

Computational Biomedical Research

Michael Schatz

August 30, 2021

Lecture I: 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/biomedicalresearch2021>

Course Discussions: <http://piazza.com>

Class Hours: Mon + Wed @ 3:00p – 3:50p in Hodson 211

Schatz Office Hours: TBD and by appointment

Das Office Hours: TBD and by appointment

Please try Piazza first!

TA: Samantha Zarate



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 2018/2020
- <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: 25% Due at 11:59pm a week later
Practice using the tools we are discussing
- 1 Exam: 30% In class (Tentatively 11/1)
Assess your performance, focusing on the methods
- 1 Class Project: 45% 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:***
 - A total of 96 hours (24×4) can be used to extend the deadline for assignments, but not the class project, without any penalty; after that time assignments will not be accepted

Course Webpage

The screenshot shows a GitHub repository page for 'schatzlab/biomedicalresearch2021'. The repository has 1 branch and 0 tags. The main branch has 9 commits. The repository is described as 'Course Materials for EN.601.452 / AS.020.415 Computational Biomedical Research & Advanced Biomedical Research'. It includes a Readme file and a CC-BY-NC-ND license. There are no releases or packages published.

Code Issues Pull requests Actions Projects Wiki Security Insights

main · 1 branch · 0 tags

Go to file Code

mschatz Merge branch 'main' of https://github.com/schatzlab/biomedical... 17 minutes ago 9 commits

assignments/assignment1 update links 17 minutes ago

policies add policies 39 minutes ago

.gitignore import schedule 1 hour ago

LICENSE initial commit 2 hours ago

README.md update schedule 1 hour ago

.config.yml Set theme jekyll-theme-cayman 21 minutes ago

README.md

EN.601.452 / AS.020.415 Computational Biomedical Research & Advanced Biomedical Research

Prof: Michael Schatz (mschatz@cs.jhu.edu)
TA: Samantha Zаратé (szarate@jhu.edu)
Class Hours: Monday + Wednesday @ 3:00 - 3:50p in Hodson 211
Schatz Office Hours: By appointment
Tentative Office Hours: TBD and by arrangement

Notifications Star Fork

About

Course Materials for EN.601.452 / AS.020.415 Computational Biomedical Research & Advanced Biomedical Research

Readme

CC-BY-NC-ND License

Releases

No releases published

Packages

No packages published

<https://github.com/schatzlab/biomedicalresearch2021>

Piazza

The screenshot shows a web browser window with the URL piazza.com/class/ksoxihnaqr2v6ggz?cid=6. The browser's address bar also displays 'EN.601.452'. The page title is 'EN.601.452'.

The Piazza interface includes a top navigation bar with links for 'Q & A', 'Resources', 'Statistics', and 'Manage Class'. On the right side of the top bar, there is a user profile for 'Michael Schatz'.

The main content area shows a note titled 'Welcome to Piazza!' by 'Instr. Welcome to Piazza!'. The note content 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 email), 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' and has a 'good note' button with a count of 0. It was updated 1 hour ago by Michael Schatz.

On the left sidebar, there is a list of recent posts:

- Instr. Welcome to Piazza! (8:55PM)
- Private Introduce Piazza to your stu... (8:53PM)
- Private Get familiar with Piazza (8:53PM)
- Private Tips & Tricks for a successf... (8:53PM)
- Welcome to Piazza! (8:53PM)

At the bottom of the page, there are sections for 'Average Response Time' (N/A) and 'Special Mentions' (There are no special mentions at this time). There are also links for 'Online Now' (1) and 'This Week' (1).

Copyright information at the bottom: Copyright © 2021 Piazza Technologies, Inc. All Rights Reserved. [Privacy Policy](#) [Copyright Policy](#) [Terms of Use](#) [Blog](#) [Report Bug](#)

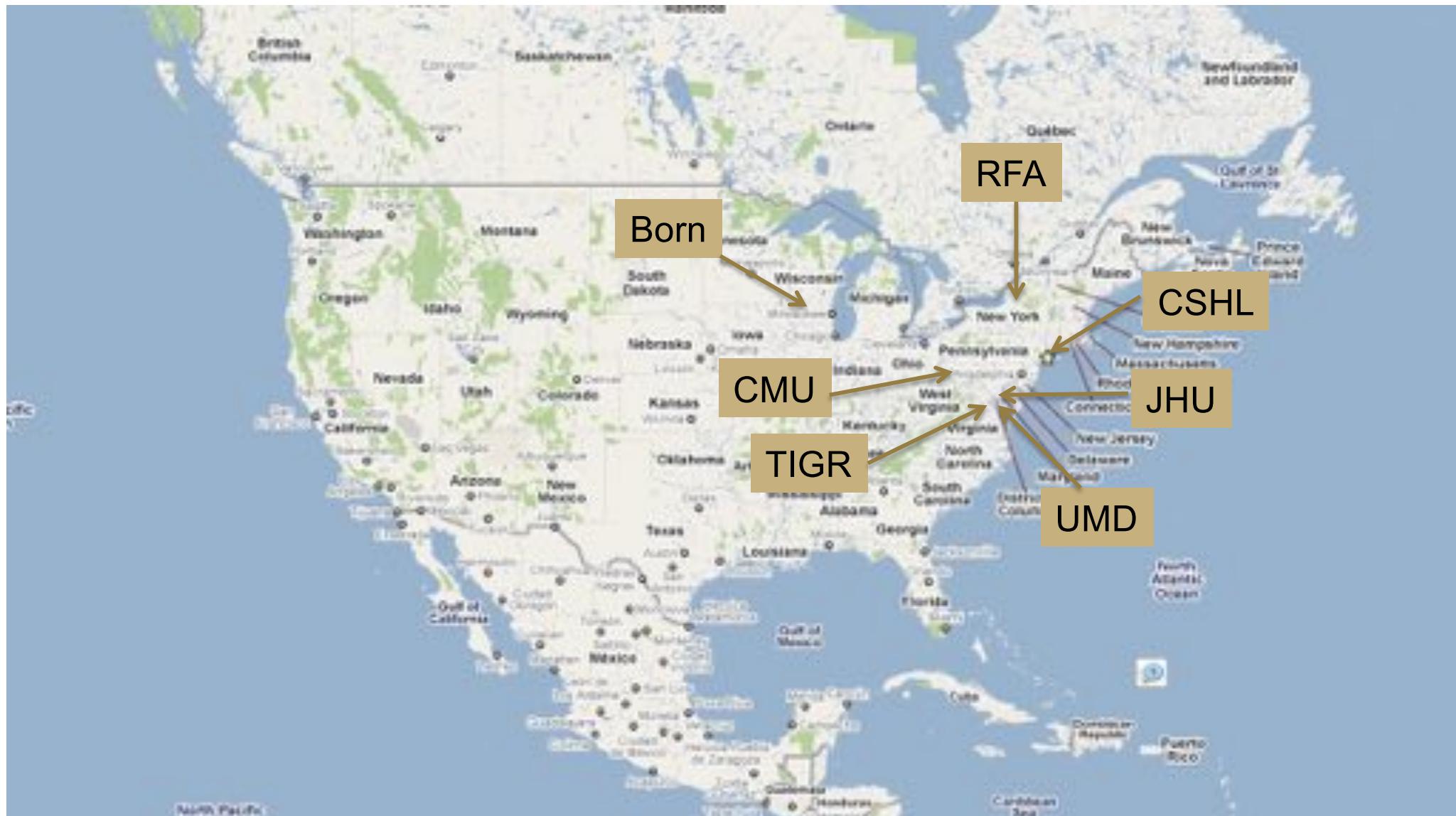
<http://piazza.com/jhu/fall2021/en601452/home>

GradeScope

The screenshot shows the GradeScope dashboard. At the top, there's a header bar with browser controls, a tab labeled "Dashboard | Gradescope", and a search bar. Below the header, the URL "gradescope.com" is visible. The main content area is titled "Your Courses". It displays two course sections: "Fall 2021" and "Spring 2021".
Fall 2021: Shows course "EN.601.452 Computational Biomedical Research" with "0 assignments".
Spring 2021: Shows course "EN.601.749 Applied Comparative Genomics" with "10 assignments".
A dashed box highlights a button labeled "Create a new course". At the bottom of the page, there's a teal footer bar with links for "Account", "Enroll in Course", and "Create Course +".

<https://www.gradescope.com/>
Entry Code:D5GDXP

A Little About Me



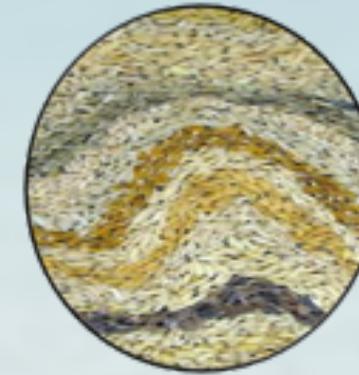
Schatzlab Overview



Human Genetics

Role of mutations
in disease

Aganezov *et al.* (2020)
Wang *et al.* (2019)



Agricultural Genomics

Genomes &
Transcriptomes

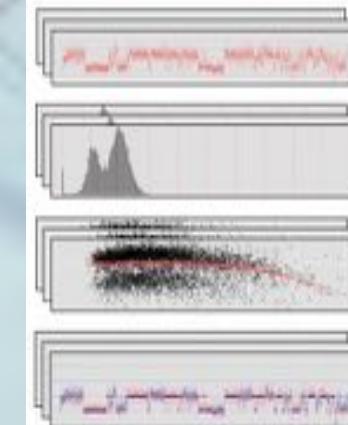
Alonge *et al.* (2020)
Soyk *et al.* (2019)



Algorithmics & Systems Research

Ultra-large scale
biocomputing

Kirsche *et al.* (2020)
Fang *et al.* (2018)



Biotechnology Development

Single Cell + Single
Molecule Sequencing

Kovaka *et al.* (2020)
Sedlazeck *et al.* (2018)

Earliest Genomics

Any Guesses?

