

# Genomic Technologies

Michael Schatz

Feb 1, 2018

Lecture 2: Applied Comparative Genomics



# 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/appliedgenomics2018>

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

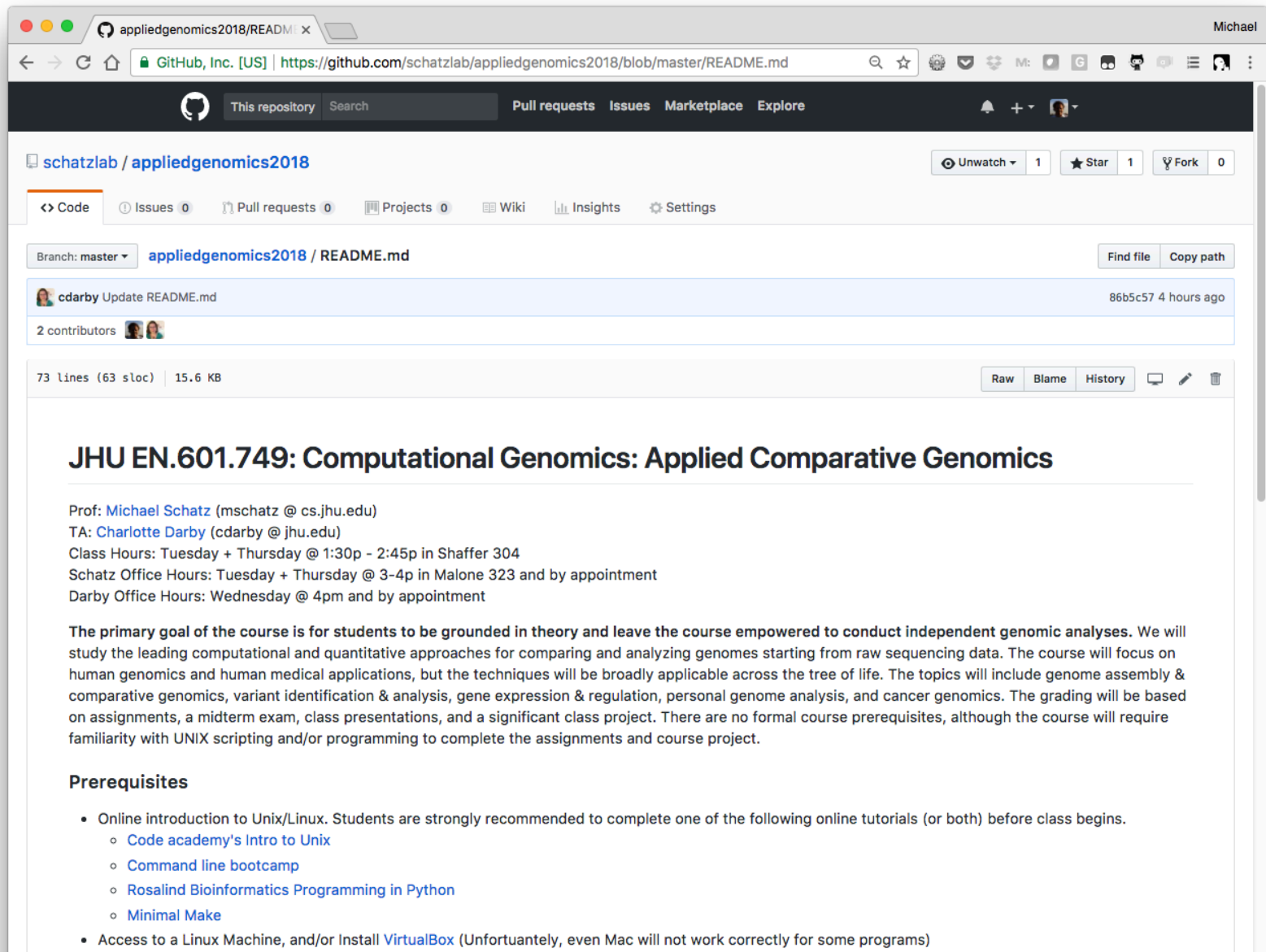
**Class Hours:** Tues + Thurs @ 1:30p – 2:45p, Shaffer 304

**Schatz Office Hours:** Tues + Thurs @ 3-4p and by appointment

**Darby Office Hours:** Wed @ 4-5 and by appointment

Please try Piazza first!

# Course Webpage



The screenshot shows a web browser displaying the GitHub repository page for 'schatzlab / appliedgenomics2018'. The page title is 'appliedgenomics2018 / README.md'. The repository has 1 star and 0 forks. The current branch is 'master'. The README file is 73 lines (63 sloc) and 15.6 KB. The commit history shows a recent update by 'cdarby' 4 hours ago. The main content of the README is for the course 'JHU EN.601.749: Computational Genomics: Applied Comparative Genomics'. The course is taught by Prof. Michael Schatz and TA Charlotte Darby. The primary goal is for students to be grounded in theory and leave the course empowered to conduct independent genomic analyses. The prerequisites include an online introduction to Unix/Linux, access to a Linux machine, and familiarity with UNIX scripting and programming.

Branch: master **appliedgenomics2018 / README.md** Find file Copy path

cdarby Update README.md 86b5c57 4 hours ago  
2 contributors

73 lines (63 sloc) | 15.6 KB Raw Blame History

## JHU EN.601.749: Computational Genomics: Applied Comparative Genomics

Prof: [Michael Schatz](mailto:mschatz@cs.jhu.edu) (mschatz@cs.jhu.edu)  
TA: [Charlotte Darby](mailto:cdarby@jhu.edu) (cdarby@jhu.edu)  
Class Hours: Tuesday + Thursday @ 1:30p - 2:45p in Shaffer 304  
Schatz Office Hours: Tuesday + Thursday @ 3-4p in Malone 323 and by appointment  
Darby Office Hours: Wednesday @ 4pm and 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.

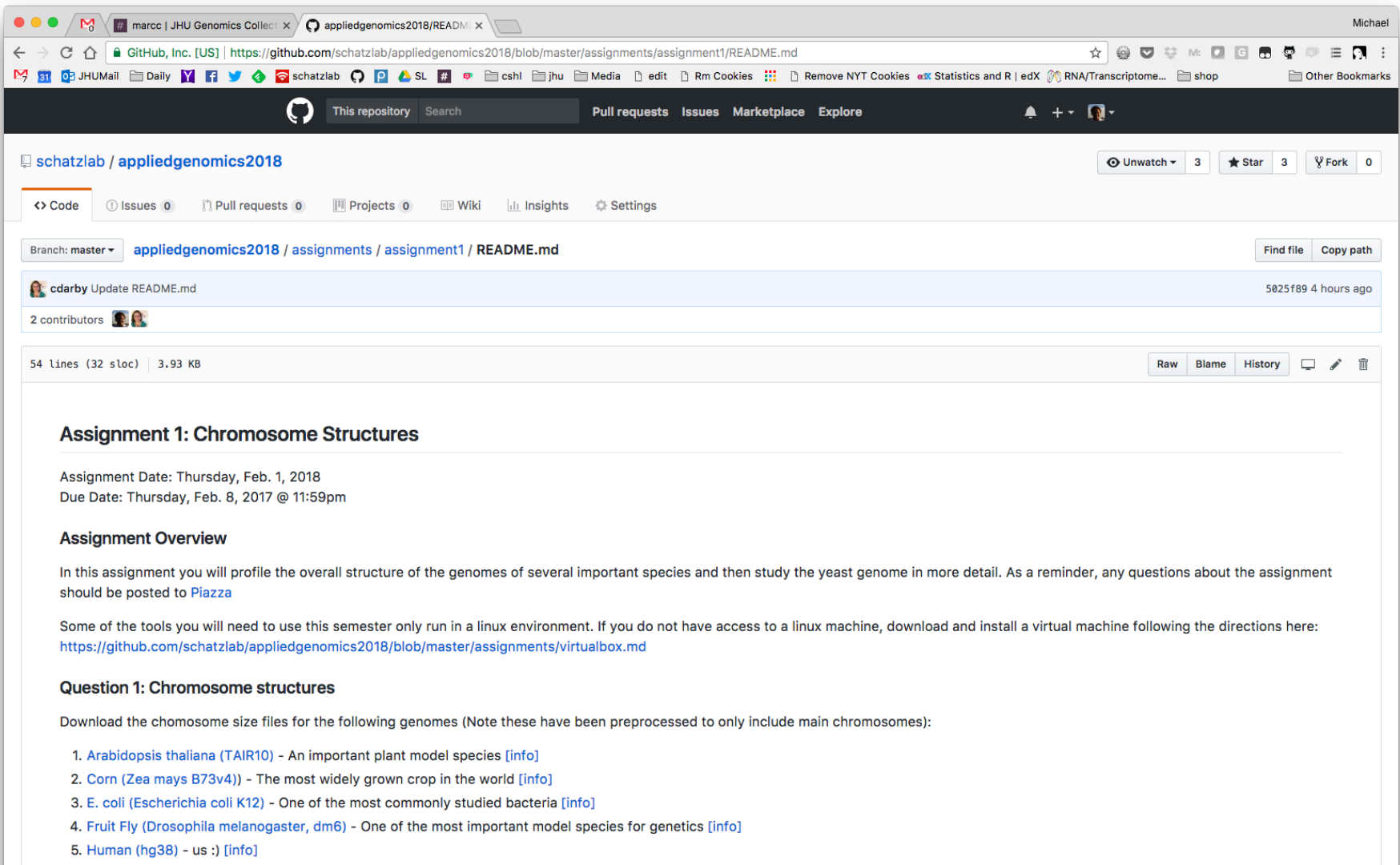
### Prerequisites

- 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](#)
  - [Command line bootcamp](#)
  - [Rosalind Bioinformatics Programming in Python](#)
  - [Minimal Make](#)
- Access to a Linux Machine, and/or Install [VirtualBox](#) (Unfortunately, even Mac will not work correctly for some programs)

