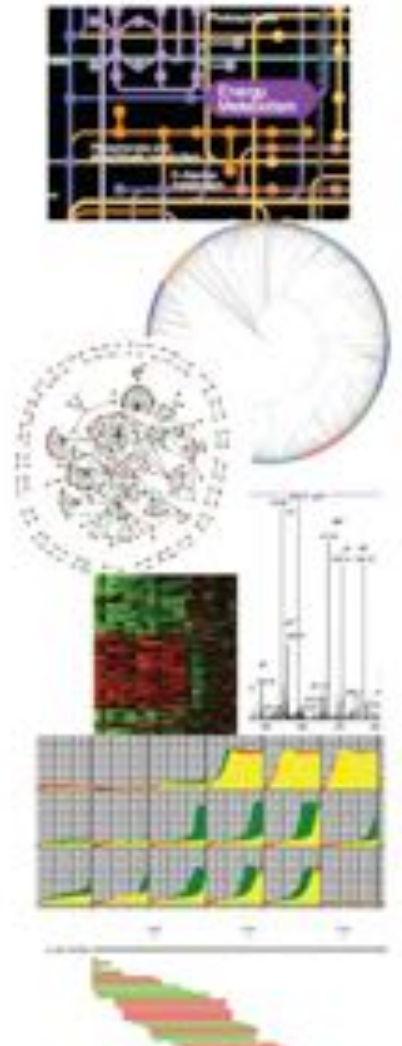
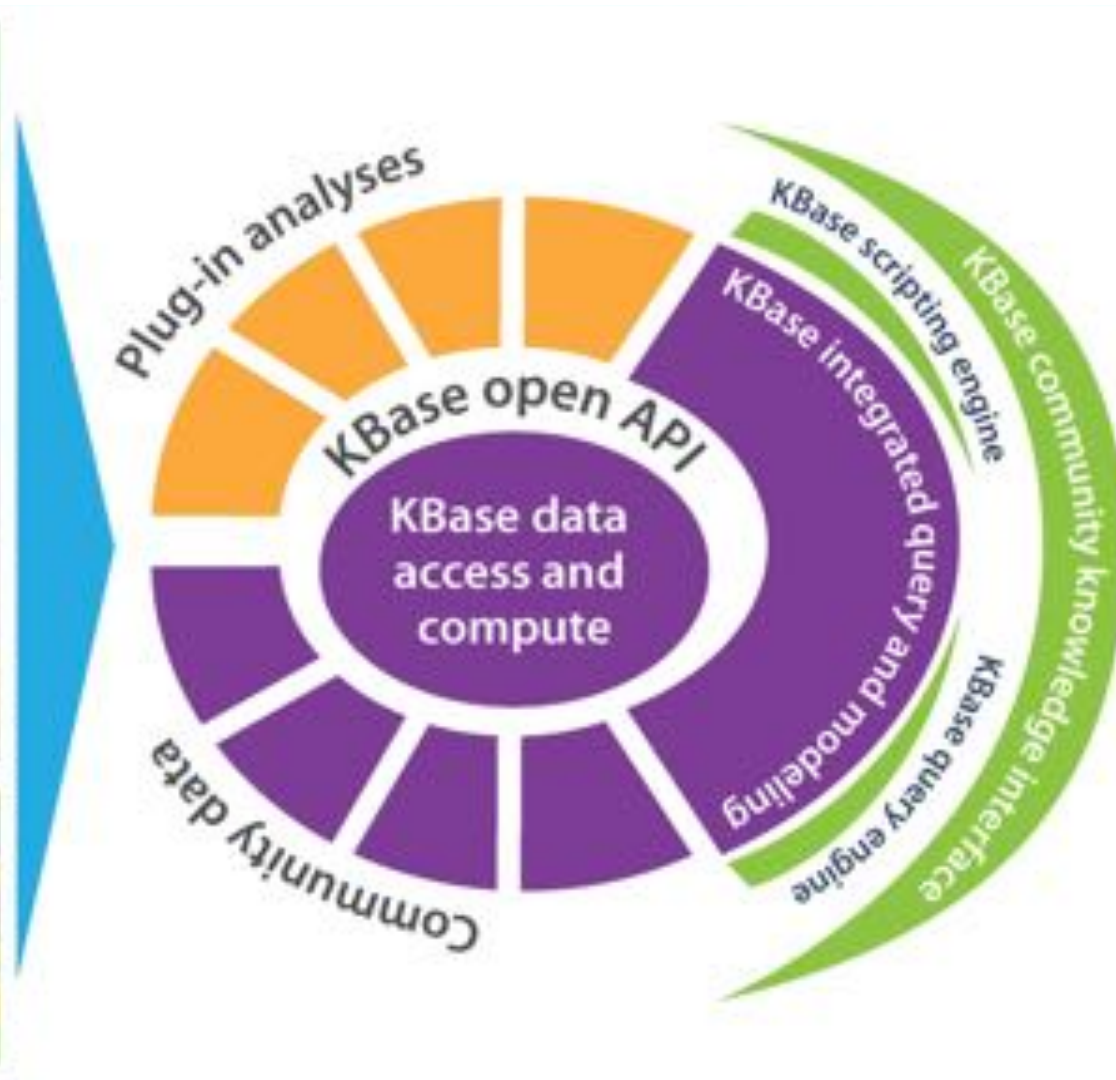
Informatics Centers 2014



