



Lecture 8:

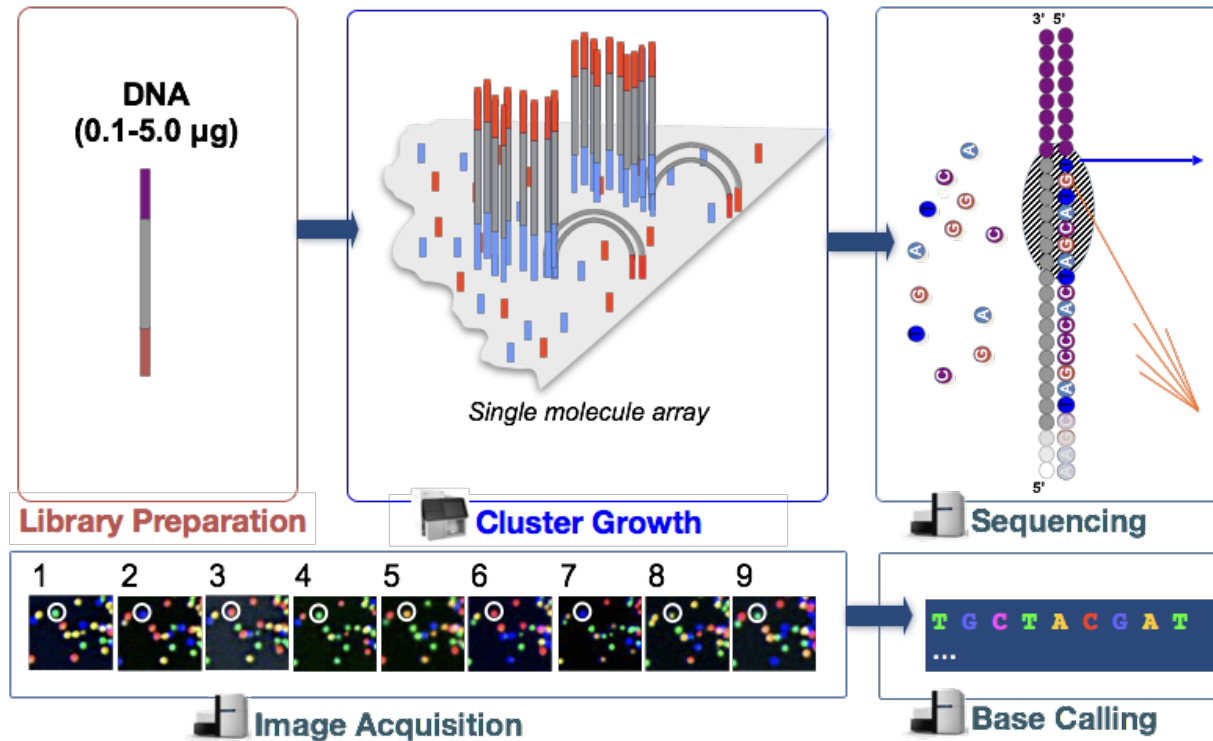
Genome assembly Stitching the pieces

Course 485

Introduction to Genomics

What is the outcome of sequencing?

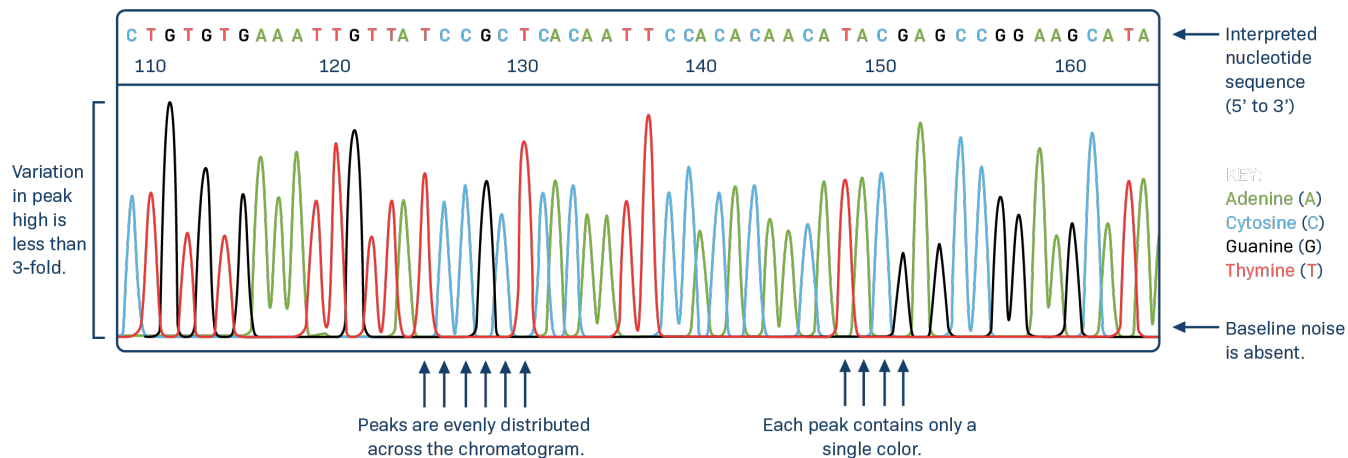
What comes out of the sequencing apparatus/machine?



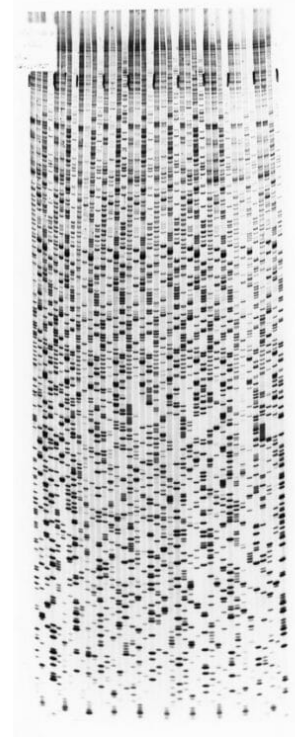
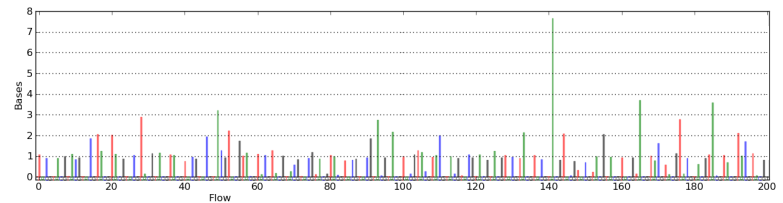
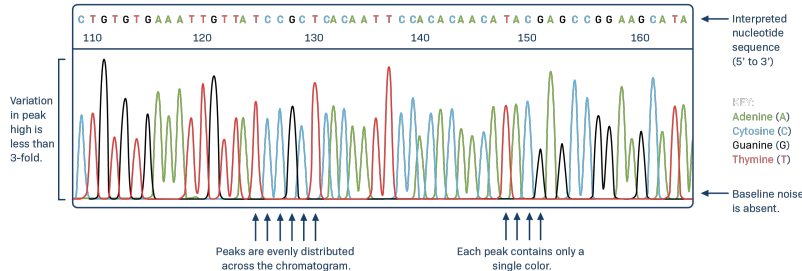
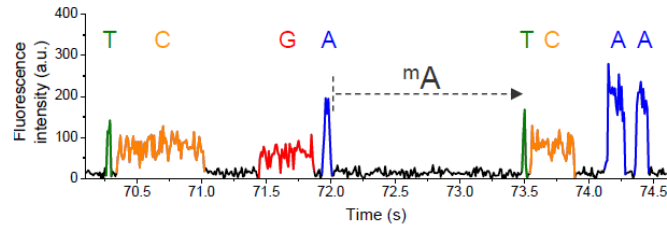
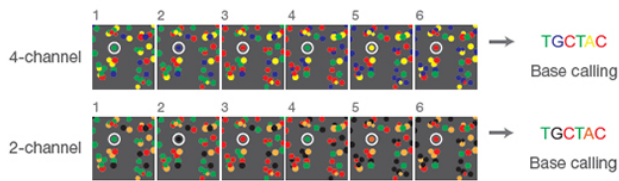
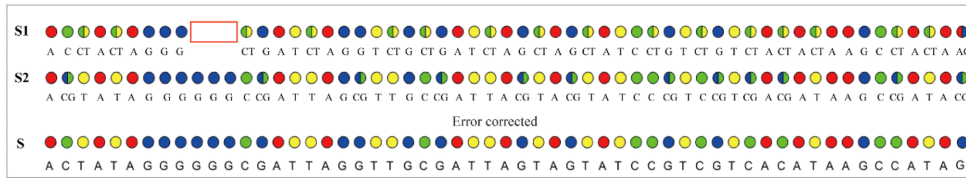
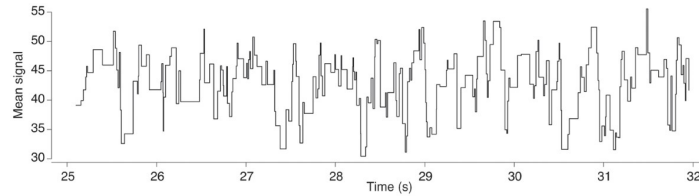
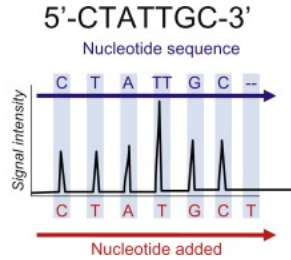
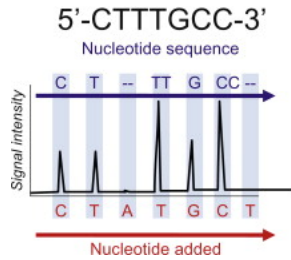
What is the outcome of sequencing?