<https://github.com/schatzlab/appliedgenomics2018>

# Assignment 1: Chromosome Structures

## Due Feb 8 @ 11:59pm



The screenshot shows a web browser displaying a GitHub repository page. The repository is named 'appliedgenomics2018' and is owned by 'schatzlab'. The current file being viewed is 'README.md' located in the path 'assignments/assignment1/'. The page shows a commit by 'cdarby' titled 'Update README.md' from 4 hours ago. The README content includes the following sections:

### Assignment 1: Chromosome Structures

Assignment Date: Thursday, Feb. 1, 2018  
Due Date: Thursday, Feb. 8, 2017 @ 11:59pm

#### Assignment Overview

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

Some of the tools you will need to use this semester only run in a linux environment. If you do not have access to a linux machine, download and install a virtual machine following the directions here: <https://github.com/schatzlab/appliedgenomics2018/blob/master/assignments/virtualbox.md>

#### 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. Corn (*Zea mays* B73v4) - The most widely grown crop in the world [[info](#)]
3. *E. coli* (*Escherichia coli* K12) - One of the most commonly studied bacteria [[info](#)]
4. Fruit Fly (*Drosophila melanogaster*, dm6) - One of the most important model species for genetics [[info](#)]
5. Human (hg38) - us :) [[info](#)]

<https://github.com/schatzlab/appliedgenomics2018>

# Piazza

The screenshot shows a web browser window with the URL <https://piazza.com/class/jcumooljtd46p7?cid=6>. The page title is "EN. 601.749". The user is logged in as "Michael".

The main content area displays a "note" section with the following text:

**Welcome!**  
Welcome to JHU EN.601.749: Computational Genomics: Applied Comparative Genomics  
Please feel free to ask any questions here! Also see the course webpage here:  
<https://github.com/schatzlab/appliedgenomics2018>  
Good luck!!  
Mike & Charlotte

Below the note is a "followup discussions" section with the text "Start a new followup discussion" and a text input field containing "Compose a new followup discussion".

At the bottom of the page, there are statistics:

Average Response Time:	Special Mentions:	Online Now	This Week:
N/A	There are no special mentions at this time.	1	4

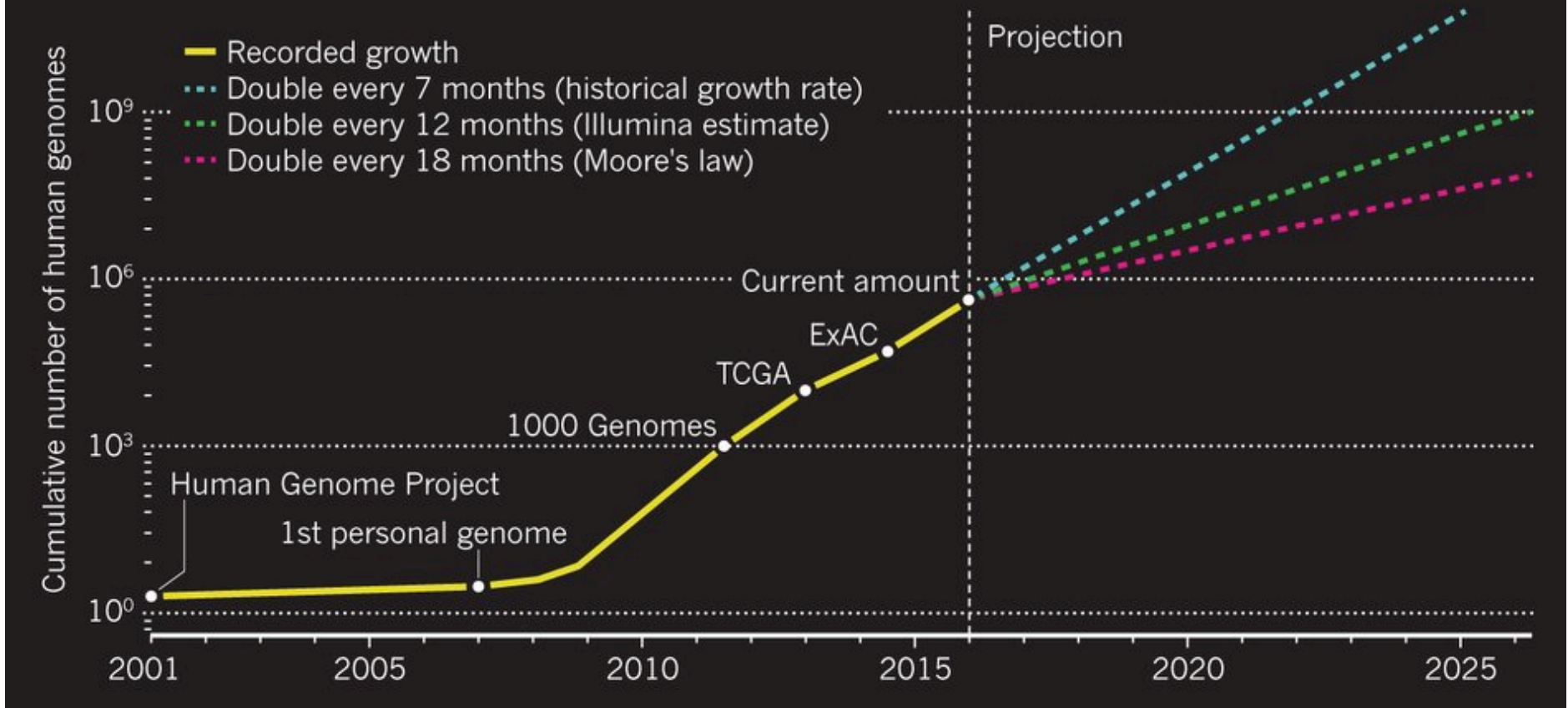
Copyright © 2018 Piazza Technologies, Inc. All Rights Reserved. Privacy Policy Copyright Policy Terms of Use Blog Report Bug!

<http://piazza.com/jhu/spring2018/en601749>

# 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?**

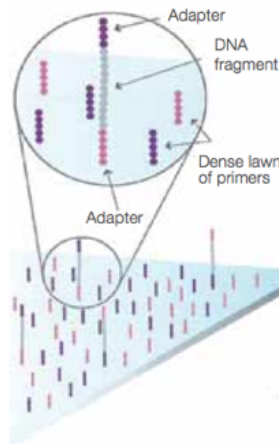
Stephens, Z, et al. (2015) PLOS Biology DOI: 10.1371/journal.pbio.1002195

# Second Generation Sequencing

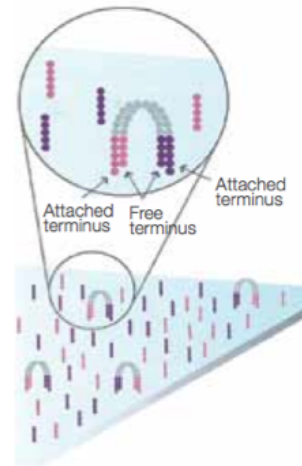


**Illumina HiSeq 2000**  
*Sequencing by Synthesis*

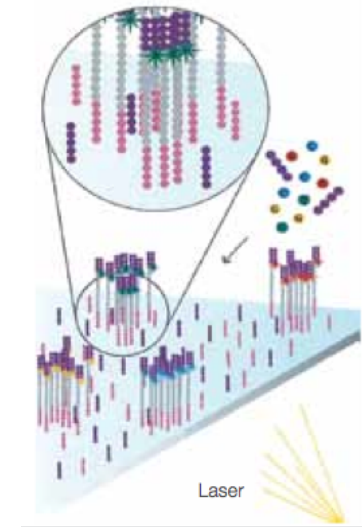
>60Gbp / day



1. Attach



2. Amplify



3. Image



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

# Question?

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

What can we do?



# Genome Hacking

## 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 Coope<sup>1</sup>, Greg A. Taylor<sup>1</sup>, Maccuire 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>1</sup>, Barry Jaquish<sup>5</sup>, Alvin Yanchuk<sup>5</sup>, Ca Brian Boyle<sup>7</sup>, Jean Bousquet<sup>7,8</sup>, Kermit Ritland<sup>6</sup>, John MacKay<sup>7,8</sup>, Jörg B. Ståhl<sup>9</sup>, 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 Computing Science, University of Alberta, Edmonton, Alberta, Canada, <sup>4</sup>Michael Smith Laboratories, University of British Columbia, BC V6T 1Z4, Canada, <sup>5</sup>British Columbia Ministry of Forests, Lands and Natural Resource Operations, BC V8W 9C2, Canada, <sup>6</sup>Department of Forest Sciences, University of British Columbia, BC V1N 1C2, 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>Department of Molecular Biology and Biochemistry, Simon Fraser University, Burnaby, BC V5A 1S6, Canada

Associate Editor: Michael Brudno

### 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 efforts. Sequencing and assembling the large and highly repetitive spruce genome through pushes the boundaries of the current technology. Here, we describe a whole-genome shotgun sequencing strategy using two Illumina sequencing platforms and an assembly approach using the ABySS software. We report a 20.8 Gb genome size and a 4.9 million scaffolds, with a scaffold N50 of 20356 bp. We demonstrate how recent improvements in the sequencing technology, especially increasing read lengths and paired end reads from longer fragments 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 through NCBI (Accession#: ALXJ2010000000 PID: PRJNA83438). <http://www.ncbi.nlm.nih.gov/bioproject/83438>.

