# Genomic Technologies

Michael Schatz

August 30, 2022

Lecture 2: Applied Comparative Genomics
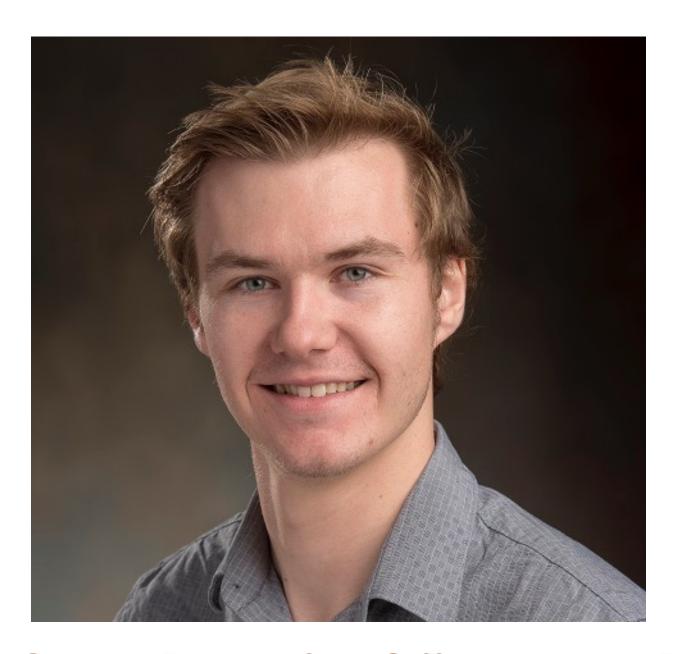
# Course Webpage
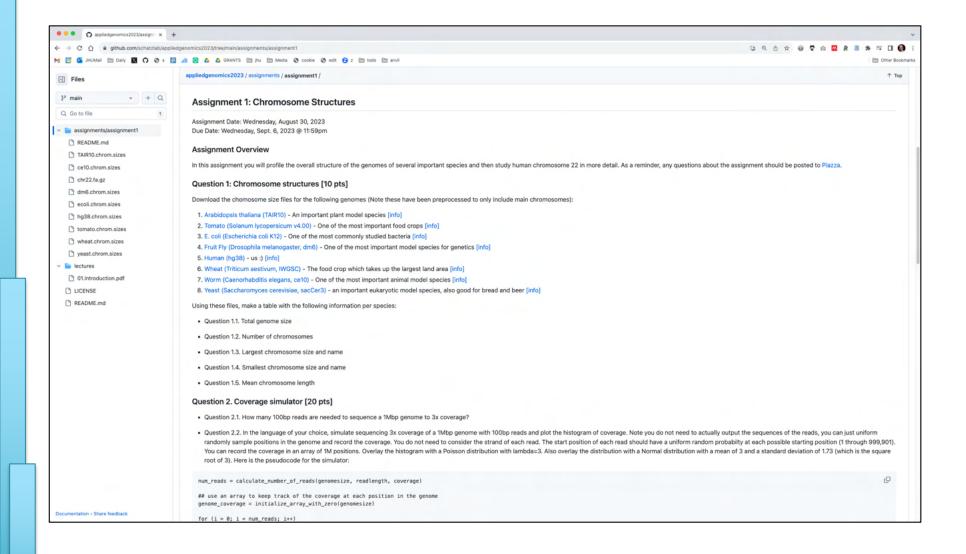


https://github.com/schatzlab/appliedgenomics2023

# TA: Alex Sweeten



# Check Piazza for Office Hours Poll

# Assignment 1

Due end of day on Sept 6 (right before midnight)

# Plotting in Python



https://matplotlib.org/

# Plotting in R / ggplot2



https://ggplot2.tidyverse.org/

# What is ChatGPT and Why Does it Work?

# ChatGPT Prompts for Academic Writing



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

# 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 is 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



*Big Data: Astronomical or Genomical?*
*Stephens, Z, et al. (2015) PLOS Biology DOI: 10.1371/journal.pbio.1002195*