Earliest Genomics



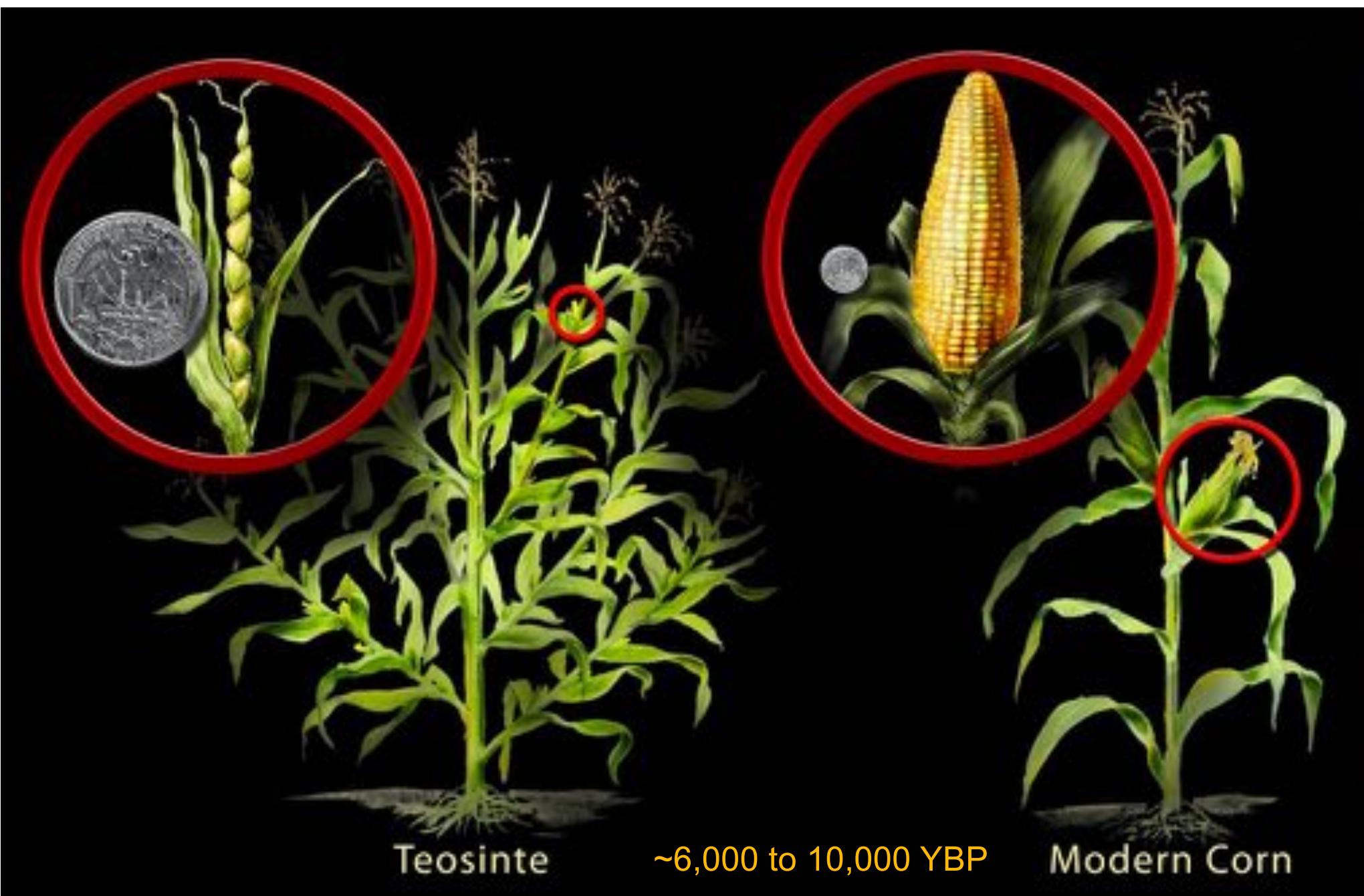
15,000 to 35,000 YBP

Earliest Genomics



~1,000 to 10,000 YBP

Earliest Genomics



Angiosperms (Flowering Plants)



~130 Ma

Discovery of Chromosomes

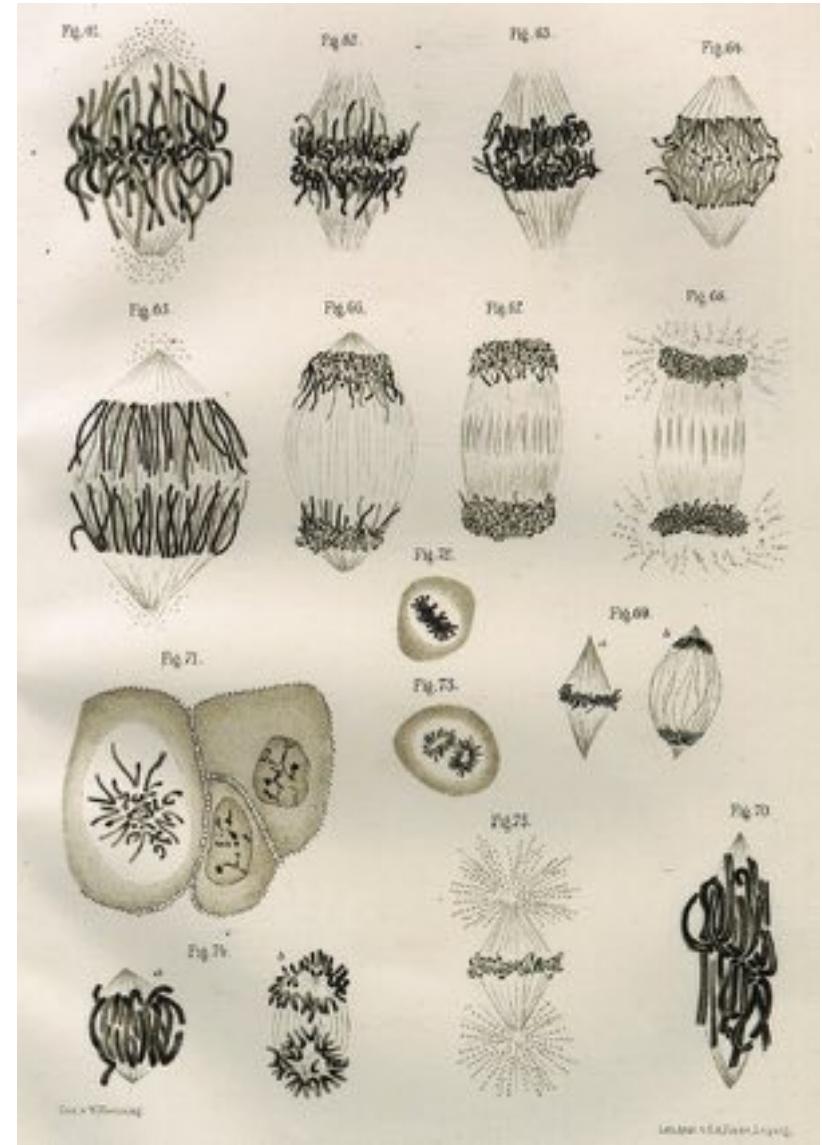
By the mid-1800s, microscopes were powerful enough to observe the presence of unusual structures called “chromosomes” that seemed to play an important role during cell division.

It was only possible to see the chromosomes unless appropriate stains were used

“Chromosome” comes from the Greek words meaning “color body”

Today, we have much higher resolution microscopes, and a much richer varieties of dies and dying techniques so that we can visualize particular sequence elements.

When you see something unexpected that you think might be interesting, give it a name



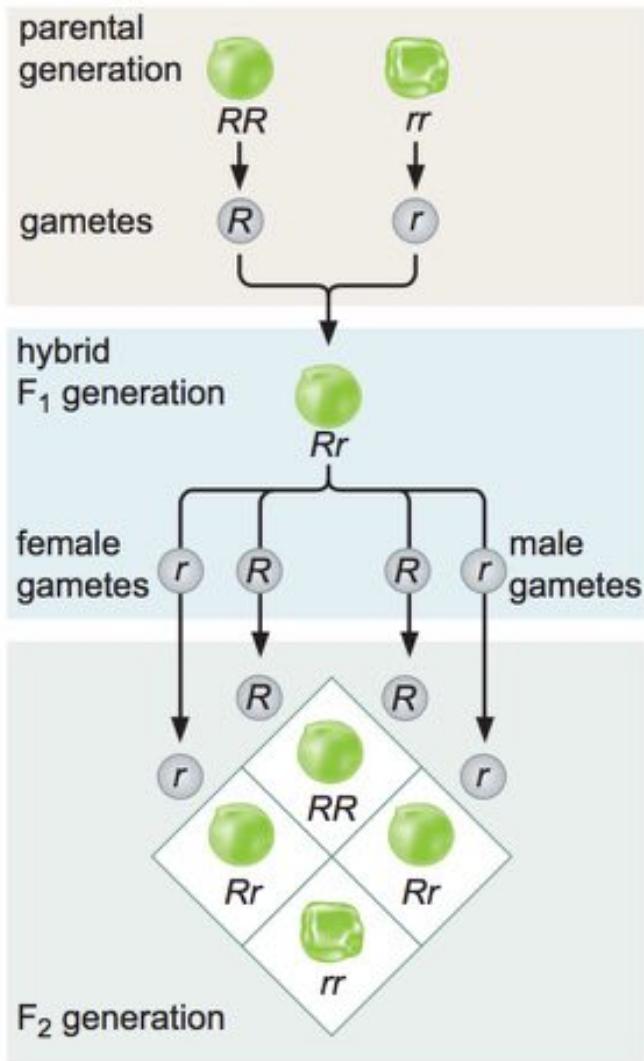
Drawing of mitosis by Walther Flemming.

Flemming, W. Zellsubstanz, Kern und Zelltheilung (F. C. W. Vogel, Leipzig, 1882).

The “first” quantitative biologist

Any Guesses?

Laws of Inheritance



Seed		Flower		Pod		Stem	
Form	Cotyledons	Color		Form	Color	Place	Size
Grey & Round	Yellow	White		Full	Yellow	Axial pods, Flowers along	Long (6-7ft)
White & Wrinkled	Green	Violet		Constricted	Green	Terminal pods, Flowers top	Short & 1ft)
	1	2	3	4	5	6	7

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

Observations of 29,000 pea plants and 7 traits

Generation	in Verhältniss gestellt:			
	A	Aa	a	$A : Aa : a$
1	1	2	1	1 : 2 : 1
2	6	4	6	3 : 2 : 3
3	28	8	28	7 : 2 : 7
4	120	16	120	15 : 2 : 15
5	496	32	496	31 : 2 : 31
n				$2^n - 1 : 2 : 2^n - 1$

Versuche über Pflanzen-Hybriden. Verh. Naturforsch (Experiments in Plant Hybridization)
 Mendel, G. (1866). Ver. Brünn 4: 3–47 (in English in 1901, J. R. Hortic. Soc. 26: 1–32).

