

Genomic Technologies

Michael Schatz

January 29, 2020

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/appliedgenomics2020>

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

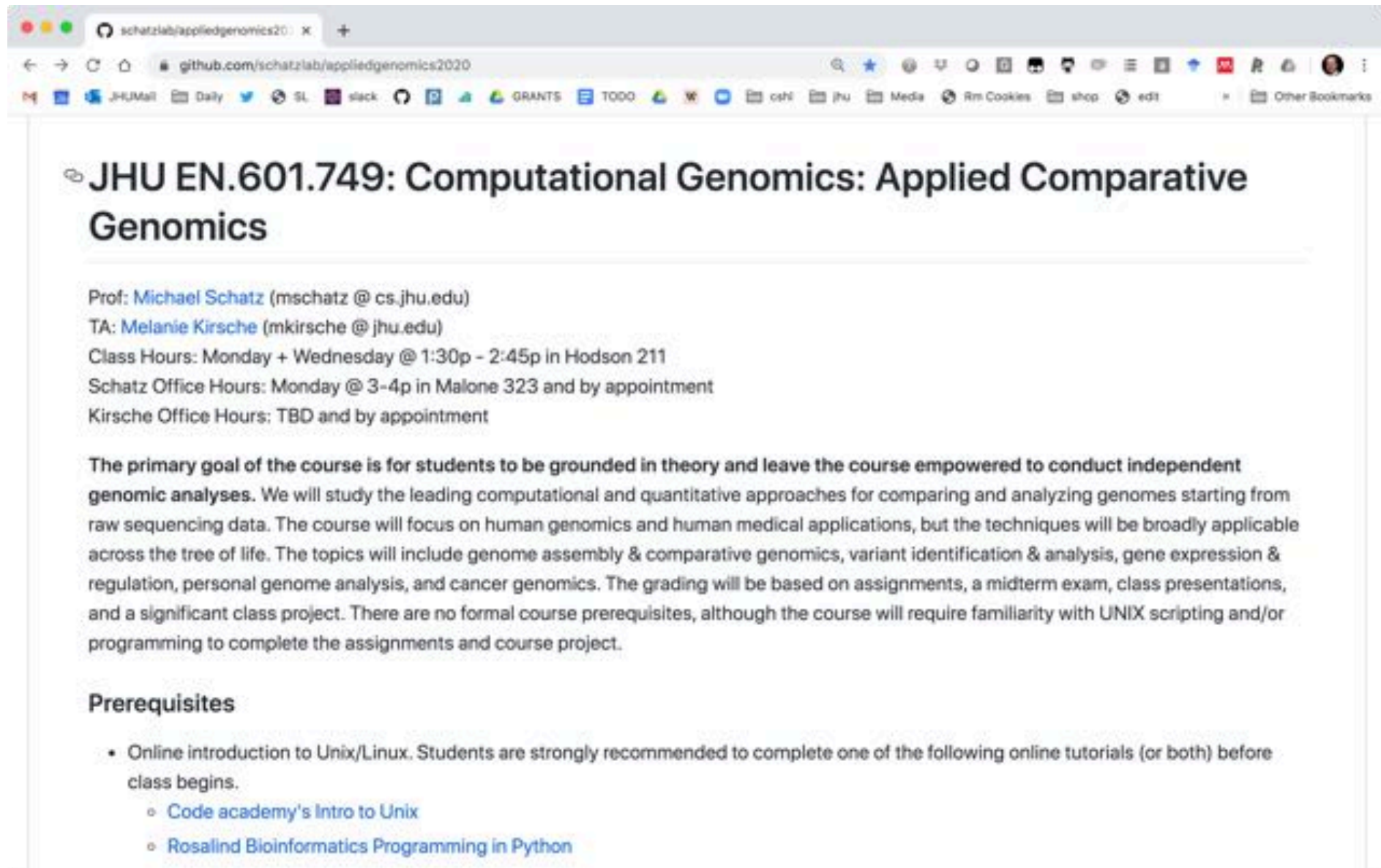
Class Hours: Mon + Wed @ 1:30p – 2:45p, Hodson 211

Schatz Office Hours: Mon @ 3-4p and by appointment

Kirsche Office Hours: TBD and by appointment

Please try Piazza first!

Course Webpage



schatzlab/appliedgenomics2020

github.com/schatzlab/appliedgenomics2020

JHU EN.601.749: Computational Genomics: Applied Comparative Genomics

Prof: [Michael Schatz](#) (mschatz@cs.jhu.edu)
TA: [Melanie Kirsche](#) (mkirsche@jhu.edu)
Class Hours: Monday + Wednesday @ 1:30p - 2:45p in Hodson 211
Schatz Office Hours: Monday @ 3-4p in Malone 323 and by appointment
Kirsche Office Hours: TBD 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.

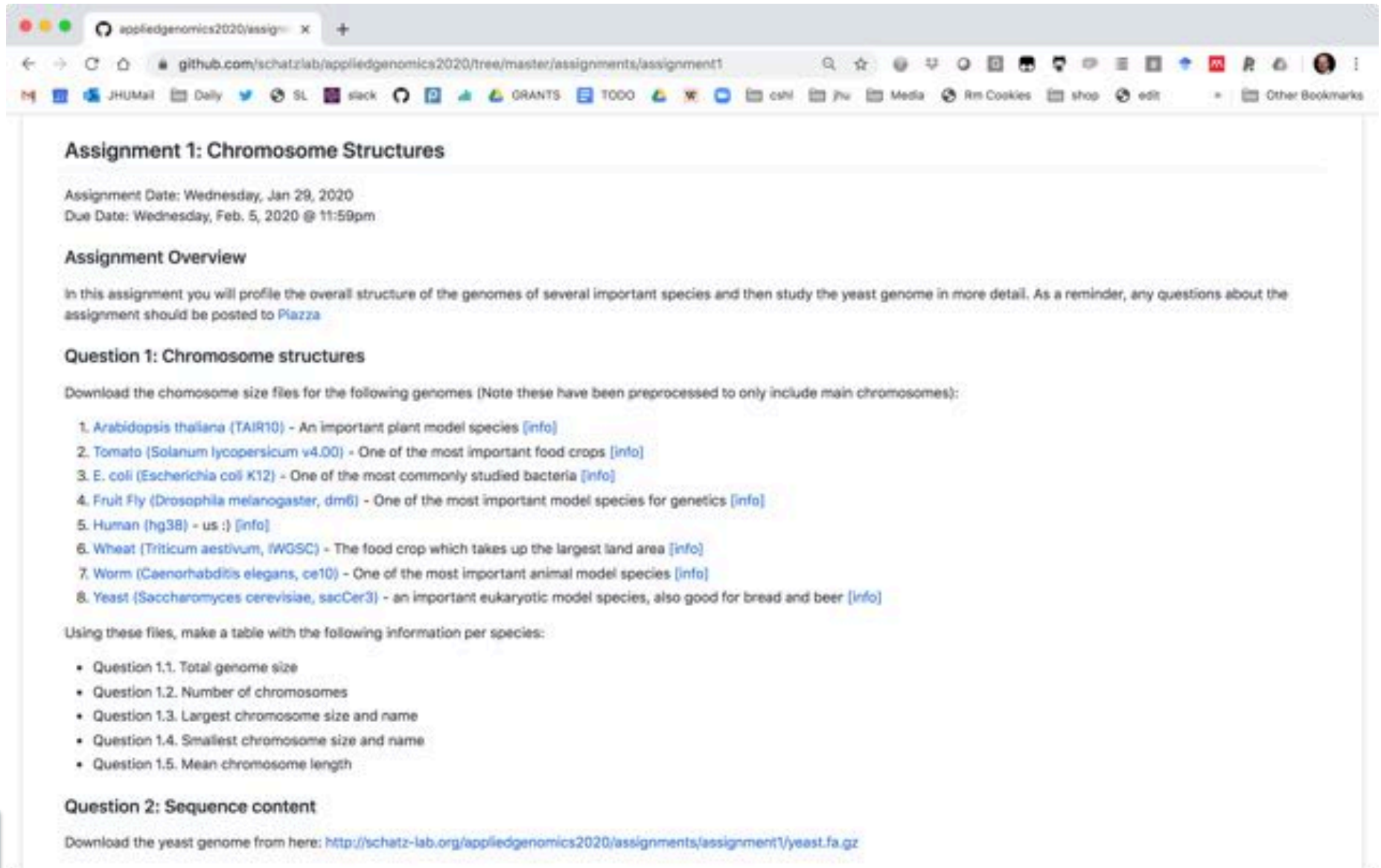
Prerequisites

- 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](#)
 - [Rosalind Bioinformatics Programming in Python](#)

