

# **Strain-Level Comparative Genomics**

Phylogenomics Crash Course  
Multiple Genome Alignment & Visualization

---

STAMPS – Day 5

July 18, 2025

# Quick Phylogenetics Refresher

---

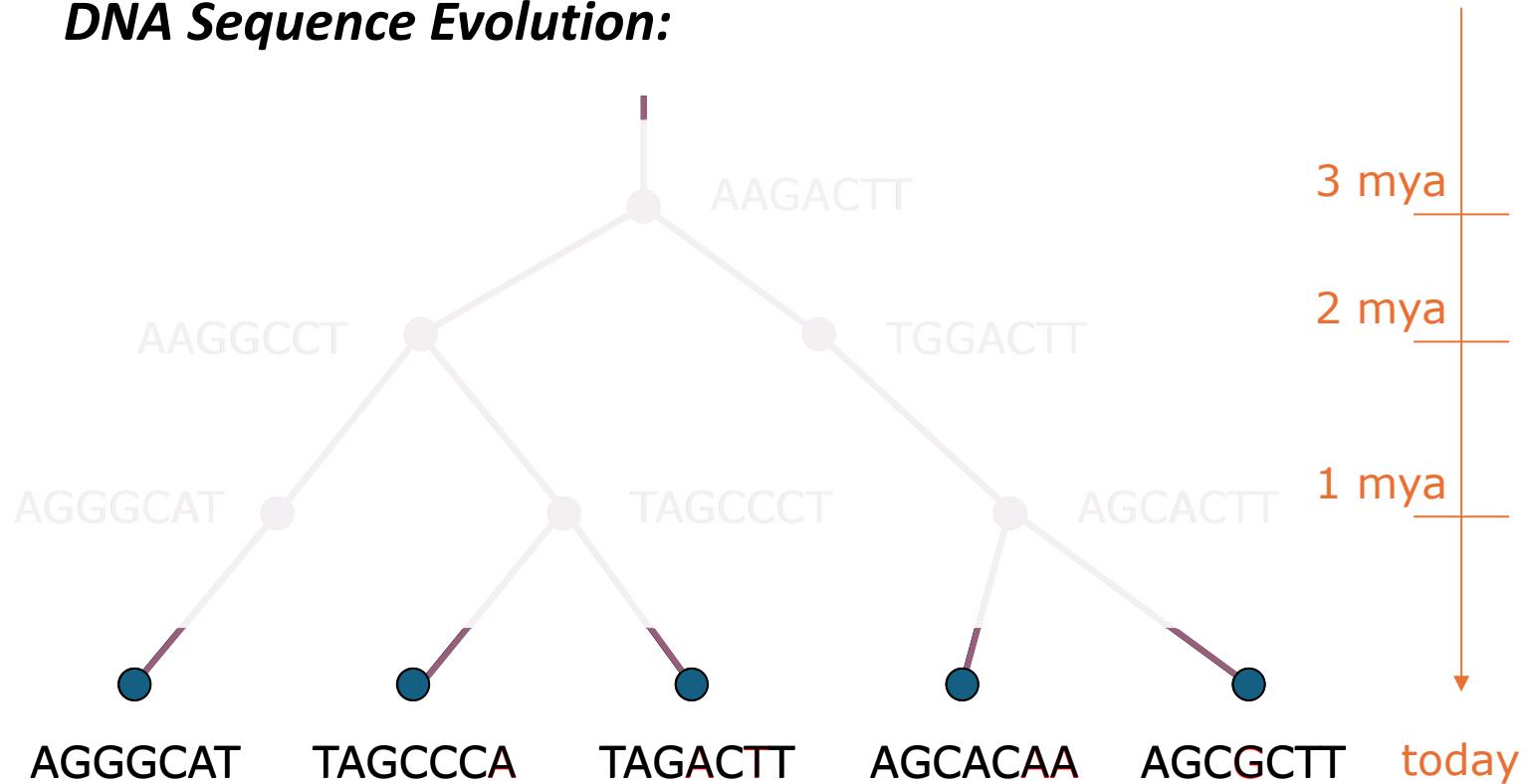
Multiple Sequence Alignment

Phylogeny Estimation

Population-level (“Species”) Tree Estimation

# Brief Intro to Molecular Phylogenetics

## DNA Sequence Evolution:



## Notes:

- insertions and deletions also occur randomly on branches (not shown)
- In statistical terms, the extant sequences are the **observed data**, while the tree shape, tree topology and mutation rates are the **unknown model parameters** to be estimated.

*Task is to estimate the tree shape from the extant sequences at the tree leaves...*

## Two Separate Problems (both NP-Hard):

1. Identify which groups of characters share a common ancestor. (Multiple Sequence Alignment)
2. Find the maximum likelihood tree and model parameters (ML Tree Estimation)

\*mya = million years ago

# Multiple Sequence Alignment: Definition & Goal

---

**Input:** Sequences from different organisms (or different loci) that evolved from a common ancestor.

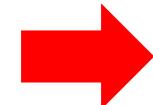
**Goal:** Align sequences so that all sets of positions having a common ancestor are grouped together.

- *Not the same as aligning short sequences (or reads) to a reference (“mapping”).*
- *Not the same as **genome** alignment*
- *Typically done before creating a phylogenetic tree...*

## Tools:

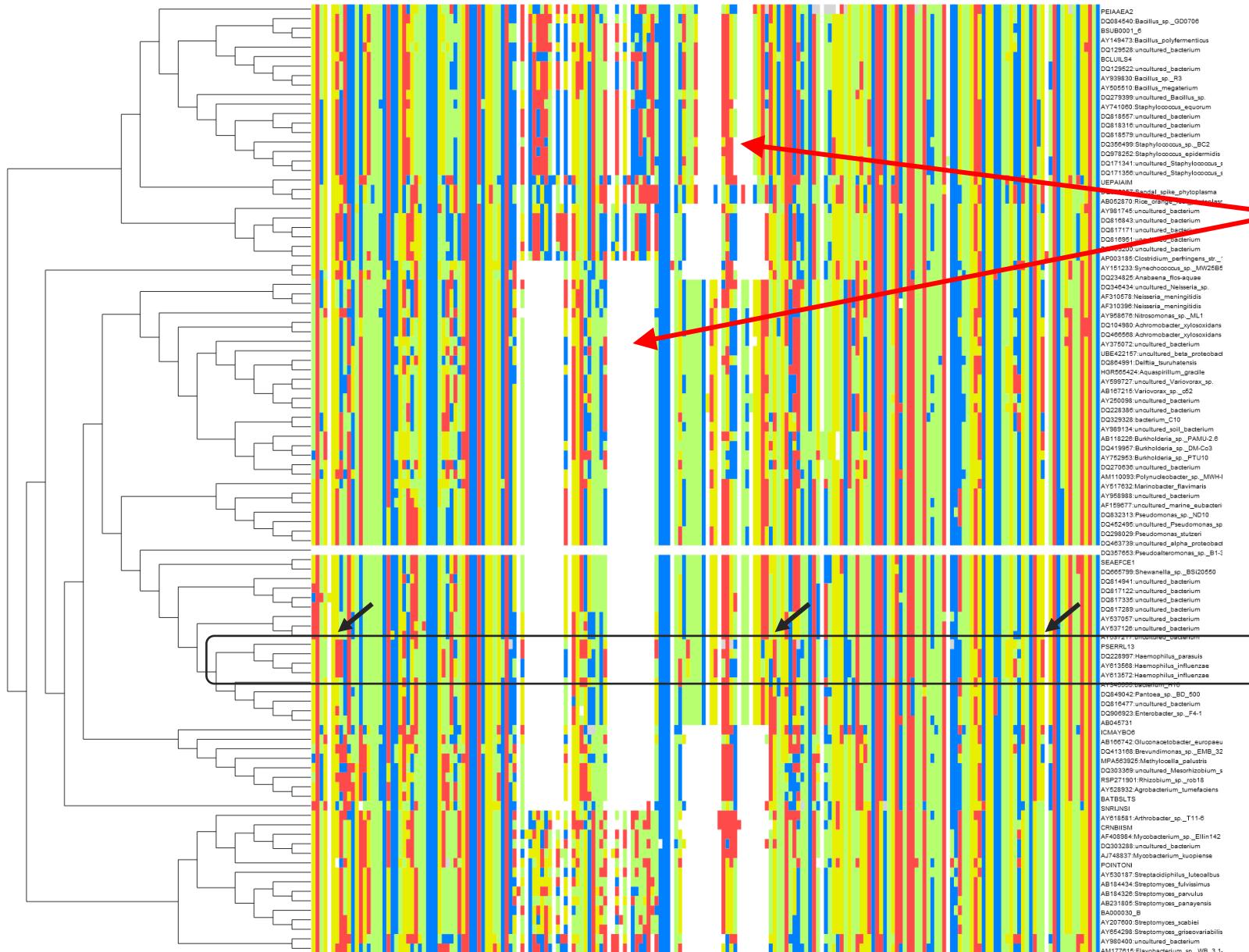
- **MAFFT**
  - Muscle
  - **PASTA**
  - ClustalW
  - DiAlign
  - BAli-Phy
- PRANK
  - T-COFFEE
  - ...et cetera

**Garbage  
Alignment**



**Garbage  
Tree**

# Multiple Sequence Alignment (Example)



- White indicates a “gap” in the alignment, a.k.a. an insertion or deletion (indel)

*Conserved mutation patterns indicate evolutionary closeness.*

*Phylogeny estimation  
algorithms use this to build a  
tree...*

# Data Properties Affecting Multiple Sequence Alignment

---

Very Approximate Order of Importance

- Avg Sequence Similarity (rate of evolution)
- # of Sequences
- Presence of highly conserved regions
- Sequence Length heterogeneity
- Gap length/frequency
- Sequence fragmentation (*not the same as heterogeneity*)
- Avg Sequence Length

