The human genome

Michael Schatz

Feb 10, 2020 Lecture 5: Computational Biomedical Research



Assignment 2: Genome Assembly Due Wednesday Feb 12 @ 11:59pm

- I. Setup Docker/Ubuntu
- 2. Initialize Tools
- 3. Download Reference Genome & Reads

4. Decode the secret message

- I. Estimate coverage, check read quality
- 2. Check kmer distribution
- 3. Assemble the reads with spades
- 4. Align to reference with MUMmer
- 5. Extract foreign sequence
- 6. dna-encode.pl -d

https://github.com/schatzlab/appliedgenomics2020/blob/mas ter/assignments/assignment2/README.md



Part I: Recap

de Bruijn Graph Construction

- $G_k = (V, E)$
 - V = Length-k sub-fragments
 - E = Directed edges between consecutive sub-fragments
 - Sub-fragments overlap by k-I words



- Overlaps between fragments are implicitly computed

de Bruijn, 1946 Idury et al., 1995 Pevzner et al., 2001

Pop Quiz 2

Assemble these reads using a de Bruijn graph approach (k=3):



Unitigging / Unipathing

- After simplification and correction, compress graph down to its non-branching initial contigs
 - Aka "unitigs", "unipaths"





Why do contigs end?

(1) End of chromosome! ⁽ⁱ⁾, (2) lack of coverage, (3) errors,
(4) heterozygosity and (5) repeats

Contig N50

Def: 50% of the genome is in contigs as large as the N50 value



Part 2: The human genome

The scale of DNA in our body is staggering.

- A typical human is comprised of roughly 40 trillion human cells (excluding trillions of bacterial cells in our gut)
- If stretched out, each haploid genome would be roughly 2 meters.
- So, each cell has 4 meters of DNA.
- 40 trillion * 4 meters = 160 trillion meters.
- 160 trillion meters / 1609.34 = 99,750,623,441 miles
- 99,750,623,441 / 92,960,000 = 1,073.05 trips to the sun.

A typical cell replicates about 100 times

160 trillion meters x 100 =

1.69123746 light years

More info

The first genetic map

Mendel's Second Law (The Law of Independent Assortment) states alleles of one gene sort into gametes independently of the alleles of another gene: *Pr(smooth/wrinkle) is independent of Pr(yellow/green)*

Morgan and Sturtevant noticed that the probability of having one trait given another was **not** always 50/50– those traits are **genetically linked**



http://www.caltech.edu/news/first-genetic-linkage-map-38798

Sturtevant realized the probabilities of co-occurrences could be explained if those alleles were arranged on a linear fashion: traits that are most commonly observed together must be locates closest together



The Linear Arrangement of Six Sex-Linked Factors in Drosophila as shown by their mode of Association Sturtevant, A. H. (1913) Journal of Experimental Zoology, 14: 43-59

Chromosome Giemsa banding (G-banding)



https://en.wikipedia.org/wiki/G banding, https://ghr.nlm.nih.gov/chromosome/1#ideogram

- Heterochromatic regions, which tend to be rich with adenine and thymine (AT-rich) DNA and relatively gene-poor, stain more darkly with Giemsa and result in G-banding
- Less condensed ("open") chromatin, which tends to be (GC-rich) and more transcriptionally active, incorporates less
 Giemsa stain, resulting in light bands in G-banding.
- Cytogenetic bands are labeled p1, p2, p3, q1, q2, q3, etc., counting from the centromere out toward the telomeres. At higher resolutions, sub-bands can be seen within the bands.
- For example, the locus for the CFTR (cystic fibrosis) gene is
 7q31.2, which indicates it is on chromosome 7, q arm, region
 3, band 1, and sub-band 2. (Say 7,q,3,1 dot 2)

The human karyotype



Bacterial Artificial Chromosomes (BACs)



- A BAC is an engineered DNA molecule used to clone DNA sequences in bacterial cells (for example, E. coli).
- BACs are often used in connection with DNA sequencing.
- Segments of an sample's DNA, ranging from 100,000 to about 300,000 base pairs, can be inserted into BACs.
- The BACs, with their inserted DNA, are then taken up by bacterial cells.
- As the bacterial cells grow and divide, they amplify the BAC DNA, which can then be isolated and used in sequencing DNA.