What comes out of the sequencing apparatus/machine?

Sequence reads are the raw data and the output of the sequencing machine



Do all sequencing methods have identical sequence reads?





```
fasta files
>sequence1
ACCCATGATTTGCGA
```

```
qual files
>sequence1
40 40 39 39 40 39 40 40 40 40 20 20 36 39 39
```

```
fastq files
@sequence1
ACCCATGATTTGCGA
+
IIHHIIIIII55EHH
```

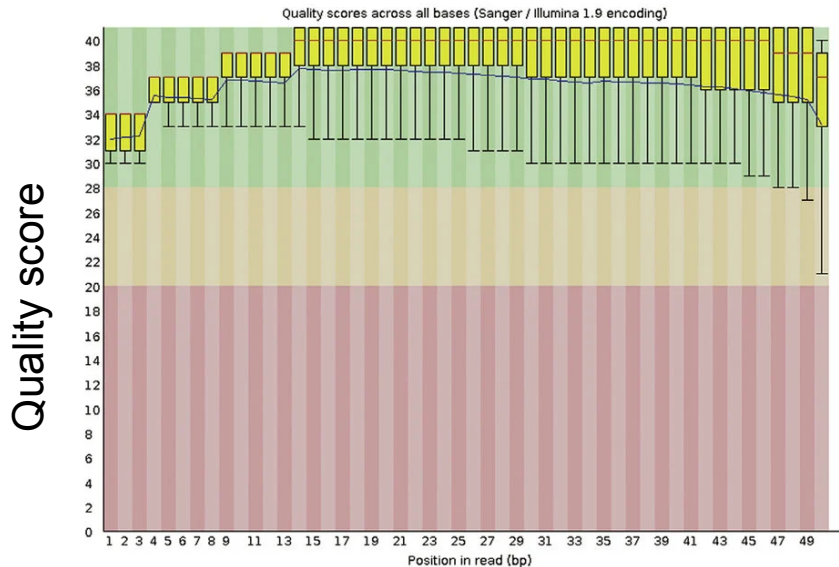
Regardless to the sequencing method used, the most important outputs are:

- 1) The identity of a nucleotide
- 2) Confidence of the call (quality)



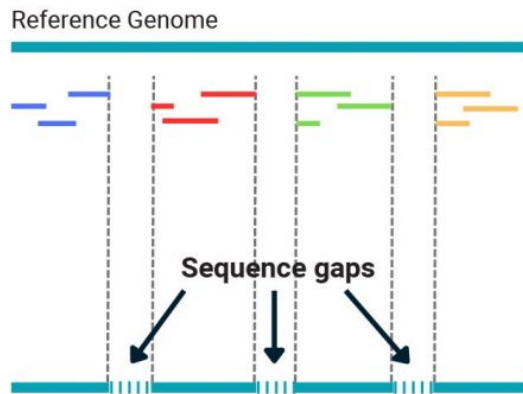
Are all nucleotides calls in a sequence read correct and trustworthy?

```
@HWI-ST193:397:D16B3ACXX:2:1101:1091:2467 1:N:0:CGATGT  
ATCACAGACAGAAGAGGATTGTACAGAGGAGCTCTTTGACTTCCTGCATC  
+  
:=:ABBDAFFDDFHIGCEEB:CEBF<+A??F3?D*?D*?B*?:?B<?)?#
```

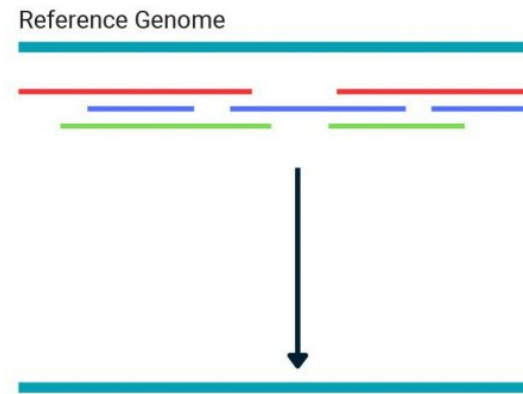


Are all sequence reads equal in length?

① SHORT READS



② LONG READS





Do you remember the average sequence read's length across sequencing methods?

B Whole-genome short read data



C Whole-genome long read data



RESEARCH ARTICLE

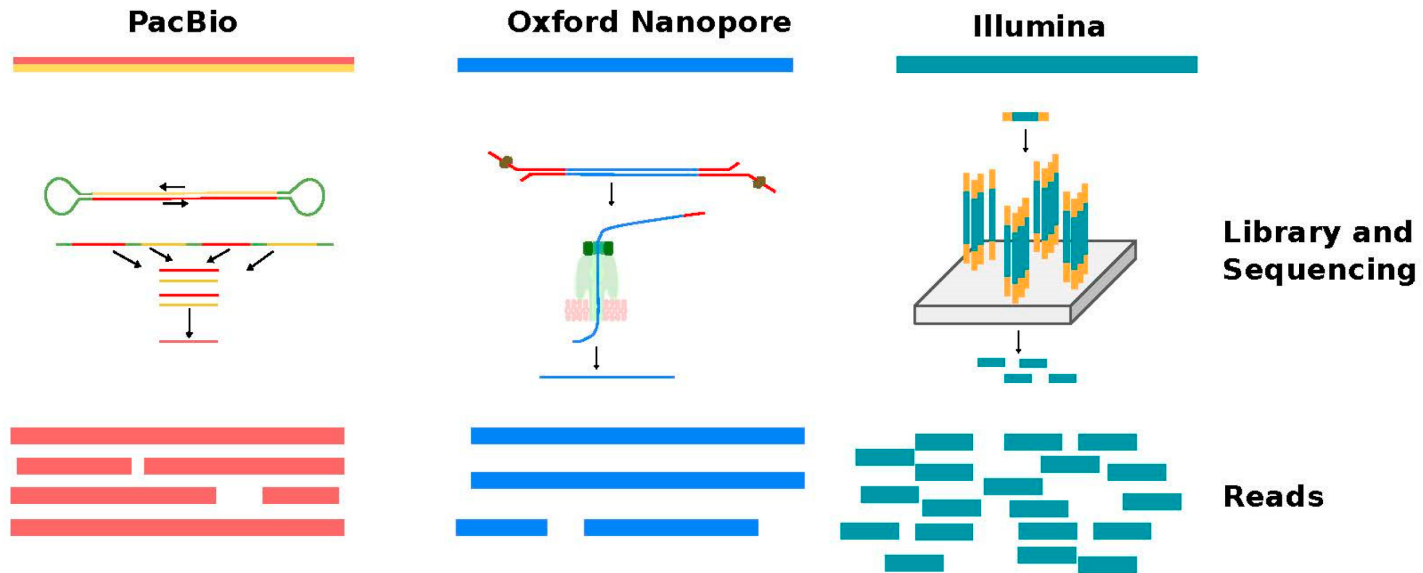
An Error Correction and DeNovo Assembly Approach for Nanopore Reads Using Short Reads

Mehdi Kchouk^{1,2} and Mourad Elloumi^{2,*}

Table 1. Characteristics of sequencing technologies.

Sequencing Technology	Average Read Length (pb)	Run Time	Advantages	Drawbacks	Error Types	Year
<i>First Generation</i>						
Sanger	~1000	~2 Hours	Long individual read length	First sequencing technology and expensive cost.	Insertion - Deletion - Substitution	1977
<i>Second Generation</i>						
Roche 454	700	24 Hours	Long read length	High error rate and expansive runs	Insertion - Deletion	2005
Illumina/Solexa	2*(36 to 100) paired-end	6 Days or 2 Days in rapid mode	Low cost and low error rate	Long run time	Substitution	2006
AB SOLID	85	8 Days	Low cost per base and low error rate	Very long run time	Substitution	2007
Helicos Biosciences	55	Depending on read length	Low error rate	Long Run time	Insertion - Deletion	2009
Ion Torrent	400	2 Hours	Short run time and Less expansive equipment	High error rate	Insertion - Deletion	2010
<i>Third Generation</i>						
Pacific Biosciences	3000 (and up to 15000)	20 Min to 4 Hours	Longest read length and fast run time	High error rate	Insertion - Deletion	2010
Oxford Nanopore (MinIon)	9545	< 6 Hours	Long read; small cost USB device	High error rate	Insertion - Deletion - Substitution	2014

