

# Genome Sequencing

Michael Schatz

August 28, 2024

Lecture 2: Applied Comparative Genomics



# Course Webpage

The screenshot shows a GitHub repository page for 'appliedgenomics2024'. The repository has 1 branch and 0 tags. The README file contains the following content:

## JHU EN.601.449/EN.601.649: Computational Genomics: Applied Comparative Genomics

Prof: Michael Schatz ([mshatz@cs.jhu.edu](mailto:mshatz@cs.jhu.edu))  
TA: Matthew Nguyen ([mnguye99@jh.edu](mailto:mnguye99@jh.edu))  
Class Hours: Monday + Wednesday @ 3:00p - 4:15p Hodson 316  
Schatz Office Hours: By appointment  
Sweeten Office Hours: TBD and by appointment

The primary goal of the course is for students to be grounded in the fundamental theory and applications to 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. A major focus will be on deep learning and machine learning to tackle these problems. 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.

### Prerequisites

- Online introduction to Unix/Linux. Students are strongly recommended to complete one of the following online tutorials (or both) before class begins.
  - [Code academy's intro to Unix](#)
  - [Rosalind Bioinformatics Programming in Python](#)
  - [Minimal Make](#)

Access to a Linux Machine, and/or install Docker (Unfortunately, even Mac will not work correctly for some)

**About**

Materials for EN.601.449/649  
Computational Genomics: Applied Comparative Genomics

- Readme
- CC0-1.0 license
- Activity
- Custom properties
- 0 stars
- 1 watching
- 0 forks
- Report repository

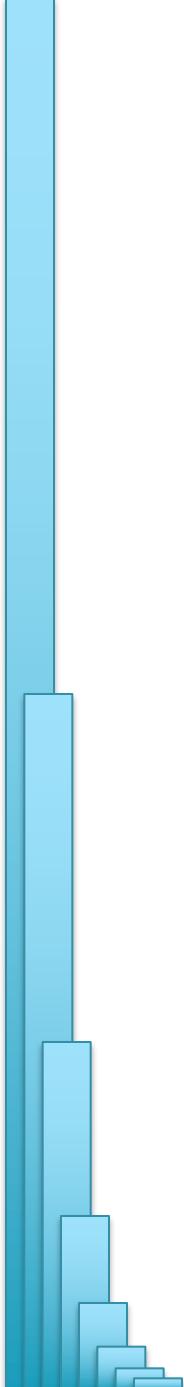
**Releases**

No releases published

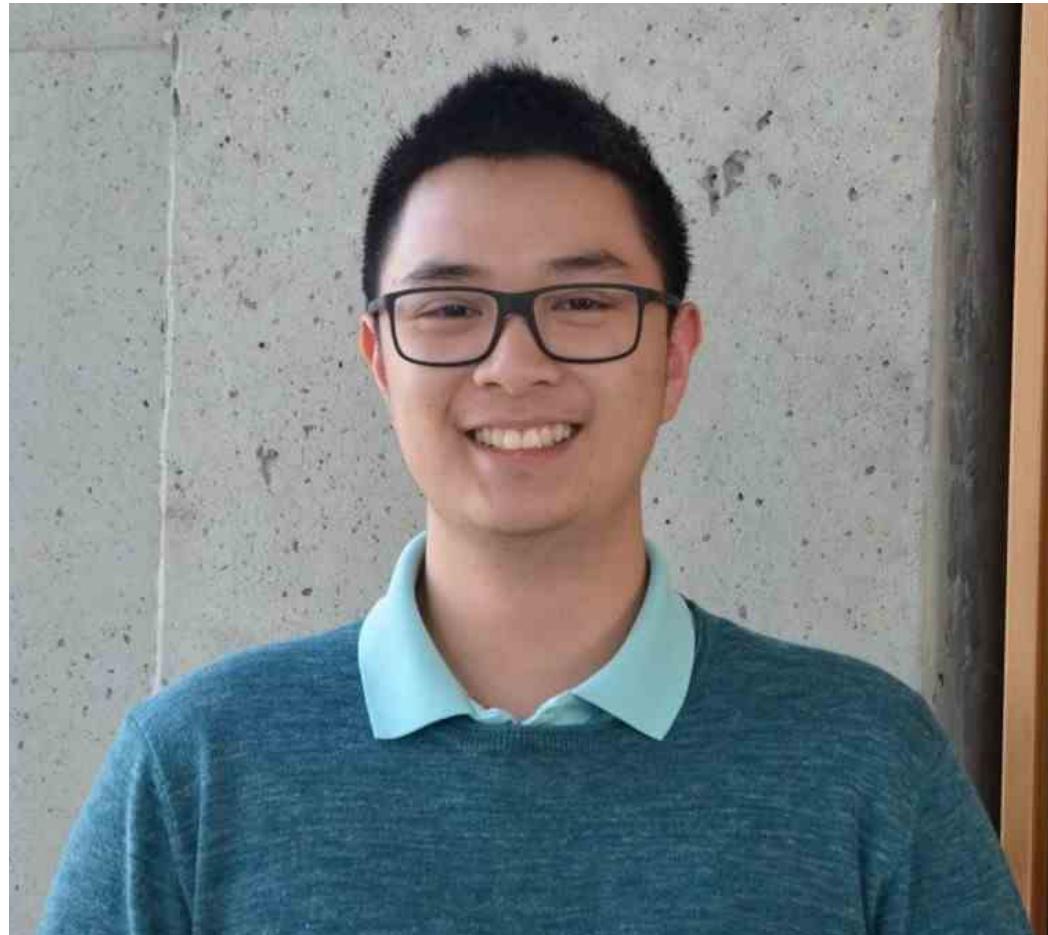
**Packages**

No packages published

<https://github.com/schatzlab/appliedgenomics2024>



# TA: Matthew Nguyen



Check Piazza for Poll

# Assignment 1

The screenshot shows a GitHub repository page for 'appliedgenomics2024' with the path 'assignments/assignment1'. The 'README.md' file is open, displaying the assignment details.

**Assignment 1: Chromosome Structures**

Assignment Date: Wednesday, August 28, 2024  
Due Date: Wednesday, Sept. 4, 2024 @ 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 [10 pts]**

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

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

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

- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes
- Question 1.3. Largest chromosome size and name
- Question 1.4. Smallest chromosome size and name
- Question 1.5. Mean chromosome length

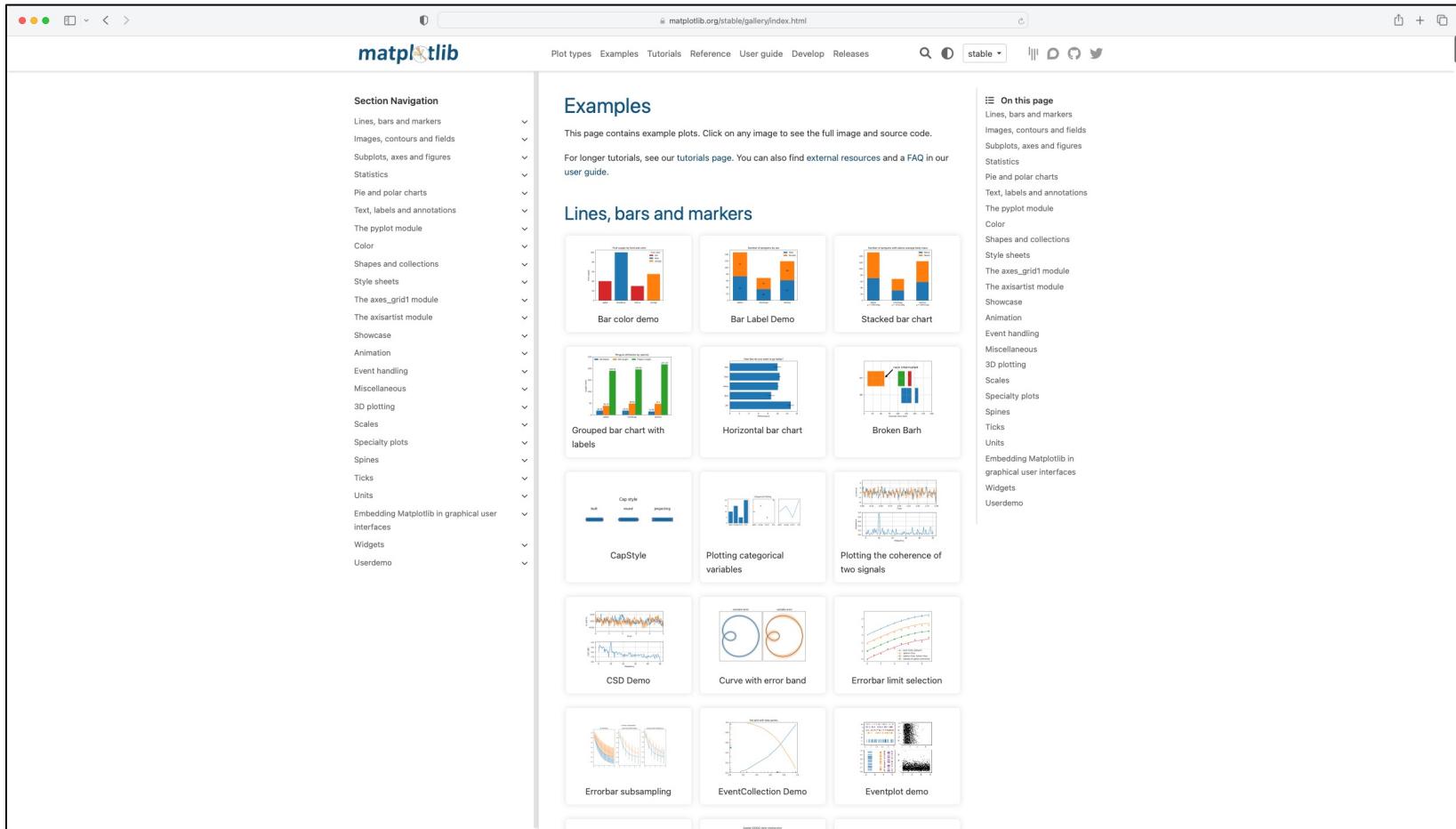
**Question 2. Coverage simulator [20 pts]**

- Question 2.1. How many 100bp reads are needed to sequence a 1Mbp genome to 3x coverage?
- Question 2.2. In the language of your choice, simulate sequencing 3x coverage of a 1Mbp genome with 100bp reads and plot the histogram of coverage. Note you do not need to actually output the sequences of the reads, you can just uniformly sample positions in the genome and record the coverage. You do not need to consider the strand of each read. The start position of each read should have a uniform random probability at each possible starting position (1 through 999,901). You can record the coverage in an

<https://github.com/schatzlab/appliedgenomics2024/tree/main/assignments/assignment1>

Due end of day on Sept 4 (right before midnight)

# Plotting in Python



The screenshot shows the official Matplotlib gallery website at [matplotlib.org/stable/gallery/index.html](https://matplotlib.org/stable/gallery/index.html). The page features a sidebar with 'Section Navigation' containing links to various plotting categories like 'Lines, bars and markers', 'Images, contours and fields', etc. The main content area is titled 'Examples' and displays a grid of thumbnail images for different plot types. The first row includes 'Bar color demo', 'Bar Label Demo', and 'Stacked bar chart'. The second row includes 'Grouped bar chart with labels', 'Horizontal bar chart', and 'Broken Barh'. The third row includes 'CapStyle', 'Plotting categorical variables', and 'Plotting the coherence of two signals'. The fourth row includes 'CSD Demo', 'Curve with error band', and 'Errorbar limit selection'. The fifth row includes 'Errorbar subsampling', 'EventCollection Demo', and 'Eventplot demo'. A sidebar on the right lists 'On this page' topics such as 'Lines, bars and markers', 'Images, contours and fields', 'Subplots, axes and figures', 'Statistics', 'Pie and polar charts', 'Text, labels and annotations', 'The pyplot module', 'Color', 'Shapes and collections', 'Style sheets', 'The axes\_grid1 module', 'The axisartist module', 'Showcase', 'Animation', 'Event handling', 'Miscellaneous', '3D plotting', 'Scales', 'Specialty plots', 'Spines', 'Ticks', 'Units', 'Embedding Matplotlib in graphical user interfaces', 'Widgets', and 'Userdemo'.

<https://matplotlib.org/>

# Plotting in R / ggplot2

# Data visualization with ggplot2 :: CHEATSHEET



## Basics

**ggplot2** is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.

To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

Complete the template below to build a graph.

```
ggplot(data = <DATA>) +  
  <GEOGRAPHICAL_FUNCTION>(<mapping> = aes(<COORDINATE_VARIABLES>),  
  stat = <STAT>, position = <POSITION>) +  
  <COORDINATE_FUNCTION> +  
  <FACET_FUNCTION> +  
  <SCALE_FUNCTION> +  
  <THEME_FUNCTION>
```

**ggplot**(data = mpg, aes(x = cyl, y = hwy)) Begins a plot that you finish by adding layers to. Add one geom per layer.

**last\_plot()** Returns the last plot.

**ggsave**("plot.png", width = 5, height = 5) Saves last plot as 5" x 5" file named "plot.png" in working directory. Matches file type to file extension.

## Geoms

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

### GRAPHICAL PRIMITIVES

```
a <- ggplot(economics, aes(date, unemploy))  
b <- ggplot(seals, aes(x = long, y = lat))
```

- a + geom\_blank()** and **a + expand\_limits()**  
Ensures limits include values across all plots.
- b + geom\_curve(aes(yend = lat + 1, xend = long + 1, curvature = 1))** - x, yend, y, yend, alpha, angle, color, curve, linetype, size
- a + geom\_path(linetype = "butt", linejoin = "round", linemethod = 1)** - x, y, alpha, color, group, linetype, size
- a + geom\_polygon(aes(alpha = 50))** - x, y, alpha, color, fill, group, subgroup, linetype, size
- b + geom\_rect(aes(xmin = long - ymin = lat, xmax = long + 1 - ymax + lat - 1), xmax, ymin, ymax, ymin, alpha, color, fill, linetype, size, weight)**
- a + geom\_ribbon(aes(ymin = unemploy - 900, ymax = unemploy + 900))** - x, ymax, ymin, alpha, color, fill, group, linetype, size

### TWO VARIABLES both continuous

```
e <- ggplot(mpg, aes(cty, hwy))
```

- e + geom\_label(aes(label = cty, nudge\_x = 1, nudge\_y = 1))** - x, y, label, alpha, angle, color, family, fontface, fill, hjust, lineheight, size, vjust
- e + geom\_point()** - x, y, alpha, color, fill, shape, size, stroke
- e + geom\_quantile()** - x, y, alpha, color, group, linetype, size, weight
- e + geom\_rug(sides = "bl")** - x, y, alpha, color, linetype, size
- e + geom\_smooth(method = lm)** - x, y, alpha, color, fill, group, linetype, size, weight
- e + geom\_text(aes(label = cty, nudge\_x = 1, nudge\_y = 1))** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

### continuous bivariate distribution

```
h <- ggplot(diamonds, aes(carat, price))
```

- h + geom\_bin2d(binwidth = c(0.25, 500))** - x, y, alpha, color, fill, linetype, size, weight
- h + geom\_density\_2d()** - x, y, alpha, color, group, linetype, size
- h + geom\_hex()** - x, y, alpha, color, fill, size

### continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

- i + geom\_area()** - x, y, alpha, color, fill, linetype, size
- i + geom\_line()** - x, y, alpha, color, group, linetype, size
- i + geom\_step(direction = "hv")** - x, y, alpha, color, group, linetype, size

### mapping error

```
df <- data.frame(prt = c("A", "B"), fit = 4.5, se = 1.2)  
j <- ggplot(df, aes(prt, fit, ymin = fit - se, ymax = fit + se))
```

- j + geom\_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size
- j + geom\_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width
- Also geom\_errorbar!().**
- j + geom\_linerange()** - x, ymin, ymax, alpha, color, group, linetype, size
- j + geom\_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size, weight

### one discrete, one continuous

```
f <- ggplot(mpg, aes(class, hwy))
```

- f + geom\_col()** - x, y, alpha, color, fill, group, linetype, size
- f + geom\_boxplot()** - x, y, lower, middle, upper, ymax, ymin, alpha, color, fill, group, linetype, shape, size, weight
- f + geom\_dotplot(binaxis = "y", stackdir = "center")** - x, y, alpha, color, fill, group
- f + geom\_violin(scale = "area")** - x, y, alpha, color, fill, group, linetype, size, weight

### one VARIABLE continuous

```
c <- ggplot(mpg, aes(hwy)); c <- ggplot(mpg)
```

- c + geom\_area(stat = "bin")** - x, y, alpha, color, fill, linetype, size
- c + geom\_density(kernel = "gaussian")** - x, y, alpha, color, fill, group, linetype, size, weight
- c + geom\_dotplot()** - x, y, alpha, color, fill
- c + geom\_freqpoly()** - x, y, alpha, color, group, linetype, size
- c + geom\_histogram(binwidth = 5)** - x, y, alpha, color, fill, linetype, size, weight
- c2 <- geom\_qq(aes(sample = hwy))** - x, y, alpha, color, fill, linetype, size, weight

### both discrete

```
g <- ggplot(diamonds, aes(cut, color))
```

- g + geom\_count()** - x, y, alpha, color, fill, shape, size, stroke
- g + geom\_jitter(height = 2, width = 2)** - x, y, alpha, color, fill, shape, size

### maps

```
data <- data.frame(murder = USArrests$Murder, state = tolower(rownames(USArrests)))  
map <- map_data("state")  
k <- ggplot(data, aes(fill = murder))
```