The first genetic map

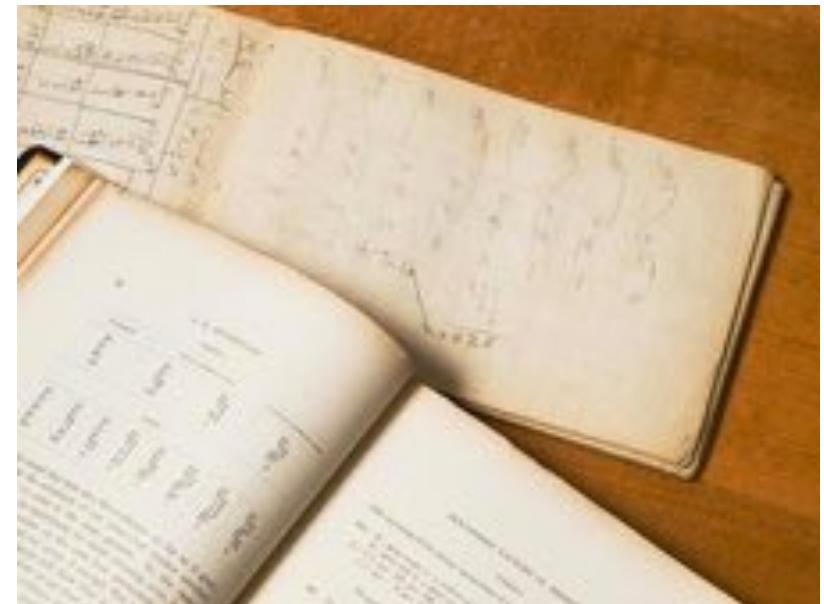
Mendel's Second Law (The Law of Independent Assortment) states alleles of one gene sort into gametes independently of the alleles of another gene: ***Pr(smooth/wrinkle) is independent of Pr(yellow/green)***

Morgan and Sturtevant noticed that the probability of having one trait given another was **not** always 50/50— those traits are ***genetically linked***

Sturtevant realized the probabilities of co-occurrences could be explained if those alleles were arranged on a linear fashion: traits that are most commonly observed together must be located closest together



The Linear Arrangement of Six Sex-Linked Factors in Drosophila as shown by their mode of Association
Sturtevant, A. H. (1913) *Journal of Experimental Zoology*, 14: 43-59



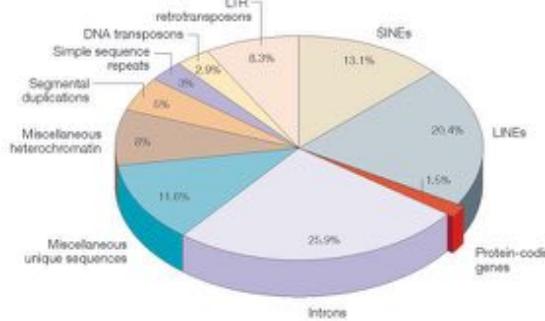
<http://www.caltech.edu/news/first-genetic-linkage-map-38798>

Jumping Genes



Previously, genes were considered to be stable entities arranged in an orderly linear pattern on chromosomes, like beads on a string

Careful breeding and cytogenetics revealed that some elements can move (cut-and-paste, DNA transposons) or copy itself (copy-and-paste, retrotransposons)



(Gregory, 2005, Nature Reviews Genetics)

(Much) later analysis revealed that nearly 50% of the human genome is composed of transposable elements, including LINE and SINE elements (long/short interspersed nuclear elements) which can occur in 100k to 1M copies

“The genome is a graveyard of ancient transposons”

The origin and behavior of mutable loci in maize.

McClintock, B. (1950) PNAS. 36(6):344–355.

Nobel Prize in Physiology or Medicine in 1983

Discovery of the Double Helix

No. 4324 April 25, 1953

NATURE

737

equipment, and to Dr. G. E. R. Dobson, and the captain and officers of H.M.S. *Discovery II* for their part in making the observations.

* Doring, F. B., Dothie, E., and Evans, W., *Phil. Mag.*, **38**, 145 (1933).

† Argent-Nigges, M. S., *Nat. Natl. Bur. Stand. Ser. C*, **29**, 109 (1948).

‡ On 424, W. R. Woods-Riley Papers in Phys. Section, Kongr. 21 (1936).

§ Elson, T. W., *Advanc. Med. Electron. Phys.* (Baltimore), **10** (1960).

MOLECULAR STRUCTURE OF NUCLEIC ACIDS

A Structure for Deoxyribose Nucleic Acid

We wish to suggest a structure for the salt of deoxyribose nucleic acid (D.N.A.). This structure has novel features which are of considerable biological interest.

A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us as an advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure is unsatisfactory for two reasons: (1) We believe that the material which gives the X-ray diagrams is the salt, not the free acid. Without the acidic hydrogen atoms it is not clear where forces would hold the structure together, especially as the negatively charged phosphates near the axis will repel each other. (2) Some of the van der Waals distances appear to be too small.

Another three-chain structure has also been suggested by Fraser (in the press). In his model the phosphates are on the outside and the bases on the inside, linked together by hydrogen bonds. This structure is described as rather ill-defined, and for this reason we shall not comment on it.

We wish to put forward a radically different structure for the salt of deoxyribose nucleic acid. This structure has two helical chains each rolled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining β -D-deoxyribose residues with $2^{\circ}, 2'$ linkages. The two chains (but not their bases) are related by a dyad perpendicular to the fibre axis. Both chains follow right-handed helices, but owing to the dyad the sequences of the atoms in the two chains run in opposite directions. Each chain loosely resembles Franklin's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The conformation of the sugar and the atoms near it is close to Pernberg's "standard configuration", the sugar being roughly perpendicular to the attached base. There

is a residue on each chain every 3-4 Å. in the α -direction. We have assumed an angle of 36° between adjacent residues in the same chain, so that the structure repeats after 10 residues on each chain, that is, after 34 Å. The distance of a phosphorus atom from the fibre axis is 19 Å. As the phosphates are on the outside, cations have easy access to them.

The structure is an open one, and its water content is rather high. At lower water contents we would expect the bases to tilt so that the structure could become more compact.

The novel feature of the structure is the manner in which the two chains are held together by the purine and pyrimidine bases. The planes of the bases are perpendicular to the fibre axis. They are joined together in pairs, a single base from one chain being hydrogen-bonded to a single base from the other chain, so that the two lie side by side with identical α -co-ordinates. One of the pair must be a purine and the other a pyrimidine for bonding to occur. The hydrogen bonds are made as follows: purine position 1 to pyrimidine position 4; purine position 6 to pyrimidine position 4.

If it is assumed that the bases only occur in the structure in the most plausible tautomeric forms (that is, with the keto rather than the enol configurations) it is found that only specific pairs of bases can bond together. These pairs are: adenine (purine) with thymine (pyrimidine), and guanine (purine) with cytosine (pyrimidine).

In other words, if an adenine forms one member of a pair, on either chain, then on these assumptions the other member must be thymine; similarly for guanine and cytosine. The sequence of bases on a single chain does not appear to be restricted in any way. However, if only specific pairs of bases can bond, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

It has been found experimentally^{2,3} that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribonucleic acid.

It is probably impossible to build this structure with a ribose sugar in place of the deoxyribose, as the extra oxygen atom would make too close a van der Waals contact.

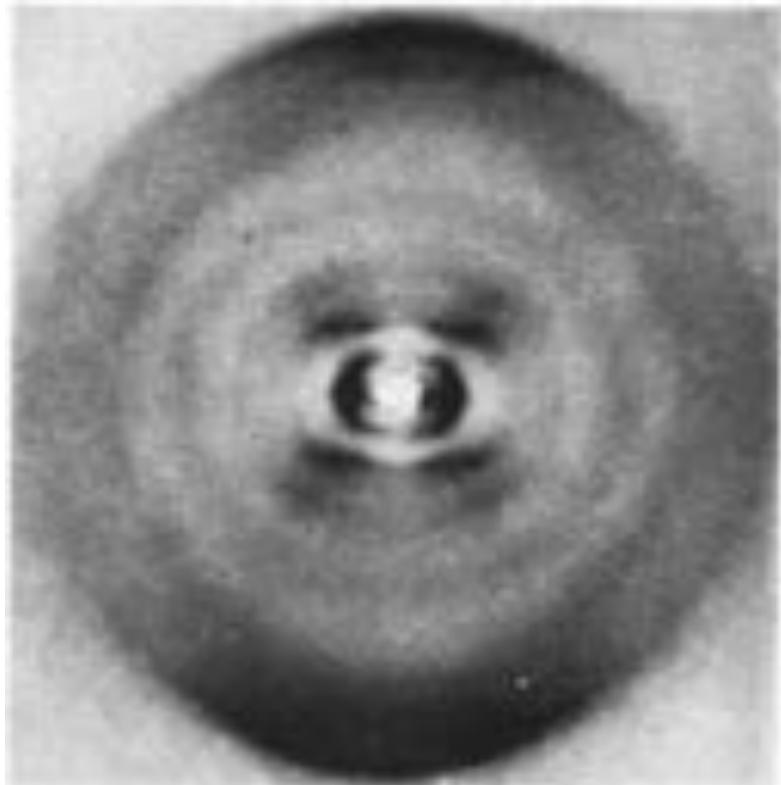
The previously published X-ray data^{4,5} on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. So far as we can tell, it is roughly compatible with the experimental data, but it must be regarded as unproved until it has been checked against more exact results. Some of these are given in the following communications. We were not aware of the details of the results presented there when we devised our structure, which rests mainly though not entirely on published experimental data and stereochemical arguments.

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material.

ACKNOWLEDGEMENTS

It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material. Full details of the structure, including the con-

This figure is greatly exaggerated. The two ribbons representing the two chains follow the horizontal axis; the vertical ends of the pairs of bases follow the vertical right-angle axis. The vertical axis marks the fibre axis.



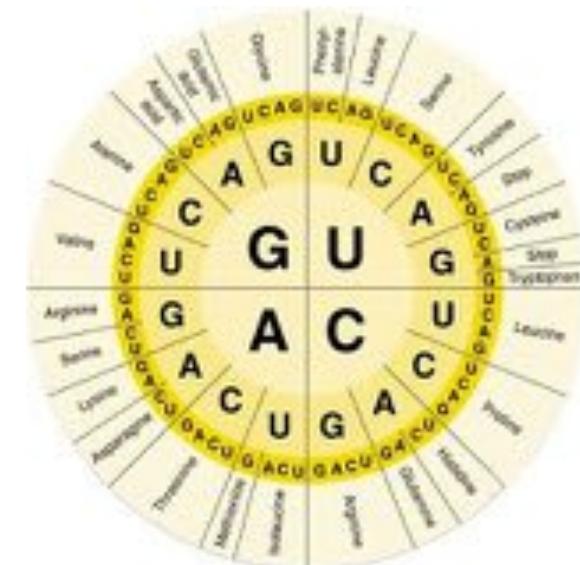
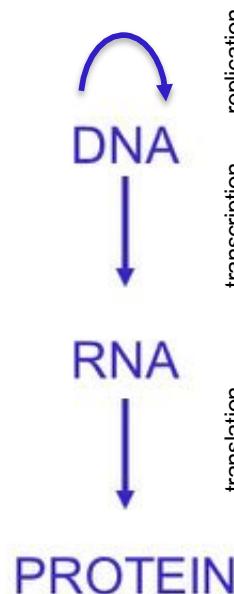
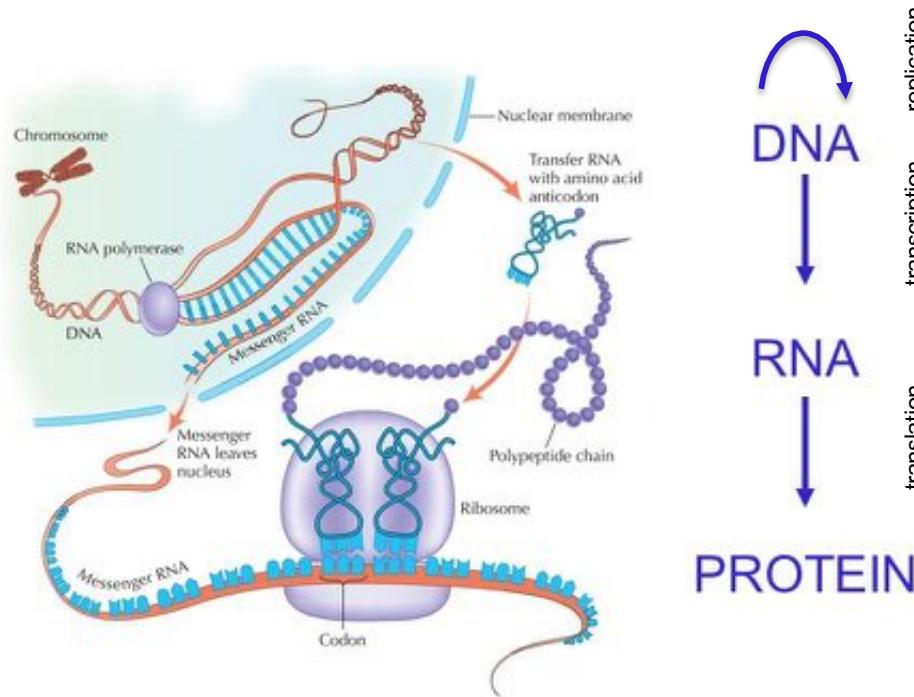
Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid

Watson JD, Crick FH (1953). *Nature* 171: 737–738.