The DNA Data Deluge

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

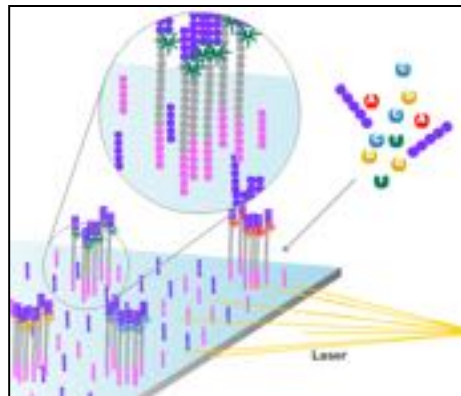
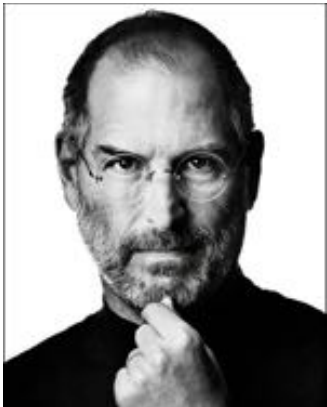
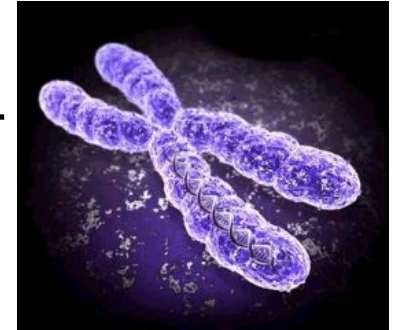
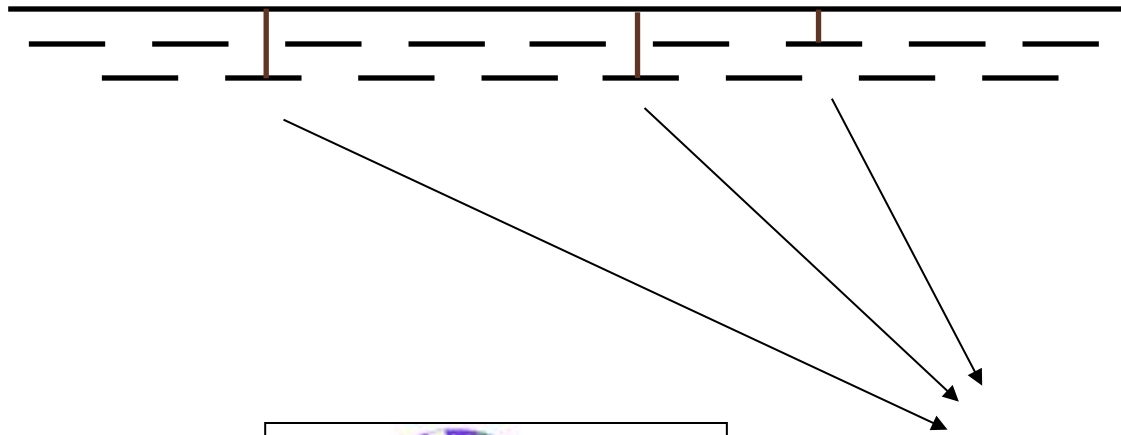
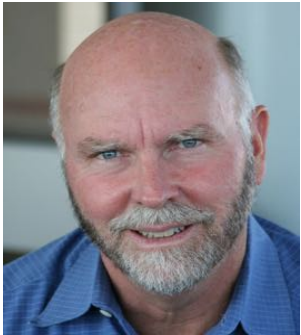
DOE Systems Biology Knowledgebase



<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!