- k + geom\_map(aes(map\_id = state), map = map) + expand\_limits(x = map\$long, y = map\$lat)** - map\_id, alpha, color, fill, linetype, size

### discrete

```
d <- ggplot(mpg, aes(fill))
```

- d + geom\_bar()** - x, alpha, color, fill, linetype, size, weight

### THREE VARIABLES

```
sealsSz <- with(seals, sqrt(delta_long^2 + delta_lat^2)); I <- ggplot(seals, aes(long, lat))
```

- I + geom\_contour(aes(z = z))** - x, y, z, alpha, color, group, linetype, size, weight
- I + geom\_contour\_filled(aes(z = z))** - x, y, alpha, color, fill, group, linetype, size, subgroup
- I + geom\_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)** - x, y, alpha, fill
- I + geom\_tile(aes(fill = z))** - x, y, alpha, color, fill, linetype, size, width

CC BY SA Posit Software, PBC • info@posit.co • posit.co • Learn more at [ggplot2.tidyverse.org](http://ggplot2.tidyverse.org) • HTML cheatsheets at [pos.it/cheatsheets](http://pos.it/cheatsheets) • ggplot2 3.4.2 • Updated: 2023-07

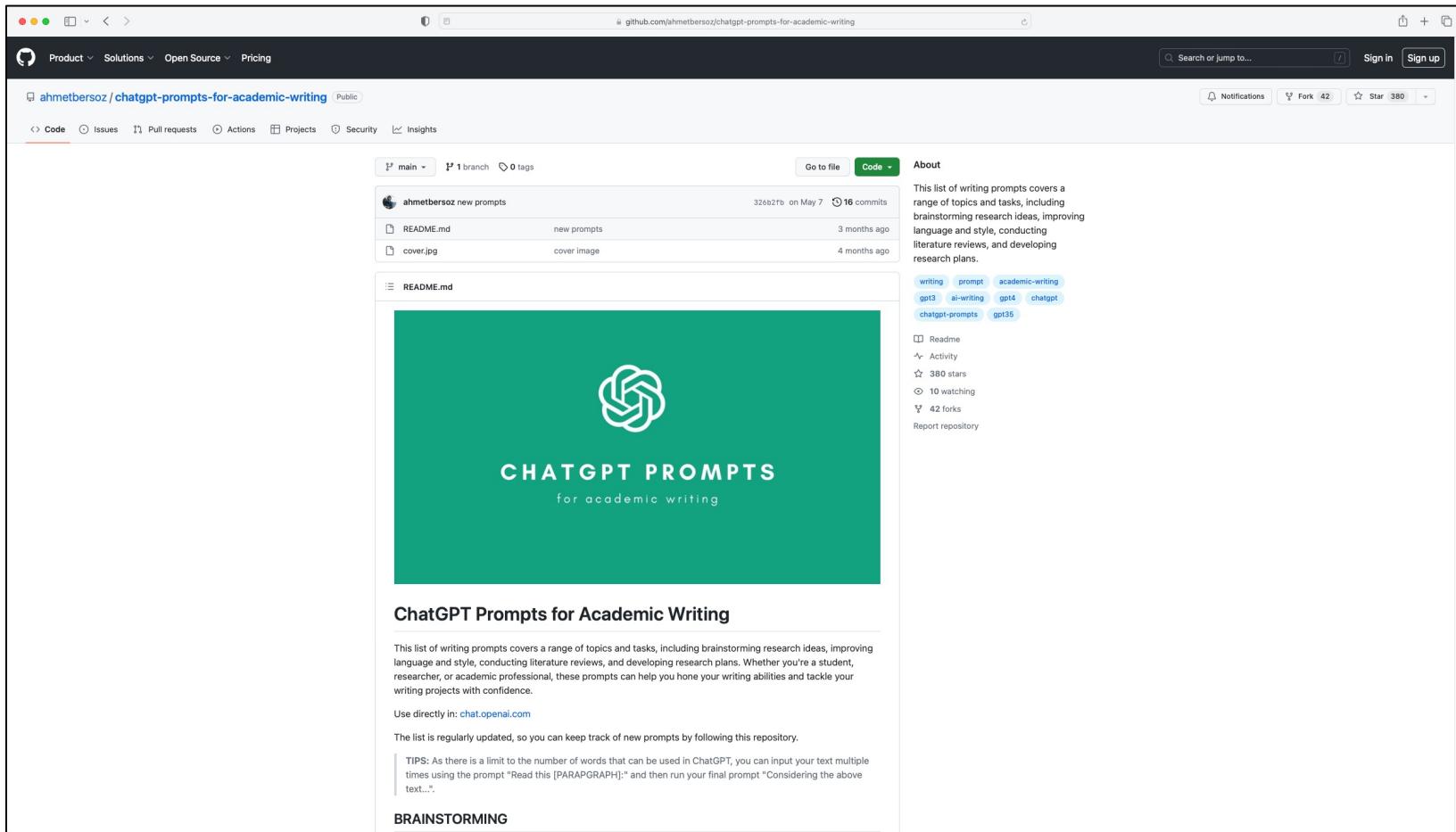
<https://ggplot2.tidyverse.org/>

# What is ChatGPT and Why Does it Work?

The screenshot shows a web browser displaying a blog post on Stephen Wolfram's website. The page has a yellow header bar with the title "STEPHEN WOLFRAM | Writings". Below the header, there are links for "RECENT" and "CATEGORIES" and a search bar. The main content area features a book cover for "What Is ChatGPT Doing ... and Why Does It Work?", which is now available in paperback and Kindle. A button labeled "Order Now" is visible. Below the book cover, there are "See also:" links to "Wolfram|Alpha as the Way to Bring Computational Knowledge Superpowers to ChatGPT" and "A discussion about the history of neural nets". The main article title is "What Is ChatGPT Doing ... and Why Does It Work?", dated February 14, 2023. The article begins with a section titled "It's Just Adding One Word at a Time" and discusses how ChatGPT generates text. To the right of the article, there is a sidebar titled "Recent Writings" featuring several thumbnail images and titles of other posts, such as "Remembering the Improbable Life of Ed Fredkin (1934–2023) and His World of Ideas and Stories" and "Introducing Chat Notebooks: Integrating LLMs into the Notebook Paradigm". There is also a "Popular Categories" sidebar with a list of topics like Artificial Intelligence, Big Picture, Companies and Business, etc.

<https://writings.stephenwolfram.com/2023/02/what-is-chatgpt-doing-and-why-does-it-work/>

# ChatGPT Prompts for Academic Writing



<https://github.com/ahmetbersoz/chatgpt-prompts-for-academic-writing>

# Discovery of the Double Helix

No. 4356 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.

<sup>1</sup>Yodkin, F. B., Germar, H., and Jeunet, W., *Phil. Mag.*, **40**, 149 (1929).

<sup>2</sup>Lionnet-Higgins, M. S., *Mon. Not. Roy. Astr. Soc., Geophys. Suppl.*, **5**, 285 (1949).

<sup>3</sup>Van Arkel, A. S., Woods Hole Papers in Phys. Oceanogr. Meteor., **11**, 131 (1948).

<sup>4</sup>Ekman, V. W., *Arktis Mat. Astron. Fysik* (Stockholm), **2** (11) (1905).

## MOLECULAR STRUCTURE OF NUCLEIC ACIDS

### A Structure for Deoxyribose Nucleic Acid

WE wish to suggest a structure for the salt  $\text{W}^-$  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 Corey<sup>1</sup>. 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 salt which gives the X-ray diagrams is the acid, not the base salt. While the acidic hydrogen atoms do not overlap, the 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 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 pairs of ester groups joining  $\beta$ -deoxyribofuranose residues with 3',5' 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 Furberg's<sup>2</sup> 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 Furberg's 'standard configuration', the sugar being roughly perpendicular to the attached base. There

This figure is purely diagrammatic. The two ribbons symbolize the two chains, and the horizontal rods the pairs of bases held together by hydrogen bonds. The vertical line marks the fibre axis.

is a residue on each chain every 3-4 Å. in the z-direction. We have assumed an angle of  $36^\circ$  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 10 Å. As the phosphates are on the outside, cations have easy access to them.

The fibre axis 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 z-coordinates. 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 1; purine position 6 to pyrimidine position 6.

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 be formed, it follows that if the sequence of bases on one chain is given, then the sequence on the other chain is automatically determined.

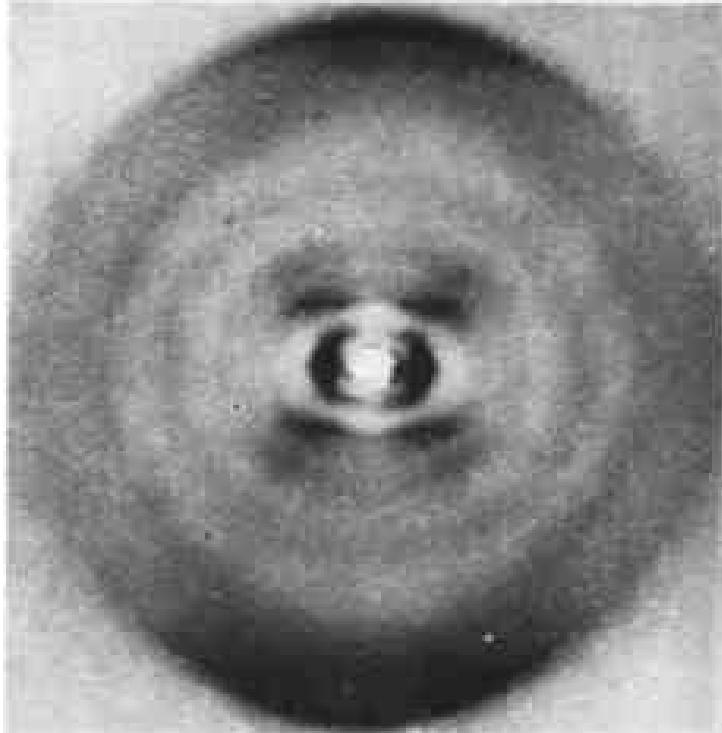
It has been found experimentally<sup>3,4</sup> that the ratio of the amounts of adenine to thymine, and the ratio of guanine to cytosine, are always very close to unity for deoxyribose nucleic acid.

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<sup>5,6</sup> on deoxyribose nucleic acid are insufficient for a rigorous test of our structure. The first of these is completely 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.

Full details of the structure, including the con-



### ACKNOWLEDGMENTS

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-

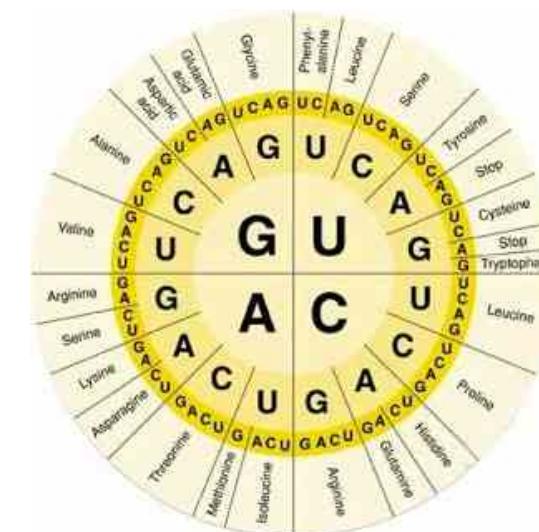
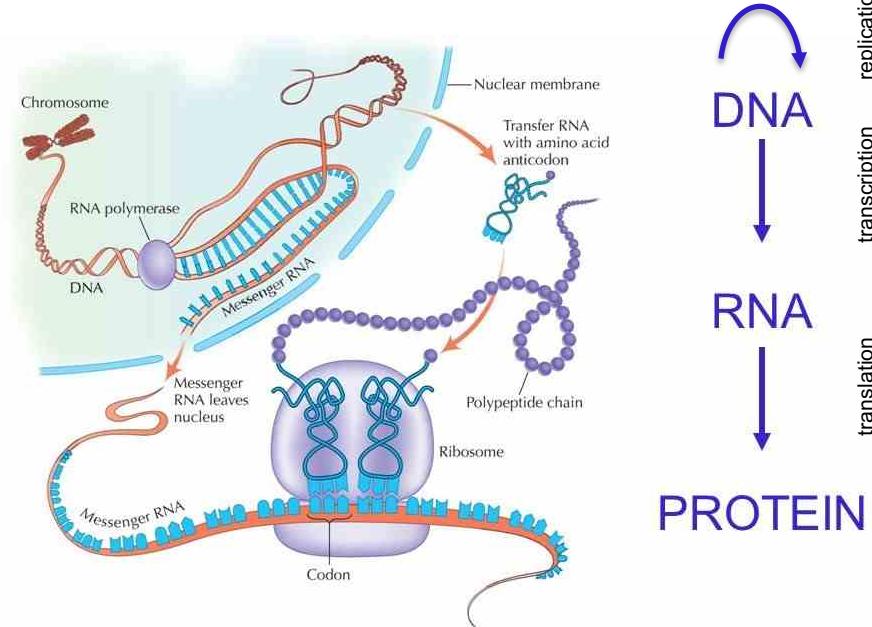
**Molecular structure of nucleic acids; a structure for deoxyribose nucleic acid**

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”

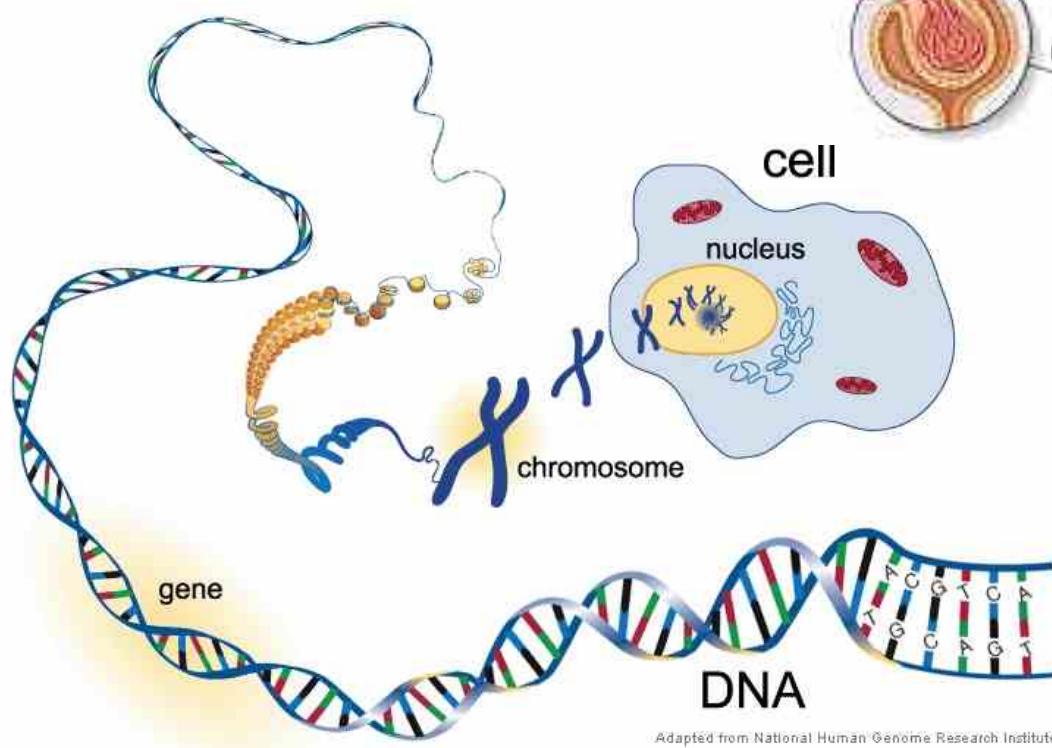
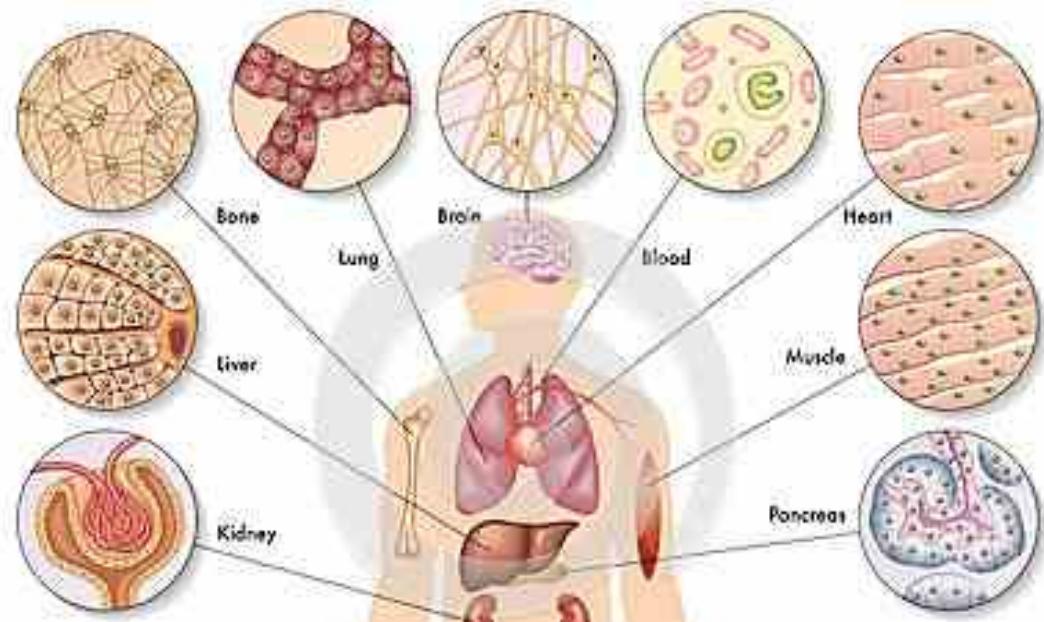