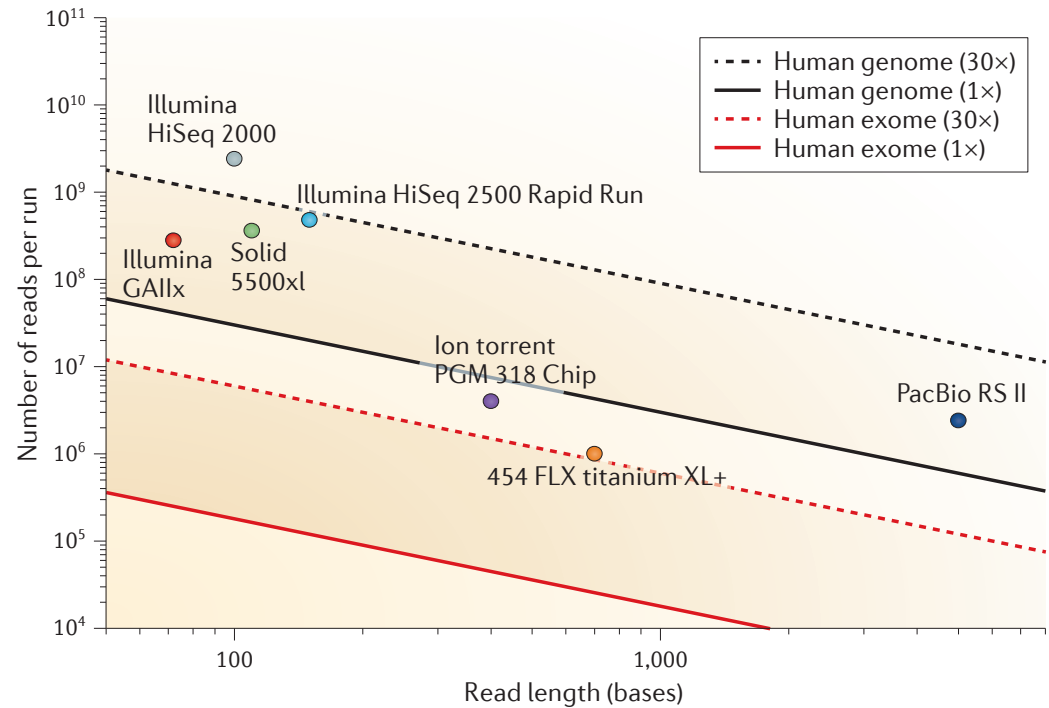
Does the number of reads resulting from sequencing a genome differs between sequencing technologies?



Does the number of reads resulting from sequencing a genome differs between sequencing technologies?

Sequencing depth and coverage: key considerations in genomic analyses

David Sims, Ian Sudbery, Nicholas E. Illott, Andreas Heger and Chris P. Ponting



Plant Communications

Research article

CellPress
Partner Journal

Pistachio genomes provide insights into nut tree domestication and ZW sex chromosome evolution



illumina[®]

495,729

100bp

50 Mb

PacBio

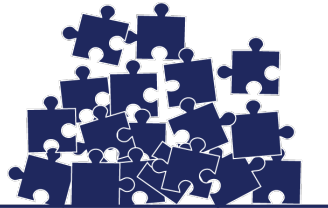
9,536

5000bp

50 Mb



What is the effect of sequence length on genome assembly?



>99.9% accuracy sequencing data for large batches and low budget per sample

Short-read sequencing

The genome is randomly fragmented in fragments of **100-600bp**

Libraries are sequenced to produce **reads up to 300bp**

NovaSeq X **MiSeq**
NextSeq **NovaSeq 6000**

VS



Continuous long sequence data with improve accesibility to AT/GC-rich regions

Long-read sequencing

The genome is randomly fragmented in fragments of **thousand of base pairs**

Libraries are sequenced to produce **reads up to 300Kbp**

Sequel IIe **Revio**
GridION



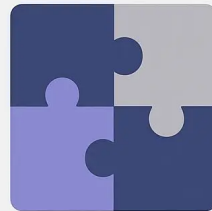
What is the total length of combined sequence reads in relation to the genome size?

Genome assembly: An analogy

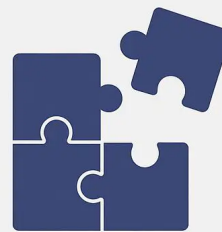




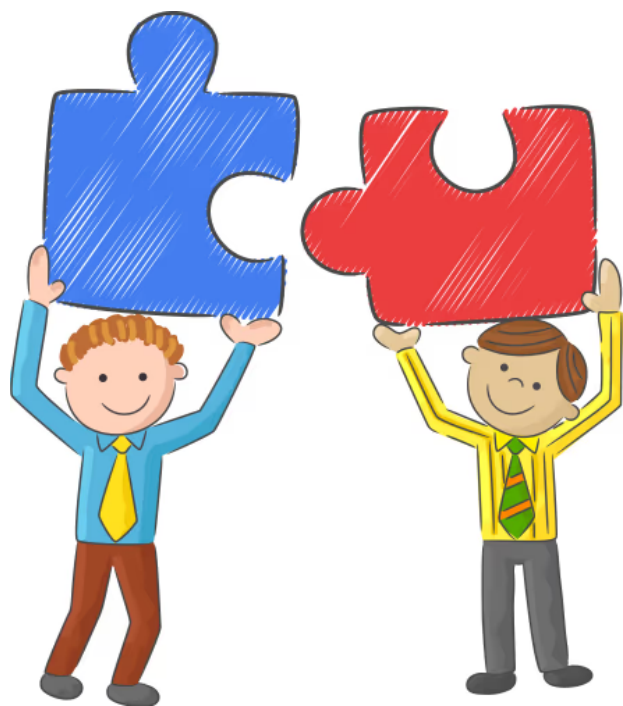
**Edge Pieces
First**



**Grouping by
Color or Pattern**



**Working on
Small Sections**





What do we need to assemble a genome under a hierarchical strategy?

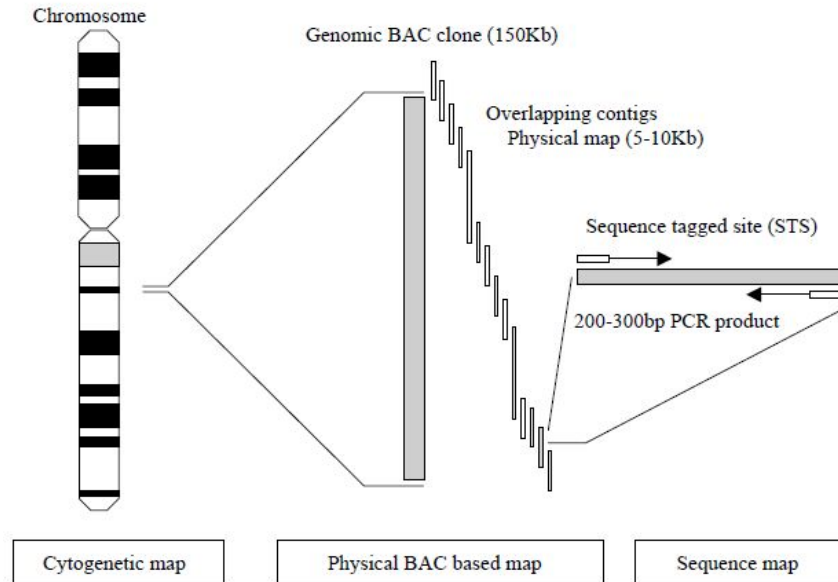
- 1) Ordered clones
- 2) Paired-end sequences of segments
- 3) Computational assembly



What do we need to assemble a genome under a whole-genome shotgun strategy?

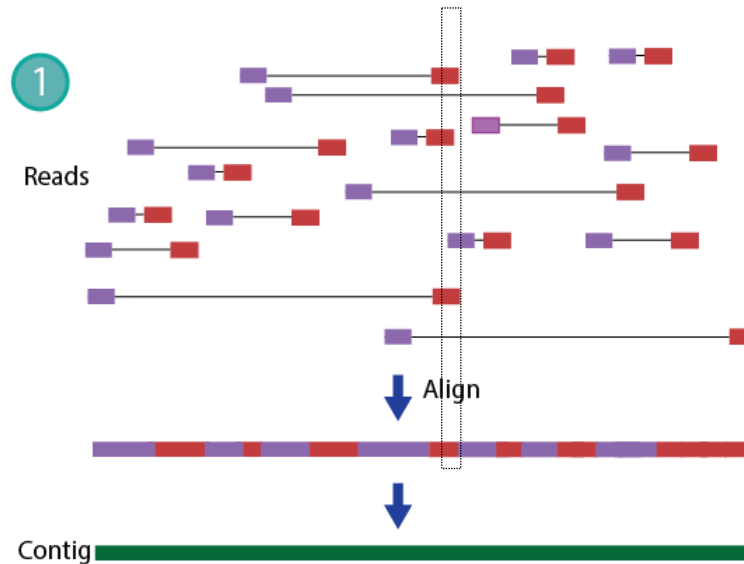
- 1) Single-end or paired-end sequences
- 2) Perform de-novo assembly

What does ordered clones and paired-end sequence provide us?



Order clones and paired-end sequences provide map to assemble large genomic fragments

How exactly can we assemble sequences or stitch sequence reads together?