Contact: [birol@bcgsc.ca](mailto:birol@bcgsc.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/2.0/>), which permits non-commercial reuse, distribution, and reproduction in any medium, provided the original work is properly cited. For all other uses, permission should be sought from Oxford University Press.

large genomes was demonstrated by (Simpson *et al.*, 2009) using human and was later used to assemble the SOA/Denovo tool (Li *et al.*, 2010), high quality results, as demonstrate *et al.*, 2012; Ladner *et al.*, 2013; Ritland *et al.*, 2013). The ABySS software (Chan *et al.*, 2011; Chiu *et al.*, 2011; Godel *et al.*, 2012; Swart *et al.*, 2012). Estimated at 20 Giga base pairs (Gb) and assembly of the genome of the pine (*Pinaceae*) family present unique generation end, those challenges include whole-genome shotgun sequencing data reduced representation resources due of the problem. On the bicla massive sequencing datasets is extremely cycles, memory usage, storage re programming implementations on co

We addressed the data representation and sequencing multiple whole-genome HiSeq 2000 and MiSeq sequencers (CA, USA). Compared with localized as building and sequencing fosmid approach of isolating ~10-kb DNA sequencing fragments in high through CA, USA), a shotgun only sequencing sequence data effectively covering 1e that can be an order of magnitude less especially substantial when sequencing In this work, we demonstrate that at this scale remains viable and pro

assemble the spruce genome, we used the ABySS algorithm (Simpson *et al.*, 2009), which captures a representation of read-to-read overlaps by a distributed de Bruijn graph and uses parallel computations to build the genome. The modular nature of the tool allowed us to execute a large number of tests to tune the message passing interface for a successful execution, train the assembly parameters for an optimal assembly and quantify the utility of long reads for large genome assemblies. To the best of our knowledge, the ABySS algorithm is unique in its ability to enable genome assemblies of this scale using whole-genome shotgun sequencing data.

### 2 METHODS

#### 2.1 Sample collection

Apical 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 mg tissue by BioS&T (<http://www.biost.com/>, Montreal, QC, Canada) using an organelle exclusion 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 for 45 s using an E210 sonicator (Covaris) and then analysed on 8% PAGE gels. The 200–300 bp (for libraries with 250 bp insert size) or 450–550 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 [5.1 (vol/vol) LaTE buffer [3 mM Tris-HCl (pH 7.5), 0.2 mM EDTA]/7.5 M ammonium acetate) and was purified using a Spin-X Filter Tube (Fisher Scientific) and ethanol precipitation. Genome libraries were prepared using a modified paired-end tag (PET) protocol supplied by Illumina Inc. This involved DNA end repair and formation of 3' adenosine overhangs using the Klenow fragment of DNA polymerase I (3'-5' exonuclease minus) and ligation to Illumina PE adapters (with 5' overhangs). Adapter-ligated products were purified on QIAquick spin columns (Qiagen) and amplified using Phusion DNA polymerase (NEB) and 10 PCR cycles with the PE primer 1.0 and 2.0 (Illumina). PCR products of the desired size range were purified from adapter ligation artifacts using 8% PAGE gels. DNA quality was assessed and quantified using an Agilent DNA 1000 series II assay (Agilent) and Nanodrop 7500 spectrophotometer (Nanodrop). DNA was subsequently diluted to 8 nM. The final concentration was confirmed using a Quant-iT dsDNA HS assay kit and Qubit fluorometer (Invitrogen).

The mate pair (MPET, a.k.a. jumping) libraries were constructed using 4 µg of genomic DNA with the Illumina Nextera Mate-Pair library construction protocol and reagent (FC-132-1001). 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 using a ZymoClean Large Fragment DNA Recovery Kit. The fragments were circularized 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 non-biotinylated molecules were washed away. The DNA fragments were then end repaired and A-tailed following the

protocol and ligated to indexed TruSeq adapters. The final library was enriched by a 10-cycle PCR and purified by AMPure bead clean-up. Library quality and size were 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 for 20 cycles at speed code 12 using a HydroShear (Marlborough, MA) equipped with a large assembly module. Fragmented DNA was loaded on a 1% agarose gel, and fragments 18 kb were extracted. Biotinylated circularization adapters Titanium Paired-end Adaptor set (454 Life Sciences/Rochester, CT) were added to ends of the gel-extracted fragments. Recombination of the ends was performed with Cre recombinase (New England Biolabs, Ipswich, MA), and linear molecules remaining were removed with Plasmid Safe (Epicentre, Madison, WI) molecules were fragmented using GS Rapid Library Nebulization (Roche, Branford, CT), and fragment end-repair for tailing was performed with the GS Rapid Library preparation Sciences/Roche, Branford, CT). TruSeq Adaptors (Illumina, CA) were ligated to the repaired/A-tailed ends. Biotinylated were enriched using Streptavidin-coupled Dynabeads (Life Technologies, Grand Island, NY) and amplified by PCR using Illumina primers.

Random bacterial artificial chromosome (BAC) sequencing was performed using DNA from the same genotype on 4. Titanium with 6 kb paired-end libraries at the Platform Genomics of the Institute for Systems and Integrative Biology (Université Laval, Québec City, QC). 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 fragmented using a HydroShear with a standard assembly module at speed code 18, 6–10 kb fragments were extracted from GS-FLX library adaptors were ligated to the repaired/A-tailed ends. GS-FLX sequencing using the titanium chemistry was performed according to manufacturer's instructions (454 Life Sciences/Rochester, CT). Sanger sequencing method was used to verify BAC sequencing data as previously described (Hamberger, Keeling *et al.*, 2010).

#### 2.3 MiSeq modification

In sequencing the spruce genome, we generated longer read lengths by modifying the MiSeq platform. The MiSeq uses a clamshell style cartridge (Supplementary Fig. S1A) to hold reagent tubes in an arrangement that allows for independent steps such as denaturation and cluster generation at each cycle. Although the MiSeq allows any read length specified in the control software, the reagent cartridge cannot be used during the run without stopping it. Increasing the read length requires increasing the quantity of the length-dependent reagent cartridge. This led to the solution of combining the length reagents of two kits into one.

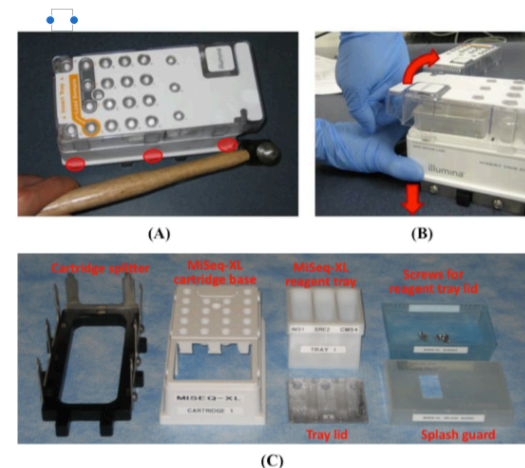
A tool was designed that opens the snap-hook latches cartridge together (Supplementary Figs S1B and S2), which allows the reagent tubes, yet allowing the cartridge to be put back together without damage to its components (Supplementary Fig. S3). The stock length-dependent reagent containers allow a run of ~650 cycles in total. To maximize the potential of the kit approach, a new reagent tray with 70 ml wells was placed in a modified clamshell base.

### Assembling the 20 Gb white spruce genome

## 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 Coope<sup>1</sup>, Greg A Taylor<sup>1</sup>, Maccuire 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>1</sup>, Barry Jaquish<sup>5</sup>, Alvin Yanchuk<sup>5</sup>, Carol Ritland<sup>6</sup>, Brian Boyle<sup>7</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>

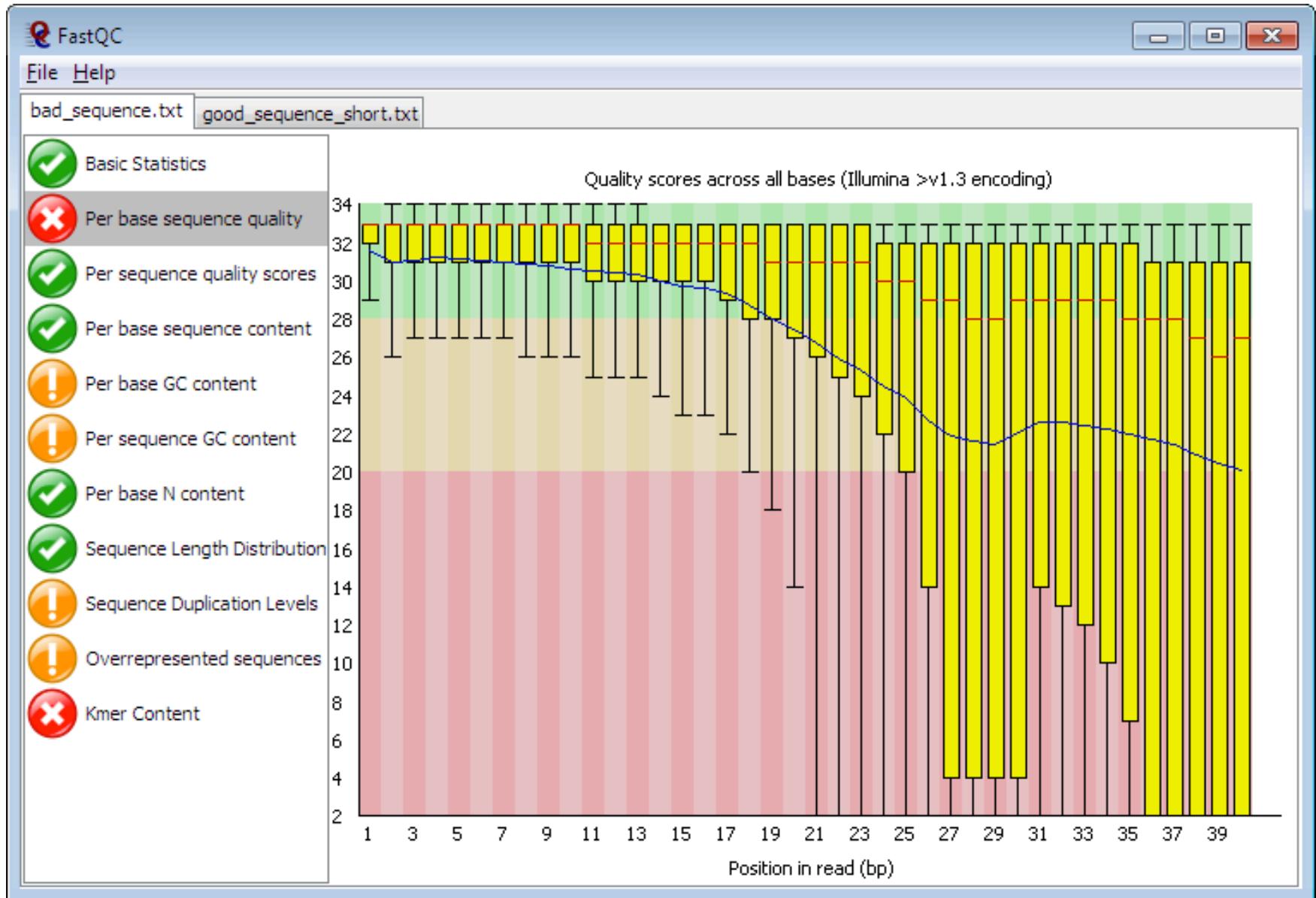
- <sup>1</sup> British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6
- <sup>2</sup> University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1
- <sup>3</sup> Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6
- <sup>4</sup> University 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.



