

Applied Comparative Genomics

Michael Schatz

January 25, 2020

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/appliedgenomics2021>

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

Class Hours: Mon + Wed @ 1:30p – 2:45p, Zoom

Schatz Office Hours: TBD and by appointment

Das Office Hours: TBD 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 2018/2020
- <https://github.com/quinlan-lab/applied-computational-genomics>
- Ben Langmead's teaching materials:
- <http://www.langmead-lab.org/teaching-materials/>

Grading Policies

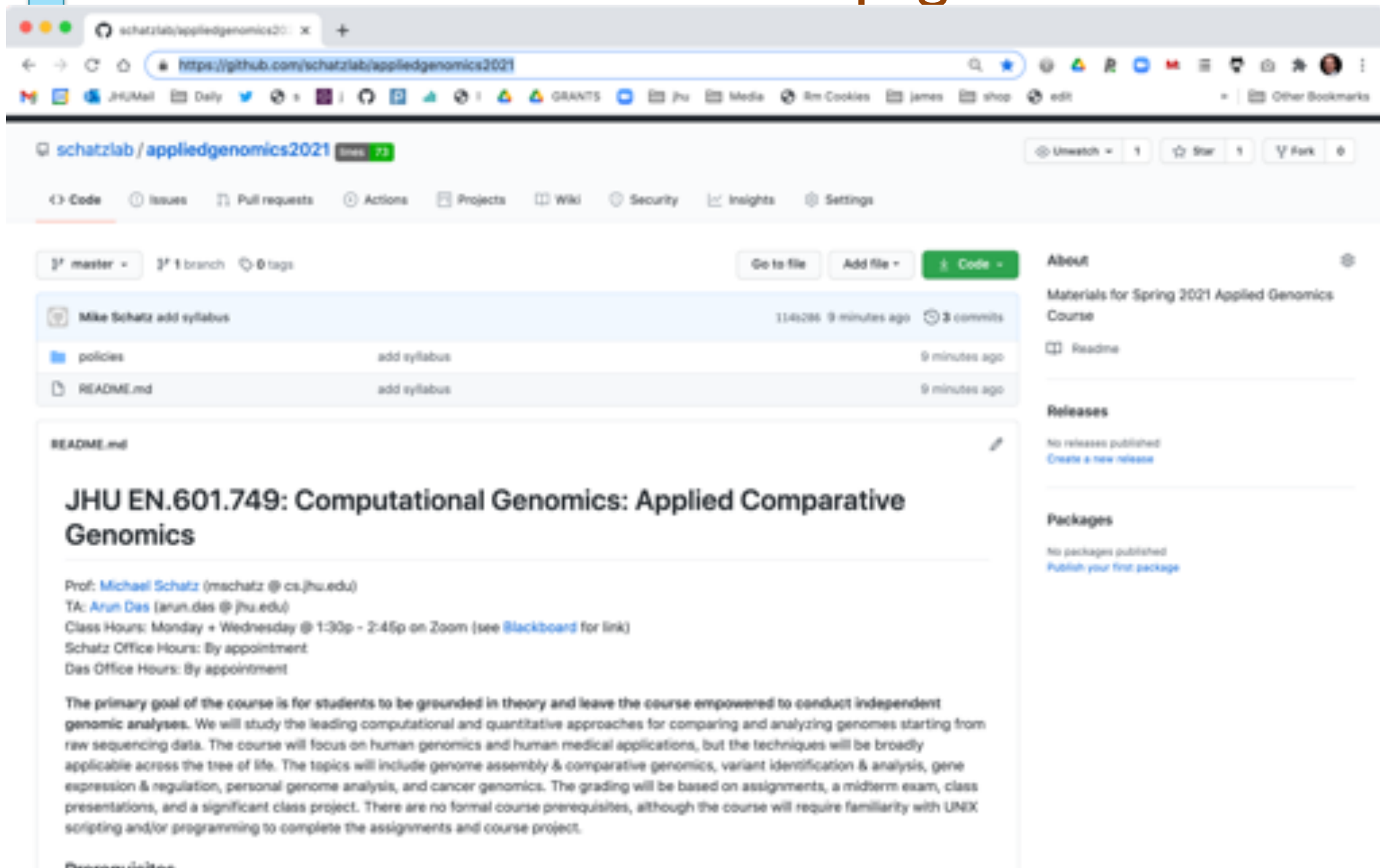
Assessments:

- 5 Assignments: 30% Due at 11:59pm a week later
Practice using the tools we are discussing
- 1 Exam: 30% In class (Tentatively 3/29)
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:***
 - A total of 96 hours (24 x 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



schatzlab / appliedgenomics2021

Unwatch 1 Star 1 Fork 0

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

master 1 branch 0 tags

Go to file Add file + Code +

Commit	Files	Time
Mike Schatz add syllabus		114s286 9 minutes ago 3 commits
	policies	add syllabus 9 minutes ago
	README.md	add syllabus 9 minutes ago

README.md

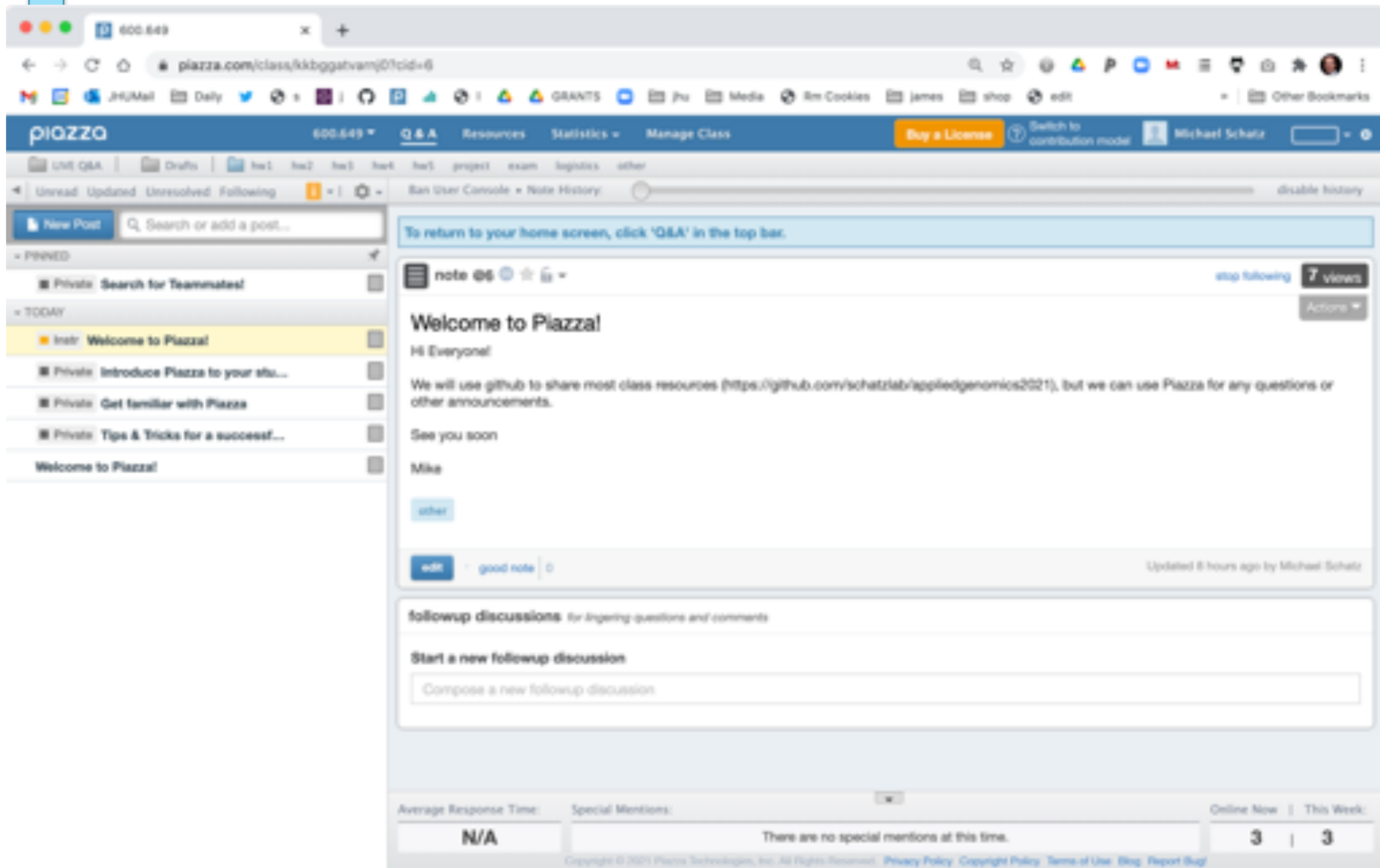
JHU EN.601.749: Computational Genomics: Applied Comparative Genomics

Prof: [Michael Schatz](#) (mschatz @ cs.jhu.edu)
TA: [Arun Das](#) (arun.das @ jhu.edu)
Class Hours: Monday + Wednesday @ 1:30p - 2:45p on Zoom (see [Blackboard](#) for link)
Schatz Office Hours: By appointment
Das Office Hours: By appointment

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. The grading will be based on assignments, a midterm exam, class presentations, and a significant class project. There are no formal course prerequisites, although the course will require familiarity with UNIX scripting and/or programming to complete the assignments and course project.

<https://github.com/schatzlab/appliedgenomics2021>

Piazza



The screenshot displays the Piazza website interface for a class page. The browser address bar shows the URL piazza.com/class/kkbggahvamj07cid+6. The page header includes the Piazza logo, the class ID 600.649, and navigation links for Q&A, Resources, Statistics, and Manage Class. A 'Buy a License' button and a 'Switch to contribution model' link are also present. The user profile 'Michael Schatz' is visible in the top right corner.

The main content area features a 'Welcome to Piazza!' message from Michael Schatz, dated 8 hours ago. The message includes a greeting, a link to the class GitHub repository, and a note about using Piazza for questions and announcements. The message is marked as 'good note' and has 0 replies.

On the left sidebar, there is a 'New Post' button and a search bar. Below the search bar, there are sections for 'Pinned' and 'Today' posts. The 'Today' section lists several posts, including 'Welcome to Piazza!' and 'Introduce Piazza to your stu...'. The 'Pinned' section lists 'Search for Teammate!'. The 'Welcome to Piazza!' post is highlighted in yellow.

At the bottom of the page, there is a 'followup discussions' section with a text input field for composing a new discussion. The footer displays the average response time as 'N/A', special mentions as 'There are no special mentions at this time.', and online status for 'Online Now' (3) and 'This Week' (3). Copyright information for Piazza Technologies, Inc. is also present.

<http://piazza.com/jhu/spring2021/600649>

GradeScope

The screenshot shows the GradeScope dashboard in a web browser. The browser's address bar displays 'gradescope.com'. The dashboard has a left sidebar with the 'gradescope' logo and a 'Your Courses' section containing a welcome message. The main content area is titled 'Your Courses' and lists two courses for 'Spring 2021' and 'Spring 2020'. Each course entry shows the course ID 'EN.601.749', the title 'Applied Comparative Genomics', and the number of assignments (0 for Spring 2021, 10 for Spring 2020). A dashed box highlights a '+ Create a new course' button. At the bottom, a teal footer bar contains an 'Account' menu, an 'Enroll in Course' button, and a 'Create Course +' button.

Dashboard | Gradescope

gradescope.com

gradescope

Your Courses

Welcome to Gradescope! Click on one of your courses to the right, or on the Account menu below.