History of the Human Genome Project



The reference human genome

"Without a doubt, this is the most important, most wondrous map ever produced by humankind."

Bíll Clínton June 26, 2000

The reference human genome

"Without a doubt, this is the most important, most wondrous map ever produced by humankind."

Bíll Clínton June 26, 2000



The Sequence of the Human Genome Venter et al. Science 291. pp 1304-1351 (2001)

Initial sequencing and analysis of the human genome International Human Genome Sequencing Consortium Nature 409, pp 860–921 (2001)

Two Human Genomes?



The Sequence of the Human Genome Venter et al. Science 291. pp 1304-1351 (2001)

(Figure from Baker (2012) Nature Methods)

Genomic DNA BAC library Organized mapped large clone contigs BAC to be sequenced Shotgun clones Shotgun ... ACCOTARATGOGCTGATCATOCTTARA TGATCATOCTTAAACCCTGTGCATCCTACTG... sequence Assembly ... ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

Hierarchical shotgun sequencing

Initial sequencing and analysis of the human genome International Human Genome Sequencing Consortium Nature 409, pp 860–921 (2001)

The Buffelo News/Sunday, March 23, 1997

ment abuse, civil disobedience

increase their authority, always at to accept the consequences." case of the people."

water and the second

by government has forgotten that reast of the people," Parlato addacting more like it's the master." to and the Lapps share an abiding non-violent civil disobedience. insist on being respectful in our of resistance," Barbara Lyn Lapp lut if we claim to care about our

rs, we must protest government inviolence has to be the watchword. said, calling civil disobedience the

is of the violent militia movement. in-violence can serve as an antigovernment oppression, he added. law is unjust or you're given an ithout moral or legal authority.

opic. But the very nature of gov-you should refuse it," Parlato said. "And,) creates a mind set that inspires if need be, you have to be brave enough

- ----

Rachel Lapp says she believes government can be good, when it costrols the aggressors in society. Instead, it too offen comes down on the side of the aggressors. who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat belt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

The Lapps and Parlato will be joined by Samuel Radford III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993. disturbance at the City Campus of Erie'. Community College.



WANTED 20 Volunteers to participate in the Human Genome Project

a very large international scientific research effort.

The goal is to decode the human beneditary information (human Nuepvine) that determines all individual traits inherated from parents. The systome of the prisect will have memendous surgest on future progress of medical science and lead to improved diagnosis and weathent of hereditary diseases.

Volunteers will receive information about the preject from the Clinical Genetics. Service at Roswell Park, and sign a consent form before participating.

No personal information will be maintained or transferred.

Voluments will provide a one-time domation of a small bland speciment. A small more tay reindurancest will be psycided to the participants for their ness and effort,

> Individuals must be at least 18 years of age. Persons who have undergone chemotherapy are not eligible.



For-invest indivendent place constant the **Clinical Casultics Service** 845-5730 (http://au.-3100.ph) March 24 - 26, UNC

Pieter de Jong, RPCI

The Buffelo News/Sunday, March 23, 1997

ment abuse, civil disobedience

increase their authority, always at to accept the consequences." case of the people."

the second second second

ly, government has forgotten that rvant of the peepic," Parlato add-acting more like it's the master." to and the Lapps share an abiding non-violent civil disobedience. insist on being respectful in our of resistance," Barbara Lyn Lapp lut if we claim to care about our

rs, we must protest government inviolence has to be the watchword,

said, calling civil disobedience the in-violence can serve as an antigovernment oppression, he added. law is unjust or you're given an ithout moral or legal authority.

opic. But the very nature of gov-you should refuse it," Parlato said. "And,) creates a mind set that inspires if need be, you have to be brave enough

- ----

Rachel Lapp says she believes government can be good, when it costrols the aggressors in society. Instead, it too often comes down on the side of the aggressors. who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat belt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

said, calling civil disobedience the The Lapps and Parlato will be joined is of the violent militia movement. by Samuel Radiord III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993, disturbance at the City Campus of Erie's Community College.