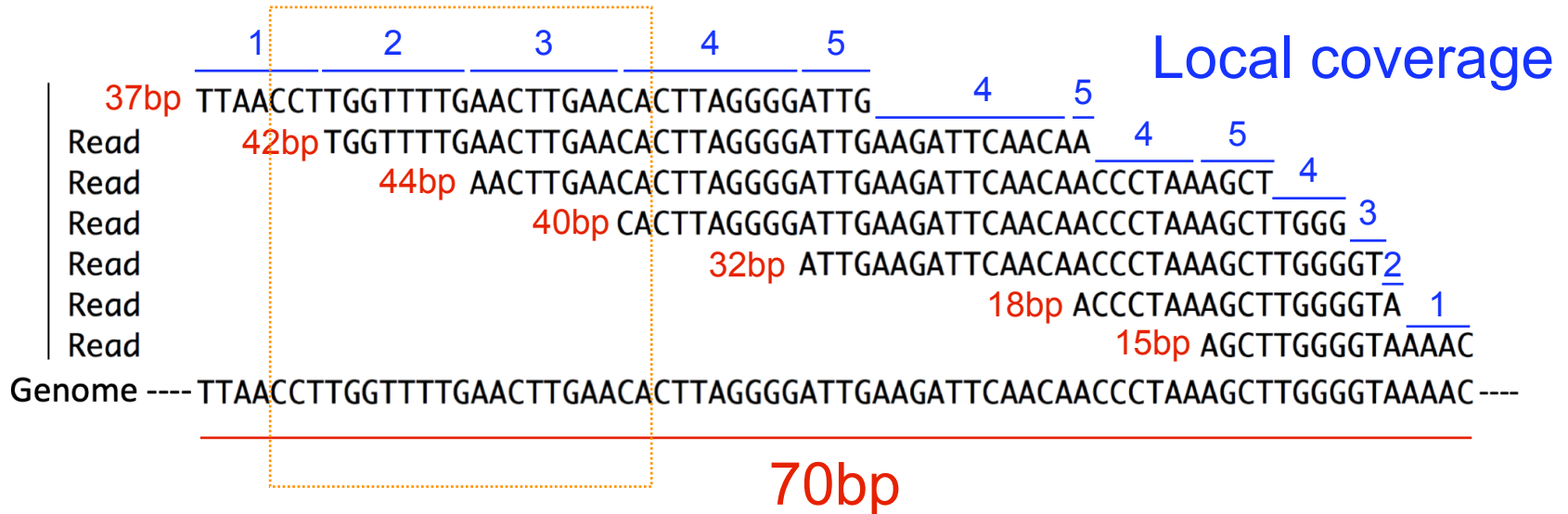
Computational assembly: looking at overlapping aligned sequence reads to get longer sequences (stitched pieces)

Sequence coverage/depth

21bp

$$1+1+1+2+2+2+2+2+2+3+3+3+3+3+3+3+3+4+4=47$$

2.24X



$$\text{Sequence coverage} = \frac{37+42+44+40+32+18+15\text{bp}}{70\text{bp}} = \frac{228\text{bp}}{70\text{bp}} = 3.26\text{X}$$

Sequence coverage/depth

The size of the genome of X species is 96Mb. Its genome was sequenced using Illumina sequencing generated a total of 480 Mb. What is the sequence coverage of X genome?

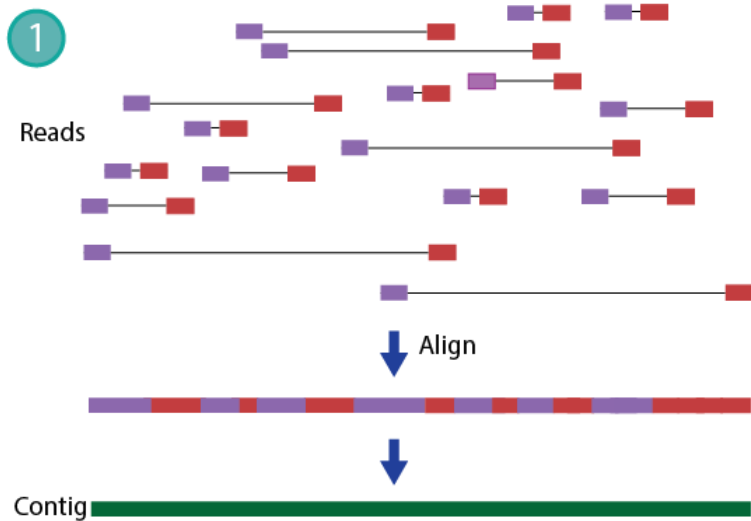
$$\text{Sequence coverage} = \frac{\text{Total sequences generated}}{\text{size of the genome}} = \frac{480\text{Mb}}{96\text{Mb}} = 5X$$



What do we need to align sequence reads and stitch them together?

- 1) A method of alignment
- 2) a test of the significance of the alignment

Information theory!



Using a method of alignment, we order and connect overlapping sequence reads.

Contigs: the result of ordering, aligning, and connecting sequence read.

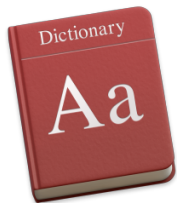
What is a contig?

con·tig·u·ous | kən'tigyəwəs |

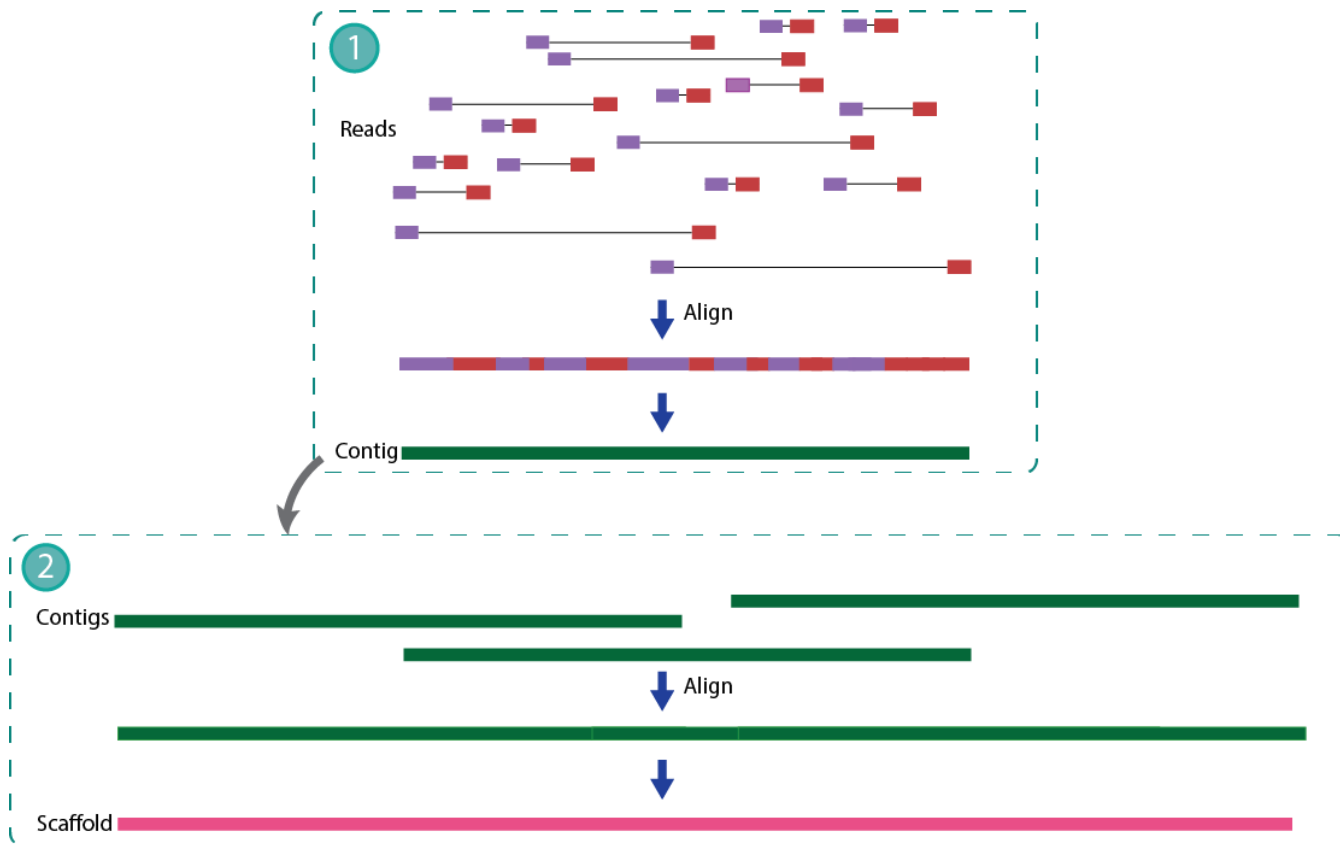
adjective

sharing a common border; touching: *the 48 contiguous states.*

- next or together in sequence: *five hundred contiguous dictionary entries.*



Scaffolds: the result of ordering, aligning, and connecting contigs

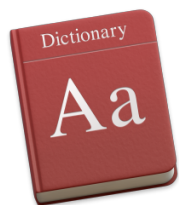


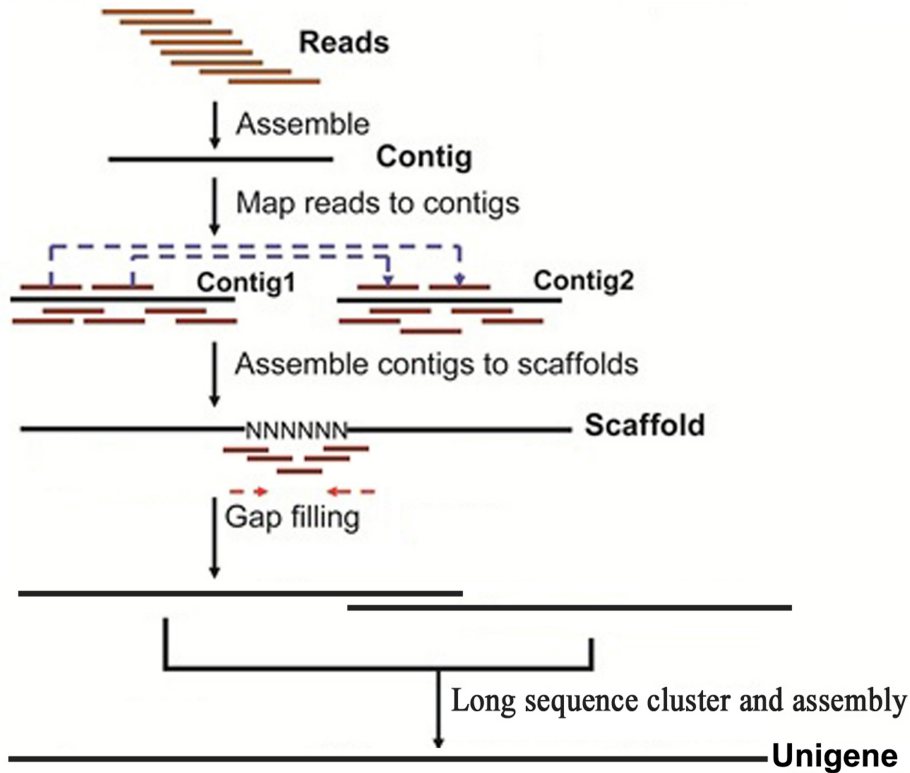
What is a scaffold?

scaf·fold | 'skafəld |

noun

- 1 a raised wooden platform used formerly for the public execution of criminals.
- 2 a structure made using scaffolding: *[as modifier] : scaffold boards.*





How do we connect disconnected contigs?