# **On Protein Synthesis**

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 a nearly 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

# Unsolved Questions in Biology

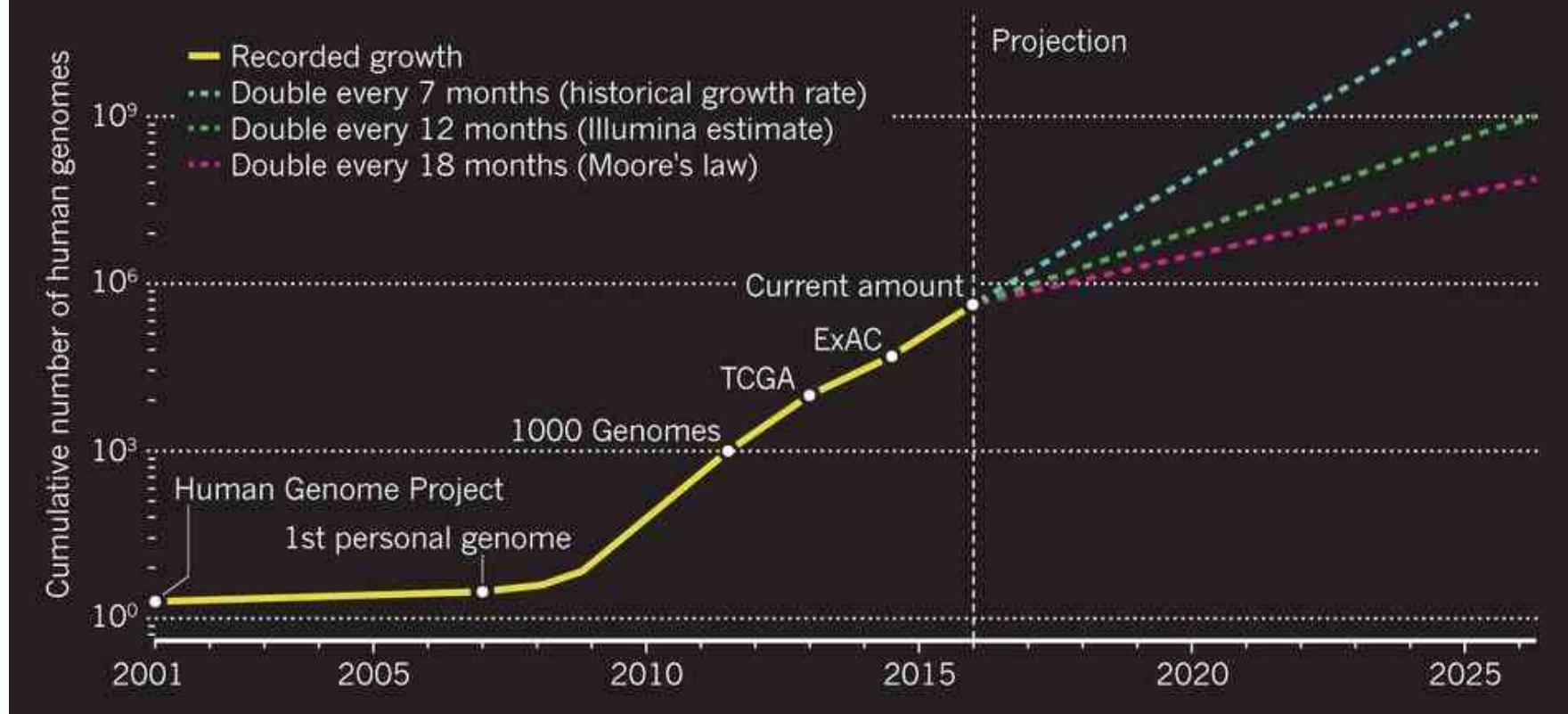
- What is your genome sequence?
- How does your genome compare to my genome?
- Where are the genes and how active are they?
- How does gene activity change during development?
- How does splicing change during development?
- How does methylation change during development?
- How does chromatin change during development?
- How does your genome folded in the cell?
- Where do proteins bind and regulate genes?
- What virus and microbes are living inside you?
- How do your mutations relate to disease?
- What drugs and treatments should we give you?
- ***Plus thousands and thousands more***



# Sequencing Capacity

## DNA SEQUENCING SOARS

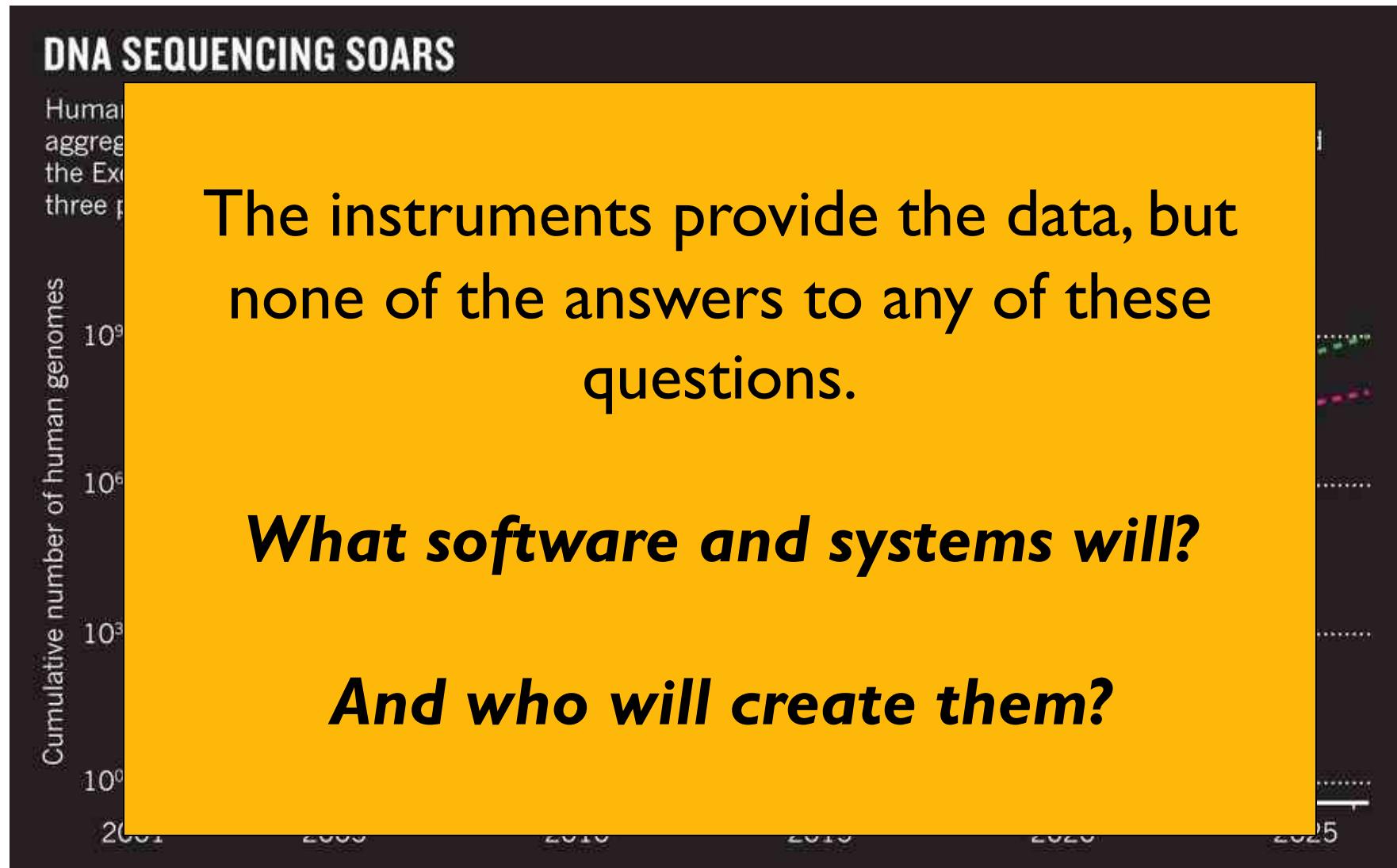
Human genomes are being sequenced at an ever-increasing rate. The 1000 Genomes Project has aggregated hundreds of genomes; The Cancer Genome Atlas (TCGA) has gathered several thousand; and the Exome Aggregation Consortium (ExAC) has sequenced more than 60,000 exomes. Dotted lines show three possible 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)

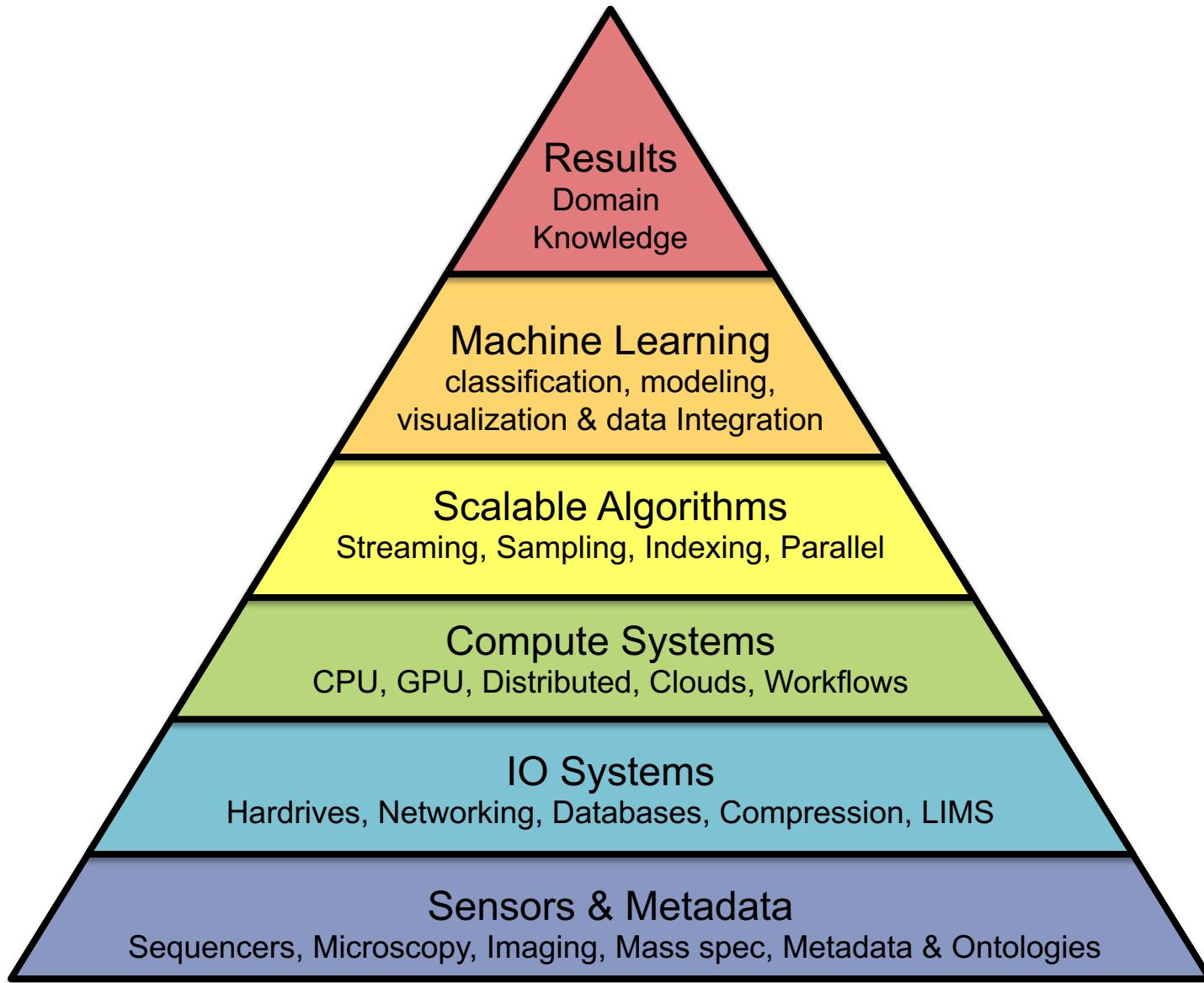
# Sequencing Capacity



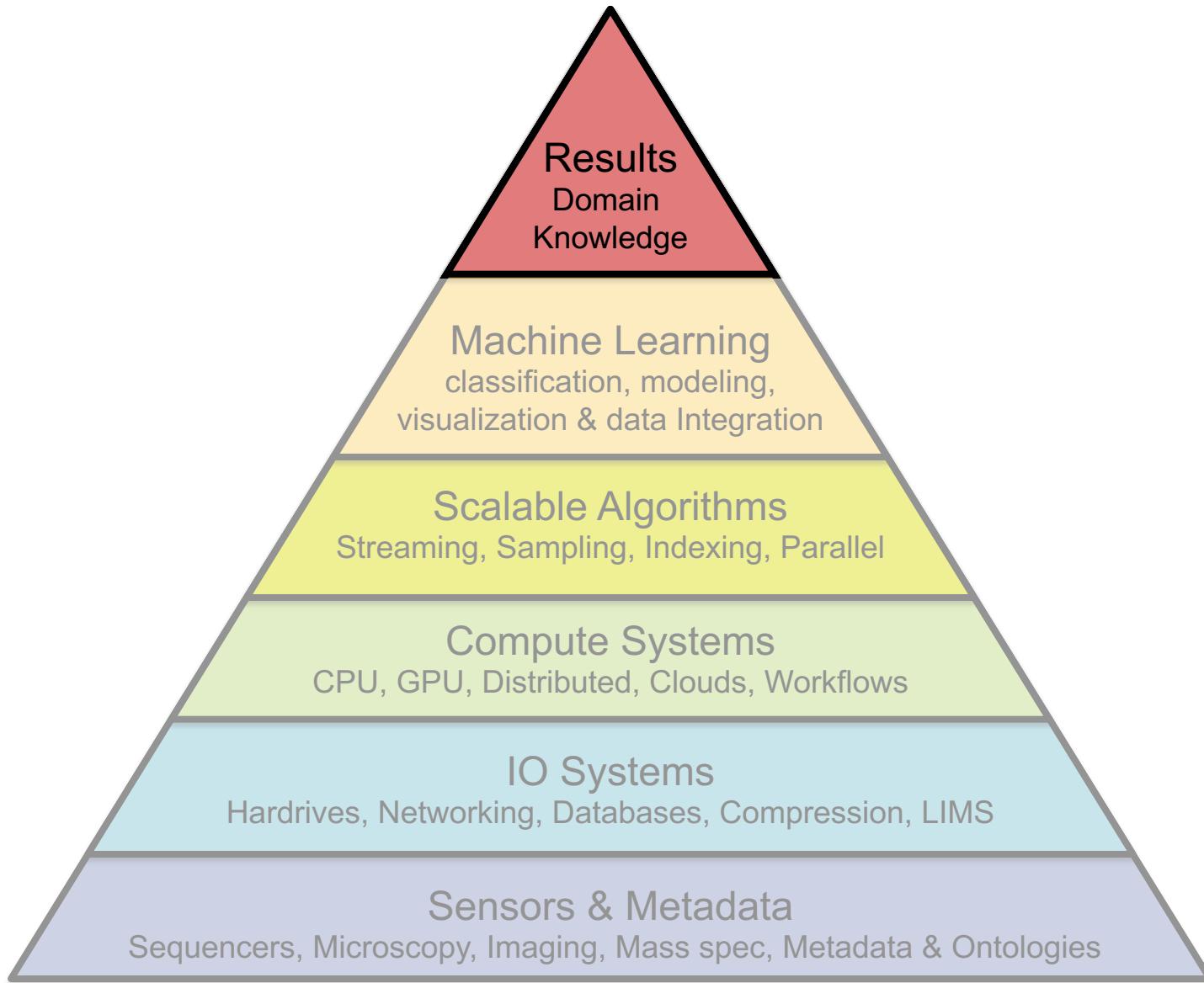
**Big Data: Astronomical or Genomical?**

