

Genomic Technologies

Michael Schatz

January 27, 2020

Lecture 2: Applied Comparative Genomics



Course Webpage

The screenshot shows a GitHub repository page for 'appliedgenomics2021'. The repository has 73 stars and 0 forks. The 'Code' tab is selected. A commit from 'Mike Schatz' was pushed 9 minutes ago, adding a syllabus. Other commits from 'policies' and 'README.md' were also made 9 minutes ago. The 'About' section describes the repository as 'Materials for Spring 2021 Applied Genomics Course'. It includes a 'Readme' link and sections for 'Releases' (no releases published) and 'Packages' (no packages published). The 'Description' section provides contact information for the professor and TA, class hours, and office hours. The 'Course Overview' section details the primary goal of the course, which is to ground students in theory and empower them to conduct independent genomic analyses using various computational and quantitative approaches.

schatzlab/appliedgenomics2021

https://github.com/schatzlab/appliedgenomics2021

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

master · 0 branches · 0 tags

Go to file Add file · Code

Mike Schatz add syllabus 1146286 · 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.

About

Materials for Spring 2021 Applied Genomics Course

Readme

Releases

No releases published

Create a new release

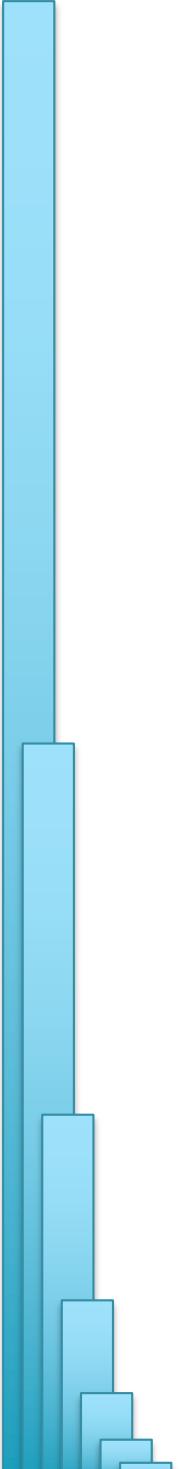
Packages

No packages published

Publish your first package

Course Overview

https://github.com/schatzlab/appliedgenomics2021



TA:Arun Das



Check Piazza for Office Hours Poll

Assignment 1

appliedgenomics2021/assignment1

github.com/schatzlab/appliedgenomics2021/tree/master/assignments/assignment1

Assignment 1: Chromosome Structures

Assignment Date: Wednesday, Jan 27, 2021
Due Date: Wednesday, Feb. 3, 2021 @ 11:59pm

Assignment Overview

In this assignment you will profile the overall structure of the genomes of several important species and then study human chromosome 22 in more detail. As a reminder, any questions about the assignment should be posted to [Piazza](#).

Question 1: Chromosome structures

Download the chromosome size files for the following genomes (Note these have been preprocessed to only include main chromosomes):

1. [Arabidopsis thaliana \(TAIR10\)](#) - An important plant model species [\[info\]](#)
2. [Tomato \(Solanum lycopersicum v4.00\)](#) - One of the most important food crops [\[info\]](#)
3. [E. coli \(Escherichia coli K12\)](#) - One of the most commonly studied bacteria [\[info\]](#)
4. [Fruit Fly \(Drosophila melanogaster, dm3\)](#) - One of the most important model species for genetics [\[info\]](#)
5. [Human \(hg38\) - us :\)](#) [\[info\]](#)
6. [Wheat \(Triticum aestivum, IWGSC\)](#) - The food crop which takes up the largest land area [\[info\]](#)
7. [Worm \(Caenorhabditis elegans, ce10\)](#) - One of the most important animal model species [\[info\]](#)
8. [Yeast \(Saccharomyces cerevisiae, sacCer3\)](#) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)

Using these files, make a table with the following information per species:

- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes



<https://github.com/schatzlab/appliedgenomics2021/tree/master/assignments/assignment1>
Due end of day on Feb 3 (right before midnight)

Unsolved Questions in Biology

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

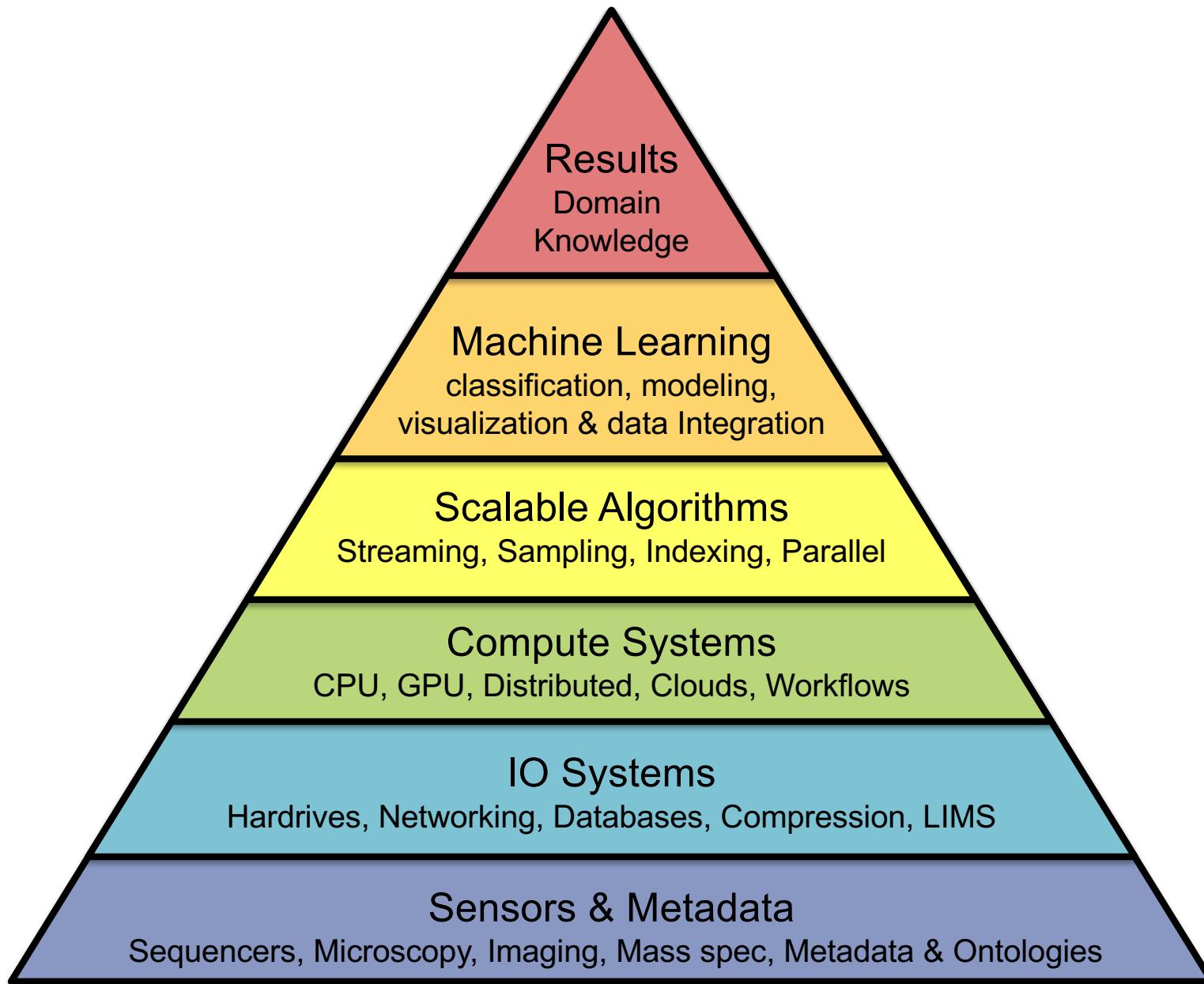
What software and systems will?

And who will create them?