# Sequencing Capacity



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?*

# Comparative Genomics Technologies



**Results**
Domain Knowledge

**Machine Learning**
classification, modeling,
visualization & data Integration

**Scalable Algorithms**
Streaming, Sampling, Indexing, Parallel

**Compute Systems**
CPU, GPU, Distributed, Clouds, Workflows

**IO Systems**
Hardrives, Networking, Databases, Compression, LIMS

**Sensors & Metadata**
Sequencers, Microscopy, Imaging, Mass spec, Metadata & Ontologies

# Comparative Genomics Technologies



**Results**
Domain Knowledge

**Machine Learning**
classification, modeling,
visualization & data Integration

**Scalable Algorithms**
Streaming, Sampling, Indexing, Parallel

**Compute Systems**
CPU, GPU, Distributed, Clouds, Workflows

**IO Systems**
Hardrives, Networking, Databases, Compression, LIMS

**Sensors & Metadata**
Sequencers, Microscopy, Imaging, Mass spec, Metadata & Ontologies

# Comparative Genomics Technologies



**Results**
Domain Knowledge

**Machine Learning**
classification, modeling, visualization & data Integration

**Scalable Algorithms**
Streaming, Sampling, Indexing, Parallel

**Compute Systems**
CPU, GPU, Distributed, Clouds, Workflows

**IO Systems**
Hardrives, Networking, Databases, Compression, LIMS

**Sensors & Metadata**
Sequencers, Microscopy, Imaging, Mass spec, Metadata & Ontologies