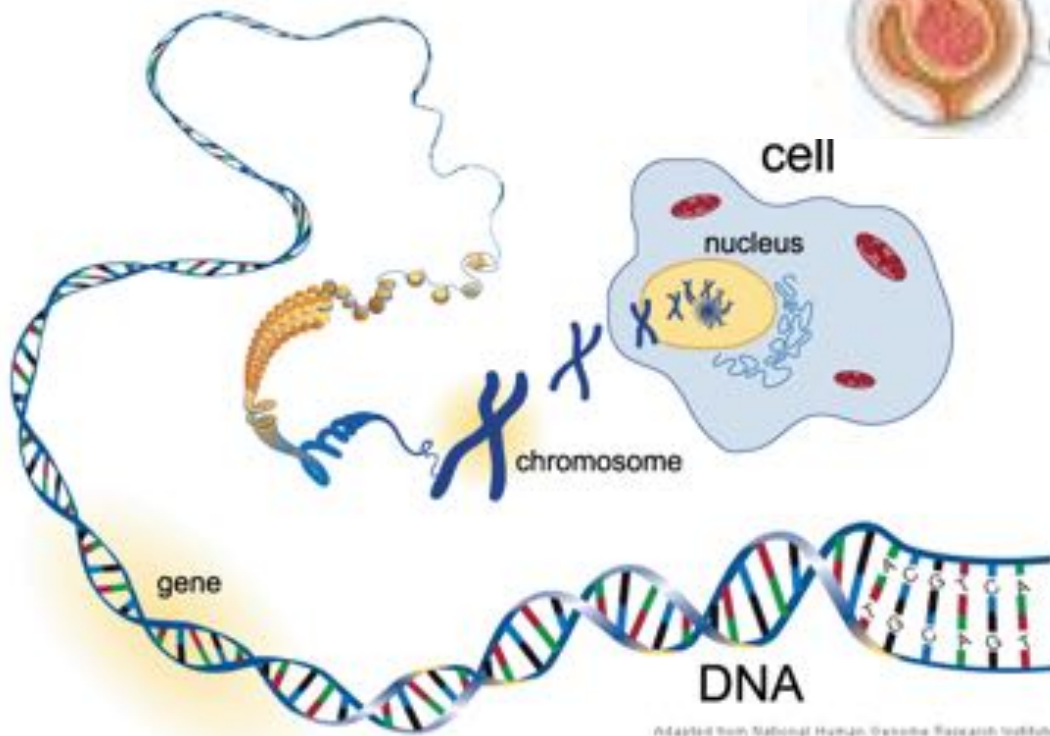


Searching for SNPs with Cloud Computing.