# FASTQC: Is my data any good?



# Paired-end and Mate-pairs

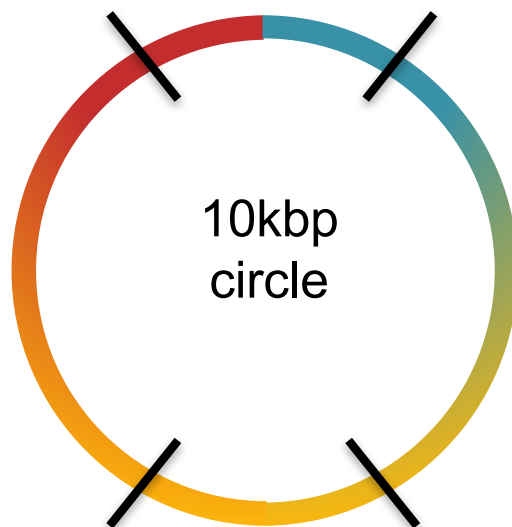
## ***Paired-end sequencing***

- 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



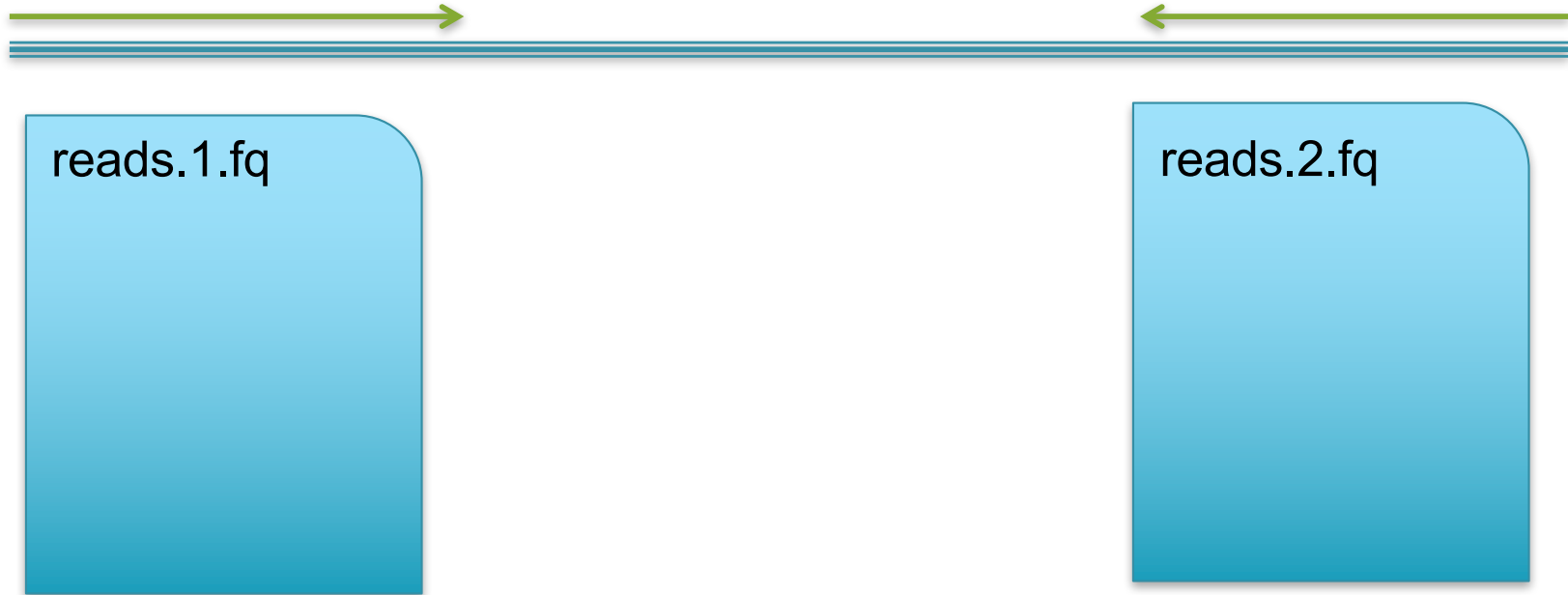
2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



# FASTQ Files

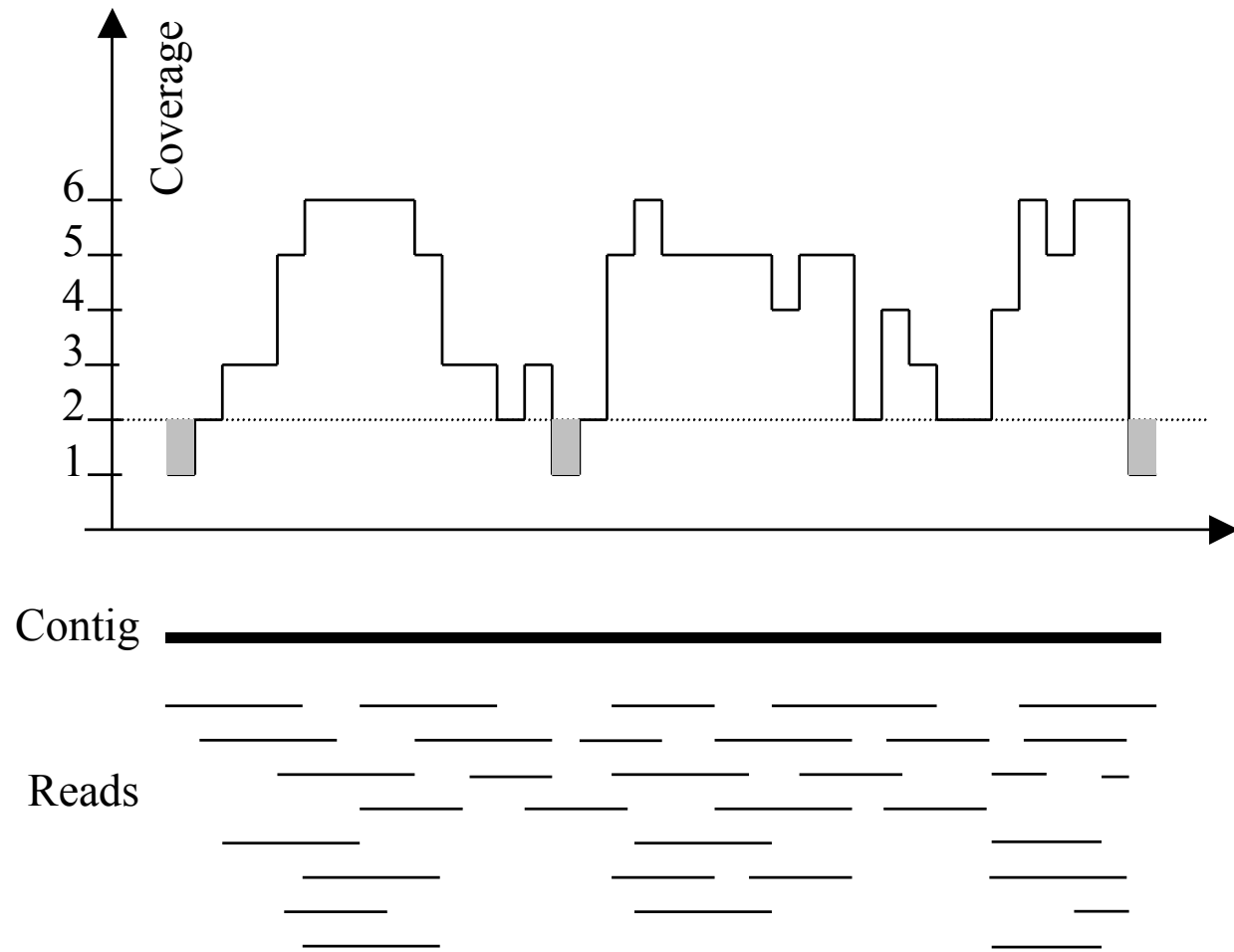


```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((( (**+))%+%+)(%+%%).1***-+*'' )**55CCF>>>>>CCCCCCC65
```

@Identifier  
Sequence  
+Separator  
Quality Values  
...