Use molecular markers that order the contigs or identify BAC clone the fills the gap



How many contigs are needed to assemble a genome?

It depends on,

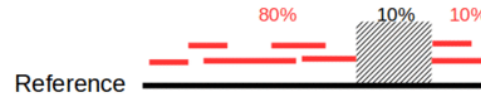
- 1) Type of the genome
- 2) Type of sequencing technology



How can we assess the quality of a genome assembly?

- 1) Completeness
- 2) Correctness
- 3) Contiguity

Genome sequence completeness



$$C = \frac{\# \text{ area covered by reads}}{\# \text{ reference area}}$$

RESEARCH ARTICLE

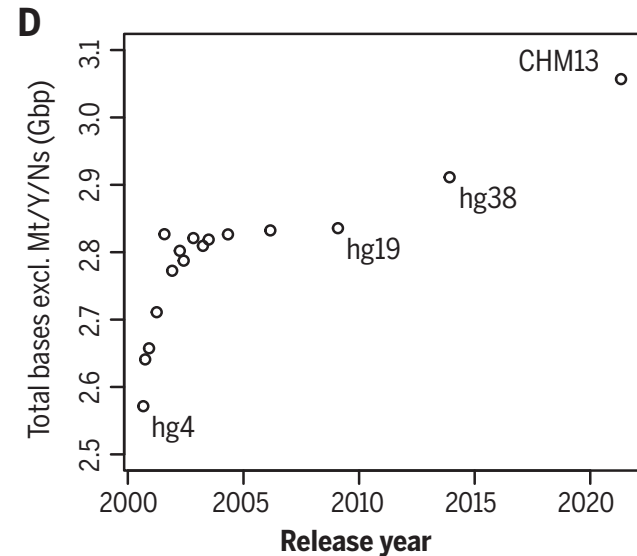
HUMAN GENOMICS

The complete sequence of a human genome

Sergey Nurk^{1†}, Sergey Koren^{1†}, Arang Rhie^{1†}, Mikko Rautiainen^{1†}, Andrey V. Bzikadze², Alla Mikheenko³, Mitchell R. Vollger⁴, Nicolas Altemose⁵, Lev Uralsky^{6,7}, Ariel Gershman⁸, Sergey Aganezov^{9‡}, Savannah J. Hoyt¹⁰, Mark Diekhans¹¹, Glennis A. Logsdon⁴, Michael Alonge⁹, Stylianos E. Antonarakis¹², Matthew Borchers¹³, Gerard G. Bouffard¹⁴, Shelise Y. Brooks¹⁴, Gina V. Caldas¹⁵, Nae-Chyun Chen⁹, Haoyu Cheng^{16,17}, Chen-Shan Chin¹⁸, William Chow¹⁹, Leonardo G. de Lima¹³, Philip C. Dishuck⁴, Richard Durbin^{19,20}, Tatiana Dvorkina³, Ian T. Fiddes²¹, Giulio Formenti^{22,23}, Robert S. Fulton²⁴, Arkarachai Fungtammasan¹⁸, Erik Garrison^{11,25}, Patrick G. S. Grady¹⁰, Tina A. Graves-Lindsay²⁶, Ira M. Hall²⁷, Nancy F. Hansen²⁸, Gabrielle A. Hartley¹⁰, Marina Haukness¹¹, Kerstin Howe¹⁹, Michael W. Hunkapiller²⁹, Chirag Jain^{1,30}, Miten Jain¹¹, Erich D. Jarvis^{22,23}, Peter Kerpedjiev³¹, Melanie Kirsche⁹, Mikhail Kolmogorov³², Jonas Korlach²⁹, Milinn Kremitzki²⁶, Heng Li^{16,17}, Valerie V. Maduro³³, Tobias Marschal³⁴, Ann M. McCartney⁴, Jennifer McDaniel³⁵, Danny E. Miller^{4,36}, James C. Mullikin^{14,28}, Eugene W. Myers³⁷, Nathan D. Olson³⁵, Benedict Paten¹¹, Paul Peluso²⁹, Pavel A. Pevzner³², David Porubsky⁴, Tamara Potapova¹³, Evgeny I. Rogae^{6,7,38,39}, Jeffrey A. Rosenfeld⁴⁰, Steven L. Salzberg^{9,41}, Valerie A. Schneider⁴², Fritz J. Sedlazeck⁴³, Kishwar Shafin¹¹, Colin J. Shew⁴⁴, Alaina Shumate⁴¹, Ying Sims¹⁹, Arian F. A. Smit⁴⁵, Daniela C. Soto⁴⁴, Ivan Sovic^{29,46}, Jessica M. Storer⁴⁵, Aaron Streets^{5,47}, Beth A. Sullivan⁴⁸, Françoise Thibaud-Nissen⁴², James Torrance¹⁹, Justin Wagner³⁵, Brian P. Walenz²⁹, Jonathan M. D. Wood¹⁹, Chunlin Xiao⁴², Stephanie M. Yan⁴⁹, Alice C. Young¹⁴, Samantha Zarate⁹, Urvashi Surti⁵⁰, Rajiv C. McCoy⁴⁹, Megan Y. Dennis⁴⁴, Ivan A. Alexandrov^{3,7,51}, Jennifer L. Gerton^{13,52}, Rachel J. O'Neill¹⁰, Winston Timp^{8,41}, Justin M. Zook³⁵, Michael C. Schatz^{9,49}, Evan E. Eichler^{4,53}, Karen H. Miga^{11,54}, Adam M. Phillippy^{1*}

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion-base pair sequence of a human genome, T2T-CHM13, that includes gapless assemblies for all chromosomes except Y, corrects errors in the prior references, and introduces nearly 200 million base pairs of sequence containing 1956 gene predictions, 99 of which are predicted to be protein coding. The completed regions include all centromeric satellite arrays, recent segmental duplications, and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies.

STATISTICS	GRCH38	T2T-CHM13	DIFFERENCE (±%)
Summary			
Assembled bases (Gbp)	2.92	3.05	+4.5
Unplaced bases (Mbp)	11.42	0	-100.0
Gap bases (Mbp)	120.31	0	-100.0
Number of contigs	949	24	-97.5
Contig NG50 (Mbp)	56.41	154.26	+173.5
Number of issues	230	46	-80.0
Issues (Mbp)	230.43	8.18	-96.5

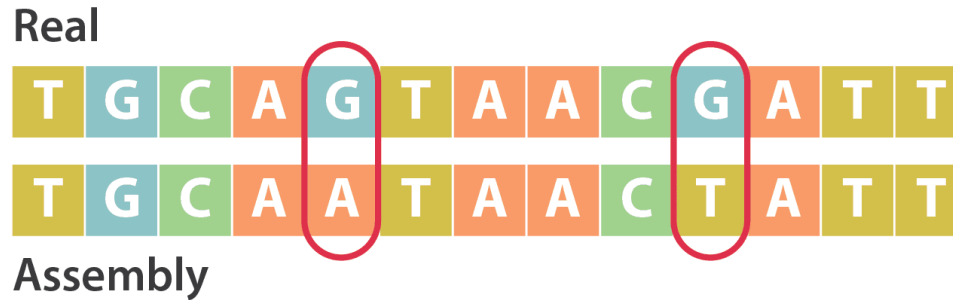


Genome sequence completeness

The genome of X species is 2.65Gb and was sequenced using pyrosequencing resulting in an assembled genome of 2.3Gb. How much of the genome was covered with sequencing?

$$\% \text{ genome sequenced} = \frac{2.30\text{Gb}}{2.65\text{Gb}} \times 100 = 86.8\%$$