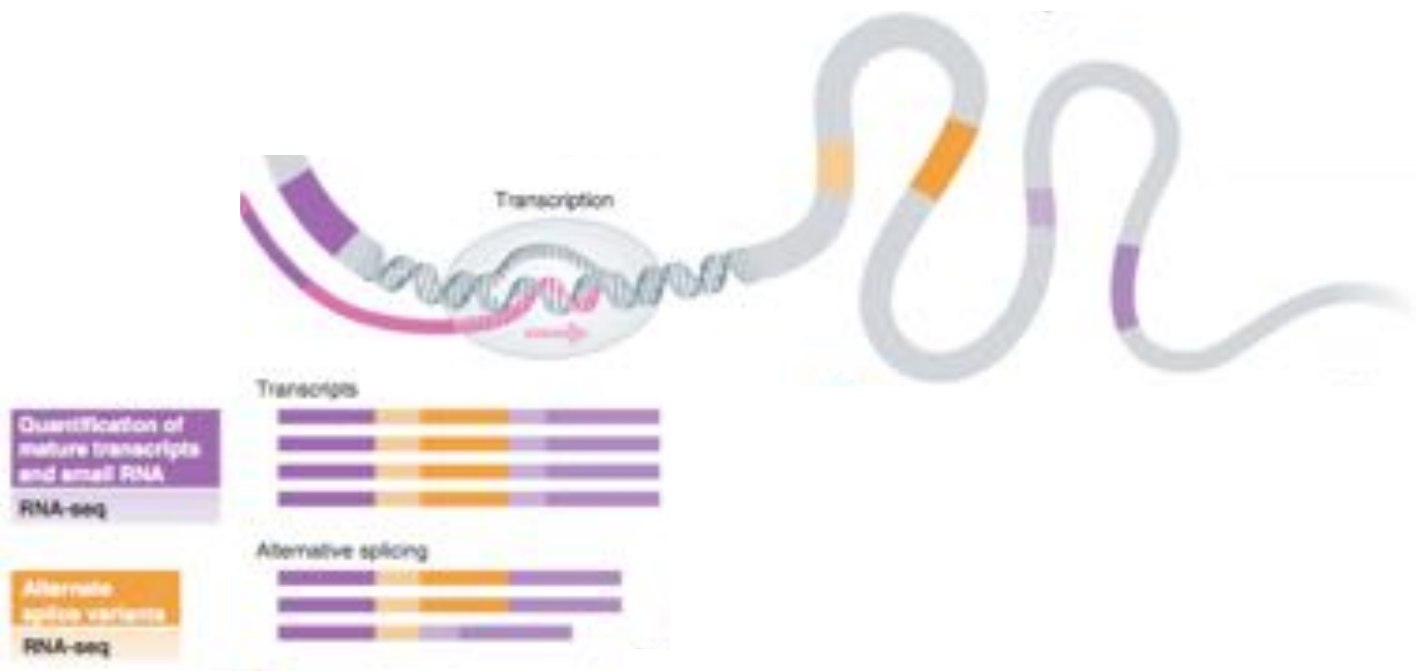
Langmead B, Schatz MC, Lin J, Pop M, Salzberg SL (2009) *Genome Biology*. **10**:R134

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



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

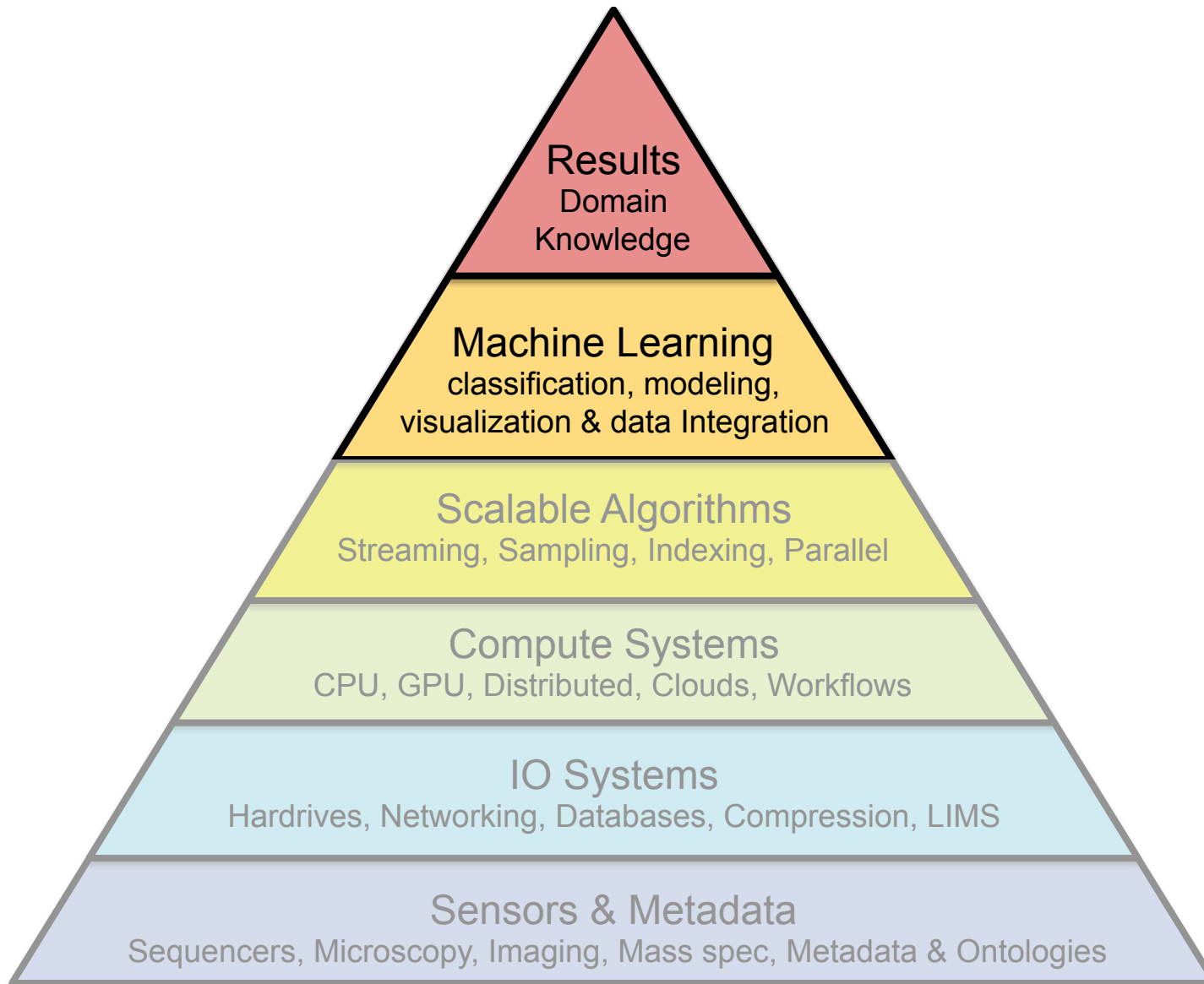
- Expect to see >1000 cores in a single computer
- Compute & IO needs to be considered together
- Rewriting efficient parallel software is complex and expensive