# 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



**Results**
Domain Knowledge

**Machine Learning**
classification, modeling, visualization & data Integration

**Scalable Algorithms**
Streaming, Sampling, Indexing, Parallel

**Compute Systems**
CPU, GPU, Distributed, Clouds, Workflows

**IO Systems**
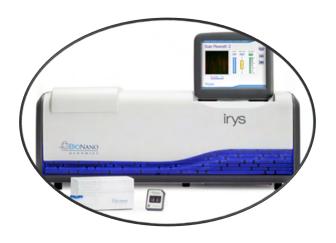Hardrives, Networking, Databases, Compression, LIMS

**Sensors & Metadata**
Sequencers, Microscopy, Imaging, Mass spec, Metadata & Ontologies

# Genomics Arsenal in the year 2023

Transcripts

Quantification of mature transcripts and small RNA

RNA-seq

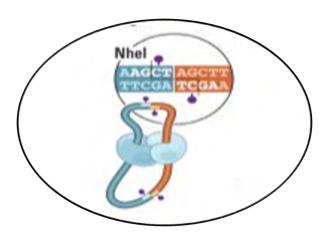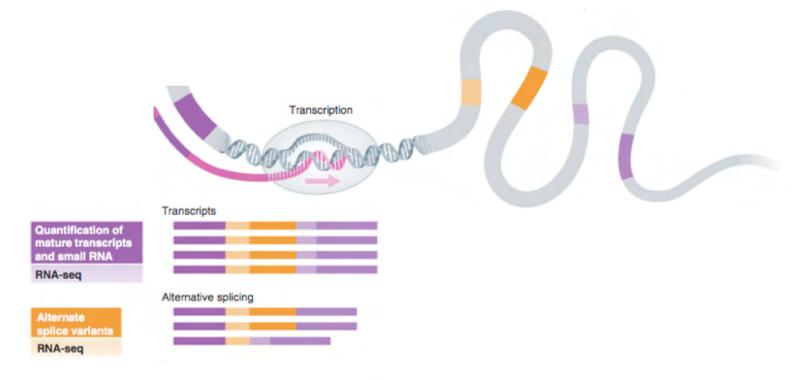Alternate splice variants
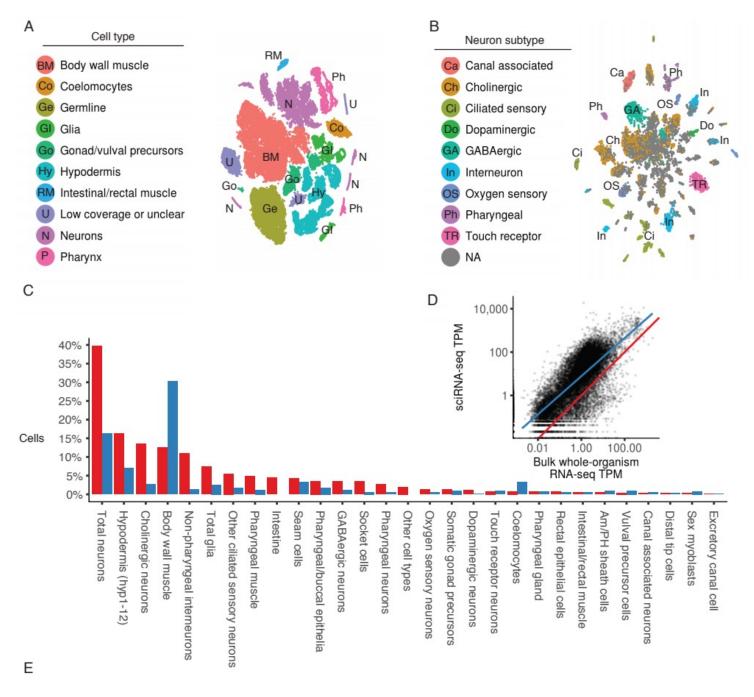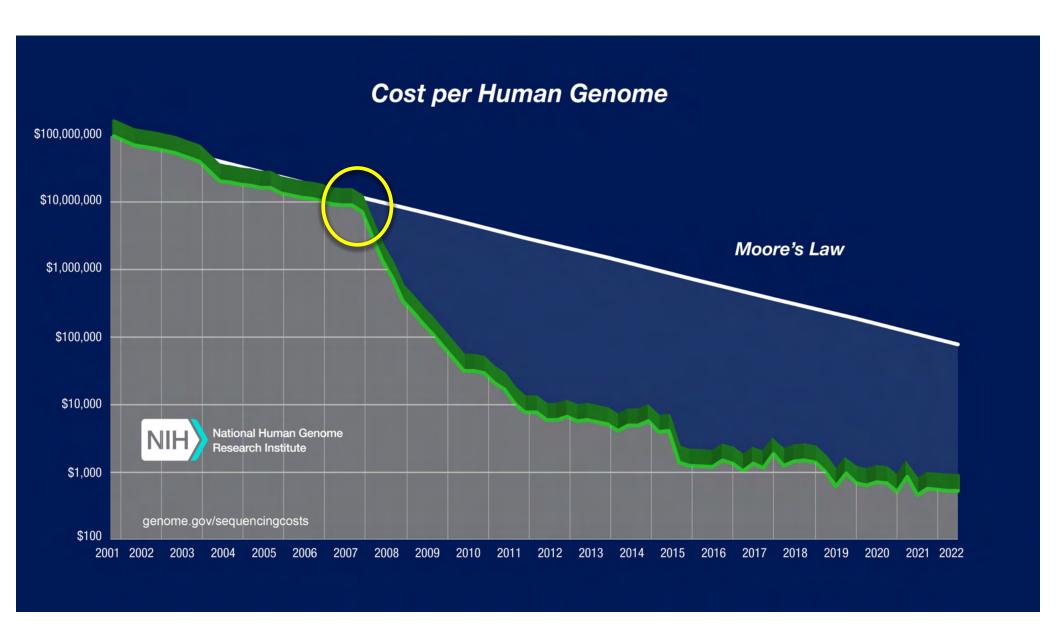
RNA-seq

Alternative splicing

Soon et al., Molecular Systems Biology, 2013

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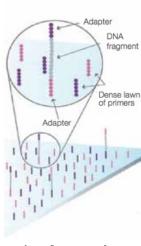
# Cost per Genome



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

# Second Generation Sequencing
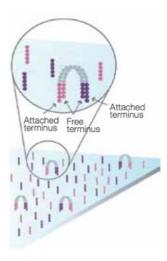


**Illumina NovaSeq 6000**
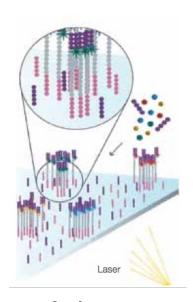*Sequencing by Synthesis*

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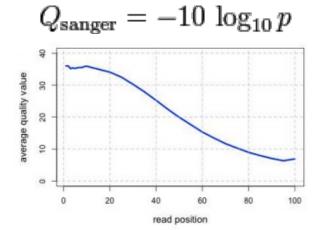
1. Attach