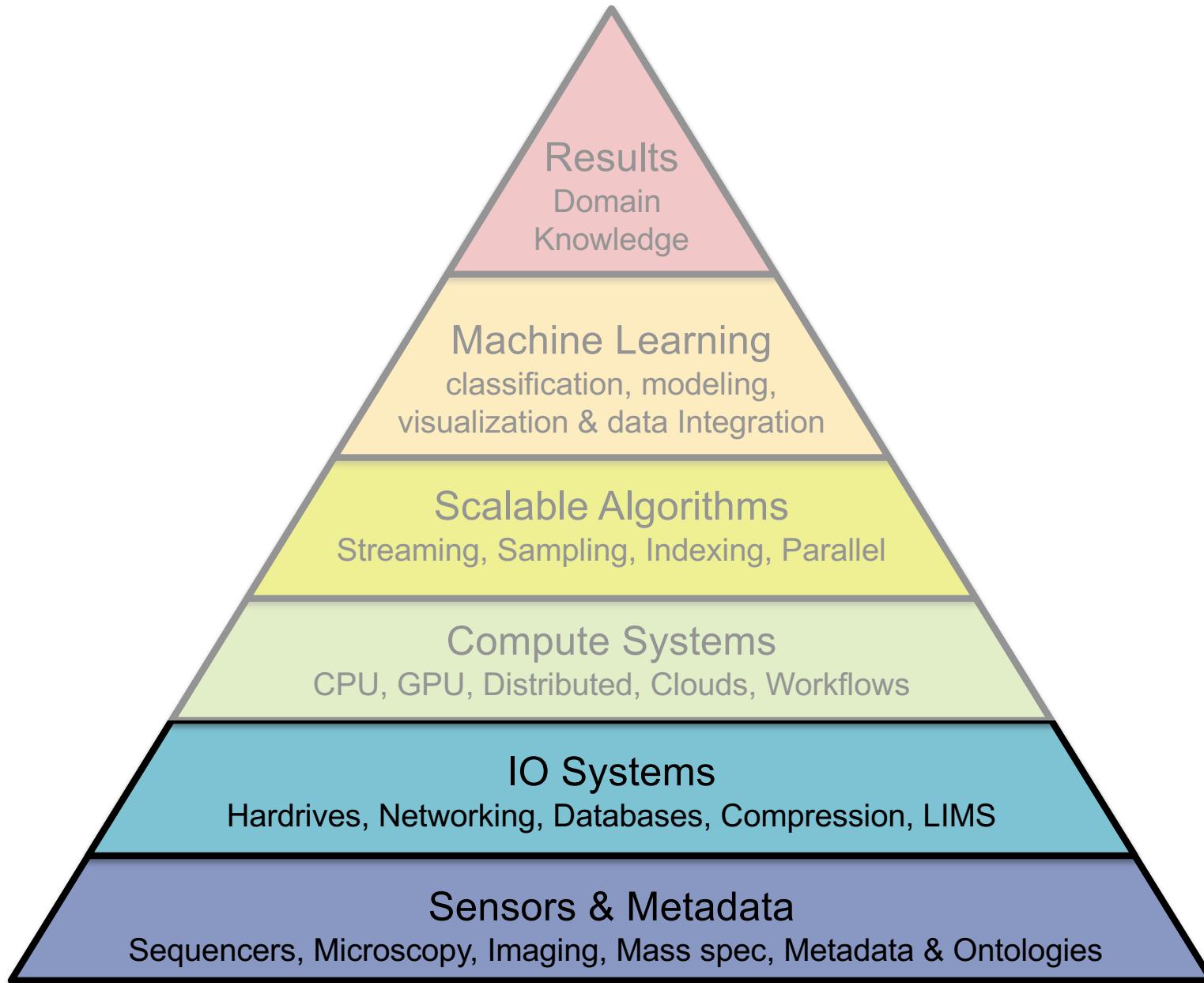
• Plus thousands and thousands more



Comparative Genomics Technologies



Comparative Genomics Technologies



Genomics Arsenal in the year 2021

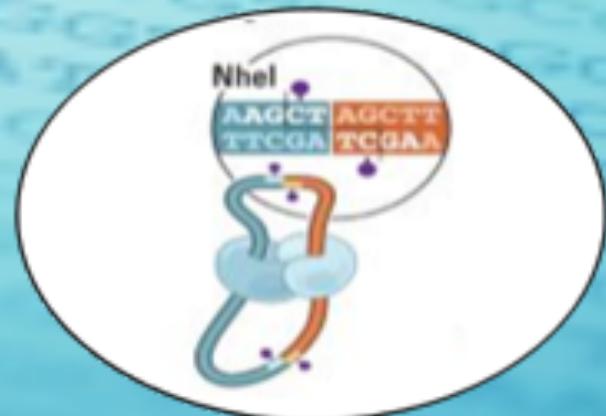
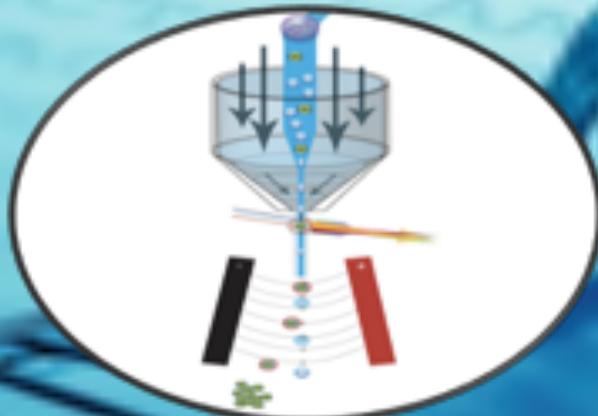
Sample Preparation