Nobel Prize in Physiology or Medicine in 1962

Central Dogma of Molecular Biology

“Once 'information' has passed into protein it cannot get out again. In more detail, the transfer of information **from nucleic acid to nucleic acid, or from nucleic acid to protein may be possible**, but transfer from protein to protein, or from protein to nucleic acid is impossible. Information means here the precise determination of sequence, either of bases in the nucleic acid or of amino acid residues in the protein”

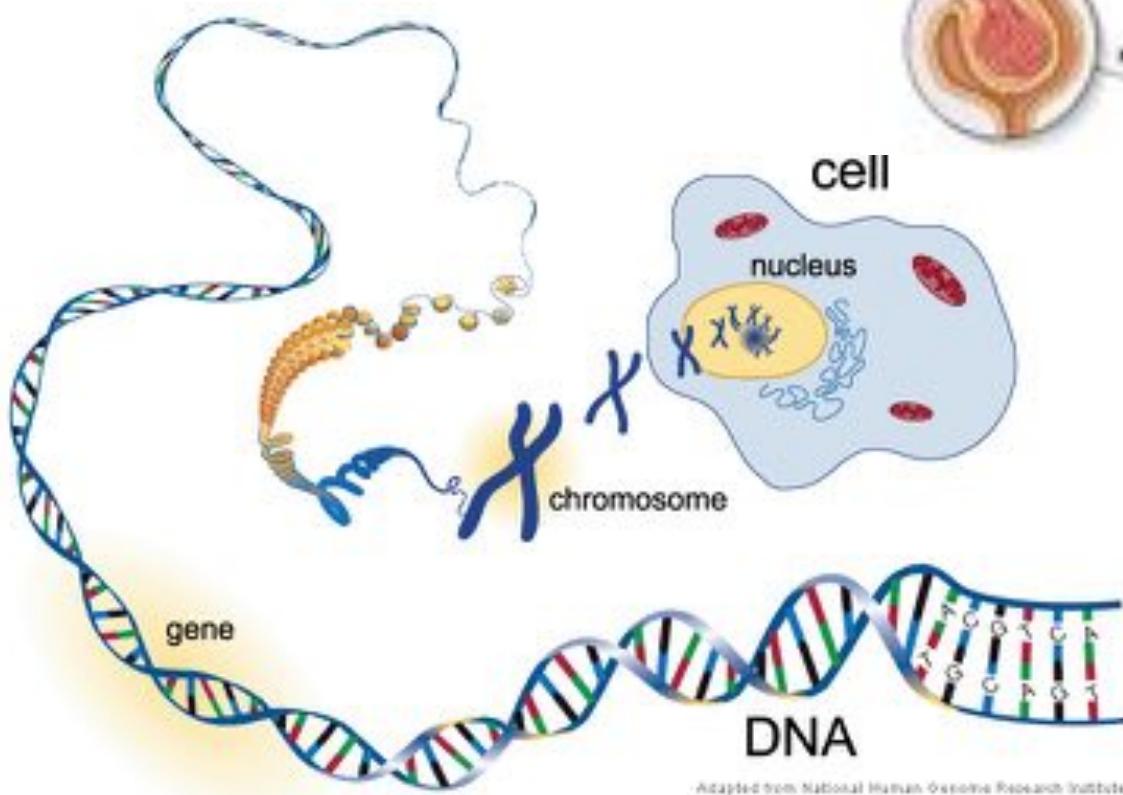


On Protein Synthesis

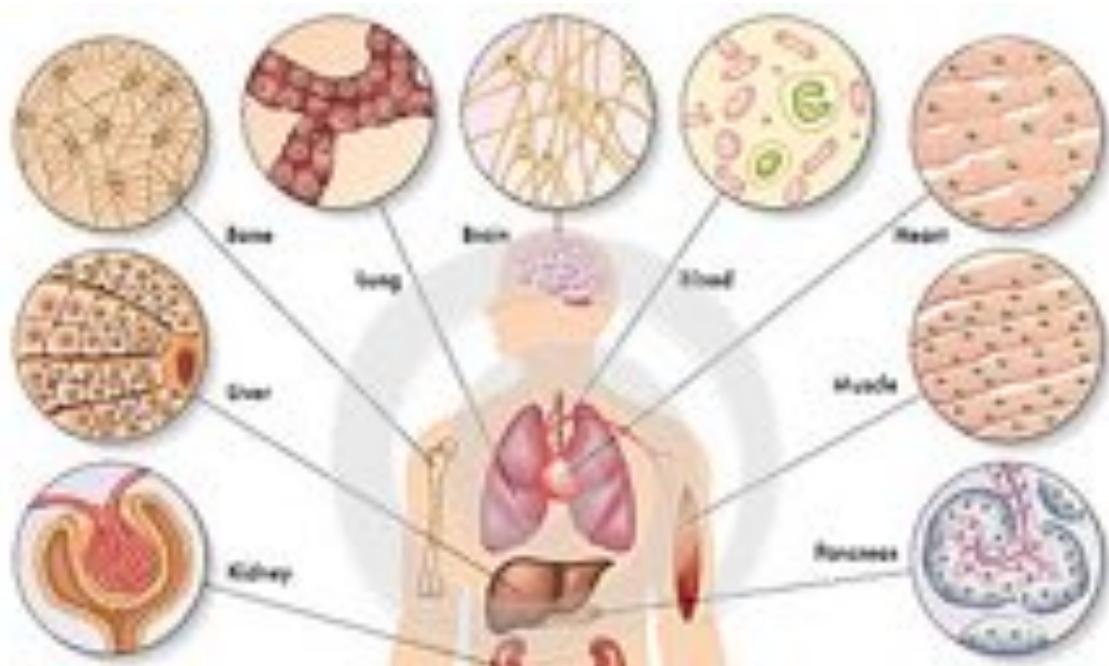
Crick, F.H.C. (1958). *Symposia of the Society for Experimental Biology* pp. 138–163.

One Genome, Many Cell Types

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



Adapted from National Human Genome Research Institute



Your body has a few hundred (thousands?) major cell types, largely defined by the gene expression patterns

Milestones in Genomics: Zeroth Generation Sequencing

Nature Vol. 265 February 24 1977 687

articles

Nucleotide sequence of bacteriophage Φ X174 DNA

F. Sanger, G. M. Air*, B. G. Barrell, N. L. Brown*, A. R. Coulson, J. C. Fiddes,
C. A. Hutchison III[†], P. M. Slocombe[‡] & M. Smith^{*}

MRC Laboratory of Molecular Biology, Hills Road, Cambridge CB2 2QE, UK

A DNA sequence for the genome of bacteriophage Φ X174 of approximately 5,375 nucleotides has been determined using the rapid and simple 'plus and minus' method. The sequence identifies many of the features responsible for the production of the proteins of the nine known genes of the organism, including initiation and termination sites for the proteins and RNAs. Two pairs of genes are coded by the same region of DNA using different reading frames.

The genome of bacteriophage Φ X174 is a single-stranded, circular DNA of approximately 5,400 nucleotides coding for nine known proteins. The order of these genes, as determined by genetic techniques¹⁻⁴, is A-B-C-D-E-J-F-G-H. Genes F, G and H code for structural proteins of the virus capsid, and gene J (as defined by sequence work) codes for a small basic protein.

1977
1st Complete Organism
Bacteriophage ϕ X174
5375 bp



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

Nucleotide sequence of bacteriophage φ X174 DNA

Sanger, F. et al. (1977) Nature. 265: 687 – 695

Nobel Prize in Chemistry in 1980

Milestones in DNA Sequencing



(TIGR/Celera, 1995-2001)

The most wondrous map...



“Without a doubt, this is the most important, most wondrous map ever produced by humankind.”