Genome sequence correctness

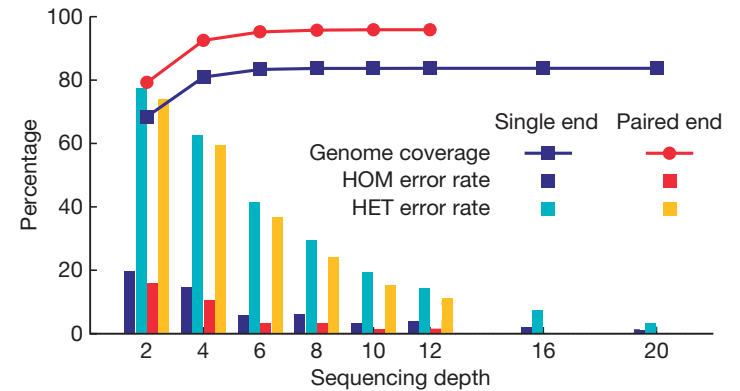


nature


Vol 456 | 6 November 2008 | doi:10.1038/nature07484

ARTICLES

The diploid genome sequence of an Asian individual



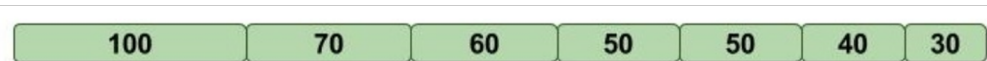
Genome sequence correctness



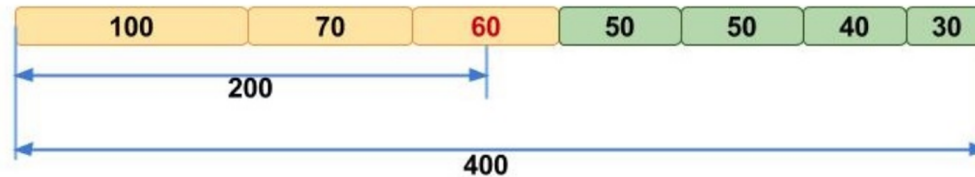
```
Read  TTAACCTTGGTTTTGAACTTGAACACTTAGGGGATTG
Read      TGGTTTTGAACTTGAACACTTAGGGGATTGAAGATTCAACAA
Read          AACTTGAACACTTAGGGGATTGAAGATTCAACAACCCTAAAGCT
Read              TACTTAGGGGATTGAAGATTCAACAACCCTAAAGCTTGGA
Read                  GTTGAAGATTCAACAACCCTAAAGCTTGGAGT
Read                      GCCCTAAAGCTTGGGGTA
Read                          AGCTTGGGGTAAAAC
Genome ---- TTAACCTTGGTTTTGAACTTGAACACTTAGGGGATTGAAGATTCAACAACCCTAAAGCTTGGRGTAAAAC ----
```

How can we assess the contiguity of a genome assembly?

- 1) Contig N50 size
- 2) Scaffold N50 size



1a. Contigs, sorted according to their lengths.



1b. Calculation of N50 using sorted contigs.

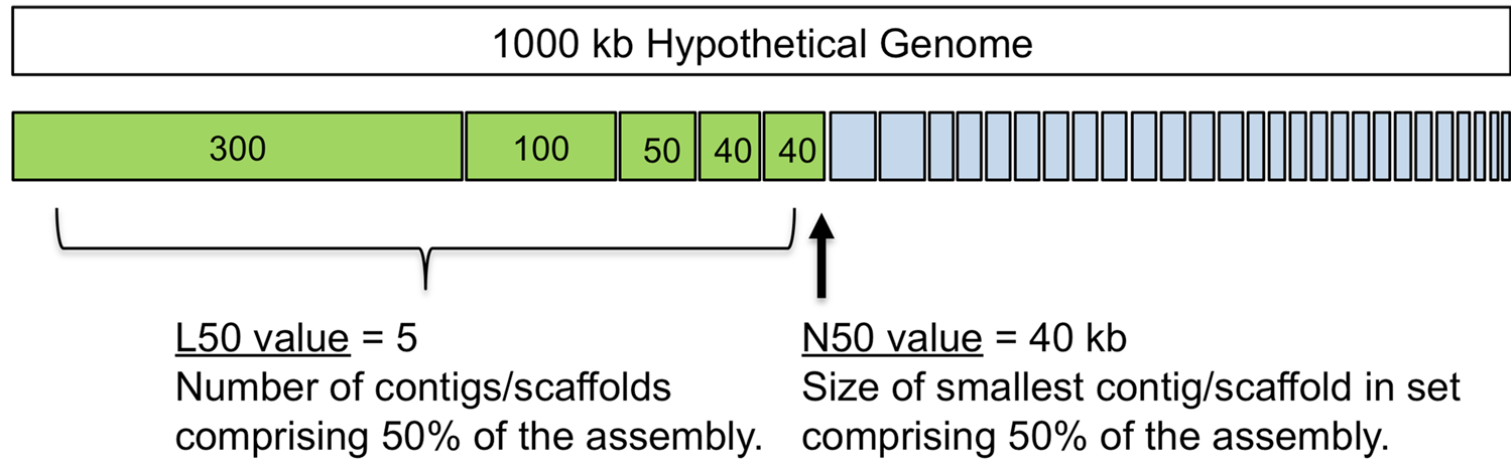
These measures are arbitrary measures and statistical parameters such as the mean. The measures are used only to compare between different assemblies.

Contig/Scaffold N50: the size of the contig/scaffold at which 50% of the genome is assembled into contigs/scaffold

Genome projects in invasion biology

Michael A. McCartney¹ · Sophie Mallez¹ · Daryl M. Goh^{2,3}

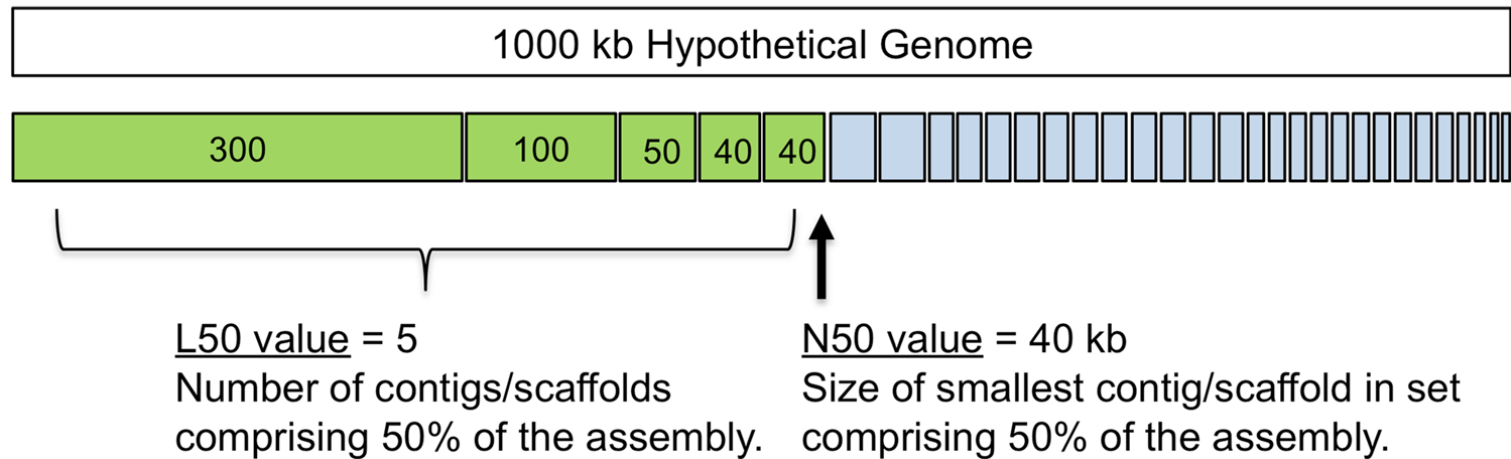
Contiguity Measures




How to calculate/find the contig N50 or L50?

- 1) Order contigs/scaffolds according to size in a descending order
- 2) Identify the number/size of the contig where 50% of the genome is assembled.

Contiguity Measures





A group of three great students (Sara, Nouf, and Hayat) have sequenced the genome of a local ant. The sequencing machine generated **4,105,387** sequence reads, which accounted to a total of **480 Mb**. Below are the details of the initial assembled contigs.

<u>Contig size</u>	<u># of Contig</u>
150 bp	6,000
200 bp	5,750
1 Kb	1,200
5 Kb	50
10 Kb	110
100 Kb	59
500 Kb	5
1 Mb	4
5 Mb	3

The three students decided to assemble the genome independently.

Sara used ALL sequences.

Nouf used sequences > 150 bp.

Hayat used sequences > 200 bp.

A group of three great students (Sara, Nouf, and Hayat) have sequenced the genome of a local ant. The sequencing machine generated **4,105,387** sequence reads, which accounted to a total of **480 Mb**. Below are the details of the initial assembled contigs.

<u>Contig size</u>	<u># of Contig</u>
150 bp	6,000
200 bp	5,750
1 Kb	1,200
5 Kb	50
10 Kb	110
100 Kb	59
500 Kb	5
1 Mb	4
5 Mb	3

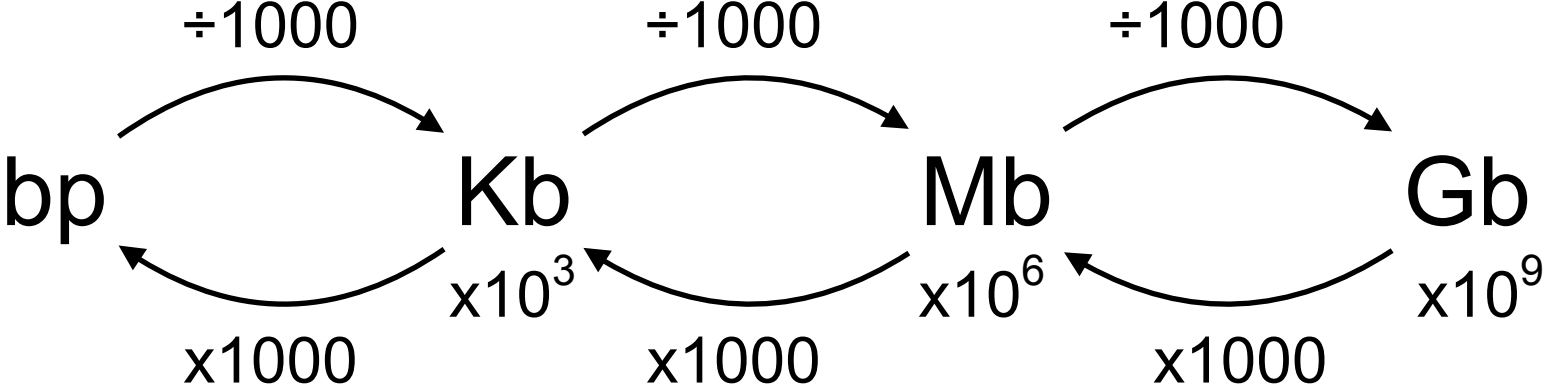
The three students decided to assemble the genome independently.