<https://github.com/schatzlab/appliedgenomics2020>

Assignment I: Chromosome Structures

Due Feb 5 @ 11:59pm



The screenshot shows a web browser displaying a GitHub repository page. The browser's address bar shows the URL `github.com/schatzlab/appliedgenomics2020/tree/master/assignments/assignment1`. The page title is "Assignment 1: Chromosome Structures". Below the title, the assignment date is "Wednesday, Jan 29, 2020" and the due date is "Wednesday, Feb. 5, 2020 @ 11:59pm". The "Assignment Overview" section states that the assignment involves profiling the overall structure of the genomes of several important species and then studying the yeast genome in more detail. It also mentions that questions should be posted to Piazza. The "Question 1: Chromosome structures" section asks the user to download chromosome size files for eight genomes and create a table with specific information for each. The list of genomes includes Arabidopsis thaliana, Tomato, E. coli, Fruit Fly, Human, Wheat, Worm, and Yeast. The "Question 2: Sequence content" section asks the user to download the yeast genome from a specific URL.

Assignment 1: Chromosome Structures

Assignment Date: Wednesday, Jan 29, 2020
Due Date: Wednesday, Feb. 5, 2020 @ 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](#)

Question 1: Chromosome structures

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

1. *Arabidopsis thaliana* (TAIR10) - An important plant model species [\[info\]](#)
2. Tomato (*Solanum lycopersicum* v4.00) - One of the most important food crops [\[info\]](#)
3. *E. coli* (*Escherichia coli* K12) - One of the most commonly studied bacteria [\[info\]](#)
4. Fruit Fly (*Drosophila melanogaster*, dm6) - One of the most important model species for genetics [\[info\]](#)
5. Human (hg38) - us :) [\[info\]](#)
6. Wheat (*Triticum aestivum*, IWGSC) - The food crop which takes up the largest land area [\[info\]](#)
7. Worm (*Caenorhabditis elegans*, ce10) - One of the most important animal model species [\[info\]](#)
8. Yeast (*Saccharomyces cerevisiae*, sacCer3) - an important eukaryotic model species, also good for bread and beer [\[info\]](#)

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

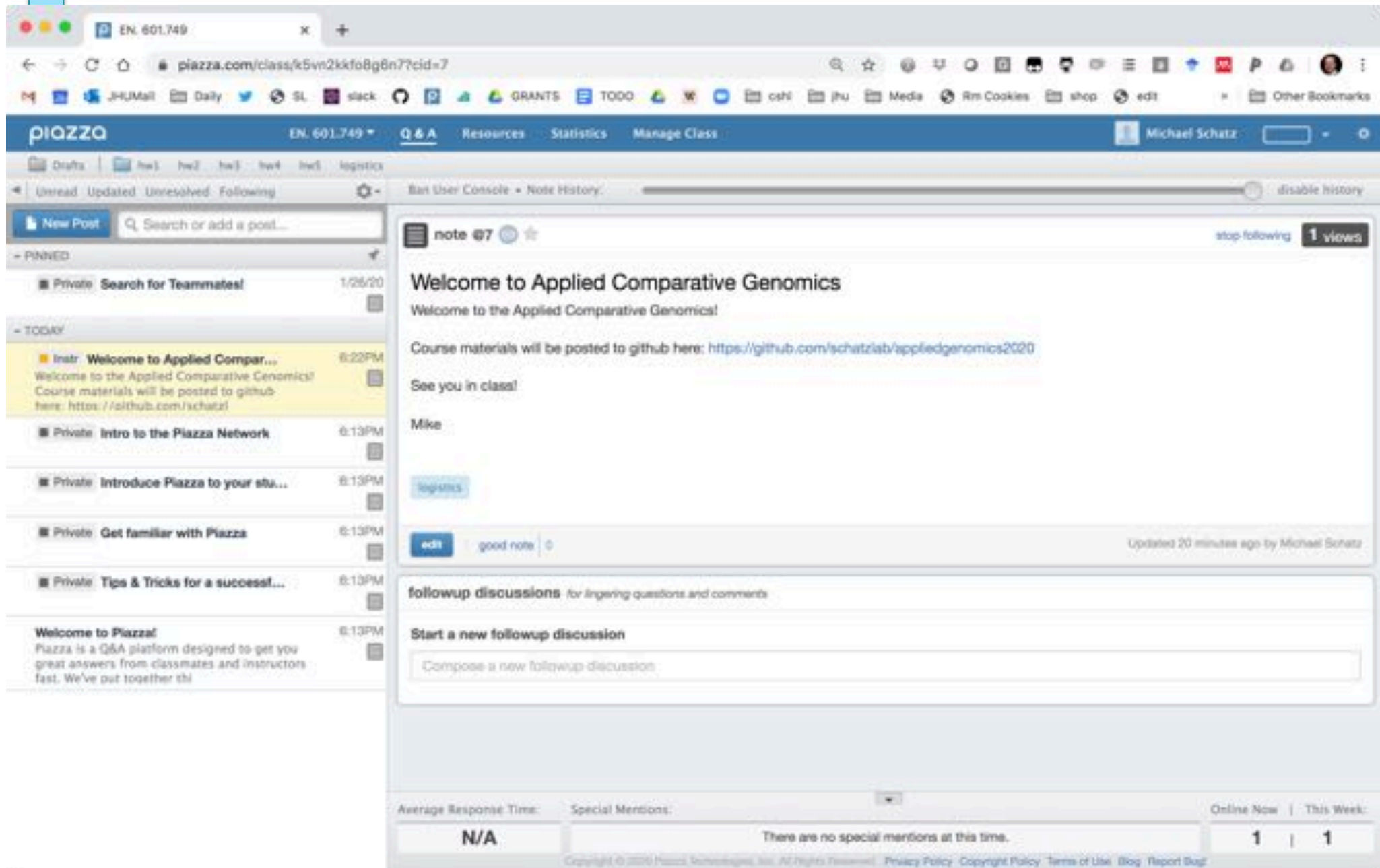
- Question 1.1. Total genome size
- Question 1.2. Number of chromosomes
- Question 1.3. Largest chromosome size and name
- Question 1.4. Smallest chromosome size and name
- Question 1.5. Mean chromosome length

Question 2: Sequence content

Download the yeast genome from here: <http://schatz-lab.org/appliedgenomics2020/assignments/assignment1/yeast.fa.gz>

<https://github.com/schatzlab/appliedgenomics2020>

Piazza



The screenshot shows the Piazza web application interface. The browser address bar displays the URL <https://piazza.com/class/k5vn2kkfo8g6n7>. The page header includes the Piazza logo, the class name 'EN. 601.749', and navigation links for 'Q & A', 'Resources', 'Statistics', and 'Manage Class'. The user profile 'Michael Schatz' is visible in the top right corner.

The main content area displays a post titled 'Welcome to Applied Comparative Genomics' by user 'note @7'. The post text reads: 'Welcome to the Applied Comparative Genomics! Course materials will be posted to github here: <https://github.com/schatzlab/appliedgenomics2020> See you in class! Mike'. The post is marked as 'logistics' and has '1 views'. It was updated 20 minutes ago by Michael Schatz.