# Typical sequencing coverage

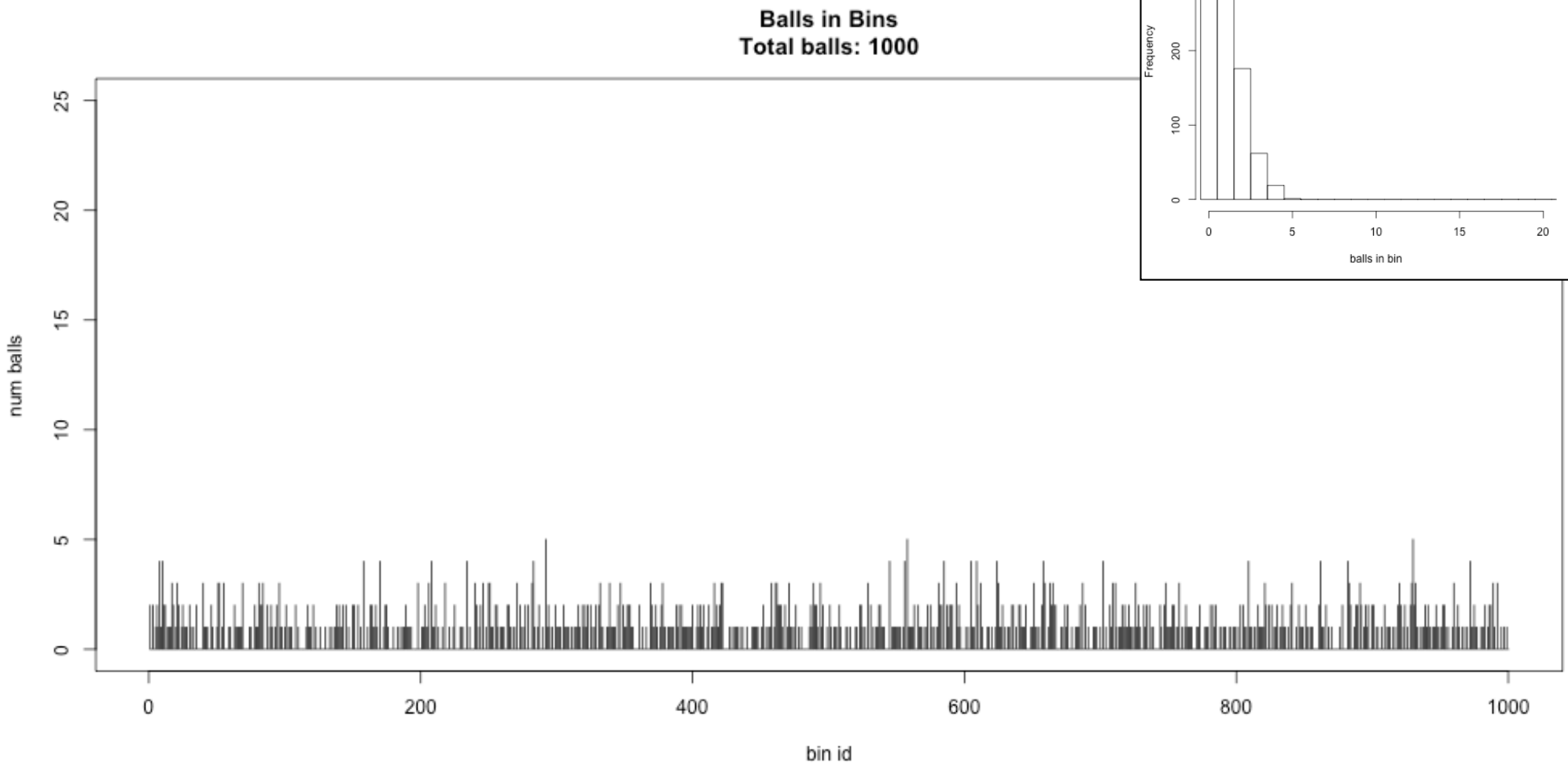


Imagine raindrops on a sidewalk

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

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

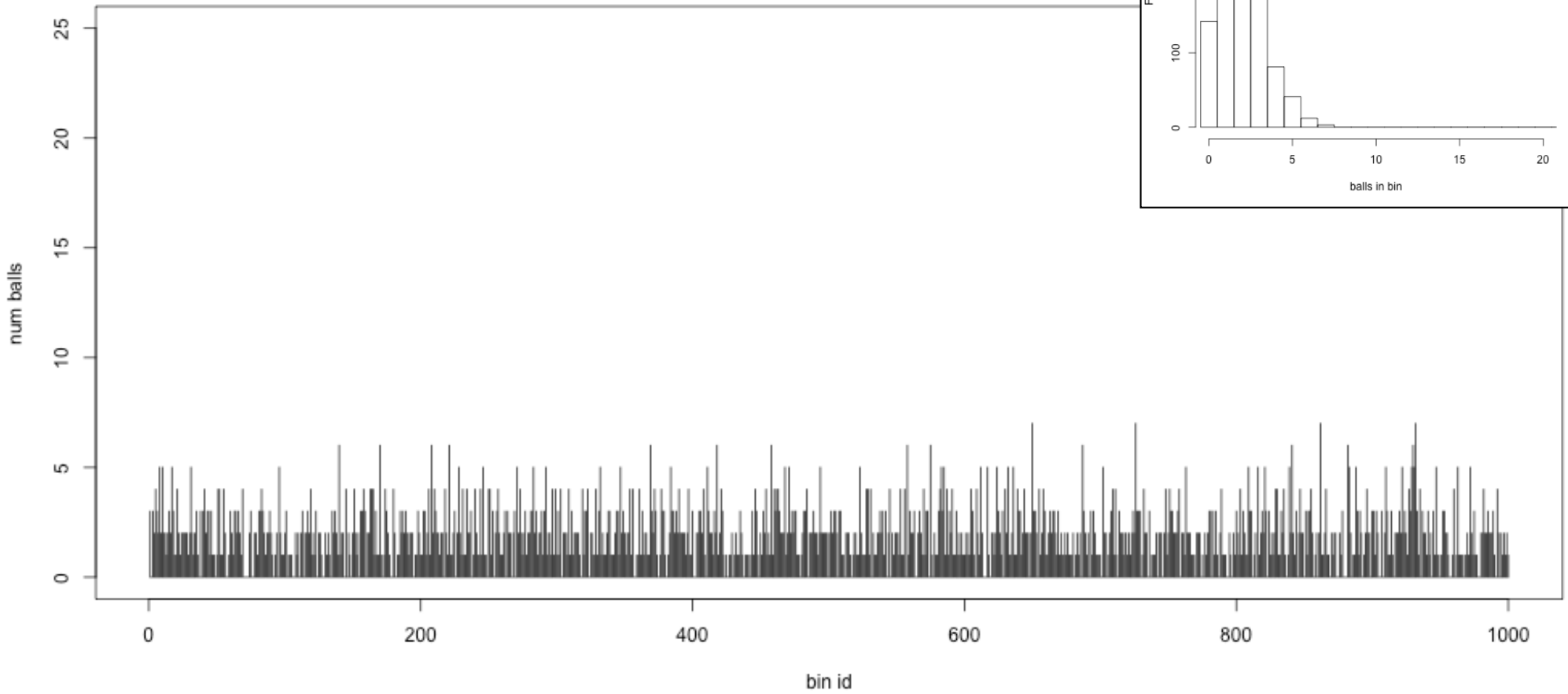
