### Genome Assembly and Alignment

Michael Schatz

Feb 8, 2018 Lecture 4: Applied Comparative Genomics



### Assignment I: Chromosome Structures Due Feb 8 @ 11:59pm

	Bilde Bildente () eit () für Lannes IV ( Lanne AV) Lanne att	Reference and a contract of the State Stat
O Thinking here	Putrepasts same Marketplace Explore	
schatzlab / appliedgenomics2018		Gimments 3 & Bar 3 Yres 1
ODM Dises # University (Passa 2 1996 L	naight Cherrys	
tent water appliedgenamics2018 ( assignments ) assignment ( BEADM	End	Fruit The College and
C velocity canadox. Miletared. and		Sectored Association
antidana IR		
14 (1994) (10 plac) 0.00 MB		the films throw G / 1
Assignment 1: Chromosome Structures		
Assignment Date: Thursday, Feb. 1, 2018 Due Date: Thursday, Feb. 8, 2017 @ 11:88pm		
Assignment Overview		
In this assignment you will profile the overall structure of the genetites of should be posted to Flacts	several important species and then study the yeast gorome in m	ore detail. As a reminder, any questions about the assignment
Some of the tools you will need to use this semester only run in a linux en https://github.com/schatzleb/spalledgenoms.2016/biob/meeter sesignm	whenment, if you do not have access to a linux machine, downloa emplyinguables, nd	at and install a virtual machine following the directions here:
Question 1: Chromosome structures		
Download the chomosome size files for the following genomes (Note the	e have been preprocessed to only include main chromosomes):	
1. Arabitopeix mariana (TABHD) - An important plant model species (in	-	
2. Corn (Zee mays 873-63) - The most widely grown crop in the world (	and a later	
3. E. col (Escherichia col K12) - One of the most commonly studied be	cteria (Info)	
4. Fruit Fly (Drootphila melanopaster, strill) - One of the most important	I model species for genetics (IMs)	
B. Human (hg38) - ue 3 (infe)		

### https://github.com/schatzlab/appliedgenomics2018

### Part I: Recap

## Two Paradigms for Assembly



Short read assemblers

- Repeats depends on word length
- Read coherency, placements lost
- Robust to high coverage



Long read assemblers

- Repeats depends on read length
- Read coherency, placements kept
- Tangled by high coverage

Assembly of Large Genomes using Second Generation Sequencing Schatz MC, Delcher AL, Salzberg SL (2010) *Genome Research*. 20:1165-1173.



Table 2 de novo human g	genome (NA18507) assemblies
-------------------------	-----------------------------

Method	Minia	C.&B.	ABySS	SOAPdenovo
Value of k chosen	27	27	27	25
Number of contigs (M)	3.49	7.69	4.35	
Longest contig (kbp)	18.6	22.0	15.9	
Contig N50 (bp)	1156	250	870	886
Sum (Gbp)	2.09	1.72	2,10	2.08
Nb of nodes/cores	1/1	1/8	21/168	1/16
Time (wall-clock, h)	23	50	15	33
Memory (sum of nodes, GB)	5.7	32	336	140

de novo human genome (NA18507) assemblies reported by our assemblier (Minia), Conway and Bromage assembler [9], ABySS [8], and SOAPdenovo [7]. Contigs shorter than 100 bp were discarded. Assemblies were made without any pairing information.



**Space-efficient and exact de Bruijn graph representation based on a Bloom filter** Chikhi and Rizk (2013) Algorithms for Molecular Biology. 8:22

## Contig N50

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



## Contig Nchart



### 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

### Part 2: Whole Genome Alignment & HW2



## Whole Genome Alignment with MUMmer

Slides Courtesy of Adam M. Phillippy NHGRI



• For two genomes, A and B, find a mapping from each position in A to its corresponding position in B



### Not so fast...

 Genome A may have insertions, deletions, translocations, inversions, duplications or SNPs with respect to B (sometimes all of the above)



## WGA visualization

- How can we visualize *whole* genome alignments?
- - A perfect alignment between A and B would completely fill the positive diagonal