Below the post, there is a section for 'followup discussions' with a prompt to 'Start a new followup discussion' and a text input field labeled 'Compose a new followup discussion'.

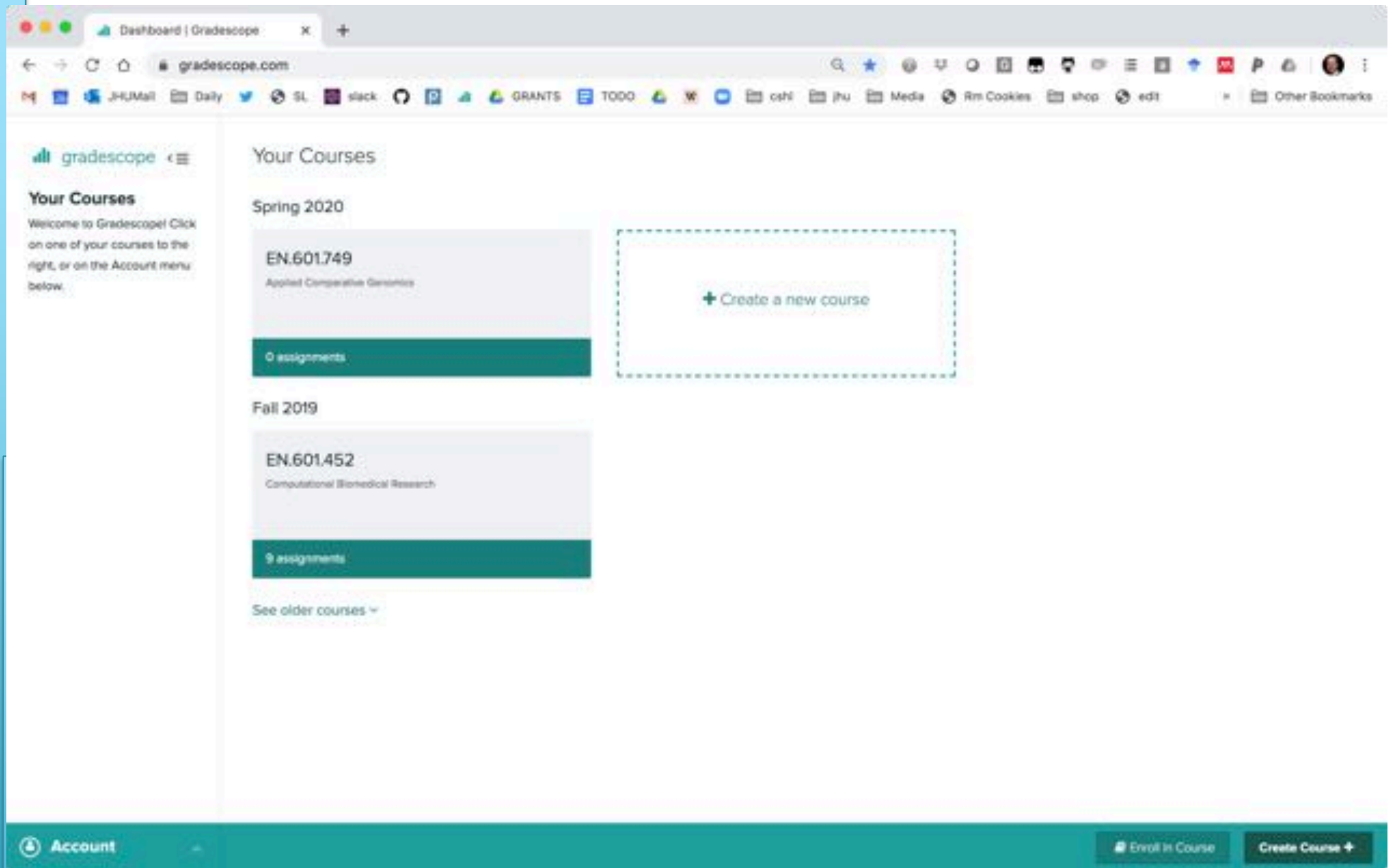
The bottom of the page shows a summary table:

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

At the very bottom, there is a copyright notice: 'Copyright © 2020 Piazza Technologies, Inc. All Rights Reserved. Privacy Policy Copyright Policy Terms of Use Blog Report Bug'.

<https://piazza.com/class/k5vn2kkfo8g6n7>

GradeScope



The screenshot displays the GradeScope dashboard in a web browser. The browser's address bar shows 'gradescope.com'. The dashboard is divided into a left sidebar and a main content area. The sidebar contains the 'gradescope' logo and a 'Your Courses' section with a welcome message. The main content area is titled 'Your Courses' and lists two courses: 'EN.601.749 Applied Comparative Genomics' for Spring 2020 and 'EN.601.452 Computational Biomedical Research' for Fall 2019. Each course card shows a green bar with the number of assignments (0 for Spring 2020, 9 for Fall 2019). A dashed box highlights a '+ Create a new course' button. At the bottom, a teal navigation bar includes an 'Account' link, an 'Enroll In Course' button, and a 'Create Course +' button.