Sara used ALL sequences.

Nouf used sequences > 150 bp.

Hayat used sequences > 200 bp.

Conversion of genome length units



What does the contigs' size categories and numbers mean?

<u>Contig size</u>	<u># of Contig</u>
150 bp	6,000
200 bp	5,750
1 Kb	1,200
5 Kb	50
10 Kb	110
100 Kb	59
500 Kb	5
1 Mb	4
5 Mb	3

What does 150bp and 200bp contigs mean?

What the most likely sequencing technology used and why?

What is the contig number for each student's genome assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750				
1 Kb	1,200				
5 Kb	50				
10 Kb	110				
100 Kb	59				
500 Kb	5	13,181			
1 Mb	4	7,181			
5 Mb	3	1,431			

What is the assembled genome size for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u>Combined length</u>
150 bp	6,000	$150\text{bp} \times 6,000 = 900,000\text{bp} = 900\text{Kb} = 0.9\text{Mb}$
200 bp	5,750	$200\text{bp} \times 5,750 = 1,150,000\text{bp} = 1,150\text{Kb} = 1.15\text{Mb}$
1 Kb	1,200	$1\text{Kb} \times 1,200 = 1,200\text{Kb} = 1.2\text{Mb}$
5 Kb	50	$5\text{Kb} \times 50 = 250\text{Kb} = 0.25\text{Mb}$
10 Kb	110	$10\text{Kb} \times 110 = 1,100\text{Kb} = 1.1\text{Mb}$
100 Kb	59	$100\text{Kb} \times 59 = 5,900\text{Kb} = 5.9\text{Mb}$
500 Kb	5	$500\text{Kb} \times 5 = 2,500\text{Kb} = 2.5\text{Mb}$
1 Mb	4	$1\text{Mb} \times 4 = 4\text{Mb}$
5 Mb	3	$5\text{Mb} \times 3 = 15\text{Mb}$

29.95Mb

31.1Mb

32Mb



What is the assembled genome size for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200				
5 Kb	50				
10 Kb	110				
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				

What is 50% the size of the assembled genome for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u>Combined length</u>
150 bp	6,000	$150\text{bp} \times 6,000 = 900,000\text{bp} = 900\text{Kb} = 0.9\text{Mb}$
200 bp	5,750	$200\text{bp} \times 5,750 = 1,150,000\text{bp} = 1,150\text{Kb} = 1.15\text{Mb}$
1 Kb	1,200	$1\text{Kb} \times 1,200 = 1,200\text{Kb} = 1.2\text{Mb}$
5 Kb	50	$5\text{Kb} \times 50 = 250\text{Kb} = 0.25\text{Mb}$
10 Kb	110	$10\text{Kb} \times 110 = 1,100\text{Kb} = 1.1\text{Mb}$
100 Kb	59	$100\text{Kb} \times 59 = 5,900\text{Kb} = 5.9\text{Mb}$
500 Kb	5	$500\text{Kb} \times 5 = 2,500\text{Kb} = 2.5\text{Mb}$
1 Mb	4	$1\text{Mb} \times 4 = 4\text{Mb}$
5 Mb	3	$5\text{Mb} \times 3 = 15\text{Mb}$

29.95Mb 31.1Mb 32Mb

50% genome = genome size x 0.5 = **14.975Mb 15.55Mb 16Mb**



What is 50% the size of the assembled genome for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50				
10 Kb	110				
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				

What is the mean contig length for each student's assembly?

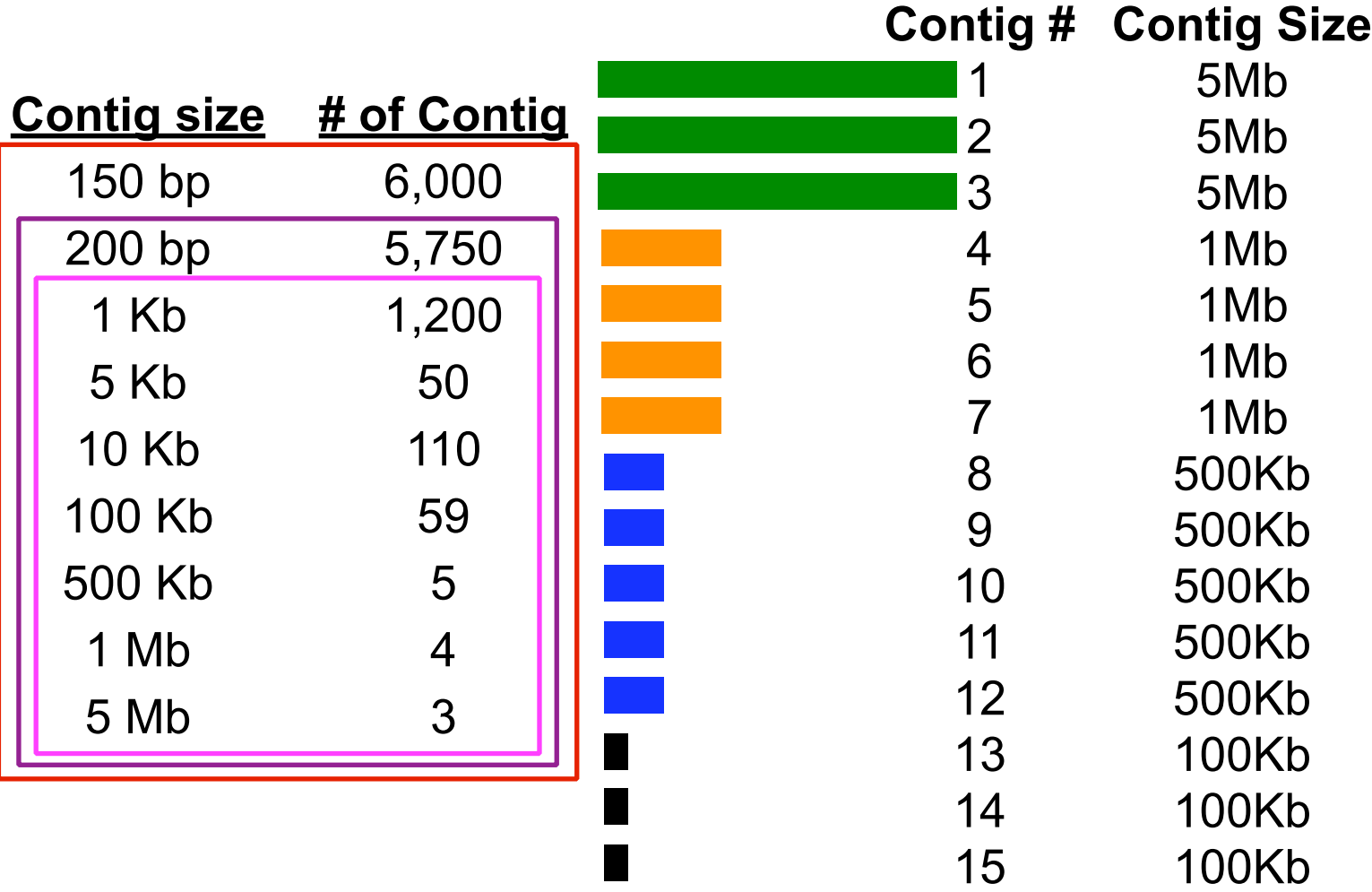
<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50				
10 Kb	110				
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				

$$\text{Mean contig length} = \frac{\text{Genome size}}{\# \text{ contigs}} = \frac{32\text{Mb}}{13,181} = 0.00242\text{Mb} = 2.42\text{Kb}$$
















What is the mean contig length for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50	Mean contig length	2.42KB	4.33Kb	20.92Kb
10 Kb	110				
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				

What is the contig N50 for each student's assembly?