Your Courses

Spring 2021

EN.601.749
Applied Comparative Genomics

0 assignments

+ Create a new course

Spring 2020

EN.601.749
Applied Comparative Genomics

10 assignments

See older courses ▾

Account

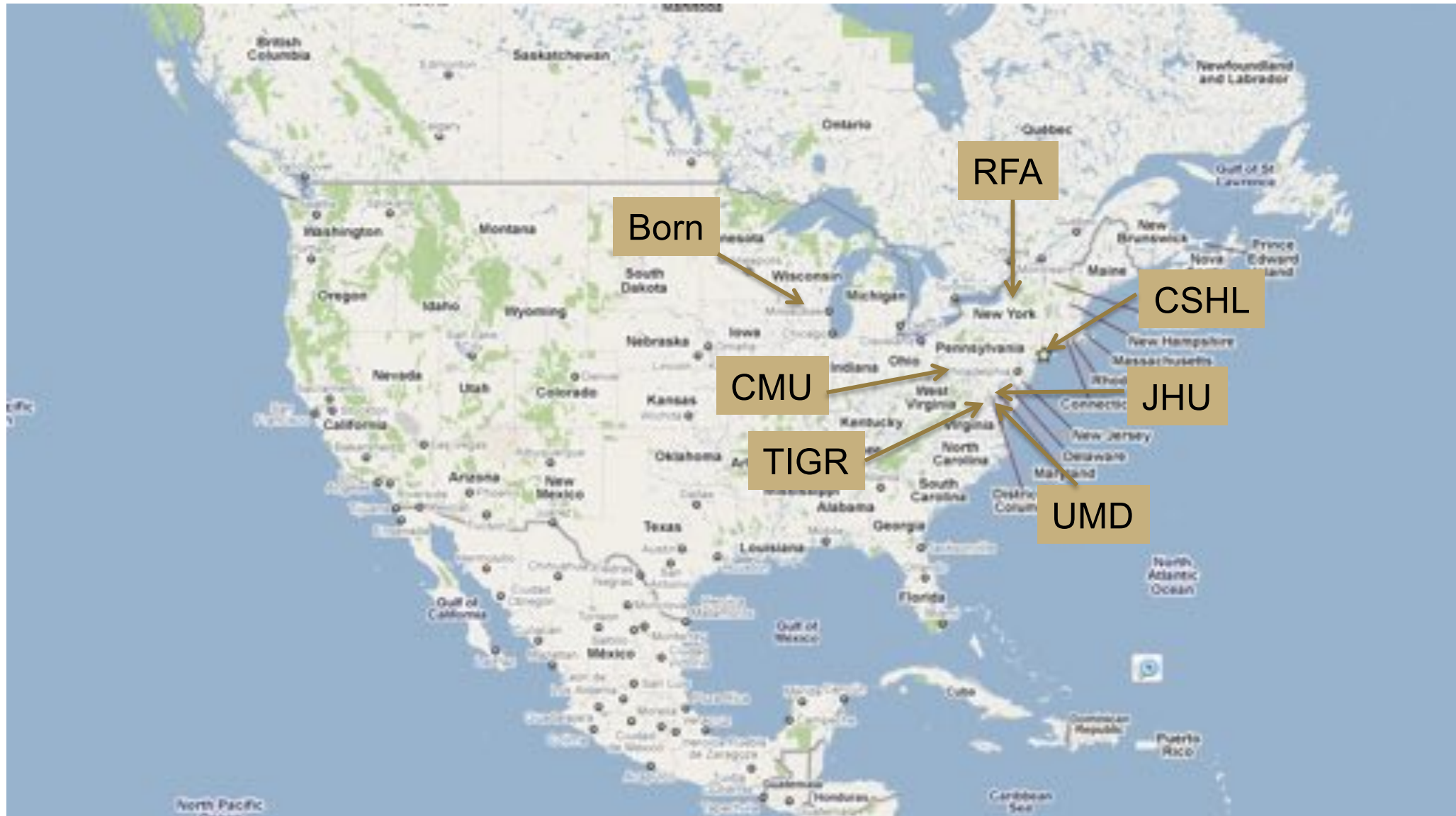
Enroll in Course

Create Course +

<https://www.gradescope.com/>

Entry Code: 86KZZP

A Little About Me



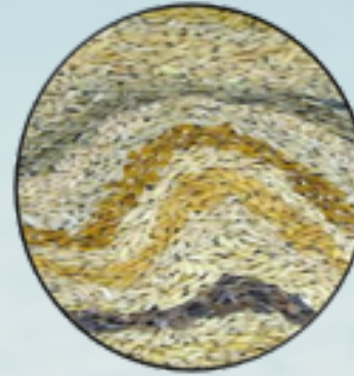
Schatzlab Overview



Human Genetics

Role of mutations
in disease

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



Agricultural Genomics

Genomes &
Transcriptomes

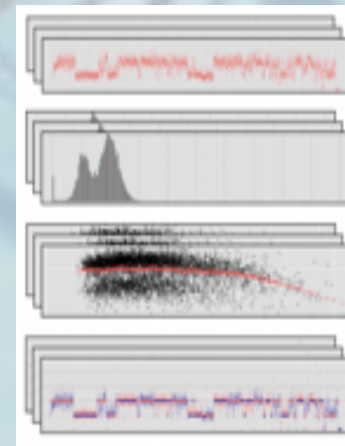
Alonge *et al.* (2020)
Soyk *et al.* (2019)



Algorithmics & Systems Research

Ultra-large scale
biocomputing

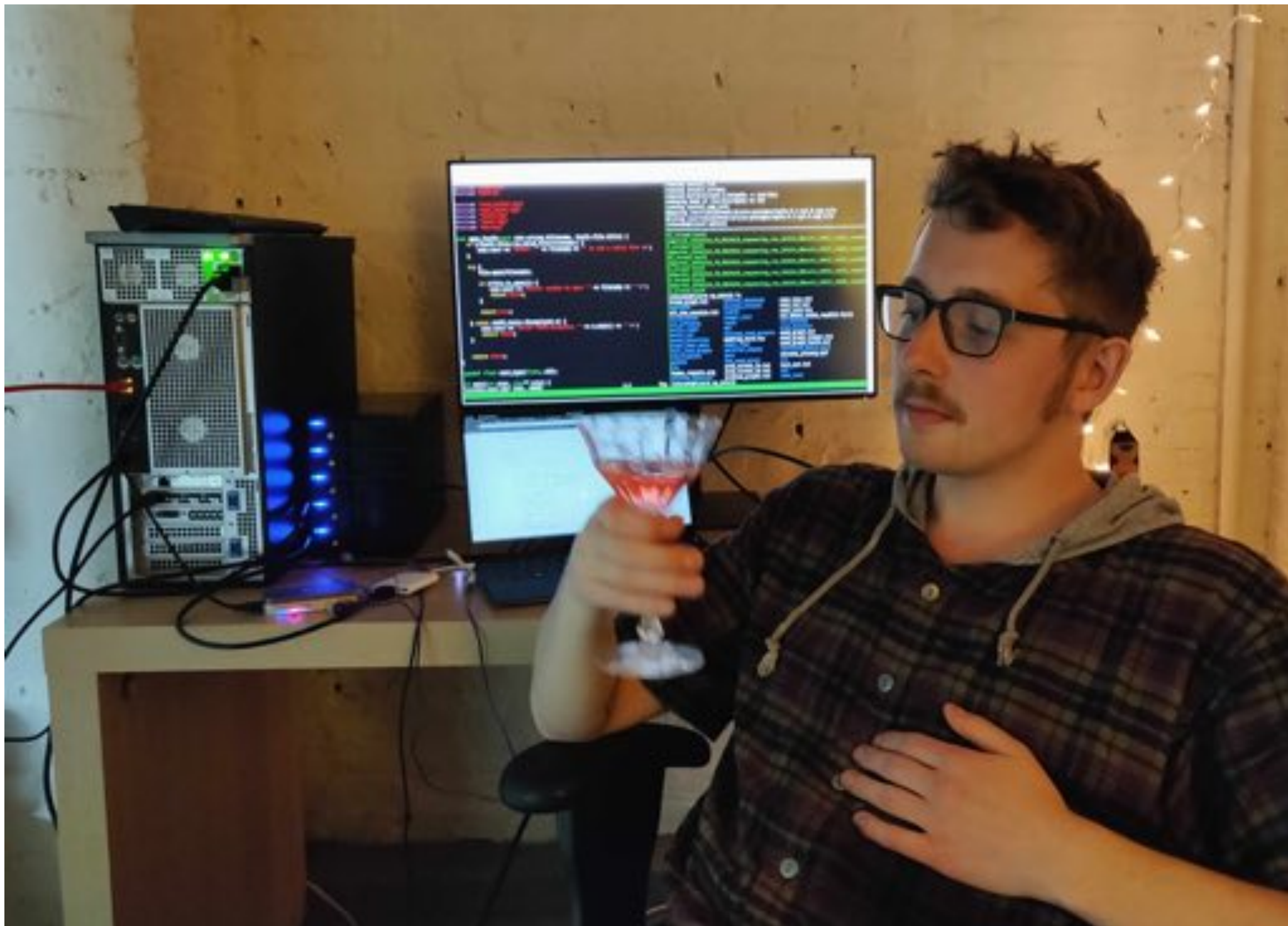
Kirsche *et al.* (2020)
Fang *et al.* (2018)



Biotechnology Development

Single Cell + Single
Molecule Sequencing

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



Targeted nanopore sequencing

nature.com/articles/s41587-020-0731-9

View all Nature Research journals

Search

My Account

Explore content

Journal information

Publish with us

Subscribe

Sign up for alerts

RSS feed

nature > nature biotechnology > articles > article

Article | Published: 30 November 2020

Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED

Sam Kovaka, Yunfan Fan, Bohan Ni, Winston Timp & Michael C. Schatz

Nature Biotechnology (2020) | Cite this article

5715 Accesses | 2 Citations | 261 Altmetric | Metrics

Abstract

Conventional targeted sequencing methods eliminate many of the benefits of nanopore sequencing, such as the ability to accurately detect structural variants or epigenetic modifications. The ReadUntil method allows nanopore devices to selectively eject reads from pores in real time, which could enable purely computational targeted sequencing. However, this requires rapid identification of on-target reads while most mapping methods require computationally intensive basecalling. We present UNCALLED (<https://github.com/skovaka/UNCALLED>), an open source mapper that rapidly matches streaming of nanopore current signals to a reference sequence. UNCALLED probabilistically

You have full access to this article via Johns Hopkins Libraries

Download PDF

Sections

Figures

References

Abstract

Main

Results

Discussion

Methods

Data availability

Code availability

References

Acknowledgements

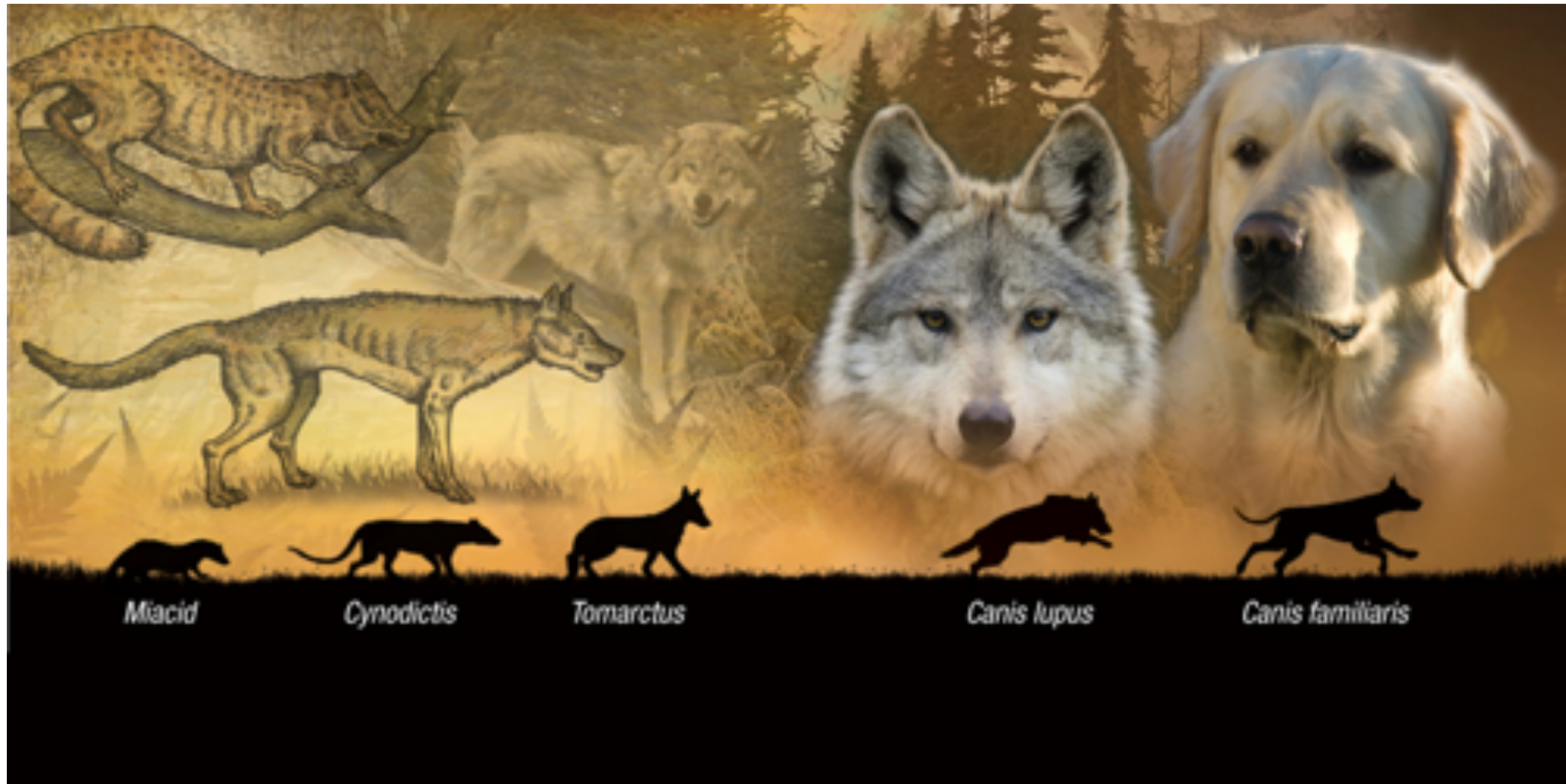
Author information

Ethics declarations

Earliest Genomics

Any Guesses?

