

Genome assembly quality control

Bioplatforms Fungi Genomics Workshop 2024
Australian National University

Rita Tam



ritahltam

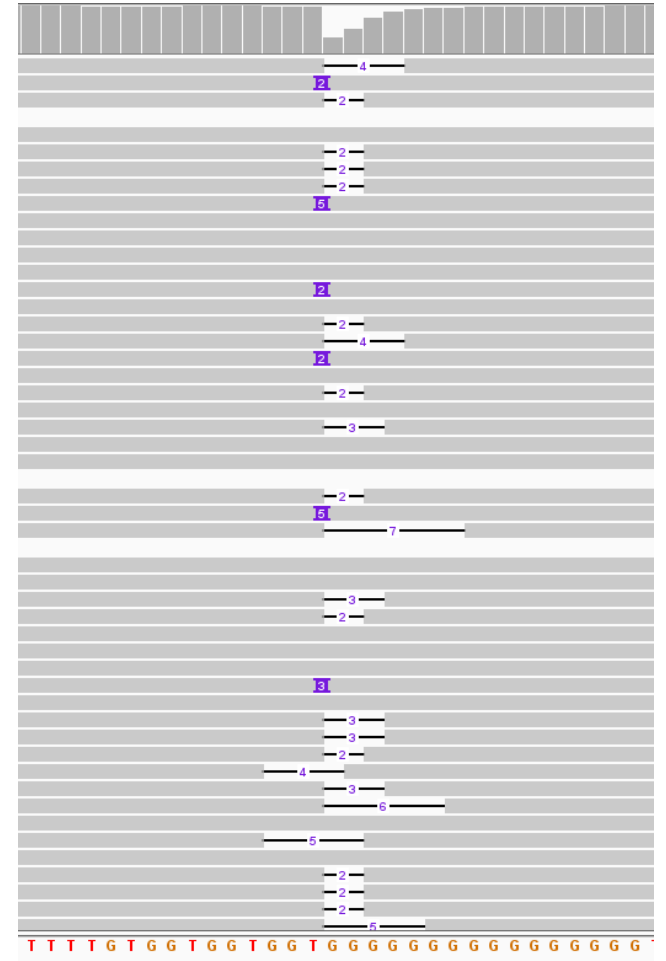


ritatam

How assembly can fail perfection

Sequencing data

- Insufficient read depth
- DNA sequencing errors (biases)
 - Lower read quality towards read 3' end
 - Prefers low-GC
 - Homopolymers (e.g. AAAAAAA or CCCCCC)
- Contaminant DNA
 - Bacteria, unexpected fungal species, host plant, etc.



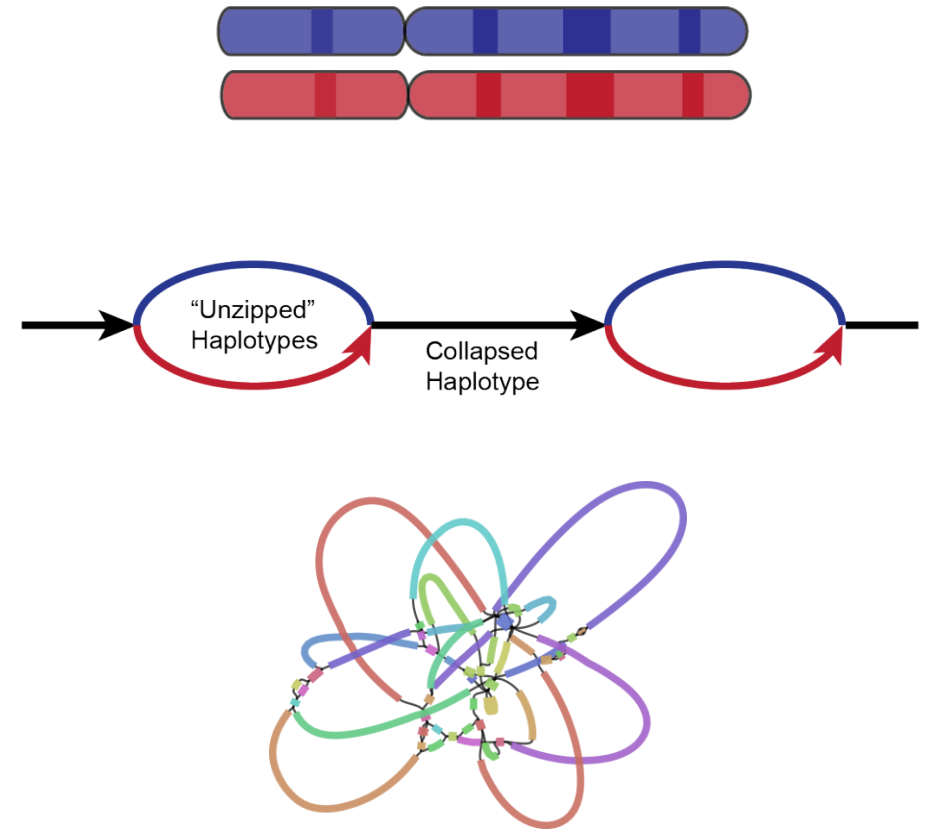
How assembly can fail perfection

Genome complexity

- Reads not long enough to span repeats
- Heterozygosity (if diploid)
e.g. confusion at highly homozygous regions

Assembly algorithms

- Tolerance for sequencing errors
- Ploidy-awareness



Why is it important to evaluate assembly quality?

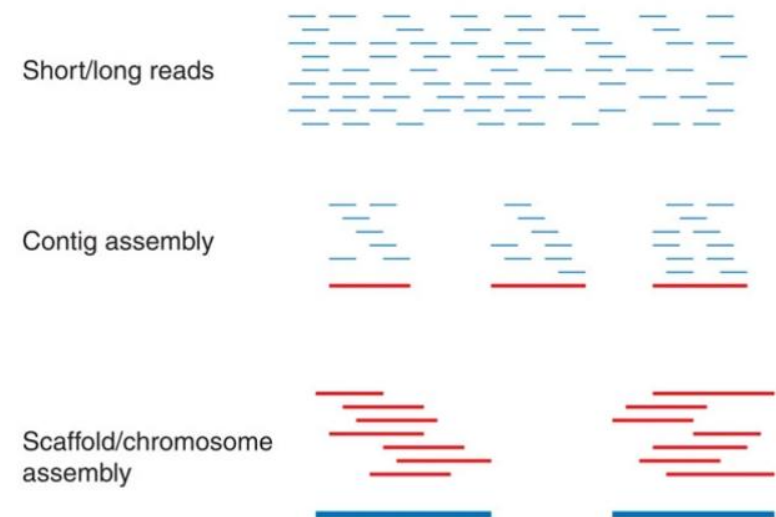
Know how trustworthy your genome assembly is!

If errors go unchecked, they can be

- propagated downstream
- misinterpreted as true biological events (e.g. a deletion)!

Users can make careful decisions

- revisit library prep?
- next steps?
 - scaffolding
 - curation
 - annotation



First – BLAST check your assembly!

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

NEWS

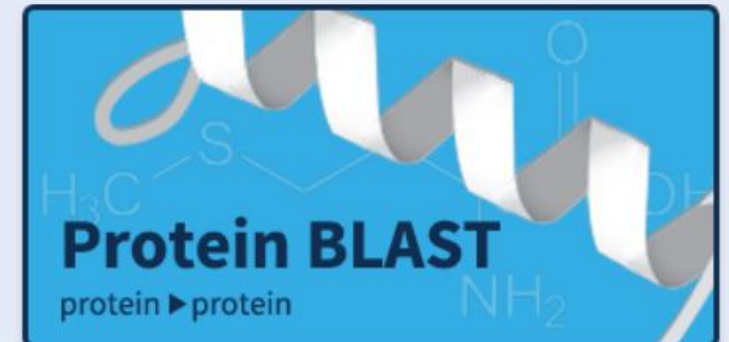
Non-interactive searches of nt switch to core_nt

Starting late September 2024 all non-interactive WebBLAST and PrimerBLAST searches of ``nt`` will

Tue, 24 Sep 2024

[More BLAST news...](#)

