

Genomic Technologies

Michael Schatz

January 30, 2019

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/appliedgenomics2019>

Course Discussions:

<http://piazza.com>

Class Hours:

Mon + Wed @ 1:30p – 2:45p, Hodson 216

Schatz Office Hours:

Wed @ 3-4p and by appointment

Kovaka Office Hours:

Wed @ 4-5p and by appointment

Please try Piazza first!

Course Webpage

The screenshot shows a web browser window with the following details:

- Title Bar:** schatzlab/appliedgenomics2019
- Address Bar:** GitHub, Inc. [US] | https://github.com/schatzlab/appliedgenomics2019
- Content Area:**
 - ## JHU EN.601.749: Computational Genomics: Applied Comparative Genomics
 - Prof: Michael Schatz (mschatz @ cs.jhu.edu)
TA: Sam Kovaka (skovaka1 @ jhu.edu)
Class Hours: Monday + Wednesday @ 1:30p - 2:45p in Hodson 216
Schatz Office Hours: Wednesday @ 3-4p in Malone 323 and by appointment
Kovaka 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
 - Command line bootcamp
 - Rosalind Bioinformatics Programming in Python
 - Minimal Make
 - Access to a Linux Machine, and/or Install VirtualBox (Unfortuantely, even Mac will not work correctly for some programs)
 - ### Course Resources:

 - Syllabus and Policies
 - Piazza Discussion Board

<https://github.com/schatzlab/appliedgenomics2019>

Assignment I: Chromosome Structures

Due Feb 6 @ 11:59pm

The screenshot shows a web browser window with the following details:

- Title Bar:** appliedgenomics2019/README
- Address Bar:** GitHub, Inc. [US] | https://github.com/schatzlab/appliedgenomics2019/blob/master/assignments/assignm...
- Toolbar:** Back, Forward, Home, Search, Favorites, etc.
- Bookmark Bar:** JHUMail, Daily, SL, P, cshl, jhu, Media, shop, edit, Rm Cookies, Remove NYT Coo..., Other Bookmarks

Content Area:

Assignment 1: Chromosome Structures

Assignment Date: Wednesday, Jan 30, 2019
Due Date: Wednesday, Feb. 6, 2019 @ 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. [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 :\)](#) [\[info\]](#)
6. [Rice \(Oryza sativa, IRGSP-1.0\)](#) - One of the most important crops in the world [\[info\]](#)
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

Piazza


A screenshot of a web browser displaying the Piazza platform for a class titled "EN. 601.749". The URL in the address bar is <https://piazza.com/class/jrlz3buim75?cid=6>. The interface includes a navigation bar with tabs for "Q & A", "Resources", "Statistics", and "Manage Class". On the left, there's a sidebar with links for "hw1" through "hw5" and "logistics". Below this is a list of posts under "LAST WEEK", with one post from "Instr Welcome" highlighted in yellow. The main content area shows a "note" titled "Welcome" by Michael Schatz, which reads: "Welcome to Applied Comparative Genomics! We will be using this system to answer any questions about homeworks, class lectures, exams, projects, and anything else. Please take a moment to look around and get used to the system." Below the note is a "followup discussions" section with a button to "Start a new followup discussion". At the bottom, there are statistics: "Average Response Time: N/A", "Special Mentions: There are no special mentions at this time.", and "Online Now | This Week: 1 | 2". The footer contains copyright information: "Copyright © 2019 Piazza Technologies, Inc. All Rights Reserved. [Privacy Policy](#) [Copyright Policy](#) [Terms of Use](#) [Blog](#) [Report Bug!](#)".

<http://piazza.com/jhu/spring2019/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 NovaSeq 6000
Sequencing by Synthesis

>3Tbp / day



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?


Illumina Quality

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



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

FASTQC: Is my data any good?



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

Paired-end and Mate-pairs


Paired-end sequencing

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



Mate-pair sequencing

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




2x100 @ ~10kbp (outies)



2x100 @ 300bp (innies)



FASTQ Files



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

@Identifier
Sequence
+Separator
Quality Values
...

Assembly, Mapping & Genotyping

Week 2/3/4

1. Split read into segments

Read
CCAGTAGCTCTCAGCCTTATTTACCCAGGCCGTGTA TACAGGCCCTGGGTTAAATAAGGCTGAGAGCTACTGG
Read (reverse complement)
Policy: extract 16 nt seed every 10 nt
Seeds
+, 0: CCAGTAGCTCTCAGCC -, 0: TACAGGCCCTGGGTTAAA
+, 10: TCAGCCTTATTTTACCC -, 10: GGTAAAATAAGGCTGTA
+, 20: TTTACCCAGGCCGTGTA -, 20: GGCTGAGAGCTACTGG

Heterozygous variant?

Homozygous variant

...CCATAG TGTCGCCCG CGGAATT TTGATAC...
...CCAT CTATIGTGCG TCGGAATT CGGTATAC...
...CCAT GGCTATG CTATCGGAAA GCGGTATA...
...CCA AGGCTATAT CCTATCGGA TTGCGGTA C...
...CCA AGGCTATAT GCCCTATCG TTTGCGGT C...
...CC AGGCTATAT GCCCTATCG AAATTTCGC ATAC...
...CC TAGGCTATA GCGCCCTA AAATTTCGC GTATAC...
...CCATAGGCTATATGCGCCCTATCGGCAATTGCGGTATAC...

2. Lookup each segment and prioritize

Seeds
+, 0: CCAGTAGCTCTCAGCC
+, 10: TCAGCCTTATTTTACCC
+, 20: TTTACCCAGGCCGTGTA →
-, 0: TACAGGCCCTGGGTTAAA
-, 10: GGTAAAATAAGGCTGTA
-, 20: GGCTGAGAGCTACTGG

Ungapped alignment with FM Index
Seed alignments (as B ranges)
{ [211, 212], [212, 214] }
{ [653, 654], [651, 653] }
{ [684, 685] }
{ }
{ }
{ [624, 625] }

3. Evaluate end-to-end match

Extension candidates
SA:684, chr12:1955
SA:624, chr2:462 → SIMD dynamic programming aligner
SA:211: chr4:762
SA:213: chr12:1935
SA:652: chr12:1945

SIMD dynamic programming aligner

SAM alignments


| | | | | |
|--|---------|---------|------|---|
| r1 | 0 | chr12 | 1936 | 0 |
| 36M | * | 0 | 0 | |
| CCAGTAGCTCTCAGCCTTATTTACCCAGGCCGTGTA | | | | |
| II | | | | |
| A5:i:0 | X5:i:-2 | XN:i:0 | | |
| XH:i:0 | XO:i:0 | XG:i:0 | | |
| NM:i:0 | MD:Z:36 | YT:Z:UU | | |
| YM:i:0 | | | | |
| ... | | | | |

- Distinguishing SNPs from sequencing error typically a likelihood test of the coverage
 - Hardest to distinguish between errors and heterozygous SNP.
 - Coverage is the most important factor!
 - Target at least 10x, 30x more reliable

Fast gapped-read alignment with Bowtie 2
Langmead & Salzberg. (2012) *Nature Methods*. 9:357-359.

The Sequence Alignment/Map format and SAMtools
Li H et al. (2009) *Bioinformatics*. 25:16 2078-9

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 $\lambda > 10$



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