*Bill Clinton
June 26, 2000*

Cost per Genome

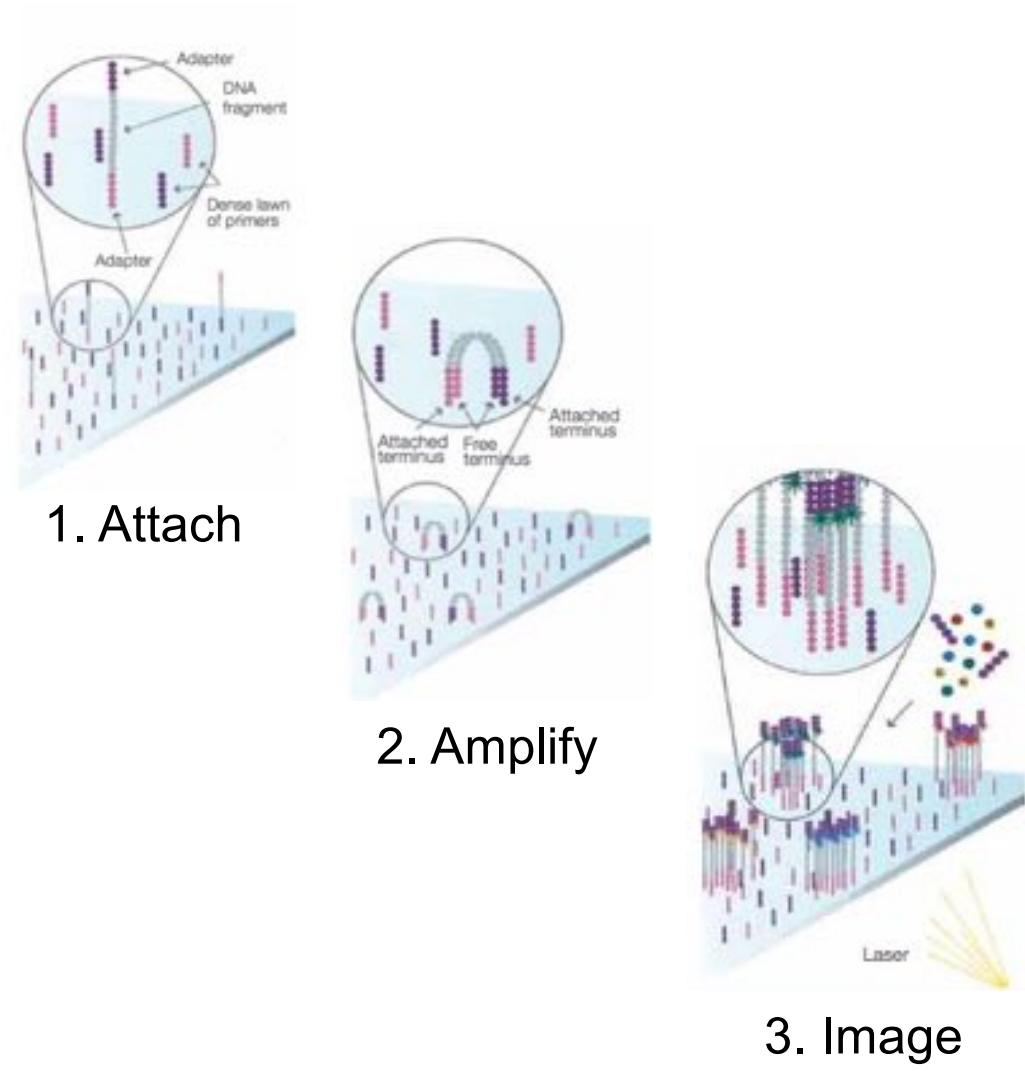


Second Generation Sequencing



Illumina NovaSeq 6000
Sequencing by Synthesis

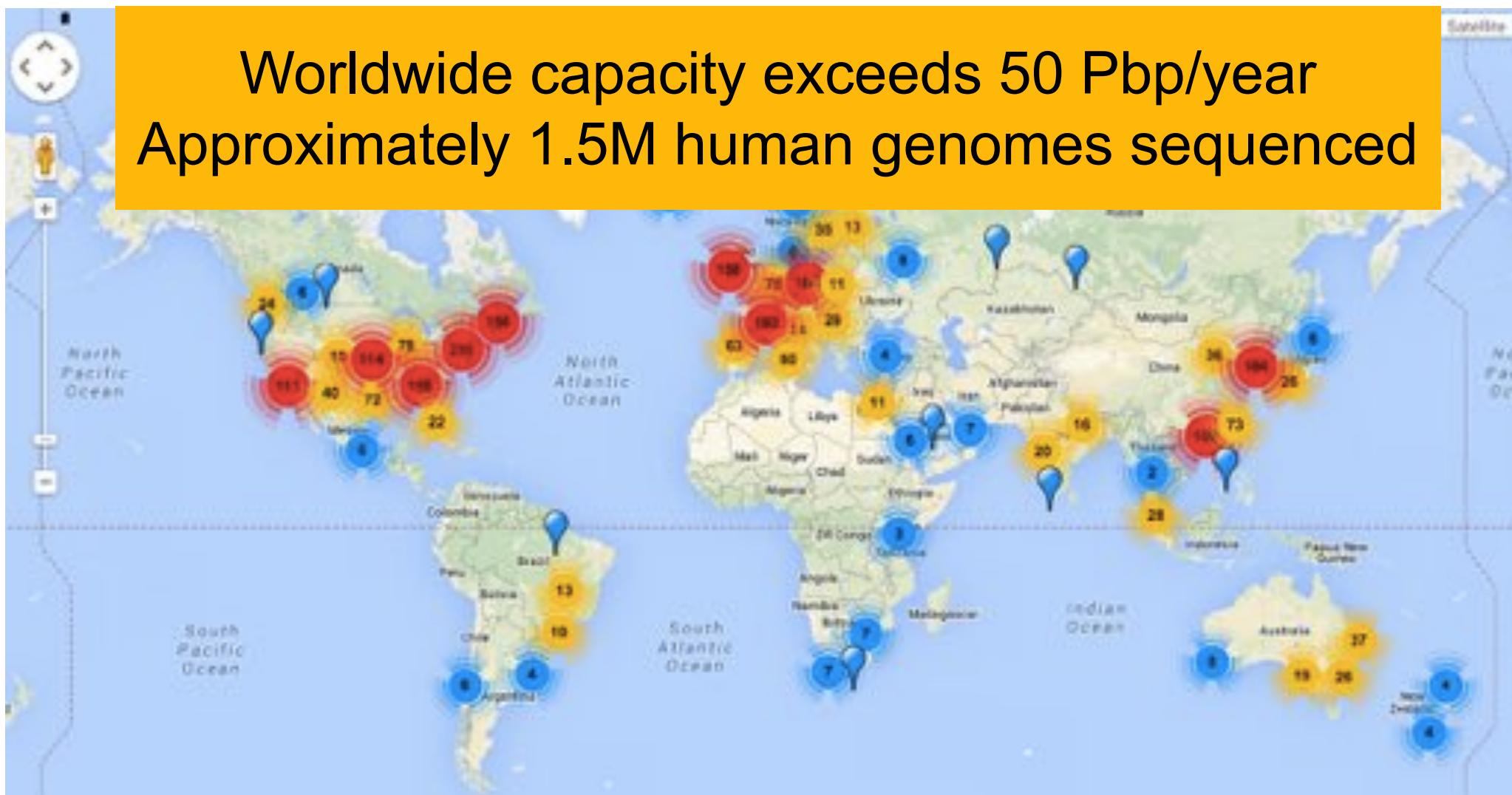
>3Tbp / day



Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Sequencing Centers

Worldwide capacity exceeds 50 Pbp/year
Approximately 1.5M 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

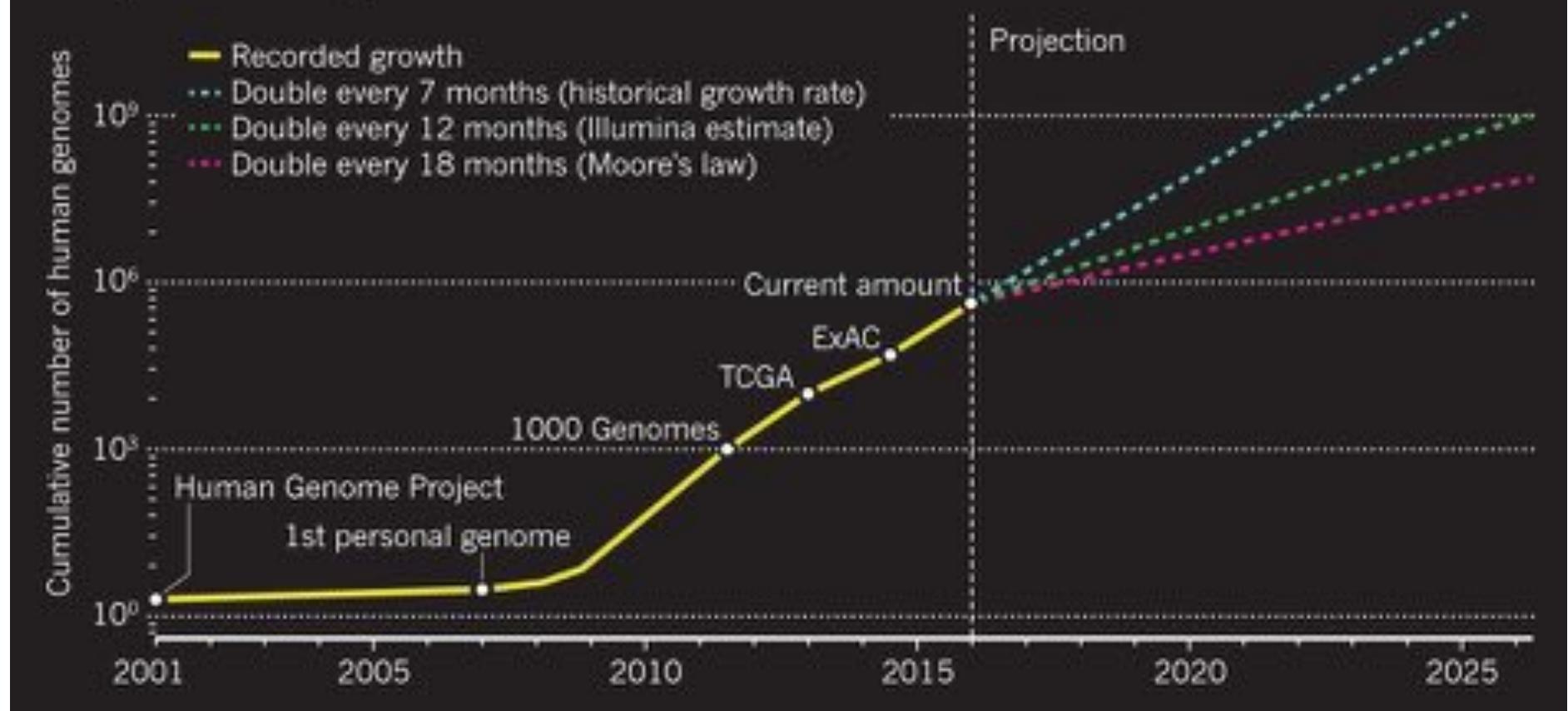


500 2 TB drives
\$50k

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)

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

Unsolved Questions in Biology

- What is your genome sequence?
 -
 -
 - The instruments provide the data, but none of the answers to any of these questions.
 -