Earliest Genomics



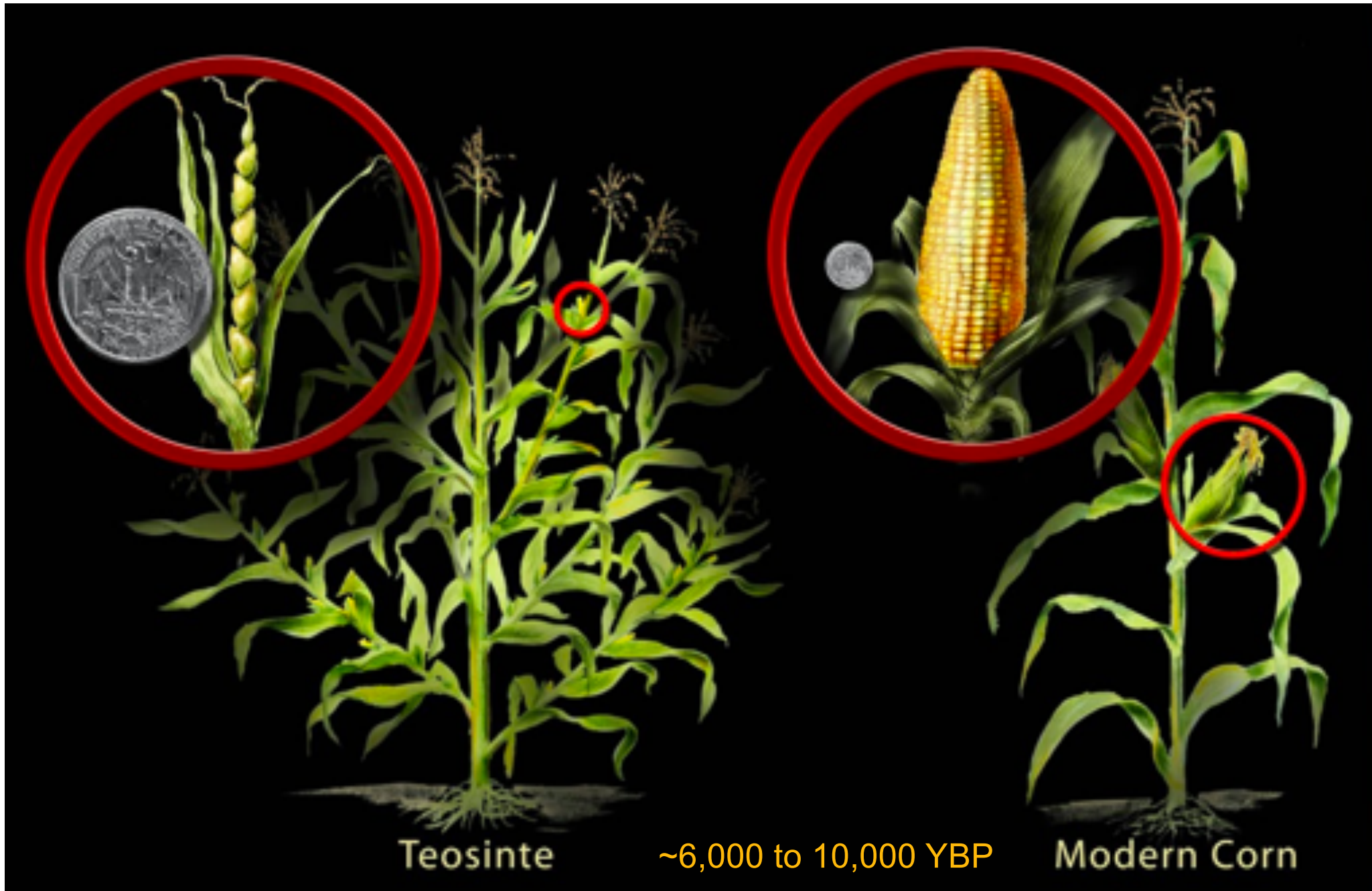
15,000 to 35,000 YBP

Earliest Genomics

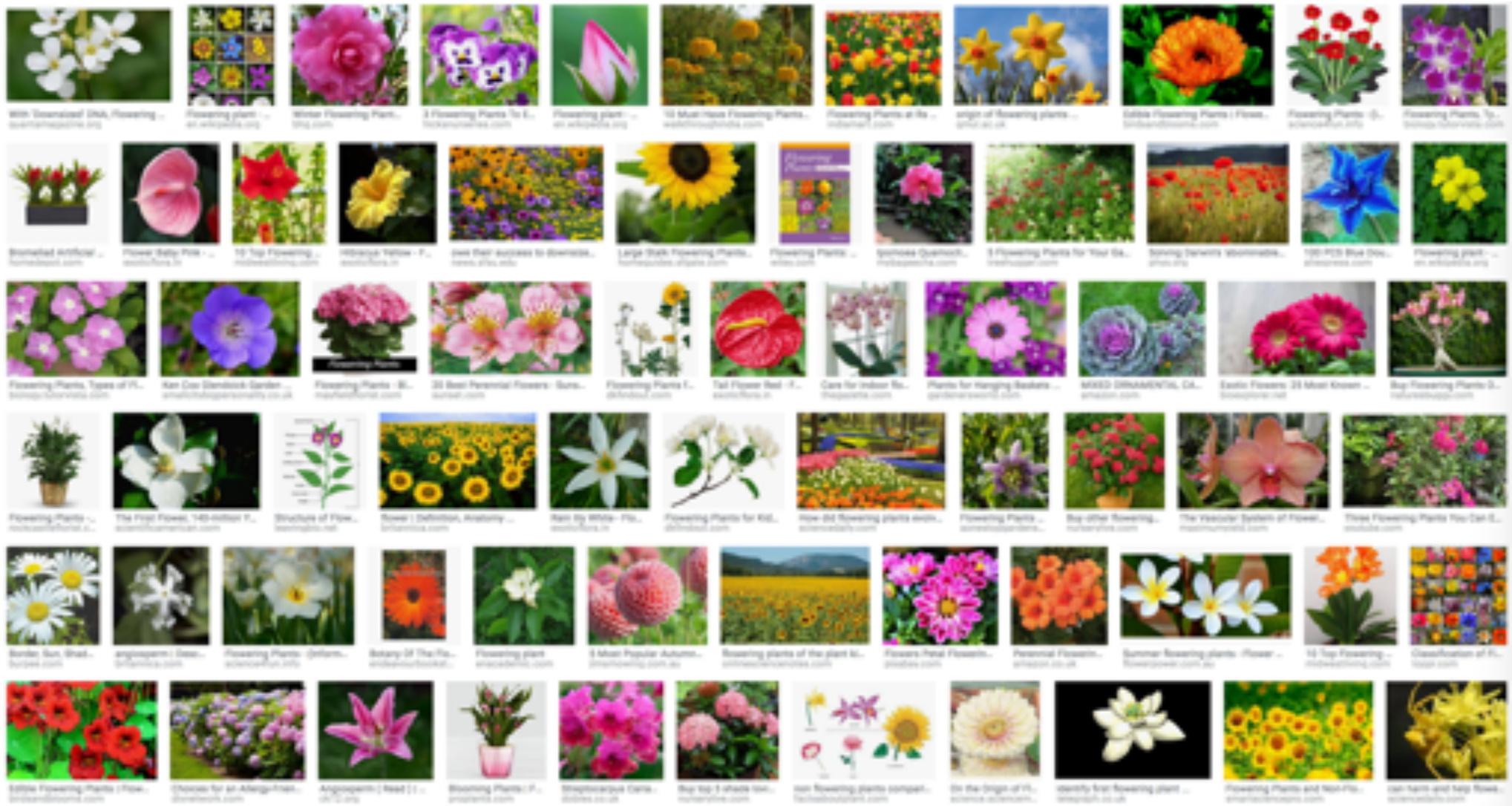


~1,000 to 10,000 YBP

Earliest Genomics



Angiosperms (Flowering Plants)



~130 Ma

Discovery of Chromosomes

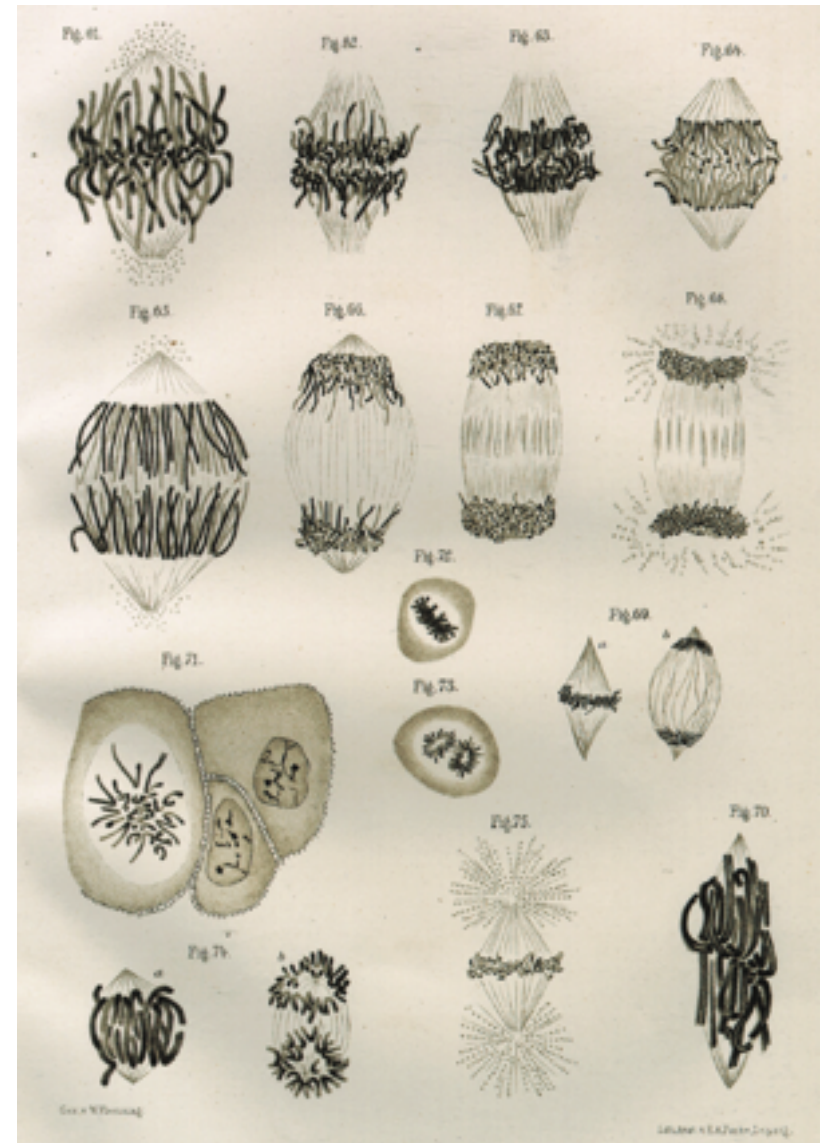
By the mid-1800s, microscopes were powerful enough to observe the presence of unusual structures called “chromosomes” that seemed to play an important role during cell division.