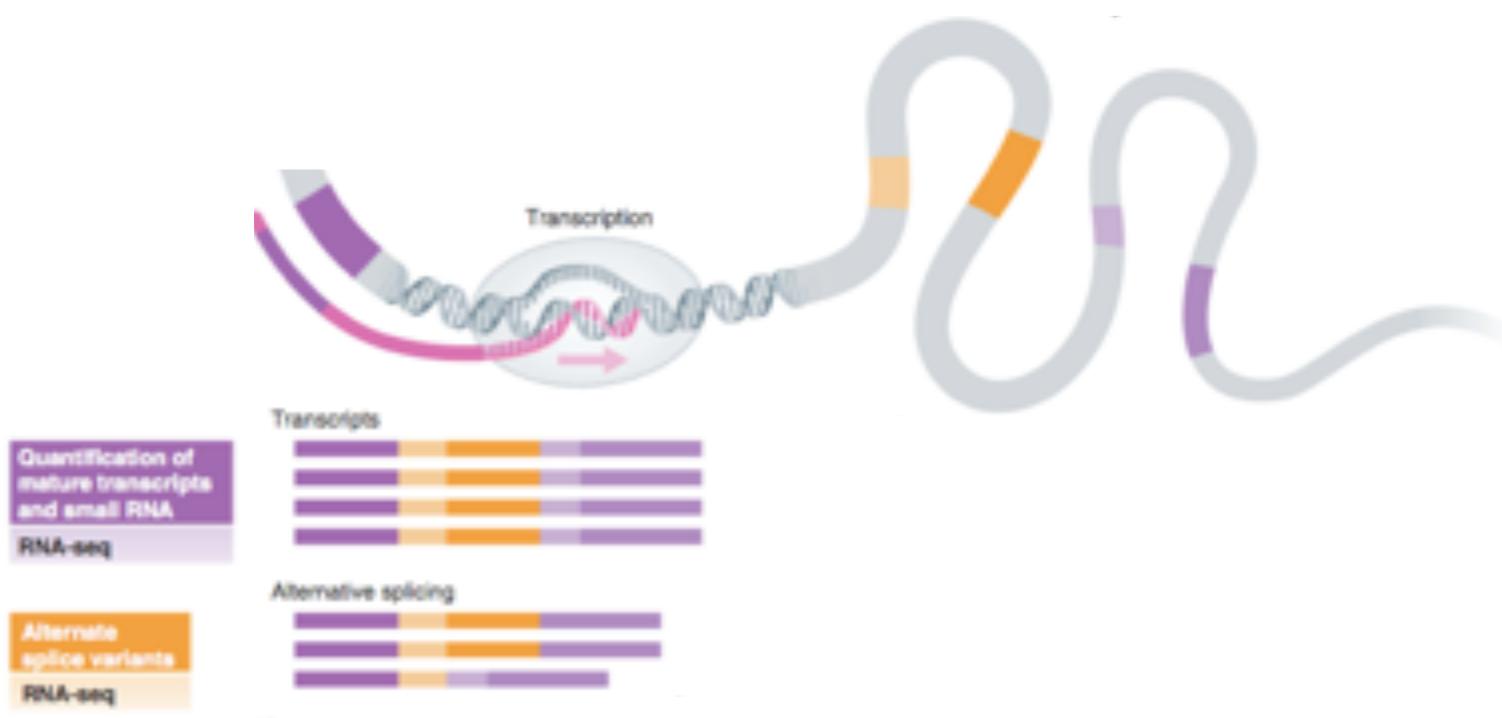


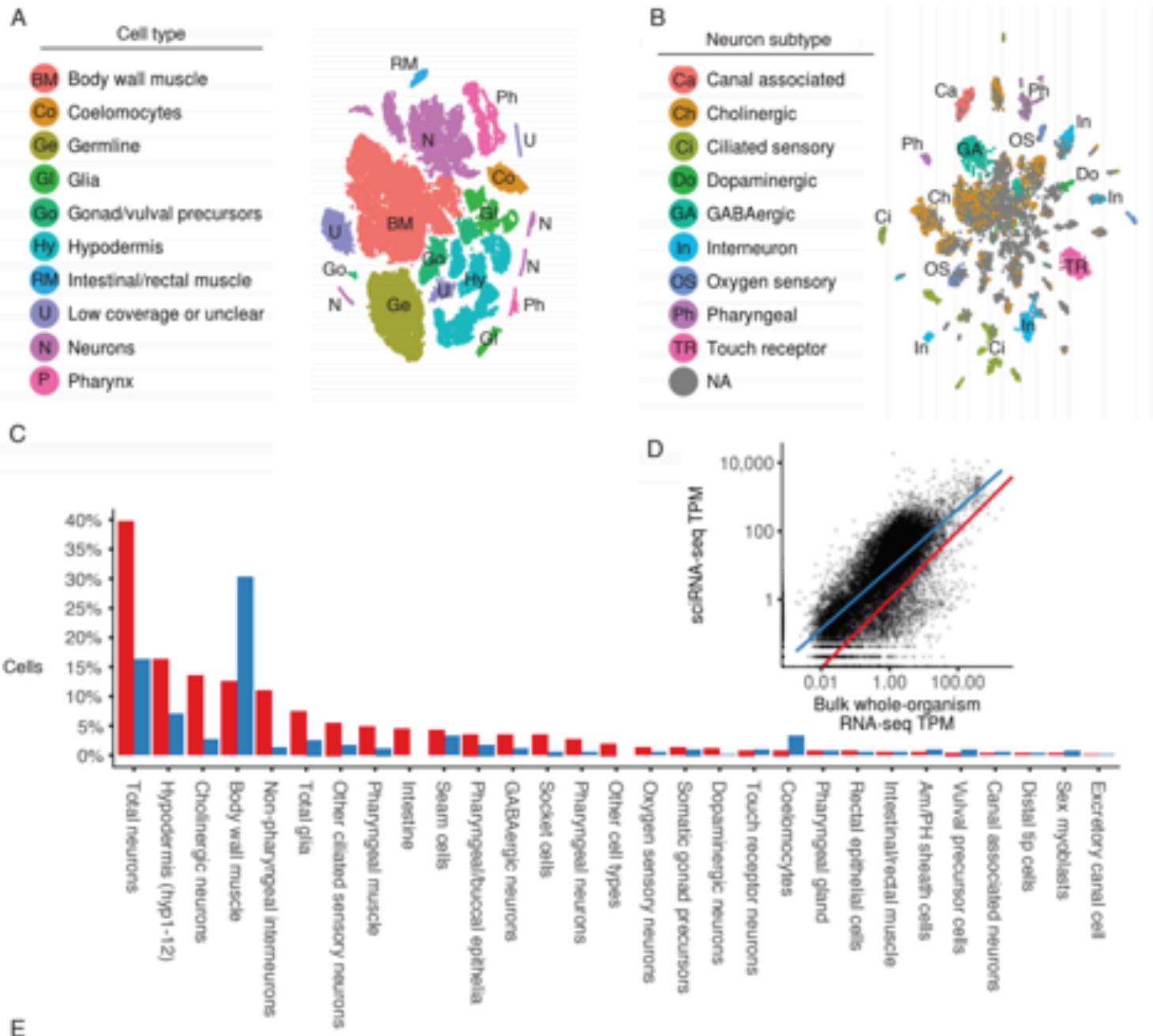
Sequencing



Chromosome Mapping





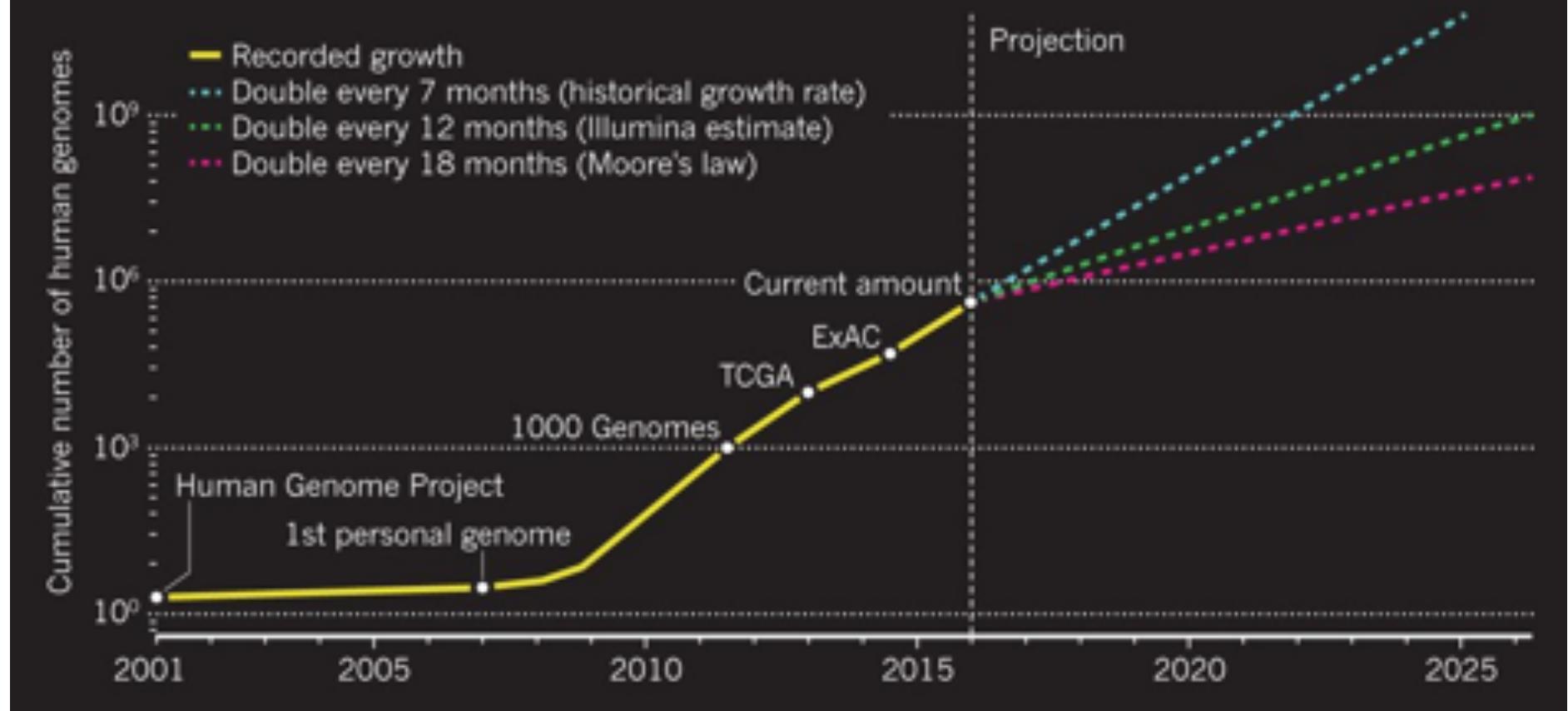


Comprehensive single-cell transcriptional profiling of a multicellular organism
 Cao, et al. (2017) Science. doi: 10.1126/science.aam8940