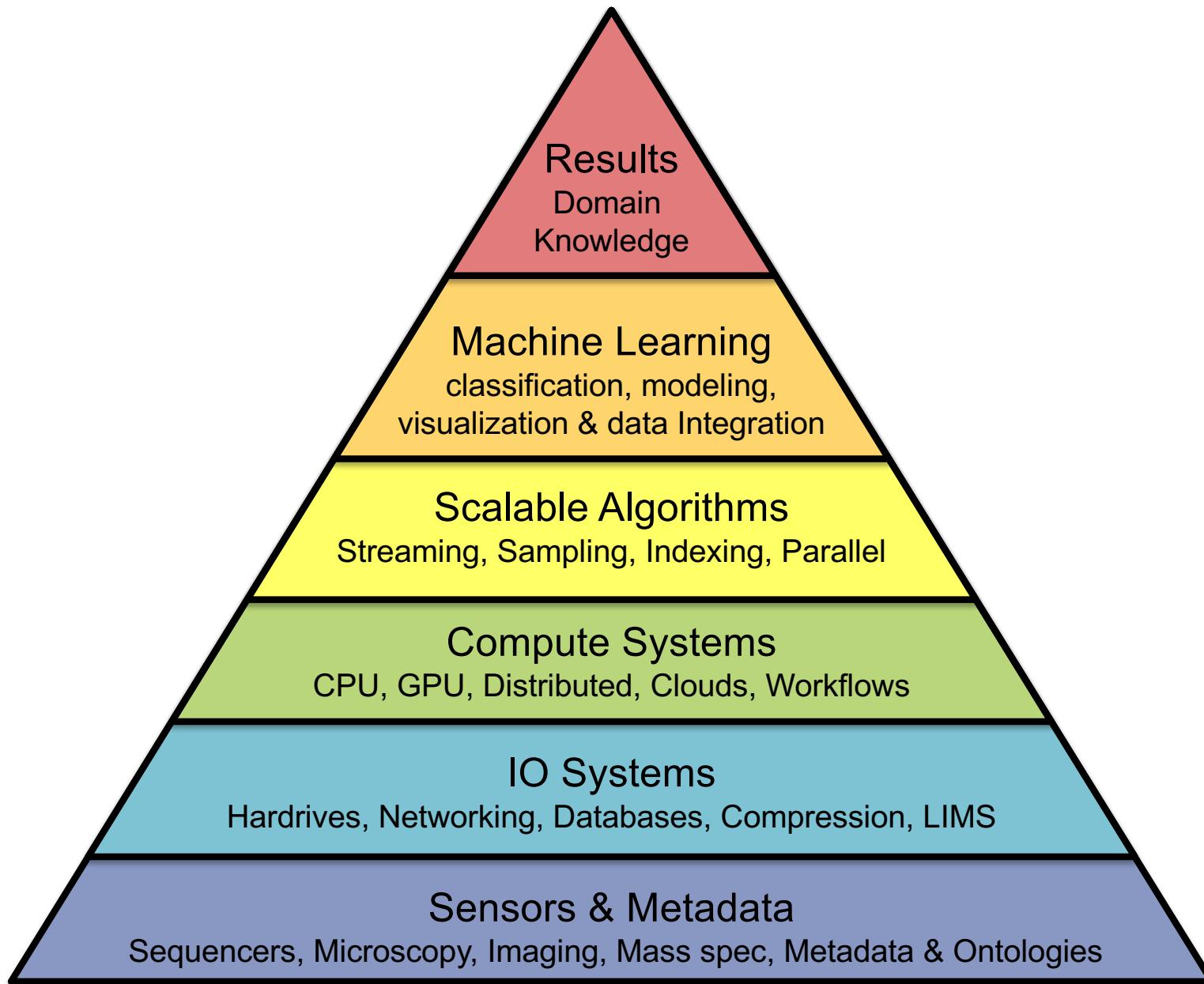
What software and systems will?

And who will create them?