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 dyes 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



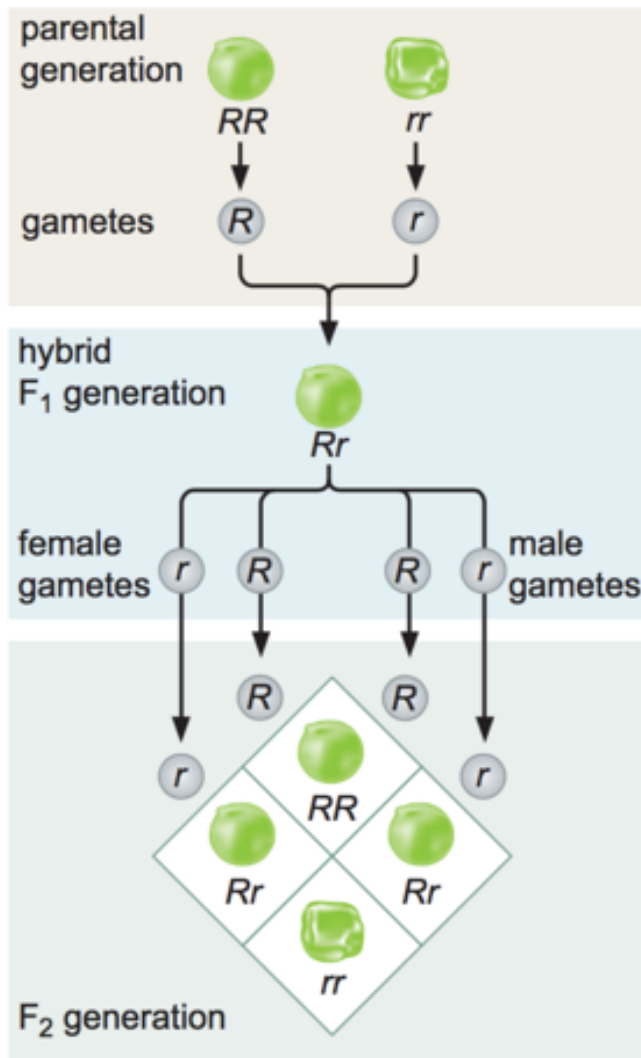
Drawing of mitosis by Walther Flemming.

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 (1-1ft)
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

				in Verhältniss gestellt:
Generation	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***



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

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

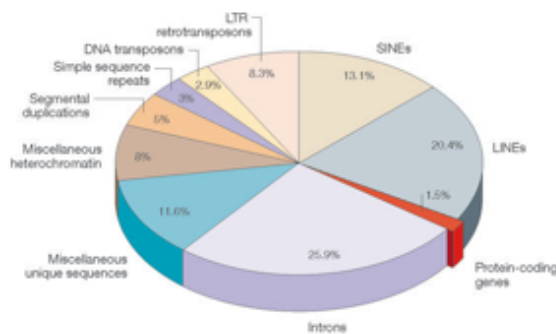
Sturtevant, A. H. (1913) *Journal of Experimental Zoology*, 14: 43-59

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. 4352 April 25, 1953 NATURE 737

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

¹Young, F. B., Gossard, E., and Jenson, W., *Phil. Mag.*, **46**, 149 (1928).

²Longuet-Higgins, M. S., *Nucl. Ac. Res. Adv. Ser., Copland, Japp., 4*, 261 (1949).

³See also, W. R., *Woods Hole Papers in Phys. Oceanogr. Meteor.*, **11** (1936).

⁴Ekman, T. W., *Adv. Nuc. Acids, Phys. (Oxford)*, **2**(1) (1953).

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.

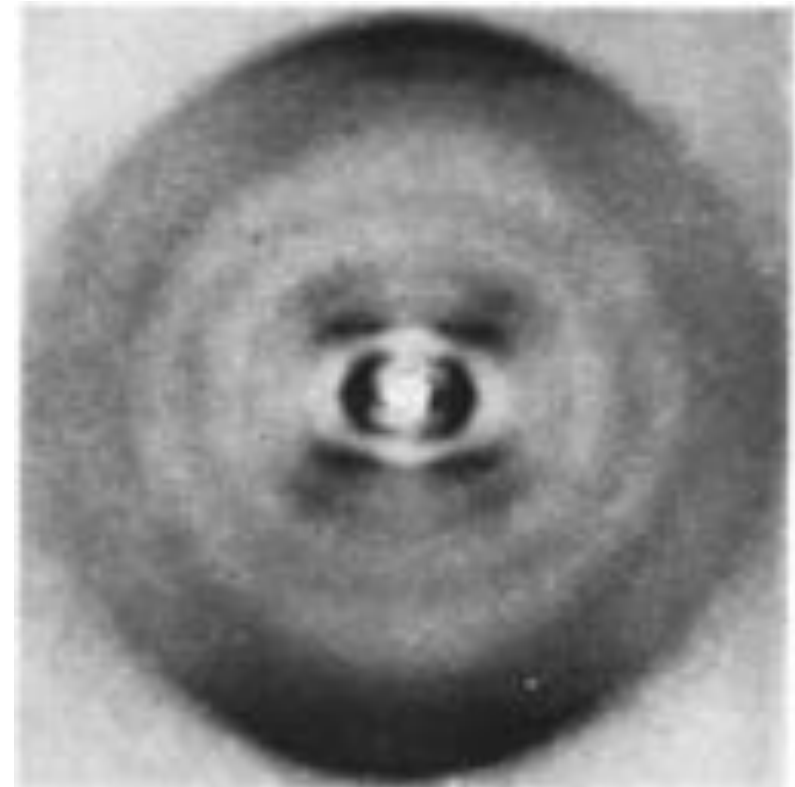
A structure for nucleic acid has already been proposed by Pauling and Corey¹. They kindly made their manuscript available to us in advance of publication. Their model consists of three intertwined chains, with the phosphates near the fibre axis, and the bases on the outside. In our opinion, this structure 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 what 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 Fawcett (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 as described is 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 coiled round the same axis (see diagram). We have made the usual chemical assumptions, namely, that each chain consists of phosphate diester groups joining 3'-phosphoryl residues with 5'/3' 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 Fawcett's model No. 1; that is, the bases are on the inside of the helix and the phosphates on the outside. The configuration of the sugar and the atoms near it is close to Fawcett's 'standard configuration', the sugar being roughly perpendicular to the attached base. There



This figure is purely diagrammatic. The two ribbons represent the two phosphate-sugar chains, and the horizontal rungs represent the pairs of bases holding the chains together. The vertical line marks the fibre axis.



CRITICAL ARGUMENTS.

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-

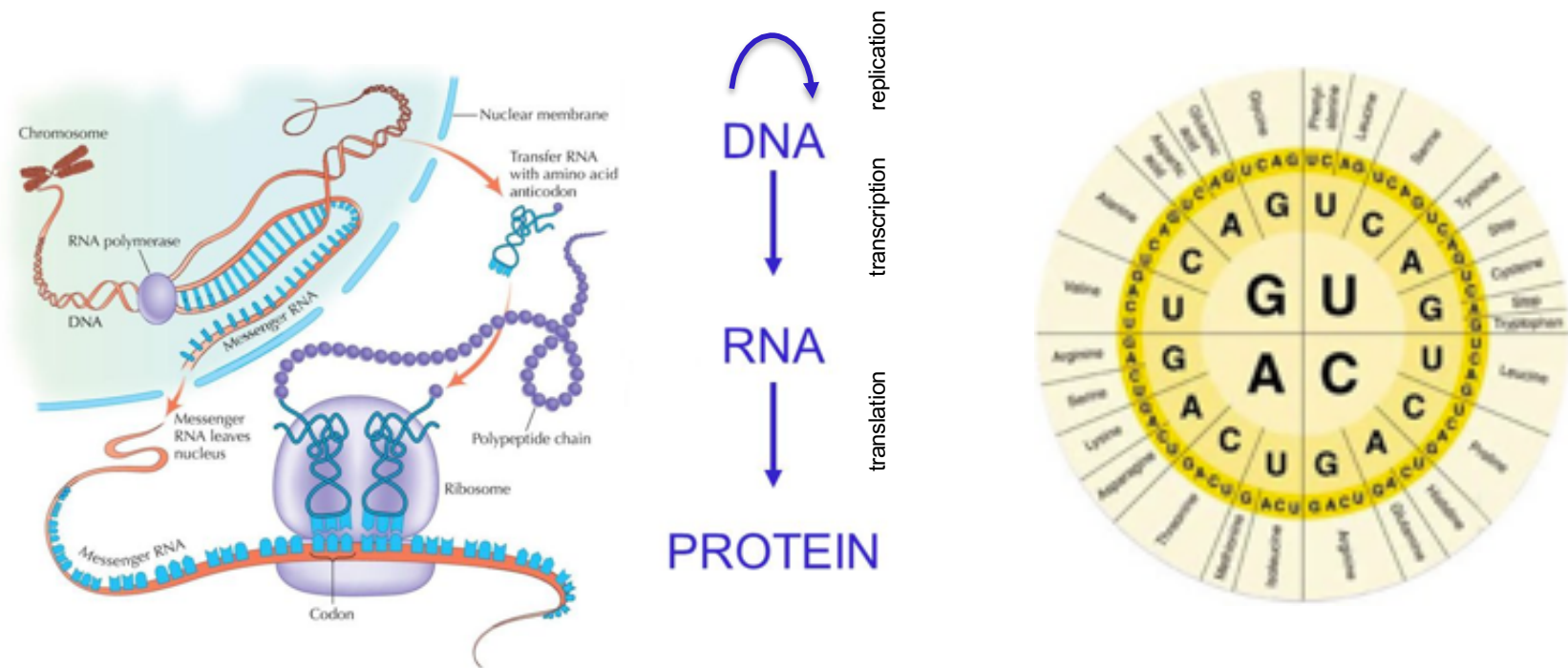
Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid

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”

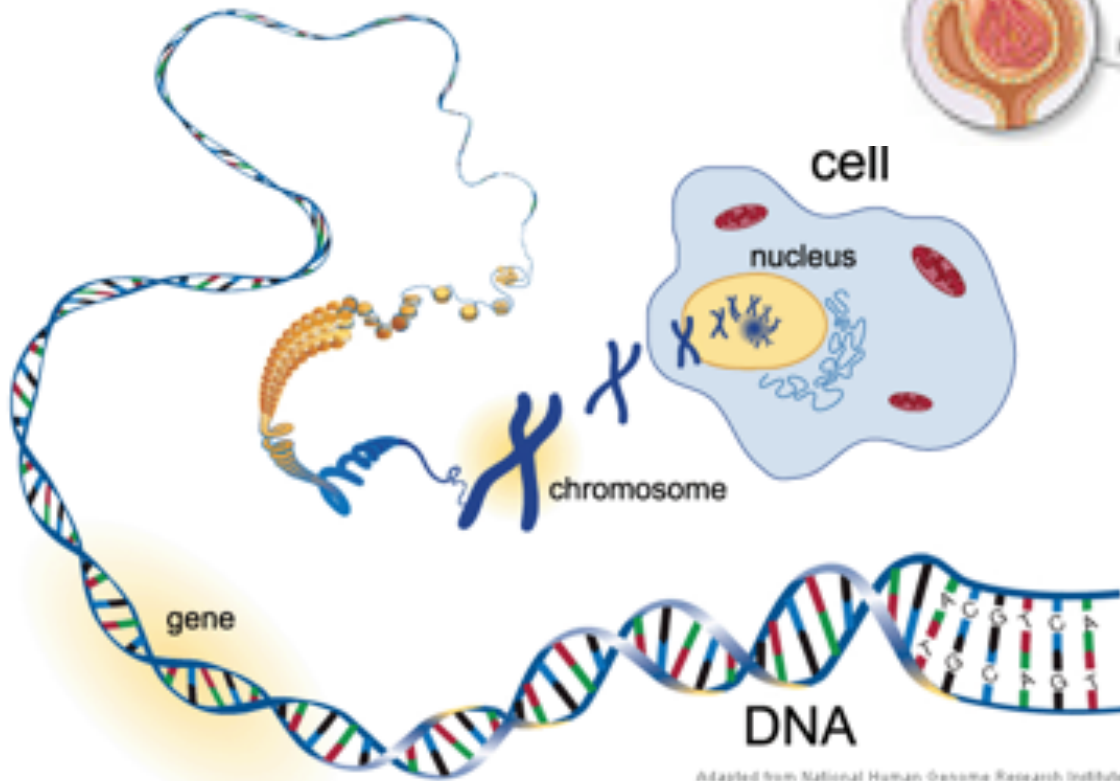


On Protein Synthesis

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.



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 2QH, 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

strand DNA of Φ X has the same sequence as the mRNA and, in certain conditions, will bind ribosomes so that a protected fragment can be isolated and sequenced. Only one major site was found. By comparison with the amino acid sequence data it was found that this ribosome binding site sequence coded for the initiation of the gene G protein¹⁴ (positions 2,362-2,413).