## SV Types



- Different structural variation types / misassemblies will be apparent by their pattern of breakpoints
- Most breakpoints will be at or near repeats
- Things quickly get complicated in real genomes

http://mummer.sf.net/manual/ AlignmentTypes.pdf



### Alignment of 2 strains of Y. pestis http://mummer.sourceforge.net/manual/

### Assignment 2: Genome Assembly Due Feb 15 @ 11:59pm

#### Assignment 2: Genome Assembly

Assignment Date: Thursday, Feb. 8, 2018 Due Date: Thursday, Feb. 15, 2018 @ 11.59pm

#### Assignment Overview

In this assignment, you are given a set of unassembled reads from a mysterious pathogen that contains a secret message encoded somepiace in the genome. The secret message will be recognizable as a novel insertion of sequence not found in the reference. Your task is to assess the quality of the reads, assemble the genome, identify, and decode the secret message. If all goes well the secret message should decode into a recognizable english text, otherwise doublecheck your coordinates and try again. As a reminder, any questions about the assignment should be posted to Plazza.

Some of the tools you will need to use only run in a linux environment. Spades, for example, will not work under Mac, even though it will compile. If you do not have access to a linux machine, download and install a visual machine following the directions have https://gthub.com/schatulab/appliedgenomics2018/b/du/master/assign/nenta/visualbox.md

#### Question 1. Coverage Analysis (10 pts)

Download the reads and reference percent from. https://github.com/schatt/sk/appliedgenomics2018/tew/master/assignments/

Note I have provided both paired-end and mate-pairs reads (see included README for details). Make sure to look at all of the reads for the coverage analysis and ioner analysis, as well as in the assembly.

- + Question 1a. How long is the reference genome? (Hint: Try samoula faile)
- Question 1b. Now many reads are provided and how long are they? Make sure to measure each file separately Drint: Try Fast(); []
- Question 1c. How much coverage do you expect to have? [Hint: A little arthmetic]
- Question 10. Plot the average quality velue across the length of the reads (Hint: Screenshot from Yax150.)

#### Question 2. Kimer Analysis (10 pts)

Use 3s(1y1)s) to count the 21-mers in the reads data. Make sure to use the "-C" flag to count cannonical kmers, otherwise your analysis will not correctly account for the fact that your reads come from either strand of DNA.

### https://github.com/schatzlab/appliedgenomics2018

### Halomonas sp. GFAJ-I





<u>Library 1: Fragment</u> Avg Read length: 100bp Insert length: 180bp <u>Library 2: Short jump</u> Avg Read length: 50bp Insert length: 2000bp

A Bacterium That Can Grow by Using Arsenic Instead of Phosphorus Wolfe-Simon et al (2010) *Science*. 332(6034)1163-1166.

## **Digital Information Storage**



Fig. S1. Schematic of DNA information storage.

Encoding/decoding algorithm implemented in dna-encode.pl from David Dooling.

### **Next-generation Digital Information Storage in DNA**

Church et al (2010) Science. 337(6102)1628

### **Mission Impossible**

- I. Setup VirtualBox
- 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/appliedgenomics2018/blob/mas ter/assignments/assignment2/README.md



### Find and decode

### nucmer -maxmatch ref.fasta \ default/ASSEMBLIES/test/final.contigs.fasta

-maxmatch	Find maximal	exact matches	(MEMs) without	repeat filtering
-p refctg	Set the outpu	ut prefix for (	delta file	

### mummerplot --layout --png out.delta

layout	Sort the	alignments	along the	diagonal
png	Create a	png of the	results	

### show-coords -rclo out.delta

-r	Sort alignments by reference position
-C	Show percent coverage
-1	Show sequence lengths