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

Cost per Genome

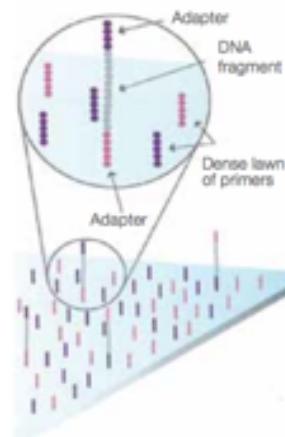


Second Generation Sequencing

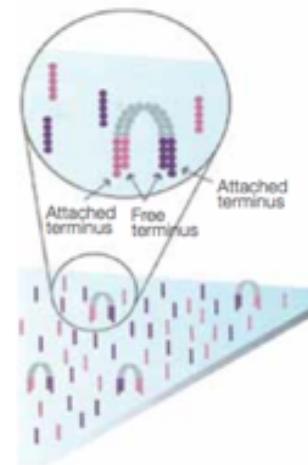


Illumina NovaSeq 6000
Sequencing by Synthesis

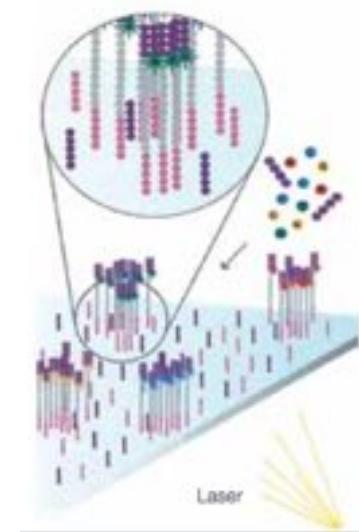
>3Tbp / day



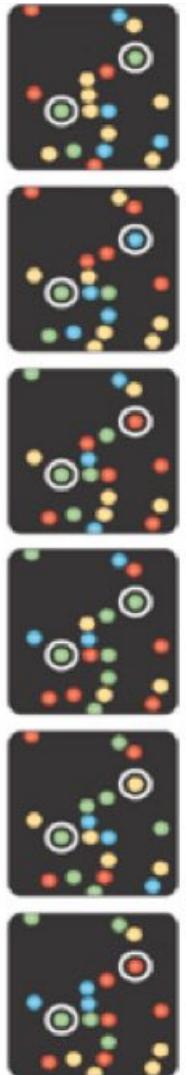
1. Attach



2. Amplify



3. Image

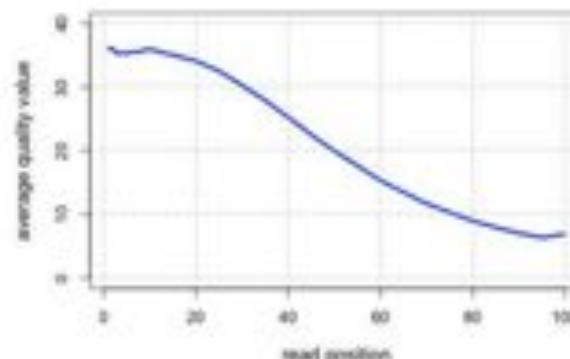


Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina Quality

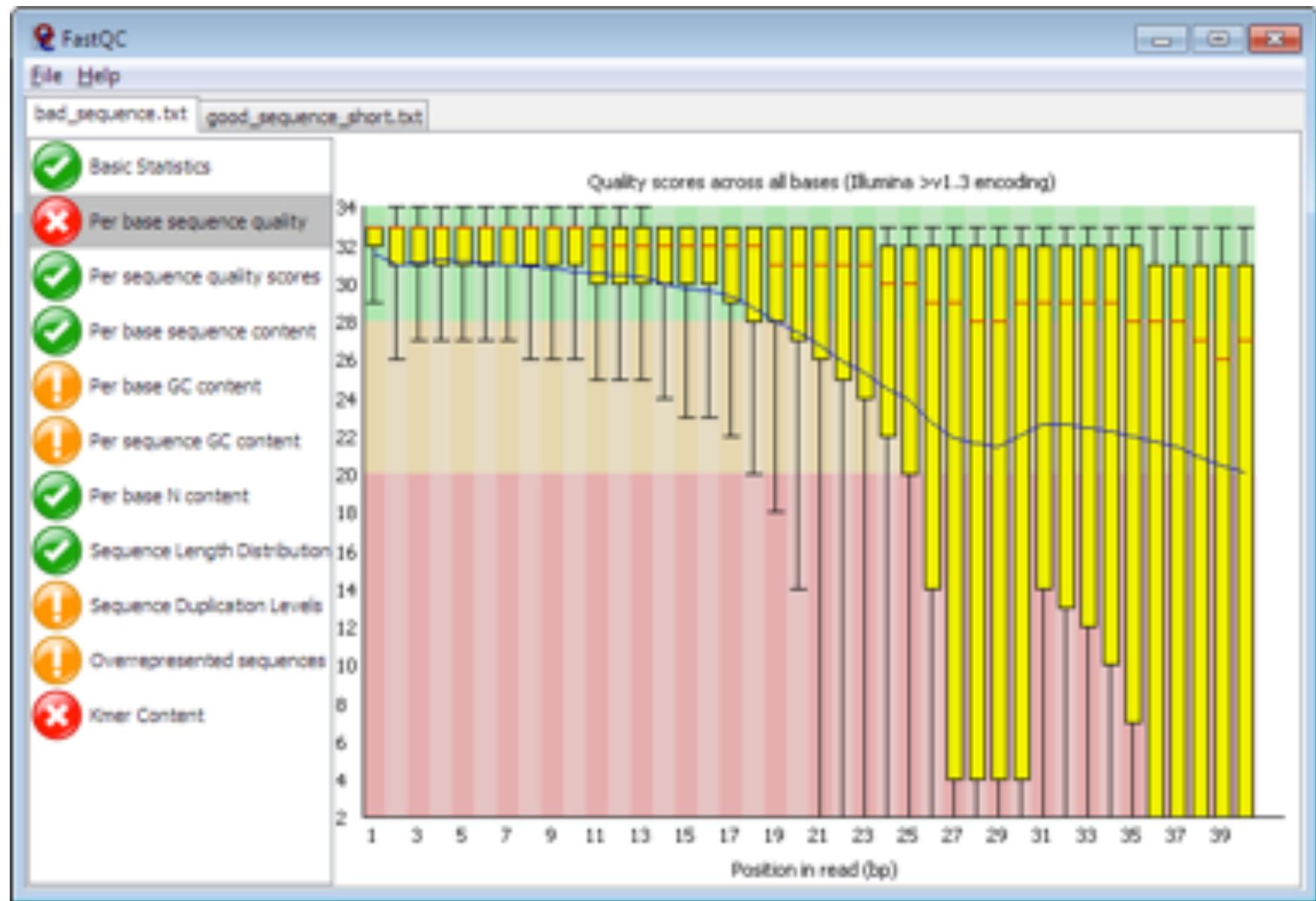
QV	p _{error}
40	1/10000
30	1/1000
20	1/100
10	1/10

$$Q_{\text{sanger}} = -10 \log_{10} p$$



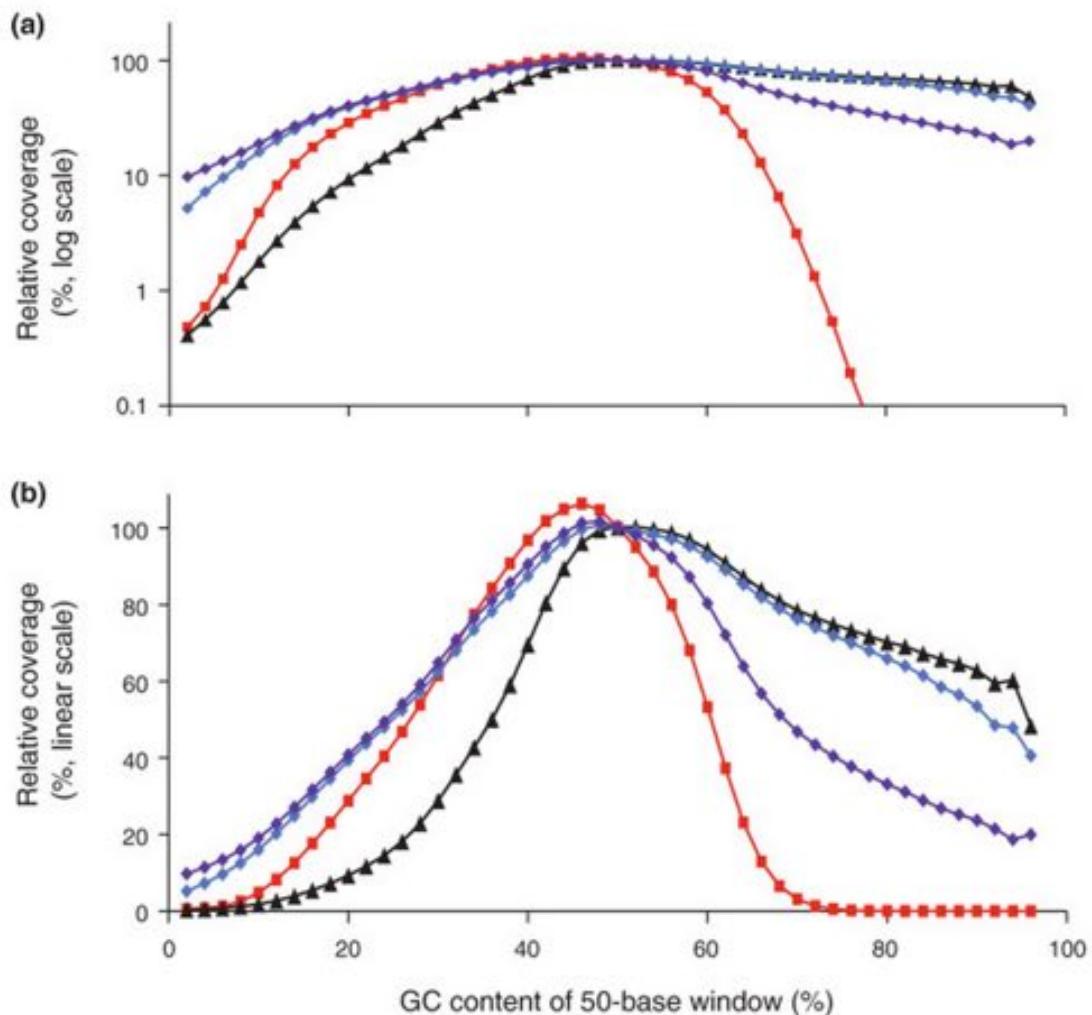
S - Sanger Phred+33, raw reads typically (0, 40)
 X - Solexa Solexa+64, raw reads typically (-5, 40)
 I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
 J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
 with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (**bold**)
 (Note: See discussion above).
 L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

