# **Genome Assembly**

Michael Schatz

Feb 3, 2020 Lecture 3: Applied Comparative Genomics



### Assignment I: Chromosome Structures Due Feb 5 @ 11:59pm

2.1	C .	è i		pithut	.004	victo	triab	Napr	plea	penor	nics2	120/	tree,h	raster	(anti-	ymen	1,445.4	gnmer	rta.			9	\$		ų	1	2	8	۲	¢	o	=		٠		R	6	0
•	4.	(Ma	8	t Owl		0	9.	-	sleck	0	Ð	4	6	GRANT	ns (	100	. 6	×	0		н. В	2 74	-	Mad	54	ø	Re l	Dook	-	10	shop	ø	efit		•	8	00++	Bookmark
A	ssig	nmi	nt	1: C	hre	mo	son	ne	Stru	Jeta	ires																											
As Du	signr e Dat	ent D a: We	ute: dner	villedr idity,	esdi Feb.	ey, Jan 5, 202	29, 0 @	202	10 19pm																													
As	isign	men	t O	verv	iew																																	
in i asi	this an	inigni ent si	nent Youk	you i the p	vill p oste	rofile t d to P	he a	vera	ill sthr	actur	e of t	ie ge	поте	is of s	even	impo	tant s	pecie	s and	then s	fudy (	the y	east (	geno	THE	in r	ore	deta	il. A	i a n	entino	ler, a	ny qu	estic	ins at	bout	the	
Q	vesti	on 1	Cł	rom	050	me s	truc	ttur	res																													
De	Download the chomosome size files for the following genomes (Note these have been preprocessed to only include main civomosomes):																																					
1	1. Arabidopsis thaliana (TAR10) - An important plant model species (info)																																					
- 3	2. Tomato (Solanum lycopersicum x4.00) - One of the most important food crops (mfo)																																					
- 3	3. E. coli (Escherishia coli K12) - One of the most commonly studied bacteria [info]																																					
- 3	4. Fruit Fly (Drosophila melanogester, dm6) - One of the most important model species for genetics (info)																																					
- 3	8, Human (hg38) - us () [info]																																					
	5. Whe	tet (T	Rica	en 14	stivu	m, Mi	280)	- T	he fo	od er	op wi	ich I	akes	up the	iarg	est lan	5 area	(into)																				
- 2	7. Wot	m (C)	iend	rhabo	Pilin e	legan	s, ce	na)	- One	t of t	he mo	at irr	porta	int ani	mal r	nodel s	pecies	(into)	ŀ																			
- 1	5. Year	M (Se	coh	vom	OPE (	benevia	lae,	sec	Cer3)	- an	impó	tant	euka	ryotic	mòđ	el spec	ies, ab	to goo	of for	bread	and b	beer	(info)															
	ing th	eso f	les,	Tuke	a tat	sie wit	h the	) fol	lawing	g info	rmati	on p	er spe	icies:																								
	Que	note	11	fotal y	peno	me siz																																
	Que	ntion	1.2.	Numb	er o	t chron	1080	amer	í																													
	Que	retion	1.3.	Large	st ch	romor	ome	i sia	e and	Inam																												
	Que	etion	1.4.	Small	est c	hromo	som	ie sl	ze an	d nar	ne i																											
	Que	stion	1.5.	Mean	chro	moso	ne k	ingt	h																													
:																																						
:	vesti	on 2	: Se	que	nce	cont	ent																															

### https://github.com/schatzlab/appliedgenomics2020

Part I: Recap & Coverage Analysis

# Second Generation Sequencing



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

# Illumina Quality





http://en.wikipedia.org/wiki/FASTQ\_format

# 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?







num balls

-----



num balls

# **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





# Normal Approximation





I want to sequence a 10Mbp genome to 24x coverage. How many 120bp reads do I need?