-o Annotate each alignment with BEGIN/END/CONTAINS

### samtools faidx default/ASSEMBLIES/test/final.contigs.fasta

Index the fasta file

### samtools faidx default/ASSEMBLIES/test/final.contigs.fasta \ contig\_XXX:YYY-ZZZ | ./dna-encode -d

### Part 3: Long Range Sequencing and Assembly

### Genomics Arsenal in the year 2018



### Assembly Complexity





### Assembly Complexity





## Assembly Complexity





The advantages of SMRT sequencing Roberts, RJ, Carneiro, MO, Schatz, MC (2013) *Genome Biology*. 14:405

### Celera Assembler aka CANU

https://github.com/marbl/canu

- I. Pre-overlap
  - Consistency checks
- 2. Trimming
  - Quality trimming & partial overlaps
- 3. Compute Overlaps
  - Find high quality overlaps
- 4. Error Correction
  - Evaluate difference in context of overlapping reads
- I. Unitigging
  - Merge consistent reads
- I. Scaffolding
  - Bundle mates, Order & Orient
- I. Finalize Data
  - Build final consensus sequences



### Overlap between two sequences



CAGTACTTGGATGCGCTGACACGTAGCTTATCCGGT...

overlap (19 bases)

overhang % identity = 18/19 % = 94.7%

**overlap** - region of similarity between regions **overhang** - un-aligned ends of the sequences

The assembler screens merges based on:

- length of overlap
- % identity in overlap region
- maximum overhang size.

[How do we compute the overlap?]

overhang (6 bases)

## Very fast overlapping

#### ARTICLES

#### manare biotechnology

#### Assembling large genomes with single-molecule sequencing and locality-sensitive hashing

Konstantin Berlin<sup>1-2,4</sup>, Sergey Konst<sup>1,4</sup>, Chen-When Chin<sup>3</sup>, James P. Dvale<sup>3</sup>, Jame M. Landolin<sup>3</sup> & Adam M. Phillippy<sup>4</sup>

Long mail, single-molecule real-time (SMRT) sequencing is realizely used to finish microbial genomes. Sut exactable assembly methods have not stated well to large genumes. We introduce the Workash Algoriterit Process (2014)P) for everyoging noise. long made using protodelistic, locally establish hashing, integrating 804P with the Colors Assembles analysis effectives grade de neue assembles el Sectematores constitue, Arabidgals thaliana, Draughila maiacegalte and a human hybridifiem mais online (CNM) has UMRI sequencing. The resulting assembles are highly continuous, include fully reached chromosome arris and close percident gaps in these reference generies. Our assembly of 2: metanogastic revealed precisedy unknown haterachronatic and televiner's transition sequences, and we poperided tow complexity sequences from CHRS that 10 gaps in the human GRCH38 reference. Using MHMP and the Celeria Assemblier, single-molecule sequencing can produce dit next near-complete substantial essentialies that are \$5,50% accurate when nempered with sualiable reference generates.

genote property. Collie energy-using property, de new secon-genetic is reconstructed from actually. An accurate reconstruction is sequencing loads of only a live bandred have pairs, which to dearbar than room common repeats. Although short made an sofficient for many adaptions, they see not sufficient for smalling studet report famiitin is other microbial or sukaryoite generate, builting to itsettend and incomplete uncertainer. Knowl whence in single millicials asparining to blockging have prosted reads loadness of hid longer that should generative methods11. Most namble, Paulle Beselement MART separately was the first committenally scalable long read indexelogy". Using a 1968 polynamic androrol in a serie mole serieguide, budia MdRT

in the assembly graph<sup>a on</sup> However, the bing reads generated by six-

gis maintails arguming correctly saffle from less accorecy (82-47%)

Pacifics11, 78-1076 MatXXX71, and unw signations have been model

to competants for miniates to the two sequences \$2.2-17, Alderight DERT segantizing may be over prove, it exhibits has segarate

ing him that previous technologies (4.7), and theoretical research

Genotes assessible in the process of reconstructing a genetic from . Thus, he overseeping the genetic at adjuster covering large. Mr. a collection of short sequencing teach and is on imaged may to any of Pacific PSCO, DRICT sequencing can be used to produce highly accuracy and continuous assembles,"122-27, including automatically