# MSA Algorithms & Software (Partial List)

Tool	Use Case	Comments
MAFFT (Katoh et al., 2002)	Single Gene MSA (small $N$ )	<ul style="list-style-type: none"><li>• Uses patterns of insertion/deletion to find optimal alignment</li><li>• Generally pretty accurate in most conditions.</li></ul>
MUSCLE (Edgar, 2004)		<ul style="list-style-type: none"><li>• Progressive alignment. Suitable for relatively high overall sequence similarity.</li></ul>
CLUSTAL (Sievers et al., 2011)		<ul style="list-style-type: none"><li>• Ideal for protein alignments with structurally important sites.</li></ul>
PASTA (Mirarab et al., 2015)	Single Gene MSA (large $N$ )	<ul style="list-style-type: none"><li>• Divide-and-conquer algorithm. Ideal for scaling alignment to large number of sequences (<math>&gt;1000</math>)</li></ul>
HMMER (Eddy, 1998)	Query sequence alignment to reference	<ul style="list-style-type: none"><li>• Represents reference alignment as HMM. Query sequence alignment performed using standard HMM algorithms.</li></ul>

I tend to tell people:

- Just run MAFFT for anything less than 500 sequences or so
- Muscle is fine too if the divergence is low...
- ...or ClustalW for AA sequences with important structural sites
- Use PASTA for over 1k sequences or if avg. %-identity is very low (high rate of evolution).

This simplistic advice is a STAMPS 2022 exclusive...

## Final MSA Points (details in appendix)

---

- MSA accuracy can drop to 0% *fast*
  - *Especially as diversity & # sequences go up.*
- Bad MSA will lead to a Star-like tree
  - *i.e. no discernable relationships between sequences*
- MSA failure can take many forms (over/under alignment)

# Tree Estimation Mechanics

---

- Given the MSA, how to estimate the Tree?
- Statistical Estimation:
  - Tree shape & branch length are ***parameters*** under a model of sequence evolution
  - Called GTR (generalized time reversible)
  - Each site (letter, residue, position, etc...) evolves I.I.D.
  - Mutation along a branch happens according to a continuous-time Markov process (“time” here = branch length).
  - Root is ***not identifiable***.
  - No indels under the model, MSA gaps treated as missing data.
- Under this model
  - Every tree shape/topology has a likelihood given the sequences (the “data”)
  - **Maximum Likelihood Tree** will be “statistically consistent” (good)
    - I.e. given enough sites, the ML tree will be the correct one with  $p \rightarrow 1$
    - Neighbor-joining, UPGMA, etc... are NOT statistically consistent (bad)
- Hence most good Tree estimators use an ML model (e.g. RAxML, FastTree, IQTree)

# MGA & Visualization for Strain Analysis

---

Digging deep...

# Introduction

---

- What does “Strain” mean?
  - Particular SNP?
  - Multiple particular SNPs?
  - Presence/Absence of certain genes?
  - Phenotype?
  - *Other?*
- How do we compare strains?
  - Multiple Genome Alignment
  - Pangenome Analysis

# Whole-Genome Alignment

---

- Idea: align specifically the *shared* (“core”) portion of several genomes.
- Use these aligned segments to identify phylogenetic relationships, etc...
- Visualize what exactly is similar and different...

## Tools:

- Parsnp
- Mauve
- SibeliaZ
- (others...)

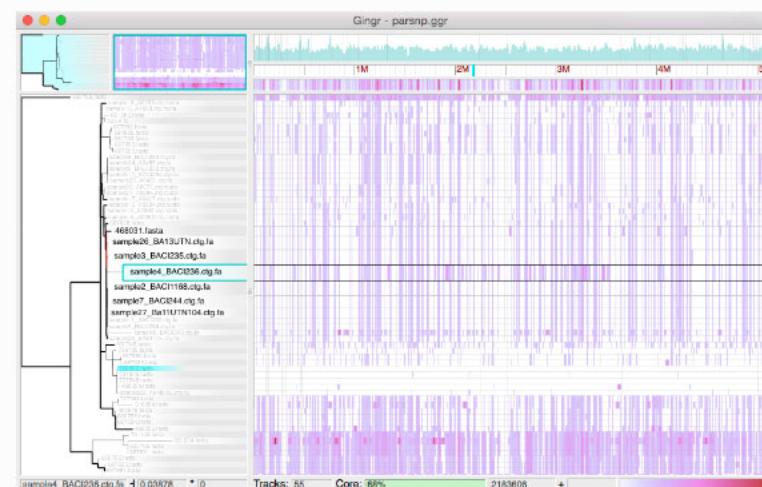
[Docs](#) » Harvest

[Edit on GitHub](#)

## Harvest



Harvest is a suite of core-genome alignment and visualization tools for quickly analyzing thousands of intraspecific microbial genomes, including variant calls, recombination detection, and phylogenetic trees.



# Whole Genome Alignment: Quick How-To with Parsnp

- Get *assembled* genomes from individual organisms
  - Isolates are nice, MAGs will do
  - Contigs are fine for this, doesn't have to be complete
  - Helps to have at least 1 high-quality, annotated reference genome
  - Useful to run QUAST to QC the assembly
- Run Parsnp:

```
contig_repo=./parsnp_contigs  
parsnp_out=./parsnp_output_13  
ref_genbank=./ref_assembly_GCF_008121495/Ref_ATCC_29149.gbff
```

```
parsnp -g $ref_genbank -d $contig_repo -p 15 -o $parsnp_out
```

Annotated Reference  
Genome (.gbff format)

Folder with 1 fasta file for each  
assembly (containing all  
contigs) ...OR...  
File with a newline-separated  
list of assembly fasta files (full  
paths)

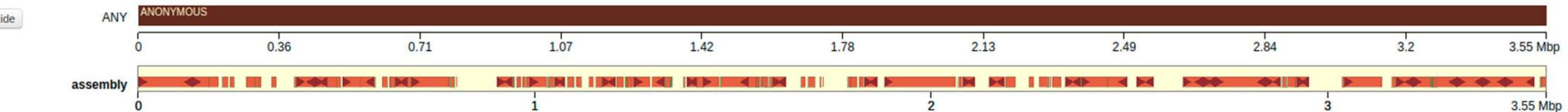
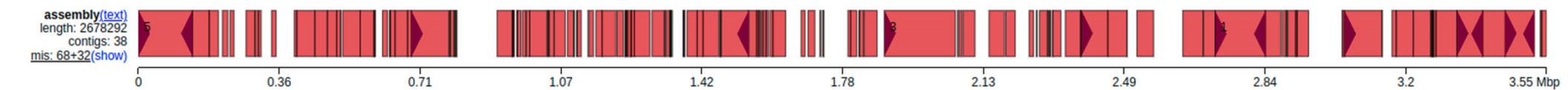
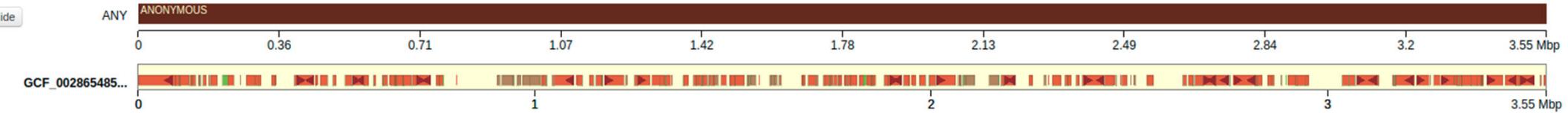
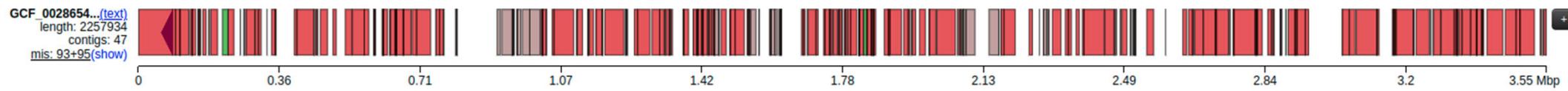
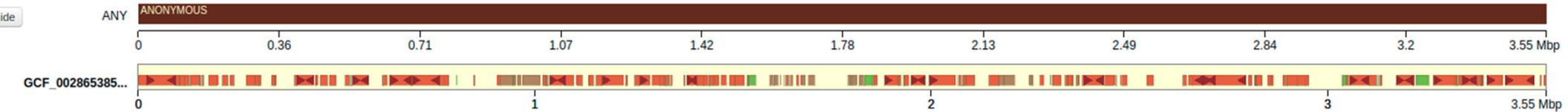
# processors

Output folder

## What can we learn?

- Assembly Quality Issues?
- Issues with Reference?

# Interlude: QC-ing an Assembly with QUAST



**Notes:**  
The top two assemblies are SPAdes assemblies done by the original authors of the R. Gnavus paper (citation later).  
The bottom is a Unicycler assembly from the same reads.