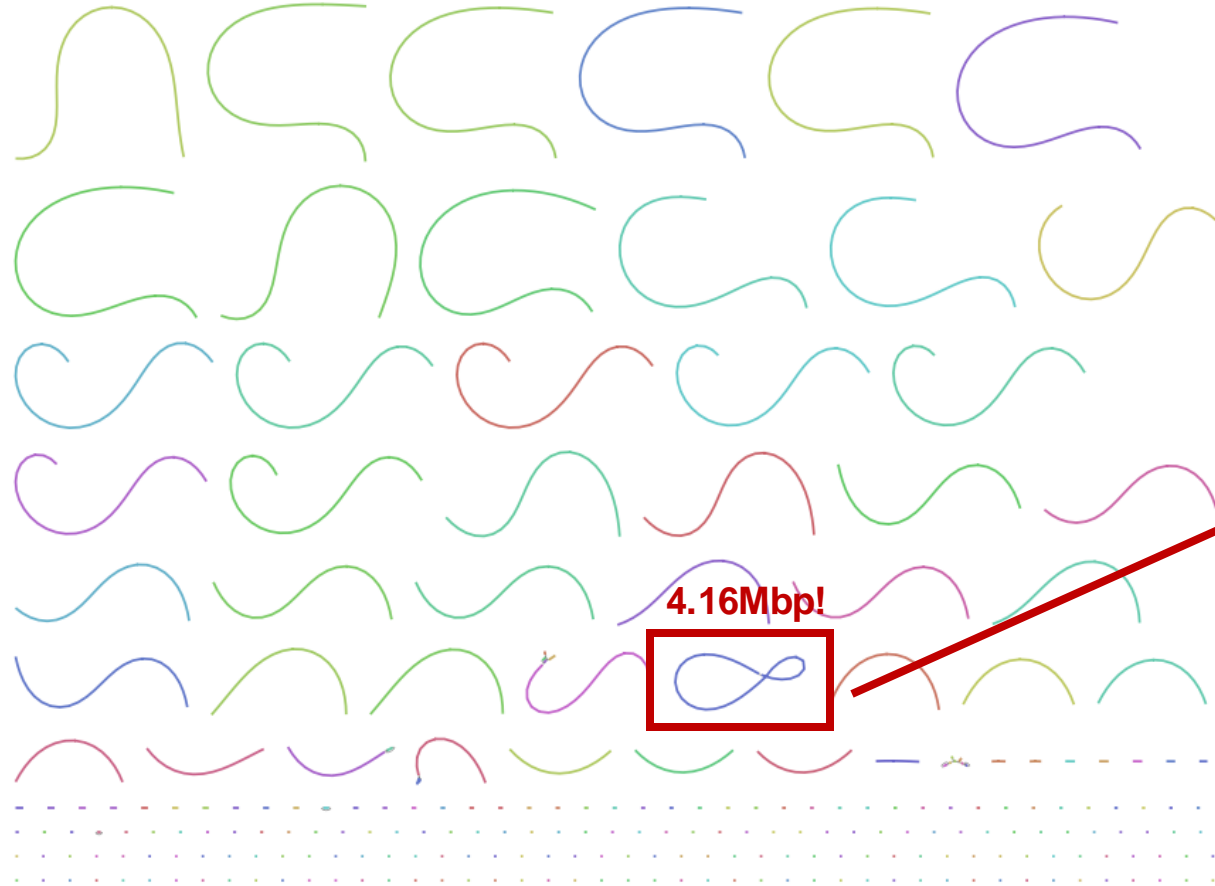
Web BLAST



First – BLAST check your assembly!



Puccinia striiformis f. sp. *tritici*
(wheat stripe rust fungi)



BLAST to detect and remove contaminant (e.g. wheat, bacteria) & mtDNA contigs.



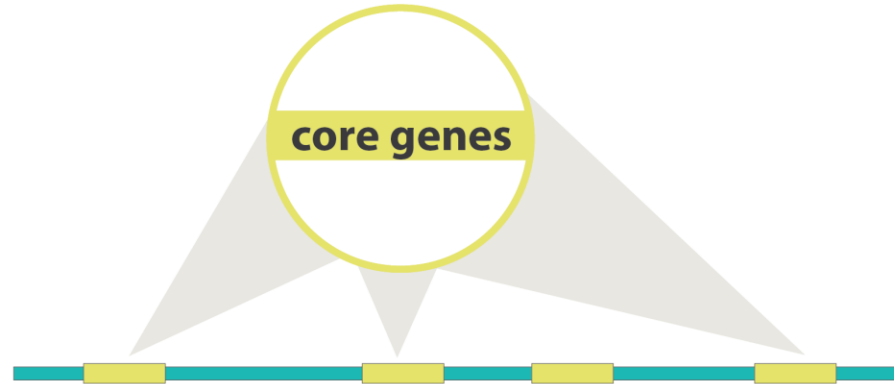
Herbaspirillum genome fully assembled along with the fungal genome!

The “3C” rules for assembly QC

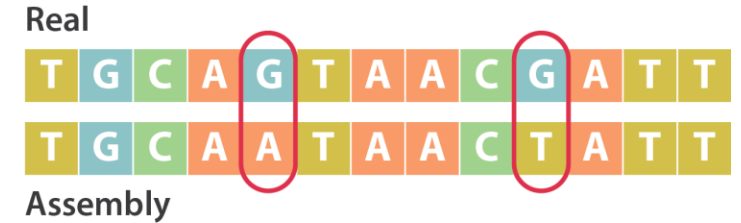
Contiguity



Completeness



Correctness

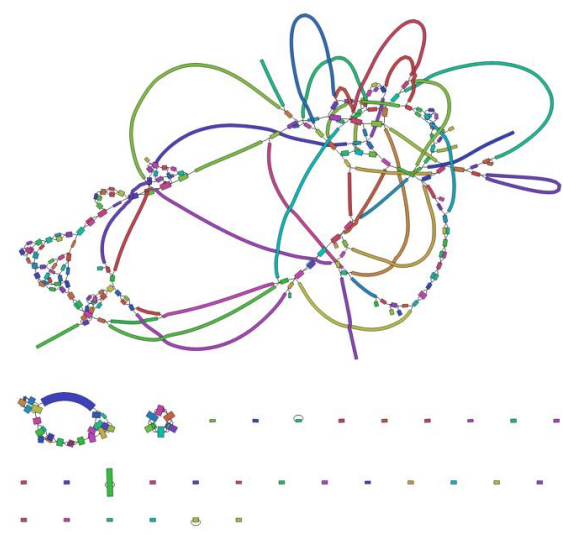


Contiguity

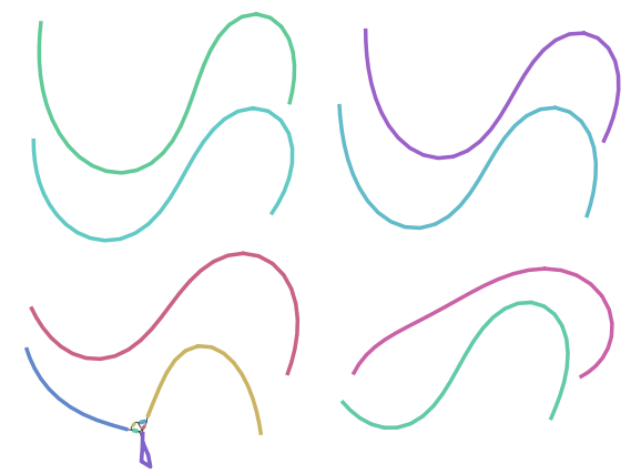
What we want:
Fewer and longer contigs,
chromosome-scale



Fragmented



Chromosome-scale



N50/L50

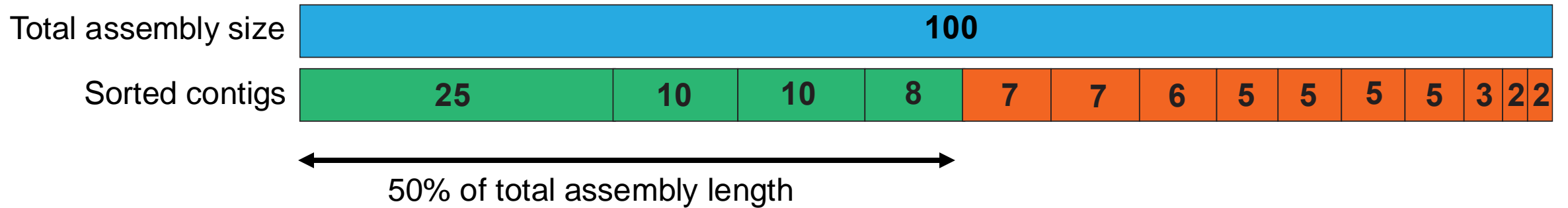
N50

Sequence length of the shortest contig at 50% of the assembly length.
Generally, the higher the better.

L50

Smallest count of contigs whose length sum make up 50% of the assembly.
Generally, the lower the better.

N50/L50



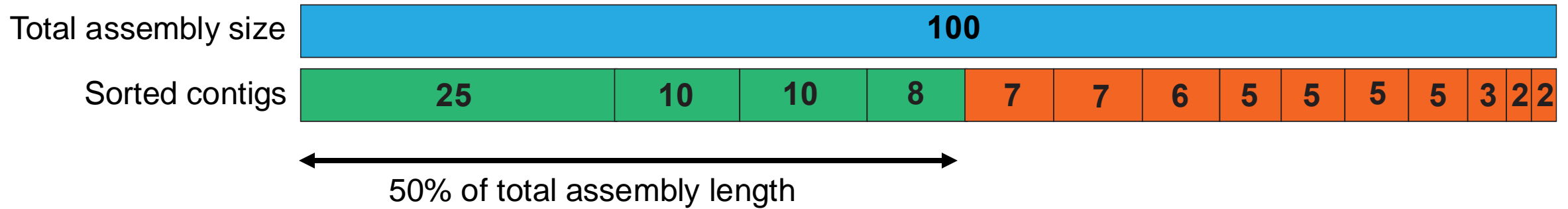
Total assembly size = 100
Sort contigs by length

50% of total length is contained within sequences of at least 8 ($25+10+10+8 = 53, \geq 50$)