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?

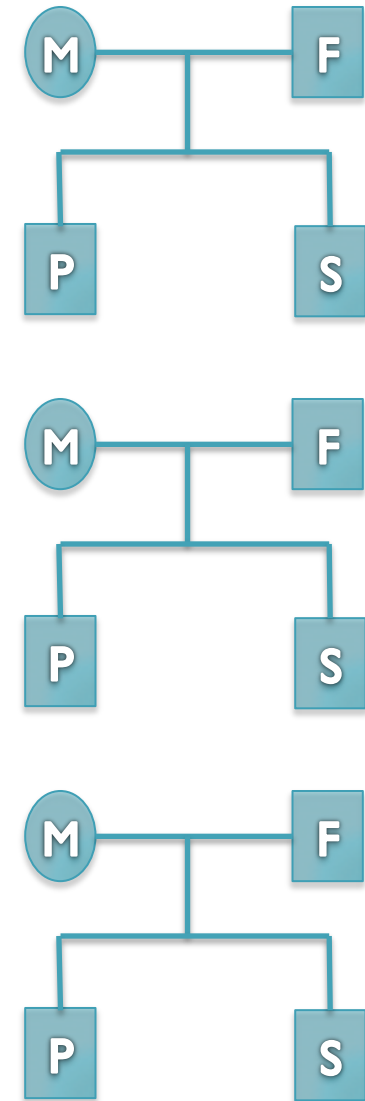
<http://www.autismspeaks.org/what-autism>