# Case Study #1: *Klebsiella pneumoniae*<sup>1</sup>

---

- 119 Carbapenam-resistant *Kp* isolates (95 from Houston Hospitals)
- Hybrid short/long-read assemblies
  - so they *should* be high-quality
- Focus of study was to show spread of two separate “clonal groups”

# Case Study #1: Observations...

## Observation #1:

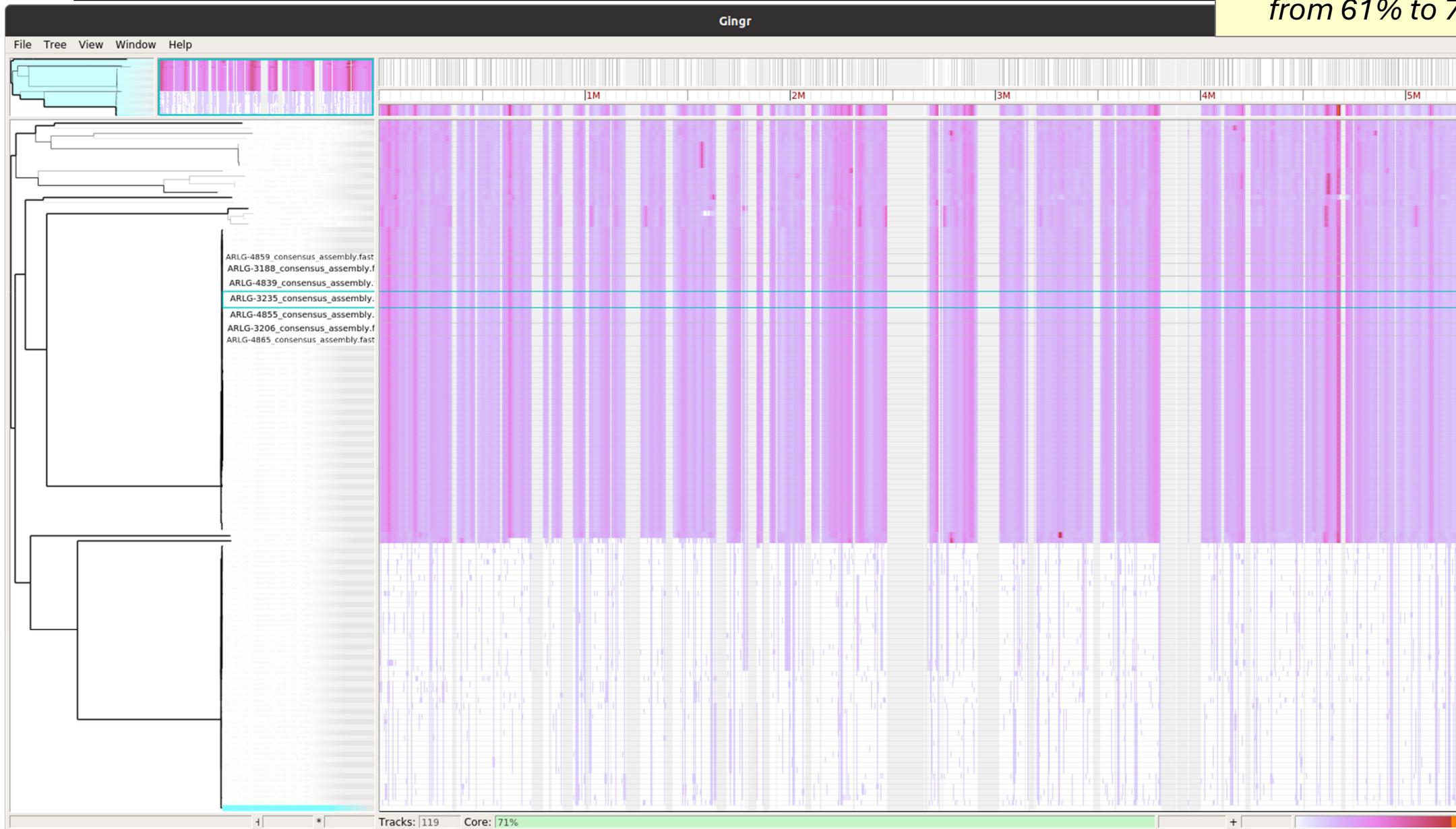
One of these genomes (ARLG-8054) is **MUCH** different from all the others.



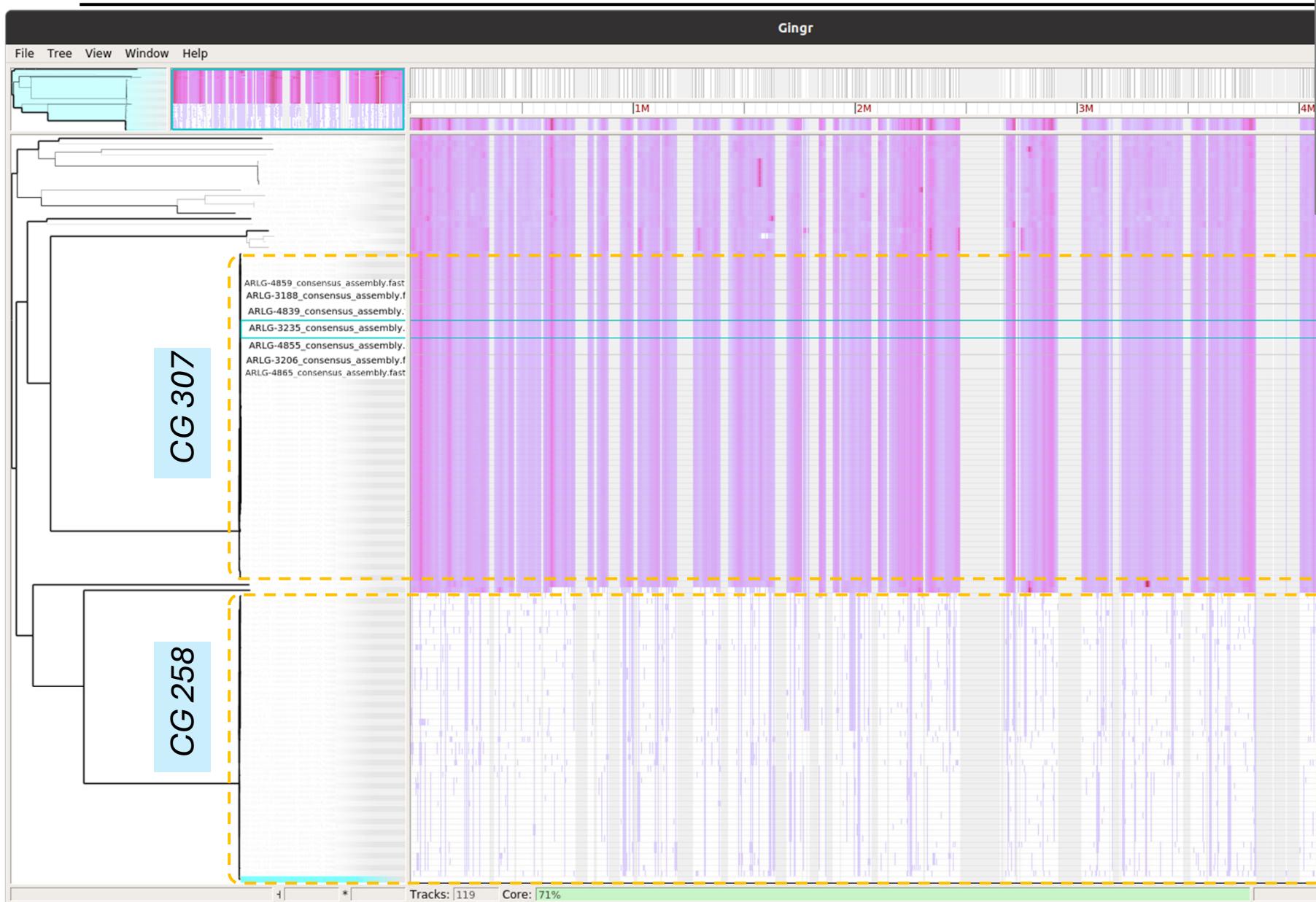
# Case Study #1: Excluding ARLG-8054

## Observation #2:

- Without 8054, core % goes from 61% to 71%



# Case Study #1: Excluding ARLG-8054



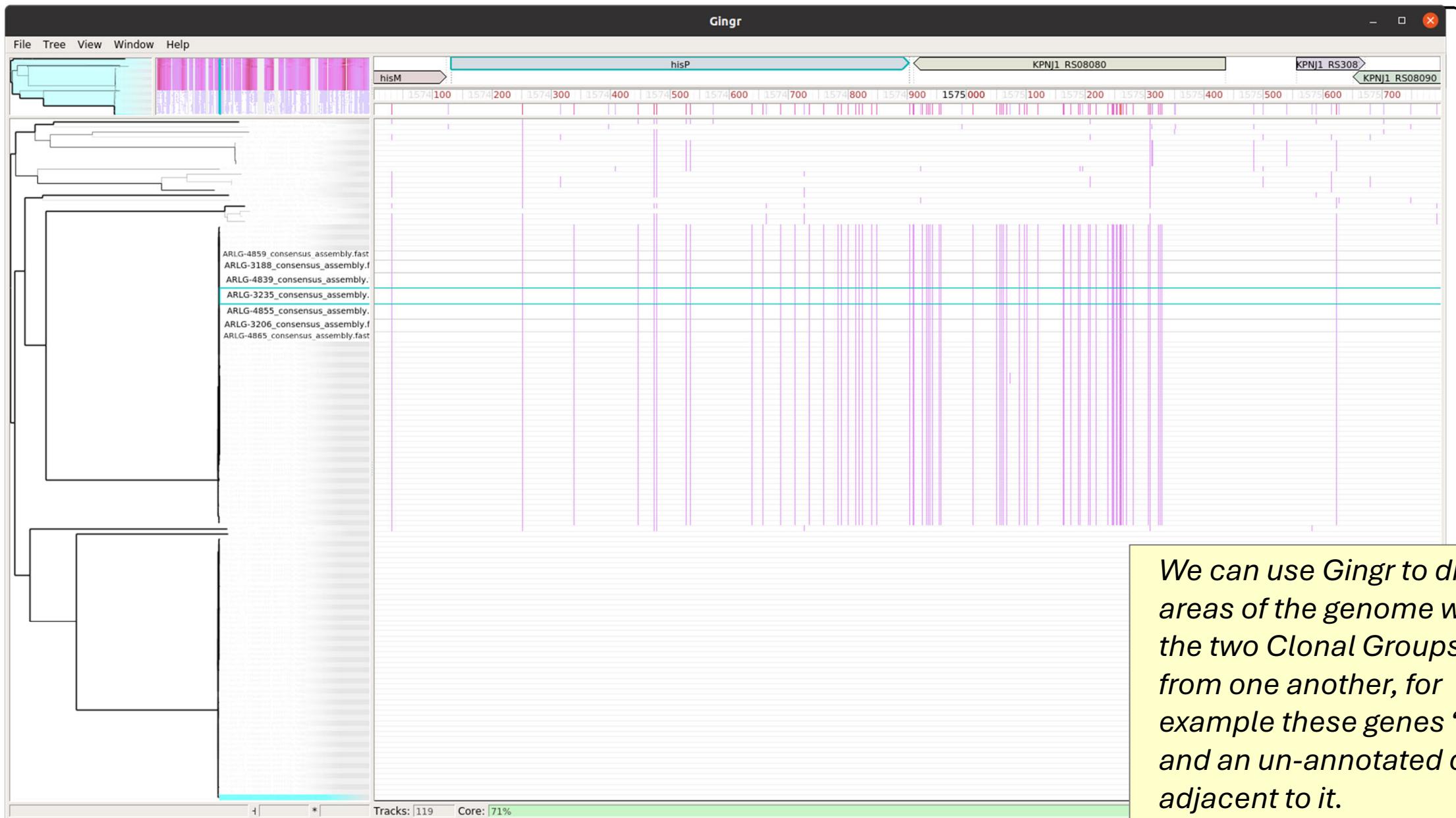
## Observation #3:

- 2 Clades with highest frequency, corresponding to clonal groups 307, 258 which were major targets of investigation in this paper.

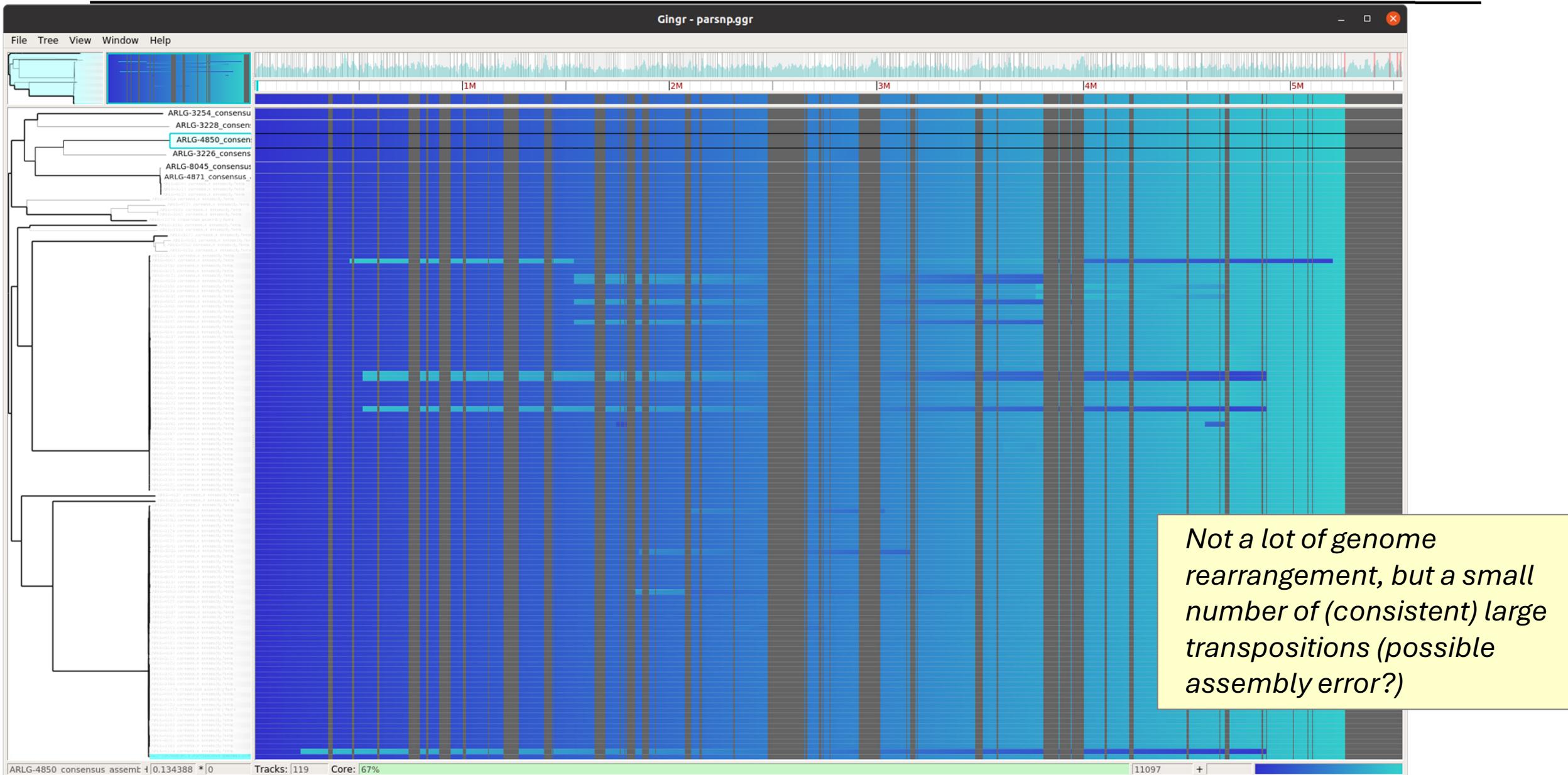
## Notes:

The reference genome here is GCF\_000598005.1, described as "strain=30660/NJST258\_1", so ostensibly ST258.

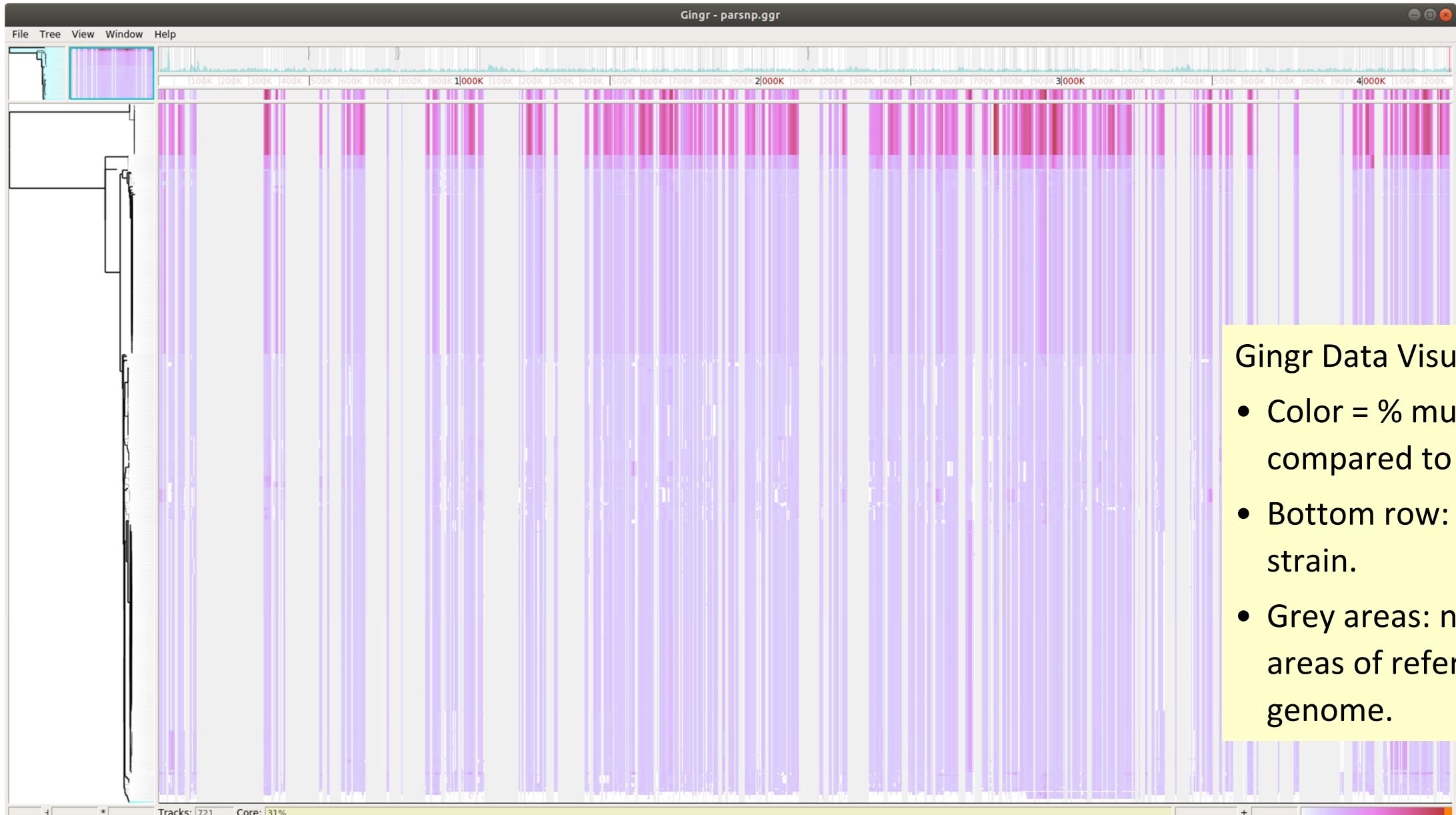
# Case Study #1: Digging In...