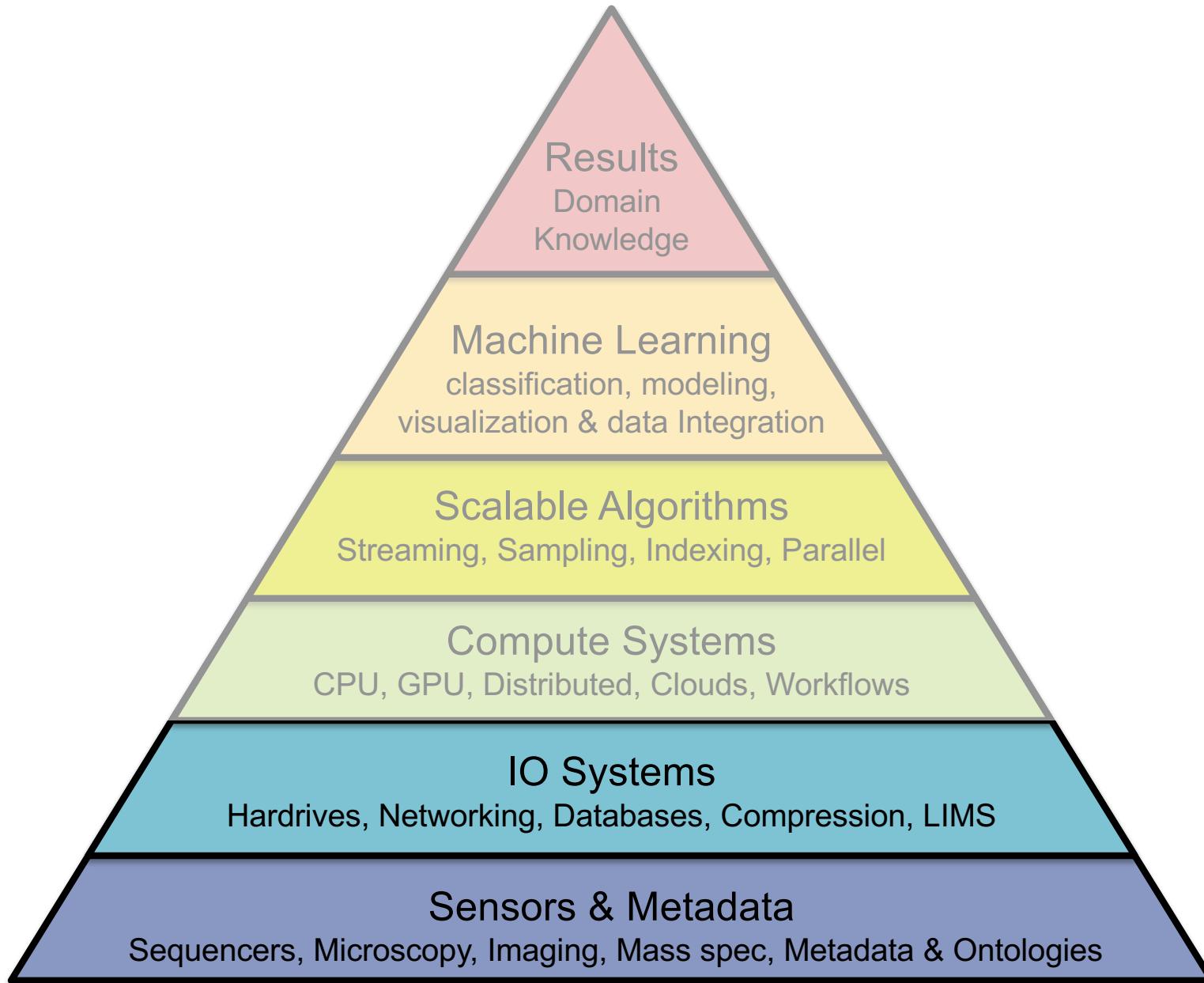
- **Plus thousands and thousands more**



Comparative Genomics Technologies



Comparative Genomics Technologies



Genomics Arsenal in the year 2021

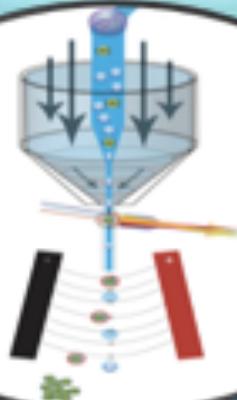
Sample Preparation

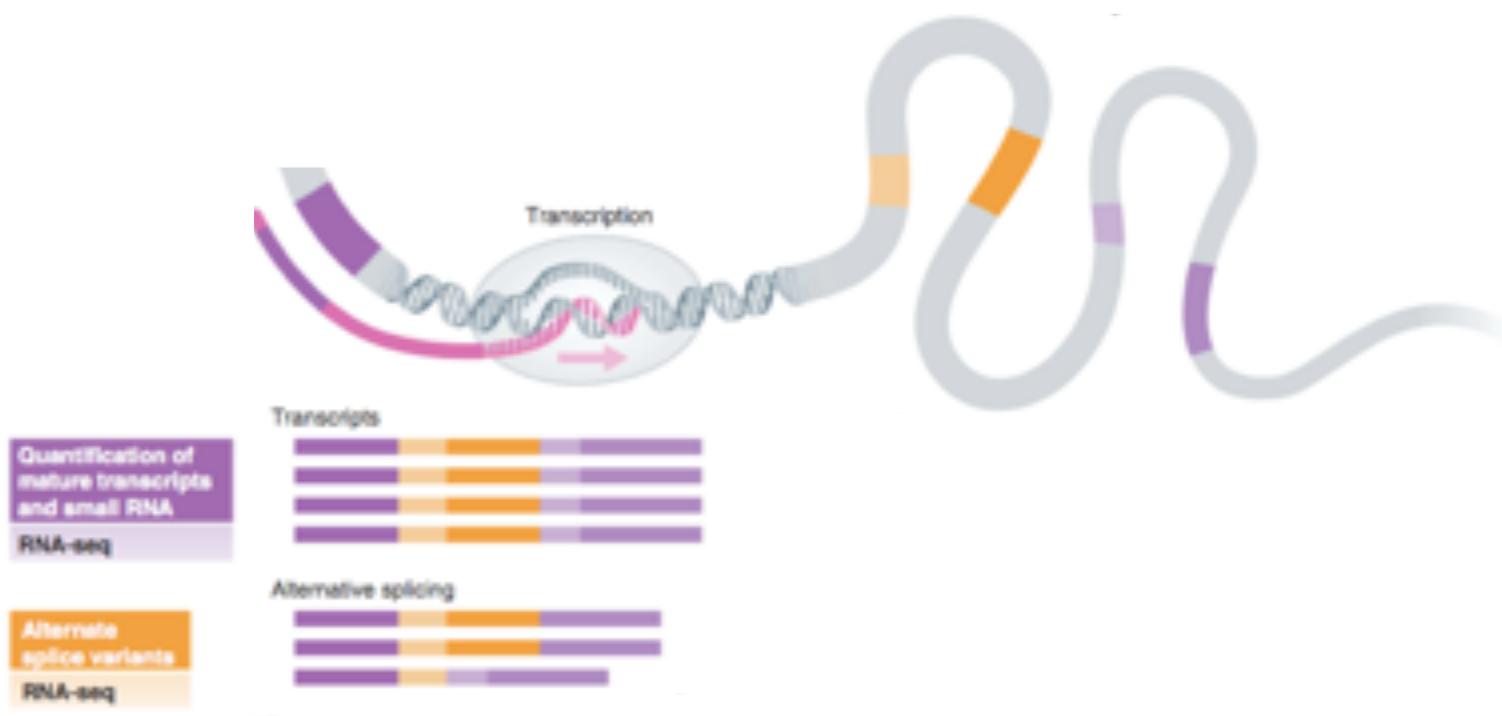


Sequencing

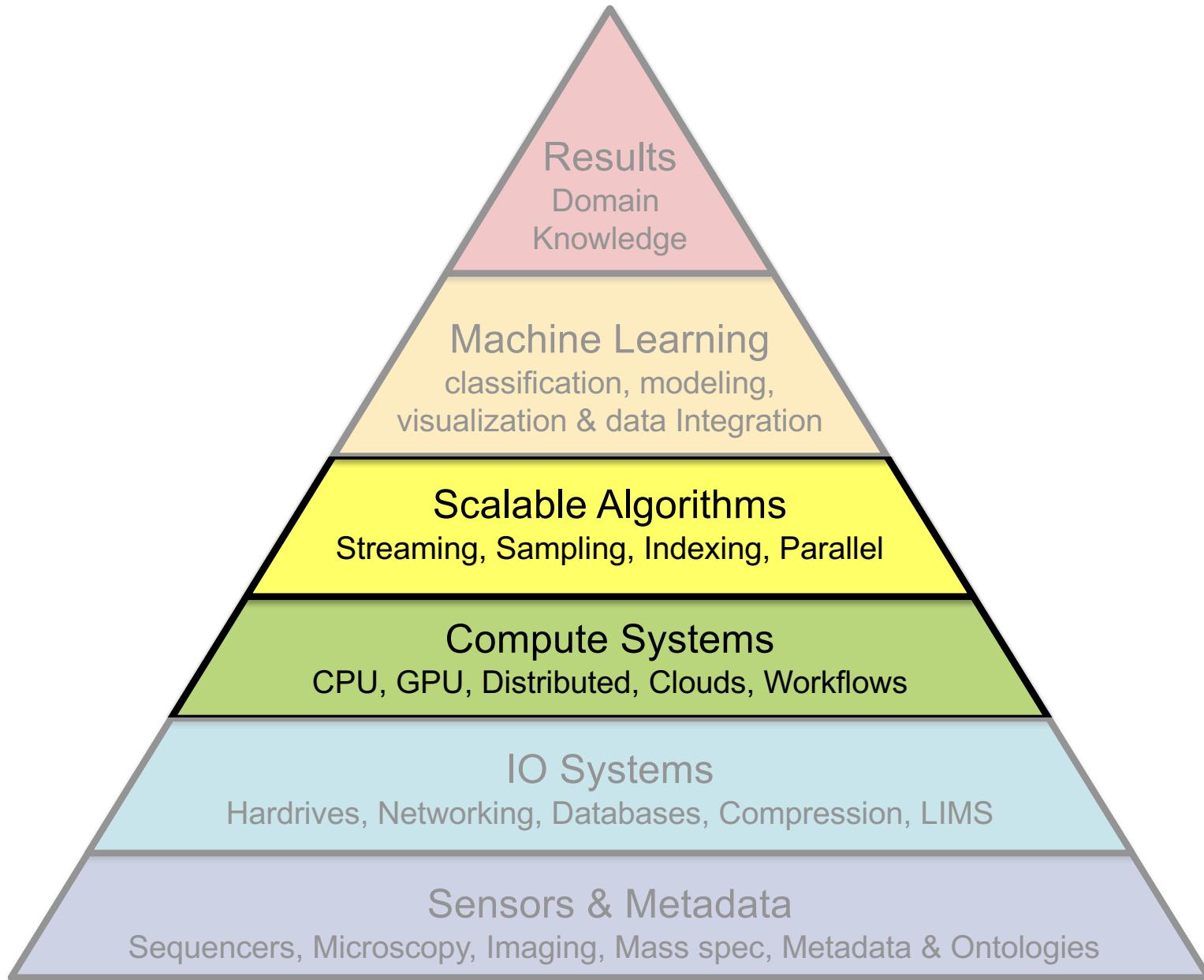


Chromosome Mapping





Comparative Genomics Technologies

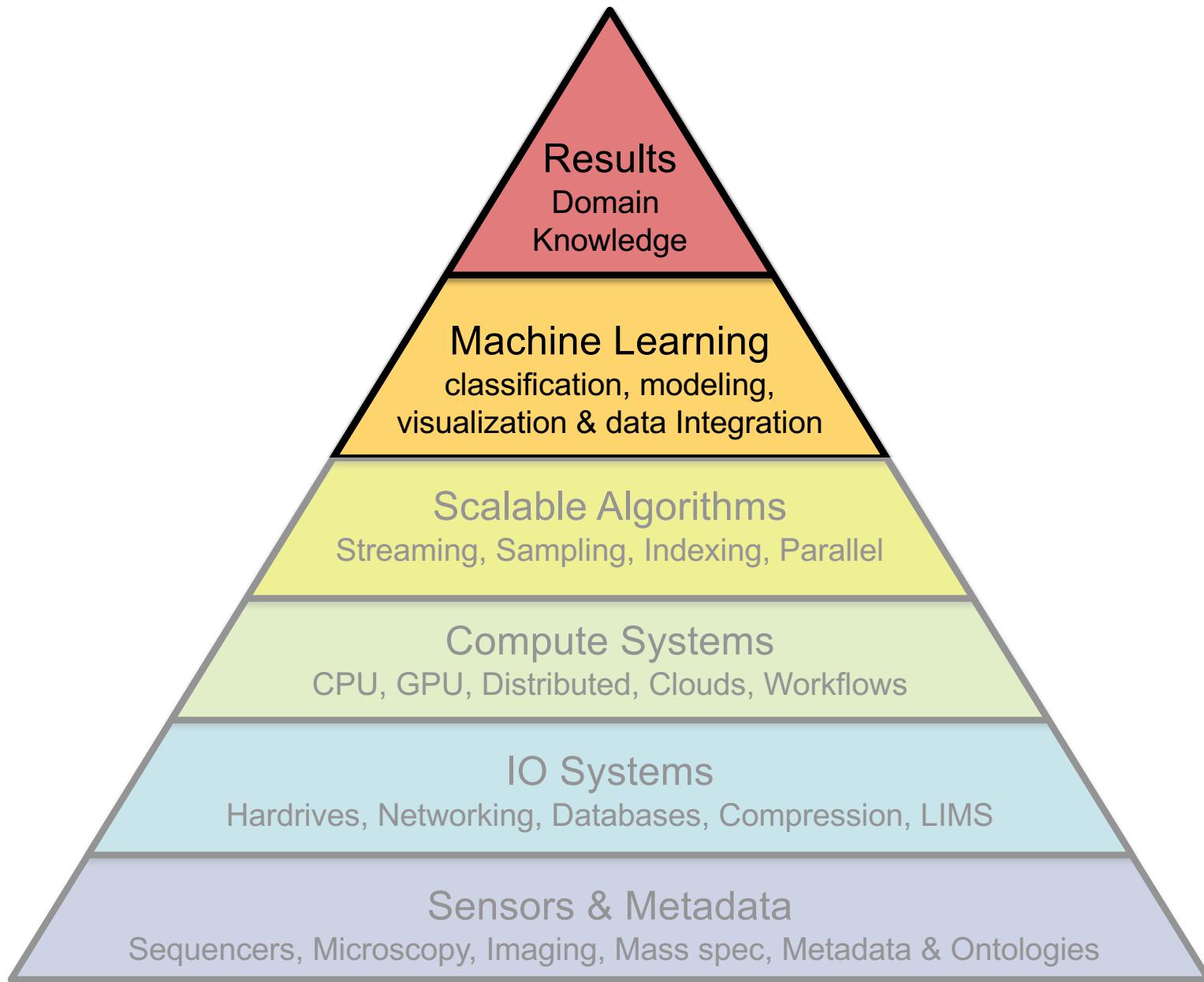


Potential Topics

- Genome assembly, whole genome alignment
- Full text indexing: Suffix Trees, Suffix Arrays, FM-index
- Dynamic Programming: Edit Distance, sequence similarity
- Read mapping & Variant identification
- Gene Finding: HMMs, Plane-sweep algorithms
- RNA-seq: mapping, assembly, quantification
- ChIP-seq: Peak finding, motif finding
- Methylation-seq: Mapping, CpG island detection
- HiC: Domain identification, scaffolding
- Chromatin state analysis: ChromHMM
- Scalable genomics: Cloud computing, scalable data structures
- Population & single cell analysis: clustering, pseudotime
- Disease analysis, cancer genomics, Metagenomics
- Deep learning in genomics



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?

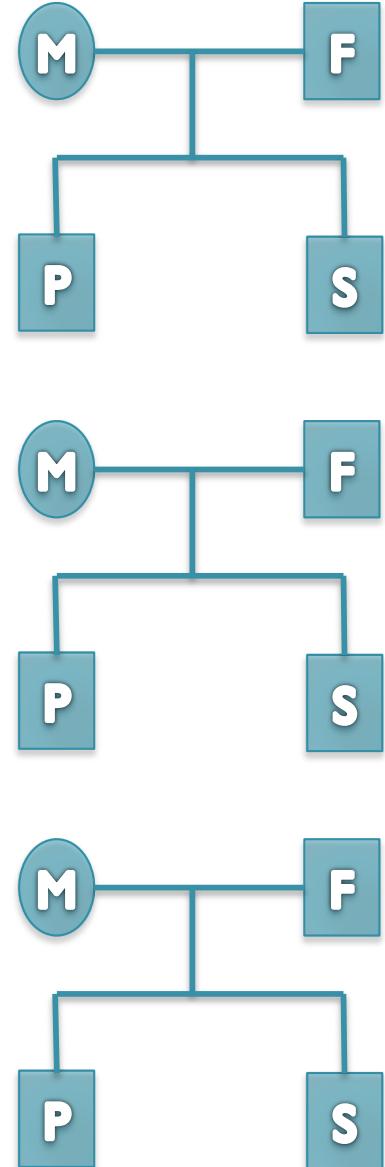
<https://autisticadvocacy.org/>

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

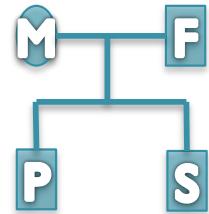
Are there any genetic variants present in affected children, that are not in their parents or unaffected siblings?



De novo mutation discovery and validation

De novo mutations:

Sequences not inherited from your parents.



Reference: . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Father(1): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Father(2): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Mother(1): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Mother(2): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Sibling(1): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Sibling(2): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

Proband(1): . . . TCAAATCCTTTAATAAAGAAGAGCTGACA . . .