# Ix sequencing



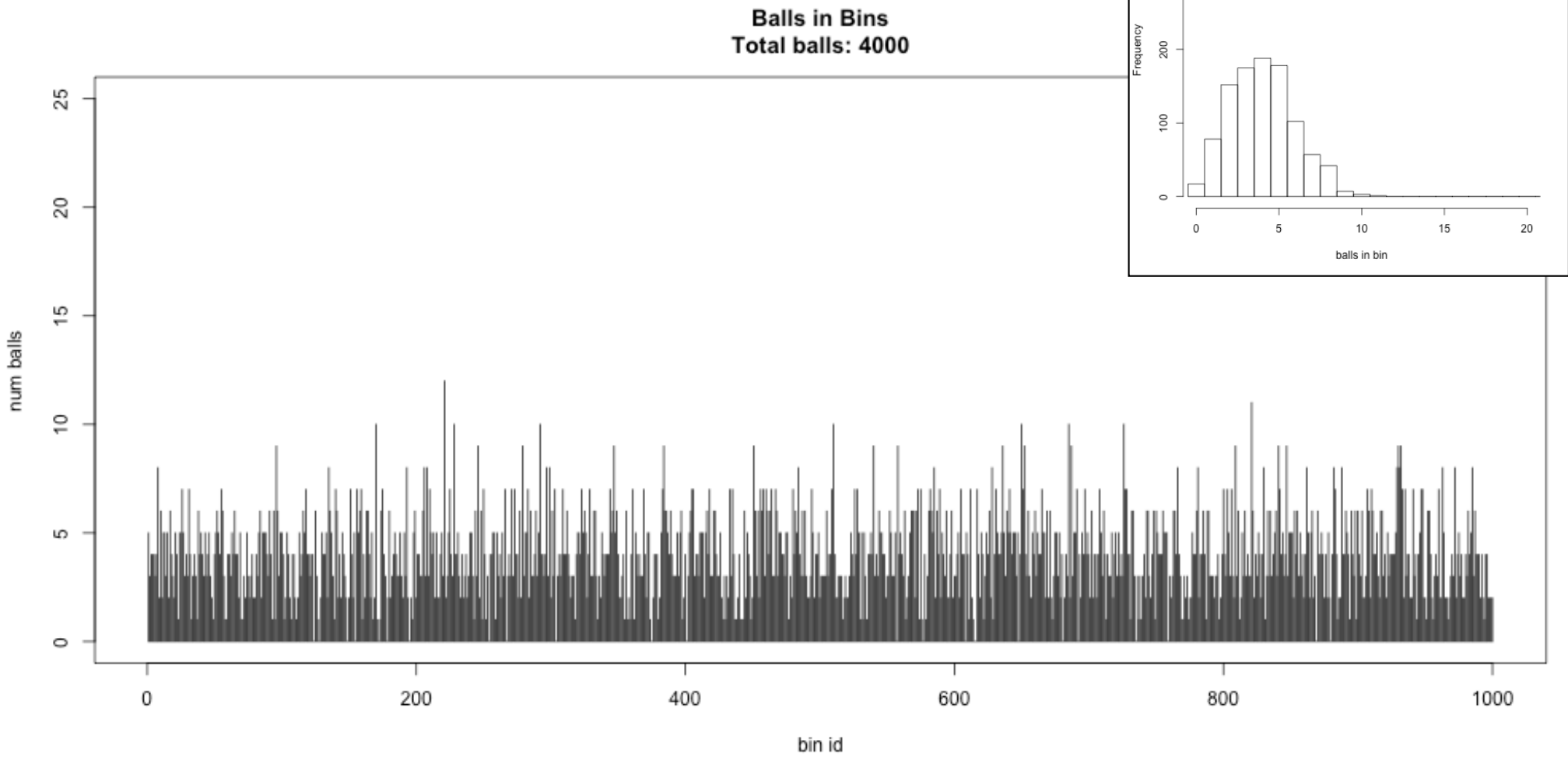


# 2x sequencing

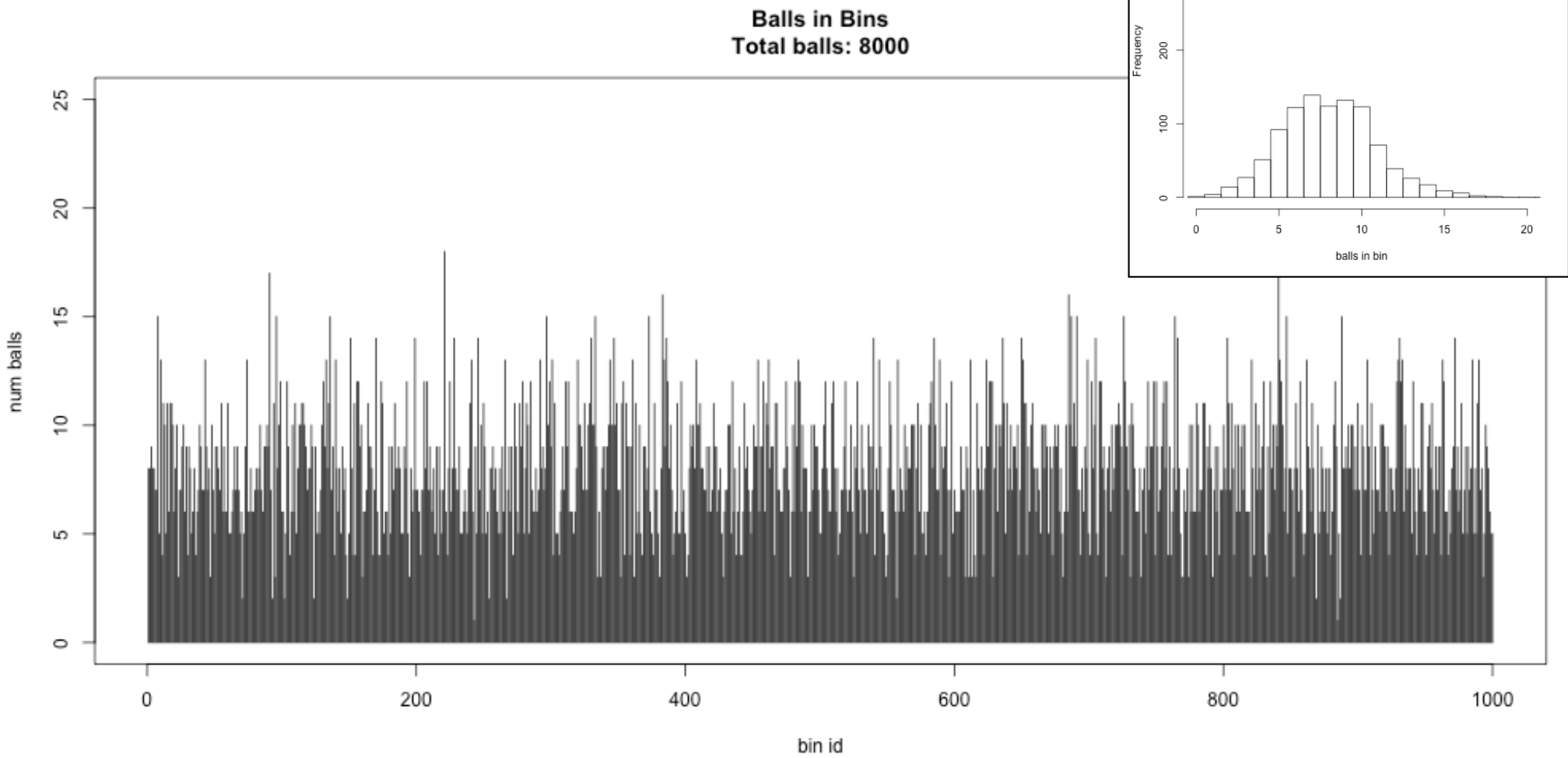
**Balls in Bins**  
Total balls: 2000



# 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.

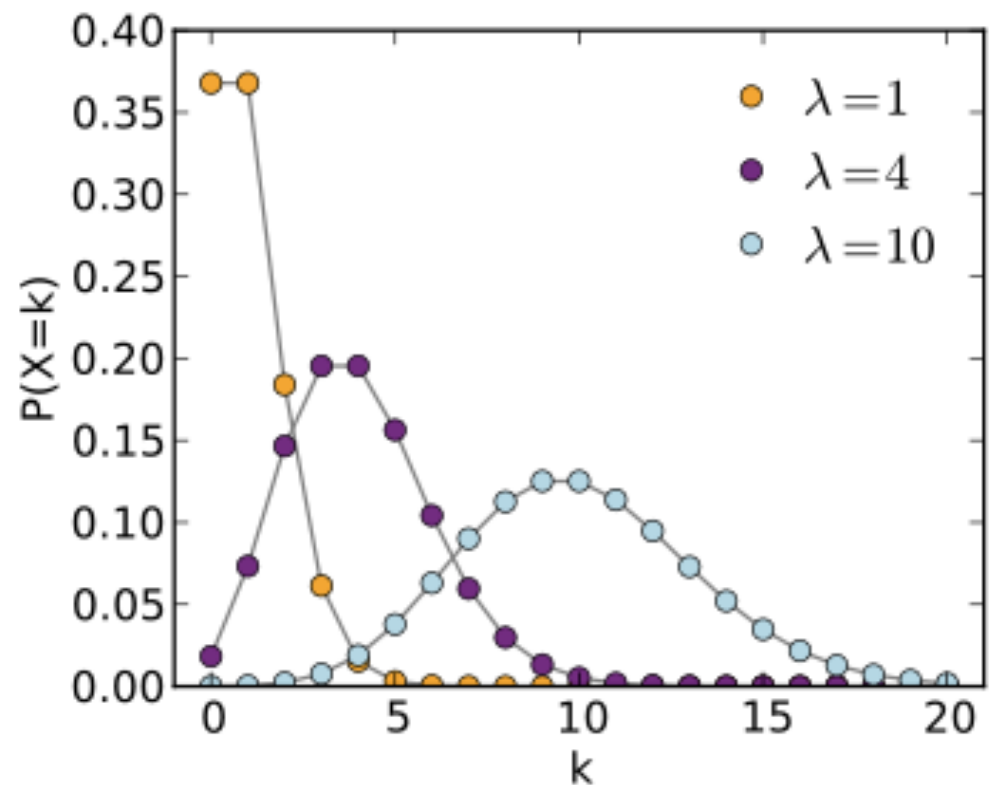
Formulation comes from the limit of the binomial equation