Dashboard | Gradescope

gradescope.com

gradescope

Your Courses

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

Your Courses

Spring 2020

EN.601.749
Applied Comparative Genomics

0 assignments

+ Create a new course

Fall 2019

EN.601.452
Computational Biomedical Research

9 assignments

See older courses

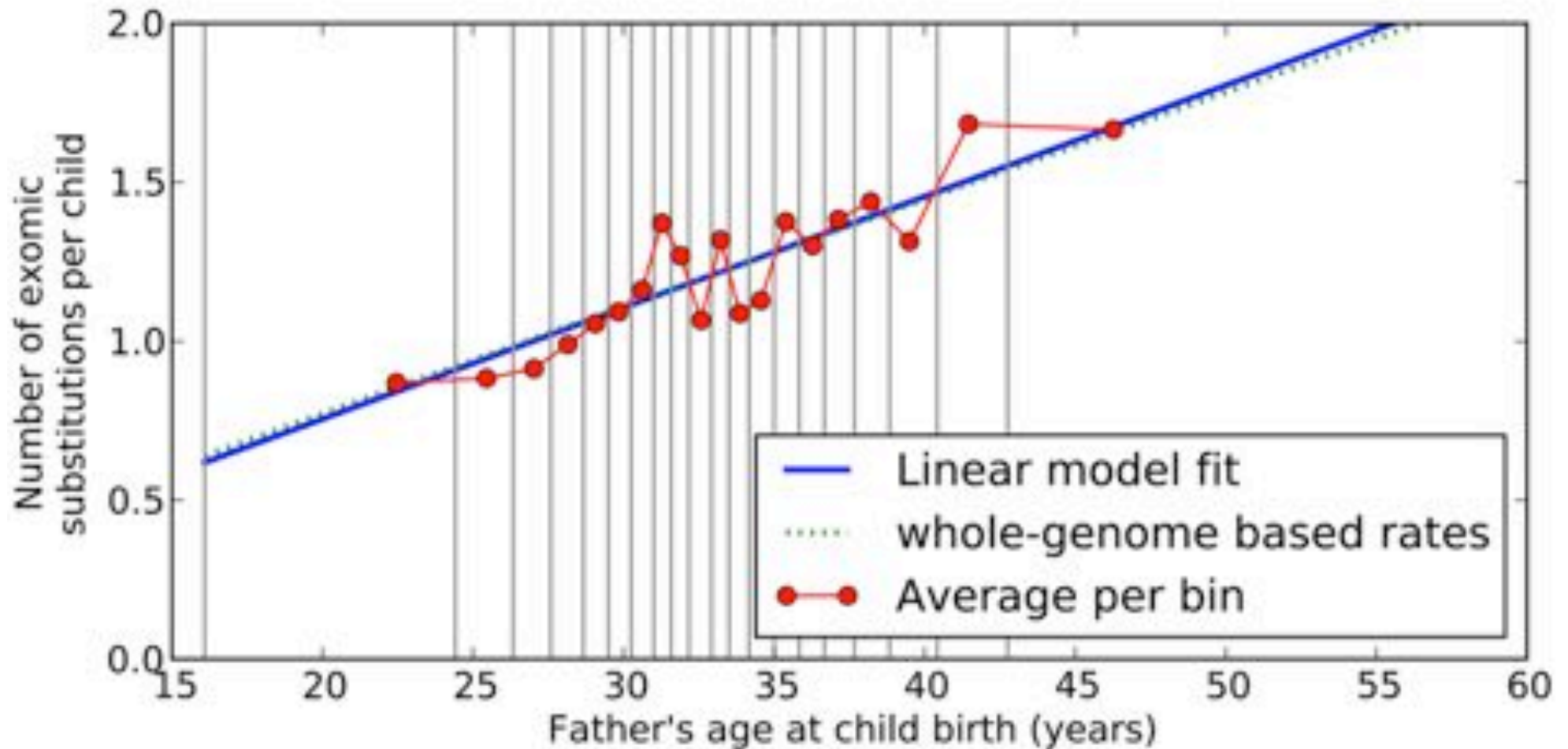
Account

Enroll In Course

Create Course +

<https://www.gradescope.com/>
Entry Code: MR652Z

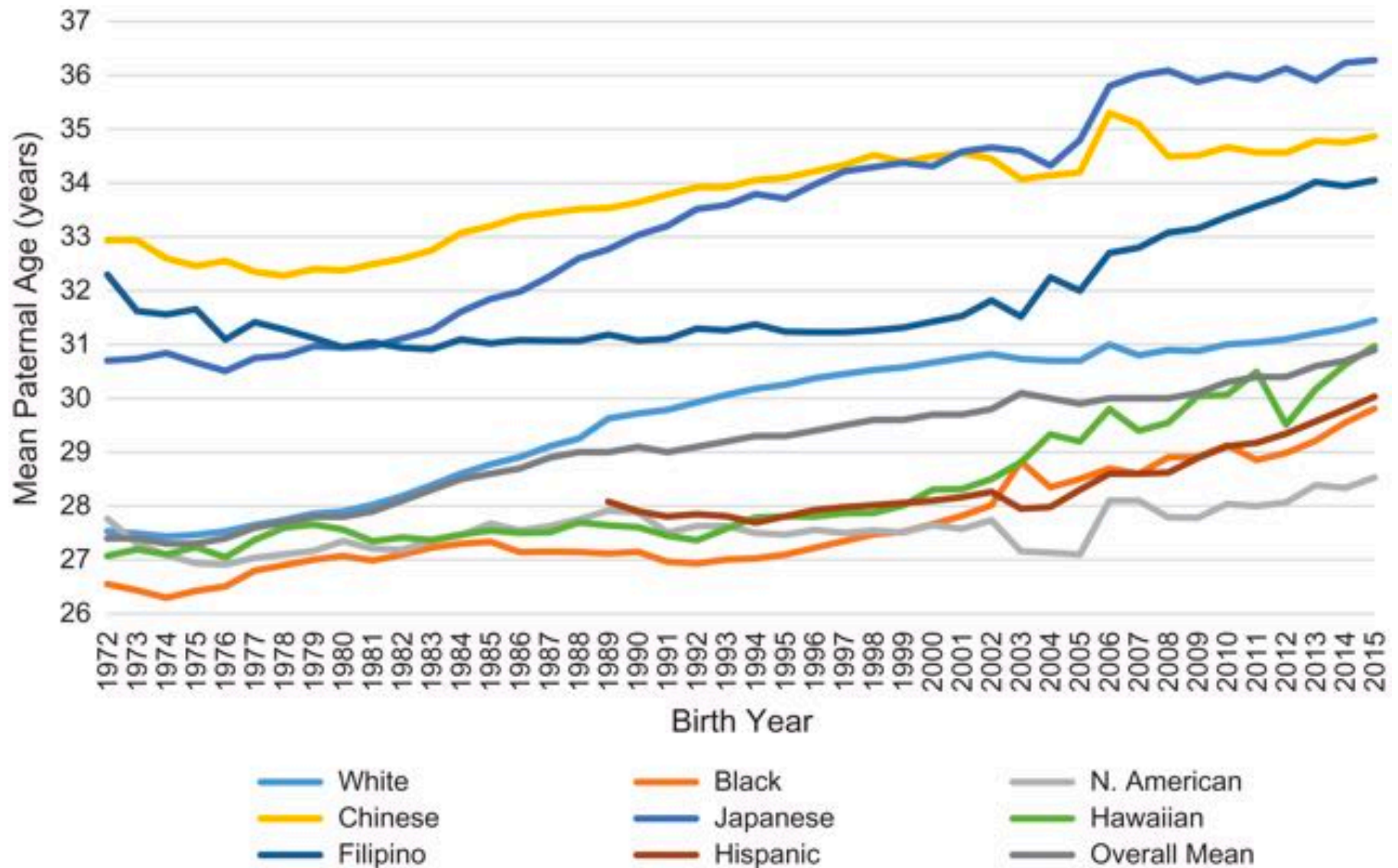
De novo Mutations in Men



The contribution of de novo coding mutations to autism spectrum disorder

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

Age of Fatherhood



The age of fathers in the USA is rising: an analysis of 168 867 480 births from 1972 to 2015

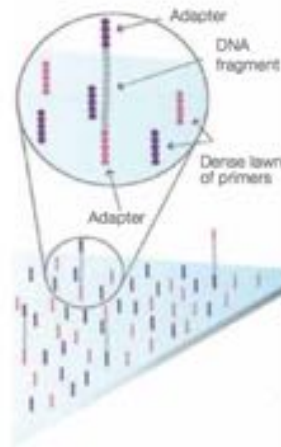
Khandwala et al (2017) *Human Reproduction*. <https://doi.org/10.1093/humrep/dex267>

Second Generation Sequencing

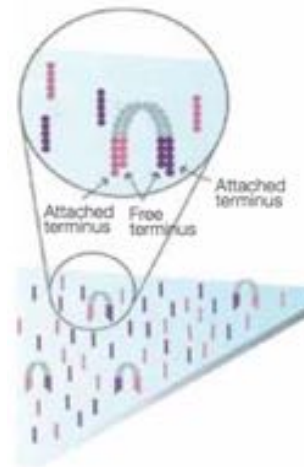


Illumina NovaSeq 6000
Sequencing by Synthesis

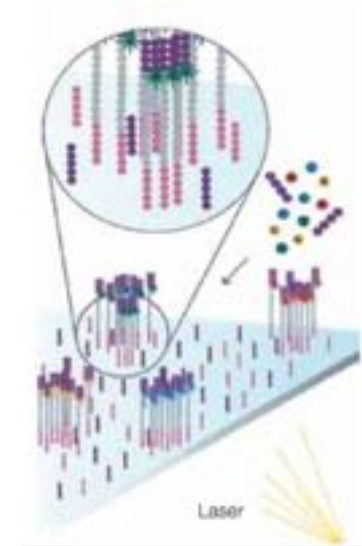
>3Tbp / day



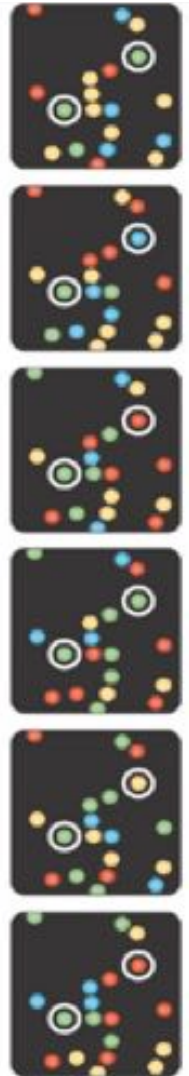
1. Attach



2. Amplify



3. Image

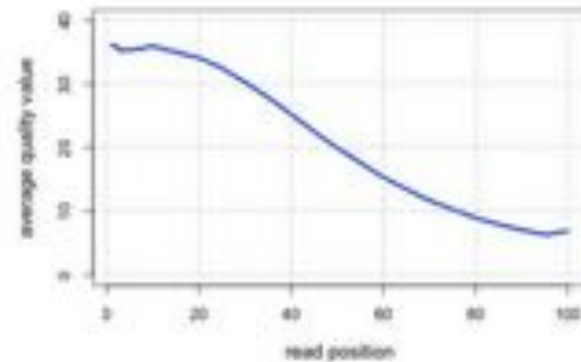


Metzker (2010) Nature Reviews Genetics 11:31-46
<https://www.youtube.com/watch?v=fCd6B5HRaZ8>