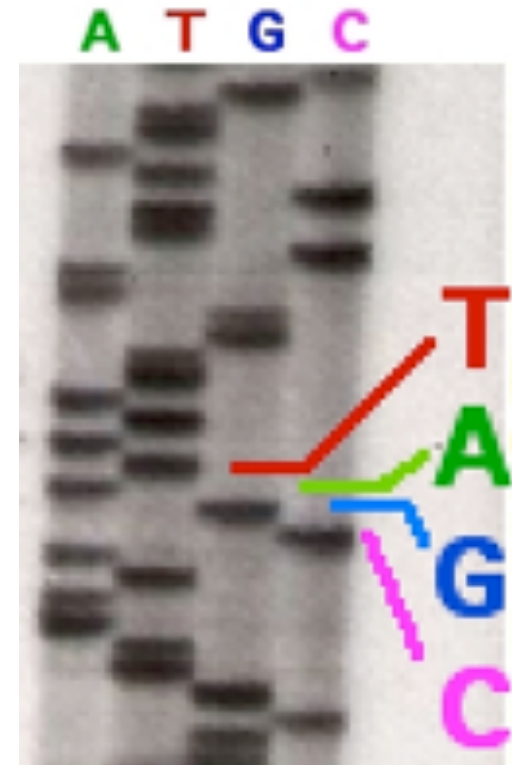
At this stage sequencing techniques using primed synthesis with DNA polymerase were being developed¹⁵ and Schott¹⁷ synthesised a decanucleotide with a sequence complementary to part of the ribosome binding site. This was used to prime into the intergenic region between the F and G genes, using DNA polymerase and ³²P-labelled triphosphates¹⁸. The ribo-substitution technique¹⁶ facilitated the sequence determination of the labelled DNA produced. This decanucleotide-primed system was also used to develop the plus and minus method⁵. Suitable synthetic primers are, however, difficult to prepare and as

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

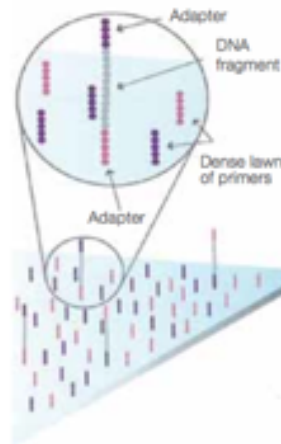


Second Generation Sequencing

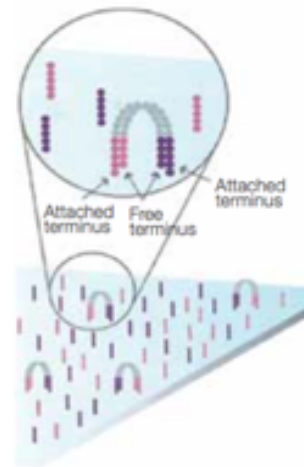


Illumina NovaSeq 6000
Sequencing by Synthesis

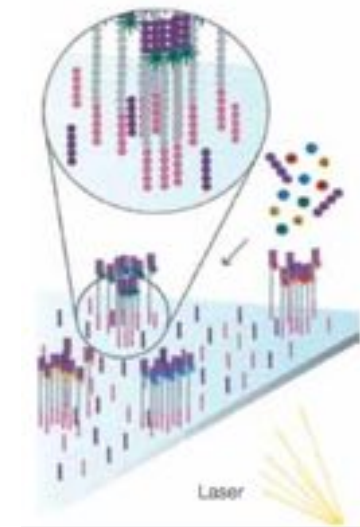
>3Tbp / day



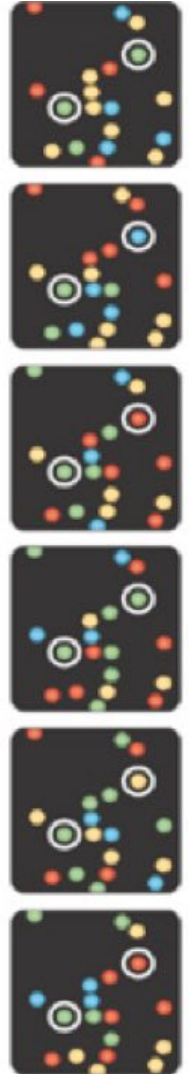
1. Attach



2. Amplify



3. Image



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

Worldwide capacity exceeds 50 Pbp/year
Approximately 1.5M human genomes sequenced

A world map illustrating the distribution of genomic sequencing capacity across various countries. The map uses a color-coded system where red and orange circles indicate higher capacity, while blue circles indicate lower capacity. Numerical values are provided for many countries, representing the capacity in Pbp/year. The map includes labels for major oceans (North Pacific, North Atlantic, South Pacific, South Atlantic, Indian) and various countries across all continents.

Country	Capacity (Pbp/year)
USA	151
Canada	34
Mexico	6
Central America	11, 114, 78, 210, 106, 40, 72, 22
South America	13, 19, 4, 8
Europe	158, 75, 101, 11, 83, 60, 183, 14, 28, 4, 11, 114, 7, 20, 16, 28, 2, 28, 104, 73, 26, 38, 5
Africa	158, 75, 101, 11, 83, 60, 183, 14, 28, 4, 11, 114, 7, 20, 16, 28, 2, 28, 104, 73, 26, 38, 5
Asia	158, 75, 101, 11, 83, 60, 183, 14, 28, 4, 11, 114, 7, 20, 16, 28, 2, 28, 104, 73, 26, 38, 5
Oceania	158, 75, 101, 11, 83, 60, 183, 14, 28, 4, 11, 114, 7, 20, 16, 28, 2, 28, 104, 73, 26, 38, 5

<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

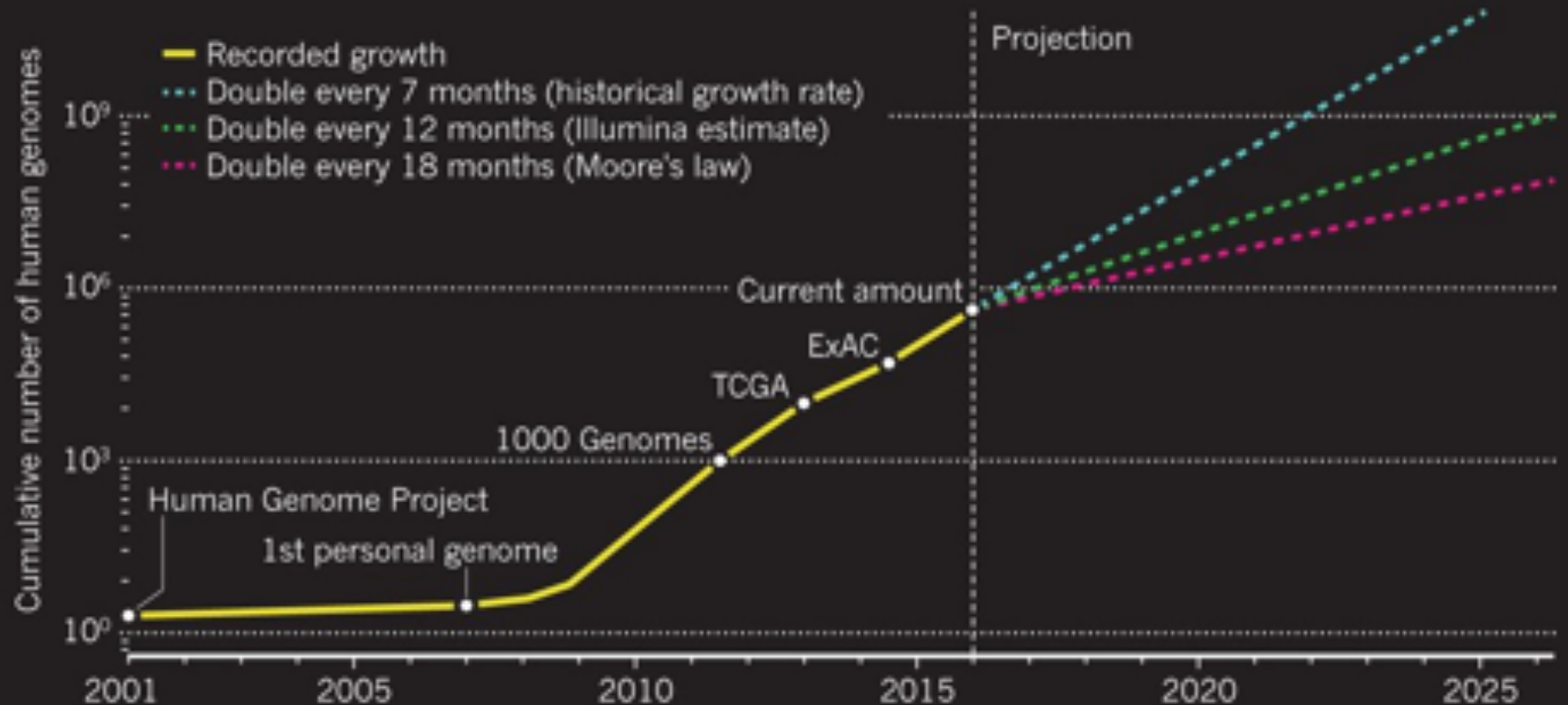


500 2 TB drives
\$100k

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.

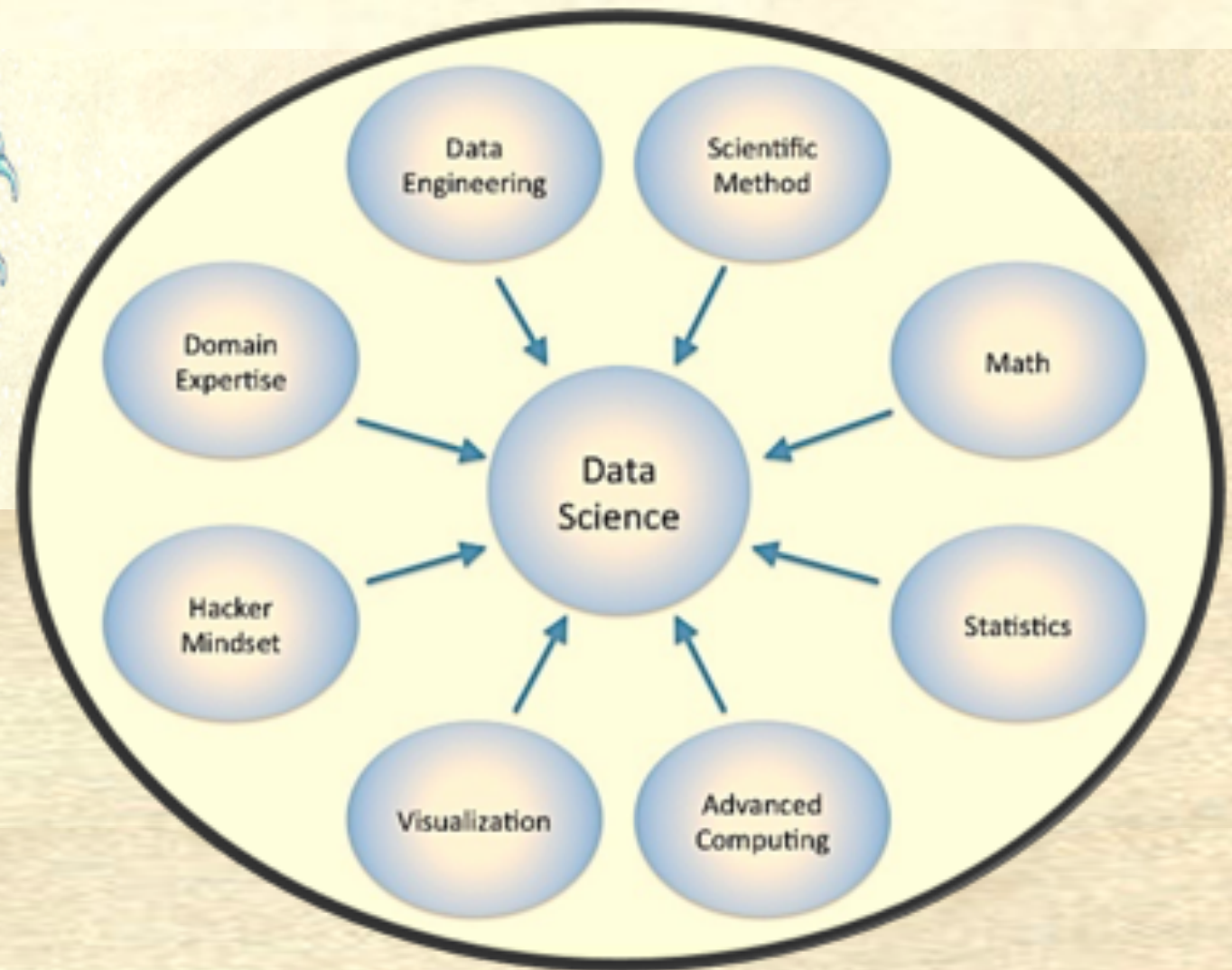
What software and systems will?