FASTQC: Is my data any good?



<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

Beware of GC Biases



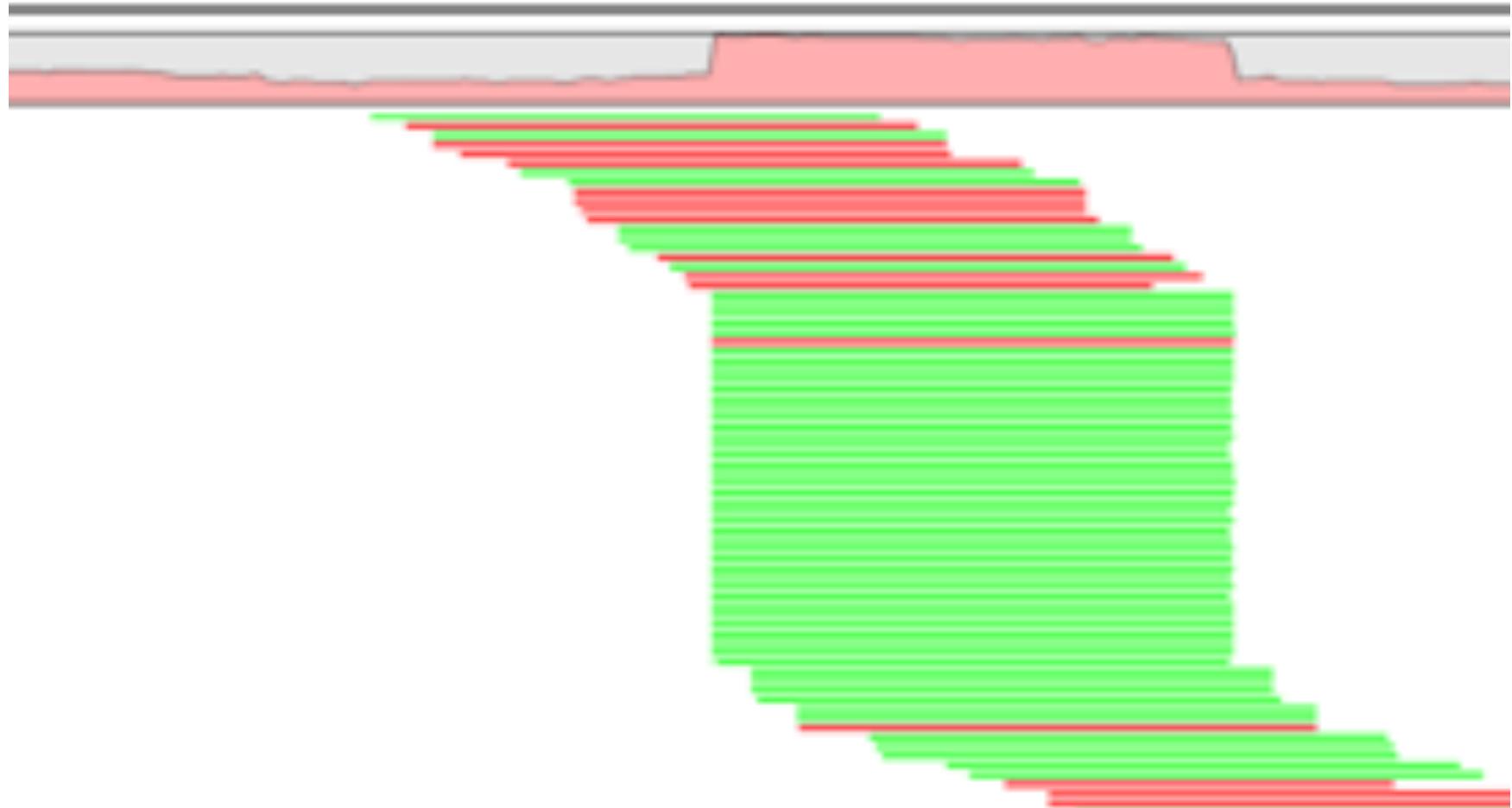
Illumina sequencing does not produce uniform coverage over the genome

- Coverage of extremely high or extremely low GC content will have reduced coverage in Illumina sequencing
- Biases primarily introduced during PCR; lower temperatures, slower heating, and fewer rounds minimize biases
- This makes it very difficult to identify variants (SNPs, CNVs, etc) in certain regions of the genome

Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries.

Aird et al. (2011) *Genome Biology*. 12:R18.

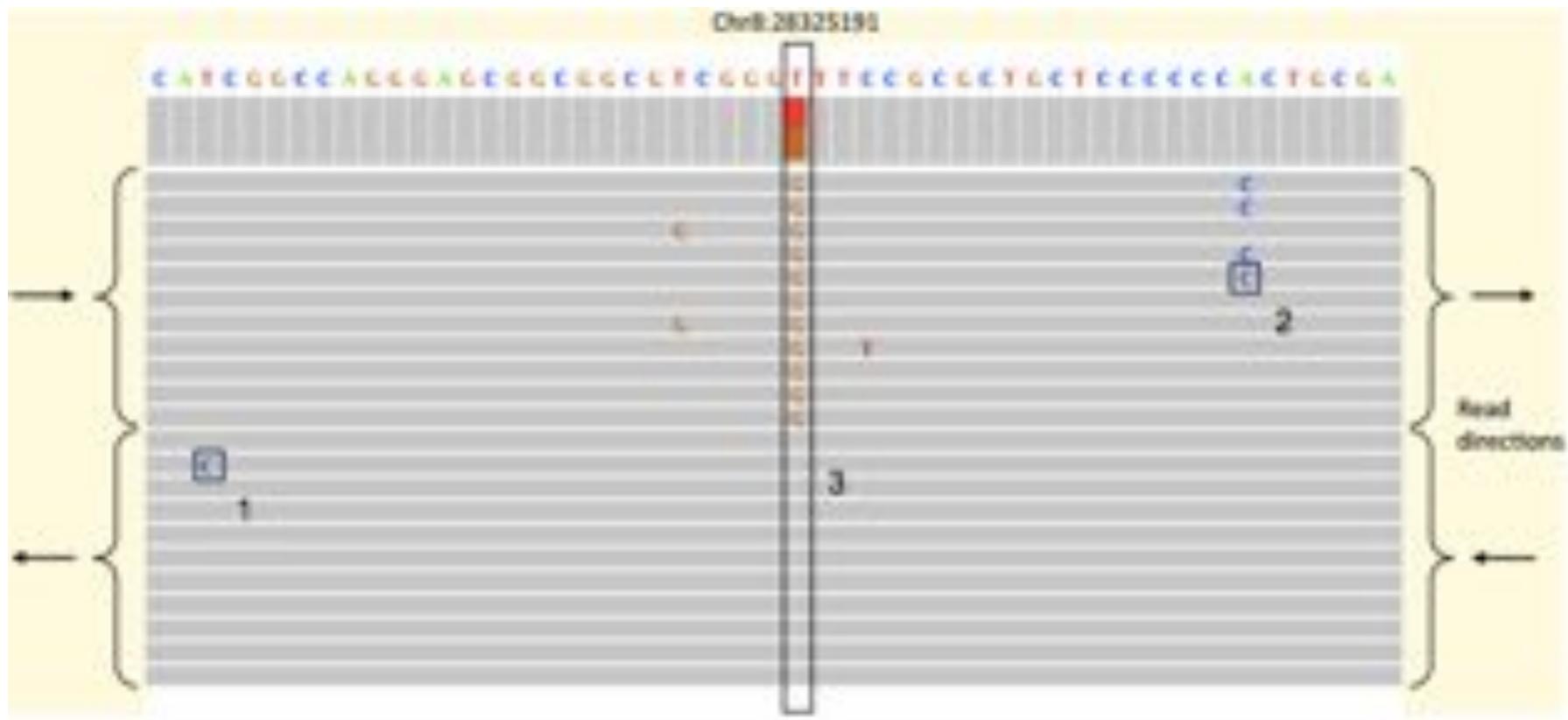
Beware of Duplicate Reads



The Sequence alignment/map (SAM) format and SAMtools.
Li et al. (2009) *Bioinformatics*. 25:2078-9