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

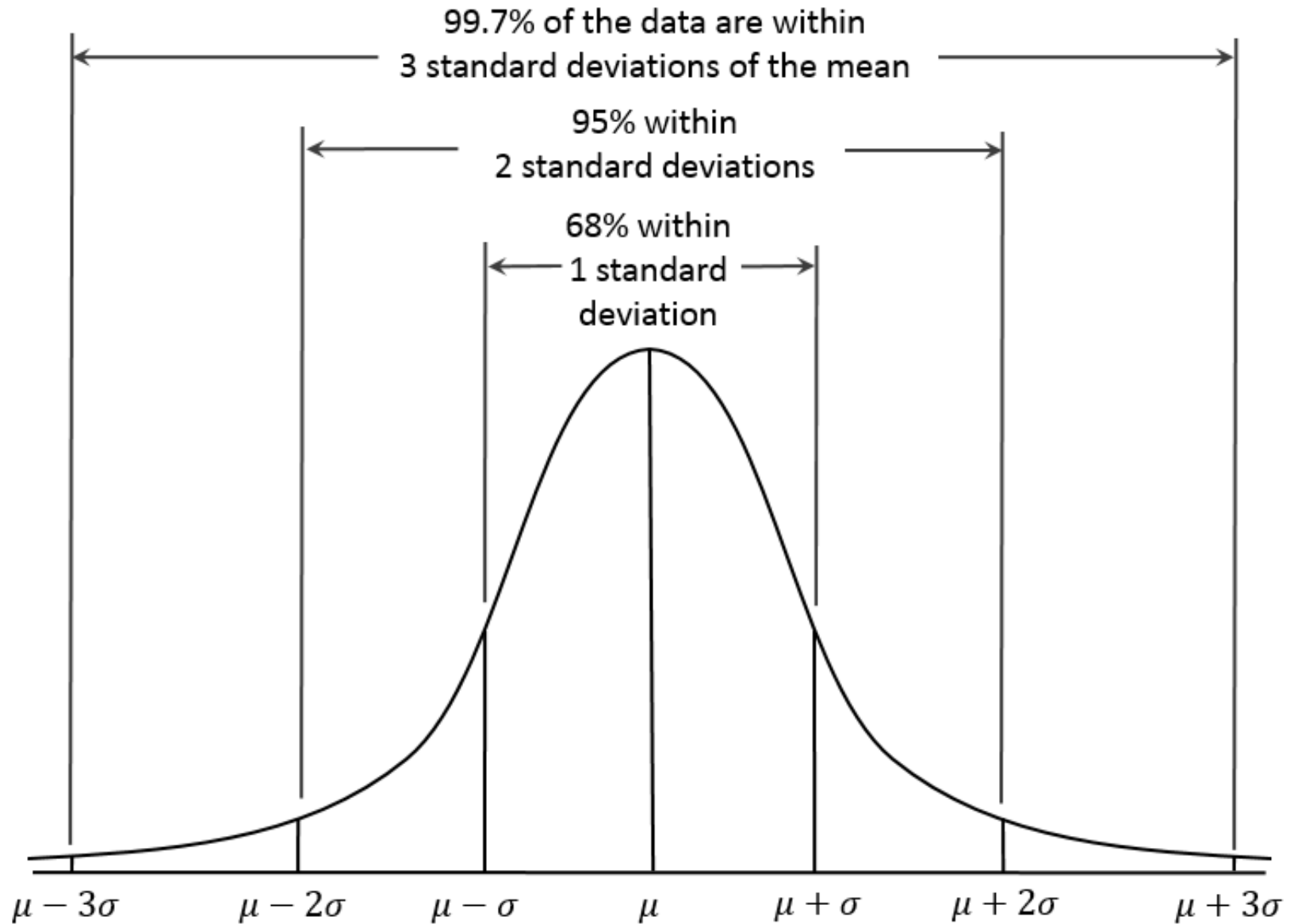
## **Key properties:**

- ***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  $10\text{Mbp} \times 24x = 240\text{Mbp}$  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  $X - 2 \cdot \sqrt{X} = 24$

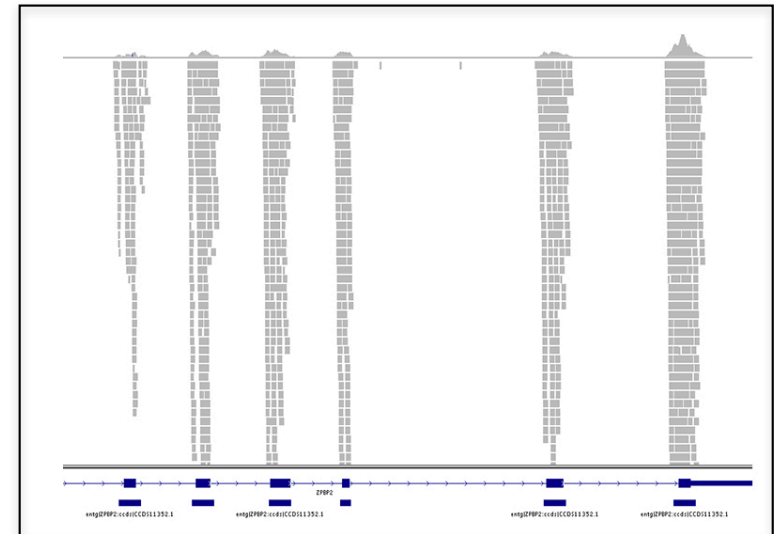
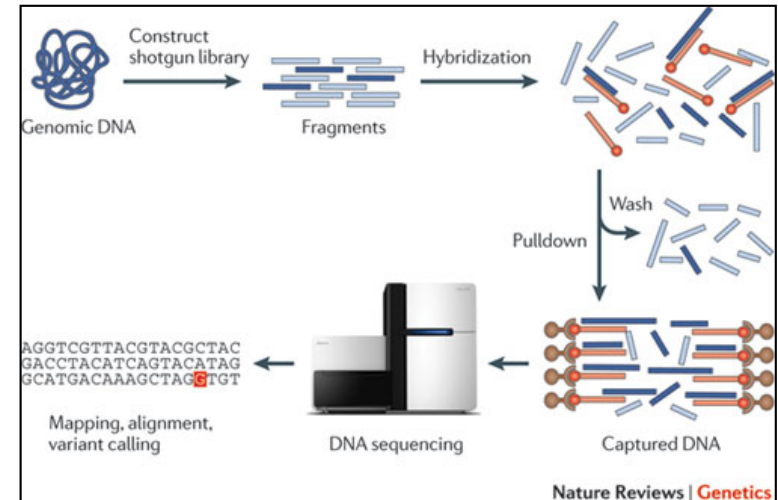
$$36 - 2 \cdot \sqrt{36} = 24$$

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

# Exome-Capture Sequencing

## Exome-capture reduces the costs of sequencing

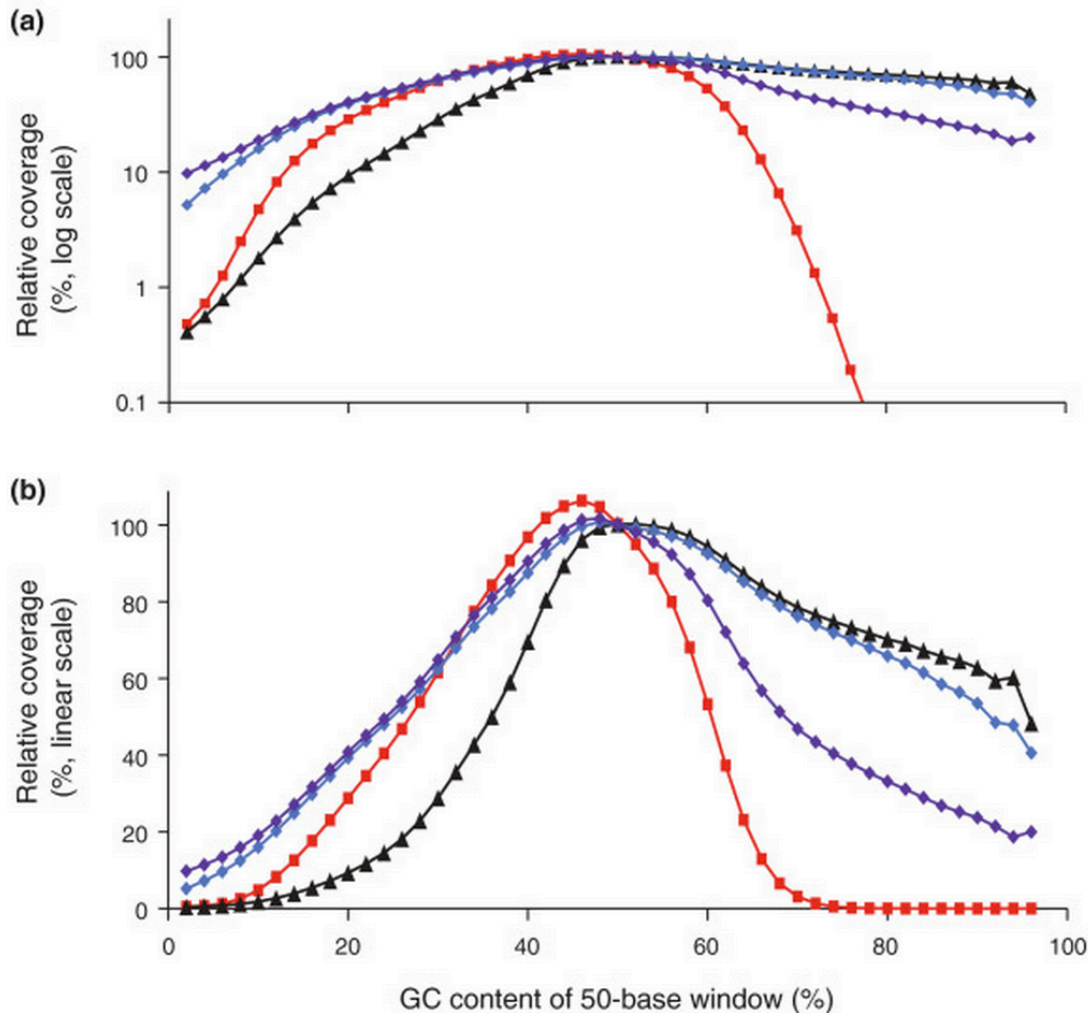
- Currently targets around 50Mbp of sequence: all exons plus flanking regions
- WGS currently costs ~\$1200 per sample, while WES currently costs ~\$300 per sample
- Coverage is highly localized around genes, although will get sparse coverage throughout rest of genome



