

Genome Sequencing & Assembly

Deb Triant

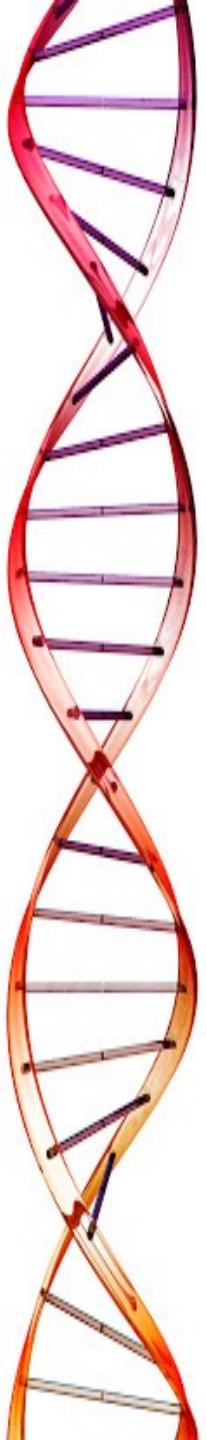
University of Virginia

Programming for Biology

Cold Spring Harbor Labs

October 2024





Lecture outline

- I. General background & theory of genome assembly
- II. Comparison of sequencing technologies
- III. Assembly quality and annotation
- IV. Bioinformatics tools for genome analysis
- V. Assembly workshop with Python programming

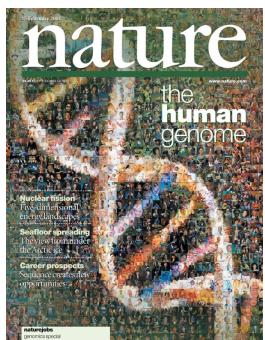
History of Genome Assembly

1977. Sanger et al. 1st Complete Organism bacteriophage 5375 bp

1995. Fleischmann et al. 1st Free Living bacteria; *Haemophilus influenzae*; TIGR Assembler. 1.8Mb

1998. *C.elegans* SC 1st Multicellular Organism BAC-by-BAC Phrap. 97Mbp

2000. *Drosophila* genome; Myers et al. 1st Large WGS Assembly Celera Assembler. 116 Mbp



Human Genome

Public: 13-year project began 1990, Dept Energy & NIH,
\$3 billion; millions of small fragments
2003 – announced as complete



Private: Craig Venter, Celera Genomics; 1998, \$300 million
Could not be patented.

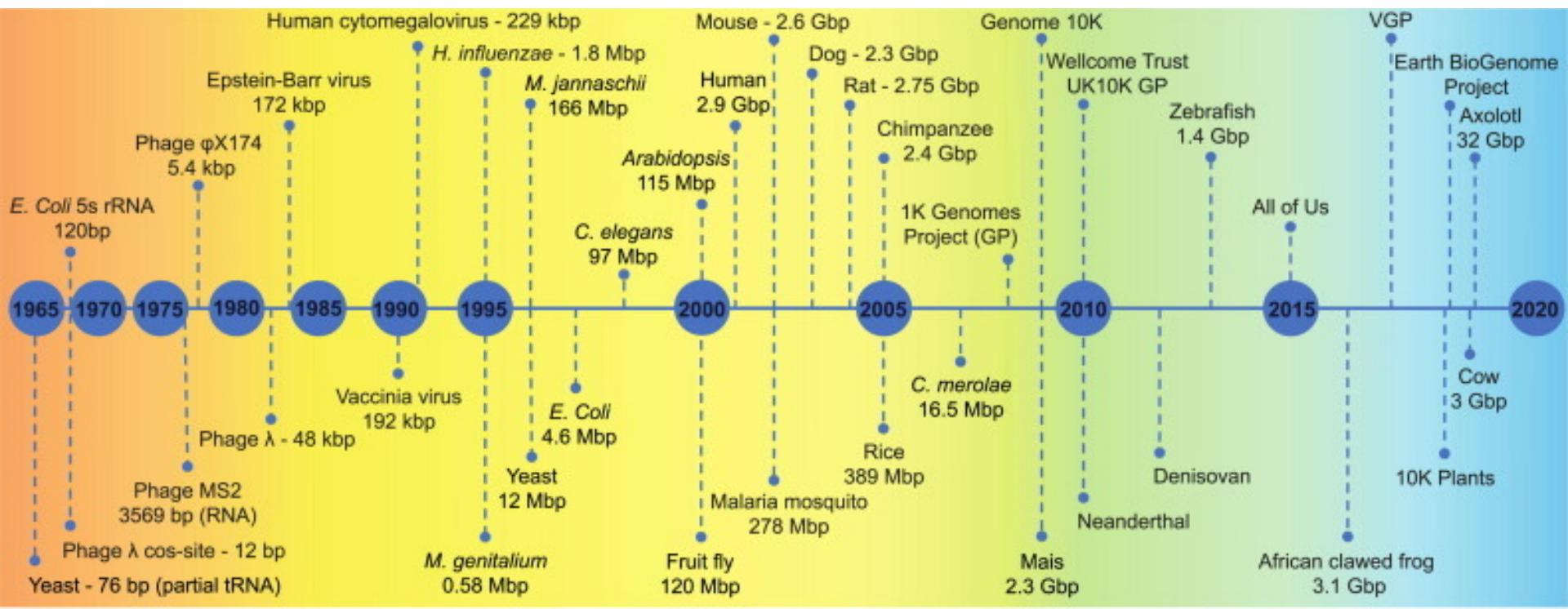
Human genome “finished” ~2003

History of Genome Assembly

- Human genome 99% finished
 - heterochromatic regions
 - centromeres, telomeres
 - tandem duplications
 - Repeats!!



History of Genome Assembly



Sanger, NGS, TGS

Giani et al. 2020

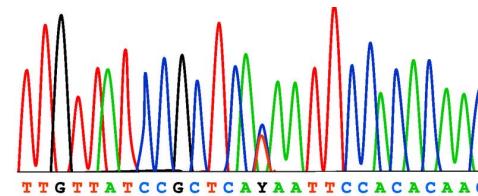
Sequencing DNA

Objective: determine sequence of nucleotides in organism or DNA molecule

- Need many copies of your sequence
- Requires fragmenting your DNA
- Sequence fragments - “reads”
- Assemble original sequence from reads
 - How many possibilities exist?

History of Genome Assembly

- First software to analyze Sanger sequencing,
Roger Staden, 1979



~800 bp reads, low error rate, costly

"If the 5' end of the sequence from one gel reading is the same as the 3' end of the sequence from another the data is said to overlap. If the overlap is of sufficient length to distinguish it from being a repeat in the sequence the two sequences must be contiguous. The data from the two gel readings can then be joined to form one longer continuous sequence."

Shredded Book Reconstruction

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

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 best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

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 best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, ...

- How can he reconstruct the text?

- Short fragments (reads) from all copies are mixed
- Proportion of reads not exact matches
- Identical reads but from different locations

It was the best of

age of wisdom, it was

best of times, it was

it was the age of

it was the age of

it was the worst of

of times, it was the

of times, it was the

of wisdom, it was the

the age of wisdom, it

the best of times, it

the worst of times, it

times, it was the age

times, it was the worst

was the age of wisdom,

was the age of foolishness,

was the best of times,

was the worst of times,

wisdom, it was the age

worst of times, it was

Join overlapping fragments

It was the best of

was the best of times,

the best of times, it

best of times, it was

of times, it was the

of times, it was the

times, it was the worst

times, it was the age

Start to join fragments based on overlaps

- How to join repetitive fragments?
- Lead to assembly misconstructions
- Build an assembly graph

Build an assembly graph

Nodes: k-1 subfragments

Edges: Directed edges between adjacent subfragments

Best match between end of one read and beginning of another

Original Fragment

Directed graph with overlaps between subfragments

It was the best of

It **was the best**

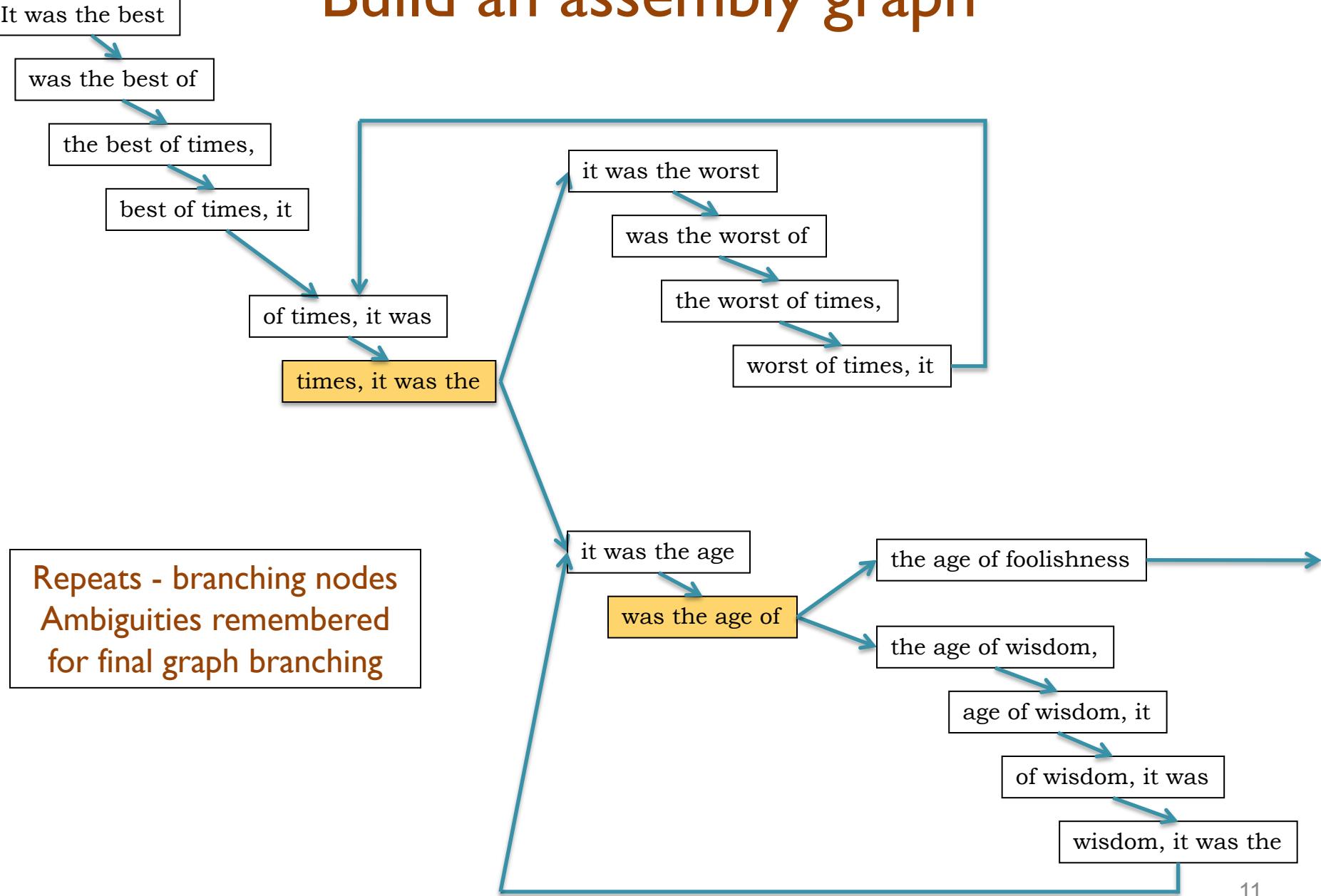


was the best of

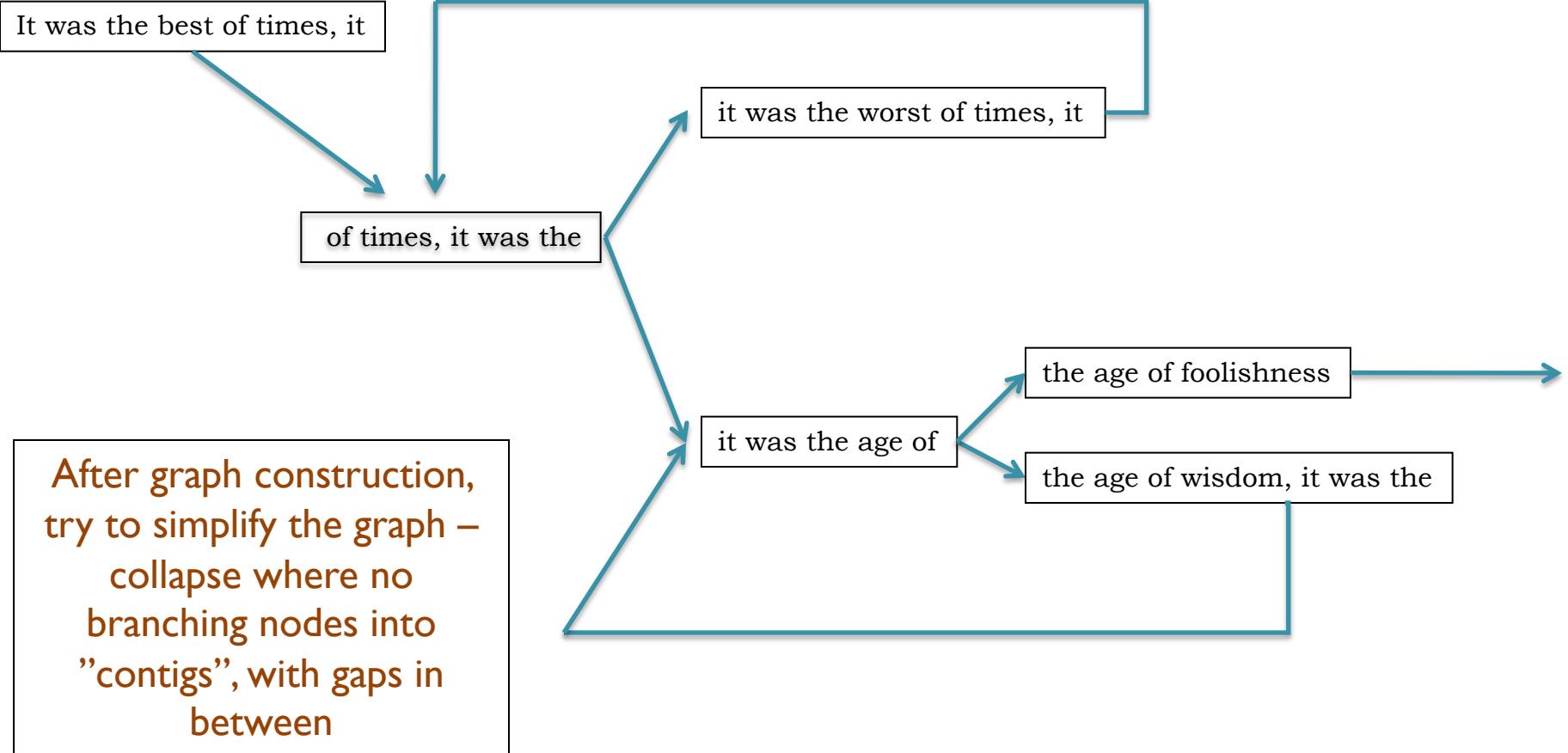
Keeps track of all possible reconstructions

- pairs of reads that share common subfragments (kmers)
- identify shortest path between overlaps