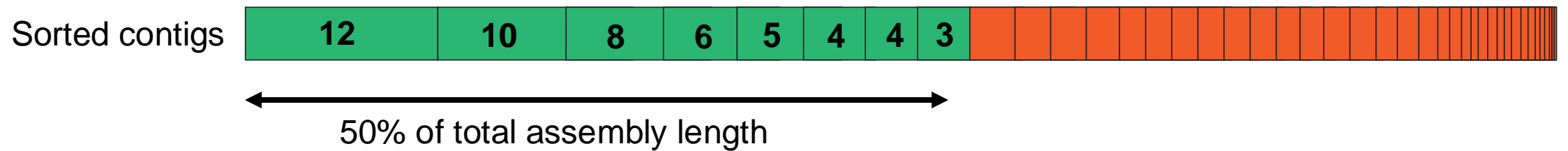
N50 = 8

L50 = 4

N50/L50



N50 = 8
L50 = 4

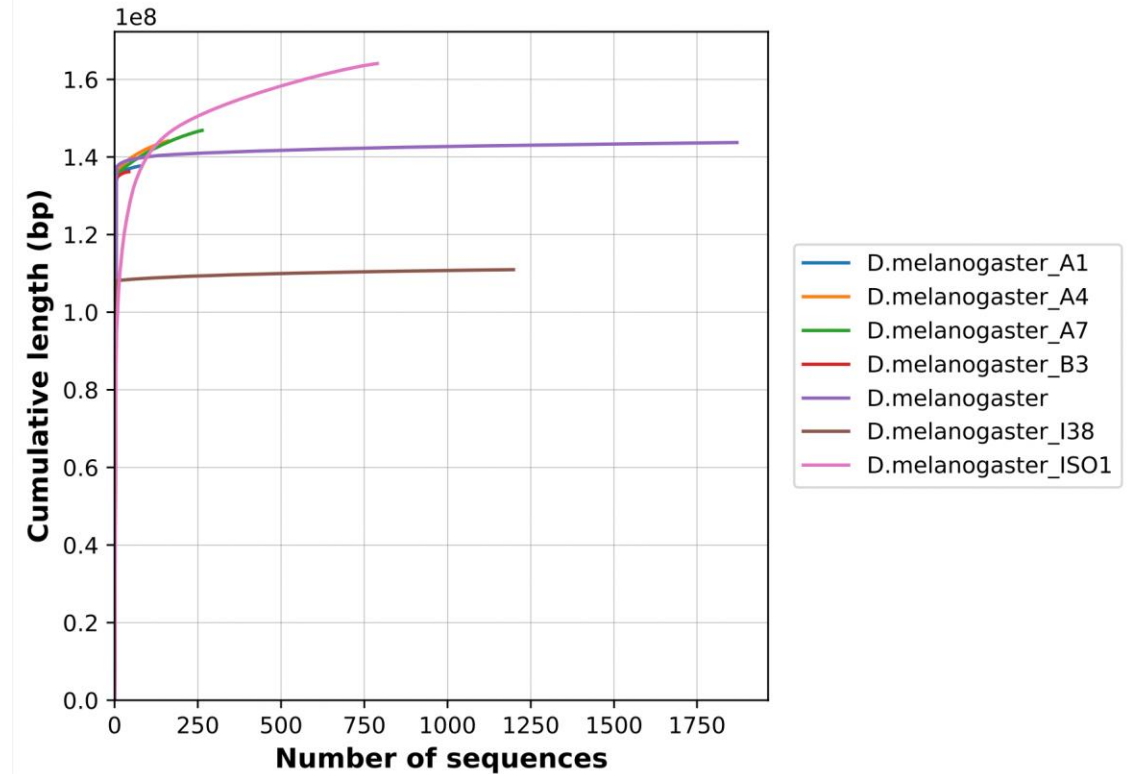
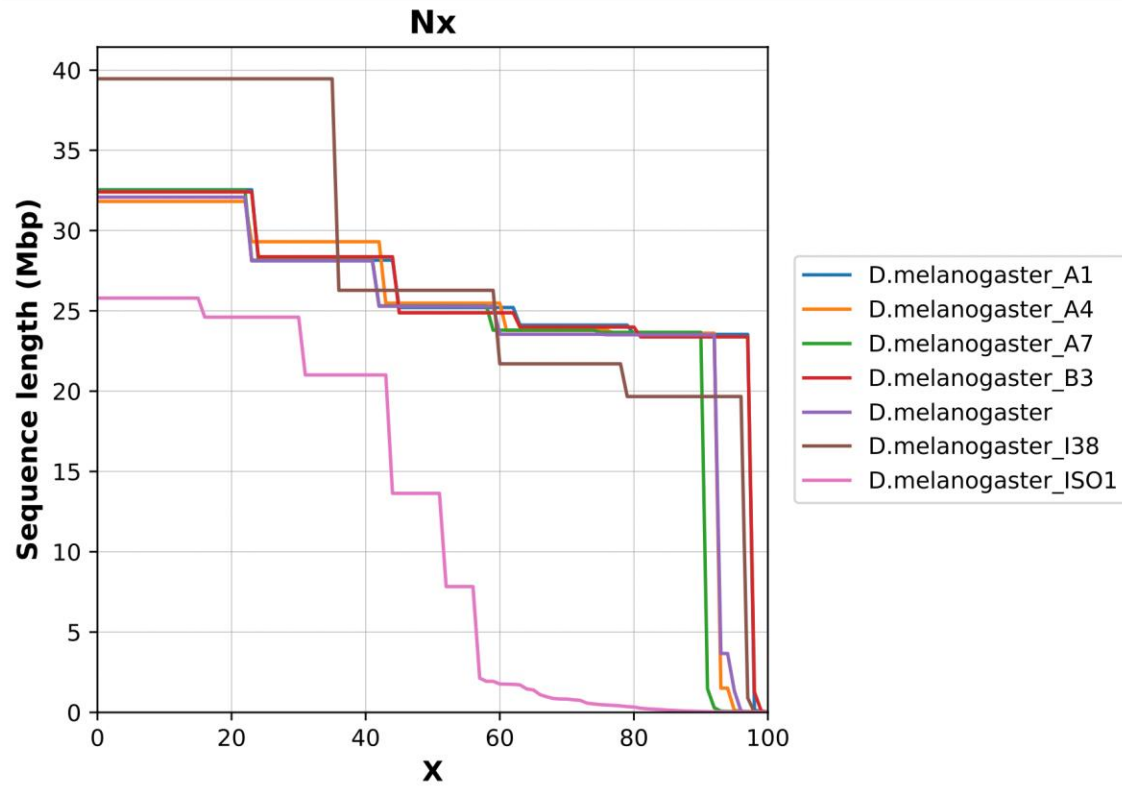


N50 = 3
L50 = 8

Nx curve

QUAST

Quality Assessment Tool for Genome Assemblies by [CAB](#)



Contiguity measured at repeats

- LTR retrotransposons are abundant in fungi (Muszewska et al. 2017)
- Assembly contiguity can be measured at repeats, e.g. LTR transposons

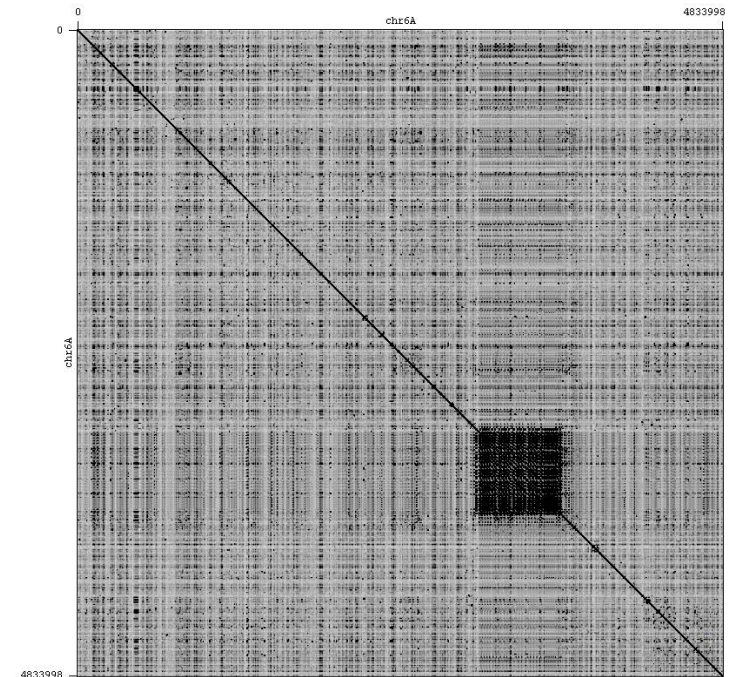
LTR Assembly Index (LAI) (Ou et al. 2018)

Category	LAI
Draft	$0 \leq \text{LAI} < 10$
Reference	$10 \leq \text{LAI} < 20$
Gold	$20 \leq \text{LAI}$

Limitations

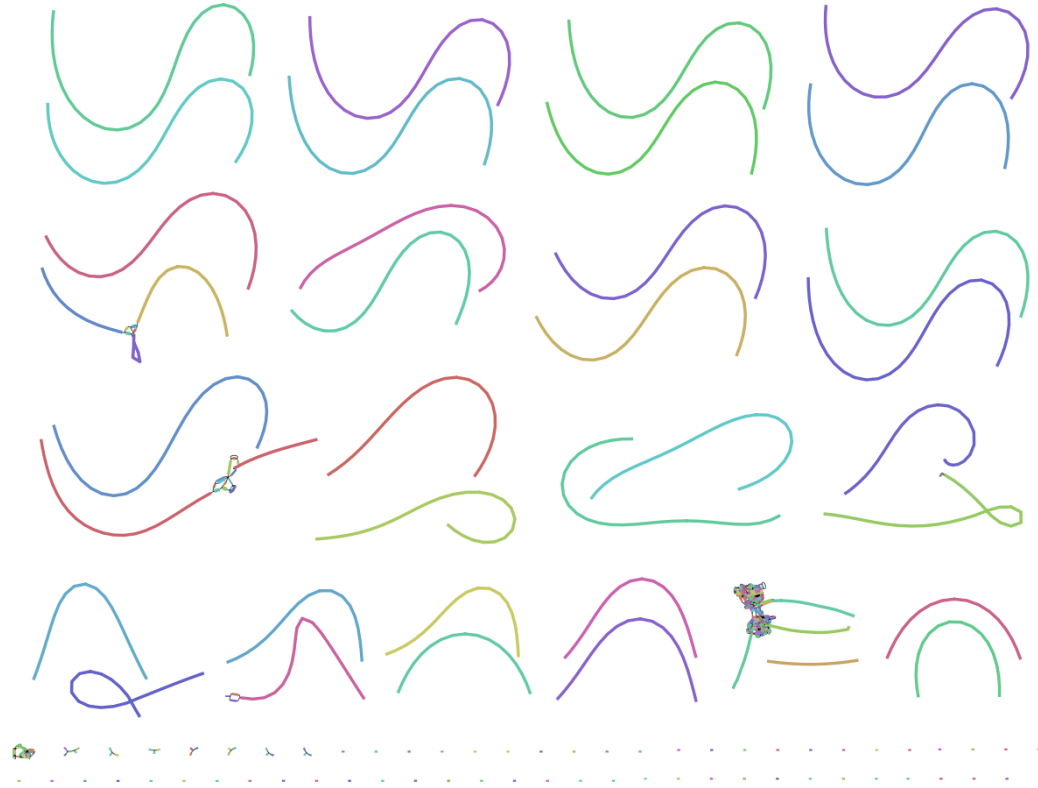
- Less suitable for species with low LTR abundance
- Intrinsic degradation rate of LTRs varies between species
- **General guide only, low LAI score \neq bad!**