> I need I0Mbp x 24x = 240Mbp of data 240Mbp / I20bp / read = 2M 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\*sqrt(X) = 24

36-2\*sqrt(36) = 24

I need 10Mbp x 36x = 360Mbp of data 360Mbp / 120bp / read = 3M 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 ~\$1000 per sample, while WES currently costs ~\$250 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

# **Adaptive Sequencing**



**Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED** Kovaka, S, Fan, Y, Ni, B, Timp, W, Schatz, MC (2020) bioRxiv doi: https://doi.org/10.1101/2020.02.03.931923

# Adaptive Sequencing



	chri 110 11 11 11 11 11 11 11 11 11 11 11 11	669 bp 613,000 bp 613,000 bp	43.310,100 to 41.310,300 to 41.310
UNCALLED DUE		and the second division of the second divisio	
UNCALLED CON	the last in the last in the last		and the second se
UNCALLED Reads	-		r
ONT WGS SVs			1
ONT WES COV	particular international products	and the second second second second	NAME AND ADDRESS OF AD
ONT WGS Reads			
PacBio SVs			
PacBio Cov	38-871	12	
PacBio Reads			
Illumina SVs		No Variants Fours	4
filumina Cov	II-III	the second se	and the second se
Illumina Reads			
Sequence = Genes		MUTYH	

**Targeted nanopore sequencing by real-time mapping of raw electrical signal with UNCALLED** Kovaka, S, Fan, Y, Ni, B, Timp, W, Schatz, MC (2020) bioRxiv doi: https://doi.org/10.1101/2020.02.03.931923

## Part 2: De novo genome assembly



# Outline

- I. Assembly theory
  - Assembly by analogy
- 2. Practical Issues
  - Coverage, read length, errors, and repeats
- 3. Next-next-gen Assembly
  - Canu: recommended for PacBio/ONT project
- 4. Whole Genome Alignment
  - MUMmer recommended

# Shredded Book Reconstruction

- Dickens accidentally shreds the first printing of <u>A Tale of Two Cities</u>
  - Text printed on 5 long spools

It wa	s the	besthof	bes <b>times</b> n	tsyas	the the	<b>rst</b> or of	times, i	it was th	e ag	<b>ge</b> bf	v <b>ivsis do</b> mį	t <b>itwas</b> l	he <b>alge</b>	age00	fi <del>fdolishn</del> ess	,
It wa	s tine	besthe	of times,	it was	the he	e worst	of times.	it was f	he the	eastev	obatrado	<b>mwat</b> s t	h <b>evas</b> et	hê fag	otioth foodishne	ess.
								,					8-		,,,,,,,	
It wa	s the	zasbiest	besimesini	ewait	vahelw	owstrof	times,eit,	, it <b>was t</b>	he ag	e of w	<b>risdom,</b> i	it was	the ag	e of	lisoolisisness	,
It wa	s tt	asbase	<b>besinge</b> sin	ies, wita	sabaha	eowstrof 1	tifntesnes	it was t	he ag	e of 7	vi <b>sdsdo</b> ,n	it, itasas	hthtgag	tef foc	ofistoolistaness	8,
It	walt ti	aæsbidset	<b>b£simes</b> in	eist, vita	wabelw	owstrof	of <b>times</b> ,	, it was t	he ag	<b>e</b> of o	fiwikahomi	it itravsats	h <b>thæ</b> ge	agfoo	fitoolistapess	,

- How can he reconstruct the text?
  - 5 copies x 138,656 words / 5 words per fragment = 138k fragments
  - The short fragments from every copy are mixed together
  - Some fragments are identical





The repeated sequence make the correct reconstruction ambiguous

It was the best of times, it was the [worst/age]

Model the assembly problem as a graph problem

How long will it take to compute the overlaps?

# 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



# de Bruijn Graph Assembly



# The full tale

... it was the best of times it was the worst of times ...
... it was the age of wisdom it was the age of foolishness ...
... it was the epoch of belief it was the epoch of incredulity ...
... it was the season of light it was the season of darkness ...
... it was the spring of hope it was the winder of despair ...