2. Amplify

3. Image

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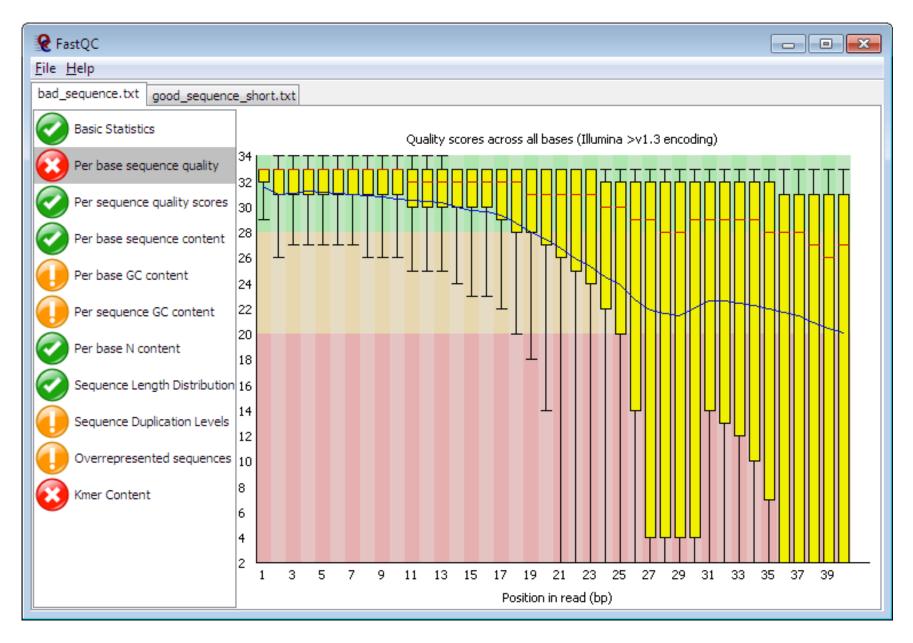https://www.youtube.com/watch?v=fCd6B5HRaZ8

# Illumina Quality

| QV | $p_{error}$ |
|---|---|
| 40 | 1/10000 |
| 30 | 1/1000 |
| 20 | 1/100 |
| 10 | 1/10 |

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

```
  SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS....................................
  ..............................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....................
  ...........................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII..............
  .................................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.................
 LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL....................................
 !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
 |                              |     |         |                                    |           |
 33                             59    64        73                                   104         126

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)
```
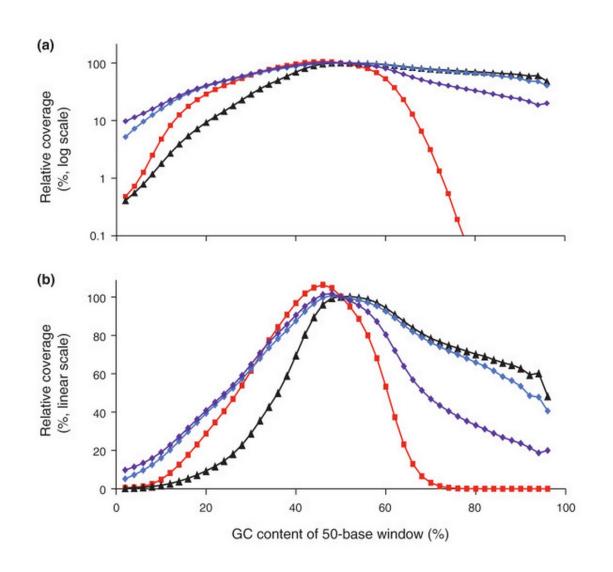
http://en.wikipedia.org/wiki/FASTQ_format