LTR retrotransposon-dominated region near *P. striiformis* f.sp. *tritici* mating type locus



Genome graph

Quick visualisation for a “feeling” of contiguity



Good sign:

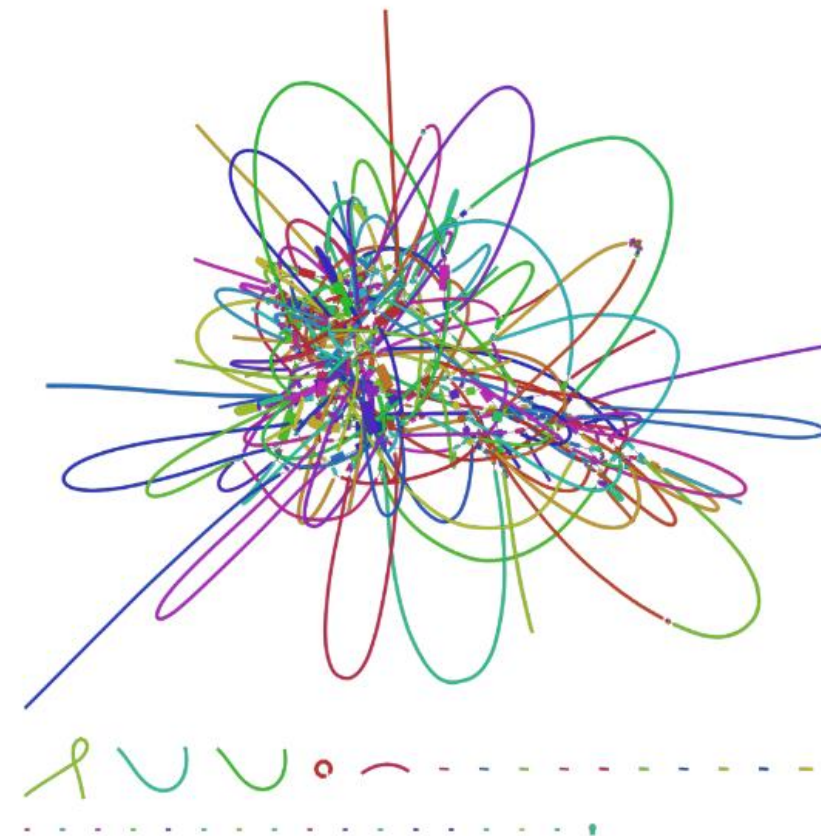
- Chromosome-scale contigs

assembly_output.fasta

assembly_output.gfa



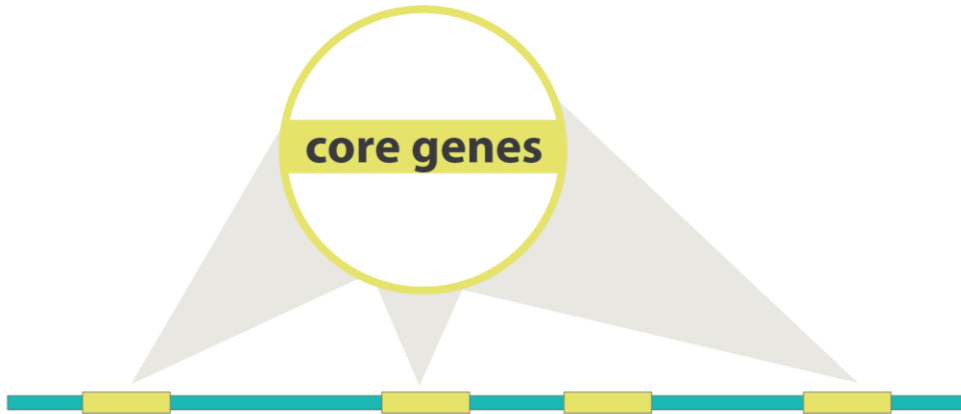
Bandage



Not so good sign:

- Chromosomes broken into small contigs
- Ambiguous contig connections

Completeness



Assembly size

Gene space completeness

- BUSCO
- Reference-based

Telomere counting

Assembly gap

Assembly size vs Expected

$$\frac{\textit{Assembly size}}{\textit{Expected genome size}}$$

Compare the assembly size to expected genome size

Very rough estimation

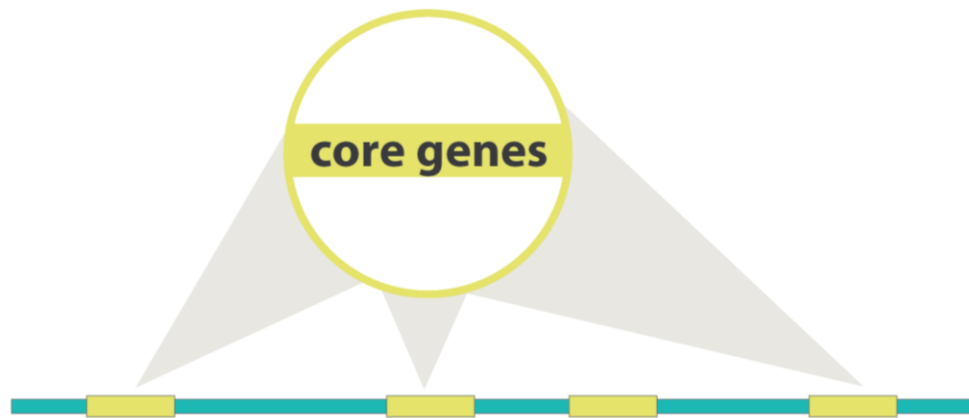
Limitations

- Empirical evidence can have inaccuracies
- Genome size can vary between species, even genotypes
- Contaminant DNA can inflate it

Gene space completeness - BUSCO

BUSCO: Benchmarking Universal Single-Copy Orthologs

Evaluate “core” gene content using a predefined set of highly conserved orthologs



BUSCO sampling space

1. High universality

Vertebrata
Mouse's orthologous groups

Arthropoda
Fly's orthologous groups

Fungi
Yeast's orthologous groups

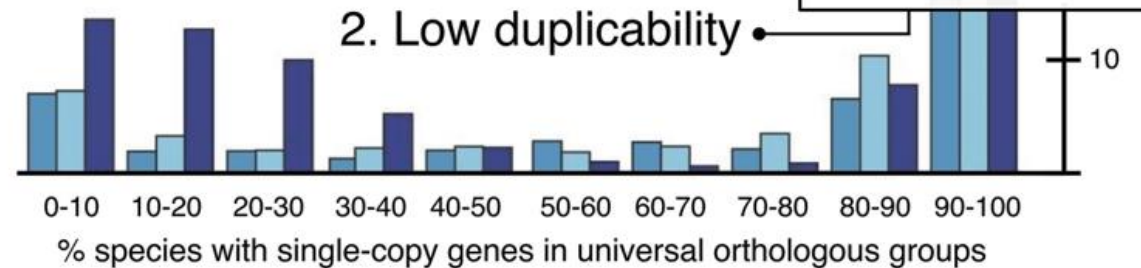
Orthologs present in
> 90% of the species
(considered as universal)

50-90%
0-50%

% universal orthologous groups

> 90% of the species
with single-copy genes

2. Low duplicability



Gene space completeness - BUSCO

```
lineages
├─ agaricales_odb10
├─ agaricomycetes_odb10
├─ archaea_odb10
├─ ascomycota_odb10
├─ bacillales_odb10
├─ bacteria_odb10
├─ endopterygota_odb10
├─ eudicots_odb10
├─ eukaryota_odb10
├─ fungi_odb10
├─ hemiptera_odb10
├─ insecta_odb10
├─ mammalia_odb10
├─ metazoa_odb10
├─ methanococcales_odb10
├─ natrialbales_odb10
├─ primates_odb10
├─ saccharomycetes_odb10
├─ solanales_odb10
```

Fungi datasets Odb10

Kingdom: [Fungi Odb10](#)

Phylum level:

[Ascomycota](#)
[Basidiomycota](#)
[Microsporidia](#)
[Mucoromycota](#)

Class level:

[Agaricomycetes](#)
[Dothideomycetes](#)
[Eurotiomycetes](#)
[Leotiomycetes](#)
[Saccharomycetes](#)
[Sordariomycetes](#)
[Tremellomycetes](#)

Order level:

[Agaricales](#)
[Boletales](#)
[Capnodiales](#)
[Chaetothyriales](#)
[Eurotiales](#)
[Glomerellales](#)
[Helotiales](#)
[Hypocreales](#)
[Mucorales](#)
[Onygenales](#)
[Pleosporales](#)
[Polyporales](#)