**Reducing assembly complexity of microbial genomes with single-molecule sequencing** Koren et al (2013) Genome Biology. **14**:R101 <u>https://doi.org/10.1186/gb-2013-14-9-r101</u>

# $\begin{array}{c} \textbf{Counting Eulerian Cycles}\\ \textbf{A} & \textbf{R} & \textbf{D} & \textbf{ARBRCRD}\\ \textbf{O} & \textbf{O} & \textbf{O} & \textbf{O} \\ \textbf{O} & \textbf{C} & \textbf{CRBRD} \end{array}$

Generally an exponential number of compatible sequences

- Value computed by application of the BEST theorem (Hutchinson, 1975)

$$\mathcal{W}(G,t) = (\det L) \left\{ \prod_{u \in V} (r_u - 1)! \right\} \left\{ \prod_{(u,v) \in E} a_{uv}! \right\}^{-1}$$
  
L = n x n matrix with  $r_u$ - $a_{uu}$  along the diagonal and  $-a_{uv}$  in entry uv  
 $r_u = d^+(u) + l$  if  $u = t$ , or  $d^+(u)$  otherwise  
 $a_{uv}$  = multiplicity of edge from u to v





It is believed 74% of the mass of the Milky Way, for example, is in the form of hydrogen atoms. The Sun contains approximately **10<sup>57</sup> atoms** of hydrogen. If you multiple the number of atoms per star (10<sup>57</sup>) times the estimated number of stars in the universe (10<sup>23</sup>), you get a value of **10<sup>80</sup> atoms** in the known universe. Nov 5, 2017



How Many Atoms Are There in the Universe? - ThoughtCo https://www.thoughtco.com/number-of-atoms-in-the-universe-603795



- Finding possible assemblies is easy!
- However, there is an astronomical genomical number of possible paths!
- Hopeless to figure out the whole genome/chromosome, figure out the parts that you can



# Contig N50

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



# Contig N50

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

### Better N50s improves the analysis in every dimension

- Better resolution of genes and flanking regulatory regions
- Better resolution of transposons and other complex sequences
- Better resolution of chromosome organization
- Better sequence for all downstream analysis

### Just be careful of N50 inflation!

- A very very very bad assembler in 1 line of bash:
- cat \*.reads.fa > genome.fa

N50 size = 3 kbp



ATTA GATT TACA TTAC

Pop Quiz I

- ATTA: ATT -> TTA
- GATT: GAT -> ATT
- TACA: TAC  $\rightarrow$  ACA
- TTAC: TTA -> TAC

Pop Quiz I

- ATTA: ATT  $\rightarrow$  TTA
- GATT: GAT  $\rightarrow$  ATT
- TACA: TAC  $\rightarrow$  ACA
- TTAC: TTA  $\rightarrow$  TAC

GAT 1 <u> ጉ</u>ጥ ACA

GATTACA

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

ACGA ACGT ATAC CGAC CGTA GACG GTAT TACG

Pop Quiz 2

ACGA ACGT ATAC CGAC CGTA GACG GTAT TACG



Pop Quiz 2





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



TACG







Pop Quiz 2











# Wow, this could double as life philosophy, too!

Michael Schatz @mike\_schatz

Replying to @ZaminIqbal @nomad421 and 4 others

Yep, very easy to find \*a\* path, very hard to find \*the\* path

11:40 AM - 22 Jan 2018





# Outline

- I. Assembly theory
  - Assembly by analogy

### 2. Practical Issues

- Coverage, read length, errors, and repeats
- 3. Next-next-gen Assembly
  - Canu: recommended for PacBio/ONT project
- 4. Whole Genome Alignment
  - MUMmer recommended

# **Assembly Applications**

Novel genomes





• Metagenomes





- Sequencing assays
  - Structural variations
  - Transcript assembly





# Next Steps

- I. Reflect on the magic and power of DNA  $\textcircled{\odot}$
- 2. Check out the course webpage
- 3. Register on Piazza
- 4. Work on Assignment I
  - I. Set up Linux, set up Virtual Machine
  - 2. Set up Dropbox for yourself!
  - 3. Get comfortable on the command line