Build an assembly graph



Building contigs

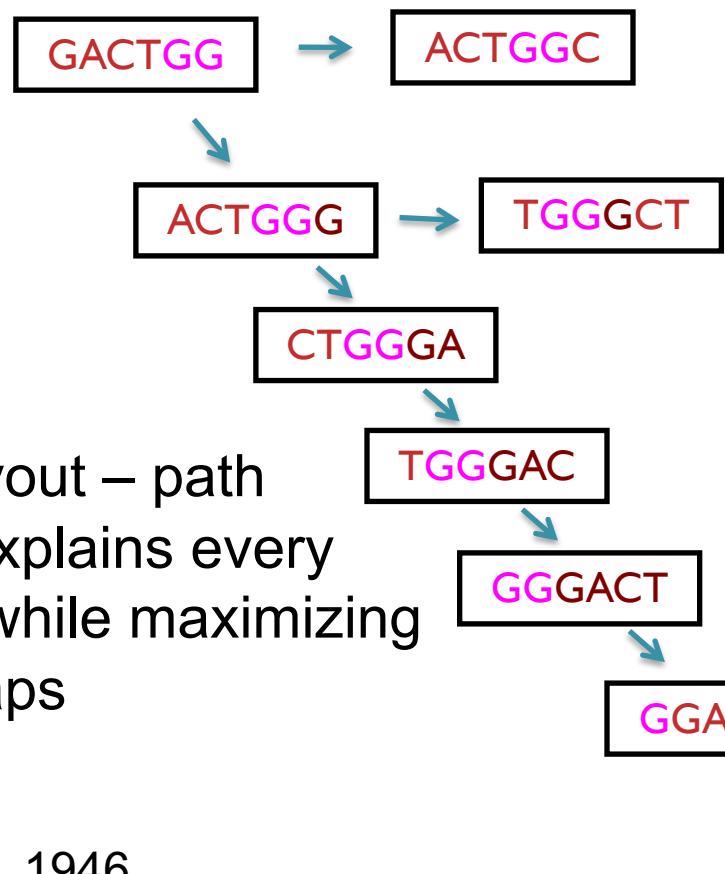


Build an assembly graph

I. Overlap

- Shredded words → k-mer

GACTGGGGACTCC



k-mer - substring of length k
 $k = 6$

2. Layout – path
that explains every
read while maximizing
overlaps

Graph Assembly

3. Consensus

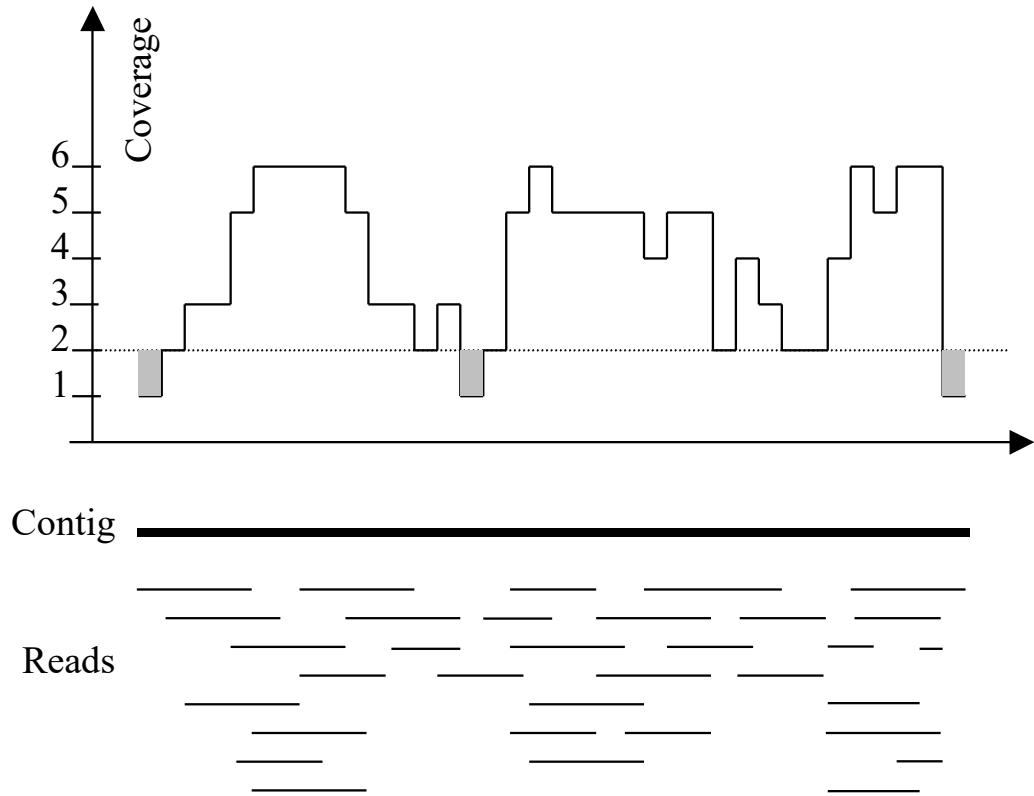
- sequence alignments to represent the same position

G TGGGACTCCGGCATTAGC TA
GA~~T~~CAGACTCCGGCATTAGCCTA
GA~~T~~CTGGGACTCCGGCATTAGCCTA
GA~~T~~CTGGG TCCGGCATT~~C~~GCCTA
GA~~T~~CTGGGACTCCGGCATTAGCCTA
GA~~T~~CTGGGACTCCGGCATTA CCTA

↓ ↓ ↓ ↓ ↓

GA~~T~~GGGACTCCGGCATT~~A~~GCCTA

Genome coverage



- Reads randomly sampled
- Want to over sample genome to improve accuracy
- Genome coverage
- Read depth at each genome position
- ~ 10x coverage with Sanger sequencing, today 30-60x

The full tale

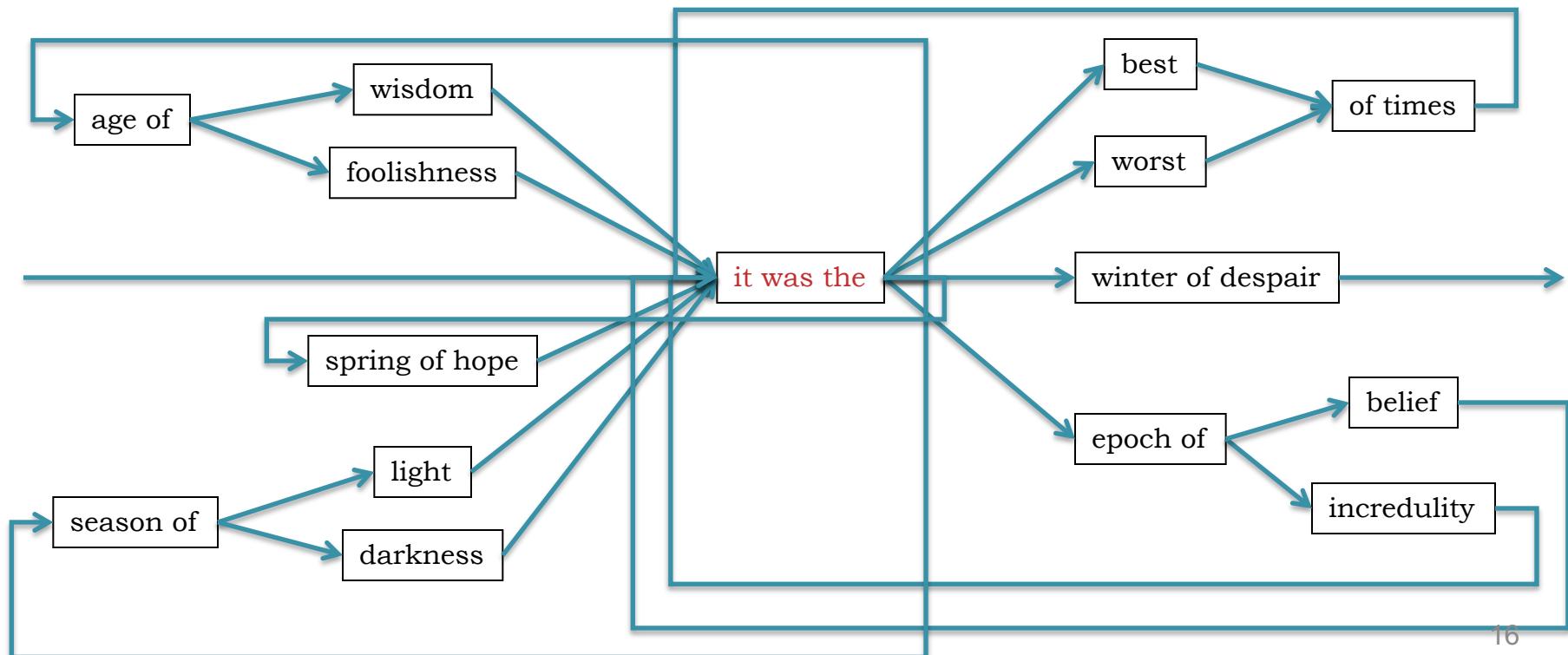
... it was the best of times it was the worst of times ...

... it was the age of wisdom it was the age of foolishness ...

... it was the epoch of belief it was the epoch of incredulity ...

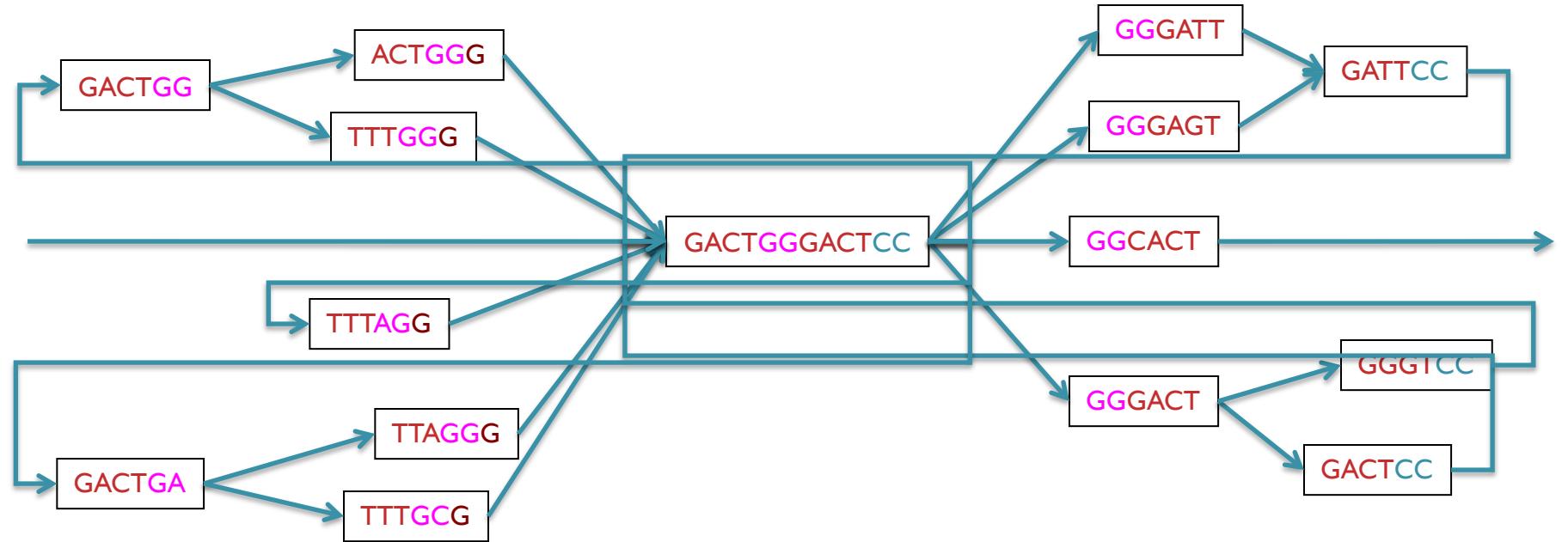
... it was the season of light it was the season of darkness ...

... it was the spring of hope it was the winter of despair ...



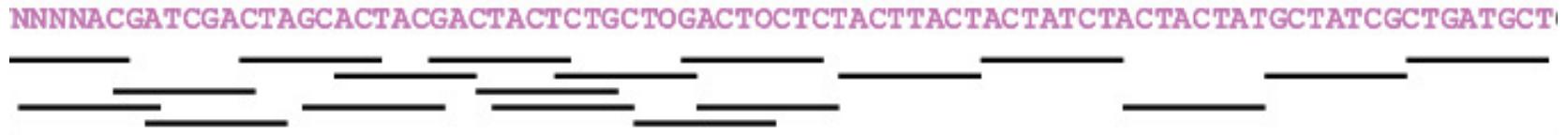
Build an assembly graph

TAGACT**GGGACTCCAG**
AAGACT**GGGACTCCGG**
GGACT**GGGAGTCCCTG**
CGTT**GGGGGT CCTTA**



Considerations for assembly projects

Coverage



- How many times has the genome been sequenced?
Number of reads that cover sequence in your assembly
Too little? Few reads lead to sequencing errors

Read 1:	CGGATTACGTGGACCATG (read length of 18)
Read 2:	ATTACGTGGACCATGAATTGCTGACA
Read 3:	ACCATGAATTGCTGACATTGTCA
Read 4:	TGAATTGCTGACATTGTCA
Depth:	1112222222233334433333333332222221

- Too much? Aim for oversampling ~30-60x.
- Raw read depth vs Mapped read depth
-driven by efficiency of alignment process

Considerations for assembly projects

Coverage



Much is driven by funding and application

- SNPs and genome rearrangements
- structural variants
- particular coding regions
- ‘complete’ genome

Depth vs Coverage

Coverage - how deeply genome is sequenced, how many times each nucleotide represented

Depth - number of reads aligned to a particular location

Calculating coverage

- Lander-Waterman Model (1988)
 - Assumes reads randomly positioned in genome & same probability of covering each region of a genome

$$\text{Coverage } c = L * N/G$$

- L = Read length
- G = Genome size
- N = Number of reads

Can be modeled by a Poisson distribution

Many coverage calculators available online

Correction by coverage.

Figure S16, evaluates the ability of high-coverage to correct for sequencing errors. With sufficient coverage, even high error rates can be compensated for.

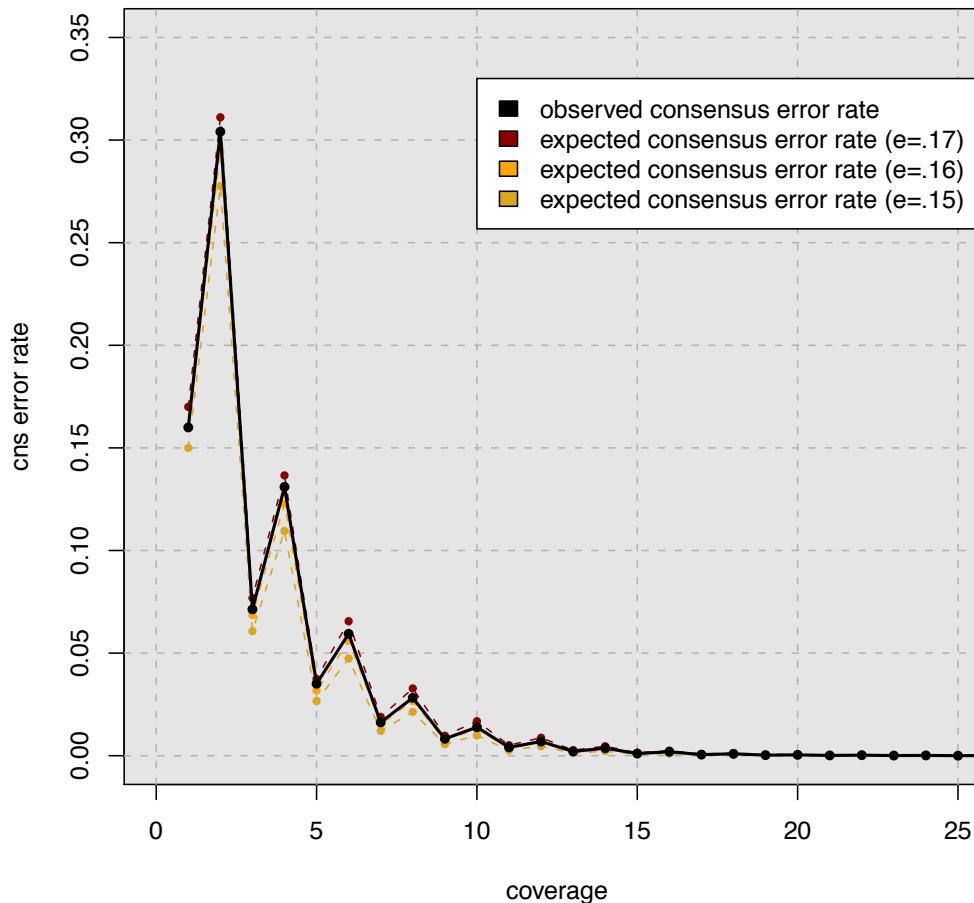
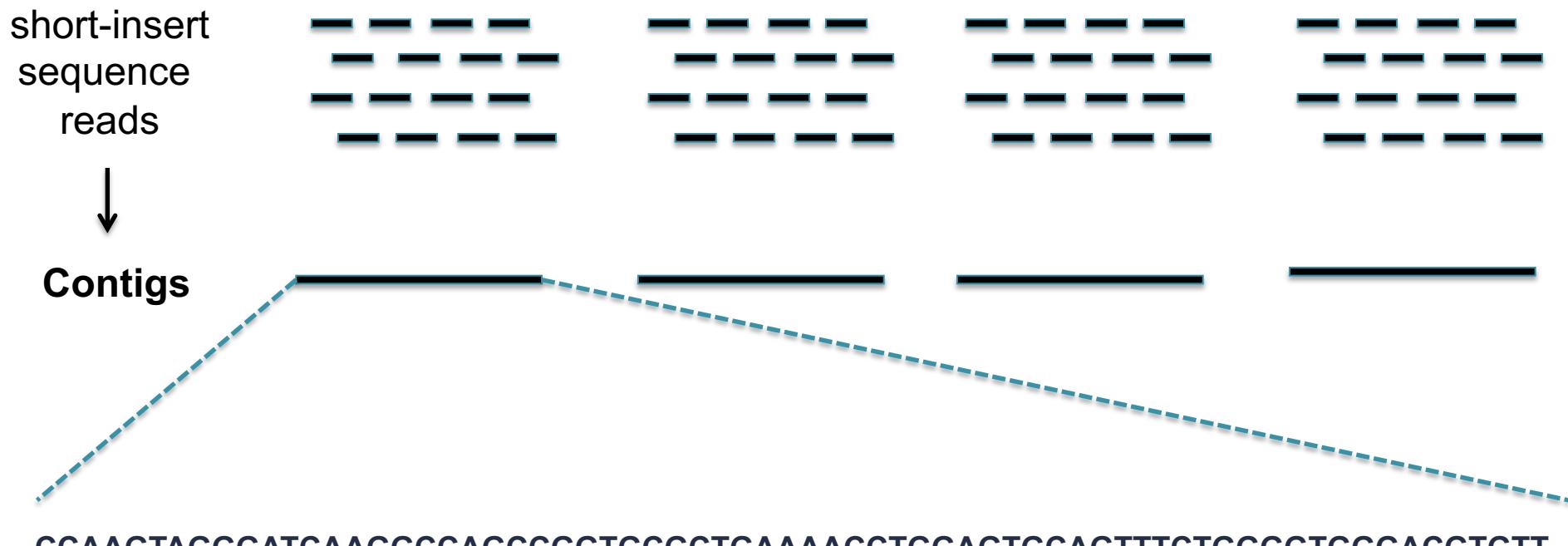


Figure S16. Coverage can overcome most random errors. 1,000 bp reads for *E. coli* K12 were simulated with random errors and the resulting consensus accuracy was measured. Even with high errors, coverage over 10X is sufficient to generate an accurate consensus. The periodic fluctuation in consensus error rate is an artifact of the tie-breaking scheme used in the consensus simulation (even numbers of reads can have ties and odd cannot).