Gene space completeness - BUSCO

BUSCO lineage dataset
for Basidiomycota

```
# BUSCO version is: 5.5.0
# The lineage dataset is: basidiomycota_odb10 (Creation date: 2024-01-08,
# number of genomes: 133, number of BUSCOs: 1764)
# Summarized benchmarking in BUSCO notation for file /media/ssd/rita/proj
# ect/104e/assembly_versions/v3.9_gapfill/v3.9.chr.haplotype-paired.fasta
# BUSCO was run in mode: euk_genome_met
# Gene predictor used: metaeuk
```

input fungal genome
assembly (.fasta)

```
***** Results: *****
```

```
C:92.6%[S:4.0%,D:88.6%],F:0.9%,M:6.5%,n:1764
1634   Complete BUSCOs (C)
  71   Complete and single-copy BUSCOs (S)
1563   Complete and duplicated BUSCOs (D)
  16   Fragmented BUSCOs (F)
  114  Missing BUSCOs (M)
1764   Total BUSCO groups searched
```

BUSCO completeness results

Gene space completeness - BUSCO

BUSCO limitations

- “Completeness” inferred from a small gene subset
- Some species can lose conserved genes and yield lower BUSCO score, even if well-assembled
- Genes with long introns are hard to detect
- Better option: annotate genes with RNA-seq evidence, then assess BUSCO

Gene space completeness – reference-based

Reference-based

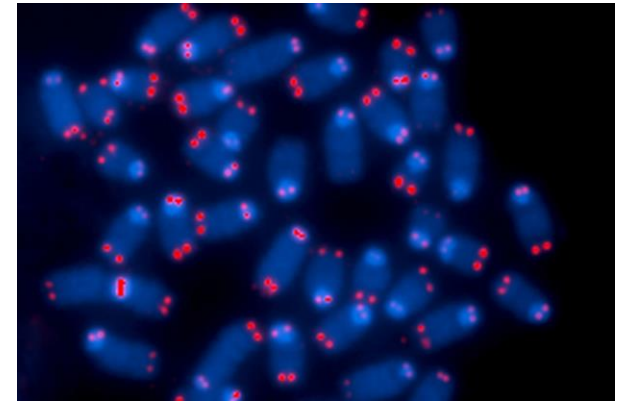
If previously annotated genes are available,
can count full and partial genes in your draft assembly

QUAST gene counting output

	assembly	genes	partial genes
draft assembly V1	verkko_duplex_assembly	31032	6
draft assembly V2	verkko_herro_assembly	31058	2

Telomeres

- “Endpoints” of a chromosome assembly
- Telomere-to-telomere assembly now highly achievable, thanks to long reads



Article | [Open access](#) | Published: 30 May 2024

Complete telomere-to-telomere genomes uncover virulence evolution conferred by chromosome fusion in oomycete plant pathogens

[Zhichao Zhang](#), [Xiaoyi Zhang](#), [Yuan Tian](#), [Liyuan Wang](#), [Jingting Cao](#), [Hui Feng](#), [Kainan Li](#), [Yan Wang](#), [Suomeng Dong](#), [Wenwu Ye](#)  & [Yuanchao Wang](#) 

RESOURCE ANNOUNCEMENT



A Telomere-to-Telomere Genome Assembly Resource of *Bipolaris sorokiniana*, the Fungal Pathogen Causing Spot Blotch and Common Root Rot Diseases in Barley and Wheat

[Yueqiang Leng](#), [Yang Du](#), [Jason Fiedler](#), [Sajeet Haridas](#), [Igor V. Grigoriev](#), and [Shaobin Zhong](#) 

Affiliations 

Published Online: 17 Jan 2024 | <https://doi.org/10.1094/PHYTOFR-08-23-0108-A>

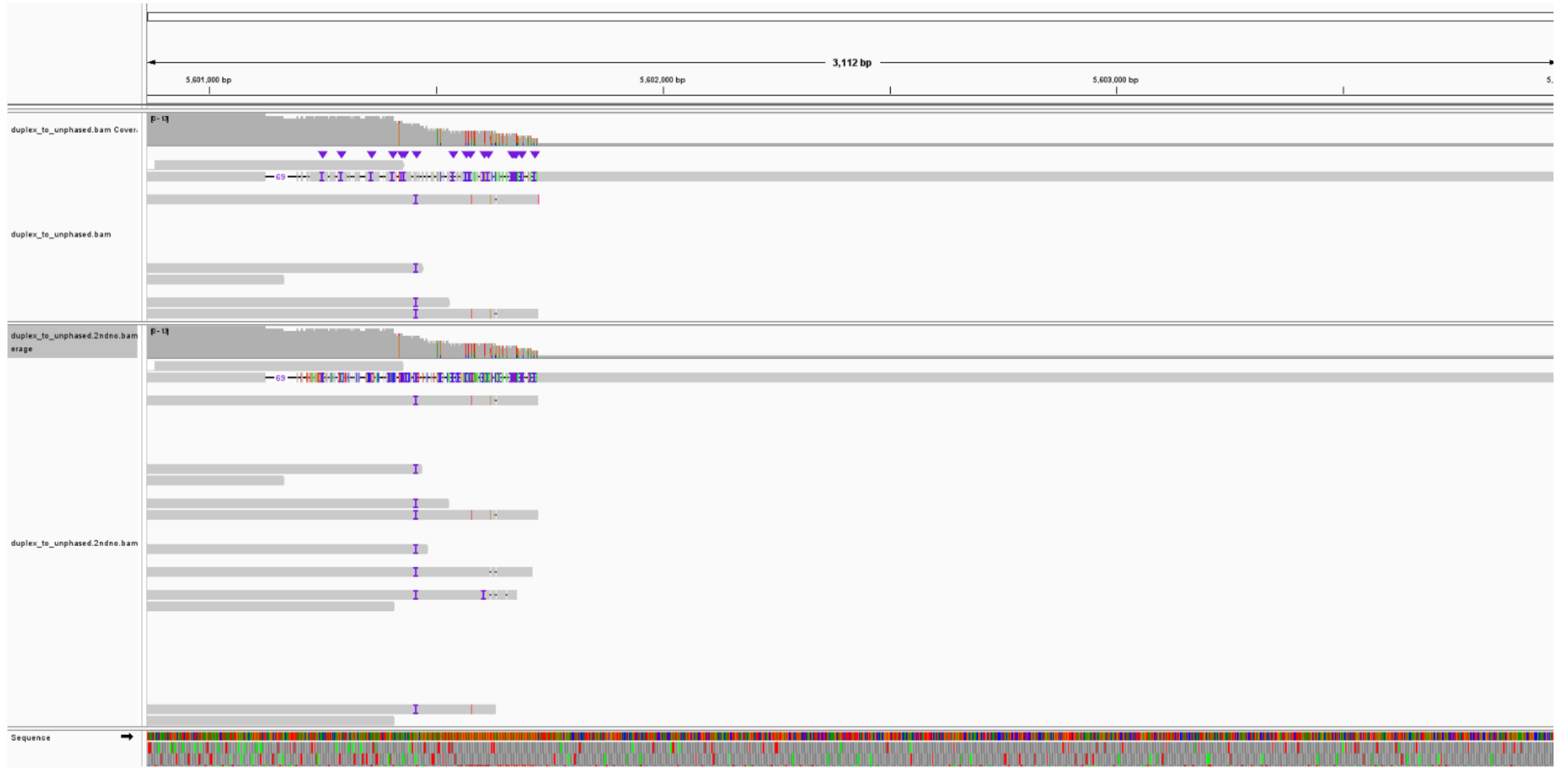
Article | [Open access](#) | Published: 14 September 2022

Telomere-to-telomere genome sequence of the model mould pathogen *Aspergillus fumigatus*

[Paul Bowyer](#) , [Andrew Currin](#), [Daniela Delneri](#)  & [Marcin G. Fraczek](#) 

[Nature Communications](#) **13**, Article number: 5394 (2022) | [Cite this article](#)

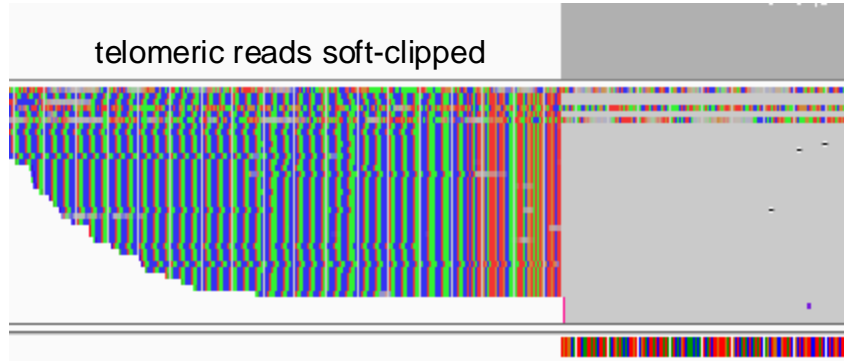
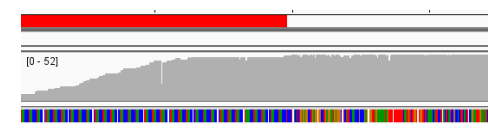
Trap: a problematic read causing extension beyond telomere



Actual telomere

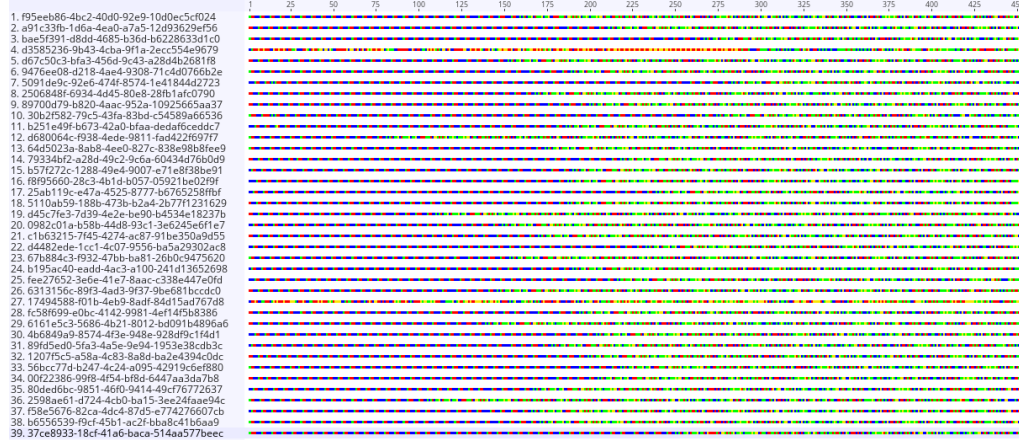
Where telomere counter wants to find telomere, but failed