Pieter de Jong, RPCI

The Buffelo News/Sunday, March 23, 1997

ment abuse, civil disobedience

increase their authority, always at to accept the consequences." case of the people."

and the second second

by government has forgotten that reast of the people," Parlato addacting more like it's the master." to and the Lapps share an abiding non-violent civil disobedience. insist on being respectful in our of resistance," Barbara Lyn Lapp lut if we claim to care about our

rs, we must protest government inviolence has to be the watchword,

said, calling civil disobedience the is of the violent militia movement. to-violence can serve as an antigovernment oppression, he added. law is unjust or you're given an ithout moral or legal authority.

ople. But the very nature of gov-you should refuse it," Parlato said. "And, creates a mind set that inspires if need be, you have to be brave enough

Rachel Lapp says she believes government can be good, when it costrols the aggressors in society. Instead, it too offen comes down on the side of the aggressors. who enforce child-protection laws, compulsory education, disclosure rules on tax forms and seat helt laws.

"We want people to see the correlation between what happened to us and what can happen to anyone when government gets out of hand," Rachel Lapp said.

The Lapps and Parlato will be joined by Samuel Radford III, a critic of public education who was arrested and pleaded guilty to reduced charges following a 1993. disturbance at the City Campus of Erie'. Community College.



Pieter de Jong, RPCI

Appendix: Identifying the ancestry of segments of the human genome reference sequence

To compare Neandertal to present-day human haplotypes for the purpose of population genetic analysis, we needed to have long haploid sequences from present-day humans that were of known ancestry. To identify such segments, we took advantage of the fact that the human reference sequence is haploid over scales of tens of kilobases, because it is comprised of a tiling-path of Bacterial Artificial Chromosomes (BACs) or other clone types that are of typical size 50-150 kb (S92). We do not know of any other substantial source of high quality human haploid sequences of the requisite size.

Determining the ancestries of the libraries in the human genome reference sequence using HAPMIX

It is crucial to know the 'ancestry' of a clone to use it in a meaningful population genetic analysis. In what follows, we define 'ancestry' as the geographic region in which a clone's ancestor lived 1,000 years ago, inferred based on its genetic proximity to other individuals from that region today. This definition allows us to classify clones from Chinese Americans as "East Asian," from European Americans as "European", and from African Americans as either "West African" or "European".

To identify the ancestries of the libraries comprising most of the human genome reference sequence, we used a list of 26,558 clones tiling the great majority of the genome, most of which we were able to assign to a library of origin. Restricting to the autosomes, we identified 21,156 clones that seemed to fall into 9 libraries based on the naming scheme: CTA (n=199), CTB (n=356), CTC (n=452), CTD (n=1,426), RPCI-1 (n=740), RPCI-3 (n=456), RPCI-4 (n=716), RPCI-5 (n=802) and RPCI-11 (n=16,009). (In a subsequent reexamination, we identified additional clones that we likely could have classified into libraries, including 953 from RPCI-11, 632 from RPCI-1, and 490 from another library RPCI-13.) The median span of the 21,156 clones we analyzed was 112 kb, and 80% are >50kb in size. About 2/3 came from a single library, RPCI-11.

- 1. RPCI-11 is an African American: RPCI-11, the individual who contributed most of the human genome reference sequence, is consistent with having African American ancestry, with 42% of the clones of confident West African ancestry and 42% of the clones of confident European ancestry, and the ancestry of the remaining clones less confidently inferred. The finding of likely African American ancestry for RPCI-11 was previously reported in a study of the ancestry of RPCI-11 clones spanning the Duffy blood group locus (S93), and here we confirm this finding, and also expand the inference to the whole genome.
- 2. CTD is an East Asian: The majority of clones from CTD, the second largest library in its contribution to the human genome sequence, is likely an East Asian. In a HAPMIX analysis with CEU (European) - CHB+JPT (East Asian) as the proposed ancestral populations, the majority of clones are of confident East Asian origin, and there is no secondary mode of confident European ancestry, as might be expected from a Latino or South Asian individual.
- 3. The remaining 7 libraries are European: The remaining libraries (CTA, CTB, CTC, RPCI-1, RPCI-3, RPCI-4 and RPCI-5) are inferred to be of European ancestry, since they all have consistent distributions of inferred clone ancestries, with the majority of clones of confident European ancestry in both our HAPMIX analyses and no secondary modes.