Stephens, Z, et al. (2015) PLOS Biology DOI: 10.1371/journal.pbio.1002195

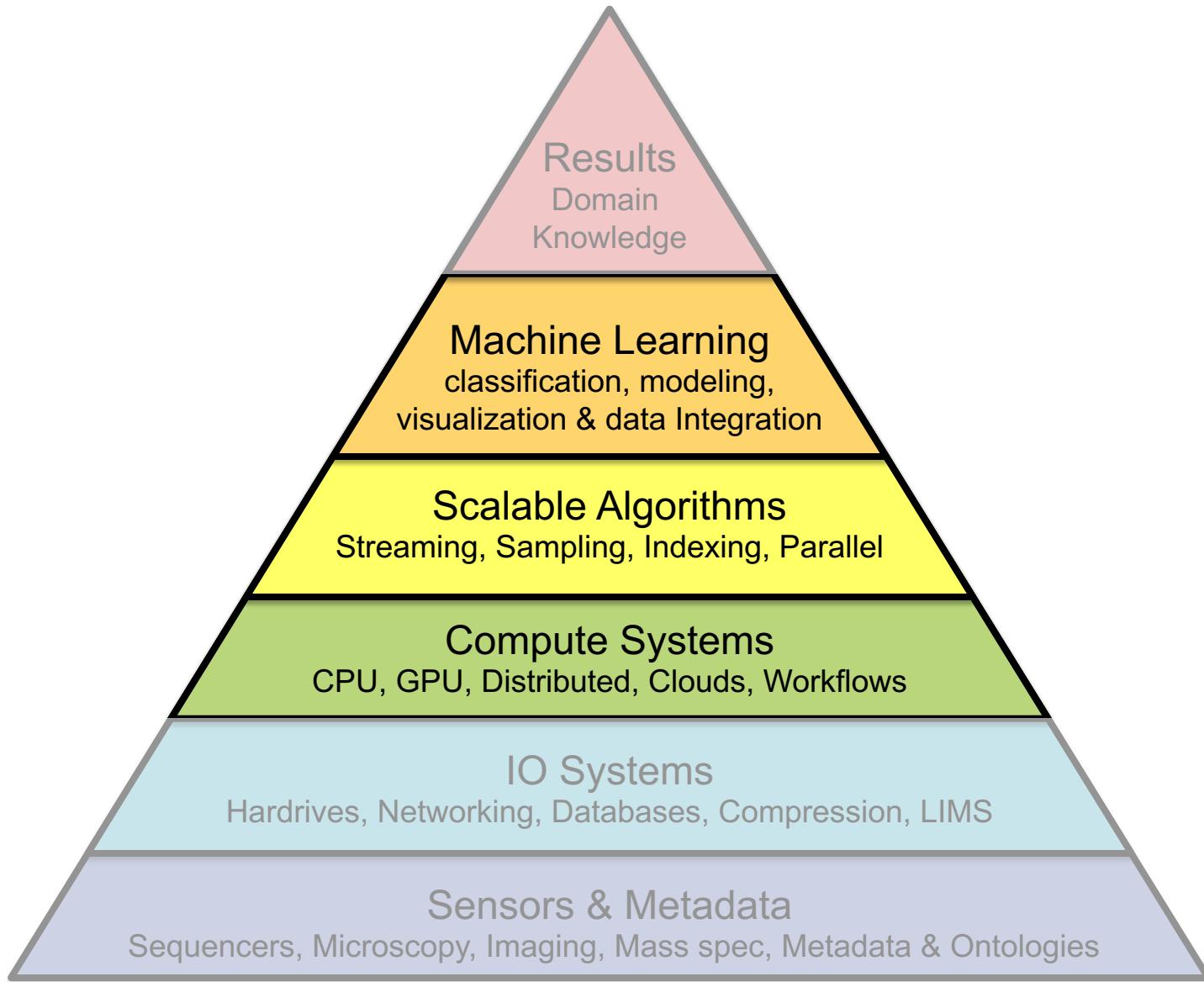
# Comparative Genomics Technologies



# Comparative Genomics Technologies



# Comparative Genomics Technologies

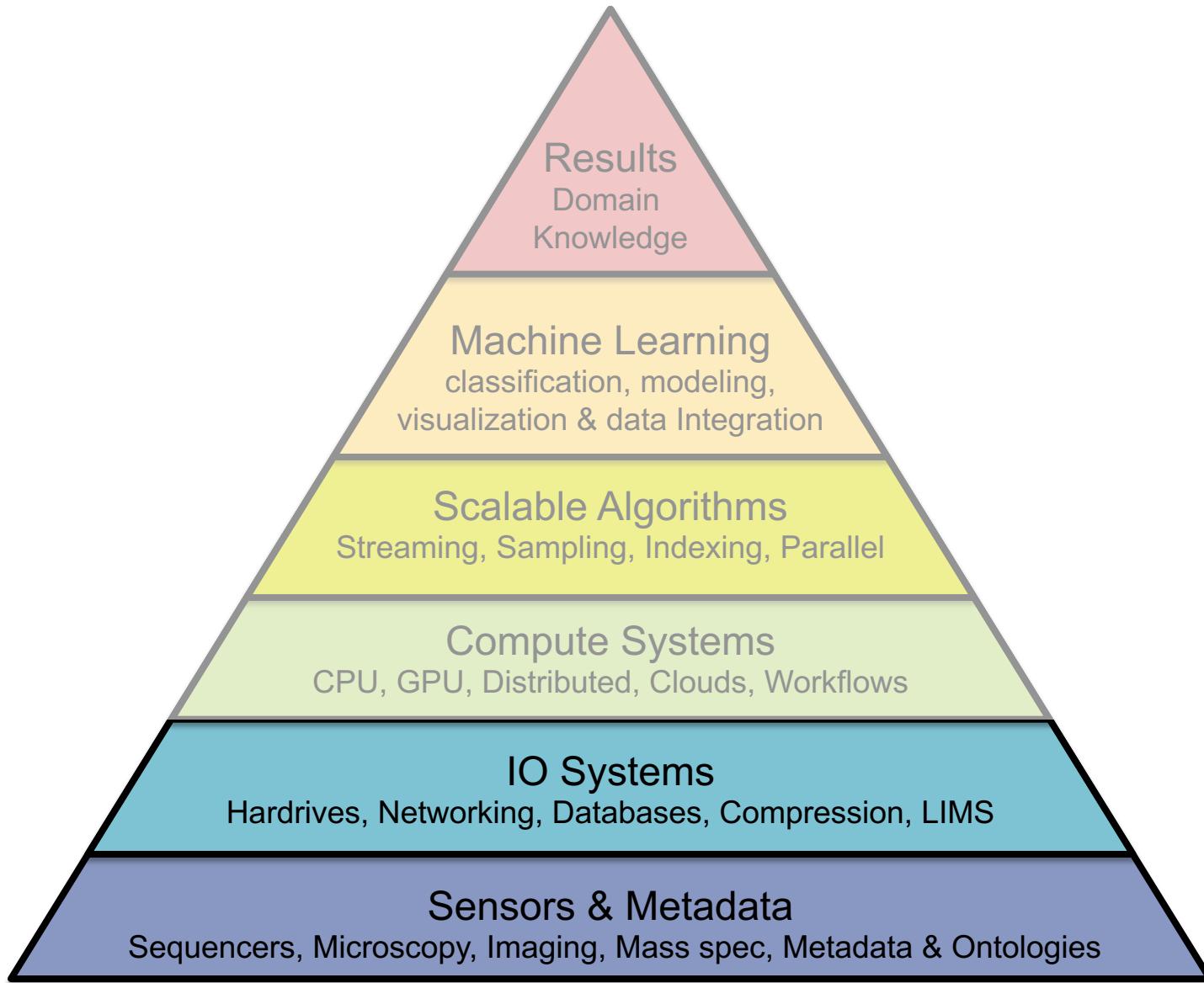


# Selected Topics

- 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



# Genomics Arsenal in the year 2024

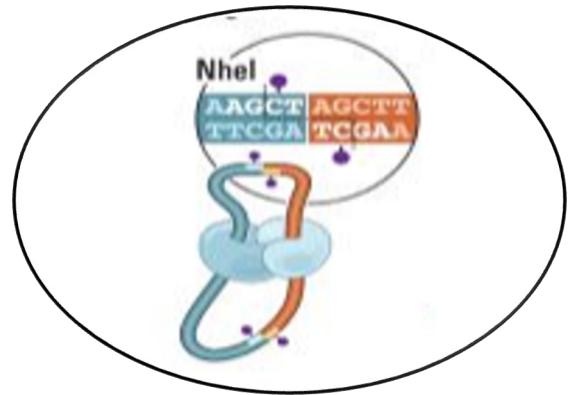
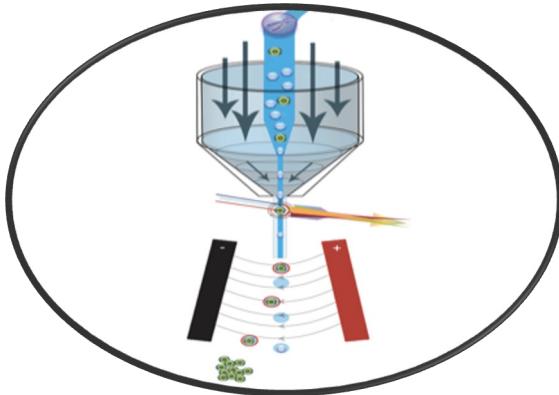
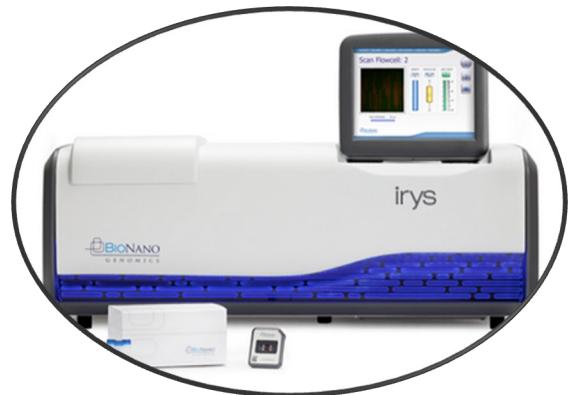
Sample Preparation