## Exome sequencing as a tool for Mendelian disease gene discovery

Bamshad et al. (2011) *Nature Reviews Genetics*. 12, 745-755

# Beware of GC Biases



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

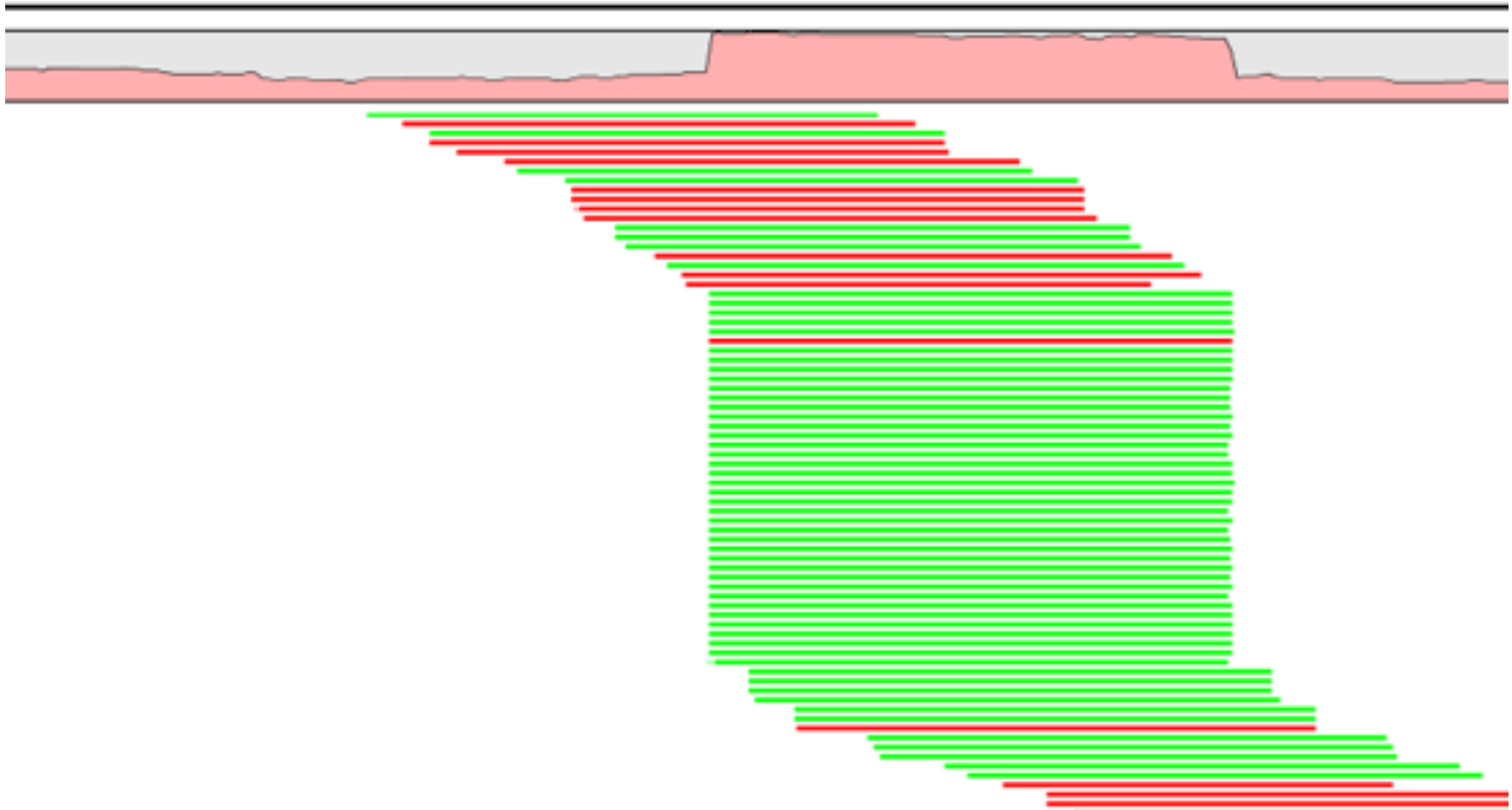
- Coverage of extremely high or extremely low GC content will have reduced coverage in Illumina sequencing
- Biases primarily introduced during PCR; lower temperatures, slower heating, and fewer rounds minimize biases
- This makes it very difficult to identify variants (SNPs, CNVs, etc) in certain regions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

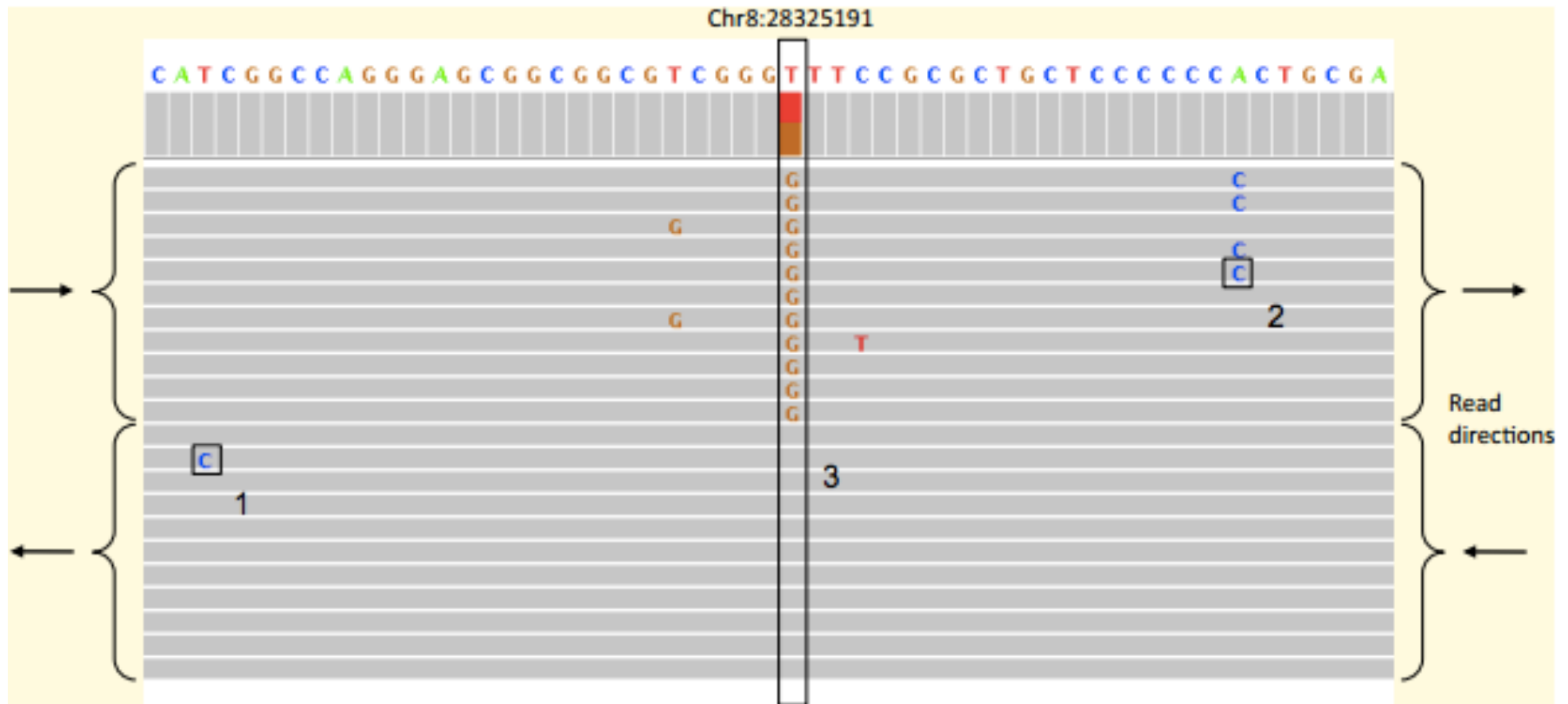


**The Sequence alignment/map (SAM) format and SAMtools.**

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

# Illumina Sequencing Summary

## Advantages:

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

## Disadvantages:

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



### Illumina HiSeq

~3 billion paired 100bp reads  
~600Gb, \$10K, 8 days  
(or “rapid run” ~90Gb in 1-2 days)

### Illumina X Ten

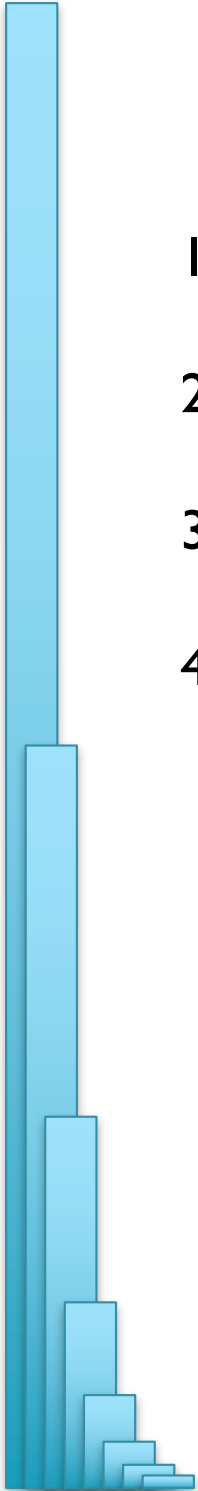
~6 billion paired 150bp reads  
1.8Tb, <3 days, ~1000 / genome(\$\$)  
(or “rapid run” ~90Gb in 1-2 days)

### Illumina NextSeq

One human genome in **<30 hours**

# 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
  2. Set up Dropbox for yourself!
  3. Get comfortable on the command line





***Welcome to Applied Comparative Genomics***  
<https://github.com/schatzlab/appliedgenomics2018>

**Questions?**