And who will create them?

- ***Plus thousands and thousands more***



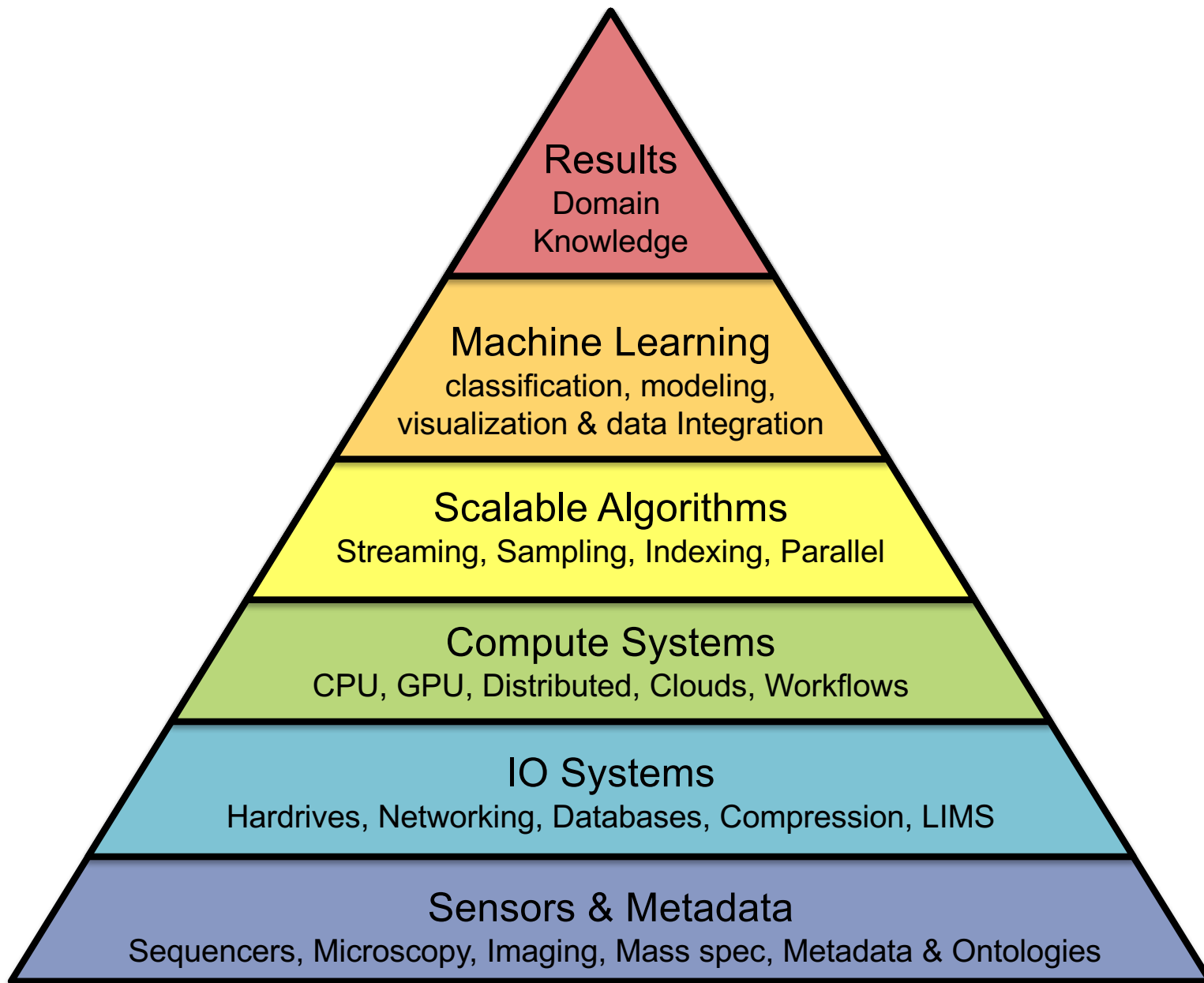
Who is a Data Scientist?



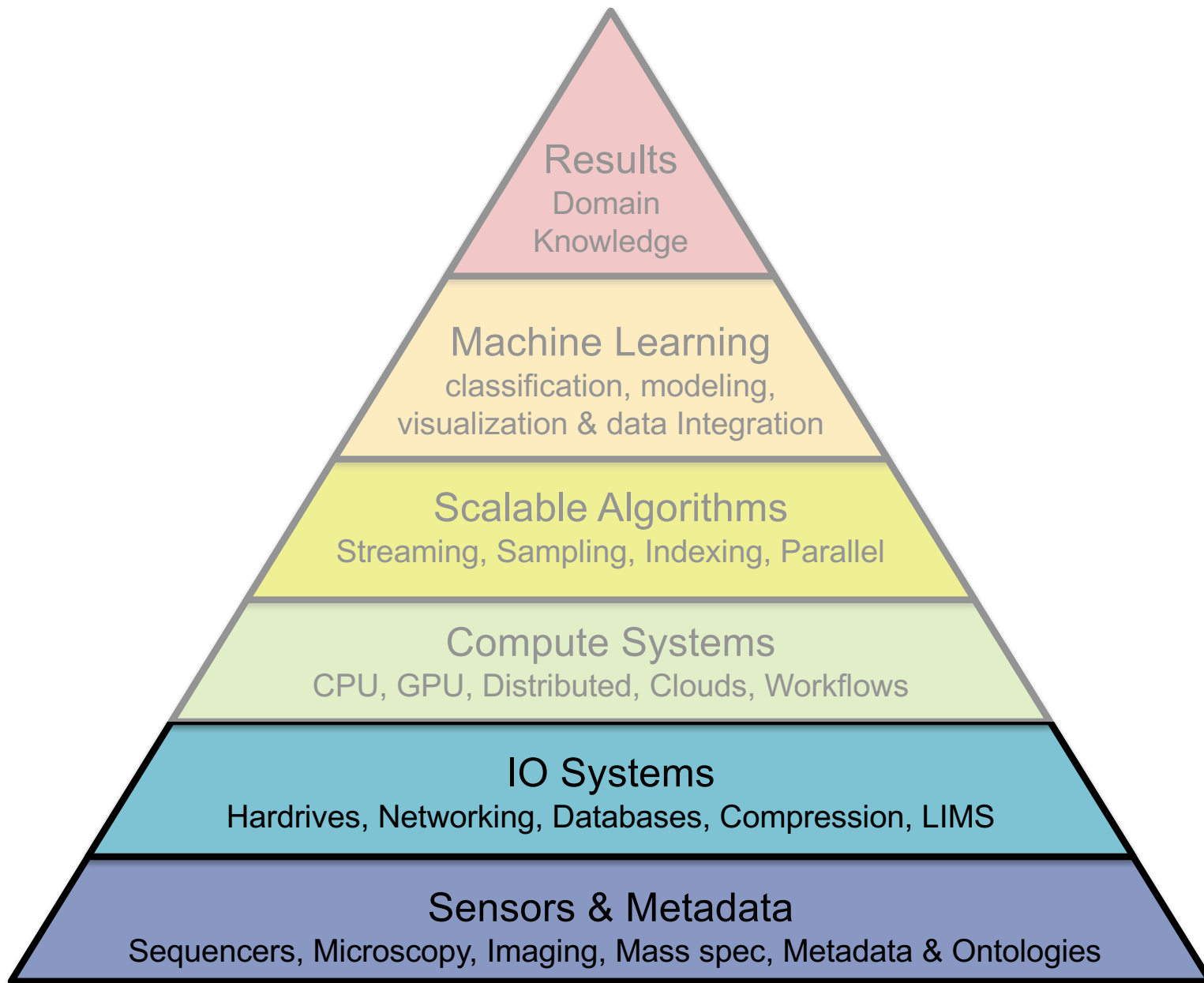
http://en.wikipedia.org/wiki/Data_science