Recover telomeres "hidden" in the reads



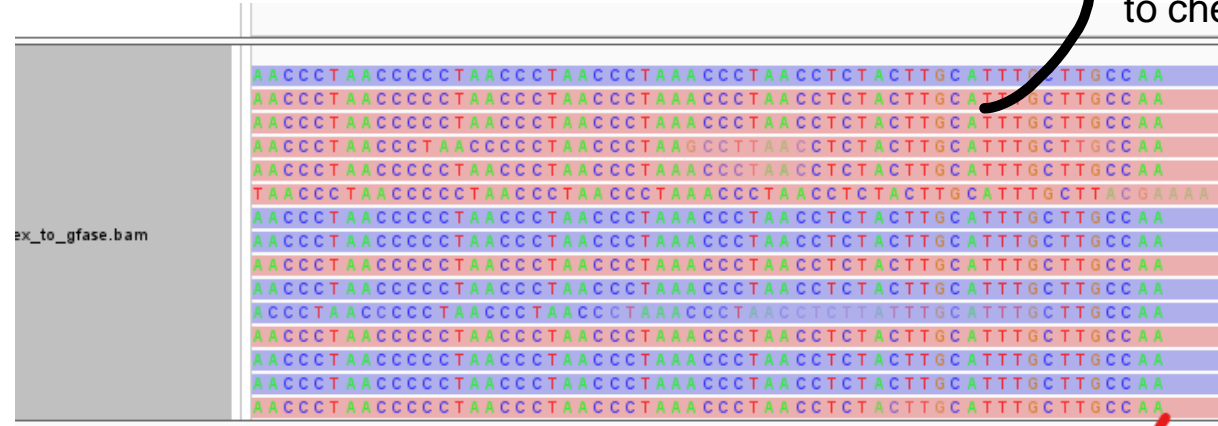
telomeric reads soft-clipped

extract soft-clipped alignments with telomeric motif "TTAGGG/CCCTAA"

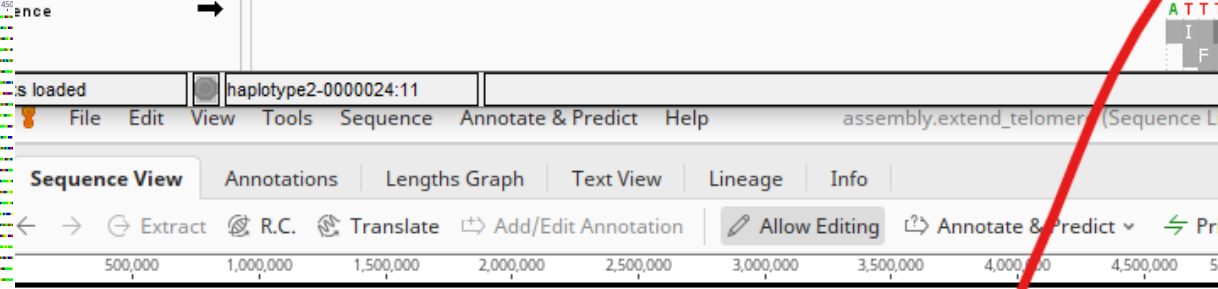


Local assembly (e.g. Flye)

telomere consensus (+ adjacent loci for anchoring)



map reads back again to check coverage



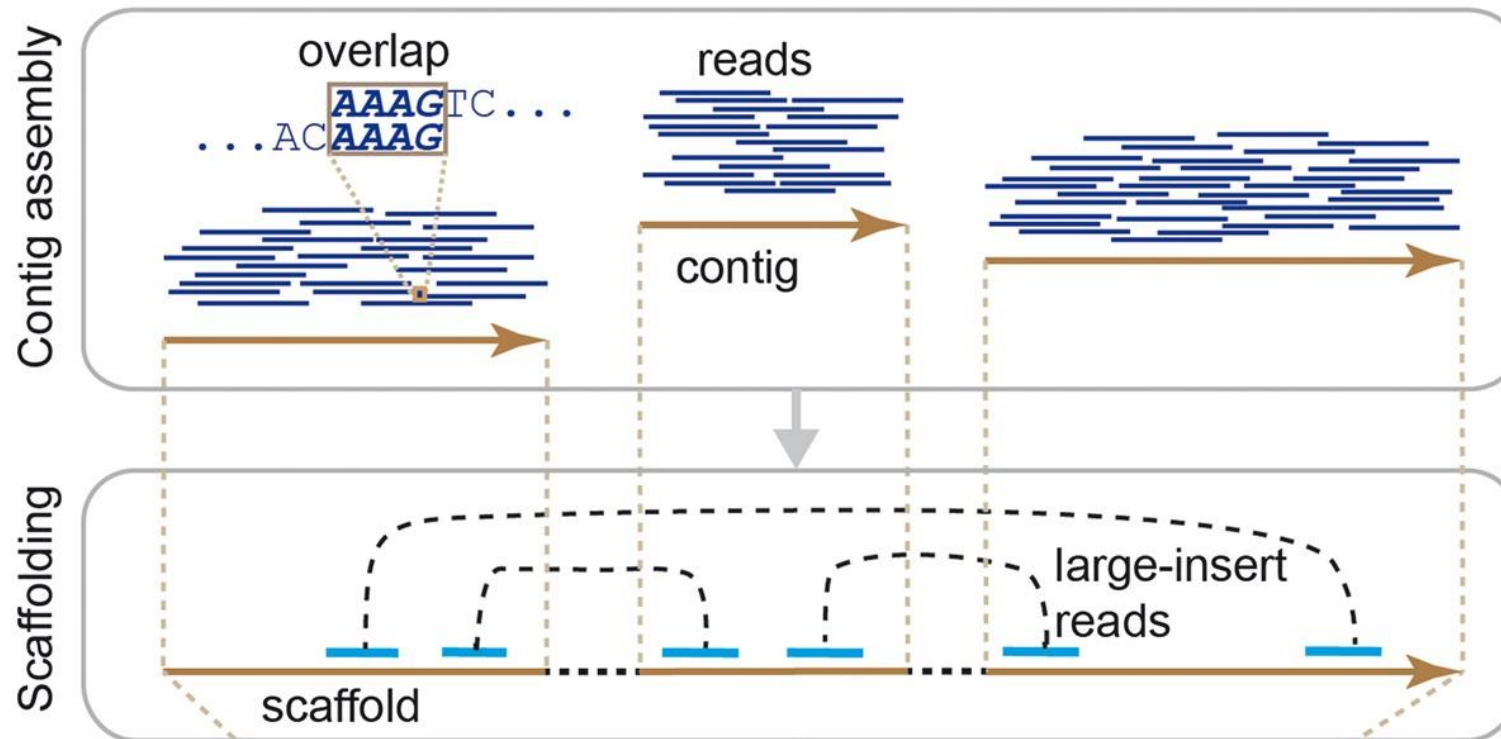
manually stitch telomeres back



Assembly gap

“unknown” nucleotides

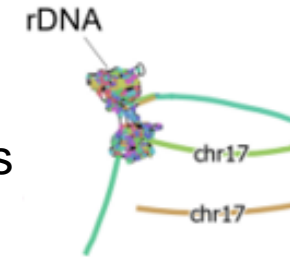
- Assembly gap (Ns) should be introduced after scaffolding (not covered in this workshop)
- But there are assemblers that might do some for you (e.g. verkko)



Gaps between contigs → filled with NNNNNN...