Burly searchildus of mitty, long tooth have bend ancionalist, but have stratist, as both the continuity and train accuracy of an assembly can addressed from a addressed al computational cost. For reample, as tobial affect the tends of all downstream analyses? Bowerer, repetitive analysis of all melanopairs from DelET made required more than expansion make asserting difficult when the repeat length cannot will diffic OV haves, when tools are didn't at the tase were used - the the read length? Most high densighput separating methods generate reportedent of some than 32 date running on a throughd core conprice chains14. Even and hasterial generate pervisedy sugared a day in assemble using the 26,45" or FB-R" assembly pipelness The printers bottlenank of king-tool assembly has been the sensi tion all surnar all alignment required to determine everlapping read pairs. For the 21 webrogaster UMET assortify, this overlapping stop conversion? \$7% of the total runtime. Does if the sciturate of heap read technologies tesperent, all parts overlapping will unusin a substantial bettimesk is overlap locitat consensus assentily. Taken together, the compositional cost and the comparatively high sequencing cross have prevented addressed application of SMET aspancing to graviture supposing her delivered peaks wath of up to 14 Mp\*. Periminary larger than 100 Mbp. The netally the reasing itercophysis of the Public readle have Deltad Nanopore aggree that vit Mp read lengths are instrument has heper to address expansion or the compupossible using NUSICO satesport organizing". Inchilorg read brights internal out of assembling larger passing his remained beyond the disationally anythic process manufally by concerning represent attachment, much of most investigation,

To aslite the computational problem of long read assembly, we present a photodilatic algorithm for officiends detecting searlaps: Inferen party long mark. MUAP uso a dimensionality reduction technique tartaed Wind balls<sup>28</sup> for credit a more compact separametries all sequencing roads. Originally developed to determine the predatity of and pages". Multials reduces a test or string to a small set of Sau shares that random server can be overcome algorithmically?" Angerptics, called a deech Minifash databas have been recoverially

Reactive of Distriction and Nuclearing's interaction of Neurosci, USAN, Neurosci, 1983, Tradition for Assances Darkstee Touches, private of Neurosci. Antiga Pan, Baylon, Alli, Sonnia Jah, Arogin, Sepris AM, Natora Biablona Antigat and Scientimation Street Families (Walker) 200. Washin Readynamic Garlieria. In: Bern Pan, Gelleria. 200. "Paga active constraint statistics in the ann. Languagement educe in adversary." tor \$1.6. Management (sector) - control

Number of August 2018 another 9 April 2018, another other 19 May 2005, multi 10 Minor 18 Min

that's do been to provide the total of a second of a serie period

а	S1: CATOGACCGACCAG					GCAGTACCGATCGT : S2 GTA CGA CGT					
		A	TG A	CC A	CC CC		AC	TC	G TC	G	
			OGA	CGA	CAG		GCA	TAC	GAT		
b	$r_1$	$\Gamma_2$	$\Gamma_3$	Γ4	1		1	$\Gamma_1$	$\Gamma_2$	$\Gamma_3$	$\Gamma_4$
	19 14 58	14 57 37	57 36 16	36 19 15	CAT ATG		GCA CAG AGT	36 18 11	19 13 54	14 56 33	57 39 28
	33 5 22	23 28 48	11 47 60	54 26 43	GAC ACC CCG	¥	ACC	49 5 22	44 48	27 47 60	49 26 43
	24 33 5 20	7 28 48 3	50 11 47 62	45 54 26 41	CGA GAC ACC CCA		CGA GAT ATC TCG	24 35 13 54	7 30 56 33	50 9 39 28	45 52 18 11
с	[5	1.	2,	15]	mir	+ n-me	rs	[5	. 1.	6,	6]
		Sketc	th (S <sub>1</sub> )	)					Sketc	th (S2)	,
d				J	(S1, S2	)=2	/4 = 0	.5			
е				S1:	CATO	GAC	CACO	AG			
				S. :	GCAG	TAC	CATC	GT			

1000

## Unitigging: Pruning the Overlap





# Questions?