Comparative Genomics Technologies



Comparative Genomics Technologies



Genomics Arsenal in the year 2021

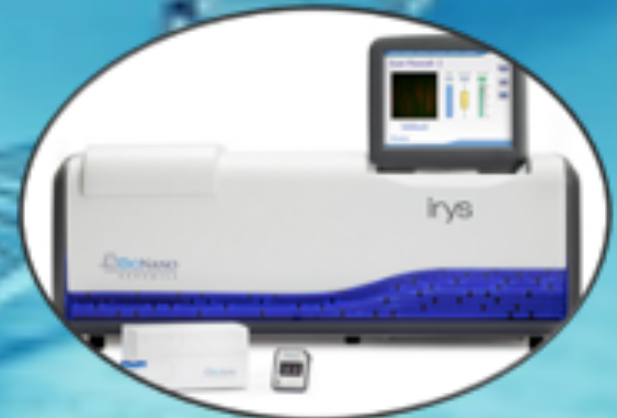
Sample Preparation

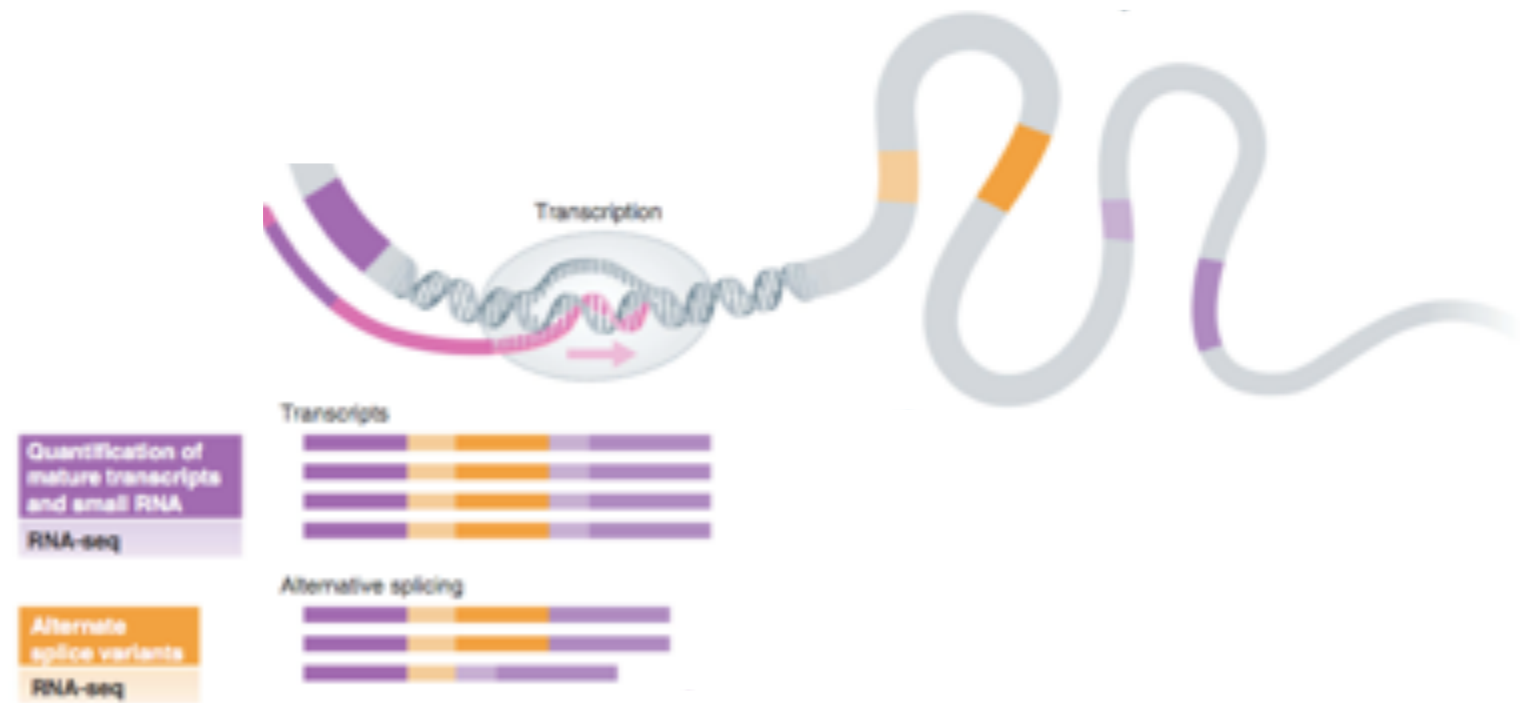


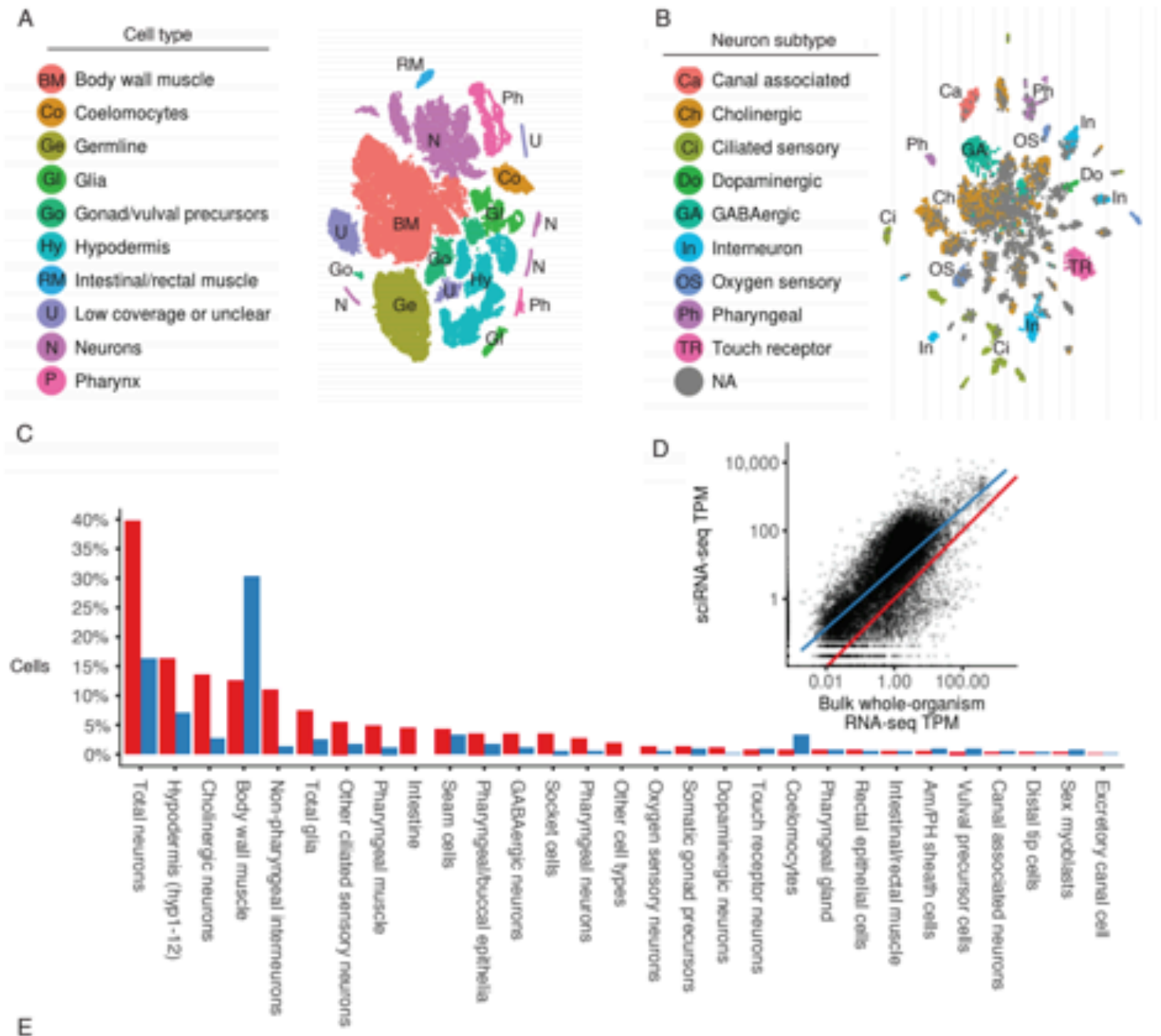
Sequencing



Chromosome Mapping

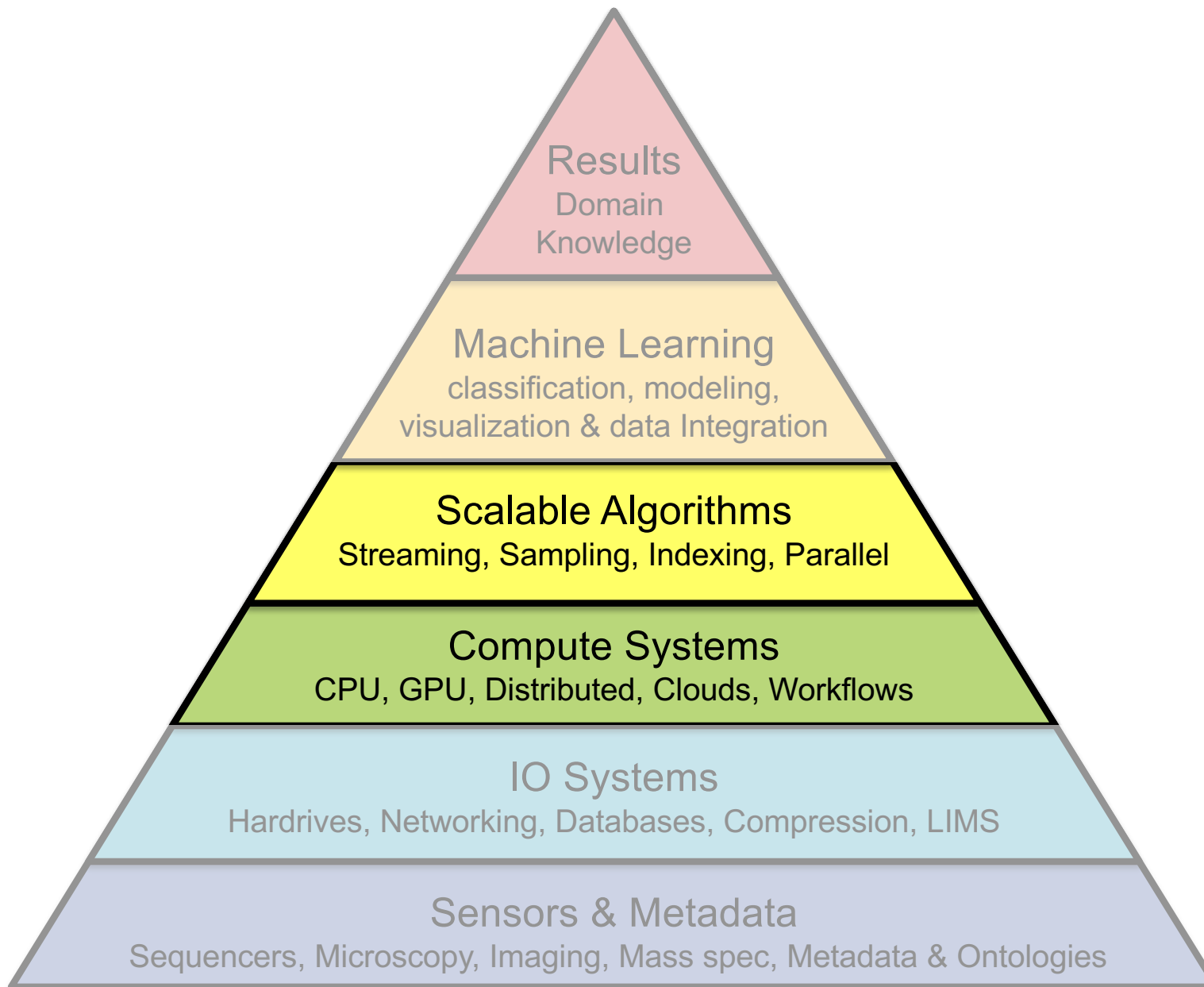






Comprehensive single-cell transcriptional profiling of a multicellular organism
 Cao, et al. (2017) *Science*. doi: 10.1126/science.aam8940

Comparative Genomics Technologies

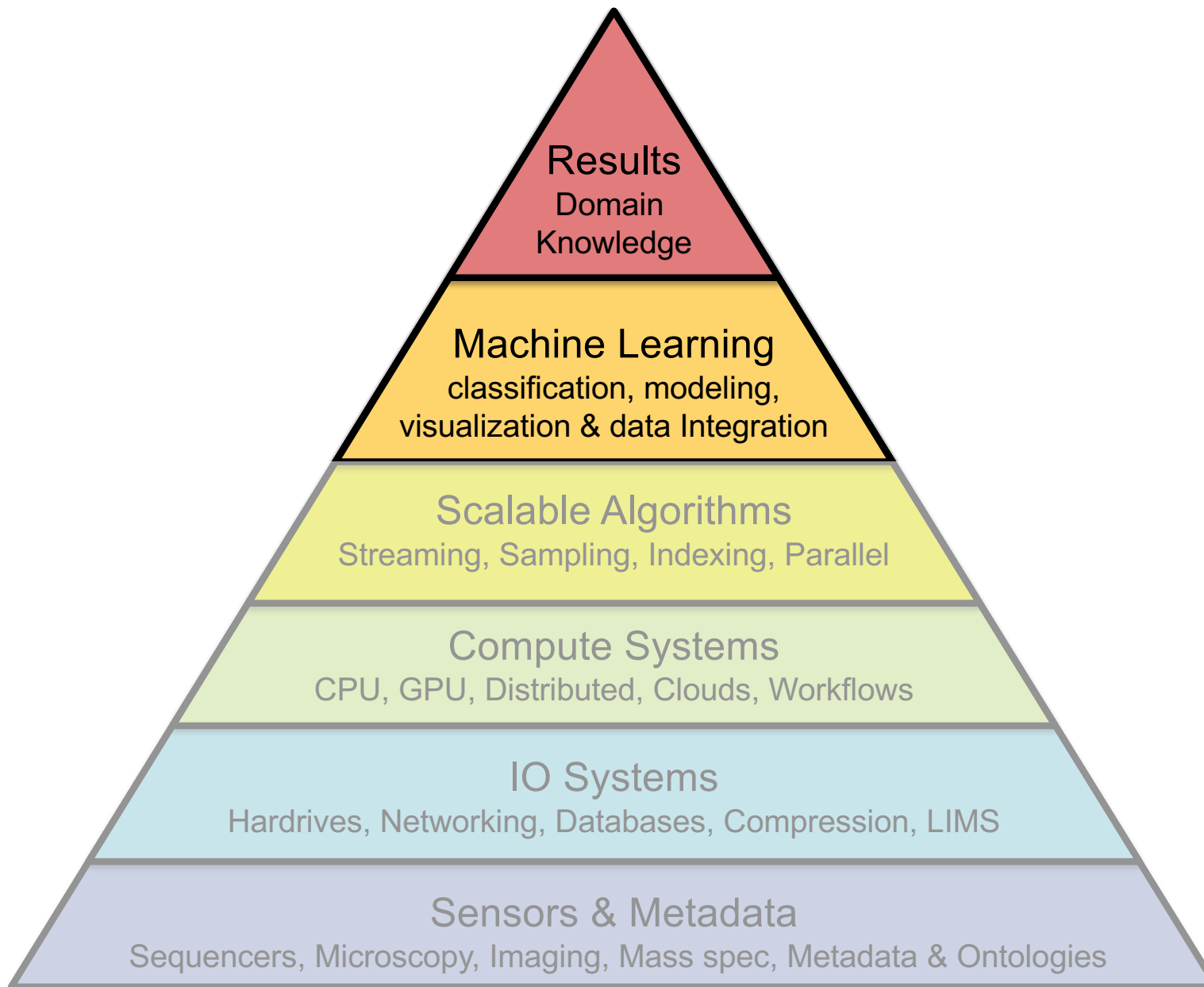


Potential Topics

- 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?

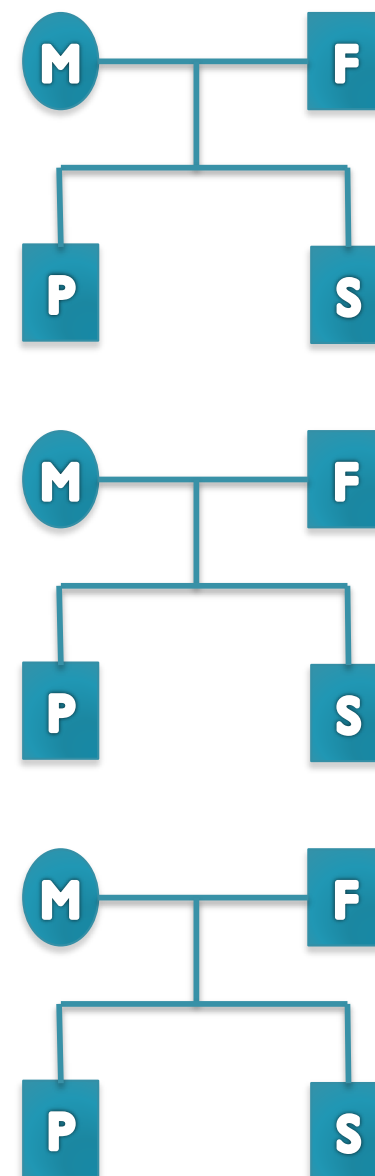
<http://www.autismspeaks.org/what-autism>

Searching for the genetic risk factors

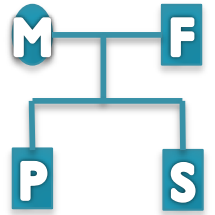
Search Strategy

- 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: . . . **TCAAATCCTTTTAATAAAGAAGAGCTGACA** . . .