A Draft Sequence of the Neandertal Genome

Green et al (2010) Science. DOI: 10.1126/science.1188021 Supplemental Note 16 (pg 145-146)

Welcome back: Michael Schatz nature methods O Logout U Cart Techniques for life scientists and chemists Search ge Advanced search Journal home > Archive > Editorial > Full Text. EDITORIAL Journal content Subscribe to Nature Methods Journal home Nature Methods 7, 331 (2010) Subscribe doi:10.1038/nmeth0510-331 Advance online publication E pluribus unum Current issue This issue If the human reference genome is to reflect more of the actual genomic diversity in Archive Table of contents humans, community participation is needed. Focuses and * Next article Supplements Article tools Please visit methagora to view and post comments on this article. Methagora blog Download PDF Method of the Year The human genome is ten years old. We acknowledge its reference assembly as an invaluable 2016 £3 Send to a friend resource essential for many purposes such as the assembly of short reads from high-CrossRef lists 11 articles citing throughput sequencing platforms into chromosome context during resequencing projects. At Multimedia the same time, we think necessary improvement of the reference genome depends on the this article Press releases willingness of the research community to provide data for the genome's less accessible Scopus lists 9 articles citing regions. this article **Journal Information** First published in 2001, the human reference genome has, since 2007, been in the hands of Export citation Guide to authors the Genome Reference Consortium (GRC) a small group of fewer than 20 scientists from the Rights and permissions European Bioinformatics Institute, the US National Center for Biotechnology Information, The Reporting checklist Sanger Institute and The Genome Center at Washington University in St. Louis, who have of Online submission committed to the improvement and completion of this reference, with very little financial naturejobs Subscribe support. New Subscription **Recruitment of Professors** and Associate Professors Renew Subscription The reference genome is now in its 19th rendition, and probably the best measure of its School of Materials Science improvement over the last ten years is the number of fragments it consists of. The very first L Paid Subscriptions and Engineering, Sun Yat-Change of Address version had ~150,000 gaps; the most recent build, GRCh37, has only around 250 gaps. sen University Sun Yat-sen University Permissions The only other publicly accessible de novo assembly of a human genome that contains For referees Faculty positions at Institut chromosome sequences is HuRef. Obtained by traditional capillary sequencing, HuRef is the franco-chinois de l'énergie diploid genome of Craig Venter. It comes in 4,500 pieces and, like any individual genome, it Contact the journal nucléaire contains many rare alleles. Institut franco-chinois de l'énergie About this site nucléaire Sun Yat-sen University GRCh37, in contrast, is a mosaic haploid genome derived from about 13 people. It still contains rare alleles, but the GRC recently decided to convert these to common haplotypes. More science jobs **Nature Research** Deciding which alleles are common and which are rare is proving challenging, and the GRC services Post a job members are collaborating with members of the 1000 Genomes project to collect enough data Authors & Referees to make these decisions. Advantiation



The human genome - basic stats

											000000000000000000000000000000000000000		0000					00	1000		Ĥ	
Į	ę	Ą	Ģ	P	ų	Ę	ļ	ļ	IDa	IIID:		Da		Dr	D	18	Q	E Da		Ş	Q	กไ

3.096 billion base pairs (haploid)

- 20,454 protein coding genes
- 226,950 coding transcripts
 (isoforms of a gene that each encode a distinct protein product)

Assembly	GRCh38.p12 (Genome Reference Consortium Human Build 38), INSDC Assembly GCA. 000001405.27, Dec 2013
Base Pairs	3,609,003,417
Golden Path Length	3,096,649,726
Annotation provider	Ensembl
Annotation method	Full genebuild
Genebuild started	Jan 2014
Genebuild released	Jul 2014
Genebuild last updated/patched	Mar 2019
Database version	97.38
Genoode version	GENCODE 31

Gene counts (Primary assembly)

Coding genes	20,454 (incl 660 readthrough)	
Non coding genes	23,940	
Small non coding genes	4,871	
Long non coding genes	16.848 (ind 302 readthrough)	
Misc non coding genes	2.221	
Pseudogenes	15,204 (incl 8 readthrough)	
Gene transcripts	226,950	

http://uswest.ensembl.org/Homo_sapiens/Location/Genome

Solely 2% of the human genome encodes proteins.



https://genome.ucsc.edu

Half of the human genome is comprised of repeats



http://www.nature.com/nrg/journal/v10/n10/pdf/nrg2640.pdf

Half of the human genome is comprised of repeats



Repetitive DNA not driven by retrotransposition (e.g., ATATATATATATATATATAT...)