Illumina Quality

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

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



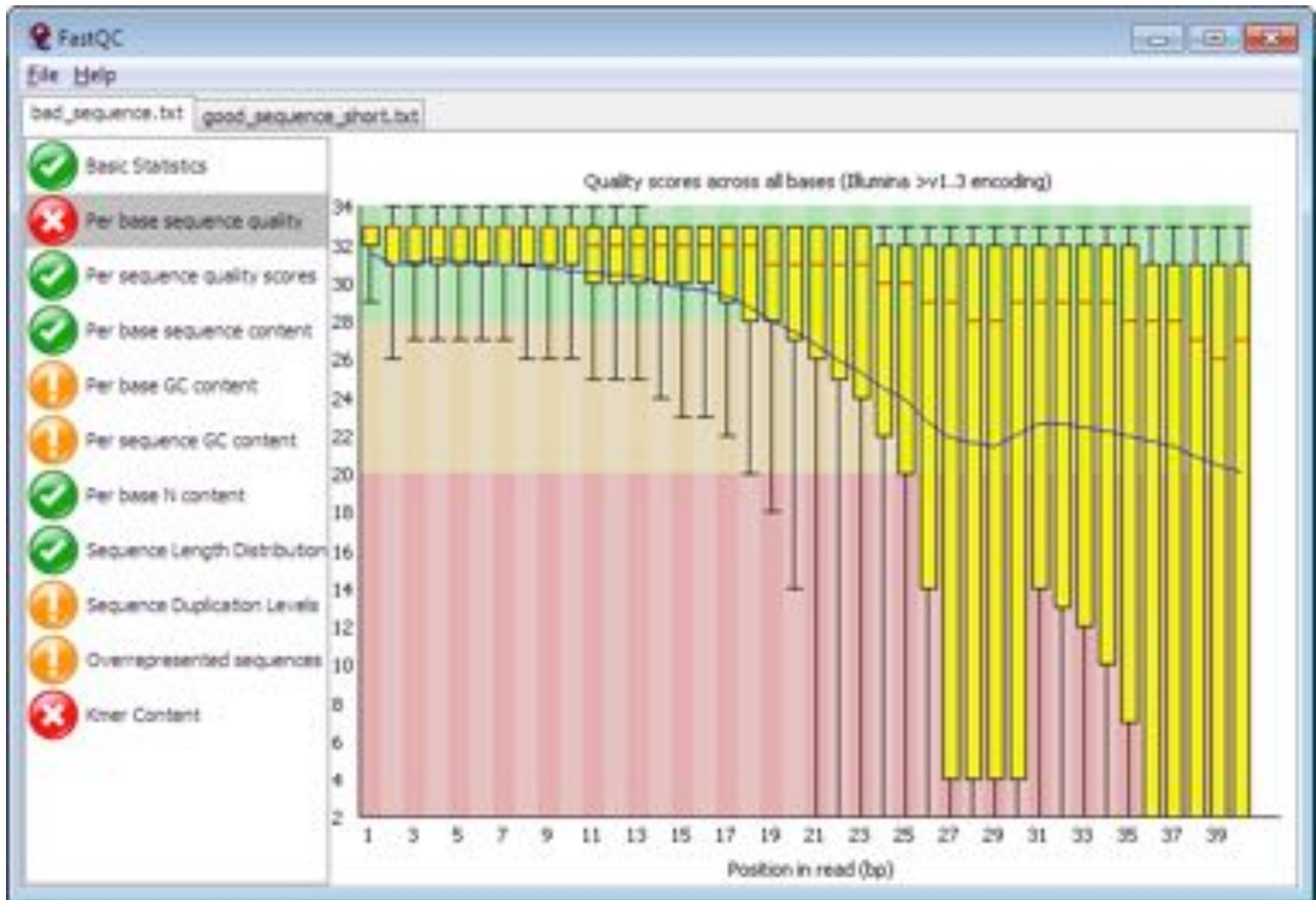
```

SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....
.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
!"#$%&'()*+,-./0123456789;:<=>?@ABCDEFGHIJKLMN
OPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|           |           |           |           |
33          59         64          73          104          126

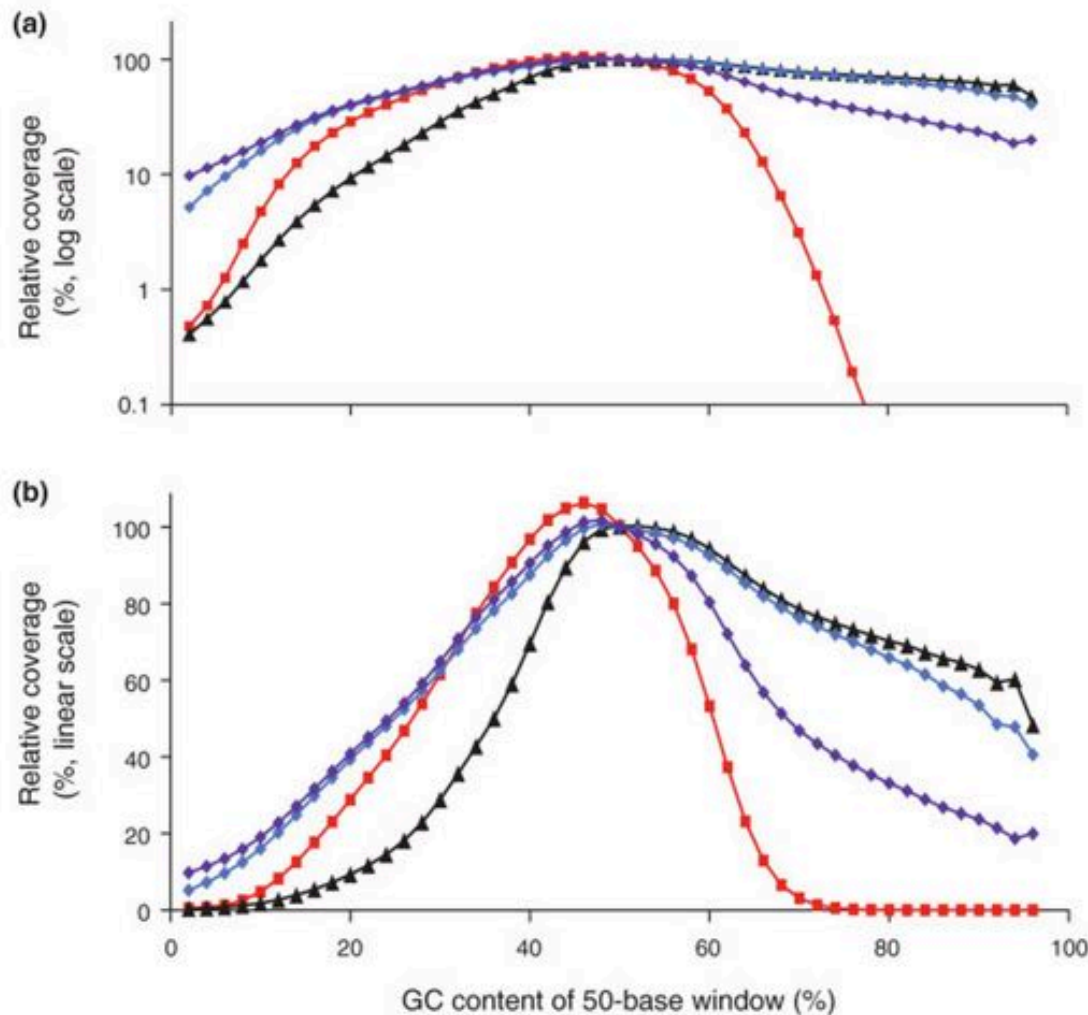
```

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)

FASTQC: Is my data any good?