# Case Study #1: Synteny Comparison



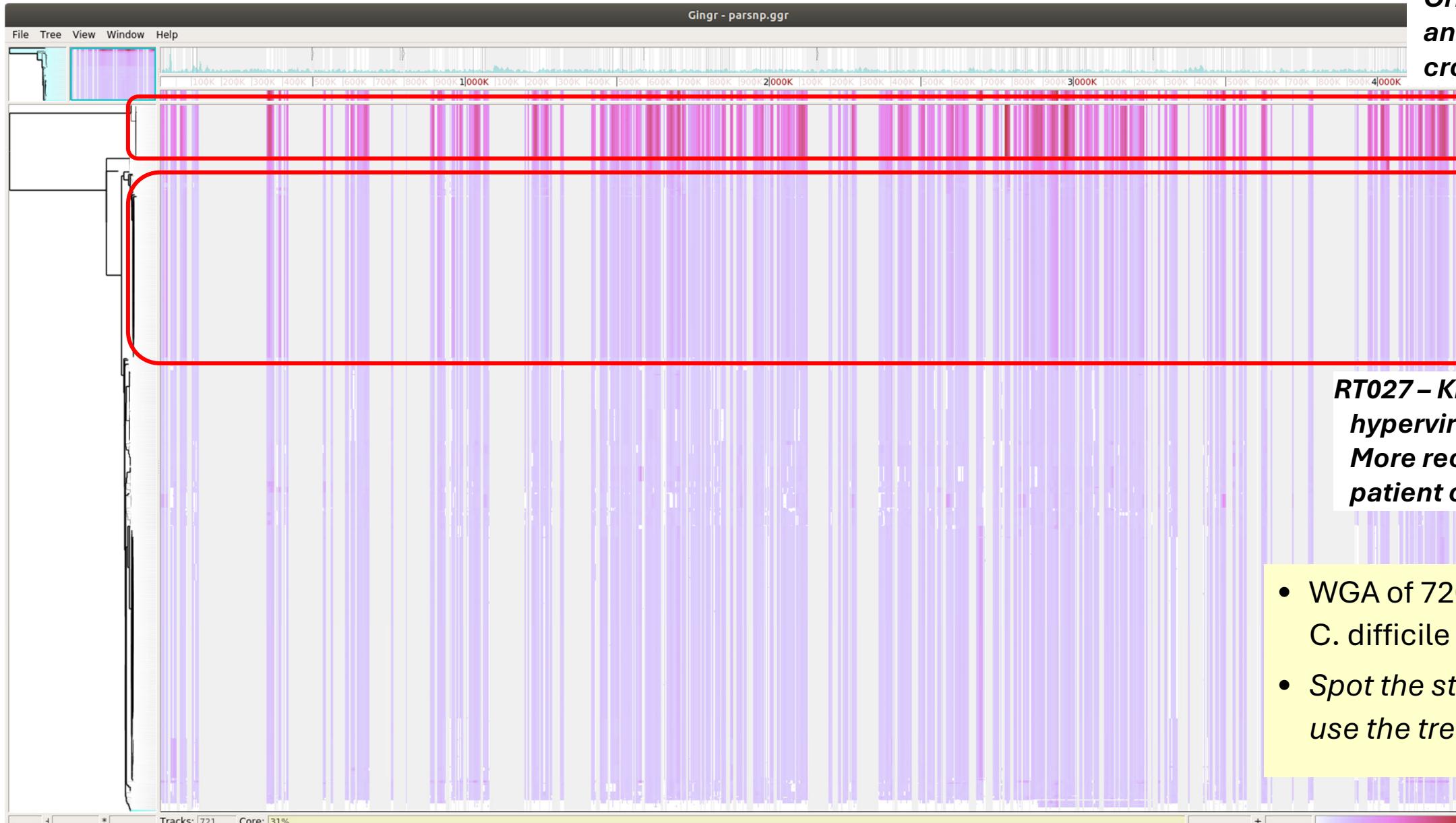
# Case-Study #2: *C. difficile* Genomes



Gingr Data Visualization:

- Color = % mutation compared to reference
- Bottom row: reference strain.
- Grey areas: non “core” areas of reference genome.

# Case-Study #2: *C. difficile* Genomes

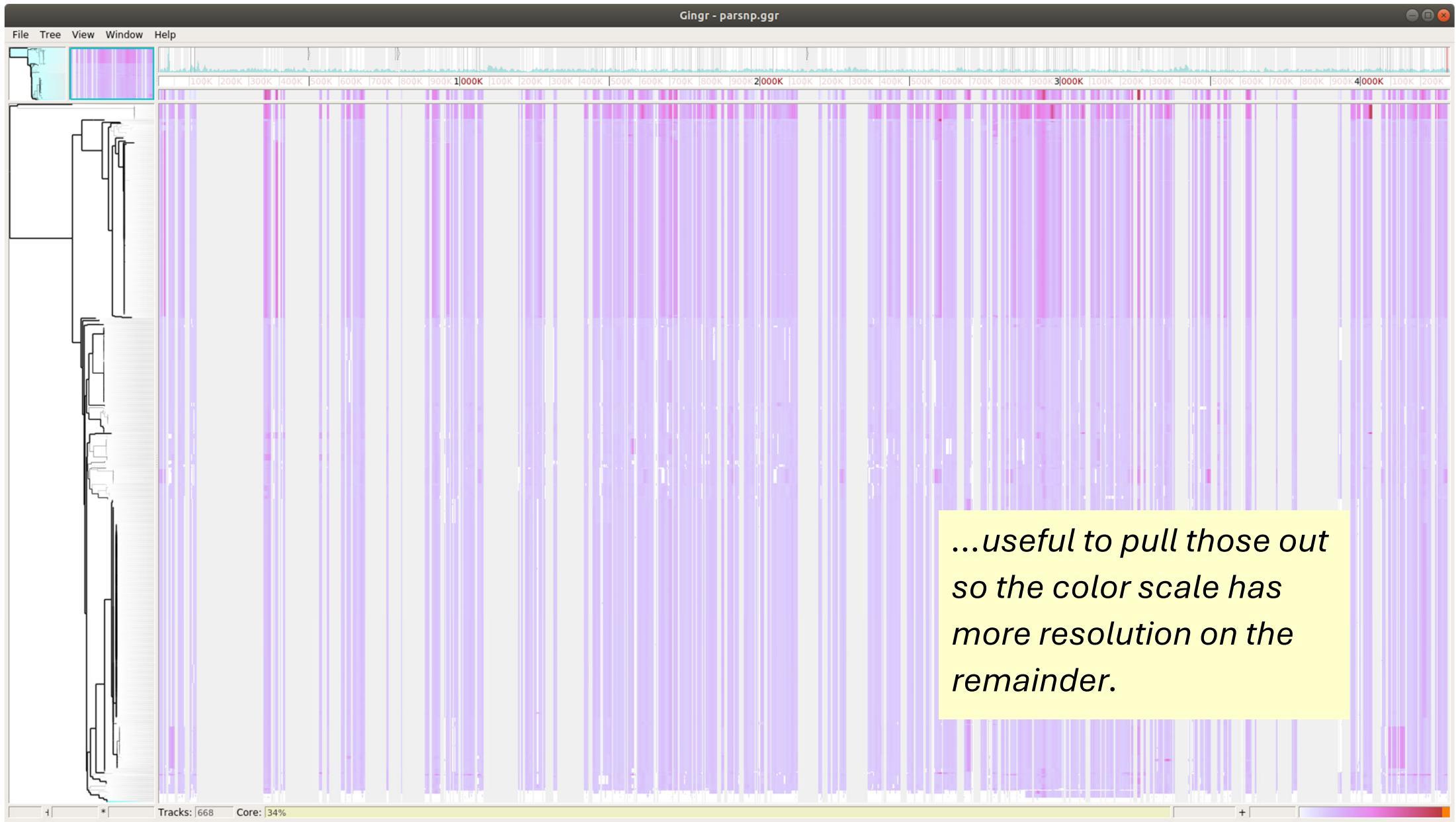


- WGA of 720 assembled *C. difficile* genomes
- Spot the strains... (hint: use the tree)