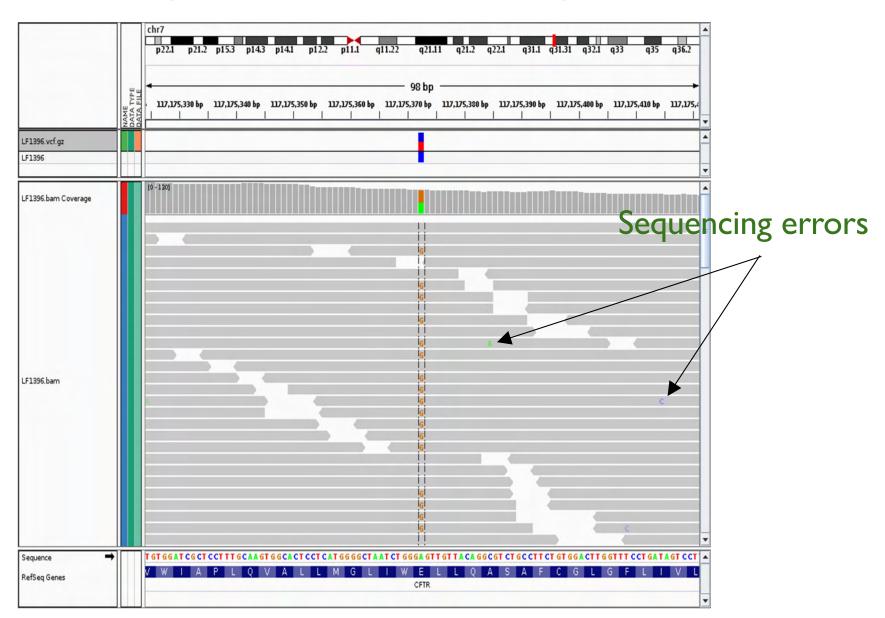# FASTQC: Is my data any good?

# Beware of GC Biases



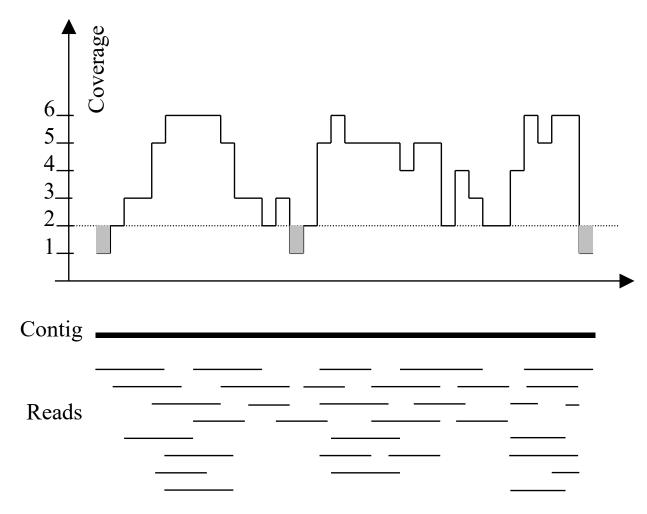**Illumina sequencing does not produce uniform coverage over the genome**

- 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)