Beware of GC Biases



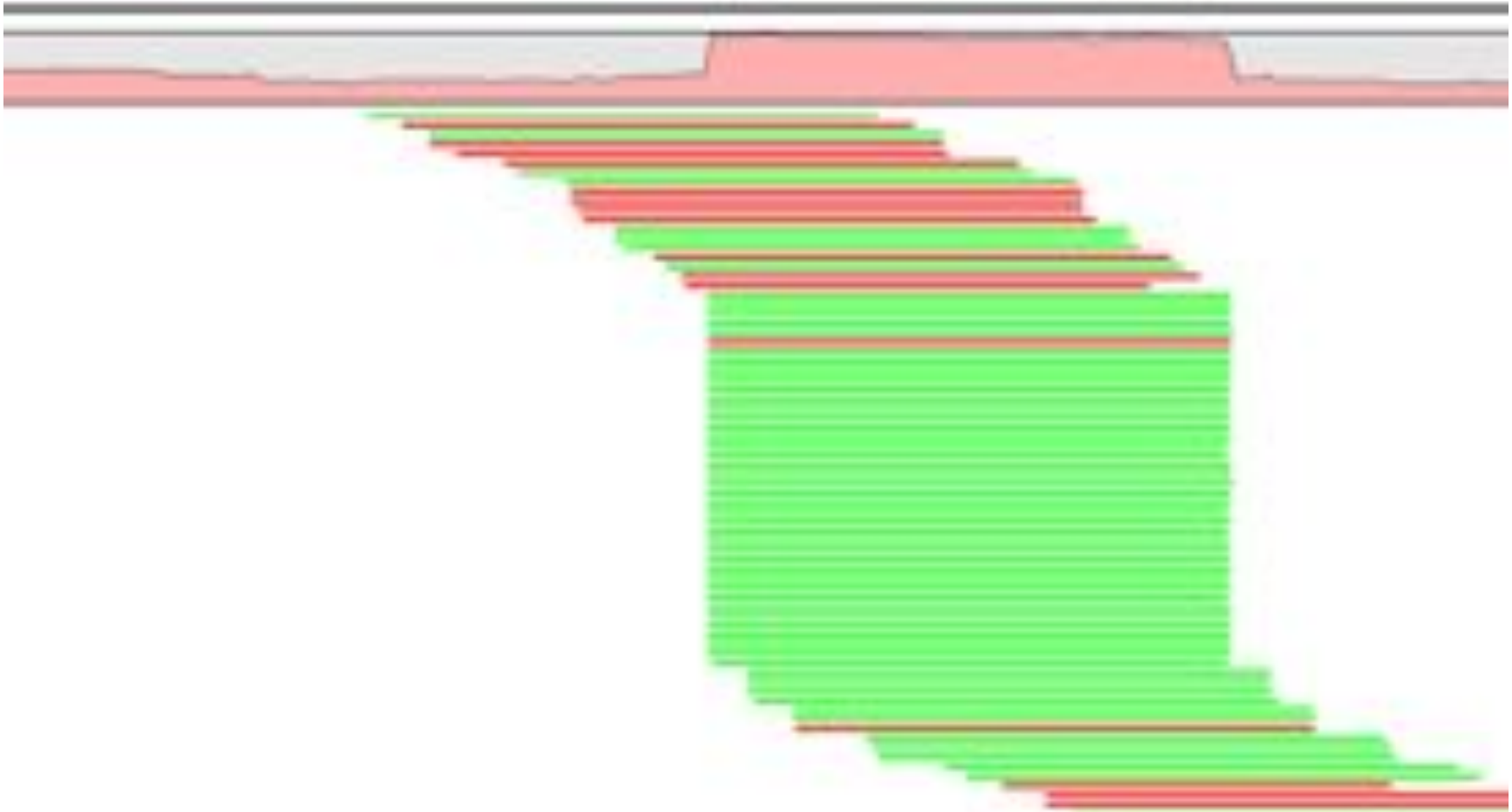
Illumina sequencing does not produce uniform coverage over the genome

- 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.

Beware of Duplicate Reads

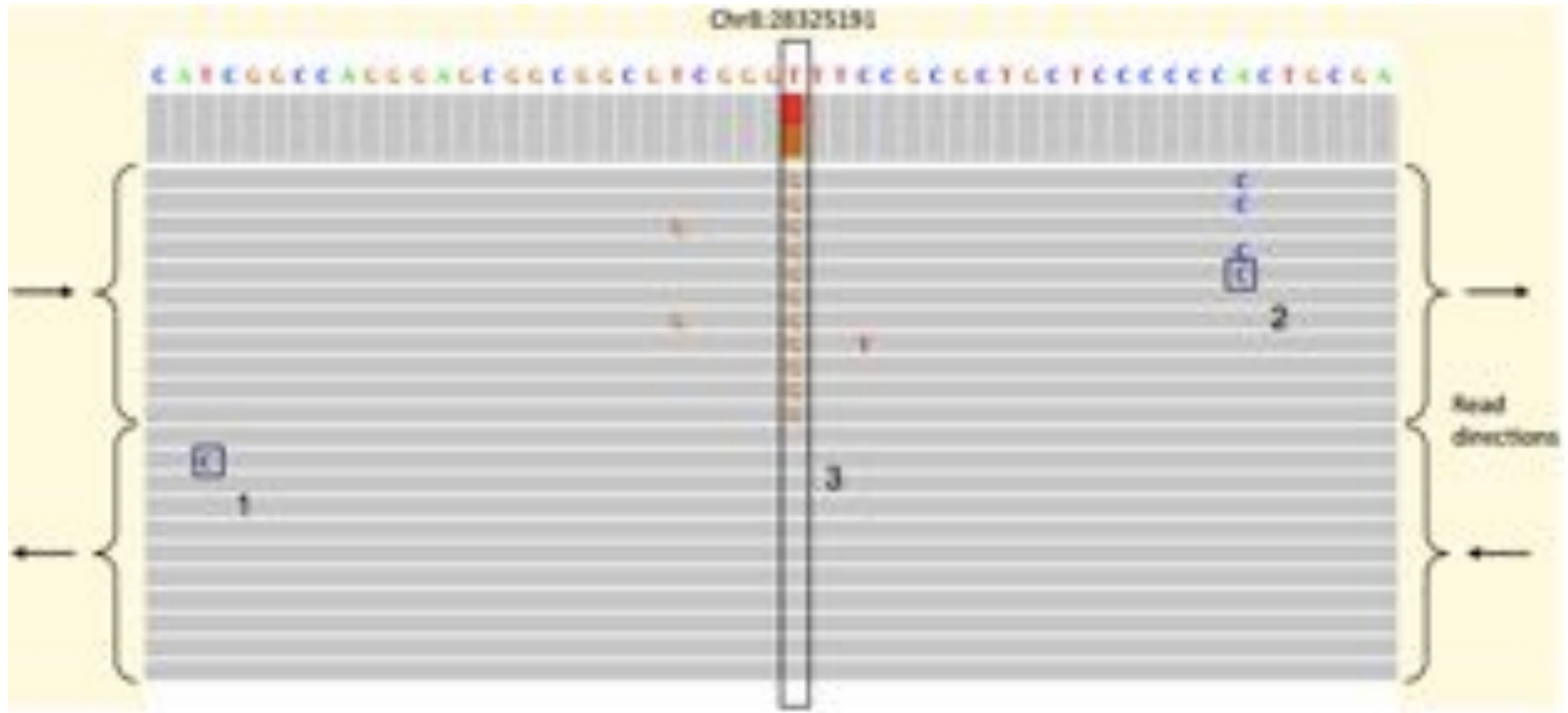


The Sequence alignment/map (SAM) format and SAMtools.