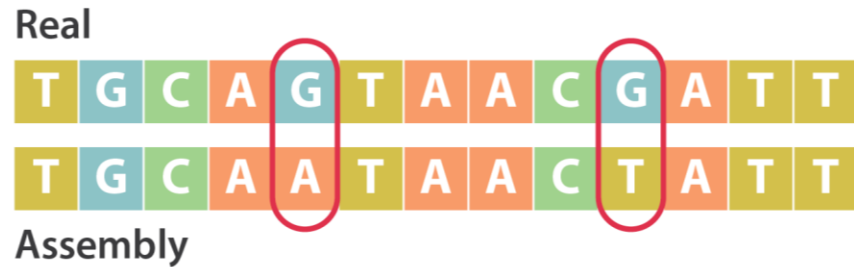
Assembly gap at ribosomal DNA array

Highly repetitive tandem repeats, identical copies



ATCCATCGCAATCGAATAGGNNNNNNNNNNNNNNNNNNNNCATTACCGACGATACACAG

Correctness



SNP/indels

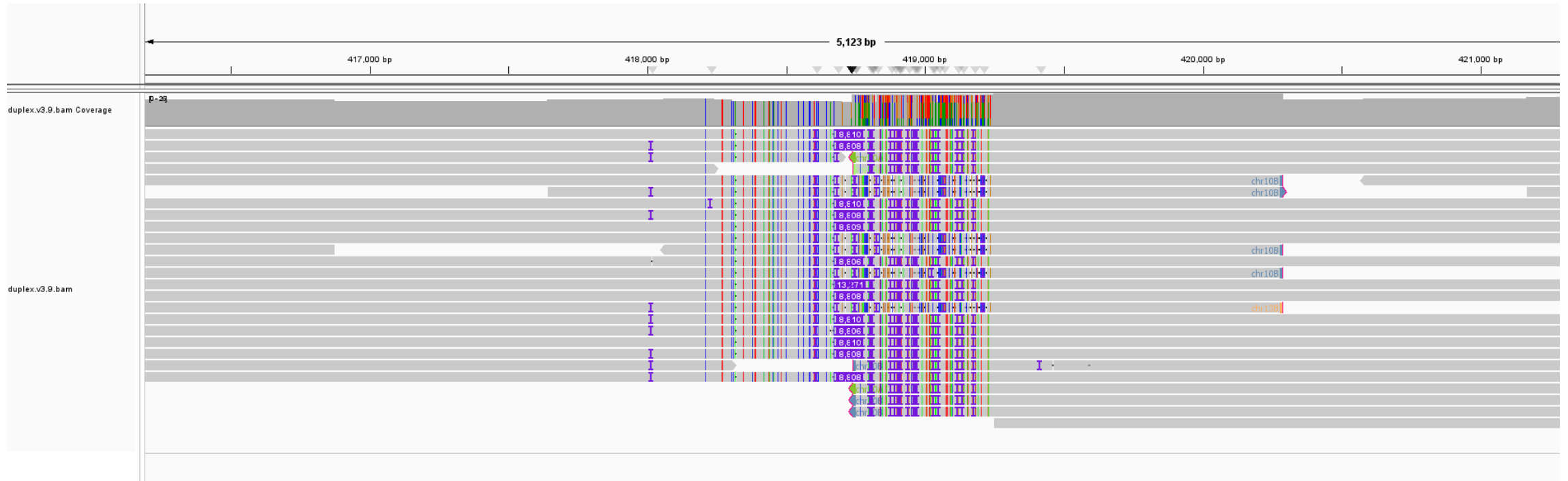
Misjoins

Mapping statistics

k-mer based consensus quality

SNP/indels

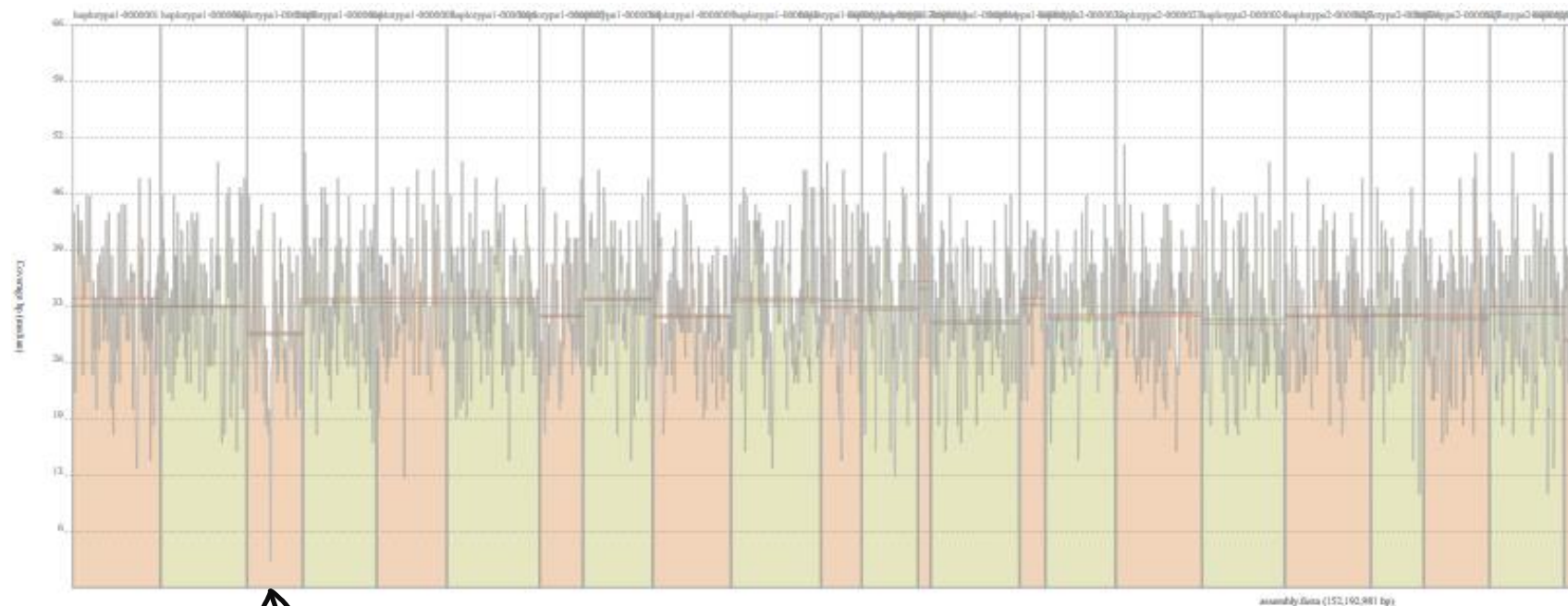
- Map reads back against the assembly
- Ideally should be completely identical
- Spot errors by calling SNP/indels



Structural errors

- Increase in coverage – unresolved repeats/duplications
- Drop in coverage – misjoins, misassembly

When ultralong reads come to rescue



Coverage drop in highly accurate ONT reads

Ultralong ONT reads filled it up!



homopolymer

[jvarkit](#)

WGSCoveragePlotter

last commit **august**

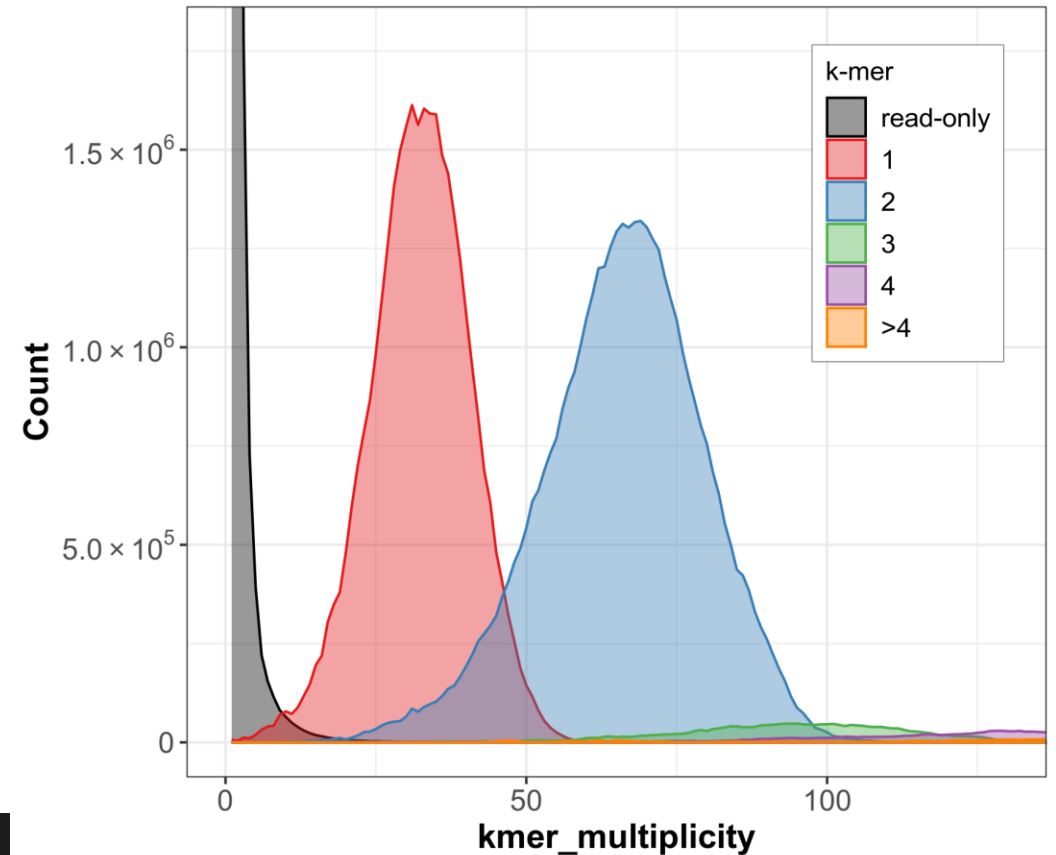
Whole genome coverage plotter

k-mer based assembly evaluation

Merqury

Find concordance between your read dataset and assembly

Reports k-merquality values (QV)



```
v3.9.chr.haplotype-paired      8597      152306908      57.3972      1.82086e-06
duplex_merqury.qv (END)
```

Equates to >99.999% accuracy

Summary

- Never assume your assembly is “error-free”
- Know how trustworthy your assembly is for any analysis
- **The “3C” rules**
 - Contiguity
 - Completeness
 - Correctness