Picard: <http://picard.sourceforge.net>

Beware of (Systematic) Errors



Identification and correction of systematic error in high-throughput sequence data
Meacham et al. (2011) *BMC Bioinformatics*. 12:451

A closer look at RNA editing.
Lior Pachter (2012) *Nature Biotechnology*. 30:246-247

Question?

We would love to generate
longer and longer reads with this technology

What can we do?

Paired-end and Mate-pairs

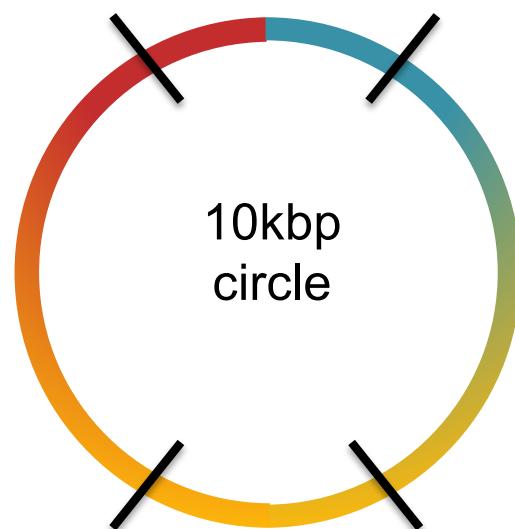
Paired-end sequencing

- Read one end of the molecule, flip, and read the other end
- Generate pair of reads separated by up to 500bp with inward orientation



Mate-pair sequencing

- Circularize long molecules (1-10kbp), shear into fragments, & sequence
- Mate failures create short paired-end reads



2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



FASTQ Files



```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
! ' * ( ( ( ****+ ) ) % % % ++ ) ( % % % ) . 1 * * * - + * ' ' ) ) * * 55CCF>>>>>cccccccc65
```

@Identifier
Sequence
+Separator
Quality Values
...

Illumina Sequencing Summary

Advantages:

- Best throughput, accuracy and read length for any 2nd gen. sequencer
- Fast & robust library preparation



Illumina HiSeq

~3 billion paired 100bp reads

~600Gb, \$10K, 8 days

(or “rapid run” ~90Gb in 1-2 days)

Illumina X Ten / NovaSeq

~6 billion paired 150bp reads

1.8Tb, <3 days, ~1000 / genome(\$\$)

(or “rapid run” ~90Gb in 1-2 days)

Illumina NextSeq

One human genome in <30 hours

Market Summary > Illumina, Inc.

NASDAQ: ILMN

+ Follow

420.35 USD **-0.65 (0.15%)** ↓

Jan 26, 3:56 PM EST · Disclaimer

1 day

5 days

1 month

6 months

YTD

1 year

5 years

Max

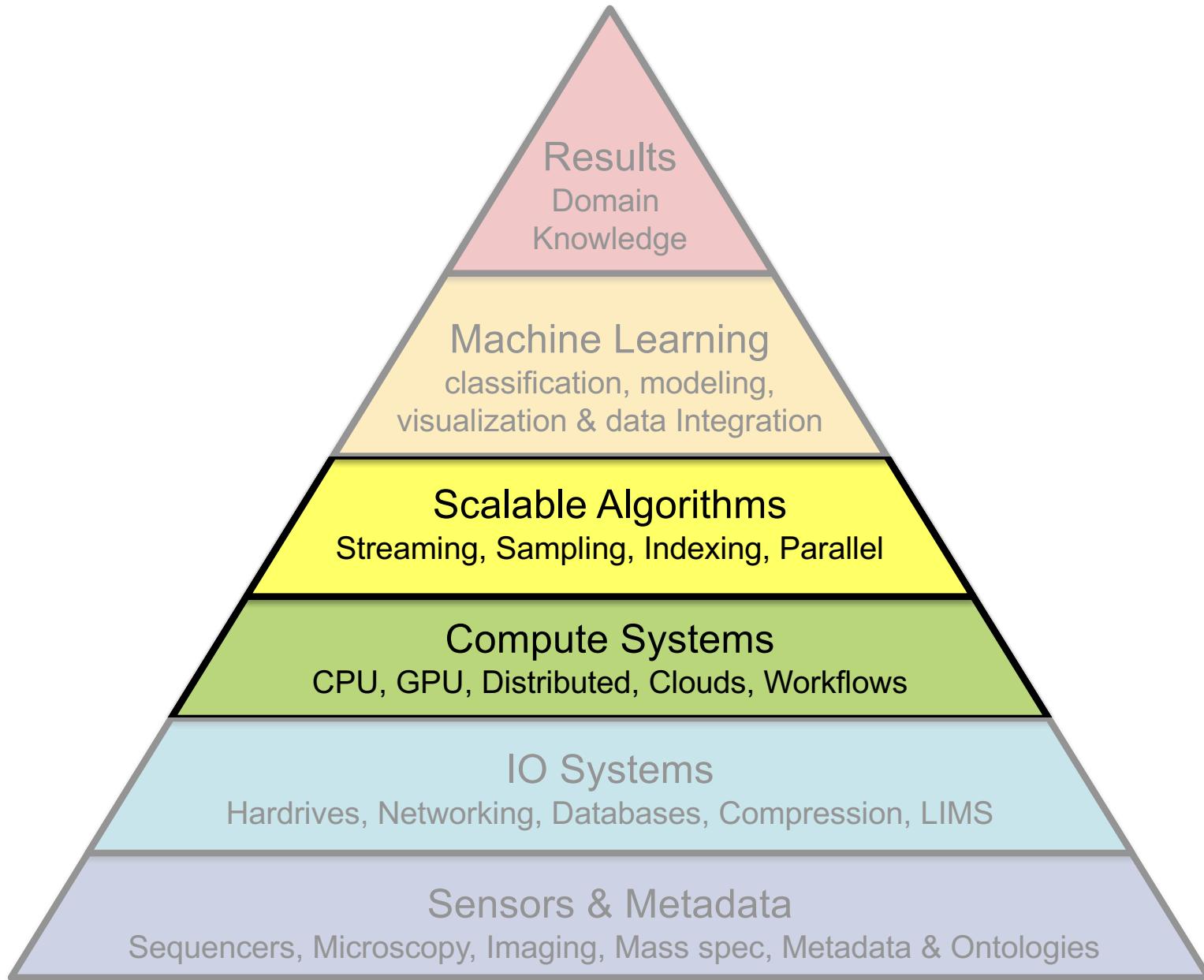


Open 418.95
High 425.92
Low 416.29
Mkt cap 61.33B
P/E ratio 97.70

Div yield -
Prev close 421.00
52-wk high 425.92
52-wk low 196.78

More about Illumina, Inc.