Searching for the genetic risk factors

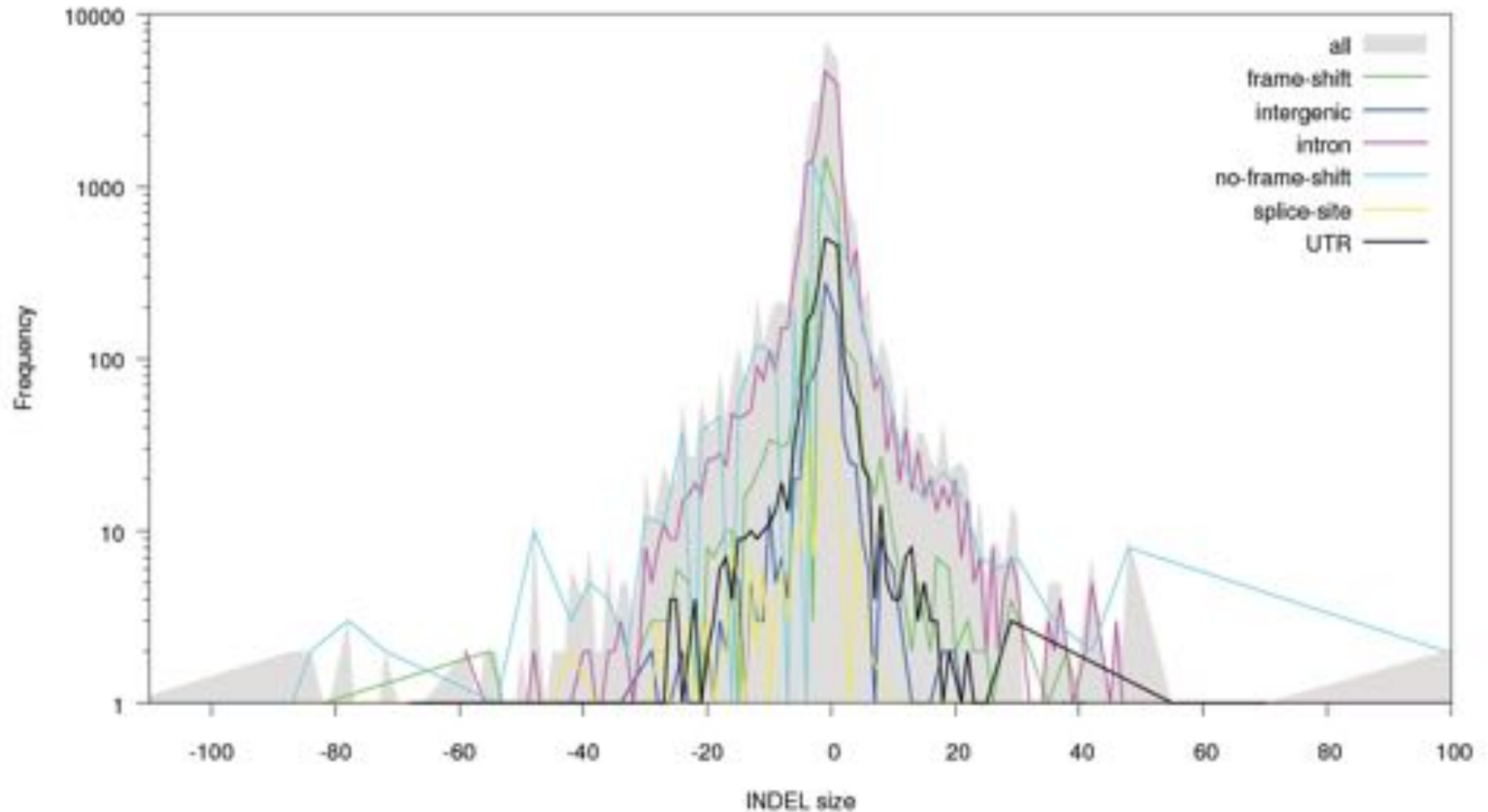
Search Strategy

- 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

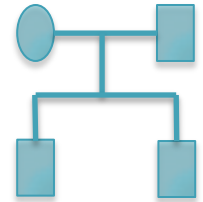


Constructed database of $>IM$ transmitted and de novo genetic mutations

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

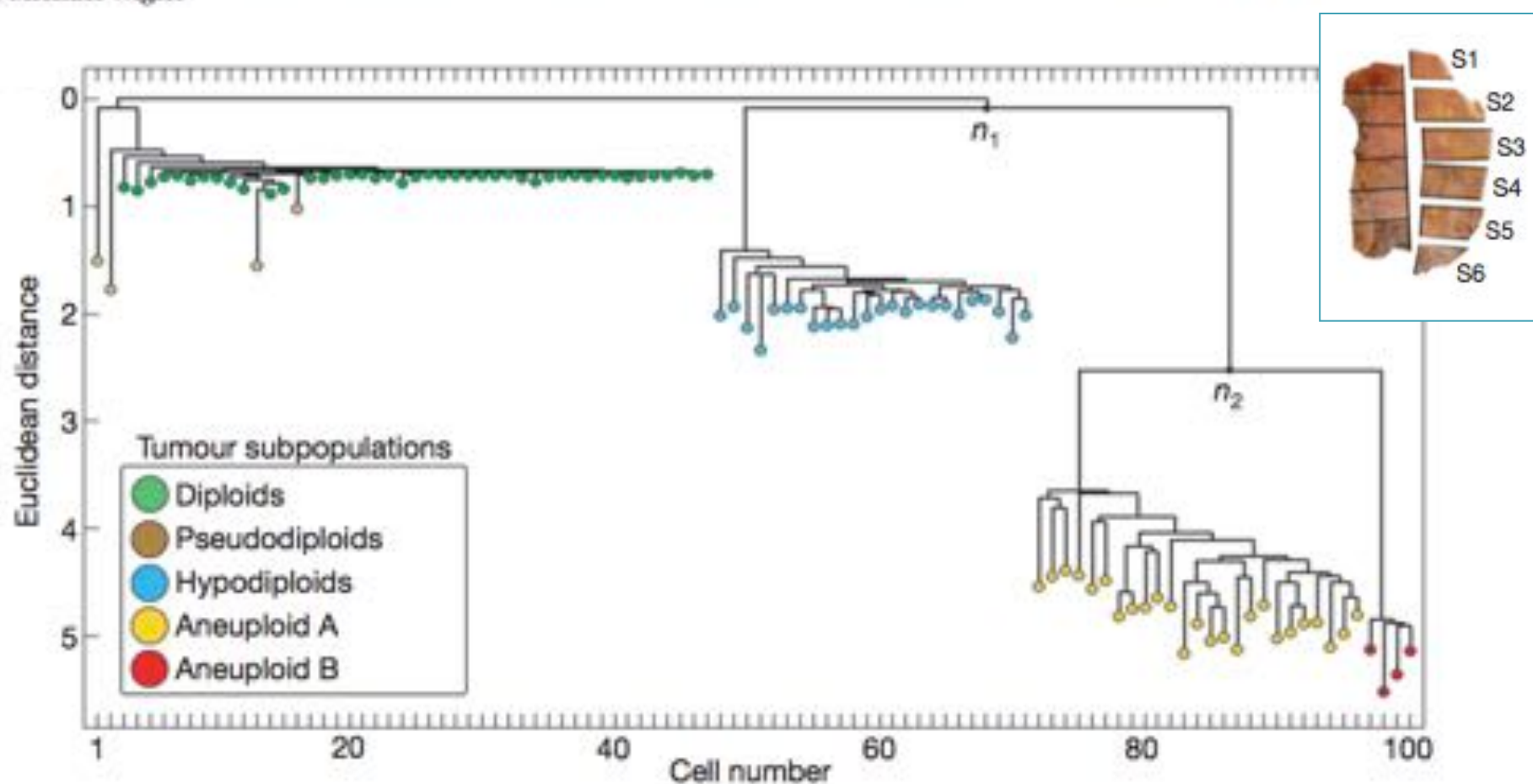
De novo Genetics of Autism

- 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^{1,2}, Jude Kendall¹, Jennifer Troge¹, Peter Andrews¹, Linda Rodgers¹, Jeanne McIndoo¹, Kerry Cook¹, Asya Stepansky¹, Dan Levy¹, Diane Esposito¹, Lakshmi Muthuswamy³, Alex Krasnitz¹, W. Richard McCombie¹, James Hicks¹ & Michael Wigler²



What makes us human?