Sequencing errors

# Typical sequencing coverage



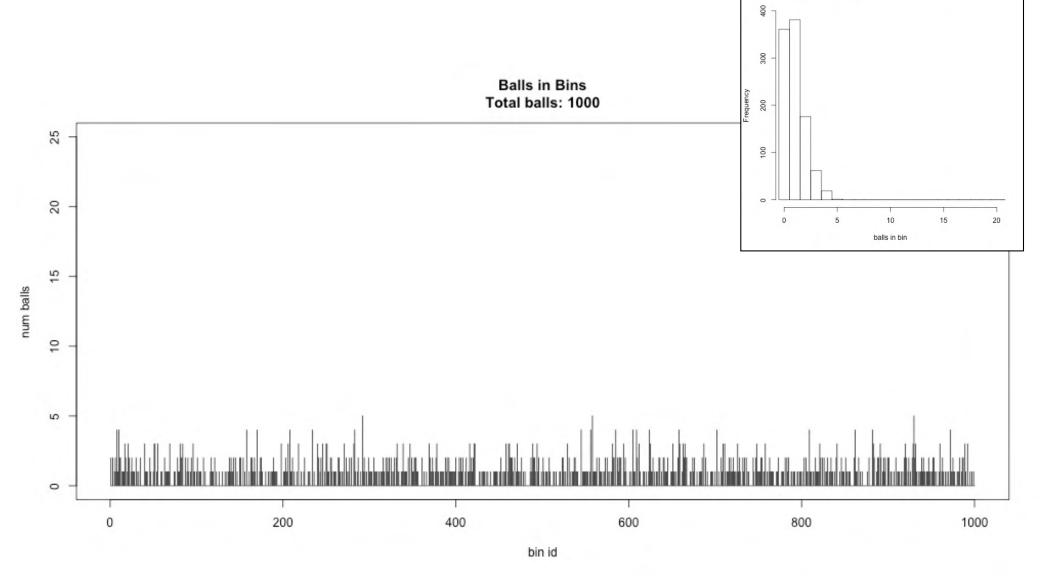Imagine raindrops on a sidewalk

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

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

# 1x sequencing



**Balls in Bins**
**Total balls: 1000**

**Histogram of balls in each bin**
**Total balls: 1000  Empty bins: 361**

# 2x sequencing
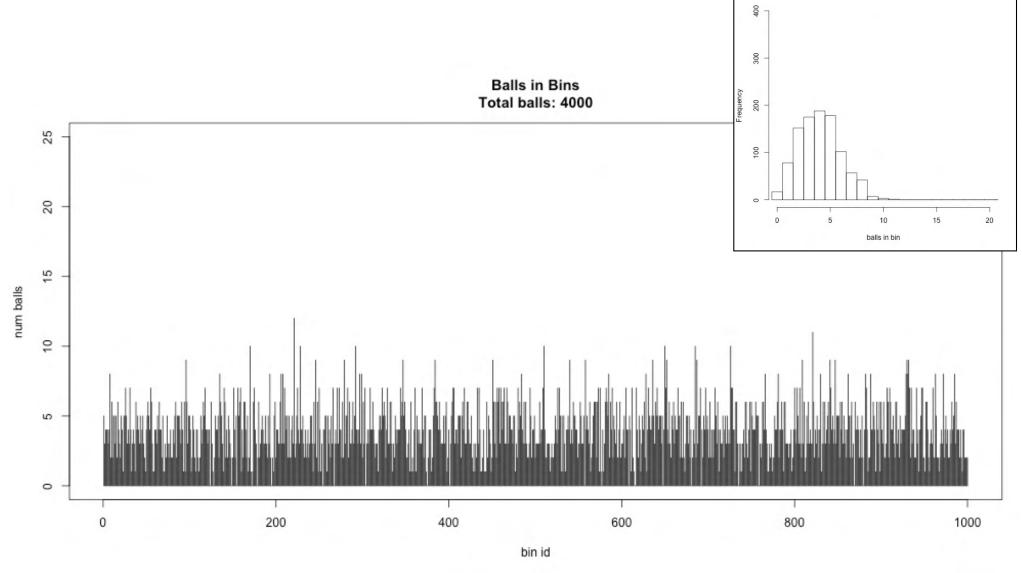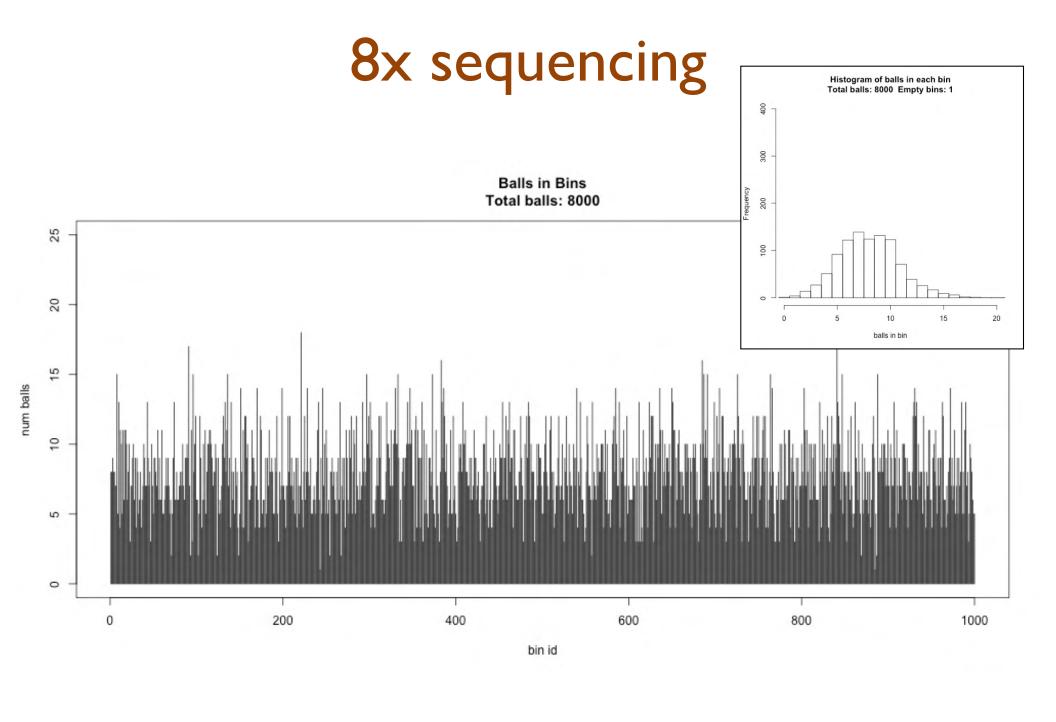
# 4x sequencing



Balls in Bins
Total balls: 4000

Histogram of balls in each bin
Total balls: 4000  Empty bins: 17

# 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.
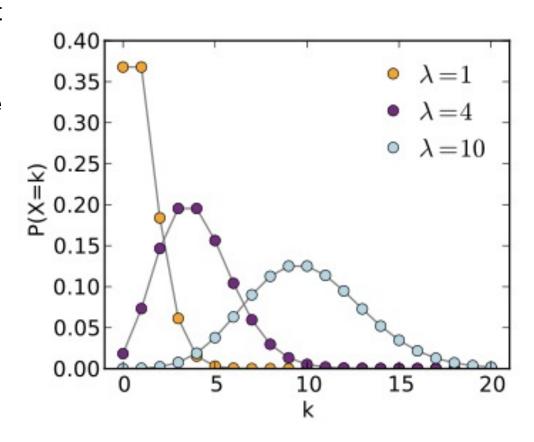
Formulation comes from the limit of the binomial equation