Li et al. (2009) *Bioinformatics*. 25:2078-9

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

Beware of (Systematic) Errors



Identification and correction of systematic error in high-throughput sequence data

Meacham et al. (2011) *BMC Bioinformatics*. 12:451

A closer look at RNA editing.

Lior Pachter (2012) *Nature Biotechnology*. 30:246-247

Question?

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

What can we do?

Illumina Hacking

BIOINFORMATICS ORIGINAL PAPER

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

Genome analysis

Advance Access publication May 22, 2013

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

Inanc Birol^{1,2,3,*}, Anthony Raymond¹, Shaun D. Jackman¹, Stephen Pleasance¹, Robin Coope¹, Greg A. Taylor¹, Macaire Man Saint Yuen¹, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Yongjun Zhao¹, Andrew J. Mungall¹, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland⁶, John MacKay^{7,8}, Jörg Bohlmann^{4,6}, Steven J.M. Jones^{1,2,9}

¹Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, BC V5Z 4S6, Canada, ²Genetics, University of British Columbia, Vancouver, BC V6H 3N1, Canada, ³School of Computer Science, University of British Columbia, Vancouver, BC V6T 1Z4, Canada, ⁴Michael Smith Laboratories, University of British Columbia, BC V6T 1Z4, Canada, ⁵British Columbia Ministry of Forests, Lands and Natural Resources, BC V8W 9C2, Canada, ⁶Department of Forest Sciences, University of British Columbia, BC V6T 1Z4, Canada, ⁷Institute for Systems and Integrative Biology, Université Laval, Québec, QC G1V 0A6, Canada, ⁸Department of Wood and Forest Sciences, Université Laval, Québec, QC G1V 0A6, Canada, ⁹Molecular 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 genomic 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 with 4.3 million scaffolds, with a scaffold N50 of 20366 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: ALX2010000000 PIR:PRUN4339], <http://www.ncbi.nlm.nih.gov/bioproject/83435>.

Contact: ibirol@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 re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact the Oxford University Press.

large genomes was demonstrated by (Simpson et al., 2009) using human and was later used to assemble the SOA/Novo tool (Li et al., 2010), high quality results, as demonstrated by (Ladner et al., 2011); Ritz has been successfully applied numerous times (Chen et al., 2011; Chu et al., 2011; Godel et al., 2012; Swart et al., 2011). Estimated at 20 Gb base pairs (GP) and assembly of the genome of the pine (*Pinus*) family present using generation end, those challenges include whole-genome shotgun sequencing data reduced representation resources due to the problem. On the basis of massive sequencing datasets is extremely high, memory usage, storage and 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 the 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 target 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 Kalamazoo Research Station of the British Columbia Ministry of Forests and Range, Vernon, British Columbia, Canada. Genomic DNA was extracted from 60 mg tissue by BioS&T (<http://www.biost.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 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% (v/v) LoTE buffer [1 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' adenine overhangs using the Klenow fragment of DNA polymerase I (5'-3' 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 Quanti-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 tagged 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 cloning. 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 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 mode. 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 Ctr 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 Nebula (Roche, Branford, CT), and fragment end-repair for tailing was performed with the GS Rapid Library preparation (Roche, Branford, CT). TruSeq Adapters (Illumina, CA) were ligated to the repaired/A-tailed ends. Biotinylated ends 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 Pacific Biosciences of the Institute for Systems and Integrative Genomics (University of British Columbia, Vancouver, BC V6T 1Z4). A single paired-end library was prepared as 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 at speed code 18, 6–10 kb fragments were extracted from GS-FLX library adapters were ligated to the repaired/A-tailed ends. GS-FLX sequencing using the titanium chemistry was according to manufacturer's instructions (454 Life Sciences, Branford, CT). Smaller sequencing method was used to do BAC sequencing data as previously described (Hamberger, Keeling et al., 2010).

2.3 MiSeq modification

In sequencing the spruce genome, we generated longer reads by modifying the MiSeq platform. The MiSeq uses a clamshell style (Supplementary Fig. S1A) to hold reagent tubes in an array by the MiSeq's supports. Most of the reagents are 5 length independent steps such as denaturation and cluster formation, the Scan, Cleavage and Incorporation mix summed at each cycle. Although the MiSeq allows any read specified in the control software, the reagent cartridge carries during the run without stopping it. Increasing the read length requires increasing the quantity of the length-dependent reagent. 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 Fig. S1B and S2), give the reagent tubes, yet allowing the cartridge to be put in without damage to its components (Supplementary Fig. S1C). The stock length-dependent reagent containers (40 ml, the stock length-dependent reagent containers) 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 (*Picea glauca*) genome from whole-genome shotgun sequencing data

Inanc Birol^{1,2,3,*}, Anthony Raymond¹, Shaun D. Jackman¹, Stephen Pleasance¹, Robin Coope¹, Greg A. Taylor¹, Macaire Man Saint Yuen¹, Christopher I. Keeling⁴, Dana Brand¹, Benjamin P. Vandervalk¹, Heather Kirk¹, Pawan Pandoh¹, Richard A. Moore¹, Yongjun Zhao¹, Andrew J. Mungall¹, Barry Jaquish⁵, Alvin Yanchuk⁶, Carol Ritland⁶, John MacKay^{7,8}, Jörg Bohlmann^{4,6}, Steven J.M. Jones^{1,2,9}

¹ British Columbia Cancer Agency, Genome Sciences Centre, Vancouver, BC V5Z 4S6

² University of British Columbia, Department of Medical Genetics, Vancouver, BC V6H 3N1

³ Simon Fraser University, School of Computing Science, Burnaby, BC V5A 1S6

⁴ University of British Columbia, Michael Smith Laboratories, Vancouver, BC V6T 1Z4

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

⁶ University of British Columbia, Department of Forest Sciences, Vancouver, BC V6T 1Z4

⁷ Université Laval, Institute for Systems and Integrative Biology, Québec, QC G1V 0A6

⁸ Université Laval, Department of Wood and Forest Sciences, Québec, QC G1V 0A6

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

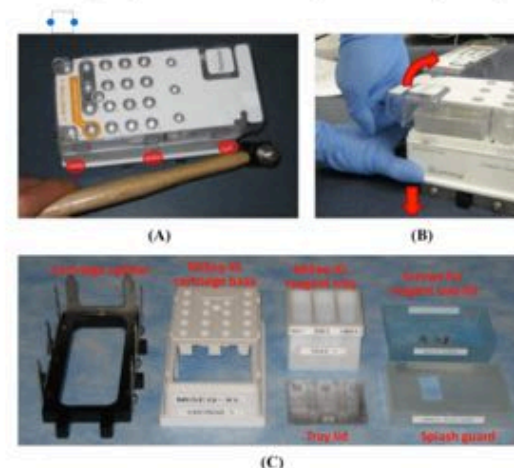


Figure S1. Modification of the MiSeq cartridge. MiSeq reagent cartridge was modified to allow for longer read lengths. (A, B) Opening of the clamshell style cartridge. (C) Contents of the modified cartridge. This was initially used to combine two PE150 kits for PE300 runs. When Illumina introduced the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