Hybrid error correction and de novo assembly of single-molecule sequencing reads.

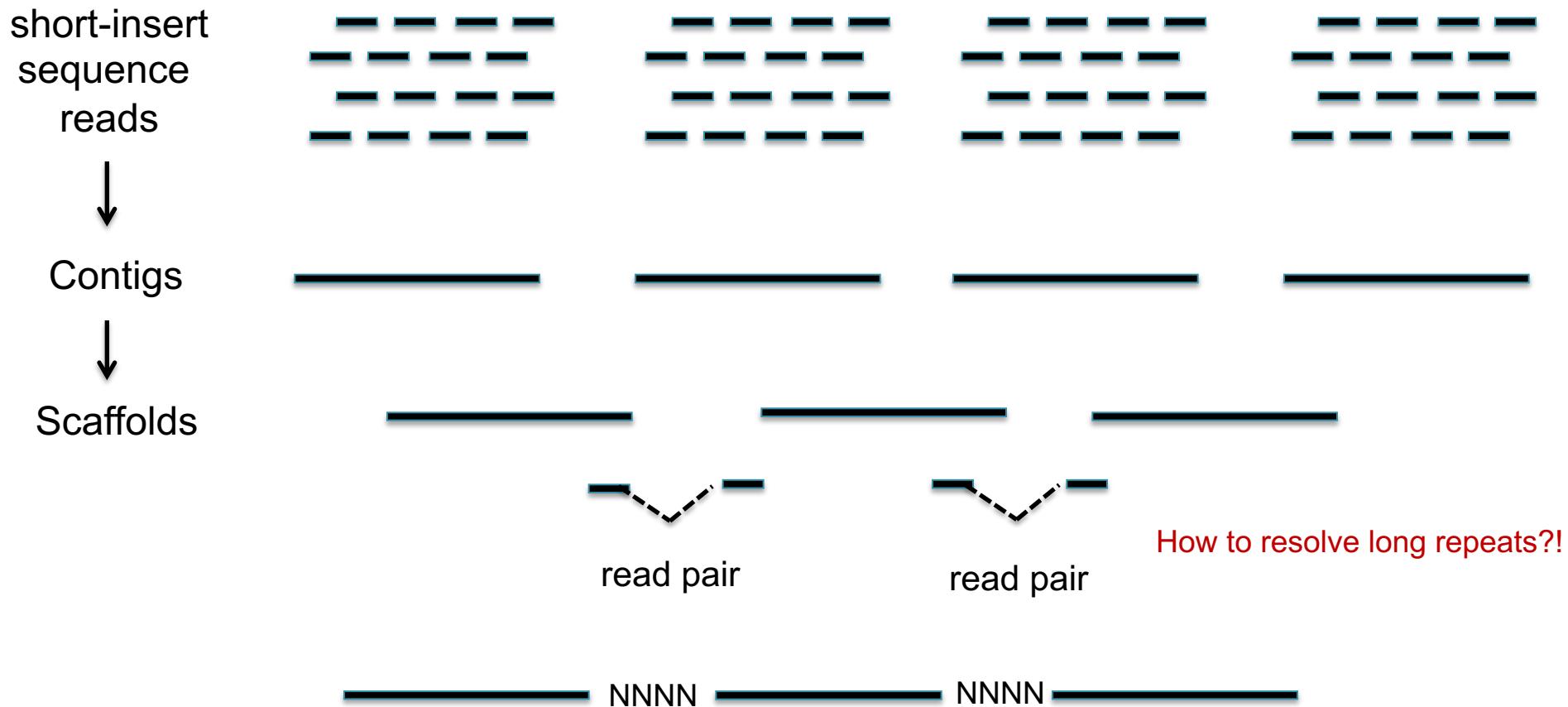
Koren et al (2012) *Nature Biotechnology*. doi:10.1038/nbt.2280

Assembly construction - Hierarchical process



Contig - contiguous sequence of nucleotides

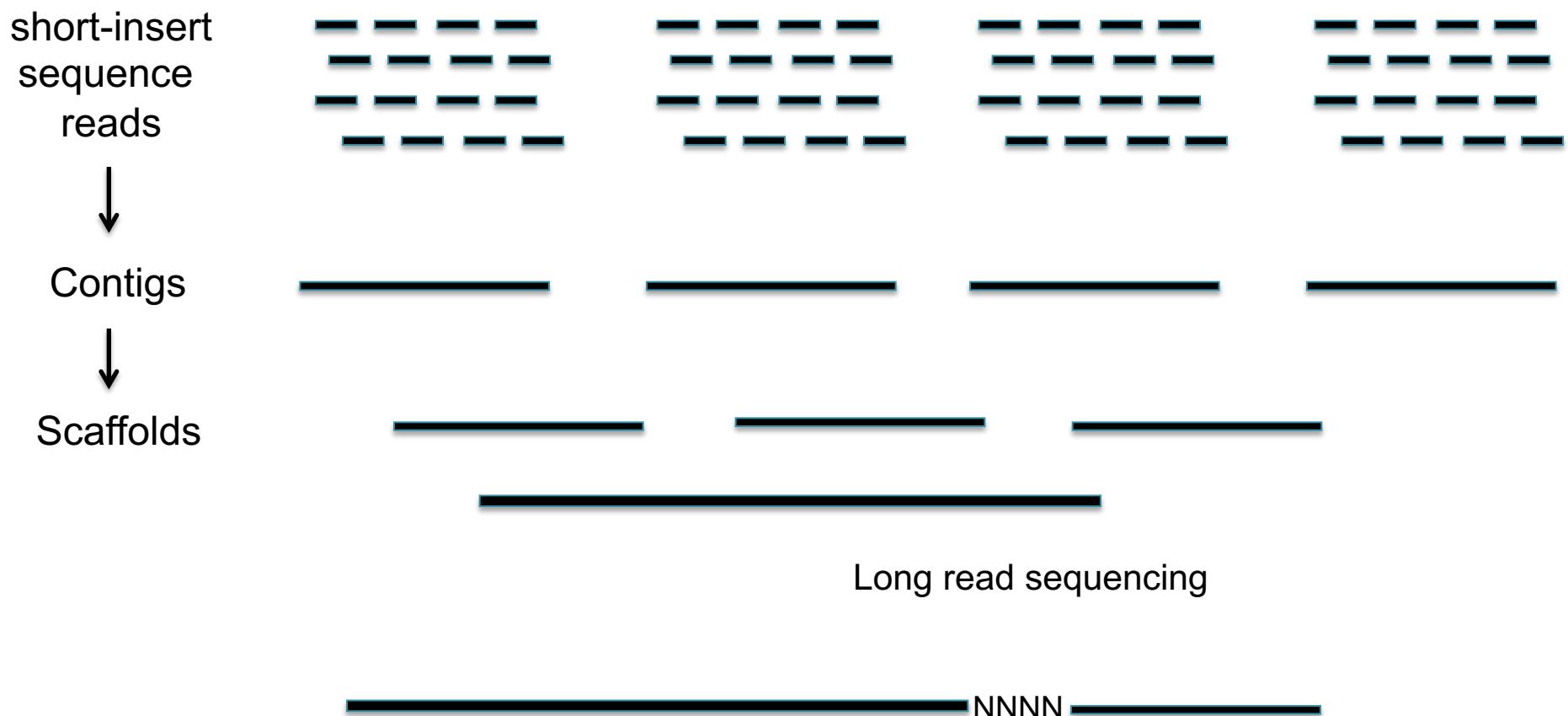
Assembly construction



Scaffold - sequence of contigs, separated by gaps - Ns are predicted gap size

```
>GTAGTATTCTAGAAAATGTTAACATAGATAGTTGTTANATCTGTTAGTGTCAGATGCTACTGAATAGTTGGAANNNNNNNNNNNNNNNNNN  
NNNNNNNTGTGAGGTTTAGCTCATGAAAGTTATGATTATTGCACCCCTACTCACAAACGAATCCCTATTCTTATCTTTNNNNNNNNNNNNNNN  
CATGTCACGGTTTATTTATTTGTGGCTGCAGAAGTCCTTGTGCTGTTAATTTGGAGTTCTCCTGTCGTATATAAGCTTCTTCTTCAGTT  
TAAATTATTTAACCTTACTATCTTCTAACATAAAATTGGAATTATCAACGAAAACATAGGNNNNNNNNNNNN  
GTCCTTATACGAAAGCTATATAG  
TGTAGGCTTCTTTNNNNNNNNNGGTGATGTTGTTAATGGTGCCTTCTGTAATCTTACTAAATCAGTTGCTGTTACTGTATAGTTG
```

Assembly construction



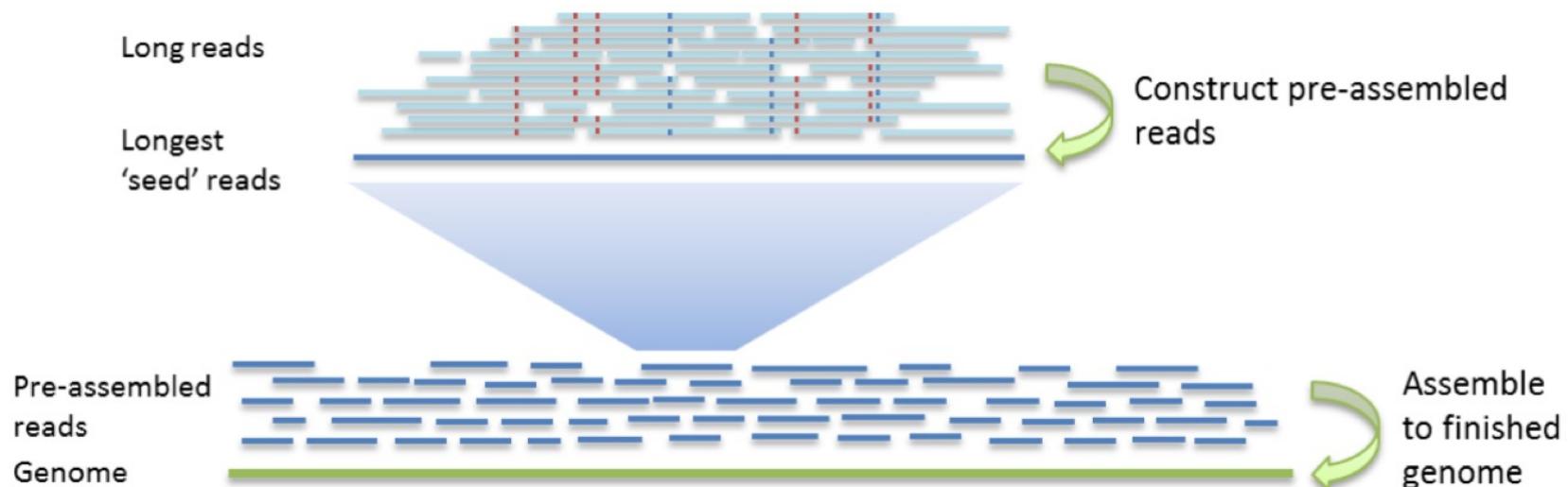
Scaffold - sequence of contigs, separated by gaps - Ns are predicted gap size

```
>GTAGTATTCTAGAAAATGTTAACATAGATAGTTGTTANATCTGTTAGTGTCAGATGCTACTGAATAGTTGGAANNNNNNNNNNNNNNNNNN  
NNNNNNNTGTGAGGTTTAGCTCATGAAAGTTATGATTATTGCACCCCTACTCACAAACGAATCCCTATTCTTATCTTTNNNNNNNNNNNNNNN  
CATGTCACGGTTATTTATTTTGTTGCGCTGCAGAAGTCCTTGTGCTGTTAATTTGGAGTTCTCCTGTCGTATATTAAAGCTTCTTCTTCAGTT  
TAAATTATTTGAACCTTACTATCTTCTAACATAAAATTGTTGAATTATCAACGAAAACATAGGNNNNNNNNNNN  
GTCCTTATACGAAAGCTATATAGTGTAGGCTTTCTTTNNNNNNNNNGGTGATGTTGTTAATGGTGCCTTCTGGTAATCTTACTAAATCAGTTGCTT  
ACTGTATAGTTG
```

Assembly construction - long reads

Still hierarchical but proceeds in two rounds

1. Seed read selection – longest reads in dataset
2. Shorter reads aligned to seed reads for consensus



Repetitive regions

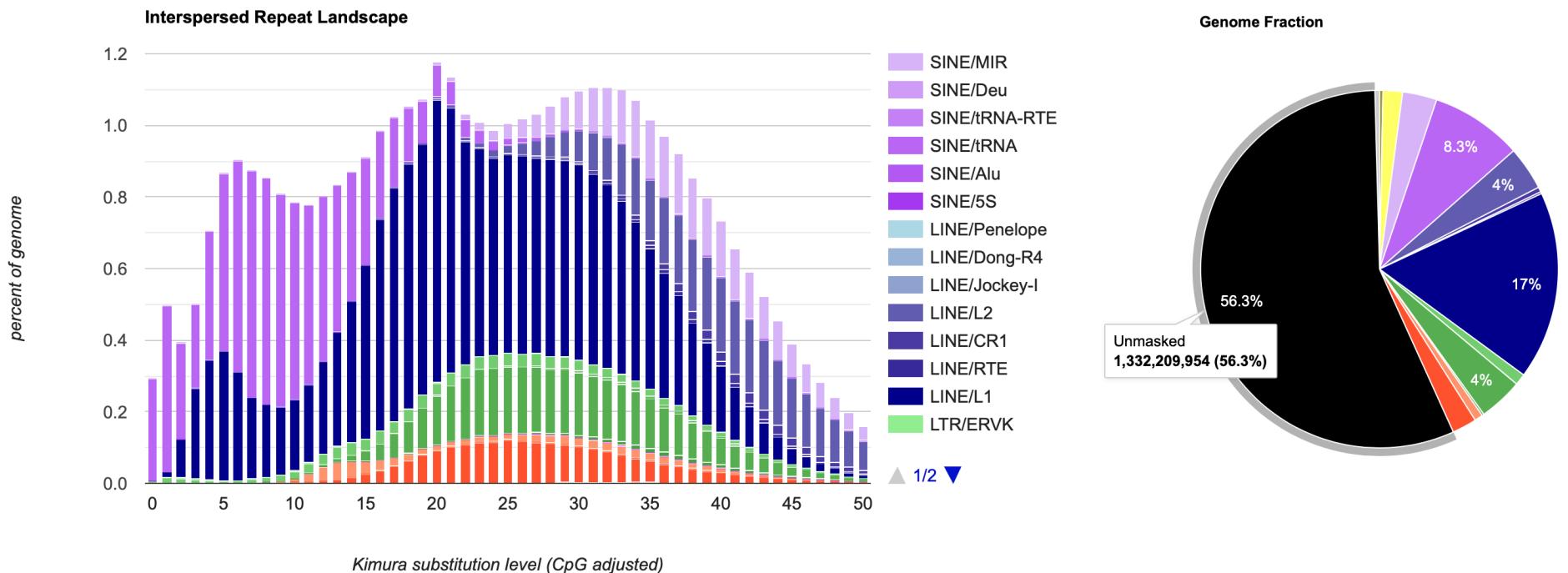
- Over 50% of mammalian genomes are repetitive
 - Large plant genomes tend to be even worse
 - *Arabidopsis* - 10% composition
 - SINEs - Short Interspersed Nuclear Elements
 - LINES - Long Interspersed Nuclear Elements
 - LTR - Long Terminal Repeats, retrotransposons
 - Segmental duplications
 - Low-complexity - Microsatellites or homopolymers