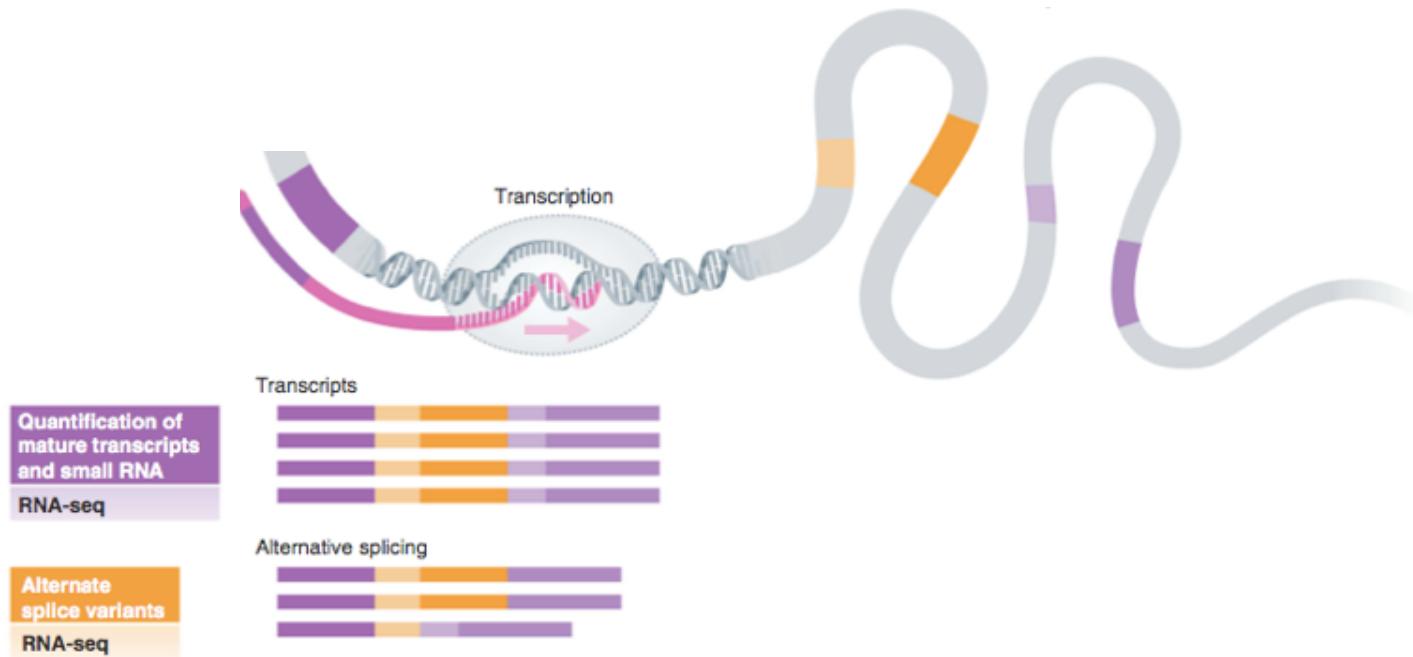


Sequencing

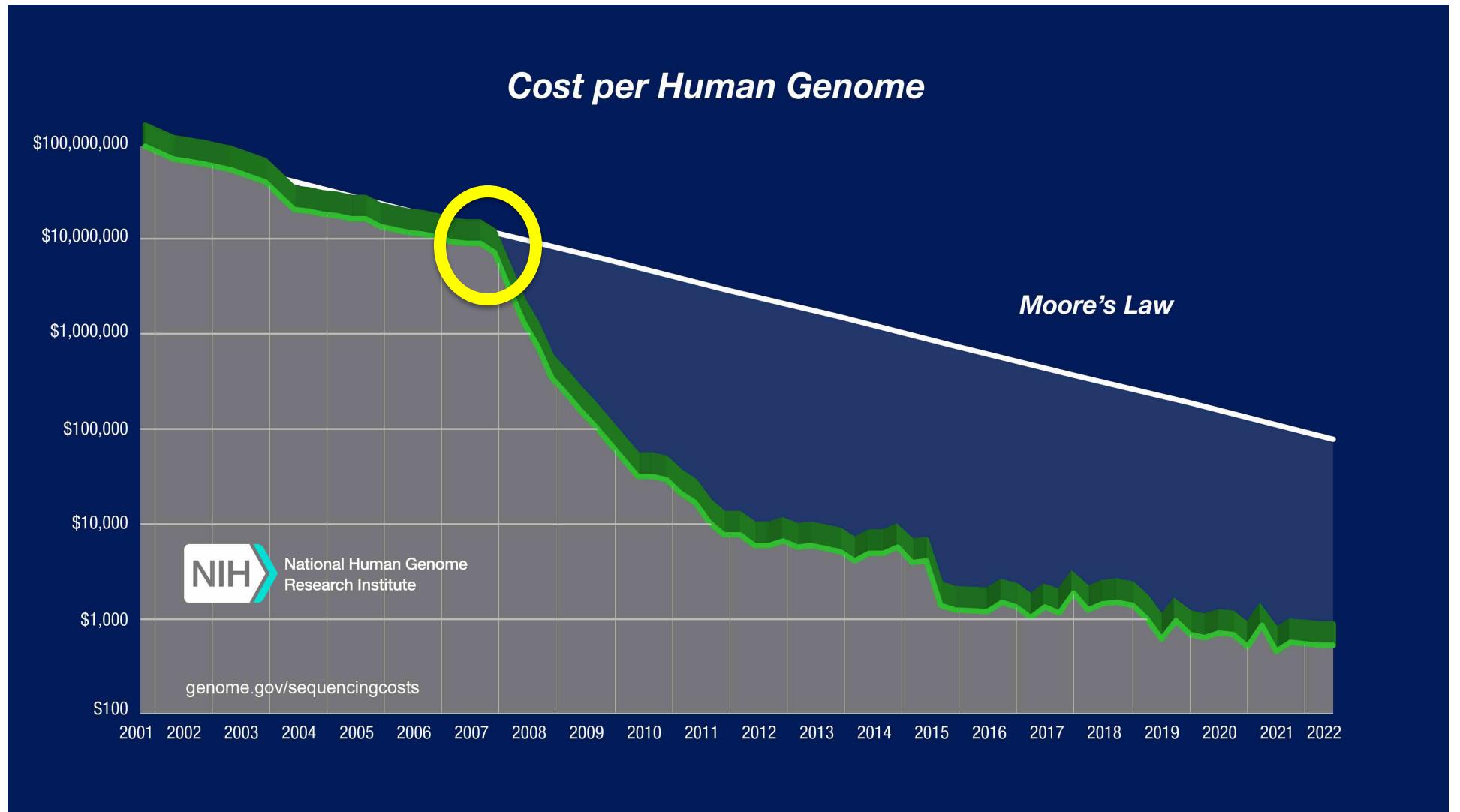


Chromosome Mapping





# Cost per Genome



<https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost>

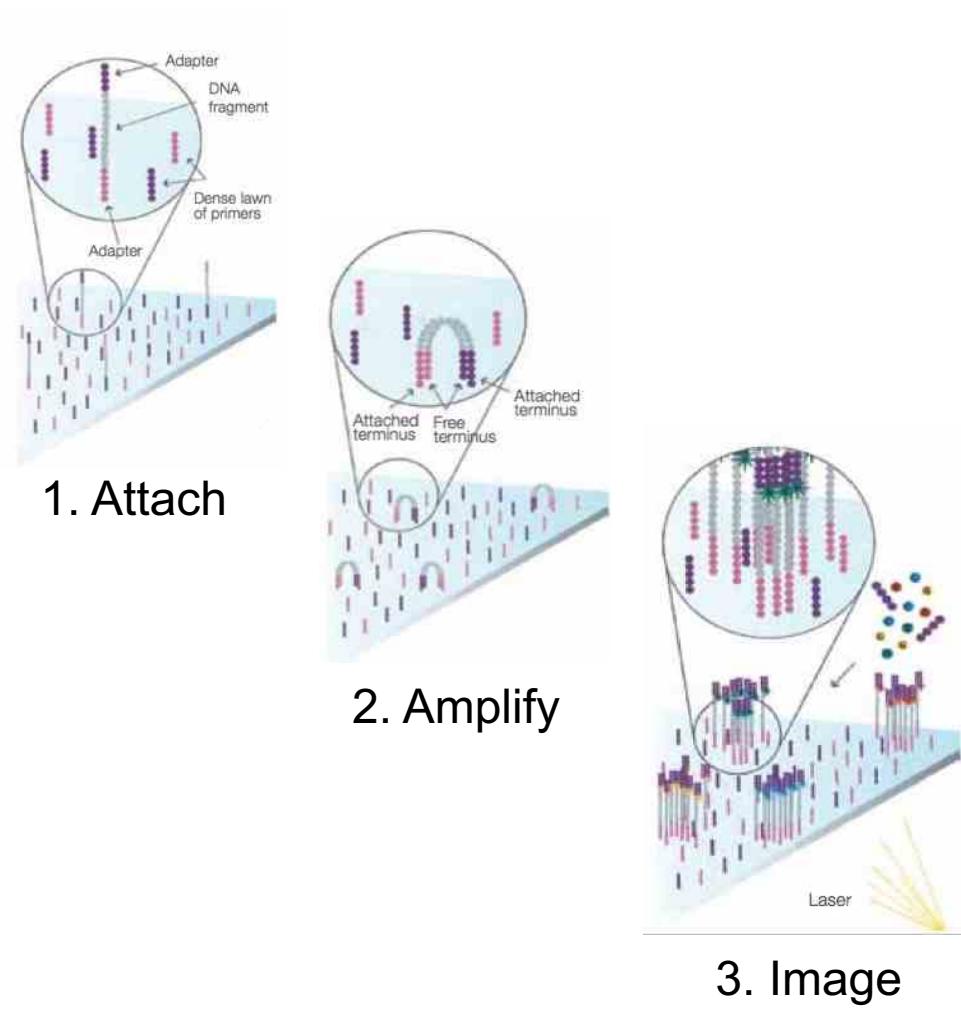
# Second Generation Sequencing



**Illumina NovaSeq 6000**

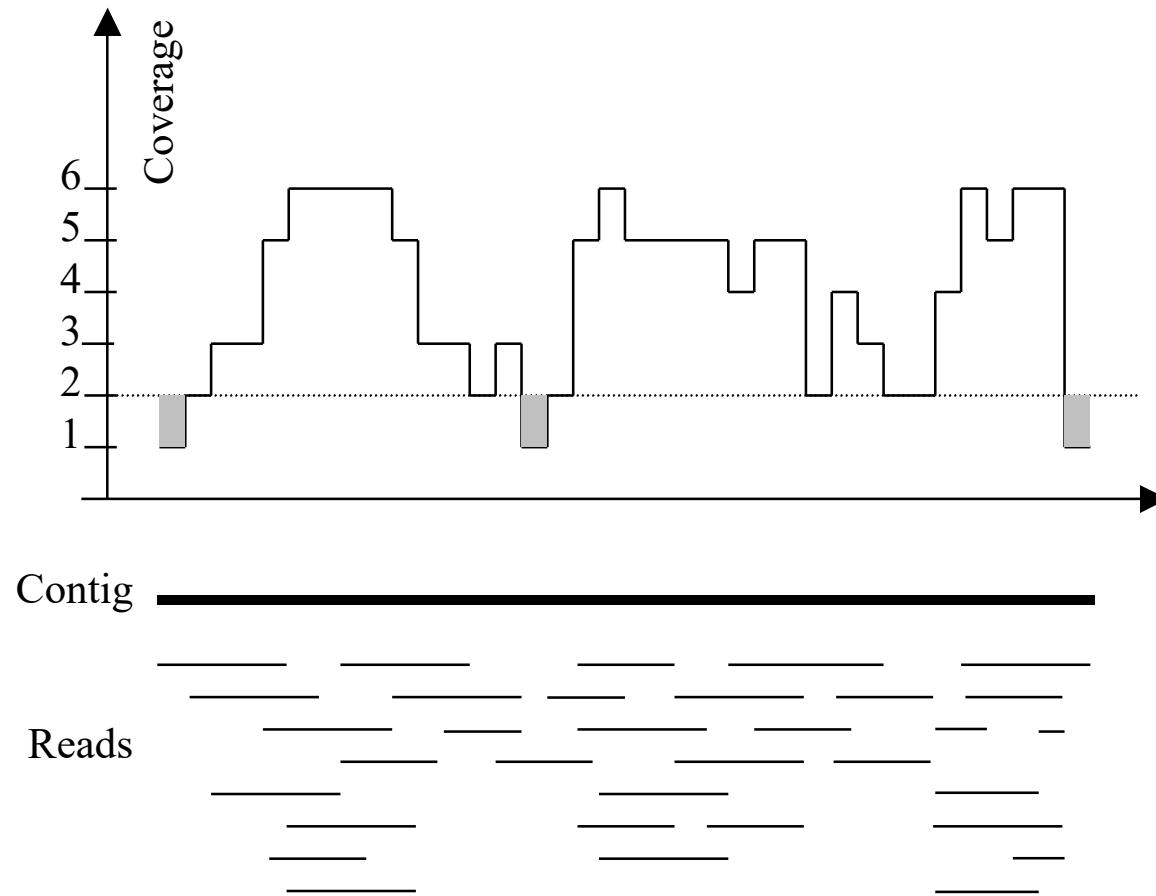
*Sequencing by Synthesis*

>3Tbp / day  
(JHU has 4 of these!)



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

# Typical sequencing coverage

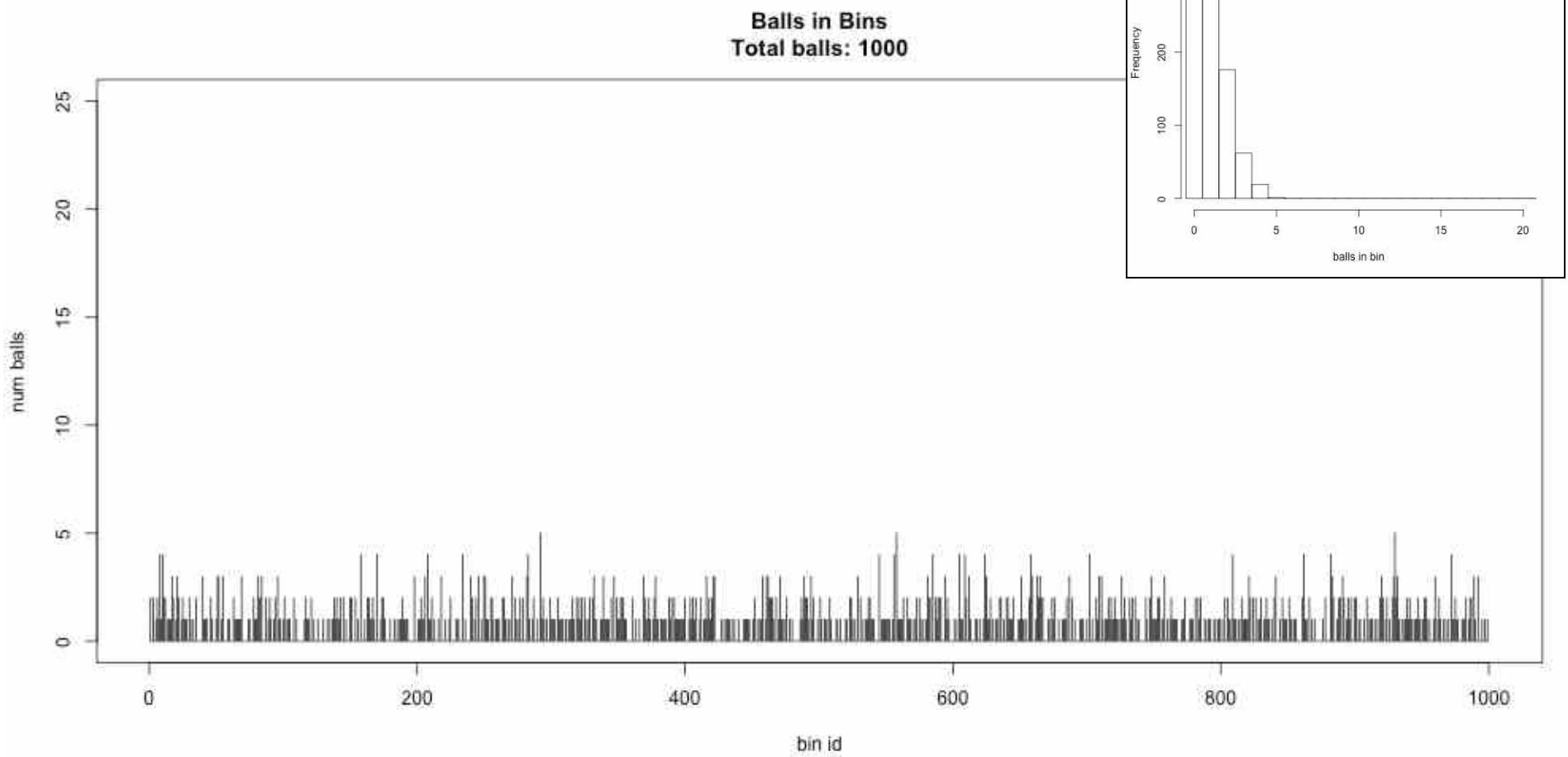


Imagine raindrops on a sidewalk

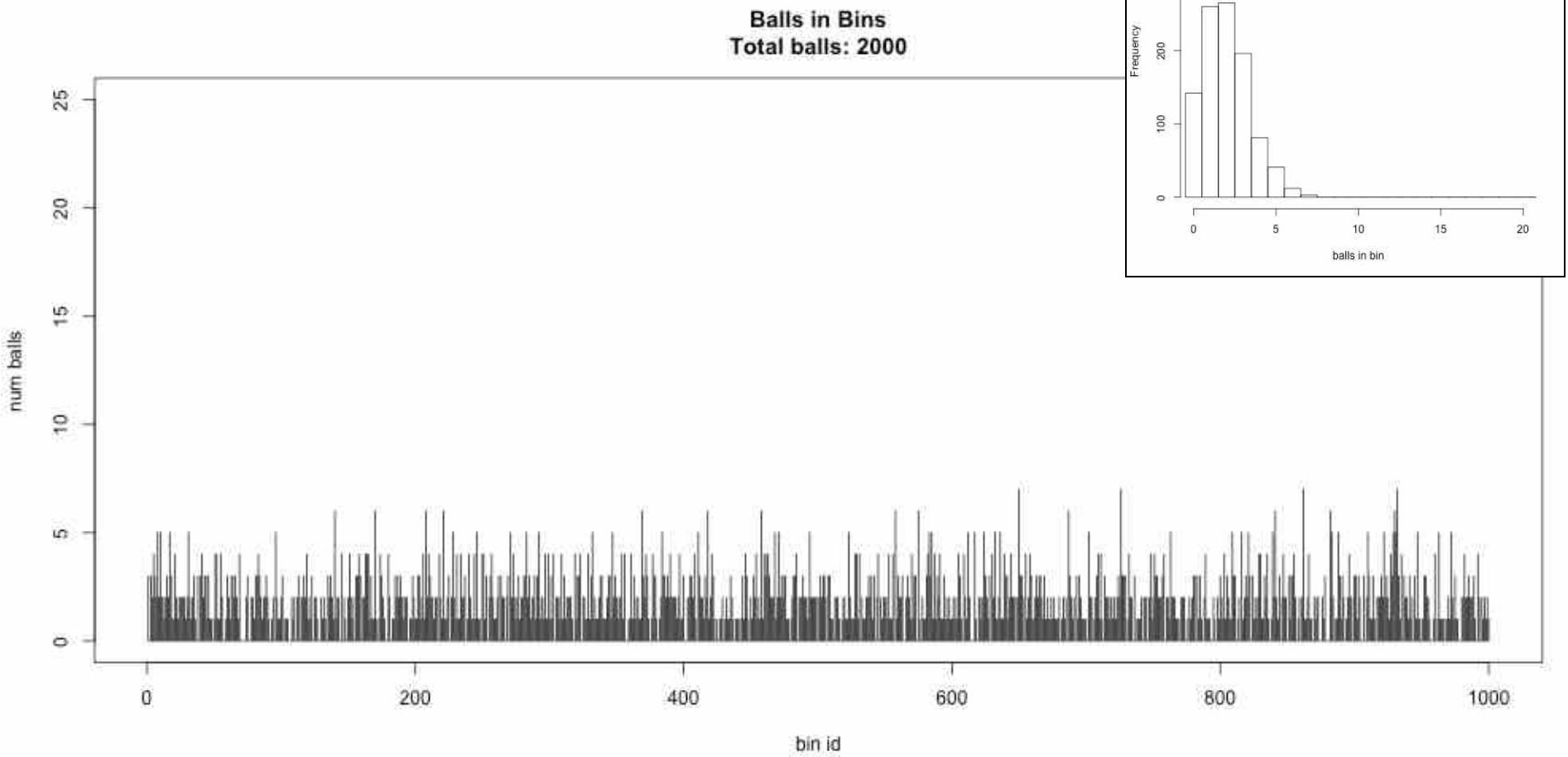
We want to cover the entire sidewalk but each drop costs \$1

If the genome is 10 Mbp, should we sequence 100k 100bp reads?