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

4bp heterozygous deletion at chr15:93524061 CHD2

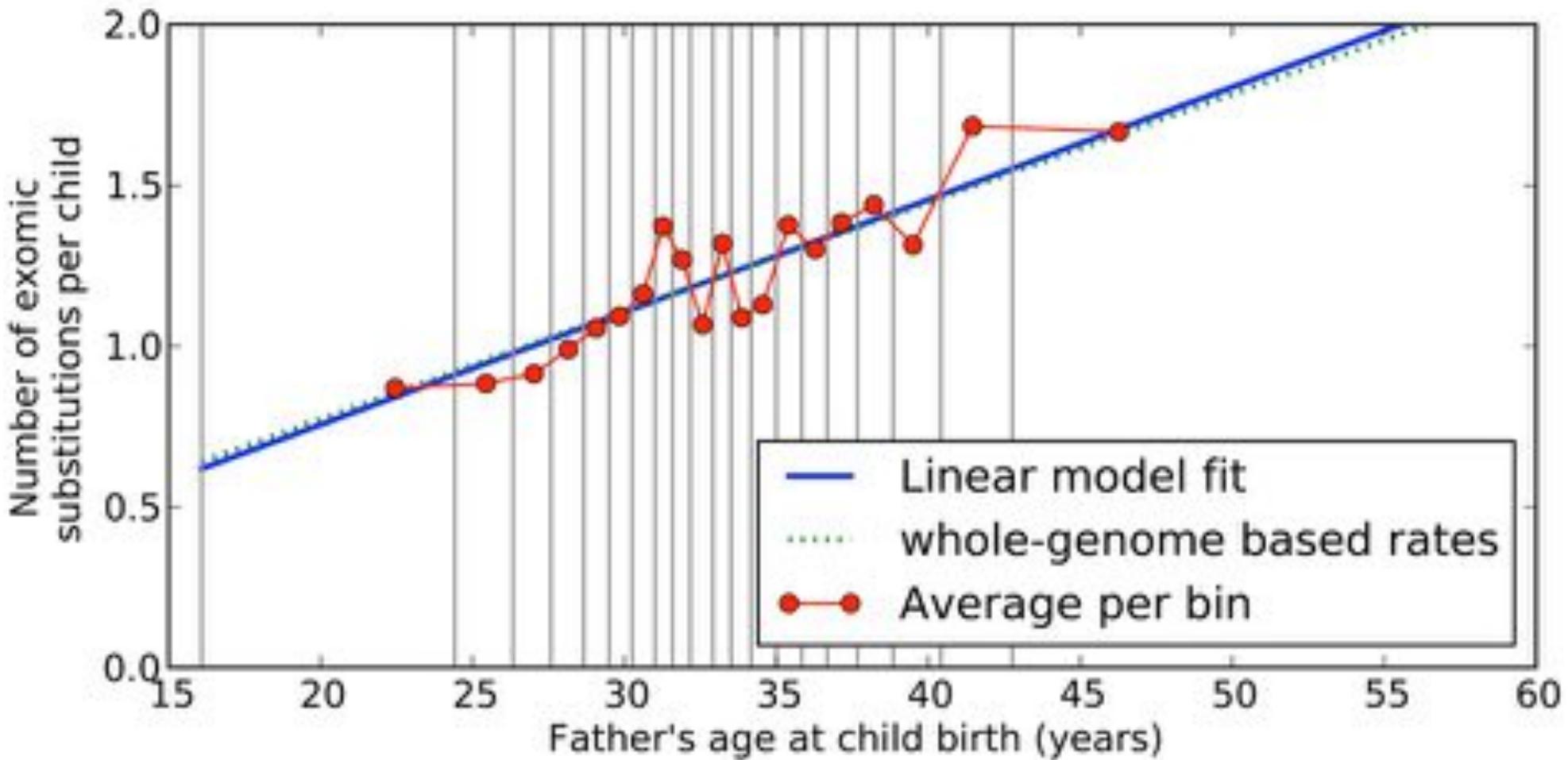
De novo Genetics of Autism

- 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.

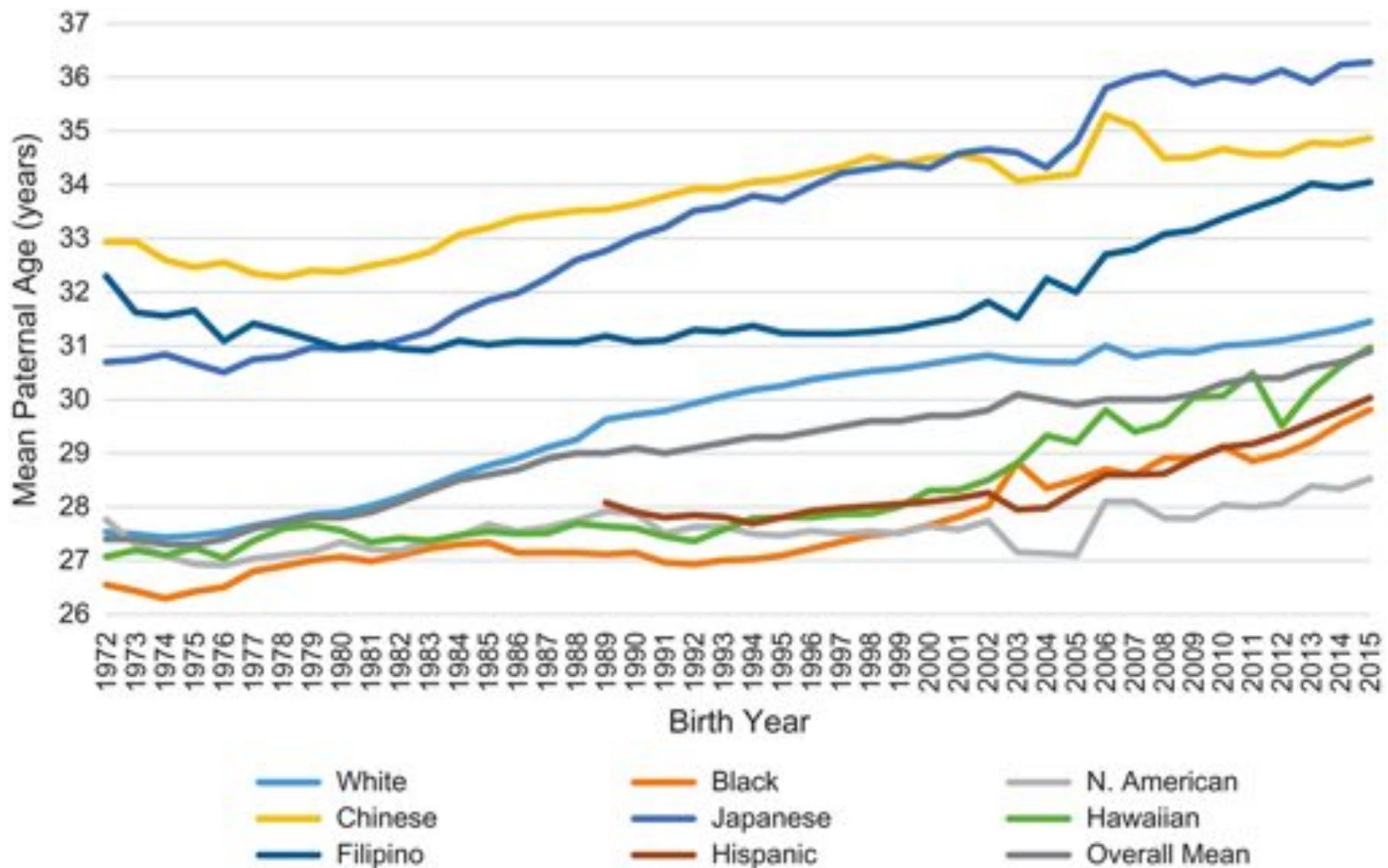
Narzisi et al (2014) Nature Methods doi:10.1038/nmeth.3069

De novo Mutations in Men



The contribution of de novo coding mutations to autism spectrum disorder
Iossifov et al (2014) *Nature*. doi:10.1038/nature13908

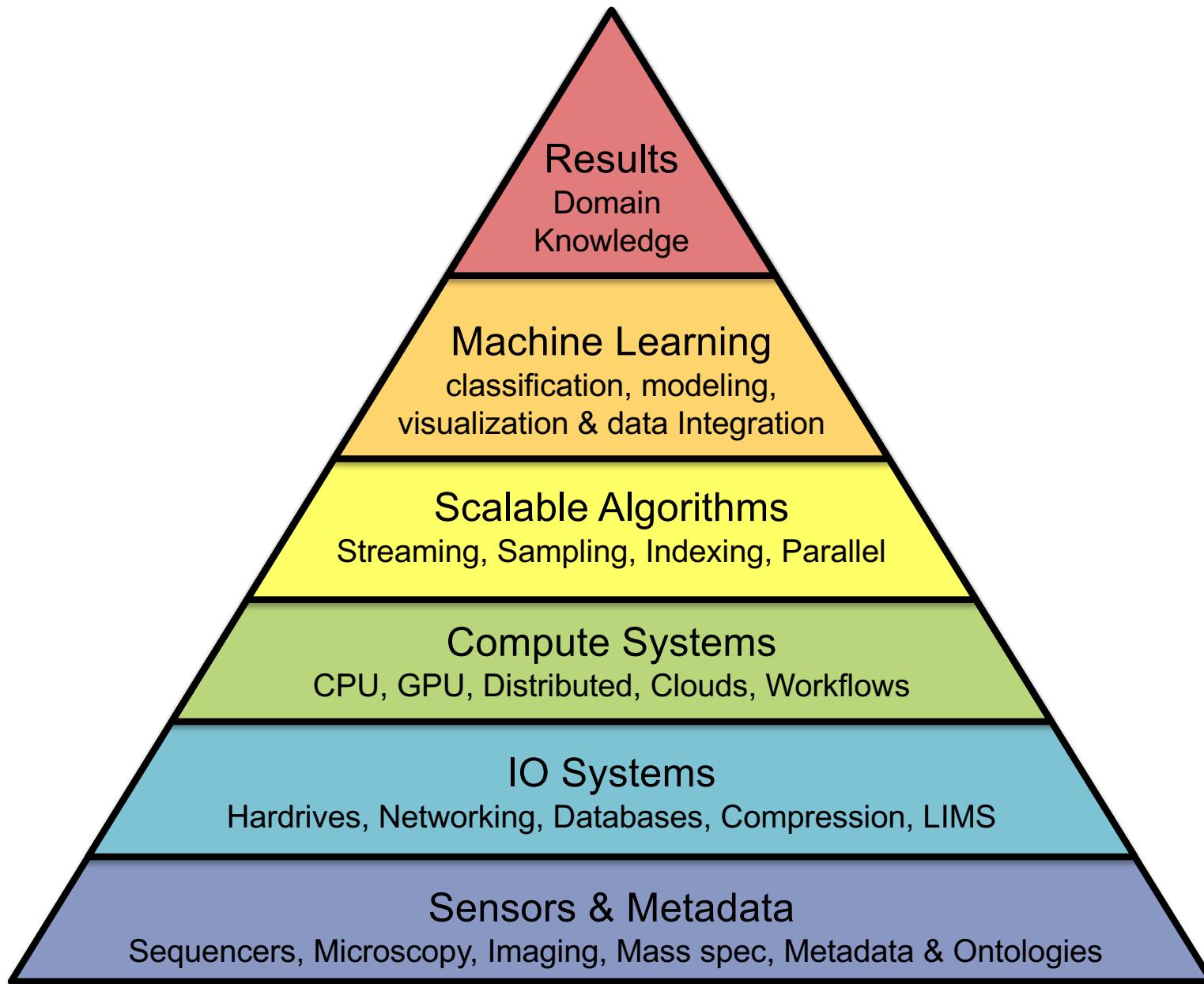
Age of Fatherhood

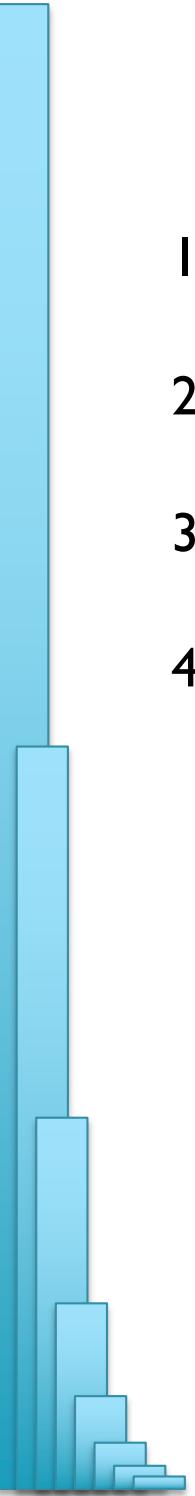


The age of fathers in the USA is rising: an analysis of 168 867 480 births from 1972 to 2015

Khandwala et al (2017) Human Reproduction. <https://doi.org/10.1093/humrep/dex267>

Comparative Genomics Technologies





Next Steps

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