GC content varies dramatically in the genome



Figure 12 Histogram of GC content of 20 kb windows in the draft genome sequence.



http://www.nature.com/nrg/journal/v10/n10/pdf/nrg2640.pdf

The human reference genome continues to change.

- Ongoing efforts to fill "gaps" and properly/thoroughly represent complex structures and loci in the genome (e.g., Major Histocompatibility Complex)
- Each improvement leads to a new genome "build". Currently on build 38.
- Experimental and computational methods provide new genome annotations
 - O New gene models, transcription factor binding sites, and loci where human individuals differ (i.e., polymorphisms)
- Therefore, the human reference genome is by no means "complete"!
- How does the same genome yield such phenotypic diversity across tissue types?
- How does the genome evolve within an individual (tissues) and among a population?

Genomics Arsenal in the Year 2020



10X Genomics Linked Reads



10X Genomics Linked Reads



(Zheng et al, 2016)

- 1. High molecular weight DNA is diluted and isolated within oil emulsion droplets
- 1. Within each droplet, short fragments are randomly amplified and tagged with barcode sequences for standard Illumina sequencing
- Reads sharing the same barcode can then be localized to the same original template molecule
- The resulting "linked reads" can be used for phasing variants or identifying SVs

https://www.youtube.com/watch?v=nk2kXM59LRM

Haplotype Phasing



c NA12878 Optimal phase variant span increases with read length

b NA12878 Optimal phase block length increases with read length



Piercing the dark matter: bioinformatics of long-range sequencing and mapping Sedlazeck et al. (2018) *Nature Reviews Genetics.* 19:329

Uncertain Future for IOX



PacBio Single Molecule Real Time Sequencing (SMRT-sequencing)



PacBio: SMRT Sequencing

Imaging of florescent phospholinked labeled nucleotides as they are incorporated by a polymerase anchored to a Zero-Mode Waveguide (ZMW).







Time

http://www.youtube.com/watch?v=v8p4ph2MAvI

Single Molecule Sequences



"Corrective Lens" for Sequencing



"Corrective Lens" for Sequencing



Consensus Accuracy and Coverage



Coverage can overcome random errors

- Dashed: error model from binomial sampling; solid: observed accuracy
- For same reason, CCS is extremely accurate when using 5+ subreads

$$CNSError = \sum_{i=\lceil c/2 \rceil}^{c} \binom{c}{i} (1-e)^{n-i}$$

"HiFi" Circular Consensus Reads

High-quality reads produced by sequencing the same molecule multiple times

Higher accuracy for low-coverage sequences like somatic variants or lowly expressed transcripts in RNAseq, more interpretable alignments, faster assembly

Limits read length, used to be very expensive but more manageable now



Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome Wenger et al (2019) Nature Biotechnology doi:10.1038/s41587-019-0217-9

Methylation Detection

- **Methylation** an epigenetic modification that can have a variety of effects, such as gene repression
- Can detect methylation from raw PacBio signal







PACB



BioWorld BioWorld MedTech BioWorld Asia Market Intelligence reports

Illumina, Pacific Biosciences abandon \$1.28 merger over FTC opposition



competition in the next generation sequencing market.

by Mark McCarty his Comments

Two players in the game sequencing space, Illumina and Placific Biosciences, have scotched their planned \$1.2 billion merger roughly two weeks after the U.S. Federal Tade Commission (FTC) posted a 5-0 vote to seek an injunction against the merger. While Illumina is consequently liable for nearly \$100 million in termination fees, it could recoup those monies under some circumstances.

The \$1.2 billion merger between Illumina inc., of San Diego, and Pacific Biosciences of California Inc., was formally announced by the two companies in November 2018, but the deal faced substantial regulatory difficulty from the outset. The FTC said in a Jan. 2 <u>statement</u> the deal would have guashed

The companies signed an extension to the deal in September 2019 to allow more time to come to terms with regulators, but that deadline was extended to the end of March 2020 in a handshake dated Dec. 18, 2019. Whether the most recent extension was a plausible effort to keep the deal together has been debated, given that the FTC had voted to oppose the merger Dec. 17, 2019, the day before the two companies agreed to give the effort one more extension.