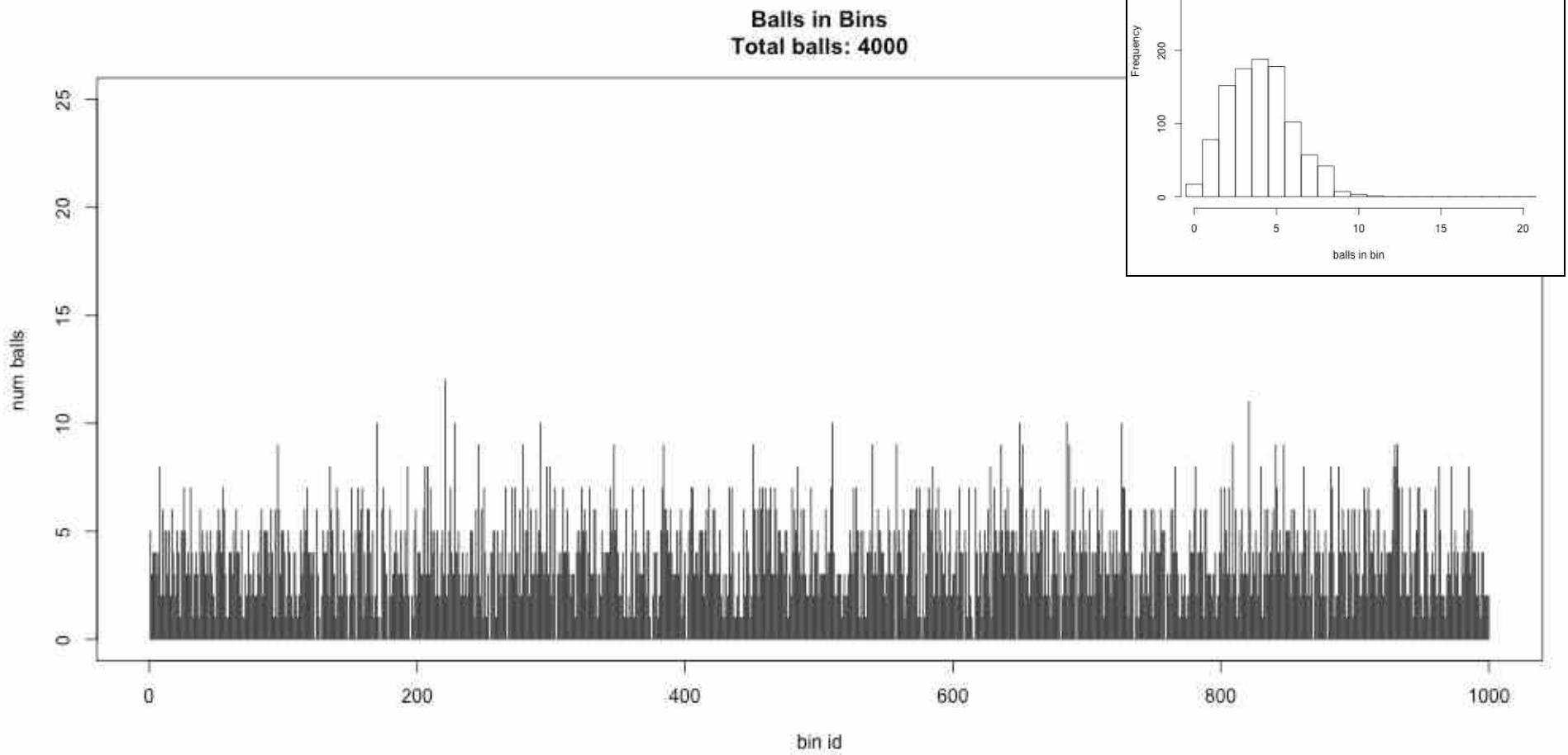
# Ix sequencing



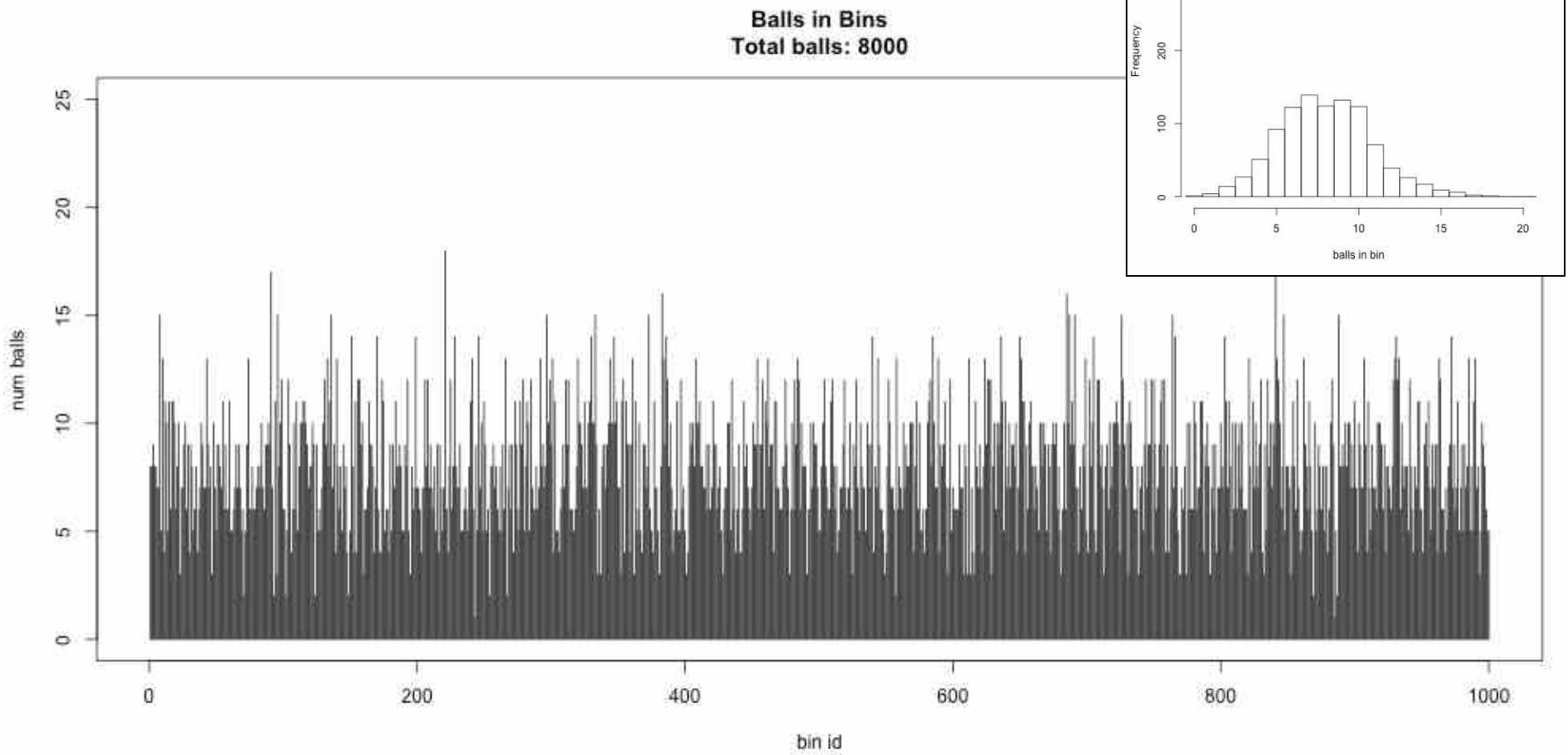
# 2x sequencing



# 4x sequencing



# 8x sequencing



# Poisson Distribution

The probability of a given number of events occurring in a fixed interval of time and/or space if these events occur with a known average rate and independently of the time since the last event.

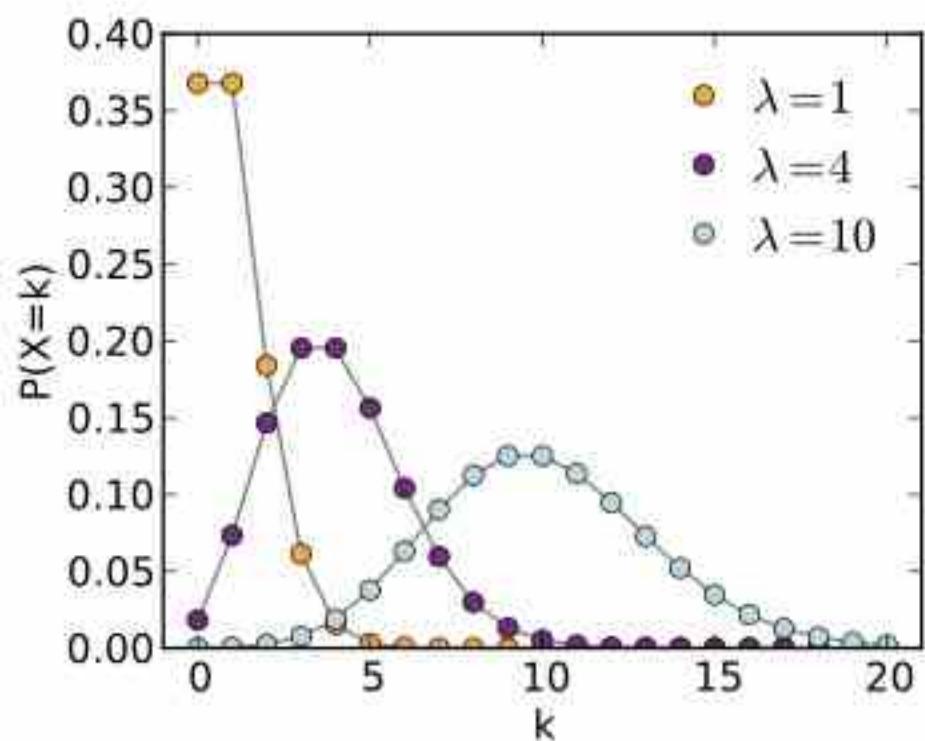
Formulation comes from the limit of the binomial equation

Resembles a normal distribution, but over the positive values, and with only a single parameter.

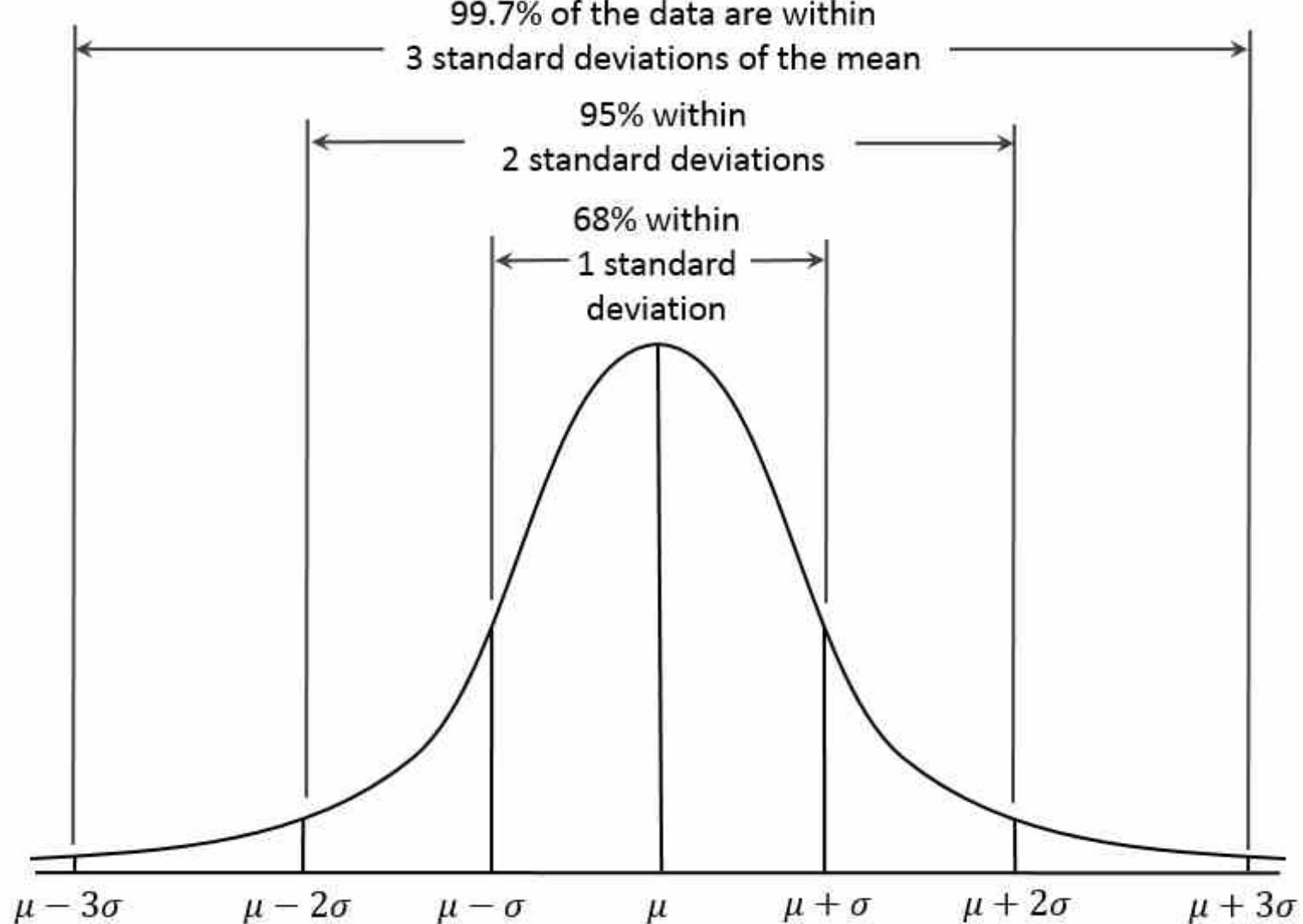
## ***Key properties:***

- ***The standard deviation is the square root of the mean.***
- ***For mean > 5, well approximated by a normal distribution***

$$P(k) = \frac{\lambda^k}{k!} e^{-\lambda}$$



# Normal Approximation



Can estimate Poisson distribution as a normal distribution when  $\lambda > 10$

# Pop Quiz!

I want to sequence a 10Mbp genome to 24x coverage.  
How many 120bp reads do I need?

I need  $10\text{Mbp} \times 24\text{x} = 240\text{Mbp}$  of data  
 $240\text{Mbp} / 120\text{bp} / \text{read} = 2\text{M reads}$

I want to sequence a 10Mbp genome so that  
>97.5% of the genome has at least 24x coverage.  
How many 120bp reads do I need?

Find X such that  $X - 2\sqrt{X} = 24$

$$36 - 2\sqrt{36} = 24$$

I need  $10\text{Mbp} \times 36\text{x} = 360\text{Mbp}$  of data  
 $360\text{Mbp} / 120\text{bp} / \text{read} = 3\text{M reads}$

# K-mers and K-mer counting

GATTACATACACATTGGATG

# K-mers and K-mer counting

GATTACATACACATTGGATG

GAT ACA ACA ATT GAT

ATT CAT CAC TTG ATG

TTA ATA ACA TGG

TAC TAC CAT GGA

## Kmers:

- Divide a string into substrings of length k
- Notice every position is covered k times
- Notice there are  $G - k + 1$  kmers from a string of length G

# K-mers and K-mer counting

GATTACATACACATTGGATG

GAT ACA ACA ATT GAT

ATT CAT CAC TTG ATG

TTA ATA ACA TGG

TAC TAC CAT GGA

## Kmers:

- Divide a string into substrings of length k
- Notice every position is covered k times
- Notice there are  $G - k + 1$  kmers from a string of length G

**Computation:** Very easy to compute, exact matches, represent 32mers in 64 bits

**Biological:** The “atomic unit” of a sequence, creates a fingerprint of a genome/read

# K-mers and K-mer counting

GATTACATACACATTGGATG

GAT ACA ACA ATT GAT  
ATT CAT CAC TTG ATG  
TTA ATA ACA TGG  
TAC TAC CAT GGA

GAT : 2 CAT : 2 ATG : 1 TGG : 1  
ACA : 3 CAC : 1 TTA : 1 TAC : 2  
ATT : 2 TTG : 1 ATA : 1 GGA : 1

# K-mers and K-mer counting

**GATTACATACACATTGGATG**

**GAT : 2   CAT : 2   ATG : 1   TGG : 1**

**ACA : 3   CAC : 1   TTA : 1   TAC : 2**

**ATT : 2   TTG : 1   ATA : 1   GGA : 1**

**1 : 7   (ATG, TGG, ...)**

**2 : 4   (GAT, CAT, ATT, TAC)**

**3 : 1   (ACA)**

# K-mers and K-mer counting

**GATTACATACACATTGGATG**

1: 7 (ATG, TGG, ...)

2: 4 (GAT, CAT, ATT, TAC)

3: 1 (ACA)

How long should k be?

# K-mers and K-mer counting

**GATTACATACACATTGGATG**

1 : 7 (**ATG**, **TGG**, ...)

2 : 4 (**GAT**, **CAT**, **ATT**, **TAC**)

3 : 1 (**ACA**)

How long should k be?

K=1 : Too short, every base is present

K=2 : Too short, every pair of bases will be present

Pick k so that  $G/(4^k) \ll 1$

$$k = \log_4(G)$$

At least 15 for human, often a bit longer

But not too long or could loose resolution

# Kmer counting applications

# Coverage Statistics

$$\text{sequencing\_coverage} = \frac{\text{total\_bases\_sequenced}}{\text{genome\_size}}$$

$$\text{genome\_size} = \frac{\text{total\_bases\_sequenced}}{\text{sequencing\_coverage}}$$

$$\text{genome\_size} = \frac{100\text{Gb}}{50x} = 2\text{Gb}$$

But how can you figure out  
the coverage without a genome?

# K-mer counting

## Kmer-ize

Read 1: GATTACA => GAT, ATT, TTA, TAC, ACA  
Read 2: TACAGAG => TAC, ACA, CAG, AGA, GAG  
Read 3: TTACAGA => TTA, TAC, ACA, CAG, AGA



GAT	ACA	ACA: 3
ATT	ACA	
TTA	ACA	
TAC	AGA	AGA: 2
ACA	AGA	
TAC	ATT	ATT: 1
ACA	CAG	CAG: 2
CAG	CAG	
AGA	GAG	GAG: 1
GAG	GAT	GAT: 1
TTA	TAC	TAC: 3
TAC	TAC	
ACA	TAC	
CAG	TTA	TTA: 2
AGA	TTA	

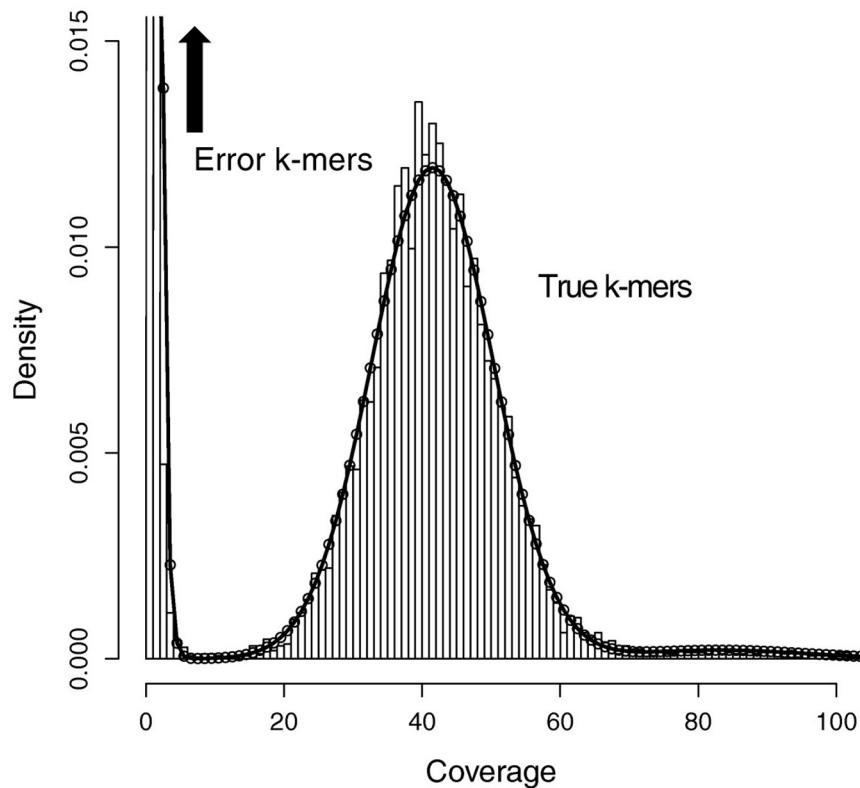
tally

sort count

From read k-mers alone, can learn something about how frequently different sequences occur (aka coverage)