Sequencing challenges

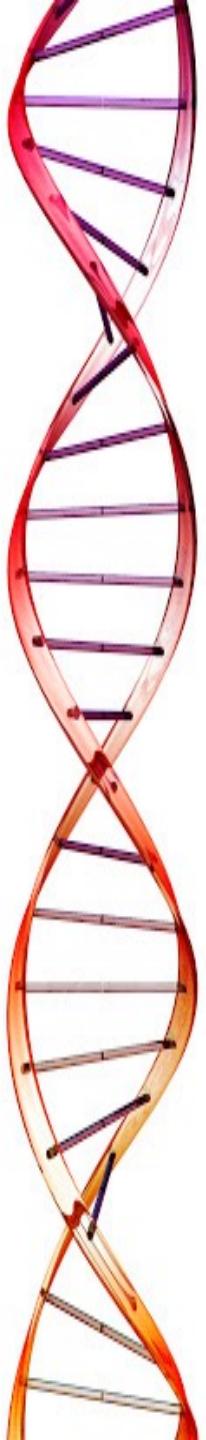
- Resolving repeats
 - Repeats longer than reads
 - Merge reads up to repeat boundaries
 - Tandem repeats often not resolved
- RepeatMasker - <http://repeatmasker.org>
screens DNA sequences for interspersed repeats and
low complexity DNA sequences

- Error correction
 - incorporation of long reads
 - correction algorithms

F. catus RepeatMasker output



NNNNNNNNNNNcaaacaatatacaaagacAAAAATTGCCACAGCAAAGACAAAGAGATAAATAAAAGGCACAAAATT



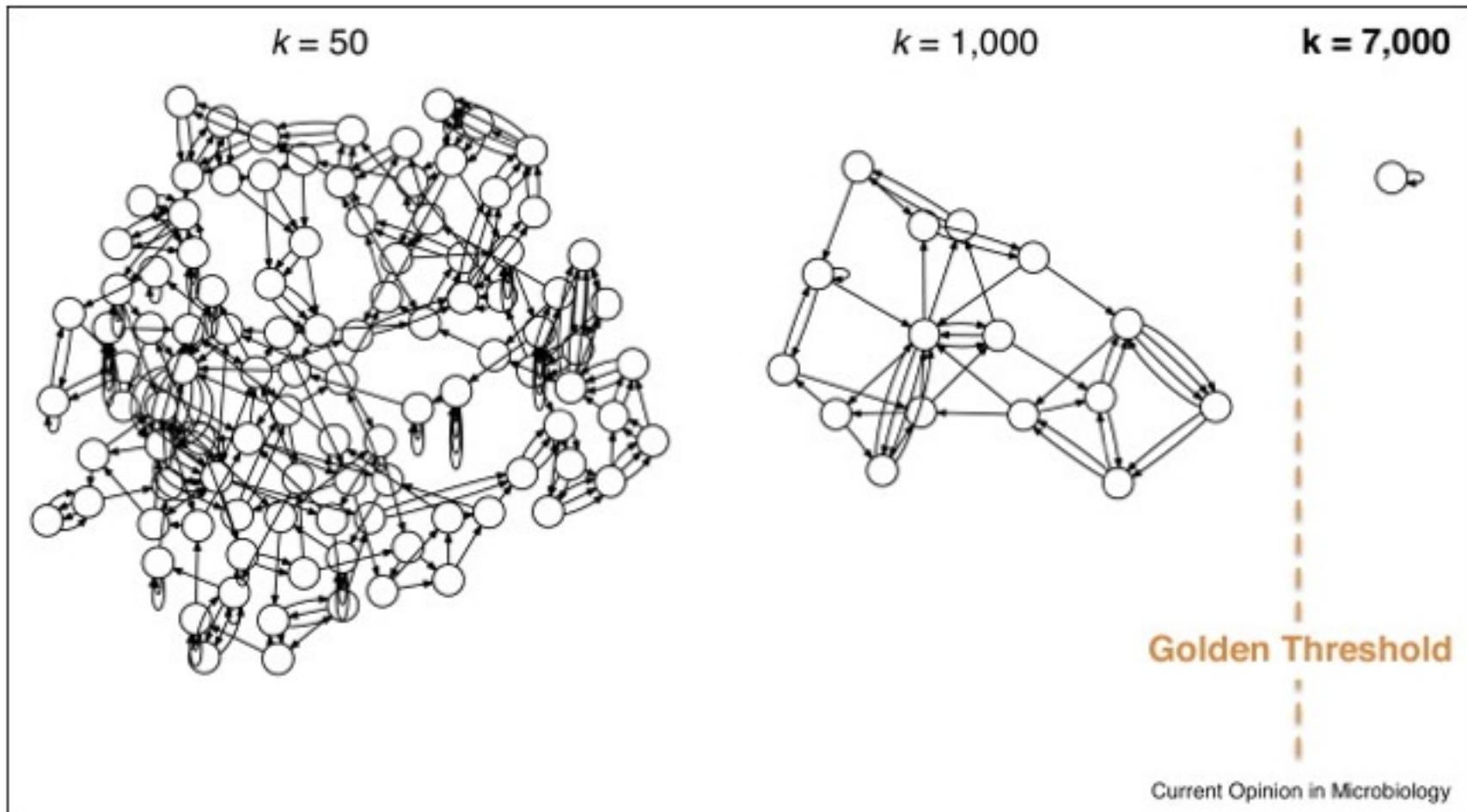
Lecture outline

- I. General background & theory of genome assembly
2. Comparison of sequencing technologies
3. Assembly quality and annotation
4. (Not explicitly numbered, implied by the sequence)
5. Assembly workshop with Python programming

Sequencing Technologies

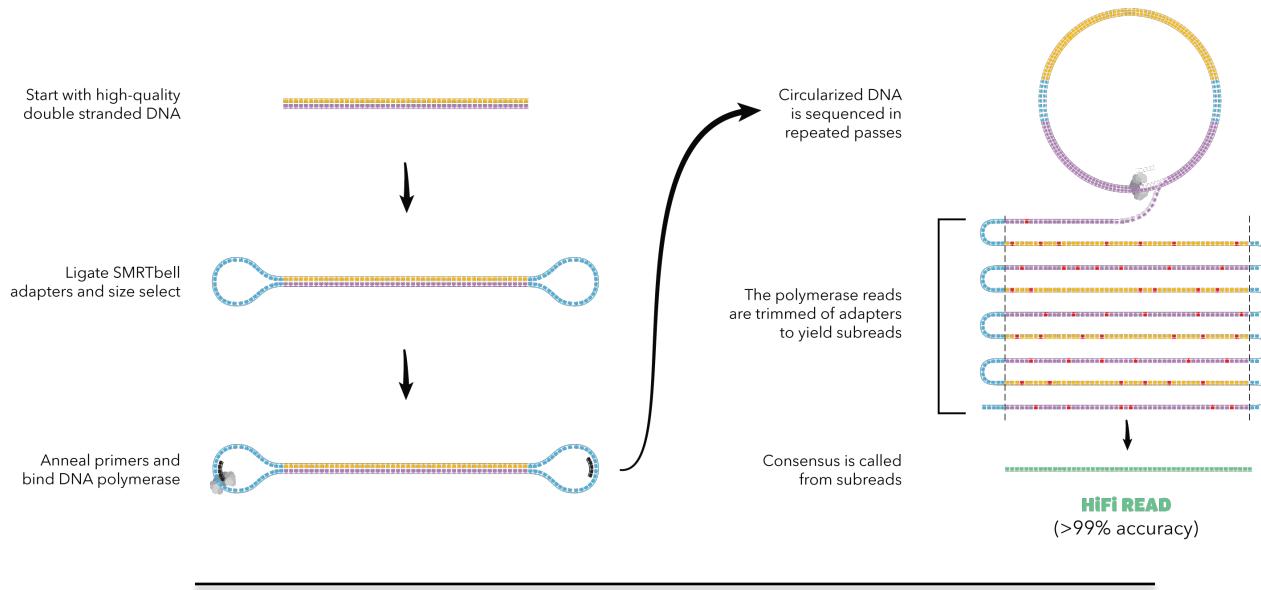
- Illumina
- PacBio
- Nanopore
- Hi-C

Read lengths



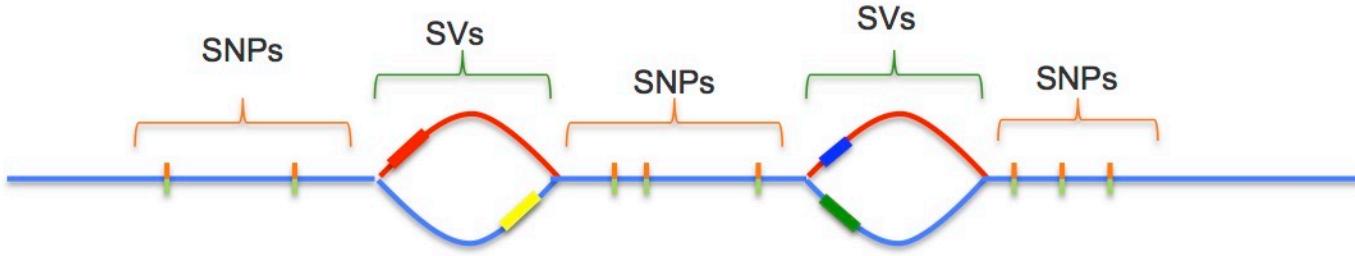
Koren & Phillippy 2015 Current Opinion in Microbiology

PacBio Hi Fi reads



- Long reads - 20kb
- 99.9% read quality (Q30) - better than some short-read contigs
- Distinguish repeats rather than spanning full repeats
- Limited length
- Lower throughput than Nanopore

Phased genome assemblies



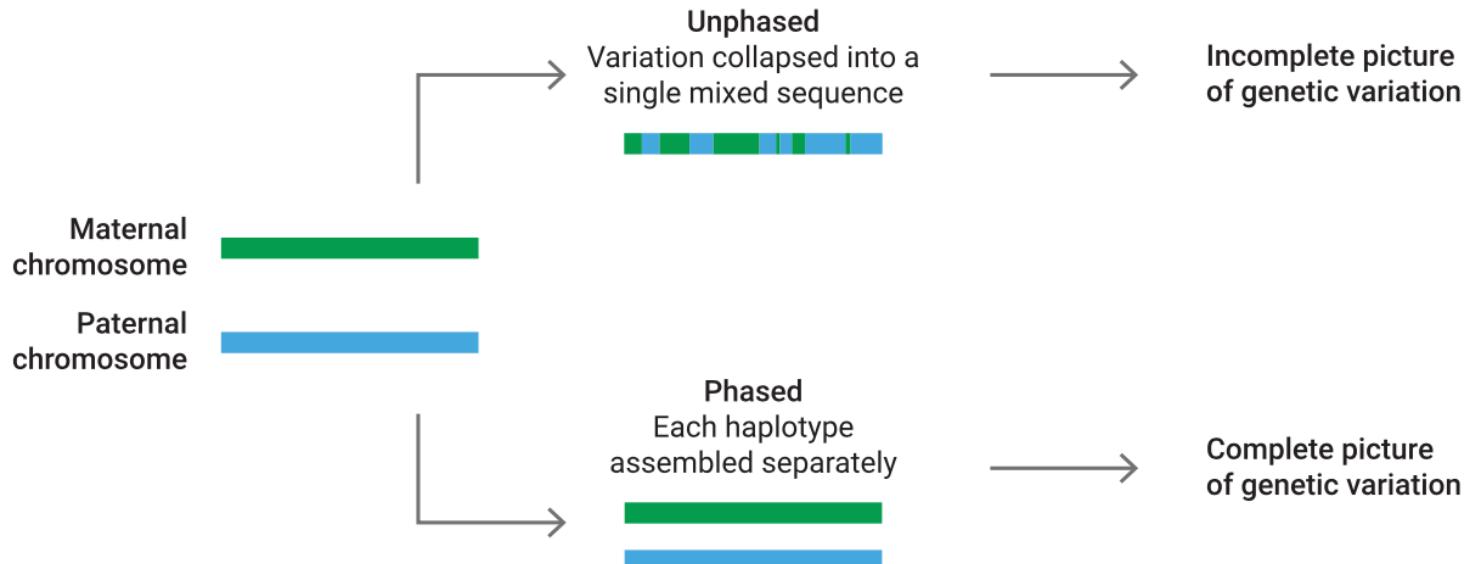
“Pseudo-haplotypes” - mixture of both types

How do we get phased genomes?

Higher heterozygosity - easier to phase.

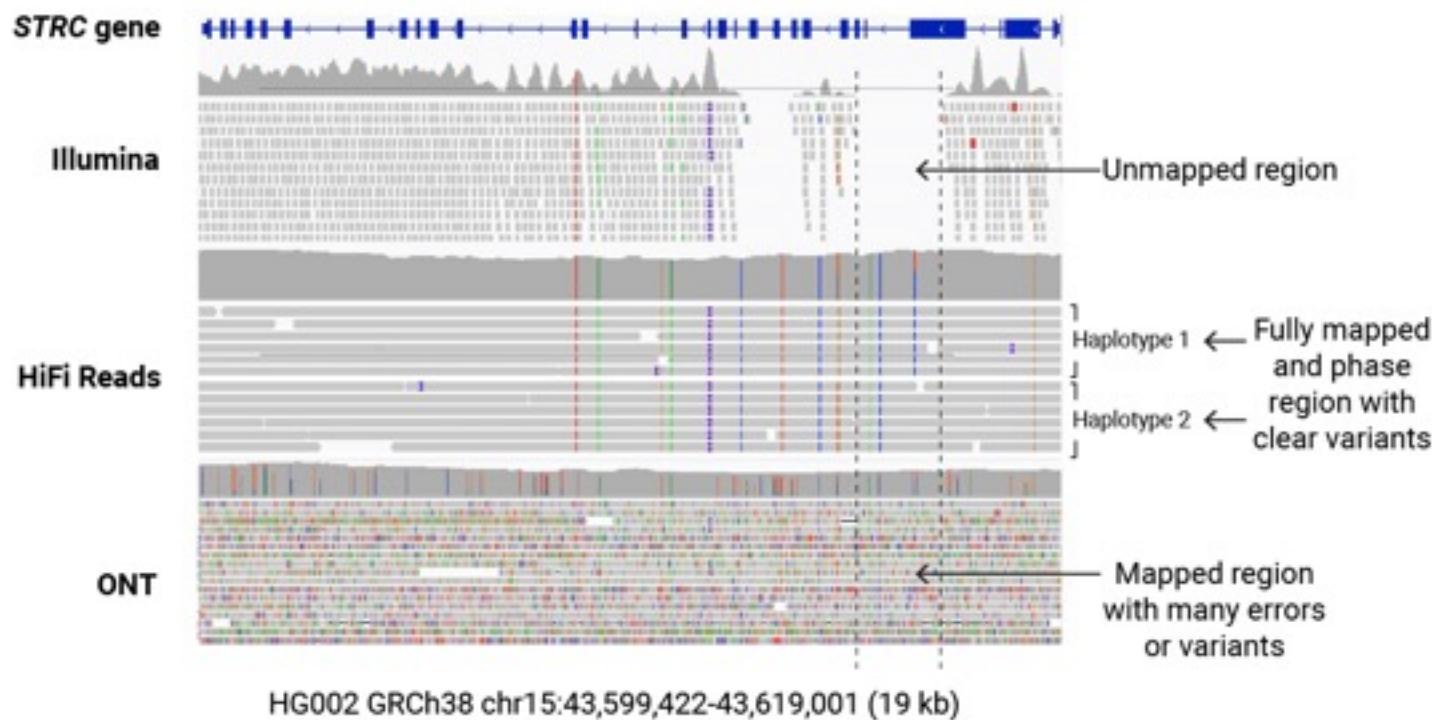
Chin et al. *Nat Meth* 2016

Phased genome assemblies



Phasing involves separating maternally and paternally inherited copies of each chromosome into haplotypes to get a complete picture of genetic variation.