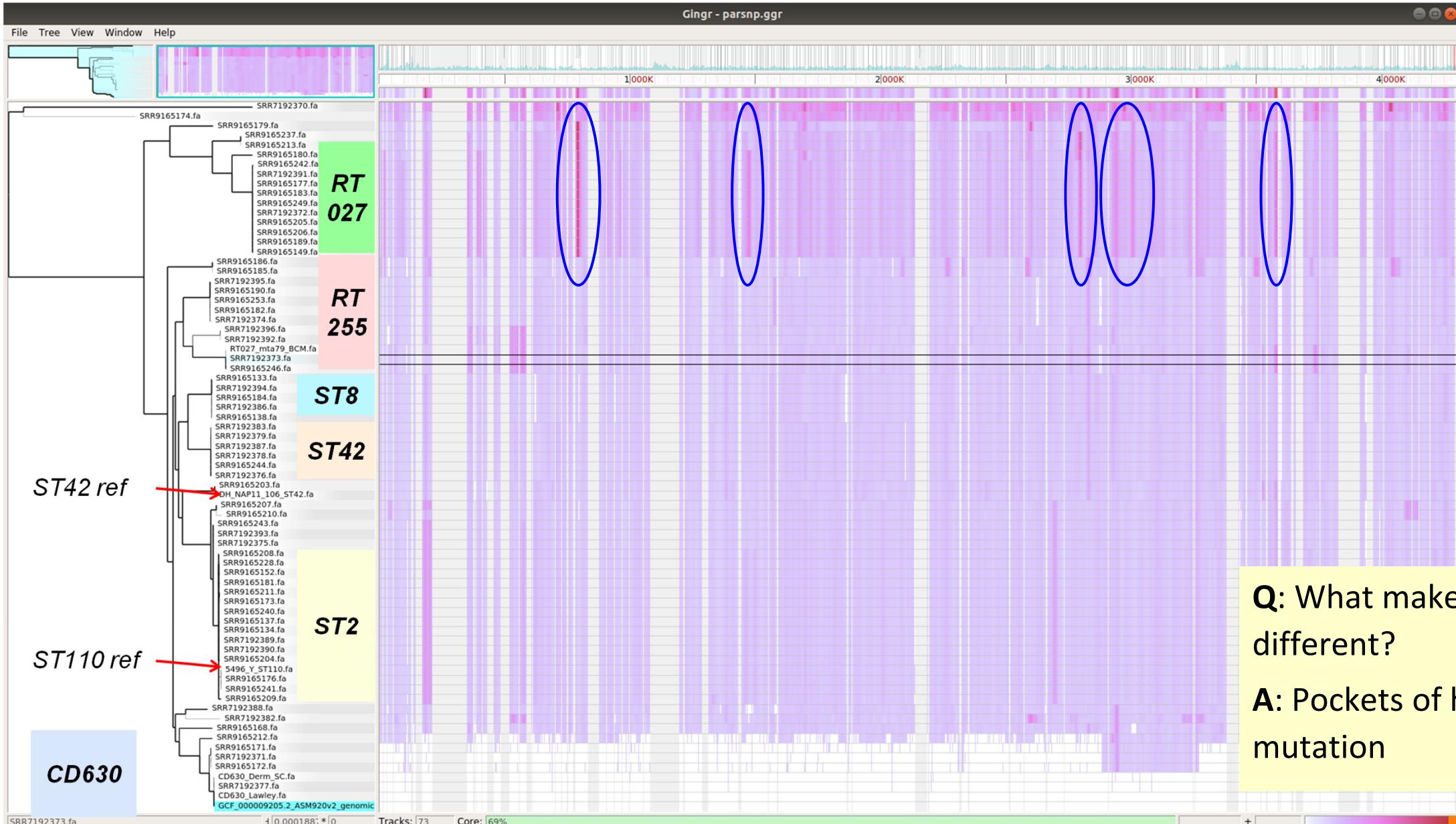
**RT078 –**  
*Originated in animal host, crossed over*

**RT027 – Known hypervirulent strain. More recurrent, nastier patient outcomes.**

# Case-Study #2: *C. difficile* Genomes (excluding RT078 samples)



# Subset of Genomes w/ST annotation

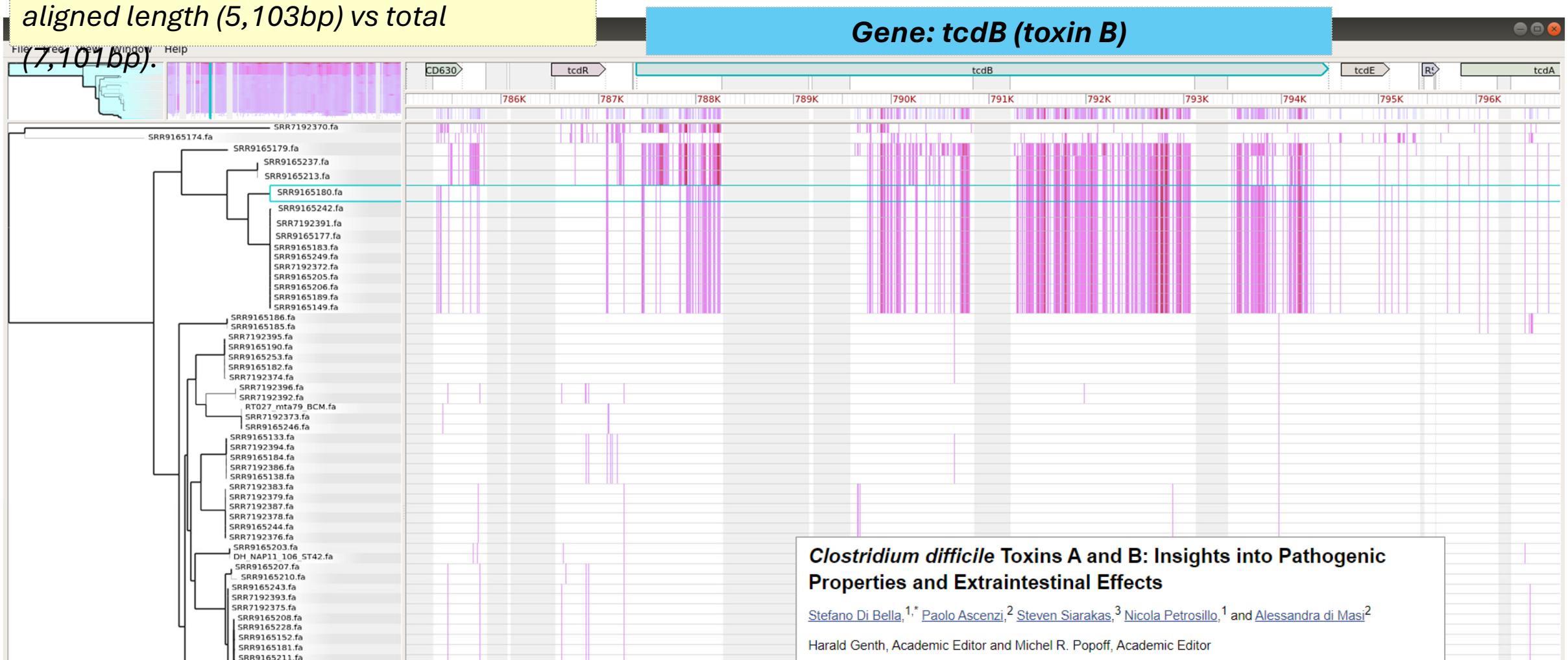


**Q:** What makes RT027 different?  
**A:** Pockets of heavy mutation

# Digging Deeper...

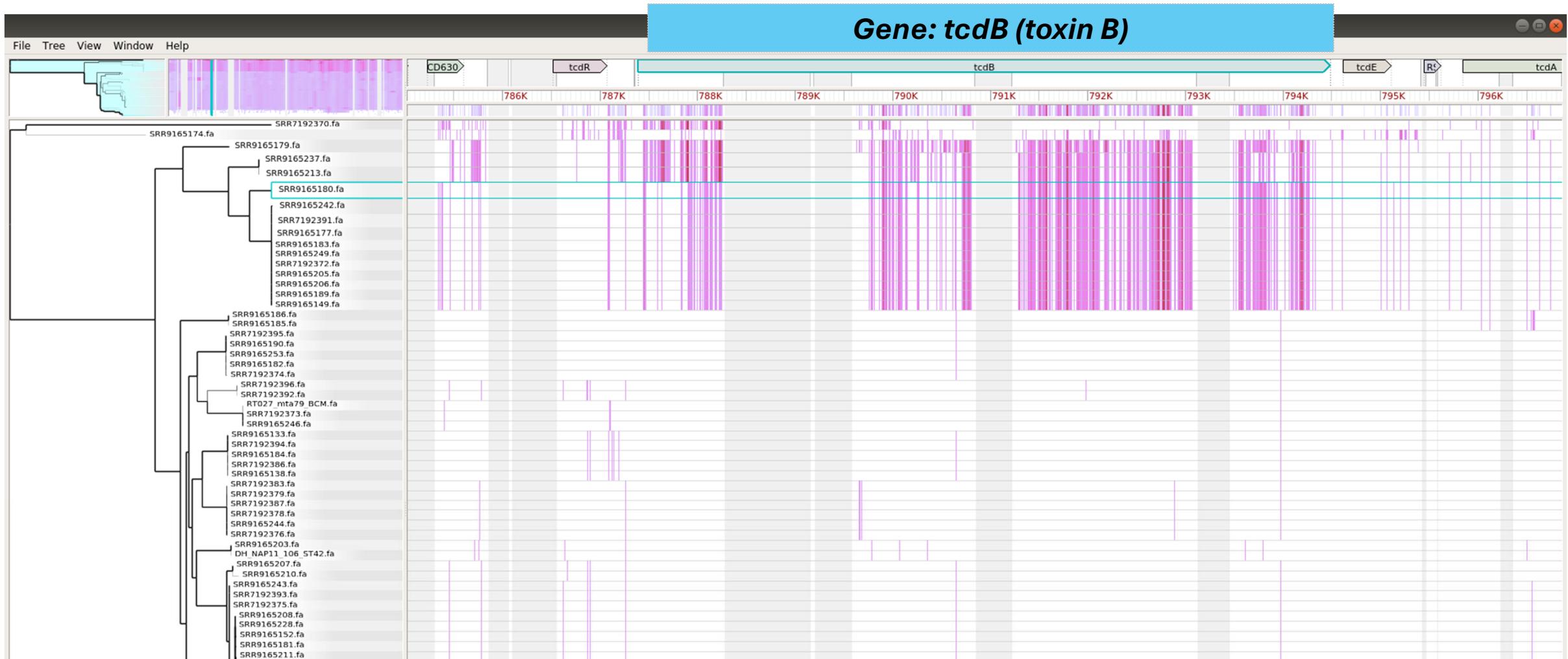
**Note:** not all of the *tcdB* gene was aligned by Parsnp, so this table represents the aligned length (5,103bp) vs total (7,101bp).

- This particular region is precisely the coding locus for Toxin B.
- RT027 carries a variant *tcdB* gene with altered function that contributes to its virulence.



# Remark...

Gene-level phylogenetic signal largely matches up with whole-genome phylogeny...