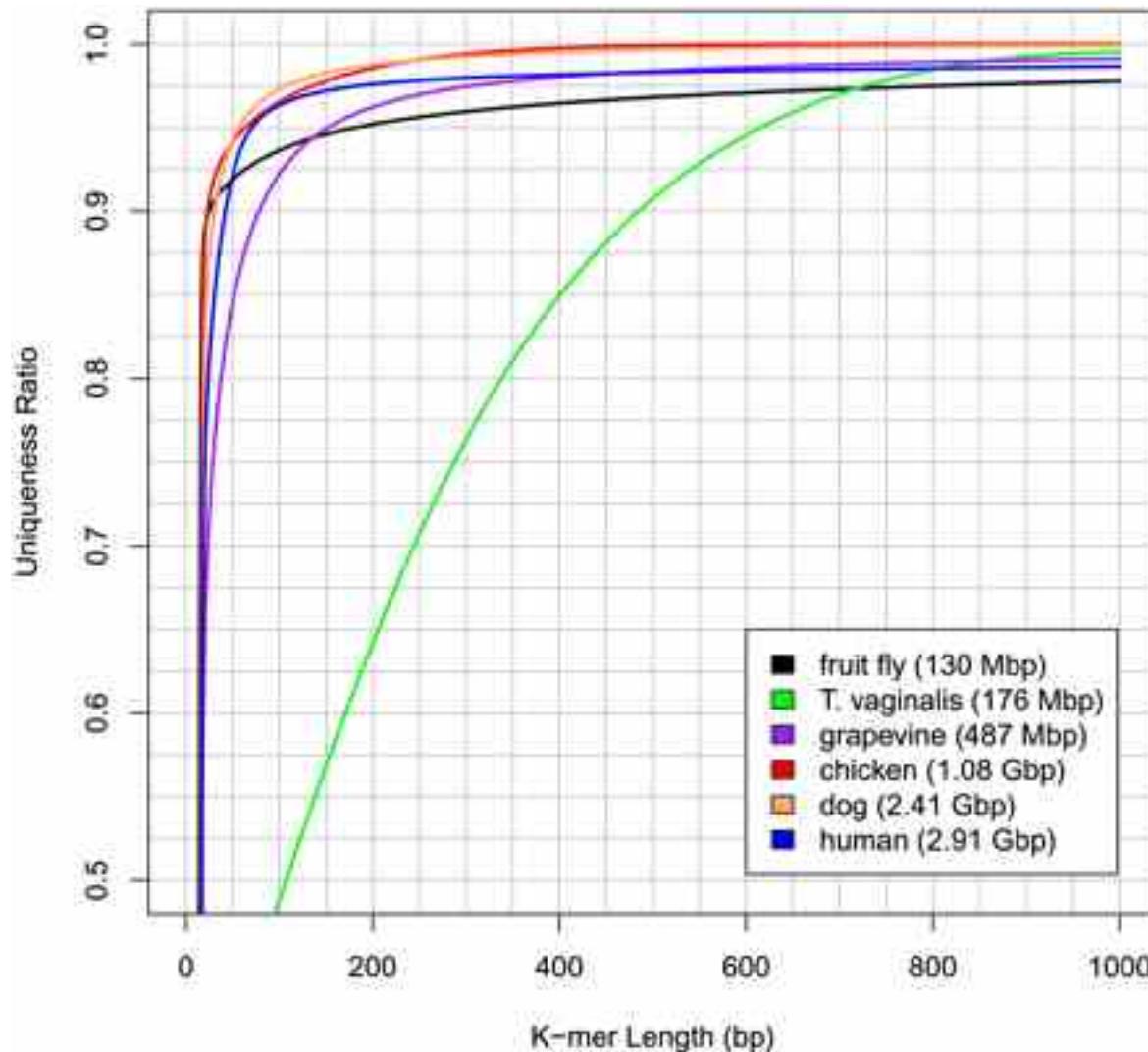
Fast to compute even over huge datasets

# K-mer counting in real genomes



- The tally of k-mer counts in real genomes reveals the coverage distribution.
- Here we sequenced 120Gb of reads from a female human (haploid human genome size is 3Gb), and indeed we see a clear peak centered at 40x coverage
- There are also many kmers that only occur <5 times. These are from errors in the reads
- There are also kmers that occur many times (>>70 times). These are repeats in the genome

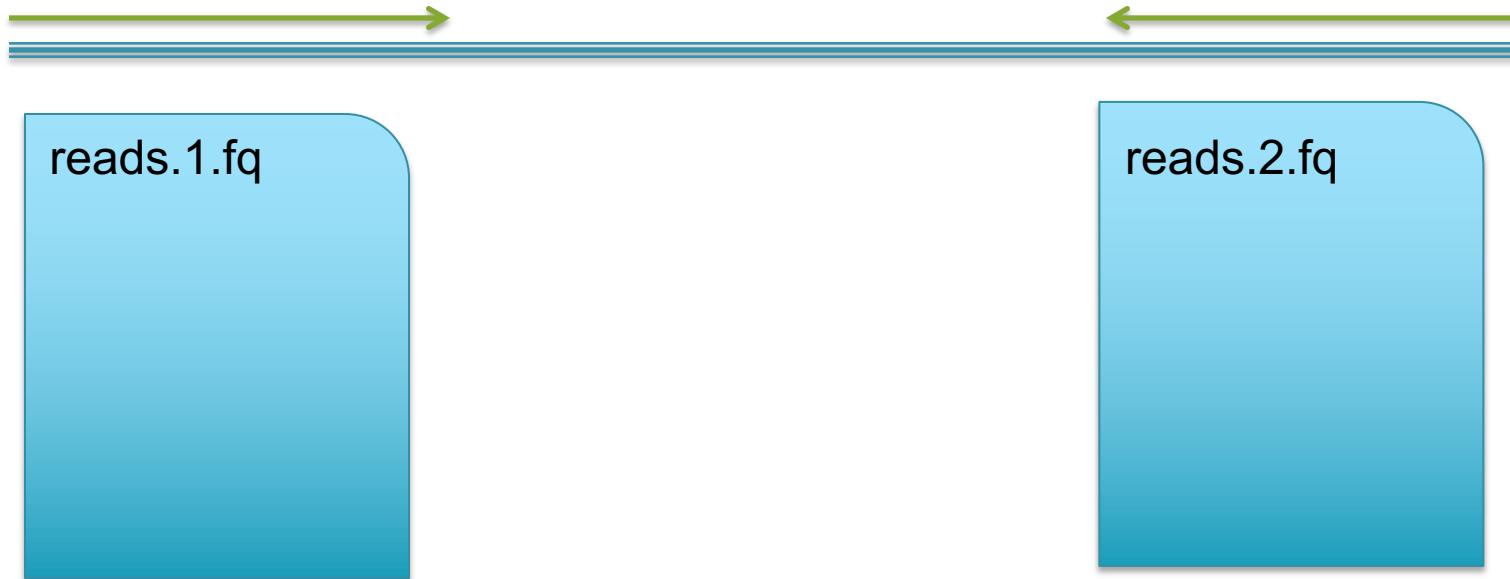
# K-mer Uniqueness



**Assembly of large genomes using second-generation sequencing**  
Schatz et al. (2010) Genome Research. doi: 10.1101/gr.101360.109

# Illumina Data Characteristics

# FASTQ Files



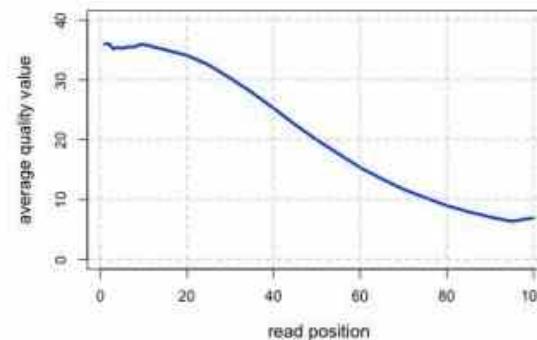
```
@SEQ_ID  
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT  
+  
!'''*((((****))%%%++)(%%%%).1***-+*'') )**55CCF>>>>>CCCCCCCC65
```

@Identifier  
Sequence  
+Separator  
Quality Values  
...

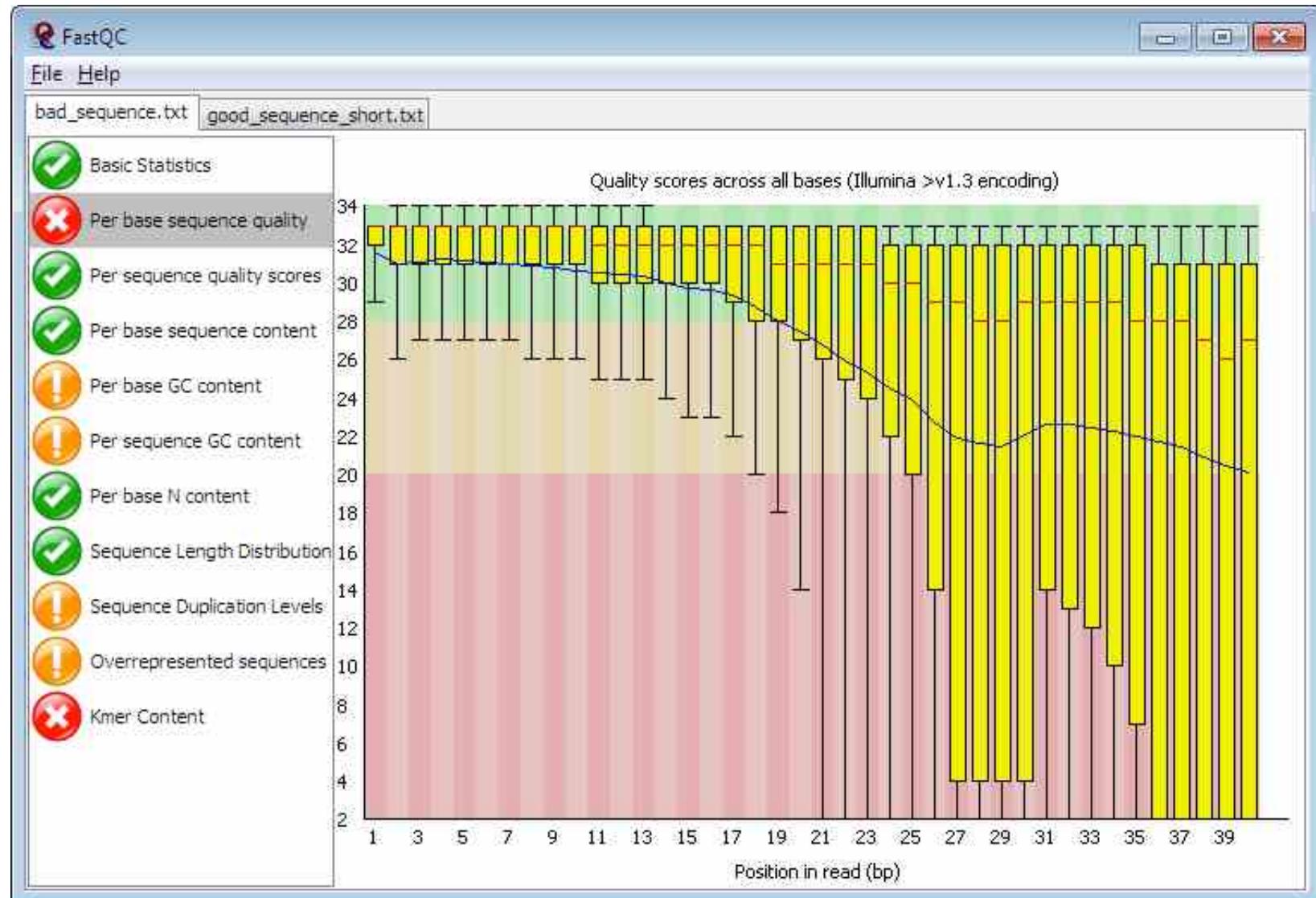
# Illumina Quality

QV	Perror
40	1/10000
30	1/1000
20	1/100
10	1/10

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

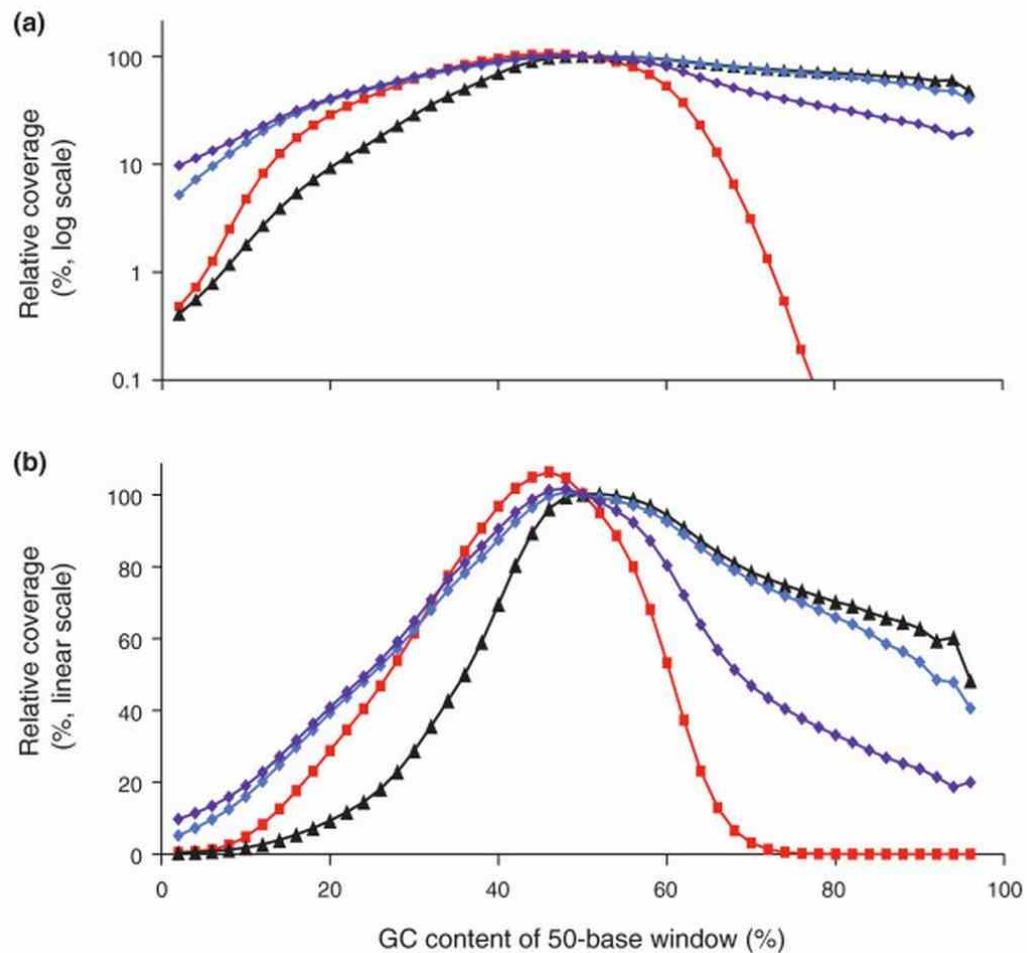


# FASTQC: Is my data any good?



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

# Beware of GC Biases

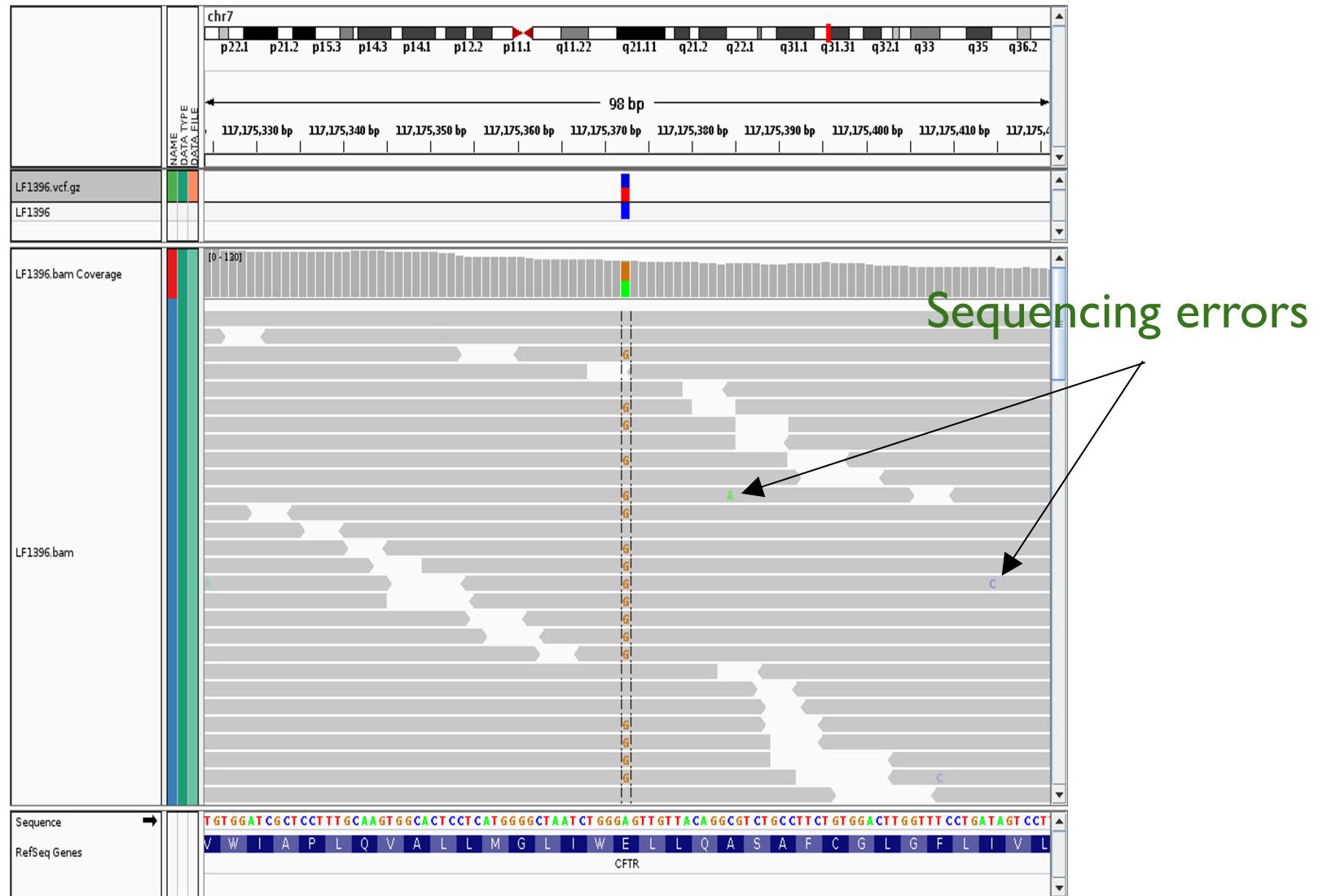


## Illumina sequencing does not produce uniform coverage over the genome

- 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.

# Sequencing errors fall out as noise (most of the time)



[https://jchoigt.files.wordpress.com/2012/07/igv\\_e217g\\_snapshot.png](https://jchoigt.files.wordpress.com/2012/07/igv_e217g_snapshot.png)

# Question?

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

What can we do?