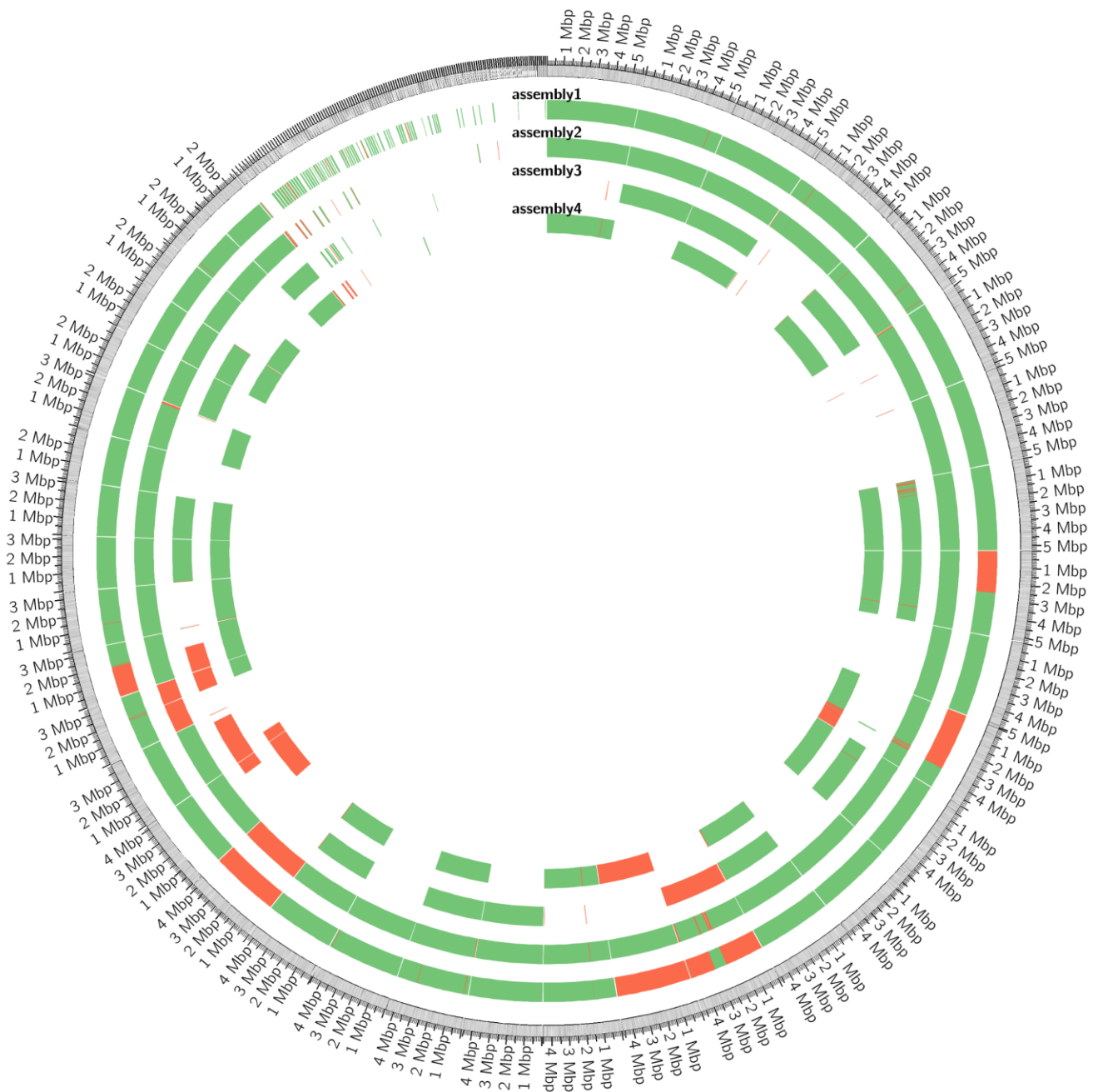
“Is my assembly ready for scaffolding?”

“Is my assembly ready for gene annotation?”

Wanna compare multiple assemblies and choose the best one?

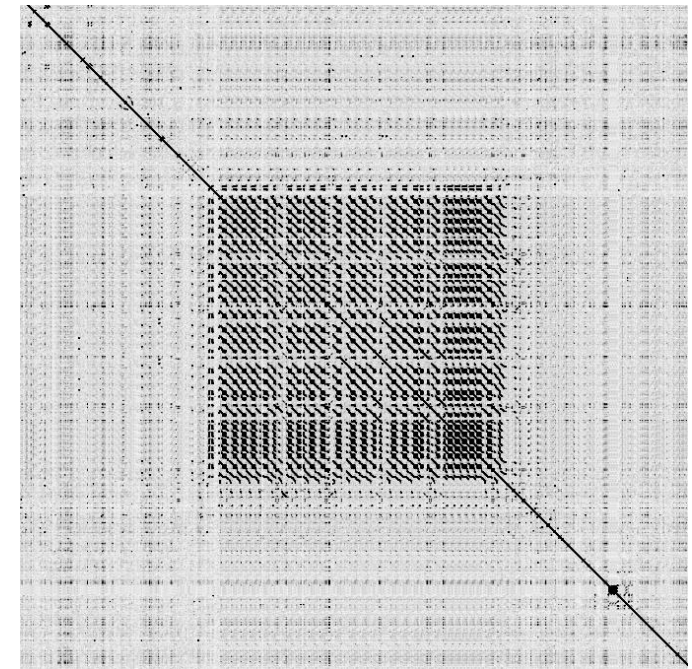
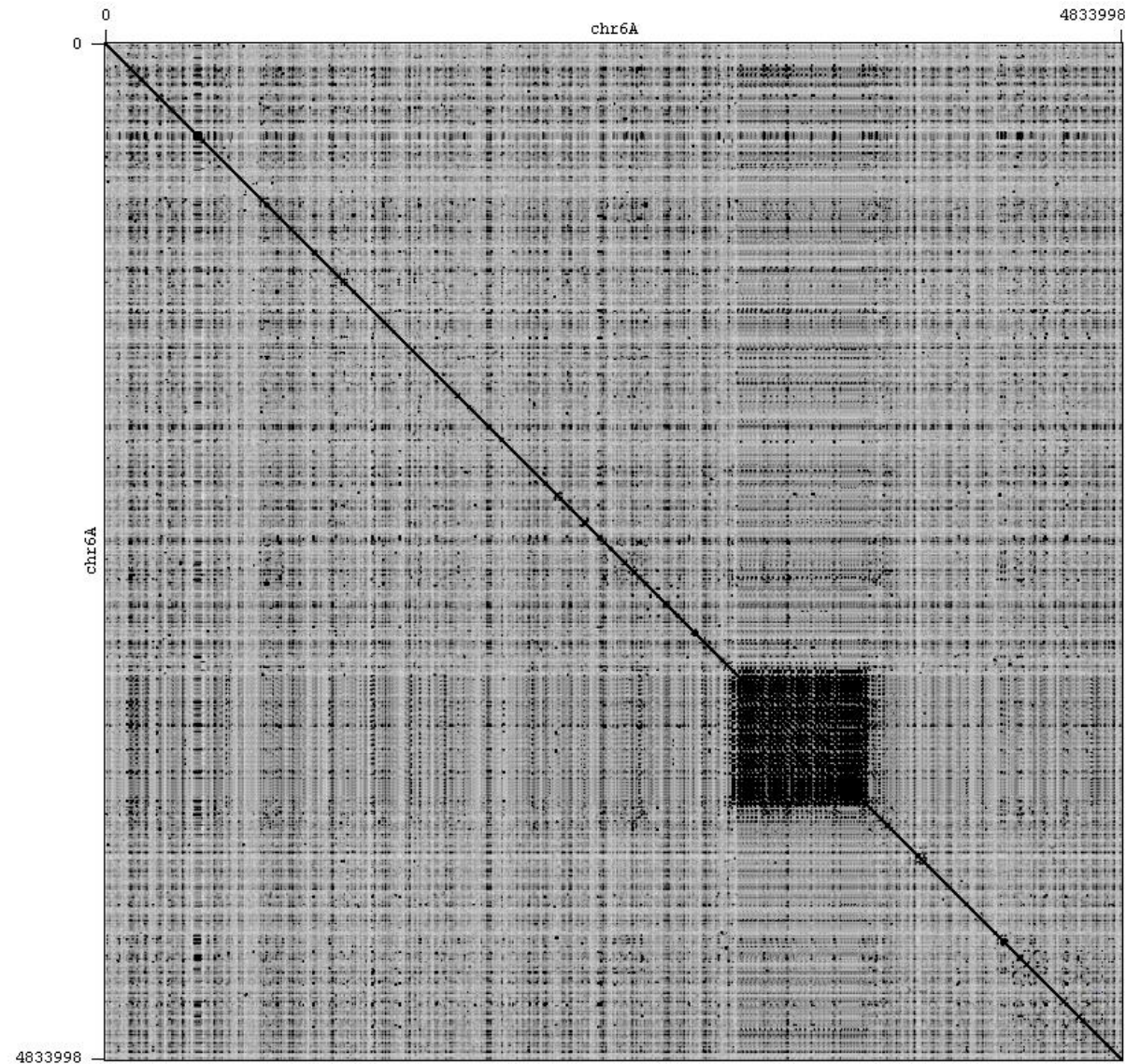
Try QUAST circo plot!

- Different read data types + assembler tool**
- 1: ONT reads + Verkko
 - 2: error-corrected ONT reads + Verkko
 - 3: ONT reads + Hifiasm (failed run)
 - 4: error-corrected ONT reads + Hifiasm (failed run)



Self alignment dotplot to find repetitive regions

Gepard, Chromeister



Last important tip

Keep track of your assembly versions as you curate it!

assembly_version_history.txt

```
raw: assembly.raw.fasta

v3.1: assembly.extend_telomere.mtDNA_lowcov_rm.fasta

v3.2: telomeres extended. telomeric motif copy number corrected for all
chromosomes.
v3.2.1: further telomere correction of v3.2 after coverage check.
extension stops when last supported by at least two reads (2 simplex, or 1
simplex 1 duplex.)

v3.3: after HiC scaffold. (v3.2.1.FINAL.fasta output from 3d-dna)
orientation swapped to make sure p arm starts first.

v3.4: replace chr5A with v2.6 manually scaffolded chr5A. telomeres
extended.

v3.5: fixed chr17B gap associated with a ~500bp GAAAA tandem repeat.

v3.5.1: fixed misassembly before the 5kbp rDNA-associated gap Ms. (Flye to
assemble HiC_scaffold_44,45,46 into a high quality local assembly. mapped
back to chr17B and replaced 89647bp before gap.)
```