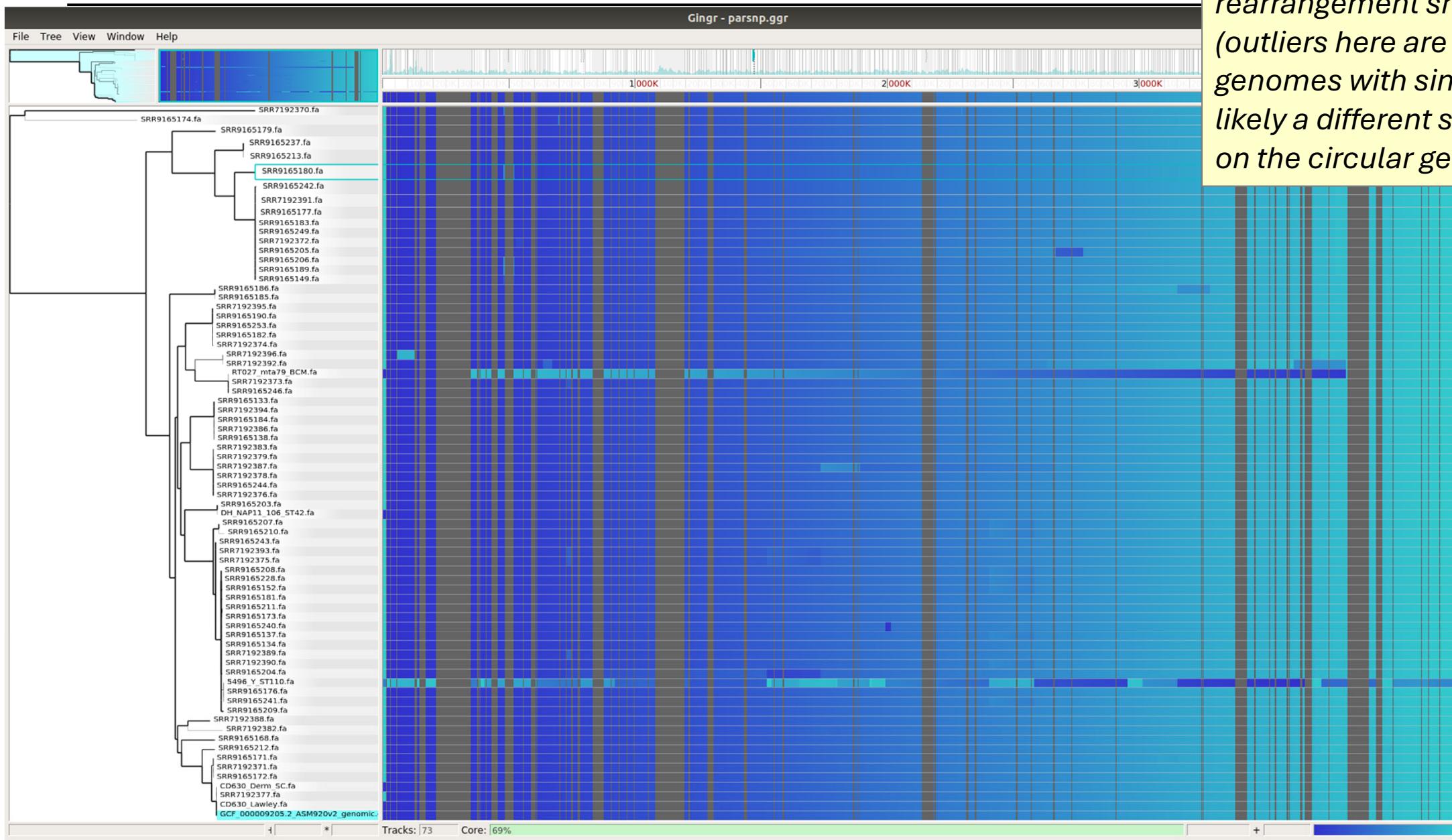


# Comparing Reference Genomes for Some Strains

**Note:** RT027 is in the top row. CD630 is a lab strain used as a common reference.



# Synteny Comparison: *C. difficile* Isolates



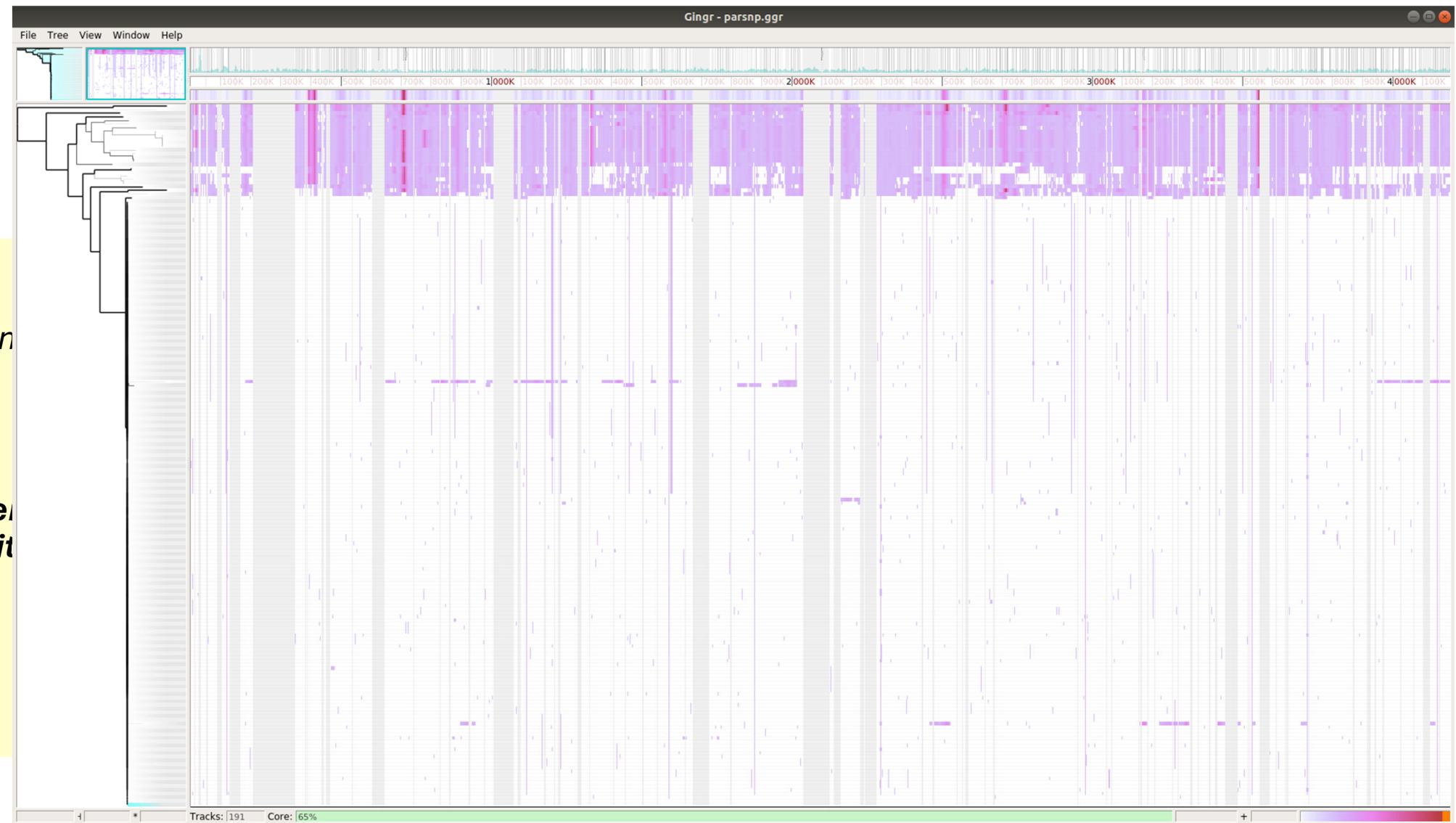
For *C. diff*, even across a huge number of isolates, very little rearrangement shows up (outliers here are reference genomes with single contig, likely a different starting point on the circular genome.)

# Alignment of RT027 isolates (and near relatives) to RT027 ref.

**Q:** Does the RT027 Reference match the genomes from the clin

**A:** ...Yes

- **Very little to see, very high match level with all RT027 isolates except 3.**

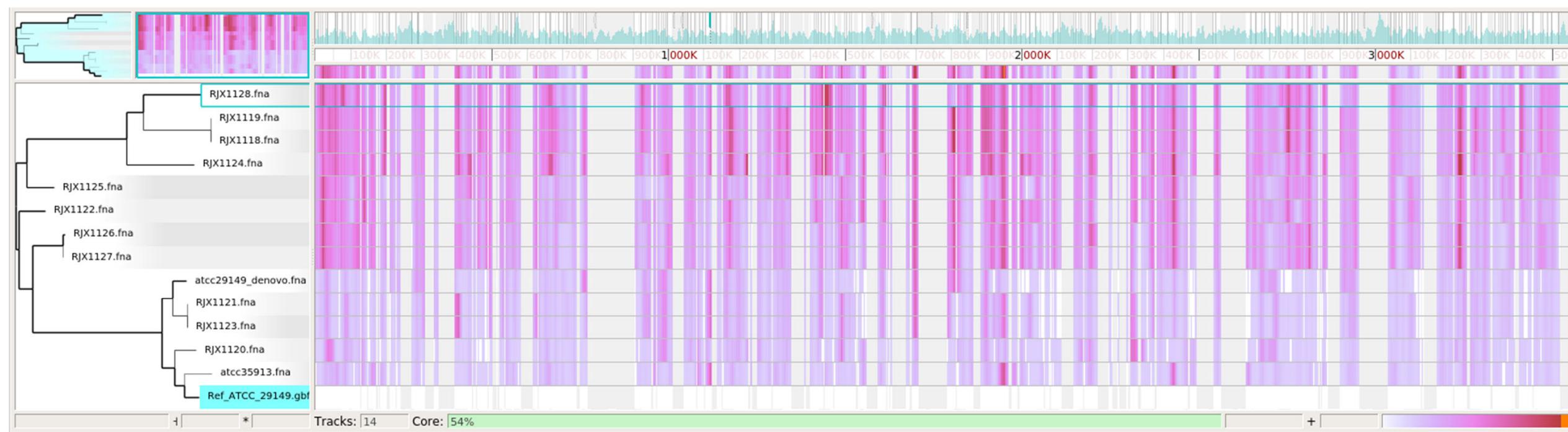


# Case Study #3: *R. gnavus* Isolates from IBD Patients

## 14 Genomes:

- Reference: ATCC 29149 (RefSeq GCF\_008121495)
- ATCC 29149 *de novo* assembly (by me)
- ATCC 35913 (GenBank GCA\_900036035)
- 12 Genomes from Hall et al. (2017) (table at right)