PacBio Hi Fi reads



HiFi reads provide the accuracy needed to call single nucleotide variants, while improving mappability and enabling phasing with no systematic bias. *STRC* gene alignments from [Genome in a Bottle \(GIAB\)](#), [HG002_NA24385_son](#). ([IGV settings](#))

Research Briefing | Published: 18 October 2024

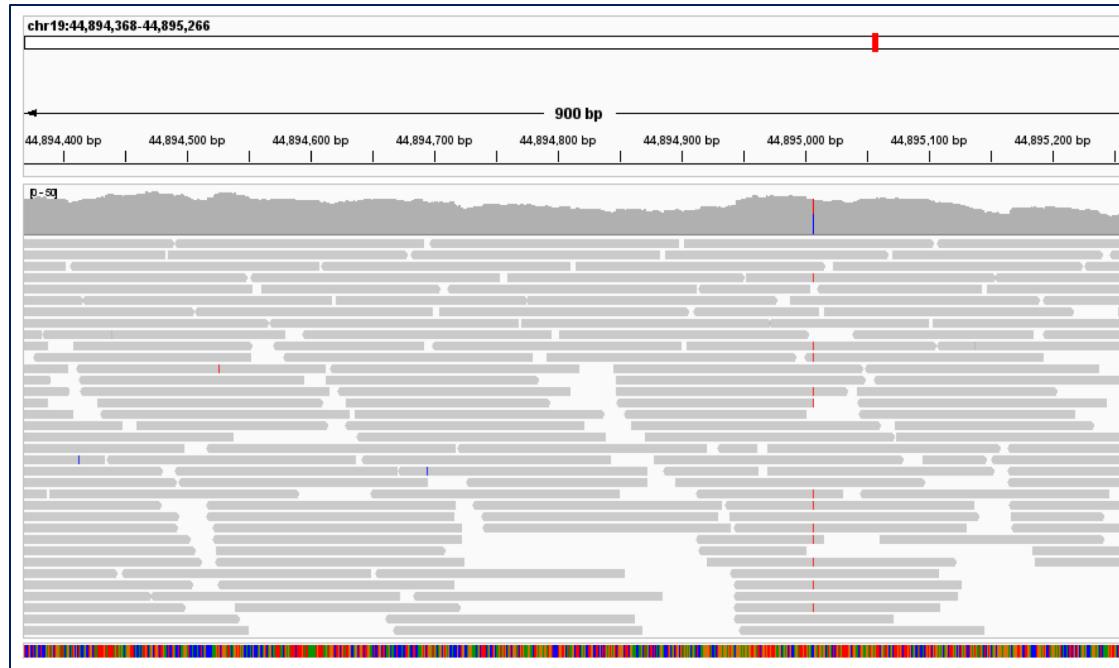
Unravelling the complex origin and breeding history of modern roses

[Nature Plants](#) (2024) | [Cite this article](#)

52 Accesses | [Metrics](#)

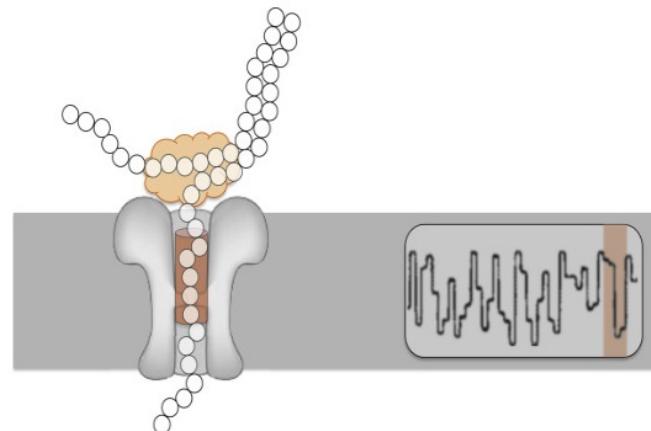
PacBio short-read sequencing

- Sequencing by binding (SBB)
- Detects light signals during binding
- Error rates 1/10,000 bases
- Cancer detection, liquid biopsies
- Tumor DNA - short fragments, variant detection



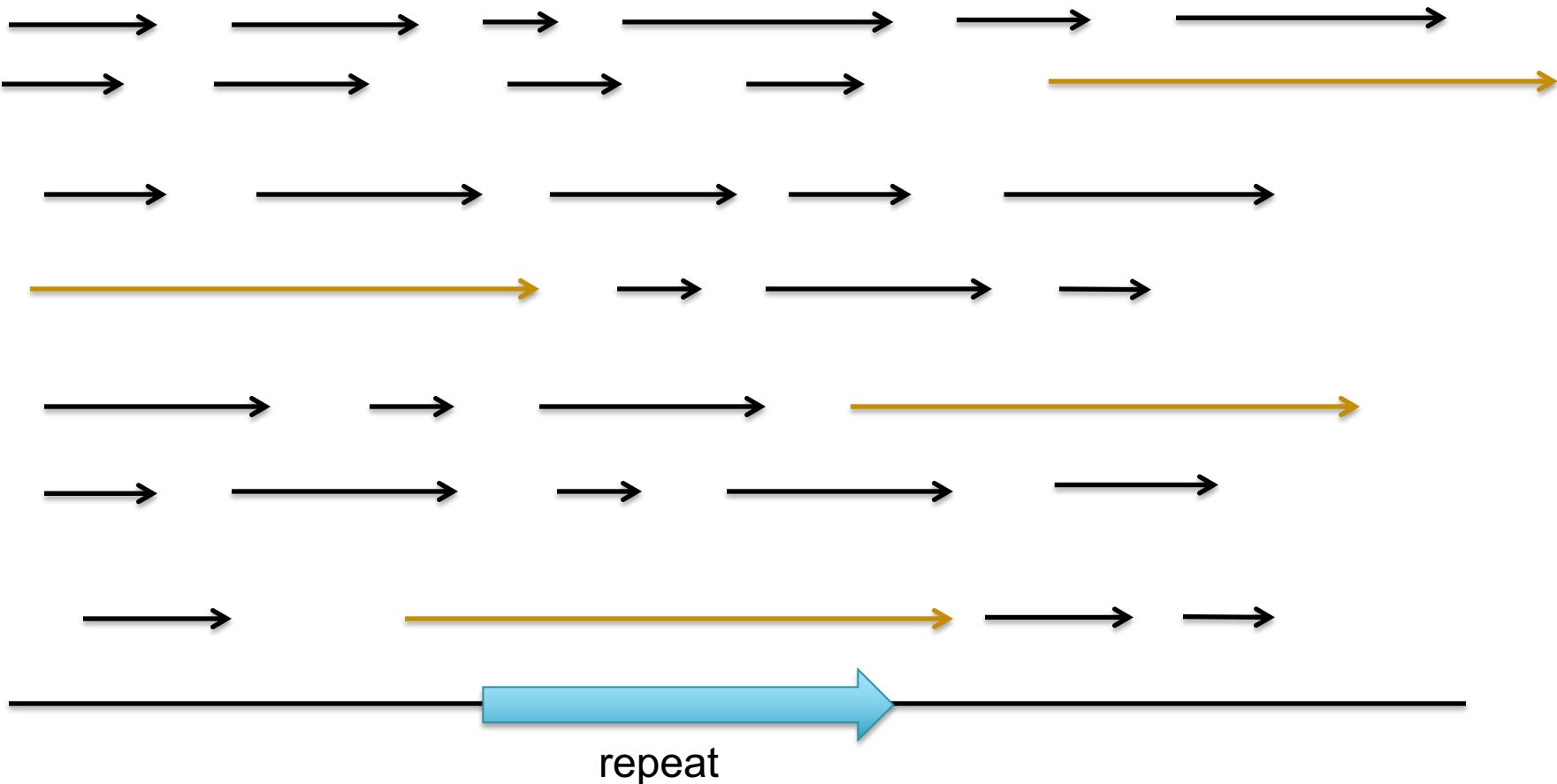
Nanopore sequencing

- “Ultra-long” reads
- Up to 1Mb read length
- 95% read quality
- Engineered
- Long lengths
- Able to span repeat regions
- Limited quality
- Took a while for technology to catch up



Nanopore sequencing

- Using long reads and coverage to detect repeats



Nanopore sequencing

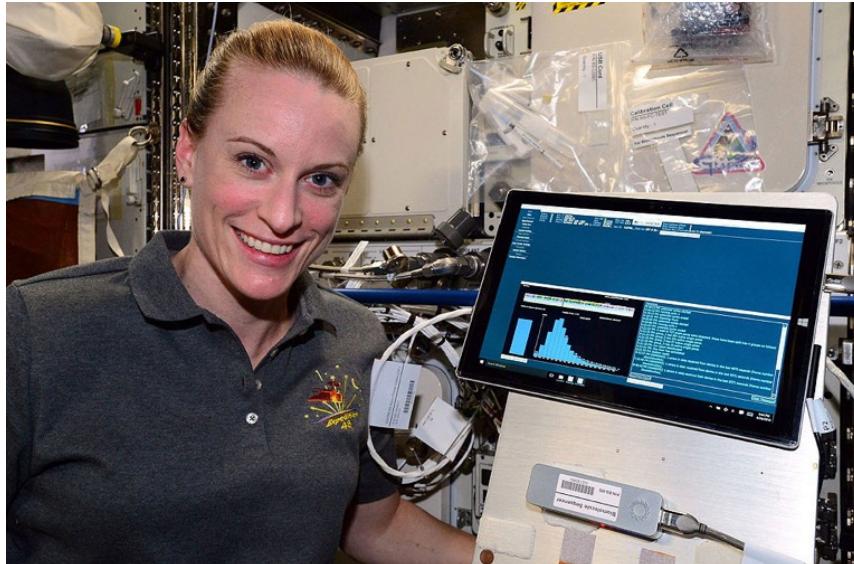
COMMENT

Open Access

Mobile real-time surveillance of Zika virus in Brazil



Nuno Rodrigues Faria¹, Ester C. Sabino², Marcio R. T. Nunes^{3,4}, Luiz Carlos Junior Alcantara⁵, Nicholas J. Loman^{6*} and Oliver G. Pybus¹



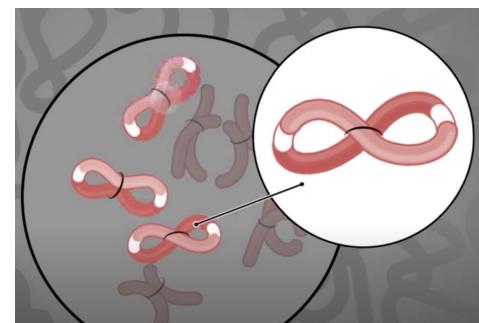
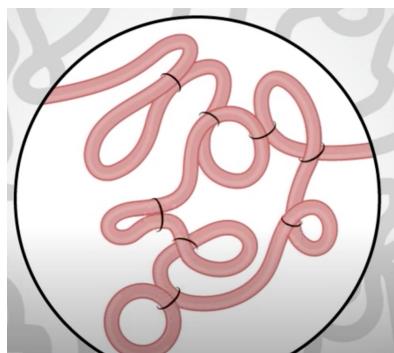
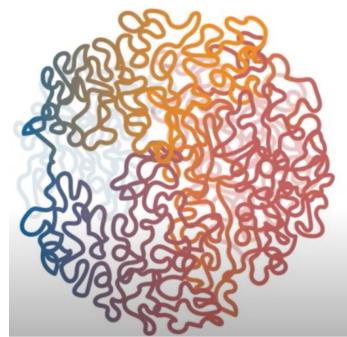
Kate Rubins, ISS



Castro-Wallace, Nature Scientific Reports, 2017

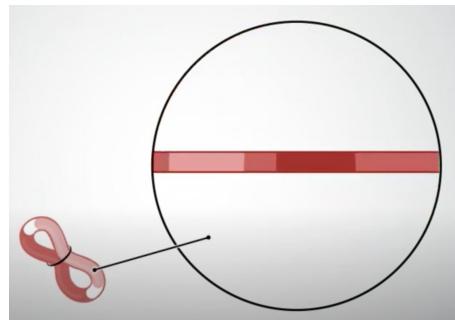
Hi-C

High throughput genomic technique to capture chromatin conformation



- Within packaged chromatin, sequences closer on the same chromosome, closer in physical space
- DNA is crosslinked to capture interactions throughout the genome
- Crosslinked DNA is fragmented and ligated

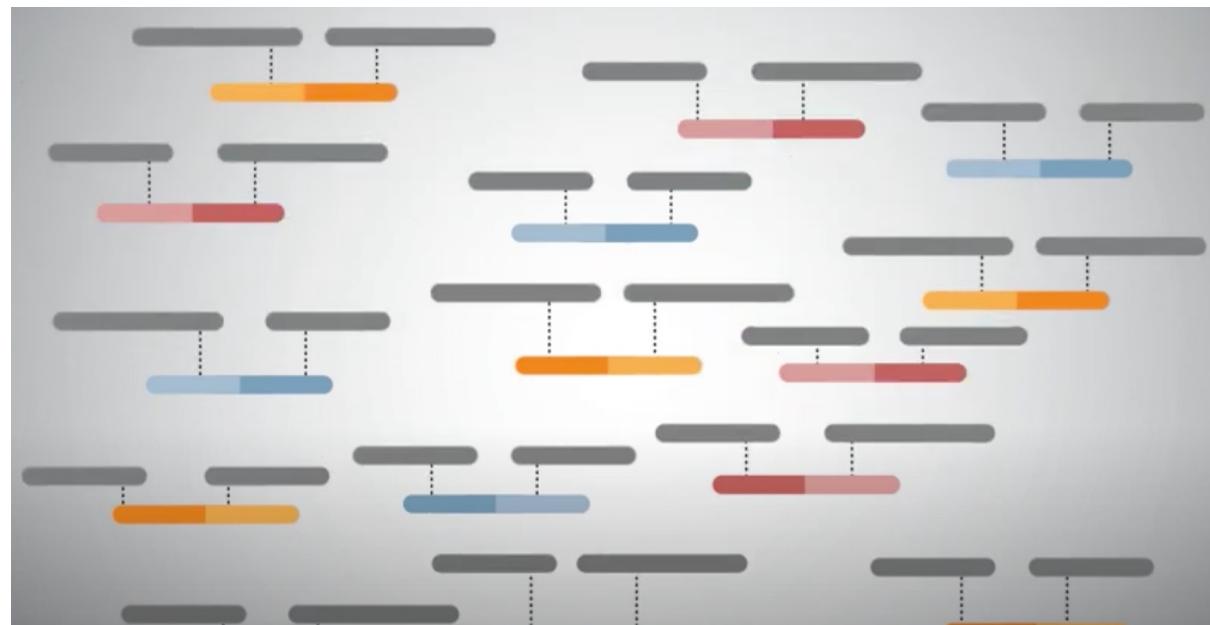
Hi-C



More frequently sequenced ligated, closer are in genome



Junctions between adjacent sequences



Ligations are paired-end sequenced and reads mapped against draft genome

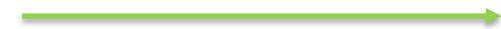
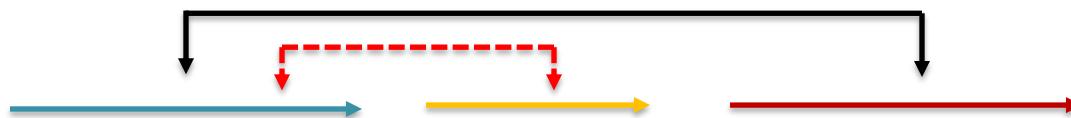
Hi-C

Draft assembly

contigs



Scaffolding



Contact map with frequency of contact: black - more contact, dotted red - less

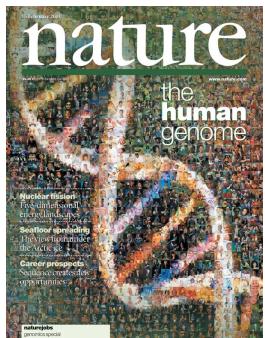
History of Genome Assembly

1977. Sanger et al. 1st Complete Organism bacteriophage 5375 bp

1995. Fleischmann et al. 1st Free Living bacteria; *Haemophilus influenzae*; TIGR Assembler. 1.8Mb

1998. *C.elegans* SC 1st Multicellular Organism BAC-by-BAC Phrap. 97Mbp

2000. *Drosophila* genome; Myers et al. 1st Large WGS Assembly Celera Assembler. 116 Mbp



Human Genome

Public: 13-year project began 1990, Dept Energy & NIH,
\$3 billion; millions of small fragments
2003 – announced as complete



Private: Craig Venter, Celera Genomics; 1998, \$300 million
Could not be patented.