Comparative Genomics Technologies

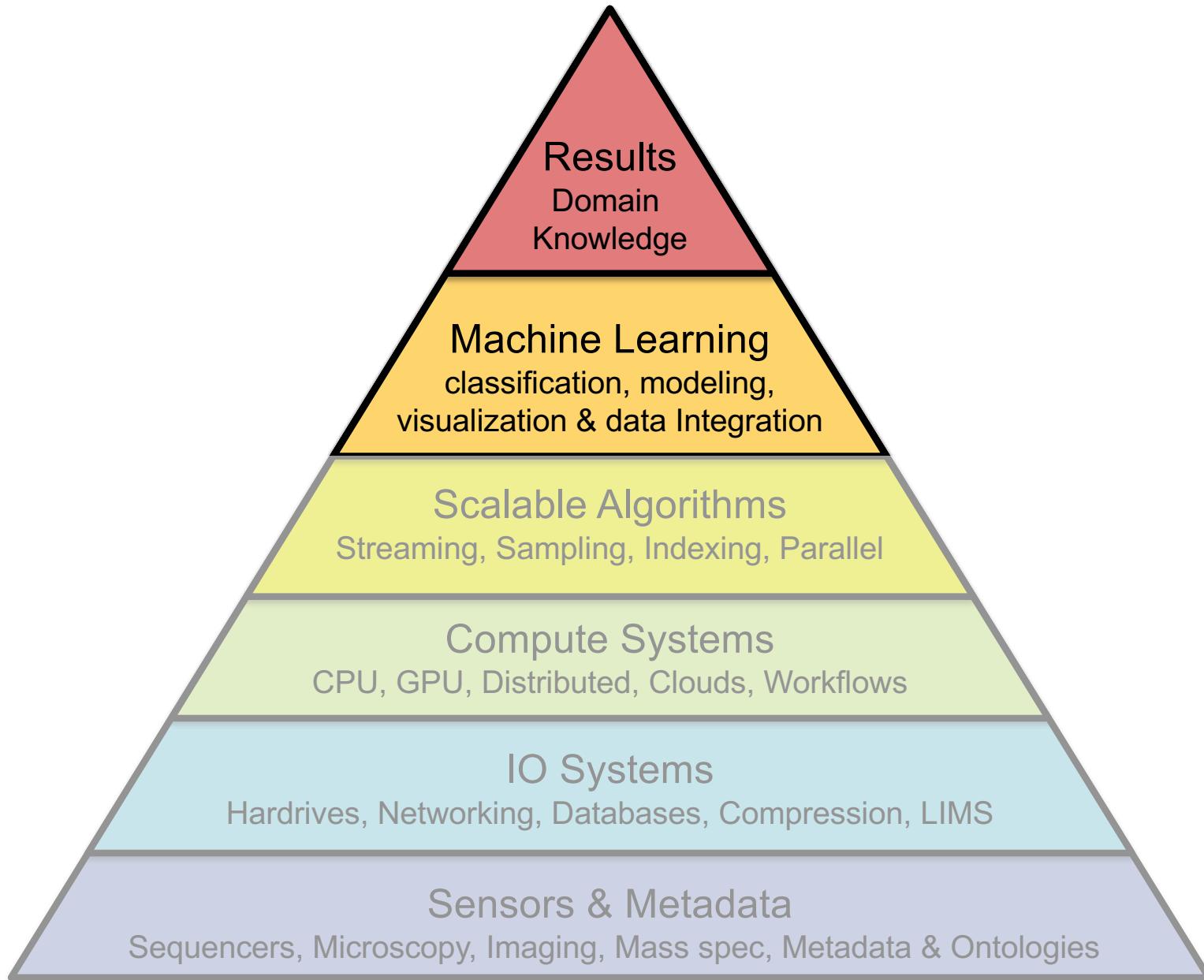


Selected Topics

- 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?

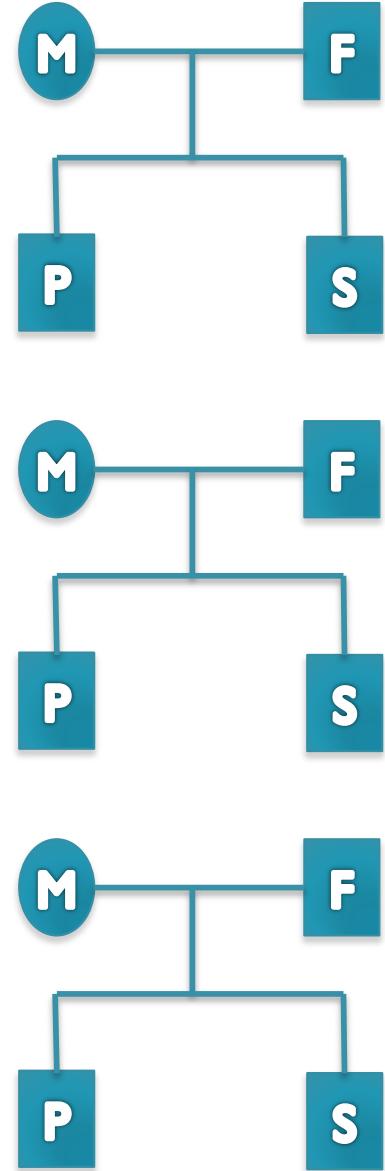
<http://www.autismspeaks.org/what-autism>

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

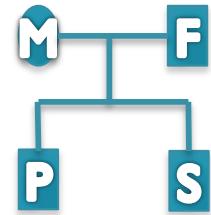
Are there any genetic variants present in affected children, that are not in their parents or unaffected siblings?



De novo mutation discovery and validation

De novo mutations:

Sequences not inherited from your parents.



Reference: . . . 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

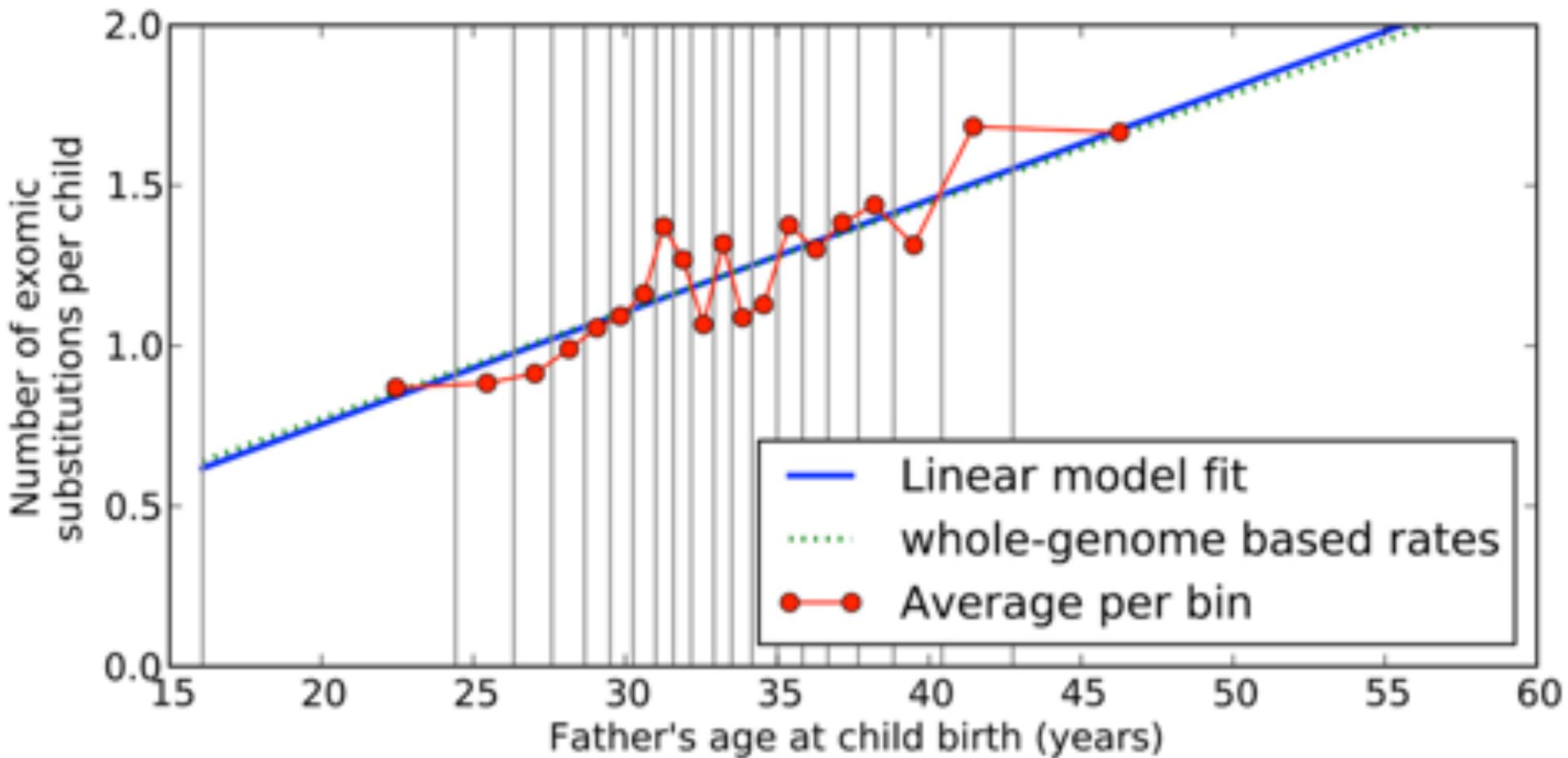
De novo Genetics of Autism

- 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.