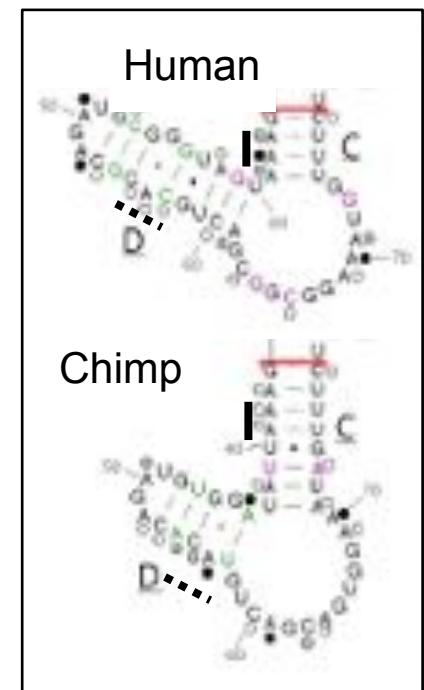
“Human Accelerated Regions”



human	TCATAGCGGTAGACCCAGTGCAGCGCGGAAATGGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA
chimp	TTATAGCGGTAGACACATGTCAGCAGTGGAAATAGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA
dog	TTATAGCGGTAGACACATGTCAGCGCGGAAACAGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA
mouse	TTATAGCGGTAGACACATGTCAGCGCGGAAATGGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA
rat	TTATAGCGGTAGACACATGTCAGCAGTGGAAATGGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA
chicken	TTATAGCGGTAGACACATGTCAGCAGTGGAAACAGTTTCTATCAAAATCAAAGTATTAGAGATTTTCCTCAAATTTCAAATTA

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)

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