Human genome “finished” ~2003

Sequencing Technologies

Combine technologies to improve assemblies

Article

Telomere-to-telomere assembly of a complete human X chromosome

<https://doi.org/10.1038/s41586-020-2547-7>

Received: 30 July 2019

Accepted: 29 May 2020

Published online: 14 July 2020

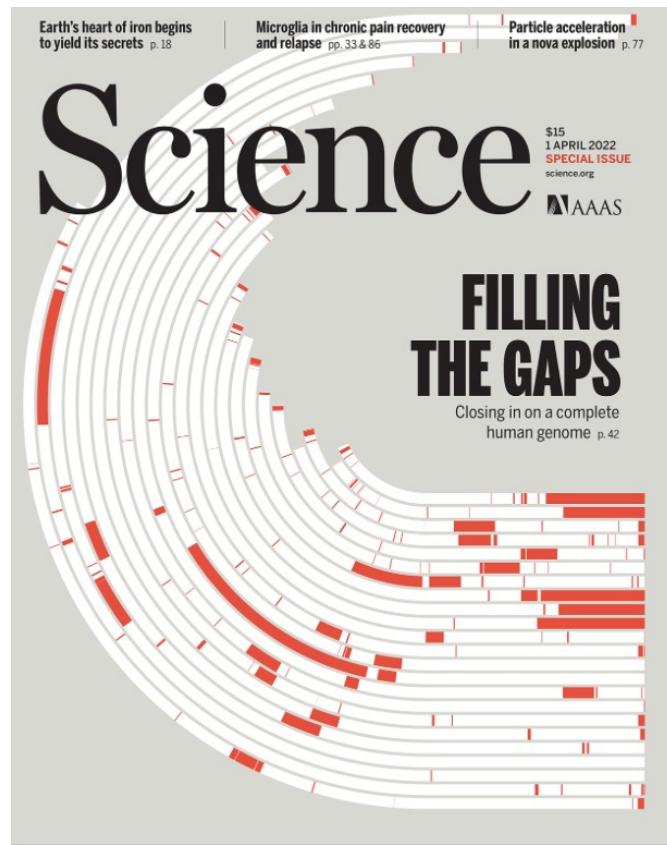
Open access

 Check for updates

Karen H. Miga^{1,24}✉, Sergey Koren^{2,24}, Arang Rhie², Mitchell R. Vollger³, Ariel Gershman⁴, Andrey Bzikadze⁵, Shelise Brooks⁶, Edmund Howe⁷, David Porubsky³, Glennis A. Logsdon³, Valerie A. Schneider⁸, Tamara Potapova⁷, Jonathan Wood⁹, William Chow⁹, Joel Armstrong¹, Jeanne Fredrickson¹⁰, Evgenia Pak¹¹, Kristof Tigyi¹, Milinn Kremitzki¹², Christopher Markovic¹², Valerie Maduro¹³, Amalia Dutra¹¹, Gerard G. Bouffard⁶, Alexander M. Chang², Nancy F. Hansen¹⁴, Amy B. Wilfert³, Françoise Thibaud-Nissen⁸, Anthony D. Schmitt¹⁵, Jon-Matthew Belton¹⁵, Siddarth Selvaraj¹⁵, Megan Y. Dennis¹⁶, Daniela C. Soto¹⁶, Ruta Sahasrabudhe¹⁷, Gulhan Kaya¹⁶, Josh Quick¹⁸, Nicholas J. Loman¹⁸, Nadine Holmes¹⁹, Matthew Loose¹⁹, Urvashi Surti²⁰, Rosa ana Risques¹⁰, Tina A. Graves Lindsay¹², Robert Fulton¹², Ira Hall¹², Benedict Paten¹, Kerstin Howe⁹, Winston Timp⁴, Alice Young⁶, James C. Mullikin⁶, Pavel A. Pevzner²¹, Jennifer L. Gerton⁷, Beth A. Sullivan²², Evan E. Eichler^{3,23} & Adam M. Phillippy²✉

Human genome completion

- 120x Nanopore
 - 70x PacBio
 - 30x PacBio HiFi
 - 50x 10X Genomics
 - 100x Illumina
 - 35x Arima Hi-C
 - BioNano optical map
 - PacBio Iso-Seq
-
- Telomere-to-Telomere
 - Centromeres resolved
 - Previous gaps filled



Article

The complete sequence of a human Y chromosome

Nature, September 2023

Genome Assembly Projects



Bhattacharya et al. 2018.
Genome Research



Axolotl salamander - 32 billion

Wheat - hexaploid, 15.3 billion bases



Anopheles mosquitos - 3 chromosomes



Article

<https://doi.org/10.1038/s42256-024-00872-0>

DNA language model GROVER learns sequence context in the human genome

Received: 31 August 2023

Melissa Sanabria¹, Jonas Hirsch¹, Pierre M. Joubert^{1,2,3} & Anna R. Poetsch^{1,4}

Accepted: 26 June 2024

So many genomes!

<https://doi.org/10.1038/s41576-024-00718-w>

Review article

Check for updates

Genome assembly in the telomere-to-telomere era

News Feature | Published: 12 January 2023

Method of the year: long-read sequencing

Vivien Marx

[Nature Methods](#) 20, 6–11 (2023) | [Cite this article](#)

Mini-Review

Toward telomere-to-telomere cat genomes for precision medicine and conservation biology

William J. Murphy^{1,2,3} and Andrew J. Harris^{1,3}¹Department of Veterinary Integrative Biosciences, ²Department of Biology, ³Interdisciplinary Program in Genetics and Genomics, Texas A&M University, College Station, Texas 77843-4458, USAArticle | [Open access](#) | Published: 29 May 2024

The complete sequence and comparative analysis of ape sex chromosomes

Reference Genome

Reference:

Multiple sequencing technologies

Complete representation

Minimal sequencing gaps

Higher cost

Conventional:

Single sequencing technology

Specific goal: genetic variants

Sequencing gaps

Lower cost

Honey bee genome



Apis mellifera

Article | Published: 26 October 2006

Insights into social insects from the genome of the honeybee *Apis mellifera*

[The Honeybee Genome Sequencing Consortium](#)

[Nature](#) **443**, 931–949 (2006) | [Cite this article](#)

Research article | [Open access](#) | Published: 08 April 2019

A hybrid de novo genome assembly of the honeybee, *Apis mellifera*, with chromosome-length scaffolds

[Andreas Wallberg](#), [Ignas Bunikis](#), [Olga Vinnere Pettersson](#), [Mai-Britt Mosbech](#), [Anna K. Childers](#), [Jay D. Evans](#), [Alexander S. Mikheyev](#), [Hugh M. Robertson](#), [Gene E. Robinson](#) & [Matthew T. Webster](#)

[BMC Genomics](#) **20**, Article number: 275 (2019) | [Cite this article](#)

Assembly Projects



The Vertebrate Genomes Project

A Collection of Research Articles from Phase I of the Vertebrate Genomes Project

nature portfolio



IRIDIAN GENOMES

nature > special

Special | 28 April 2021

Vertebrate Genomes Project



Pan genomes

> [Cell.](#) 2020 Jul 9;182(1):162-176.e13. doi: 10.1016/j.cell.2020.05.023. Epub 2020 Jun 17.

Pan-Genome of Wild and Cultivated Soybeans

Yucheng Liu ¹, Hui long Du ², Pengcheng Li ³, Yanting Shen ⁴, Hua Peng ², Shulin Liu ⁵,
Guo-An Zhou ⁵, Haikuan Zhang ³, Zhi Liu ¹, Miao Shi ³, Xuehui Huang ⁶, Yan Li ⁷, Min Zhang ⁵,
Zheng Wang ⁵, Baoge Zhu ⁵, Bin Han ⁷, Chengzhi Liang ⁸, Zhixi Tian ⁹

Affiliations + expand

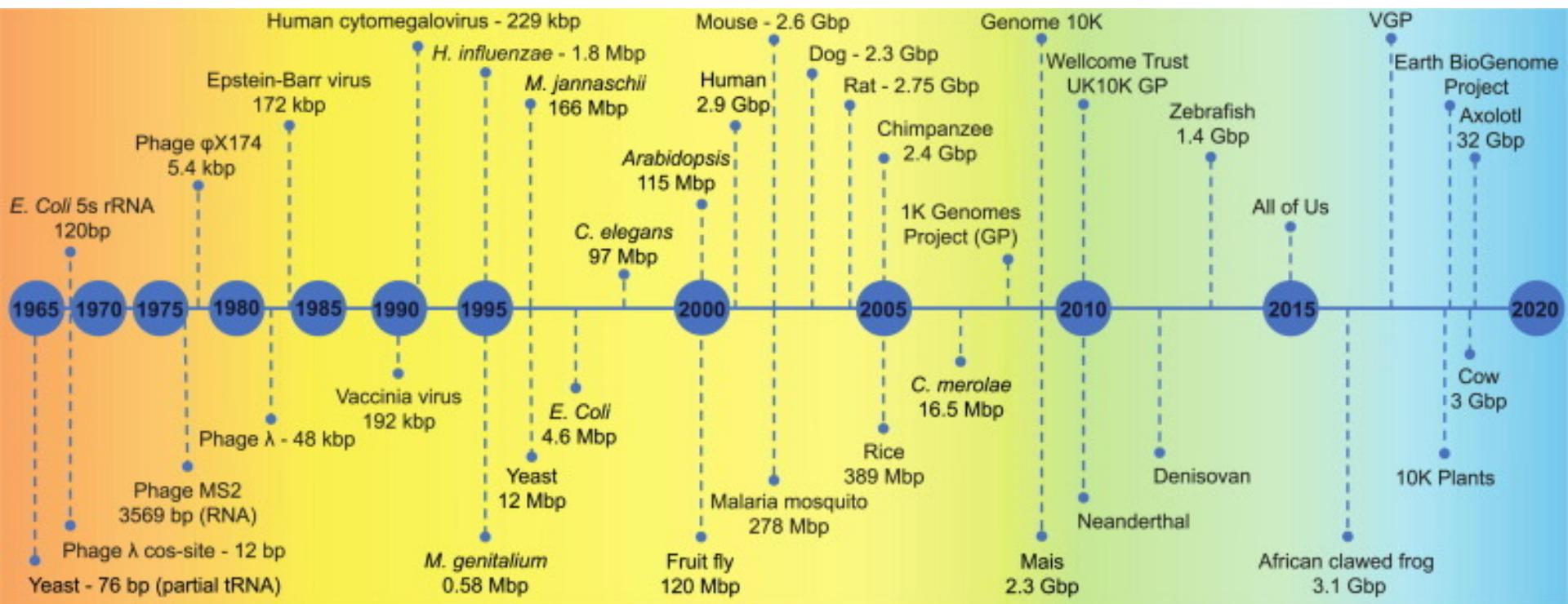
PMID: 32553274 DOI: [10.1016/j.cell.2020.05.023](#)

26 denovo assemblies

3 previously reported assemblies

Searching for variants that were undetectable in single reference genome

History of Genome Assembly



Giani et al. 2020

History of Human Genome Assembly

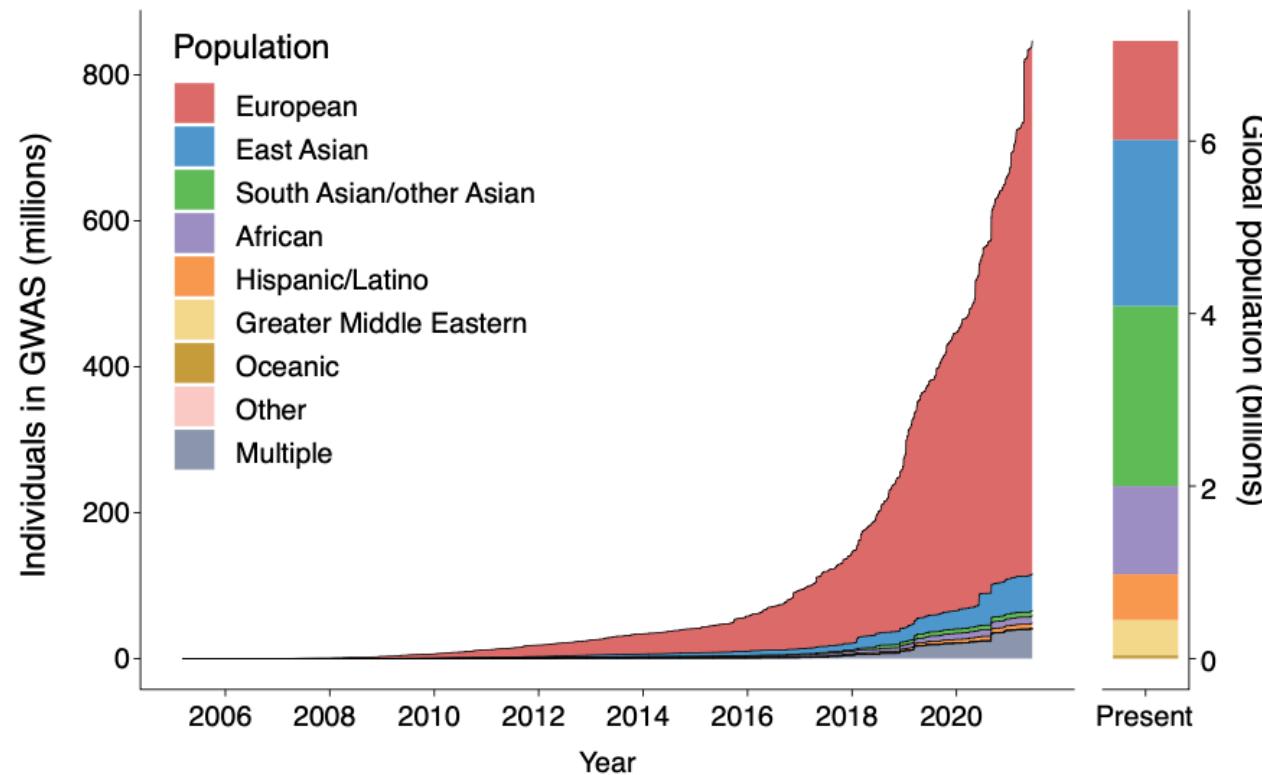


Fig. 1 | The proportion of samples from individuals cumulatively reported by the GWAS Catalog¹ as of 8 July 2021.

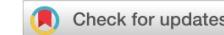
Fatumo, Nature Medicine, 2021

History of Human Genome Assembly

nature
medicine

SERIES | PERSPECTIVE

<https://doi.org/10.1038/s41591-021-01672-4>



A roadmap to increase diversity in genomic studies

Segun Fatumo^{1,2}✉, Tinashe Chikowore^{3,4}, Ananyo Choudhury³, Muhammad Ayub^{1,5},
Alicia R. Martin^{6,7} and Karoline Kuchenbaecker^{5,8}

Article

The GenomeAsia 100K Project enables genetic discoveries across Asia

<https://doi.org/10.1038/s41586-019-1793-z>

Received: 29 January 2018

Accepted: 11 October 2019

Published online: 4 December 2019

Open access

GenomeAsia100K Consortium*

The underrepresentation of non-Europeans in human genetic studies so far has limited the diversity of individuals in genomic datasets and led to reduced medical relevance for a large proportion of the world's population. Population-specific reference genome datasets as well as genome-wide association studies in diverse populations are needed to address this issue. Here we describe the pilot phase of the GenomeAsia 100K Project. This includes a whole-genome sequencing reference dataset from 1,739 individuals of 219 population groups and 64 countries across Asia. We catalogue genetic variation, population structure, disease associations and founder effects. We also explore the use of this dataset in imputation, to facilitate genetic studies in populations across Asia and worldwide.

Million-person U.S. study of genes and health stumbles over including Native American groups

By Jocelyn Kaiser | May. 29, 2019, 1:40 PM

History of Human Genome Assembly



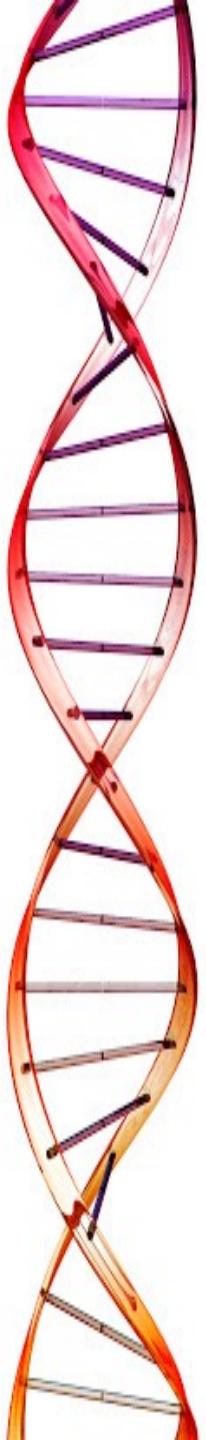
Christian Happi at Redeemer's University in Ede, Nigeria, plans to sequence human genomes.

Sequenced 426 genomes with African ancestry and discovered 3 million unknown variants

Choudhury et al. Nature 2020

Sequence three million genomes across Africa

Wonkam, Nature 2021



Lecture outline

- I. General background & theory of genome assembly
2. Comparison of sequencing technologies
3. Assembly quality and annotation
4. (Not explicitly numbered, implied by the sequence)
5. Assembly workshop with Python programming

Data Formats

HOW STANDARDS PROLIFERATE:
(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC)

SITUATION:
THERE ARE
14 COMPETING
STANDARDS.

14?! RIDICULOUS!
WE NEED TO DEVELOP
ONE UNIVERSAL STANDARD
THAT COVERS EVERYONE'S
USE CASES.



Soon:

SITUATION:
THERE ARE
15 COMPETING
STANDARDS.