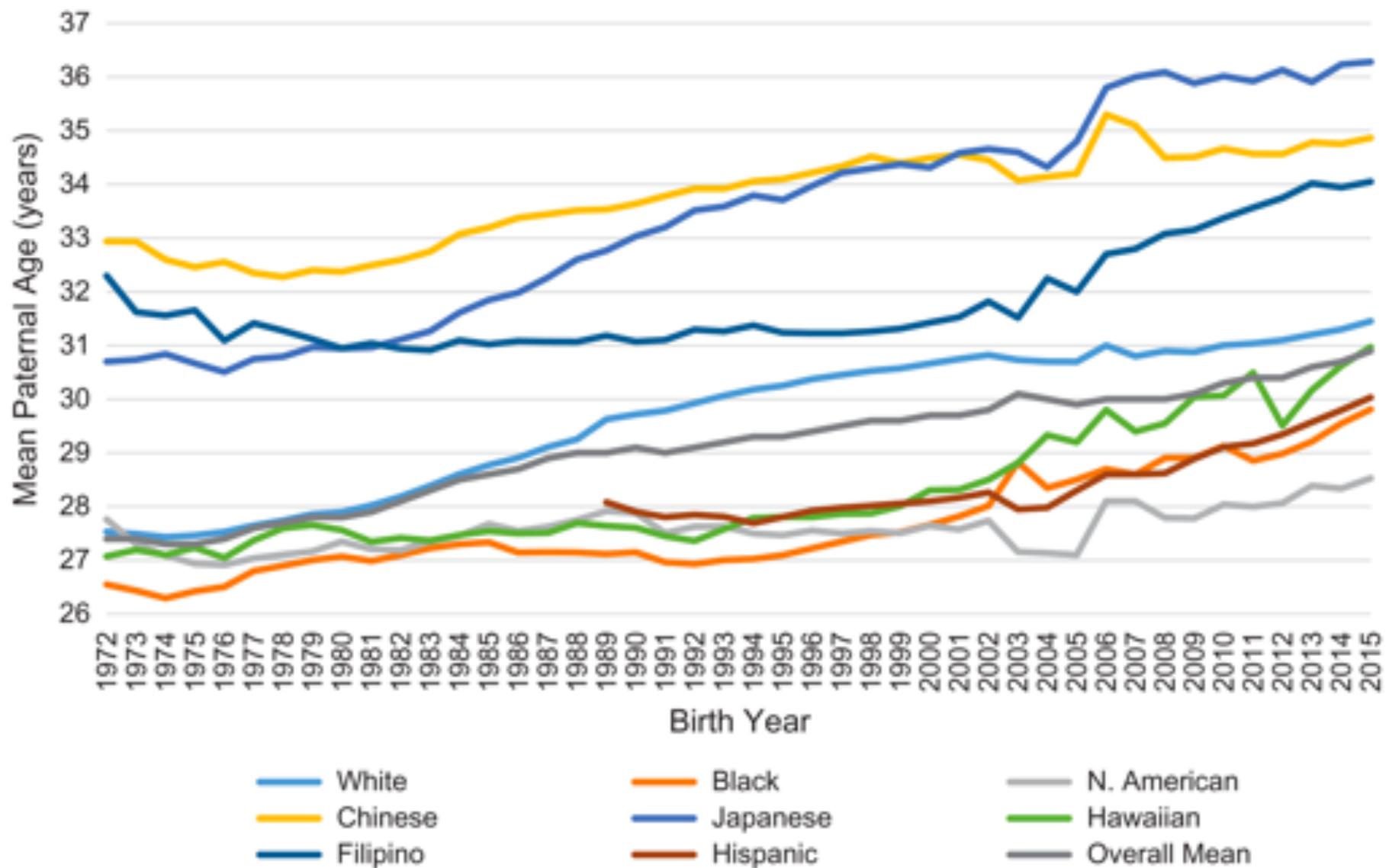
Narzisi et al (2014) Nature Methods doi:10.1038/nmeth.3069

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

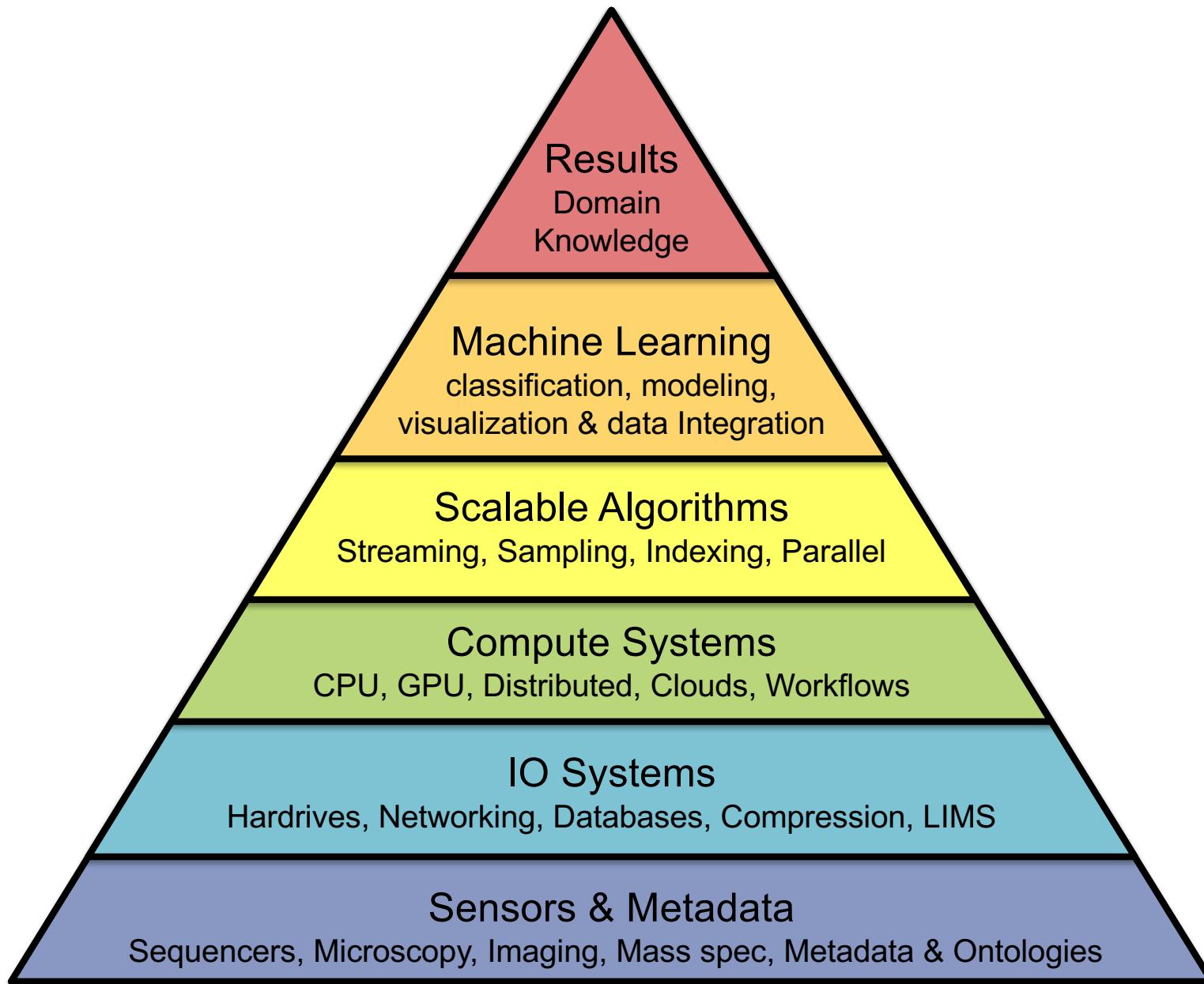
Age of Fatherhood



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. Work on Assignment I
 1. Set up Linux, set up Virtual Machine, set up Ubuntu
 2. Set up Dropbox for yourself!
 3. Get comfortable on the command line