Father(1): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Father(2): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Mother(1): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Mother(2): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Sibling(1): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Sibling(2): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

Proband(1): . . . TCAAATCCTTTTAATAAAGAAGAGCTGACA . . .

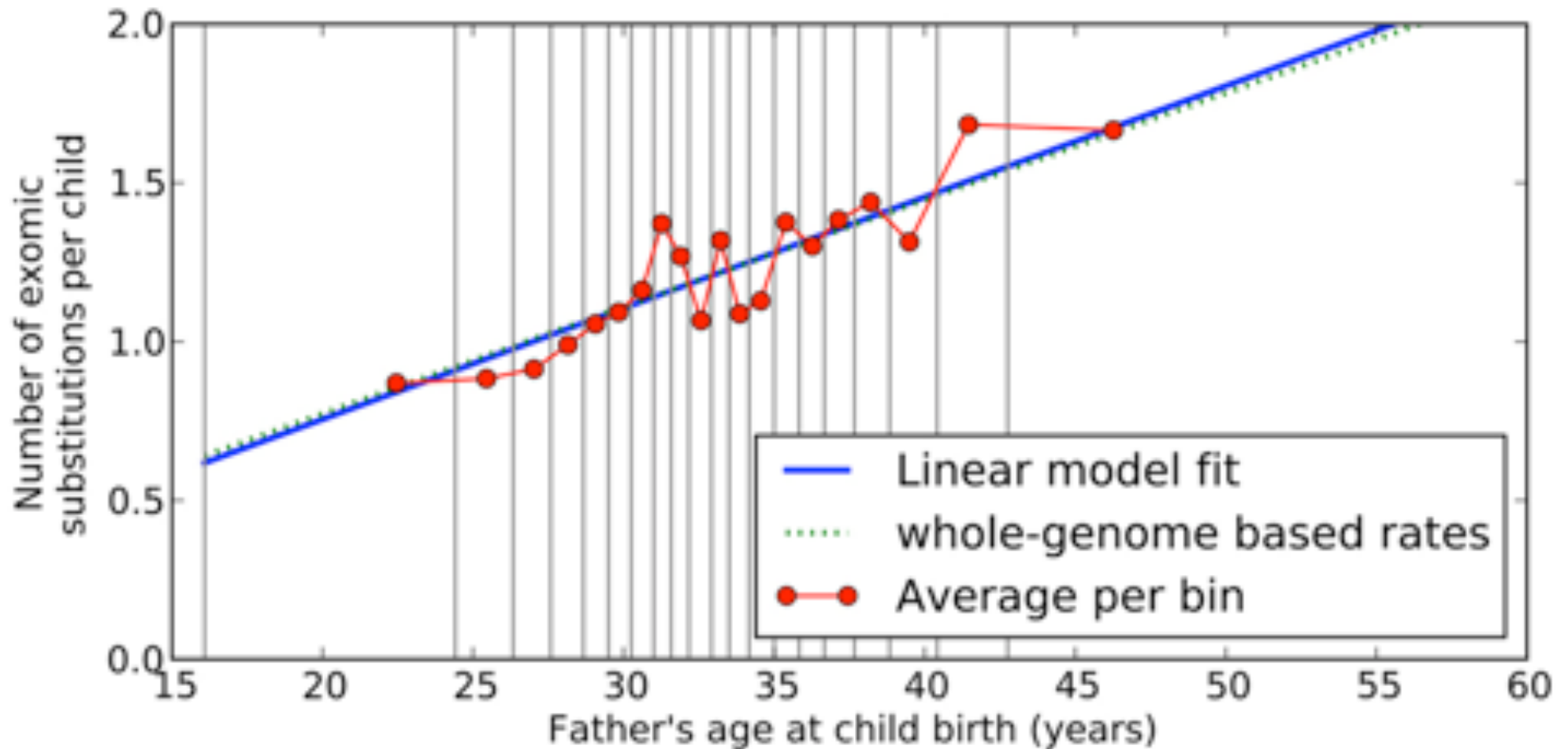
Proband(2): . . . TCAAATCCTTTTAAT****AAGAGCTGACA . . .

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

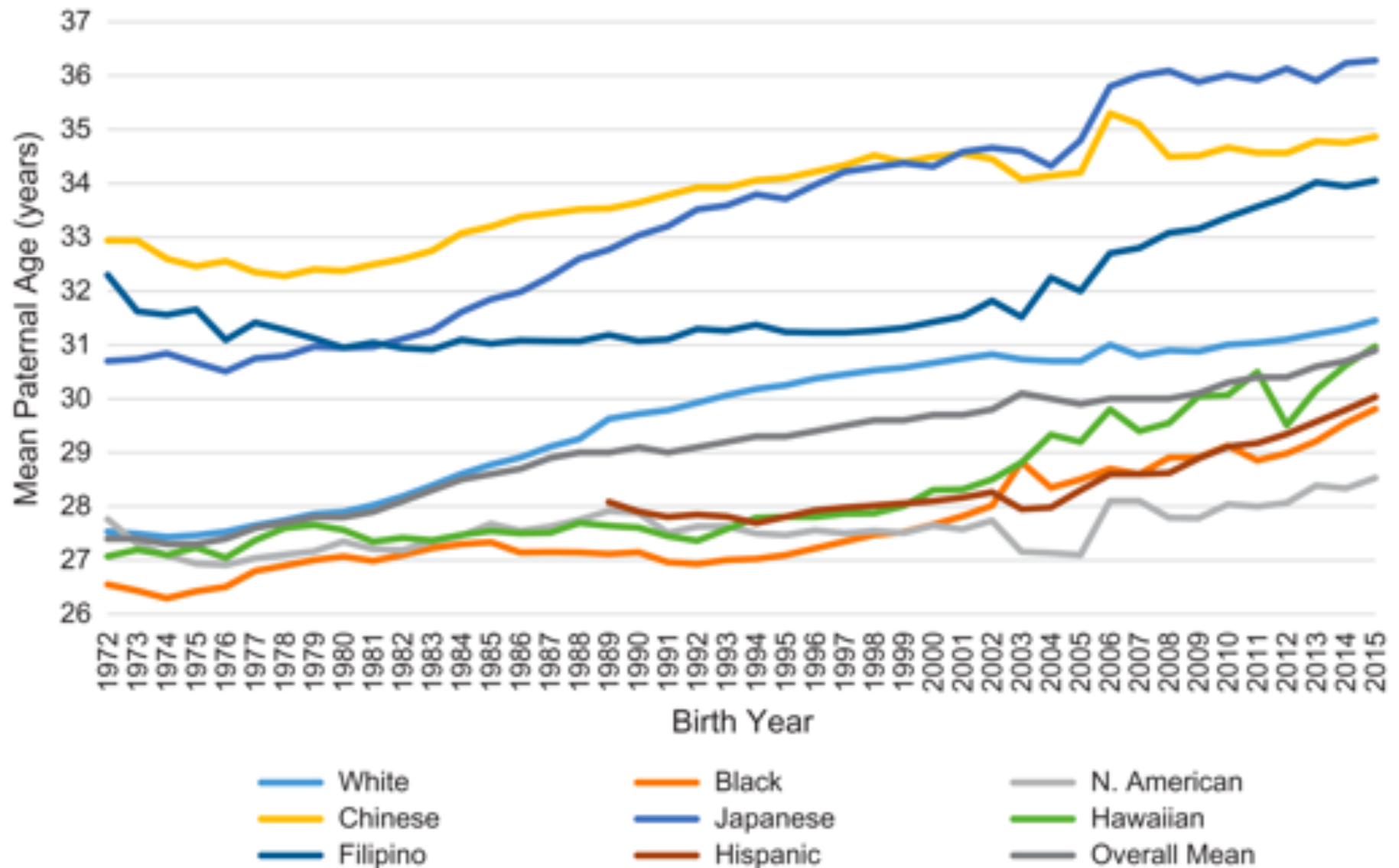
De novo Mutations in Men



The contribution of de novo coding mutations to autism spectrum disorder

lossifov *et al* (2014) *Nature*. doi:10.1038/nature13908

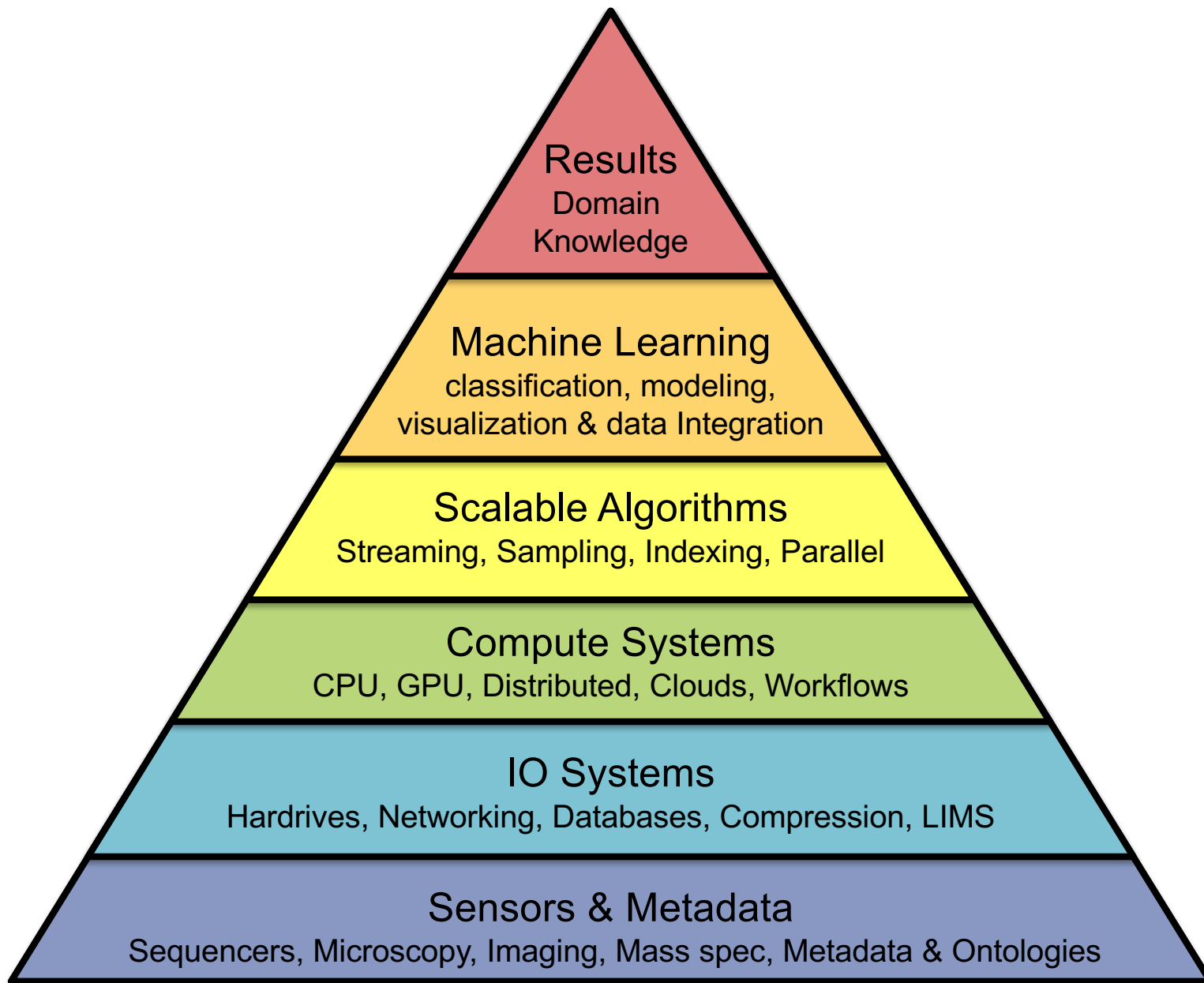
Age of Fatherhood



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