Assembly analysis

Base calling, quality control, trimming

- Most data returned in FASTQ format with quality scores included

```
@SEQ_ID          ← id  
GATTTGGGGTTCAAAGCAGTATCGATCAAATA ← sequence  
+                ← description line  
! ' ' * ( ( ( ***+ ) ) Ⓛ Ⓛ Ⓛ ++ ) ( Ⓛ Ⓛ Ⓛ Ⓛ ) . 1 ***- ← base qualities
```

FASTQ

```
@M00747:32:00000000-A16RG:1:1112:15153:29246 1:N:0:1
TCGATCGAGTAACTCGCTGCTGTCAGACTGGTTTGGTCGACTATTGTTCAGTCGCAAGAAT
ATTGTGTCCAGTCGACTGAATTCTGCTGTACGGCCACGGCGGATGCACGGTACAGCAGGCTCAG
ACGGATTAAACTGTT
+
5=9=9<=9 , -5@<<55> , 6+8AC>EE . 88AE9CDD7>+7 . CC9CD+++5@=-FCCA@EF@+**+*-
55--AA---AA-5A<9C+3+<9) 4++=E====<D94) 00=9) ) ) 2@624 (/(/2/-
(. (6;9((((((. (. ' ((6-66<6(///
@M00747:32:00000000-A16RG:1:1112:15536:29246 1:N:0:1
GTAAAATTGAGGTAAATTGTGCGGAATTAGCAATACCGTTTTATTATCACCGGATATCTATT
TGCTGTACGGCCAAGGAGGATGTACGGTACAGCAGGTGCGAACTCACTCCGACGCTCAAGTCAGTGAC
TTAATGATAAGCGTG
+
?????<BBBBBBB5<?BFFFFFFECHEFFFECCFF?9AAC>7@FHHHHHHFG?EAFFG@EEEDEHHDGHHC
BDFFGDFHF)<CCD@F , +3=CFBDFHBD++??DBDEEEDE:):CBEEEBCE68>?) ) 5?**0?:AE*A
*0//:/*:*:**:*.0)
@M00747:32:00000000-A16RG:1:1112:15513:29246 1:N:0:1
GCTAGTCTTGTGTTAGTTATGTTGCATGTTGTAACGGATTCAAACATAGGTGTTGTTCT
TTTATGGTTGTACAATTGGCCCTAACGCCCTACACTTACTGTTGTTCTTTATGGTACGACAT
TTGAGTGGTGGTTGA
+
```

FASTQ

```
@M00747:32:00000000-A16RG:1:1112:15153:29246 1:N:0:1
TCGATCGAGTAACTCGCTGCTGTCAGACTGGTTTGGTCGACTATTGTTCAGTCGCAAGAAT
ATTGTGTCCAGTCGACTGAATTCTGCTGTACGGCCACGGCGGATGCACGGTACAGCAGGCTCAG
ACGGATTAAACTGTT
+
5=9=9<=9 , -5@<<55> , 6+8AC>EE . 88AE9CDD7>+7 . CC9CD+++5@=-FCCA@EF@+**+*-
55--AA---AA-5A<9C+3+<9) 4++=E====<D94) 00=9) ) ) 2@624 (/(/2/-
(. (6;9((((((. (. ' ((6-66<6(///
@M00747:32:00000000-A16RG:1:1112:15536:29246 1:N:0:1
GTAAAATTGAGGTAAATTGTGCGGAATTAGCAATACCGTTTTATTATCACCGGATATCTATT
TGCTGTACGGCCAAGGAGGATGTACGGTACAGCAGGTGCGAACTCACTCCGACGCTCAAGTCAGTGAC
TTAATGATAAGCGTG
+
?????<BBBBBBB5<?BFFFFFFECHEFFFECCFF?9AAC>7@FHHHHHHFG?EAFFG@EEEDEHHDGHHC
BDFFGDFHF)<CCD@F , +3=CFBDFHBD++??DBDEEEDE:):CBEEEBCE68>?) ) 5?**0?:AE*A
*0//:/*:*:**:*.0)
@M00747:32:00000000-A16RG:1:1112:15513:29246 1:N:0:1
GCTAGTCTTGTGTTAGTTATGTTGCATGTTGTAACGGATTCAAACATAGGTGTTGTTCT
TTTATGGTTGTACAATTGGCCCTAACGCCCTACACTTACTGTTGTTCTTTATGGTACGACAT
TTGAGTGGTGGTTGA
+
```

Assembly Quality Scores

Calculating Phred Quality Scores - Base calling accuracy

Q scores are defined as a property that is logarithmically related to the base calling error probabilities (P)².

$$Q = -10 \log_{10} P$$

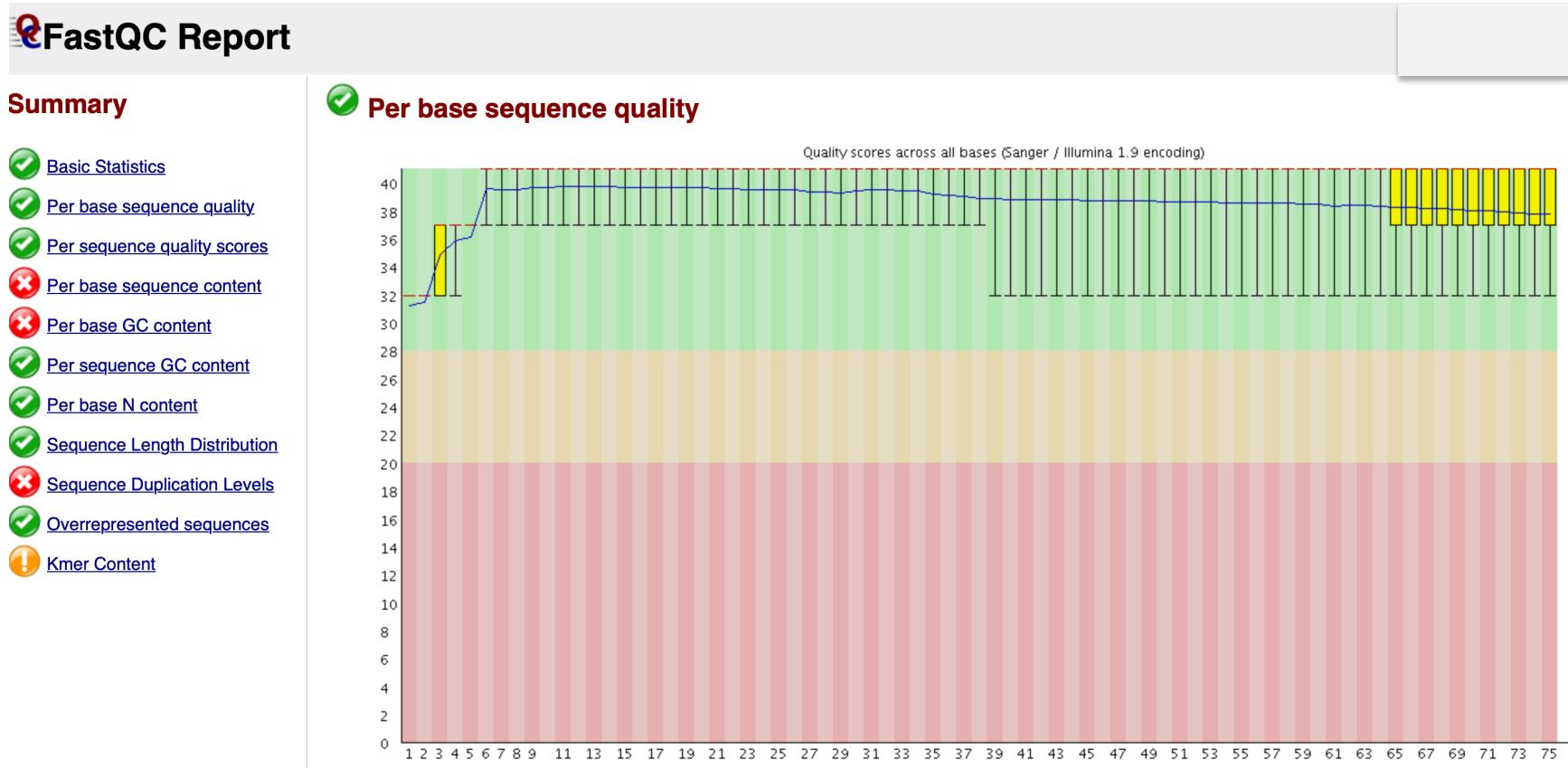
Q - sequencing quality score of a given base Q

P - probability of base call being wrong

Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

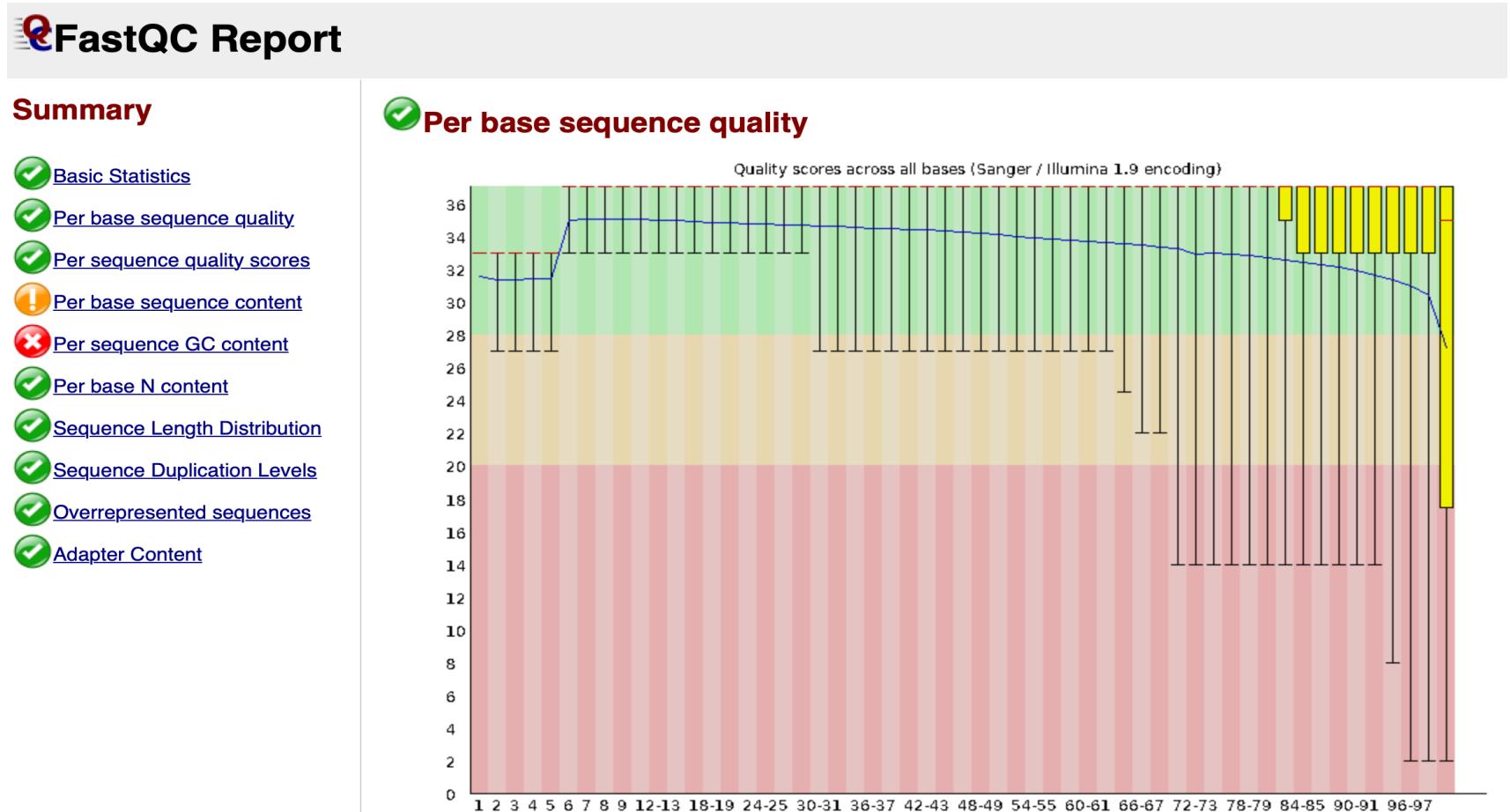
Assembly analysis

Checking quality of reads - FASTQC



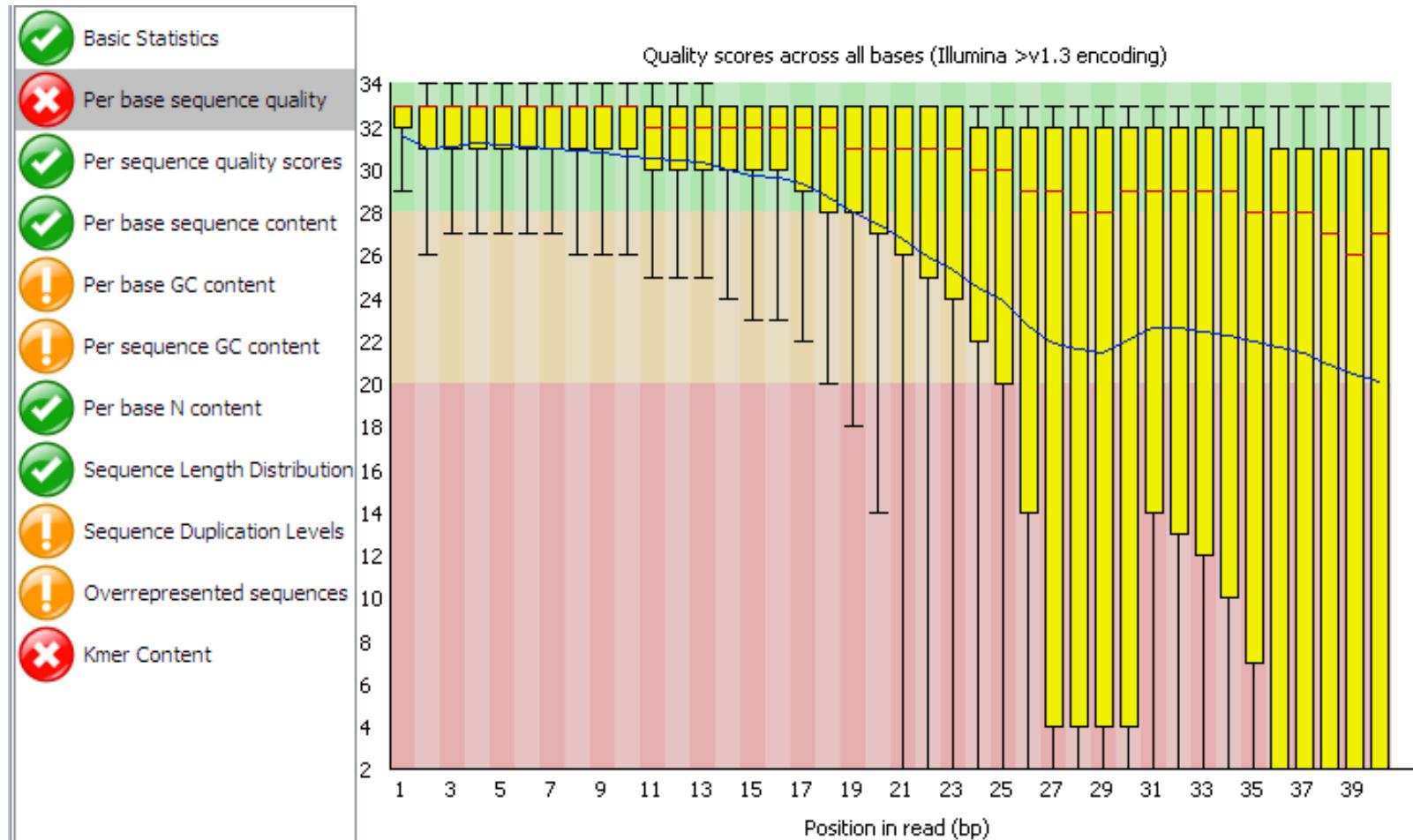
Assembly analysis

Checking quality of reads - FASTQC



Assembly analysis

Checking quality of reads - FASTQC



Sequence trimming

- Trim reads for adapters and quality
- Adapter trimming depends on technology
- Quality trimming, Q20 common cut-off

Trimmomatic - Illumina data, Bolger et al. 2014

<https://github.com/usadellab/Trimmomatic>

DynamicTrim – SolexaQA - quality trimming

<https://solexaqa.sourceforge.net>

AfterQC – filter, trim, QC, Chen et al.