What is the contig N50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>		Contig #	Contig Size	Sum
			1	5Mb	5Mb
			2	5Mb	10Mb
			3	5Mb	15Mb
			4	1Mb	16Mb
			5	1Mb	17Mb
			6	1Mb	18Mb
			7	1Mb	19Mb
			8	500Kb	19.5Mb
			9	500Kb	20Mb
			10	500Kb	20.5Mb
			11	500Kb	21Mb
			12	500Kb	21.5Mb
			13	100Kb	21.6Mb
			14	100Kb	
			15	100Kb	

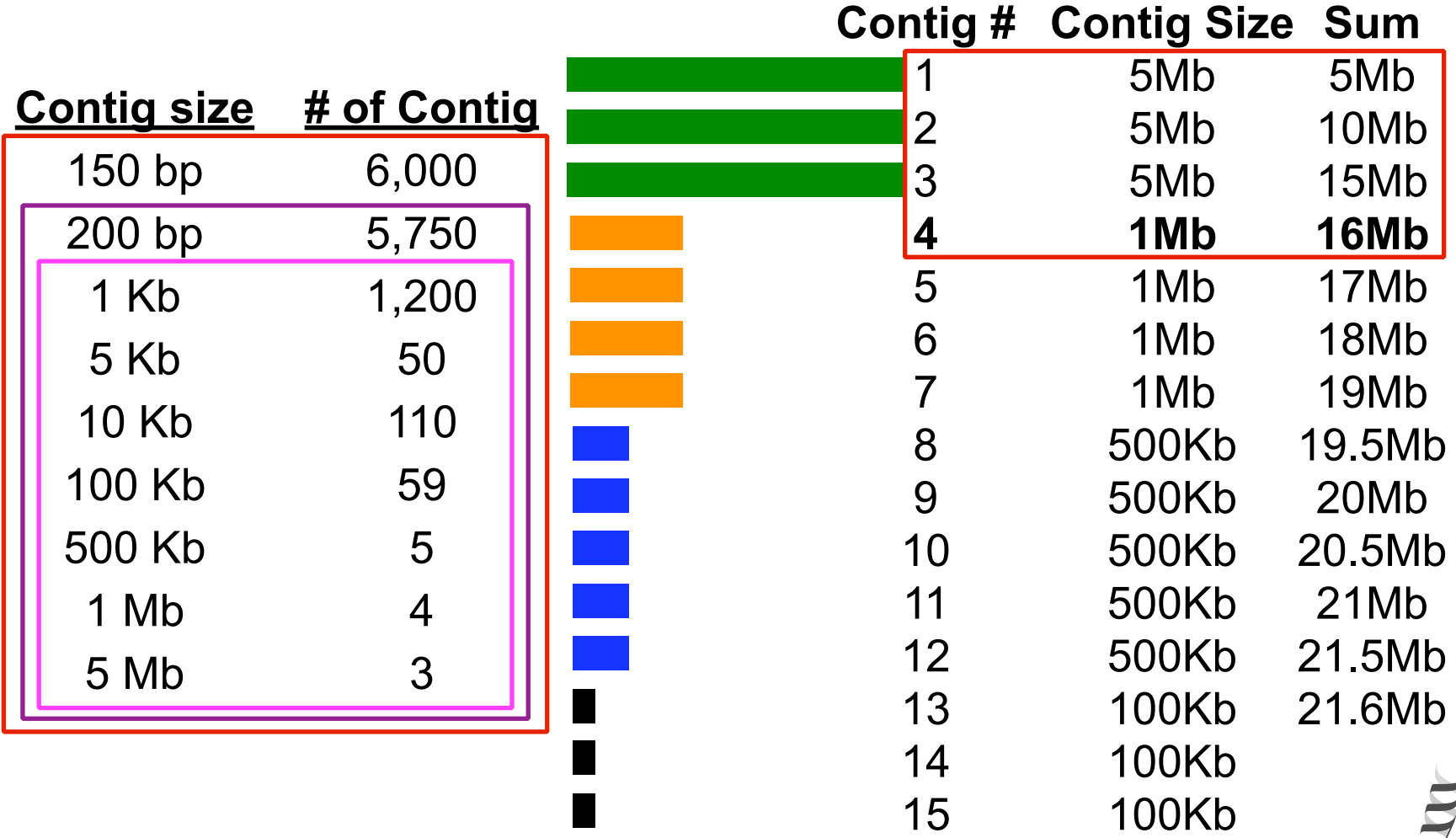
150 bp	6,000
200 bp	5,750
1 Kb	1,200
5 Kb	50
10 Kb	110
100 Kb	59
500 Kb	5
1 Mb	4
5 Mb	3


















What is the contig N50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750				
1 Kb	1,200				
5 Kb	50				
10 Kb	110				
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				
		Genome size	32Mb	31.1Mb	29.95Mb
		50% Genome	16Mb	15.55Mb	14.975Mb
		Mean contig length	2.42KB	4.33Kb	20.92Kb

What is the contig N50 for each student's assembly?


















What is the contig N50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>		Contig #	Contig Size	Sum
150 bp	6,000		1	5Mb	5Mb
200 bp	5,750		2	5Mb	10Mb
1 Kb	1,200		3	5Mb	15Mb
5 Kb	50		4	1Mb	16Mb
10 Kb	110		5	1Mb	17Mb
100 Kb	59		6	1Mb	18Mb
500 Kb	5		7	1Mb	19Mb
1 Mb	4		8	500Kb	19.5Mb
5 Mb	3		9	500Kb	20Mb
			10	500Kb	20.5Mb
			11	500Kb	21Mb
			12	500Kb	21.5Mb
			13	100Kb	21.6Mb
			14	100Kb	
			15	100Kb	



What is the contig N50 for each student's assembly?
















<u>Contig size</u>	<u># of Contig</u>		Contig #	Contig Size	Sum
150 bp	6,000		1	5Mb	5Mb
200 bp	5,750		2	5Mb	10Mb
1 Kb	1,200		3	5Mb	15Mb
5 Kb	50		4	1Mb	16Mb
10 Kb	110		5	1Mb	17Mb
100 Kb	59		6	1Mb	18Mb
500 Kb	5		7	1Mb	19Mb
1 Mb	4		8	500Kb	19.5Mb
5 Mb	3		9	500Kb	20Mb
			10	500Kb	20.5Mb
			11	500Kb	21Mb
			12	500Kb	21.5Mb
			13	100Kb	21.6Mb
			14	100Kb	
			15	100Kb	



What is the contig N50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50	Mean contig length	2.42KB	4.33Kb	20.92Kb
10 Kb	110	N50 contig	1Mb	1Mb	5Mb
100 Kb	59				
500 Kb	5				
1 Mb	4				
5 Mb	3				

What is the contig L50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>		Contig #	Contig Size	Sum
			1	5Mb	5Mb
			2	5Mb	10Mb
			3	5Mb	15Mb
			4	1Mb	16Mb
			5	1Mb	17Mb
			6	1Mb	18Mb
			7	1Mb	19Mb
			8	500Kb	19.5Mb
			9	500Kb	20Mb
			10	500Kb	20.5Mb
			11	500Kb	21Mb
			12	500Kb	21.5Mb
			13	100Kb	21.6Mb
			14	100Kb	
			15	100Kb	



What is the contig L50 for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000	Genome size	13,181	7,181	1,431
200 bp	5,750	50% Genome	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	Mean contig length	16Mb	15.55Mb	14.975Mb
5 Kb	50	N50 contig	2.42KB	4.33Kb	20.92Kb
10 Kb	110	L50 contig	1Mb	1Mb	5Mb
100 Kb	59		4	4	3
500 Kb	5				
1 Mb	4				
5 Mb	3				

What is the sequence coverage for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000	Genome size	13,181	7,181	1,431
200 bp	5,750	50% Genome	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	Mean contig length	16Mb	15.55Mb	14.975Mb
5 Kb	50	N50 contig	2.42KB	4.33Kb	20.92Kb
10 Kb	110	L50 contig	1Mb	1Mb	5Mb
100 Kb	59	~Coverage	4	4	3
500 Kb	5				
1 Mb	4				
5 Mb	3				
$\frac{480\text{Mb}}{32\text{Mb}} = 15x$					



What is the sequence coverage for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000		13,181	7,181	1,431
200 bp	5,750	Genome size	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50	Mean contig length	2.42KB	4.33Kb	20.92Kb
10 Kb	110	N50 contig	1Mb	1Mb	5Mb
100 Kb	59	L50 contig	4	4	3
500 Kb	5				
1 Mb	4				
5 Mb	3				
$\frac{480\text{Mb}-0.9\text{Mb}}{31.1\text{Mb}} = 15.4x$		~Coverage	15x	15.4x	



What is the sequence coverage for each student's assembly?

<u>Contig size</u>	<u># of Contig</u>	<u># Contigs</u>	Sara	Nouf	Hayat
150 bp	6,000	Genome size	13,181	7,181	1,431
200 bp	5,750		32Mb	31.1Mb	29.95Mb
1 Kb	1,200	50% Genome	16Mb	15.55Mb	14.975Mb
5 Kb	50		2.42KB	4.33Kb	20.92Kb
10 Kb	110	Mean contig length	1Mb	1Mb	5Mb
100 Kb	59		4	4	3
500 Kb	5	N50 contig	4	4	3
1 Mb	4		4	4	3
5 Mb	3	L50 contig	4	4	3
			4	4	3
$\frac{480\text{Mb}-2.05\text{Mb}}{29.95\text{Mb}} = \sim 12x$		~Coverage	15x	15.4x	12x



Who has the best assembly and why?

<u>Contig size</u>	<u># of Contig</u>	# Contigs	Sara	Nouf	Hayat
150 bp	6,000	Genome size	13,181	7,181	1,431
200 bp	5,750	50% Genome	32Mb	31.1Mb	29.95Mb
1 Kb	1,200	Mean contig length	16Mb	15.55Mb	14.975Mb
5 Kb	50	N50 contig	2.42KB	4.33Kb	20.92Kb
10 Kb	110	L50 contig	1Mb	1Mb	5Mb
100 Kb	59	~Coverage	4	4	3
500 Kb	5		15x	15.4x	12x
1 Mb	4				
5 Mb	3				



Disclaimer

Figures, photos, and graphs in my lectures are collected using google searches. I do not claim to have personally produced the material (except for some). I do cite only articles or books used. I thank all owners of the visual aid that I use and apologize for not citing each individual item. If anybody finds the inclusion of their material in my lectures a violation of their copy rights, please contact me via email.

hhalhaddad@gmail.com