# Illumina Hacking

**BIOINFORMATICS** ORIGINAL PAPER

Vol. 29 no. 12, 2013, pages 1492–1497  
doi:10.1093/bioinformatics/btt178

Genome analysis

Assembling the 20 Gb white spruce (*Picea glauca*) genome from whole-genome shotgun sequencing data

Inanc Birol<sup>1,2,3,\*</sup>, Anthony Raymond<sup>1</sup>, Shaun D. Jackman<sup>1</sup>, Stephen Pleasance<sup>1</sup>, Robin Coop<sup>1</sup>, Greg A. Taylor<sup>1</sup>, Macaire Man Saint Yuen<sup>4</sup>, Christopher I. Keeling<sup>4</sup>, Dana Brand<sup>1</sup>, Benjamin P. Vandervalk<sup>1</sup>, Heather Kirk<sup>1</sup>, Pawan Pandoh<sup>1</sup>, Richard A. Moore<sup>1</sup>, Yongjun Zhao<sup>1</sup>, Andrew J. Mungall<sup>5</sup>, Barry Jaquish<sup>5</sup>, Alvin Yanchuk<sup>5</sup>, Brian Boyle<sup>6</sup>, Jean Bousquet<sup>7,8</sup>, Kermit Ritland<sup>6</sup>, John MacKay<sup>7,8</sup>, Jörg E. Steven J.M. Jones<sup>1,2,9</sup>

<sup>1</sup>Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, BC V5Z 4S6, Canada; <sup>2</sup>Genetics, University of British Columbia, Vancouver, BC V6H 3N1, Canada; <sup>3</sup>School of Computer Science, University, Burnaby, BC V5A 1S6, Canada; <sup>4</sup>Michael Smith Laboratories, University of British Columbia, Vancouver, BC V6T 1Z4, Canada; <sup>5</sup>British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 0C2, Canada; <sup>6</sup>Department of Forest Sciences, University of British Columbia, 124, Canada; <sup>7</sup>Institute for Systems and Integrative Biology, Université Laval, Québec, QC G1V 0A6, Canada; <sup>8</sup>Department of Wood and Forest Sciences, Université Laval, Québec, QC G1V 0A6, Canada; <sup>9</sup>Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC V5A 1S6, Canada

Associate Editor: Michael Brum

**ABSTRACT**  
White spruce (*Picea glauca*) is a dominant conifer of the boreal forests of North America, and providing genomics resources for this commercially valuable tree will help improve forest management and conservation. We report the assembly of the largest genome ever assembled by repetitive genome shotgun through pushes the boundaries of the current technology. Here, we describe a whole-genome shotgun sequencing strategy using Illumina sequencing platforms and an assembly approach using the ABYSS software. We report a 20.0 giga base pairs draft genome in 4.9 million scaffolds, with a scaffold N50 of 20358 bp. We demonstrate how recent improvements in the sequencing technology, especially increasing read lengths and paired end reads from long fragment have a major impact on the assembly contiguity. We also note that scalable bioinformatics tools are instrumental in providing rapid draft assemblies.

**Availability:** The *Picea glauca* genome sequencing and assembly data are available at NCBI (Accession: ALW2010000000000 PDB: PRJNA83495, <http://www.ncbi.nlm.nih.gov/bioproject/83495>).  
Contact: [isrc@bcgc.ca](mailto:isrc@bcgc.ca)

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

Received on March 20, 2013; revised on April 10, 2013; accepted on April 11, 2013

**1 INTRODUCTION**  
The assembly of short reads to develop genomic resources for non-model species remains an active area of development (Schatz *et al.*, 2012). The feasibility of the approach and its scalability to

\*To whom correspondence should be addressed.

© The Author 2013. Published by Oxford University Press.  
This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is cited.

large genome was demonstrated by Simpson *et al.* (2009) using human and was later applied to assemble the 1.5 Gb genome of the SO4-dendrite tool (Li *et al.*, 2010). High quality results, as demonstrated by Chan *et al.*, 2011; Chu *et al.*, 2011; Godel *et al.*, 2012; Swart *et al.*, 2012. Estimated at 20 giga base pairs (Gbp) and assembly of the genome of the pine (*Pinaceae*) family present unique challenges, these challenges include whole-genome shotgun sequencing due to reduced representation resources of the genome. On the biotic side, we have demonstrated that using long cycle memory usage, storage and programming implementations on our

We addressed the data representation and sequencing multiple whole-genome HiSeq 2000 and MiSeq sequencers (Illumina, CA, USA). Compared to building and sequencing fosmid approach of isolating ~10-kb DNA sequencing fragments in high throughput (CA, USA), a shotgun only sequencing sequence data effectively covering the that can be an order of magnitude less especially substantial when sequencing this work, we demonstrate that at this scale remains viable and pro-

## 2 METHODS

### 2.1 Sample collection

Apricot shoot tissues were collected in April 2006 from a single white spruce (*Picea glauca*, genotype PG29) tree at the Kalamalka Research Station of the British Columbia Ministry of Forests and Ranges, Vernon, British Columbia, Canada. Genomic DNA was extracted from 60 gm tissue by Bio&Tech ([www.biotech.com](http://www.biotech.com)), Montreal, QC, Canada) using an organic extraction method yielding 300 µg of high quality purified nuclear DNA.

### 2.2 Library preparation and sequencing

DNA quality was assessed by spectrophotometry and gel electrophoresis before library construction. DNA was sheared to 45 µm using an E210 Shear仪 (Covaris) and then pooled on 8% PAGE gels. The 200 bp (for libraries with 200 bp insert size) and 500 bp (for libraries with 500 bp insert size) DNA size fractions were excised and eluted from the gel slices overnight at 4°C in 300 µl of elution buffer [5] [vol/vol] LoTE buffer (3 mM Tris-HCl (pH 7.5), 0.2 mM EDTA]/5 M ammonium acetate) and purified using a Spin-X Filter Tube (Fisher Scientific) ethanol precipitation. Gel-purified DNA samples were prepared for a modified paired-end (PE) protocol suggested by Illumina Inc. This involved DNA end repair and formation of 3' adhesive overhangs using the Klenow fragment of DNA polymerase I (3'-exonuclease minus) and ligation to Illumina PE adaptors (with 5' overhangs). Adapter-ligated products were purified on QIAquick spin columns (Qiagen) and amplified using Fusion DNA polymerase (NEB) and 10 PCR cycles with a PE primer 10 (Illumina). PCR products of the desired size range were purified using a filter column (QIAquick PCR Purification Kit, Qiagen).

The mate pair (MPET, a.k.a. long) libraries were constructed using 4 µg of sheared DNA with a Illumina PE100 kit. The PE100 construction protocol and reagent (FC-123-101). The genomic DNA sample was simultaneously fragmented and tagged with a biotin containing mate pair junction adapter, which left a short sequence gap in the fragmented DNA. The gap was filled by a strand displacement reaction using a polymerase to ensure that all fragments were flush and ready for circularization. After an AMPure Bead cleanup, size selection was done on a 0.6% agarose gel to excise 6-9 kb and 9-13 kb fractions, which were purified and combined. The mate pair junction adapter was then ligated to the ends of the fragments. A stock solution of the mate pair junction adapter was circulated by ligation, followed by a digestion to remove any linear molecules and left circularized DNA for shearing. The sheared DNA fragments that contain the biotinylated junction adapter (mate pair fragments) were purified by means of binding to streptavidin magnetic beads, and the unwanted unbiotinylated molecules were washed away. The DNA fragments were then end-repaired and A-tailed following the

protocol and ligated to indexed TruSeq adaptors. The final library was enriched by a 10-cycle PCR and purified by AMPure bead cleanup. Library quality and size was assessed by Agilent DNA 1000 series II assay and KAPA Library Quantification protocol. The two fractions were pooled for sequencing paired end 100 bp using Illumina HiSeq2000.

The construction of the 12 kb mate pair libraries was achieved by a hybrid 454/Illumina procedure. Briefly, 50 µg of genomic DNA was fragmented and 20 cycles of strand code 12 using a Hydrodash Marfham. The library was then used to construct a Mate pair library. Genomic DNA was loaded on a 1% agarose gel, and fragments <18 kb were extracted. Biotinylated circularization adaptors (TruSeq Adapters (Illumina/Roche, Brunsch, CT) and TruSeq Adapters (Illumina/Roche, Brunsch, CT) were ligated to the repaired A-tailed ends. Biotinylated adaptors were enriched using Streptavidin-coupled Dynabeads (Life Technologies, Grand Island, NY) and amplified by PCR using Illumina primers.

Genomic bacterial artificial chromosome (BAC) seq performed using DNA from the same genotype on a 454 platform with 6 kb paired-end libraries at the PlatForm Génomique of the Institute for Systems and Integrative Biology, Laval, Quebec, QC, Canada. A single paired-end library was prepared on a pool of 15 BACs (equimolar concentrations) earlier in the text with the following modifications: 15 µg of fragmented DNA using a Hydrodash with a standard assembly at 8 cycles, 18 °C, 18 °C, 10-kb fragmentation, 0.2 mM EDTA/5 M ammonium acetate) and was purified using a Spin-X Filter Tube (Fisher Scientific) ethanol precipitation. Gel-purified DNA samples were prepared for a modified paired-end (PE) protocol suggested by Illumina Inc. This involved DNA end repair and formation of 3' adhesive overhangs using the Klenow fragment of DNA polymerase I (3'-exonuclease minus) and ligation to Illumina PE adaptors (with 5' overhangs). Adapter-ligated products were purified on QIAquick spin columns (Qiagen) and amplified using Fusion DNA polymerase (NEB) and 10 PCR cycles with a PE primer 10 (Illumina). PCR products of the desired size range were purified using a filter column (QIAquick PCR Purification Kit, Qiagen).

The mate pair (MPET, a.k.a. long) libraries were constructed using 4 µg of sheared DNA with a Illumina PE100 kit. The PE100 construction protocol and reagent (FC-123-101). The genomic DNA sample was simultaneously fragmented and tagged with a biotin containing mate pair junction adapter, which left a short sequence gap in the fragmented DNA. The gap was filled by a strand displacement reaction using a polymerase to ensure that all fragments were flush and ready for circularization. After an AMPure Bead cleanup, size selection was done on a 0.6% agarose gel to excise 6-9 kb and 9-13 kb fractions, which were purified and combined. The mate pair junction adapter was then ligated to the ends of the fragments. A stock solution of the mate pair junction adapter was circulated by ligation, followed by a digestion to remove any linear molecules and left circularized DNA for shearing. The sheared DNA fragments that contain the biotinylated junction adapter (mate pair fragments) were purified by means of binding to streptavidin magnetic beads, and the unwanted unbiotinylated molecules were washed away. The DNA fragments were then end-repaired and A-tailed following the

protocol of two kits into one.

A tool was designed that opens the snap-lock latches

### Assembling the 20 Gb white spruce genome

Inanc Birol<sup>1,2,3,\*</sup>, Anthony Raymond<sup>1</sup>, Shaun D Jackman<sup>1</sup>, Stephen Pleasance<sup>1</sup>, Robin Coop<sup>1</sup>, Greg A Taylor<sup>1</sup>, Macaire Man Saint Yuen<sup>4</sup>, Christopher I Keeling<sup>4</sup>, Dana Brand<sup>1</sup>, Benjamin P. Vandervalk<sup>1</sup>, Heather Kirk<sup>1</sup>, Pawan Pandoh<sup>1</sup>, Richard A Moore<sup>1</sup>, Yongjun Zhao<sup>1</sup>, Andrew J. Mungall<sup>5</sup>, Barry Jaquish<sup>5</sup>, Alvin Yanchuk<sup>5</sup>, Carol Ritland<sup>4,6</sup>, Brian Boyle<sup>6</sup>, Jean Bousquet<sup>7,8</sup>, Kermit Ritland<sup>6</sup>, John MacKay<sup>7,8</sup>, Jörg Bohlmann<sup>4,6</sup>, Steven JM Jones<sup>1,2,9</sup>

<sup>1</sup>British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6

<sup>2</sup>University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1

<sup>3</sup>Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6

<sup>4</sup>University of British Columbia, Michael Smith Laboratories, Vancouver, BC V6T 1Z4

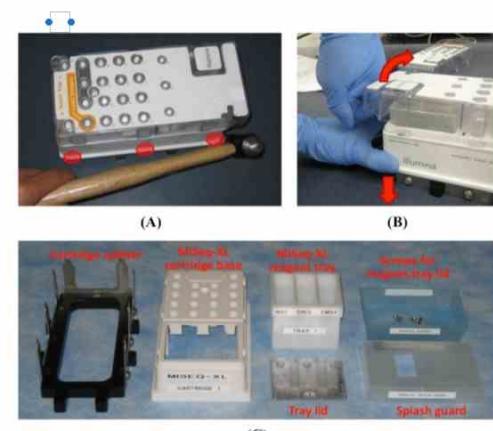
<sup>5</sup>British Columbia Ministry of Forests, Lands and Natural Resource Operations, Victoria, BC V8W 9C2

<sup>6</sup>University of British Columbia, Department of Forest Sciences, Vancouver, BC V6T 1Z4

<sup>7</sup>Université Laval, Institute for Systems and Integrative Biology, Québec, QC G1V 0A6

<sup>8</sup>Université Laval, Department of Wood and Forest Sciences, Québec, QC G1V 0A6

<sup>9</sup>Simon Fraser University, Department of Molecular Biology and Biochemistry, Burnaby, BC V5A 1S6



**Figure S1. Modification of the MiSeq cartridge.** MiSeq reagent cartridge was modified to allow for longer read lengths. (A, B) Opening of the clamshell style cartridge. (C) Contents of the modified cartridge. This was initially used to combine two PE150 kits for PE300 runs. When Illumina introduced the P250 kit, the same apparatus was used to enable PE500 runs.

# Paired-end and Mate-pairs

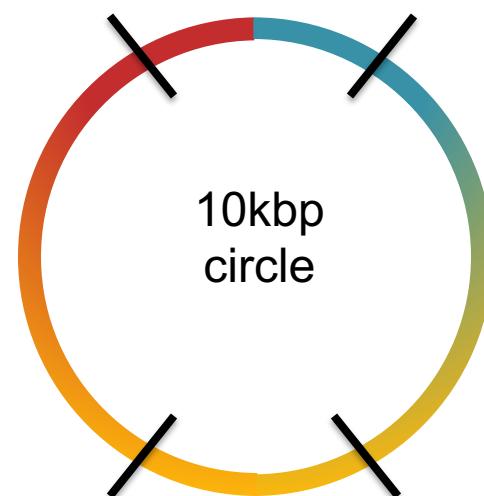
## Paired-end sequencing

- 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)



# 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)

## Disadvantages:

- Inherent limits to read length (practically, 150bp)
- Some runs are error prone
- Requires amplification, sequences a population of molecules

### Illumina NextSeq

One human genome in <30 hours

### Illumina NovaSeq

48 Genomes / run  
>8000 genomes per year  
3Tb, <3 days, ~\$300 / genome

Google ilmn

All News Finance Images Videos Shopping Web More Tools

Illumina Inc NASDAQ: ILMN

Market Summary > Illumina Inc

130.42 USD +111.39 (585.34%) ↑ all time

Closed: Aug 26, 5:20 PM EDT • Disclaimer  
After hours 130.42 0.00 (0.00%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max

Open 131.68 Mkt cap 20.78B CDP.score C  
High 132.72 P/E ratio - 52-wk high 166.65  
Low 130.15 Div yield - 52-wk low 86.50

Feedback More about Illumina Inc →

Explore more

Related Following

Moderna Inc 81.66 USD ↓ 0.95%  
Pacific Biosciences of California Inc 1.60 USD ↑ 2.56%  
Amgen Inc 326.78 USD ↓ 0.61%  
Oxford Nanopore Technologies PLC 121.04 GBX ↑ 0.37%

Disclaimer

About

[illumina.com](#)

Illumina, Inc. is an American biotechnology company, headquartered in San Diego, California. Incorporated on April 1, 1998, Illumina develops, manufactures, and markets integrated systems for the analysis of genetic variation and biological function. [Wikipedia](#)

CEO: Jacob Thayesen (Sep 25, 2023–)  
Founded: 1998

News

Google ilmn

All News Finance Images Videos Shopping Web More Tools

Illumina Inc NASDAQ: ILMN

Market Summary > Illumina Inc

130.42 USD +111.39 (585.34%) ↑ all time

Closed: Aug 26, 5:20 PM EDT • Disclaimer  
After hours 130.42 0.00 (0.00%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max

Open 131.68 Mkt cap 20.78B CDP.score C  
High 132.72 P/E ratio - 52-wk high 166.65  
Low 130.15 Div yield - 52-wk low 86.50

More about Illumina Inc → Feedback

News

www.nytimes.com/2021/07/23/science/human-genome-complete.html

SCIENCE The New York Times GIVE THE TIMES Account

MATTER

## Scientists Finish the Human Genome at Last

The complete genome uncovered more than 100 new genes that are probably functional, and many new variants that may be linked to diseases.

Share full article

A century ago, scientists knew that genes were spread across 23 pairs of chromosomes. But pinpointing any single gene and deciphering its sequence was a struggle that could have consumed a career. Michael Abbey/Science Source

By Carl Zimmer

Published July 23, 2021 Updated July 26, 2021

Sign up for Science Times Get stories that capture the wonders of nature, the cosmos and the human body. [Get it sent to your inbox.](#)

Two decades after the draft sequence of the human genome was