<https://github.com/OpenGene/AfterQC>

Contaminant removal

**Genome assemblies contain *only* genomic sequences from target organism

Contamination removal step often needed

Terminating contamination: large-scale search identifies more than 2,000,000 contaminated entries in GenBank

[Martin Steinegger](#) & [Steven L. Salzberg](#)

[Genome Biology](#) 21, Article number: 115 (2020)

NCBI FCB - Foreign Contamination Screen

<https://github.com/ncbi/fcs#readme>

ContFree-NGS - <https://github.com/labbces/ContFree-NGS>

Kraken2 - <https://github.com/DerrickWood/kraken2>

HoCoRT - <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-023-05492-w>

Genome assembly correction

Pilon protocol

Evaluate alignment pileups

TAATGGGGGCGGTGCCATATCATGAGA
TAATGGGGG**G**CGGTGCCATATC**AT**GAGA
TAATGGGGG**.**CGGTGCCATATC**TA**GAGA
TAATGGGGG**G**CGGTGCCATATC**AT**GAGA

Scan read coverage and alignment discrepancies

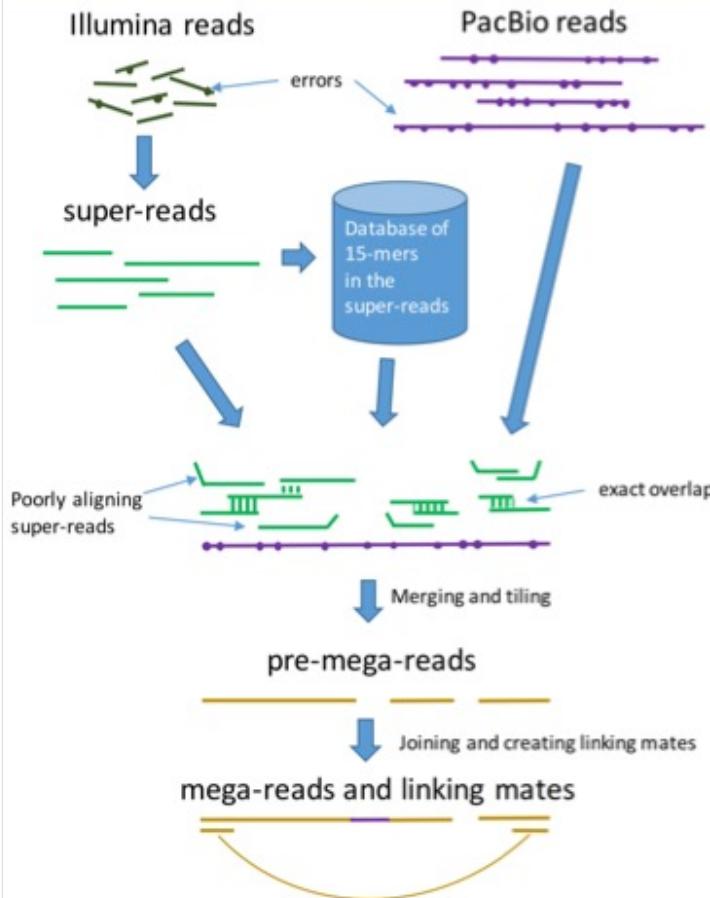


Reassemble across gaps and discordant regions



Pilon: genome assembly improvement

Walker et al. PLoS ONE. 2014



MaSuRCA mega-reads algorithm

Zimin et al. Genome Res. 2017

Assessing genome quality

- Map raw reads back to assembled genome
 - mapping back uniformly
- SAM/BAM file
- BEDTools - to retrieve coverage statistics



<https://github.com/arq5x/bedtools2>

File formatting

- FASTA
- FASTQ
 - quality scores
- SAM/BAM
 - Developed for NGS data
 - Sequence Alignment Map
 - Stores alignment information

FASTA

```
>NP_000552.2 Human glutathione transferase M1 (GSTM1)
MPMILGYWDIRGLAHAIRLLYEYTDSSYEEKKYTMGDAPDYDRSQWLNEKFKLGLDFPNLPYLIDGAH
KITQSNAILCYIARKHNLCGETEEEKIRVDILENQTMDNHMQLGMICYNPEFEKLKPYLEELPEKLK
LYSEFLGKRPWFAGNKITFVDFLVYDVL_DLHRIFEPKCLDAFPNLKDFISRFEGLEKISAYMKSSRFL
PRPVFSKMAVWGNK
```

SAM/BAM Sequence Alignment Map

- Alignment file - provides context for raw data
 - Eleven columns, tab delimited
 - One alignment record per line
- SAM is plain-text (human readable)
- BAM is a binary format
- SAMTools - suite of utilities for SAM/BAM files
- Picard - tools for sequencing data

samtools: <http://samtools.sourceforge.net>

Picard: <https://broadinstitute.github.io/picard/>

SAM/BAM Sequence Alignment Map

```
D4ZHLFP1:53:D2386ACXX:6:2115:17945:68812      0      Mle_000001      18      42      108M *      0      0
TCCCCCTGCATGTGCCGTGGCTGGATGCCATGCTCCATGCAGTATAGCTCCCAGCATGAGTTACCGATCTGGACACCTGCTTG
GCCAAGATGTACTGAGATGCAT
C@CFDFFFHHGHHFGGBFEGGDGGGGEHGGGGJJJIIIGIIB9BFBBFHGGHICEAGHGEGEDHIGEEDBECCACBDDC@CCDBCDD<
?2+4>@4>>CCCCAA@@    AS:i:-5    XN:i:0    XM:i:1    XO:i:0    XG:i:0    NM:i:1    MD:Z:0A107
YT:Z:UU

D4ZHLFP1:53:D2386ACXX:7:2110:5214:83081 0      Mle_000001      18      42      108M *      0      0
TCCCCCTGCATGTGCCGTGGCTGGATGCCATGCTCCATGCAGTATAGCTCCCAGCATGAGTTACCGATCTGGACACCTGCTGGCAA
GATGTACTGAGATGCAT
CCCFFFFFHHHHHGGEGIJIIIGJFHJJJJIIJJIIGIJIJFHJJIIJJFIIIIIIIIJJHHFFFCEEEEDDDDDDDDDDDDD
BDCDDEEEEDDDDDDDDD    AS:i:-5    XN:i:0    XM:i:1    XO:i:0    XG:i:0    NM:i:1    MD:Z:0A107
YT:Z:UU

D4ZHLFP1:53:D2386ACXX:7:2206:9985:31556 0      Mle_000001      18      42      108M *      0      0
TCCCCCTGCATGTGCCGTGGCTGGATGCCATGCTCCATGCAGTATAGCTCCCAGCATGAGTTACCGATCTGGACACCTGCTGGCAA
GATGTACTGAGATGCAT
CCCFFFFFHHHHJJJJIHJJIIIIJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
DDCD@@CDCCDDCDCDC    AS:i:-5    XN:i:0    XM:i:1    XO:i:0    XG:i:0    NM:i:1    MD:Z:0A107
YT:Z:UU
```

Helpful site for looking up SAM flag: <https://broadinstitute.github.io/picard/explain-flags.html>

Genome annotation

- Describe genetic and genomic features in raw sequence.
 - genetics and genomic features, genes, repetitive elements, mobile elements, genome duplications
 - gene prediction, ORF searches, repeat region identification, homology searches
 - most genome projects use automated methods with some manual curation
 - community motivated annotation projects

Genome Assembly

- Recommended to use multiple assemblers with different parameters to assess results – don't run with just default!
- How to assess our results?
 - Number of contigs
 - Longest contig
 - N50 - largest length for which 50% of all nucleotides are contained in the contigs of at least that length
 - L50 - number of contigs that are as long or longer than the N50 value

N50 size

If we place our contigs from largest to smallest on the genome, 50% of the genome in contigs as long as or larger than N50 value

Example: 1 Mbp genome



1000



N50 size = 30 kbp

(300k+100k+45k+45k+30k = 520k >= 500kbps)

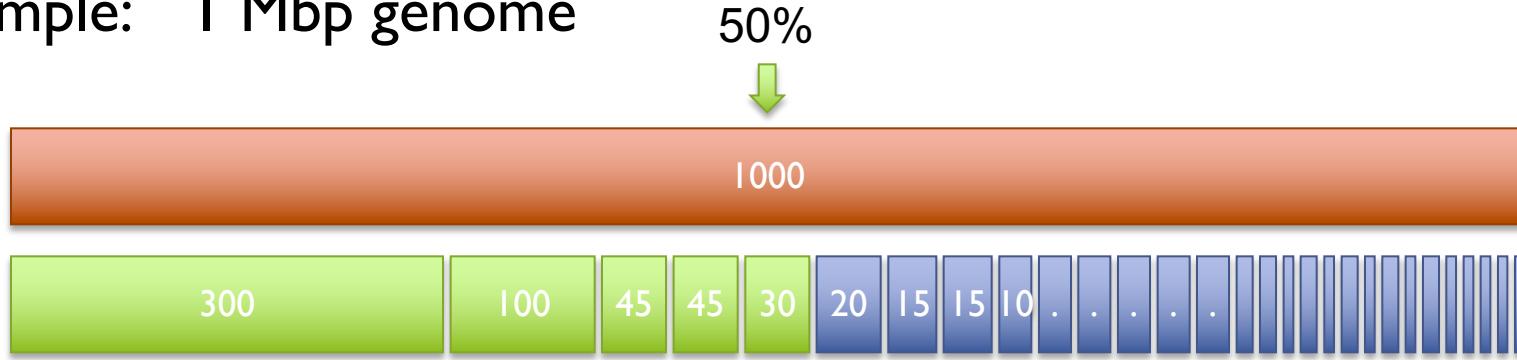
A greater N50 is usually a sign of assembly improvement

- Comparable with genomes of similar size
 - Genome composition can bias comparisons
 - Low L50 vs High N50

L50 size

Number of contigs that are as long or longer than the N50 value

Example: 1 Mbp genome



N50 size = 30 kbp

$$(300k + 100k + 45k + 45k + 30k = 520k \geq 500\text{kbp})$$

L50 - number of contigs that sum to N50 length

L50 = how many?

- Low L50 vs High N50
 - longer sequences and fewer of them...in theory
 - lower stringency can inflate N50

L50/N50 size

"N50 length" (L50) first defined:

articles

Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

* A partial list of authors appears on the opposite page. Affiliations are listed at the end of the paper.

The human genome holds an extraordinary trove of information about human development, physiology, medicine and evolution. Here we report the results of an international collaboration to produce and make freely available a draft sequence of the human genome. We also present an initial analysis of the data, describing some of the insights that can be gleaned from the sequence.

International Human Genome Sequencing Consortium, 2001. *Nature*, 409(6822), p.860.

L50/N50 size

Box 1

Genome glossary

Sequence

Raw sequence Individual unassembled sequence reads, produced by sequencing of clones containing DNA inserts.

Paired-end sequence Raw sequence obtained from both ends of a cloned insert in any vector, such as a plasmid or bacterial artificial chromosome.

Finished sequence Complete sequence of a clone or genome, with an accuracy of at least 99.99% and no gaps.

Coverage (or depth) The average number of times a nucleotide is represented by a high-quality base in a collection of random raw sequence. Operationally, a ‘high-quality base’ is defined as one with an accuracy of at least 99% (corresponding to a PHRED score of at least 20).

Full shotgun coverage The coverage in random raw sequence needed from a large-insert clone to ensure that it is ready for finishing; this varies among centres but is typically 8–10-fold. Clones with full shotgun coverage can usually be assembled with only a handful of gaps per 100 kb.

Half shotgun coverage Half the amount of full shotgun coverage

Sequenced-clone contigs

Contigs produced by merging overlapping sequenced clones.

Sequenced-clone-contig scaffolds Scaffolds produced by joining sequenced-clone contigs on the basis of linking information.

Draft genome sequence The sequence produced by combining the information from the individual sequenced clones (by creating merged sequence contigs and then employing linking information to create scaffolds) and positioning the sequence along the physical map of the chromosomes.

N50 length A measure of the contig length (or scaffold length) containing a ‘typical’ nucleotide. Specifically, it is the maximum length L such that 50% of all nucleotides lie in contigs (or scaffolds) of size at least L .

Computer programs and databases

PHRED A widely used computer program that analyses raw sequence to produce a ‘base call’ with an associated ‘quality score’ for each position in the sequence. A PHRED quality score of X corresponds to an error probability of approximately $10^{-X/10}$. Thus, a PHRED quality score of 30 corresponds to 99.9% accuracy for the base call in the raw read.

Genome assembly quality

- How to assess our results?

Alignments:

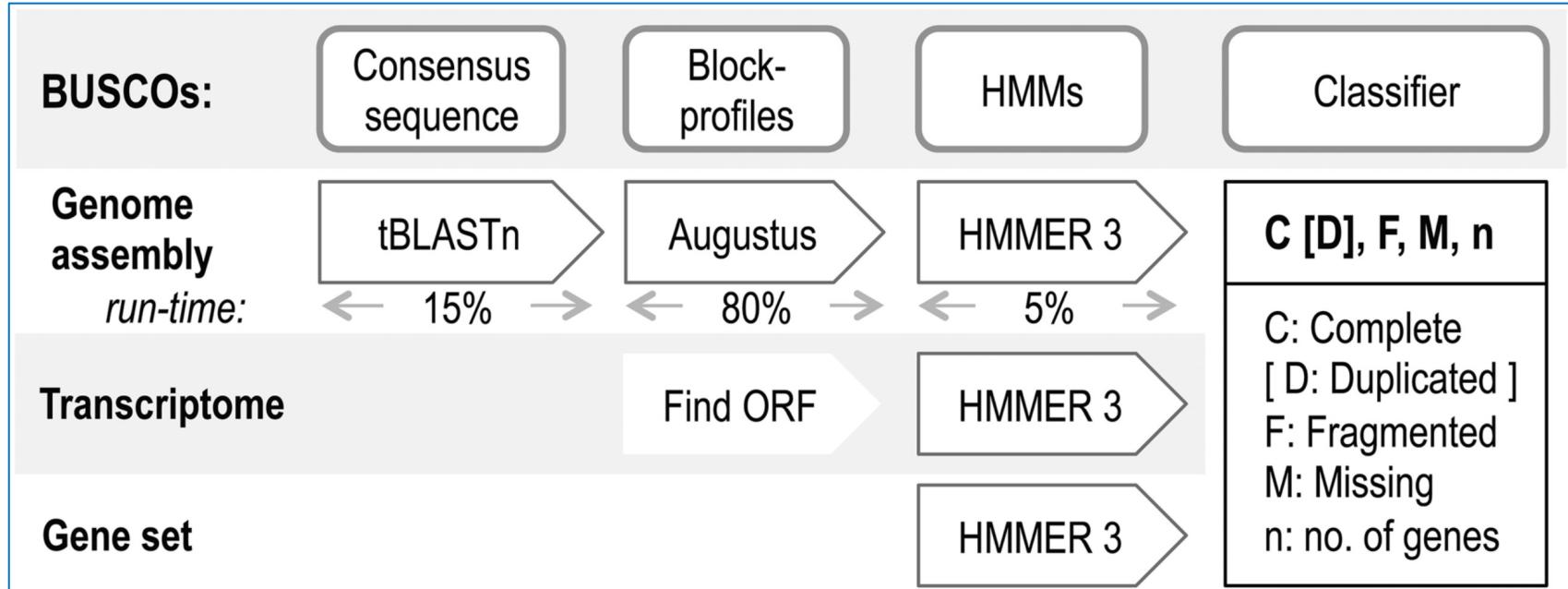
- compare to a reference genome
- align reads to your assembled genome
- assess repetitive regions
- Call SNPs

Check for completeness:

- annotation, blast against reference gene set
- Busco - Simao et al. 2015. Bioinformatics

BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs

Bioinformatics. 2015;31(19):3210-3212. doi:10.1093/bioinformatics/btv351



BUSCO assessment workflow and relative run-times

Quality of genome vs. completeness

Orthology vs Paralogy

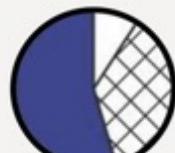
- Homolog - share a common ancestor
- Ortholog - diverged after a speciation event from ancestral gene
- Paralog - diverged after a duplication event within the same lineage

You are either homologous (share a common ancestor) or not. What varies is ability to detect homology.

- Evolution of gene families
- Species tree constructions
- genome evolution

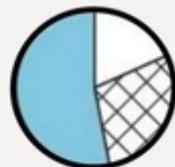
BUSCO sampling space

1. High universality



Vertebrata

Mouse's
orthologous groups



Arthropoda

Fly's
orthologous groups



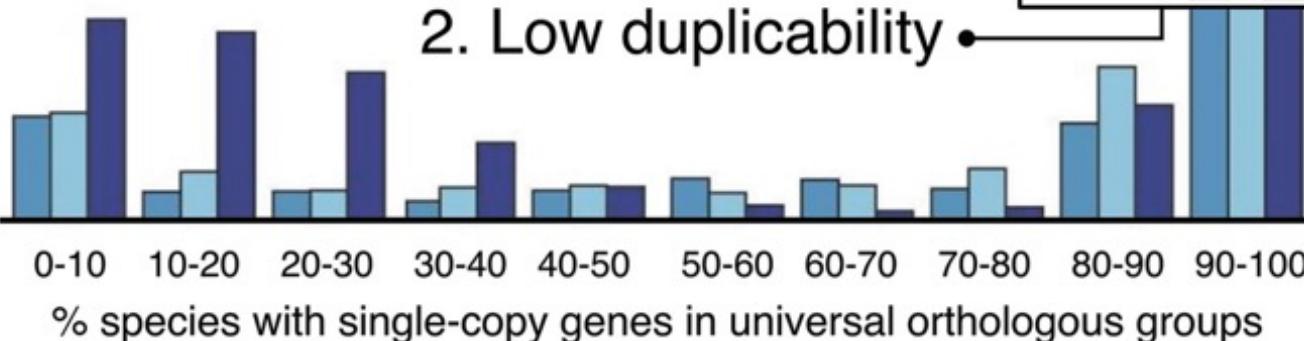
Fungi

Yeast's
orthologous groups

Orthologs present in
> 90% of the species
(considered as universal)

50-90%
0-50%

> 90% of the species
with single-copy genes



2. Low duplicability

Meta-Analysis > Mol Ecol Resour. 2021 Jul;21(5):1416-1421. doi: 10.1111/1755-0998.13364.

Epub 2021 Mar 9.

Assessing genome assembly quality prior to downstream analysis: N₅₀ versus BUSCO

April A Jauhal ¹ ², Richard D Newcomb ²

Affiliations + expand

PMID: 33629477 DOI: [10.1111/1755-0998.13364](https://doi.org/10.1111/1755-0998.13364)

- High N₅₀ reliably predicted a complete BUSCO score
- High BUSCO scores from assemblies with low N₅₀

Genome assembly quality

- Do not take first version. Try correcting/polishing your genome
- Always distrust your data!
 - Go back and reassess your genome
 - plotting, quality-control
 - Number of contigs/scaffolds change?
 - L50/N50 go up or down?
 - How much is fragmented?
 - Always new technologies and improvements

What should we expect from an assembly?

- Annotation of assembly
- Comparison to closely related genomes
- Gene content
- Percent repetitive
- Chromosomes
- Another estimate
 - Flow cytometry
 - kmer distribution

Future of genome assembly

Genome download and stats

Genome assembly workshop

- Genome assembly with PacBio data using Canu assembler
- Python programming exercise