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



99.7% of the data are within 3 standard deviations of the mean

95% within 2 standard deviations

68% within 1 standard deviation

$\mu - 3\sigma$  $\mu - 2\sigma$  $\mu - \sigma$  $\mu$  $\mu + \sigma$  $\mu + 2\sigma$  $\mu + 3\sigma$

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

# 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

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)
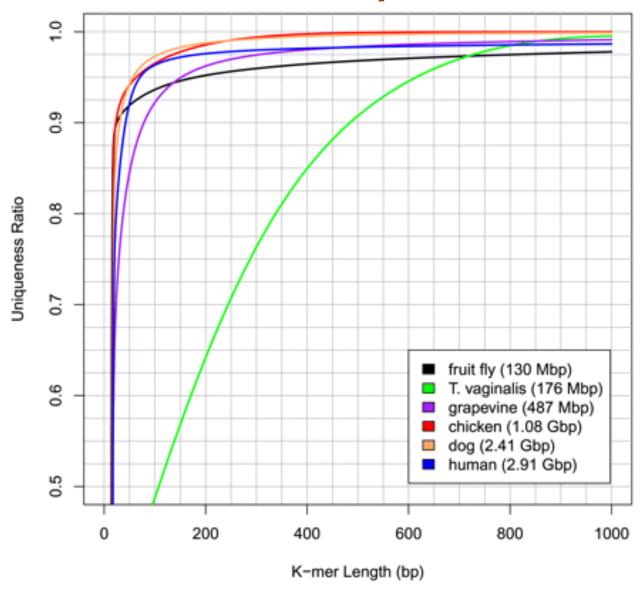
How long should k be?

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) << 1
k = log_4 (G)
At least 15 for human, often a bit longer
But not too long or could loose resolution

# K-mer Uniqueness



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