Popular Stories



Thiel calls for improving research grant, regulatory processes to enhance scientific invosation

Subscribe

Q,



Sign In



Insulet to launch wearable insulin pump this year, works with Descorn ICGM for closed loop





Noticela



Novavax developing nanoparticle vaccine for Wuhan coronavirus ButWorld

Oxford Nanopore Technologies (ONT)



Nanopore Sequencing

Sequences DNA/RNA by measuring changes in ionic current as nucleotide strand passes through a pore



nanoporetech.com/applications/dna-nanopore-sequencing



Translation of raw signal into basepairs



Translation of raw signal into basepairs

Early basecallers began by estimating k-mer boundaries using "events", which were then input to an HMM

Modern basecalers use neural networks directly on raw signal



(Based on probability of event matches)

ONT releases k-mer models with expected current distribution of every k-mer



DNA Base-Calling from a Nanopore Using a Viterbi Algorithm Timp et al. (2012) *Biophysical Journal*



Certain k-mers can be eliminated based on possible transitions



DNA Base-Calling from a Nanopore Using a Viterbi Algorithm Timp et al. (2012) *Biophysical Journal*



Final sequence determined by most probable k-mers



"DNA Base-Calling from a Nanopore Using a Viterbi Algorithm" Timp et al. (2012) *Biophysical Journal*

Basecaller/Pore Timeline

Development of both pore chemistry and basecalling algorithms is responsible for improvement in accuracy



From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy Rang *et al* (2018) *Genome Biology.* https://doi.org/10.1186/s13059-018-1462-9

New Pore Chemistries

ONT is developing alternate pore chemistries to improve accuracy, particularly for homopolymers





https://vimeo.com/272526835

DNA Modification Detection

Like PacBio, ONT can detect methylation from raw signal

• Or any other modification that changes ionic current



Piercing the dark matter: bioinformatics of long-range sequencing and mapping Sedlazeck et al. (2018) *Nature Reviews Genetics.* 19:329

Direct RNA-seq

Standard RNA sequencing (RNAseq) requires creation of complementary DNA (cDNA)

ONT recently introduced direct RNA sequencing

Allows detection of RNA modifications, and potentially secondary structure



Nanopore native RNA sequencing of a human poly(A) transcriptome

Workman et al. Nature Methods. 16:1297–1305



Oxford Nanopore sets sights on IPO

🛗 4th April 2019 👗 Callum Cyrus

The Oxford University genetic sequencing spinout is reportedly mulling an IPO that would provide exits to investors including commercialisation firm IP Group.

Oxford Nanopore Technologies, a UK-based genetic sequencing technology developer spun out from University of Oxford, is considering floating its shares in an initial public offering (IPO), The Telegraph has reported. Founded in 2005, Oxford Nanopore has developed real-time DNA and RNA sequencing technology that offers biological analyses at a relatively low cost. It has applications...

Recent Long Read Assemblies



Lee et al (2016) *bioRxiv* doi: http://dx.doi.org/10.1101/048603



First Telomere-to-Telomere Human Chromosome



Telomere-to-telomere assembly of a complete human X chromosome Miga et al. (2019) bioRxiv. https://doi.org/10.1101/735928

Assembly Summary



Assembly quality depends on

- I. Coverage: low coverage is mathematically hopeless
- 2. Repeat composition: high repeat content is challenging
- 3. Read length: longer reads help resolve repeats
- 4. Error rate: errors reduce coverage, obscure true overlaps
- Assembly is a hierarchical, starting from individual reads, build high confidence contigs/unitigs, incorporate the mates to build scaffolds
 - Extensive error correction is the key to getting the best assembly possible from a given data set
- Watch out for collapsed repeats & other misassemblies
 - Globally/Locally reassemble data from scratch with better parameters & stitch the 2 assemblies together

Next Steps

- I. Reflect on the magic and power of DNA \odot
- 2. Check out the course webpage
- 3. Register on Piazza & GradeScope
- 4. Work on HW2