Paired-end and Mate-pairs

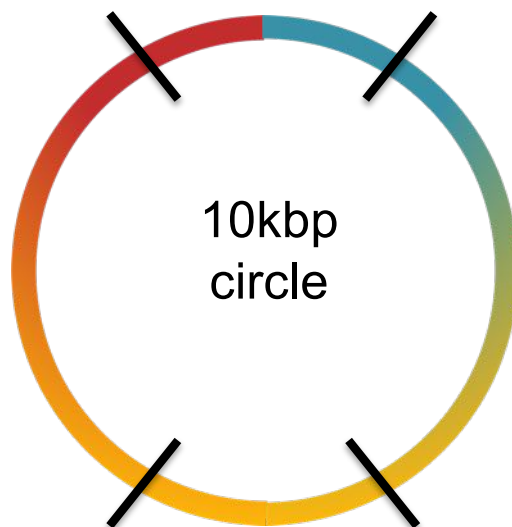
Paired-end sequencing

- 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



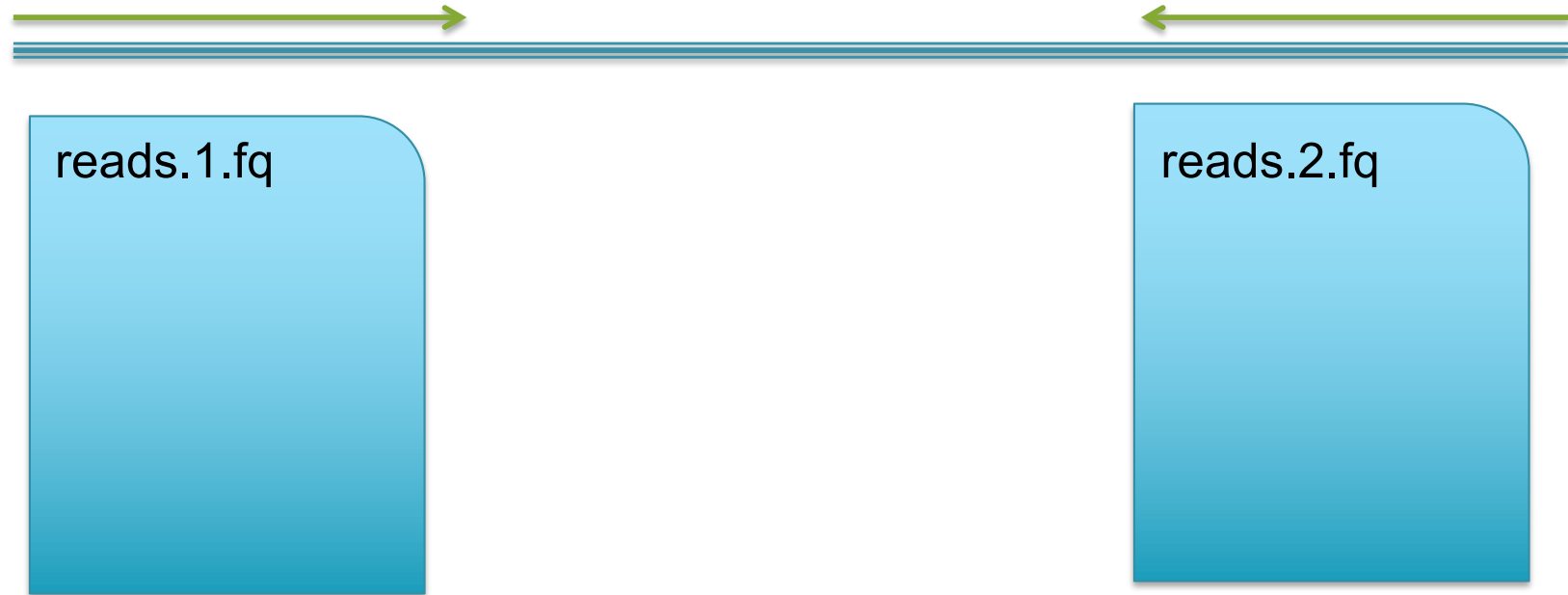
2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



FASTQ Files



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

@Identifier
Sequence
+Separator
Quality Values
...

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)

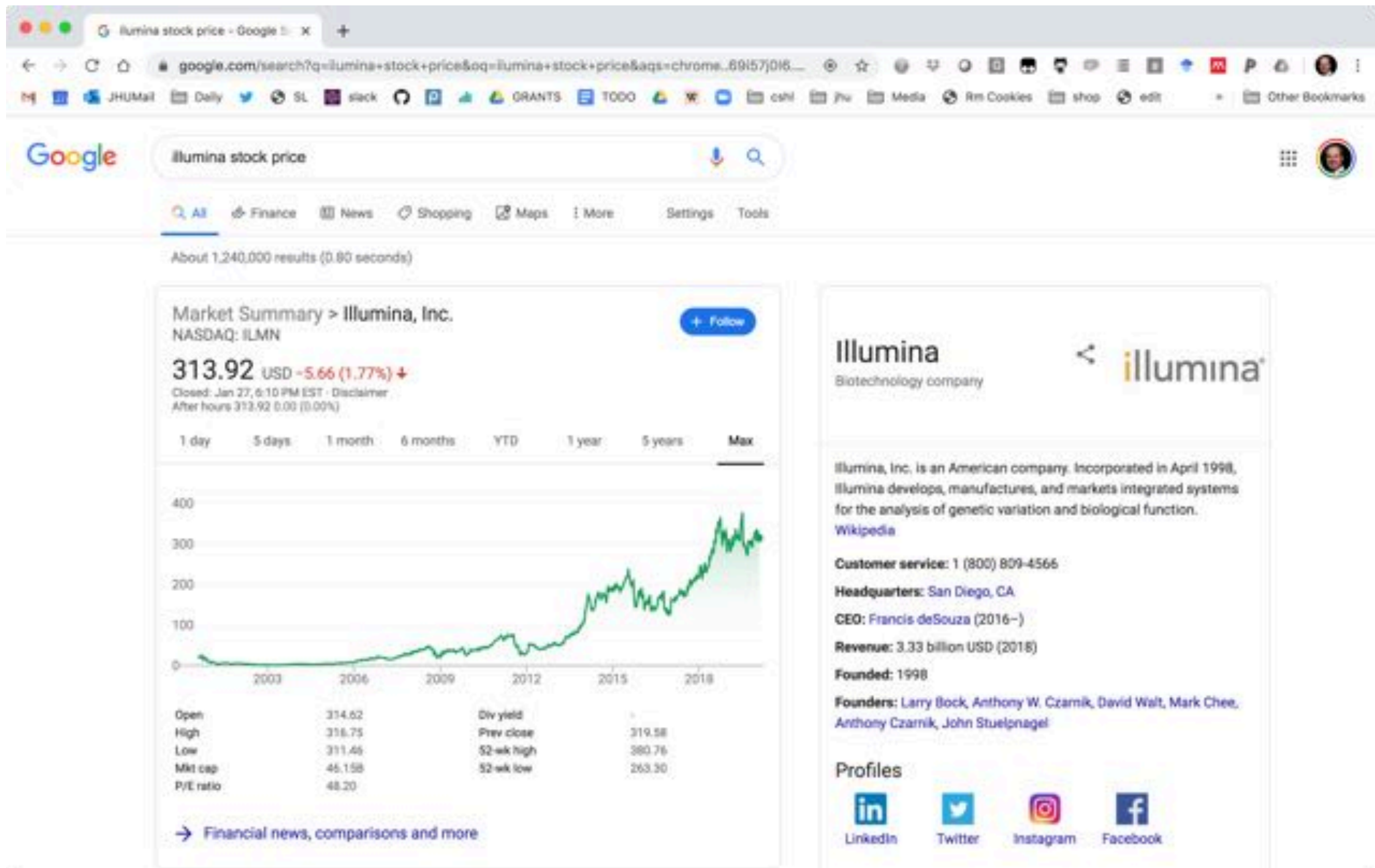
Illumina X Ten

~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**

ILMN





Next Steps

1. Reflect on the magic and power of DNA 😊
2. Check out the course webpage
3. Register on Piazza
4. Work on Assignment I
 1. Set up Linux, set up Virtual Machine, set up Ubuntu
 2. Set up Dropbox for yourself!
 3. Get comfortable on the command line