RJX1118*	<a href="#">Stool from infant treated with antibiotics</a>
RJX1119*	<a href="#">Stool from infant treated with antibiotics</a>
RJX1120*	Biopsy from IBD patient
RJX1121*	Biopsy from IBD patient
RJX1122*	Biopsy from IBD patient
RJX1123*	Biopsy from IBD patient
RJX1124*	Biopsy from IBD patient
RJX1125*	Biopsy from IBD patient
RJX1126*	Biopsy from IBD patient
RJX1127*	Biopsy from IBD patient
RJX1128*	Stool from IBD patient



# *R. gnavus* Isolates from IBD Patients

**Game 1 : Spot the 2 Genomes from Infant Stool (non-IBD)**



# *R. gnavus* Isolates from IBD Patients

**Game 1 : Spot the 2 Genomes from Infant Stool (non-IBD)**

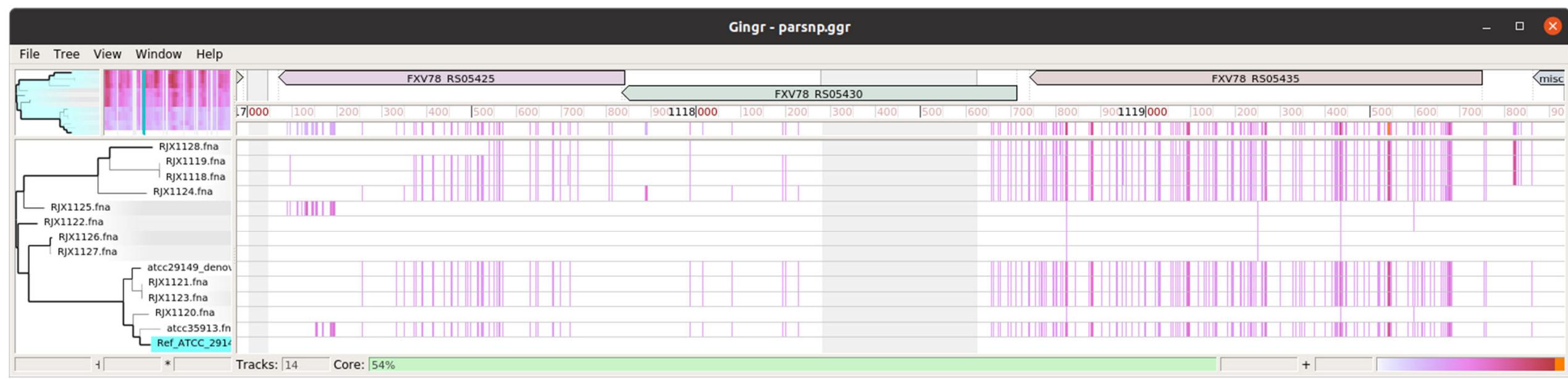
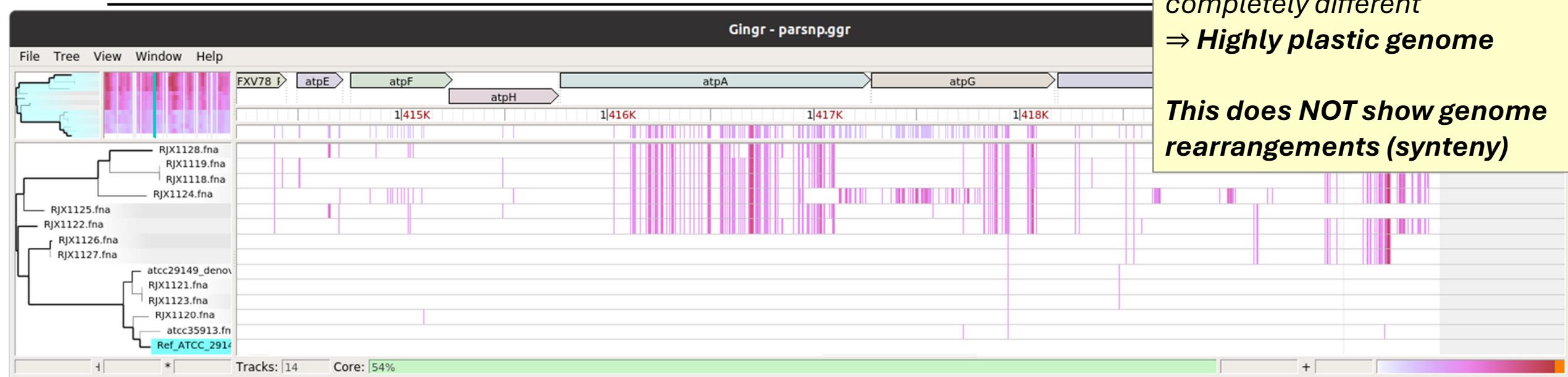
**Game 2 : Spot the 2<sup>nd</sup> ATCC 29149 genome (supposedly the same as the reference)**



# *R. gnavus* strain-level phylogenetic signal is a mess

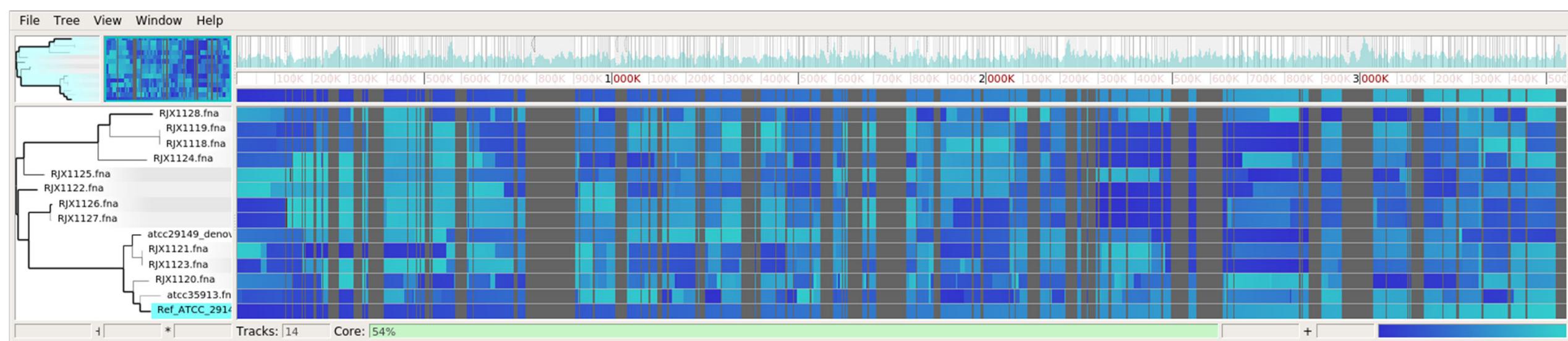
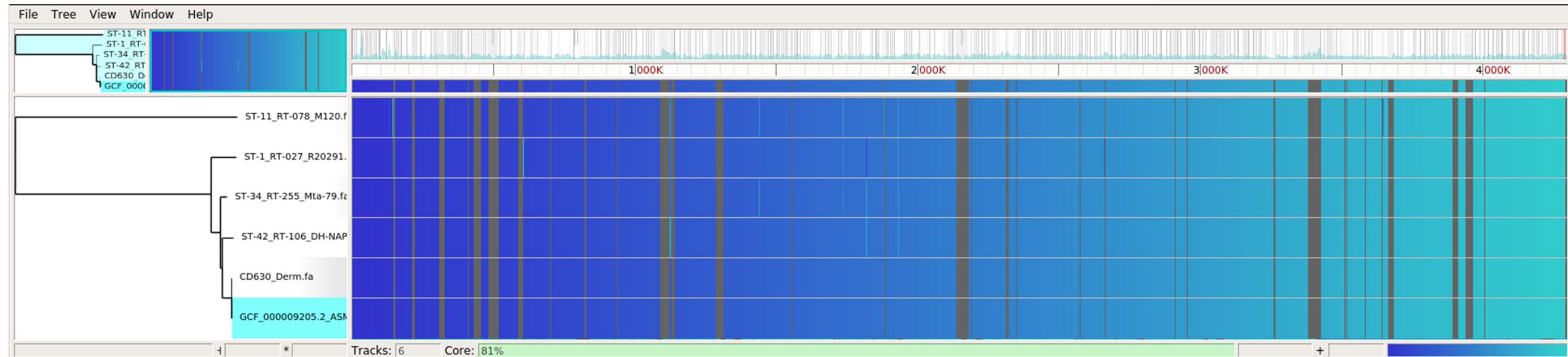
Depending on the operon, the phylogenetic appearance is completely different  
**⇒ Highly plastic genome**

This does NOT show genome rearrangements (synteny)



# Synteny Comparison: *R. gnavus* & *C. difficile*

These two organisms have very different types of genome plasticity.



# Conclusions

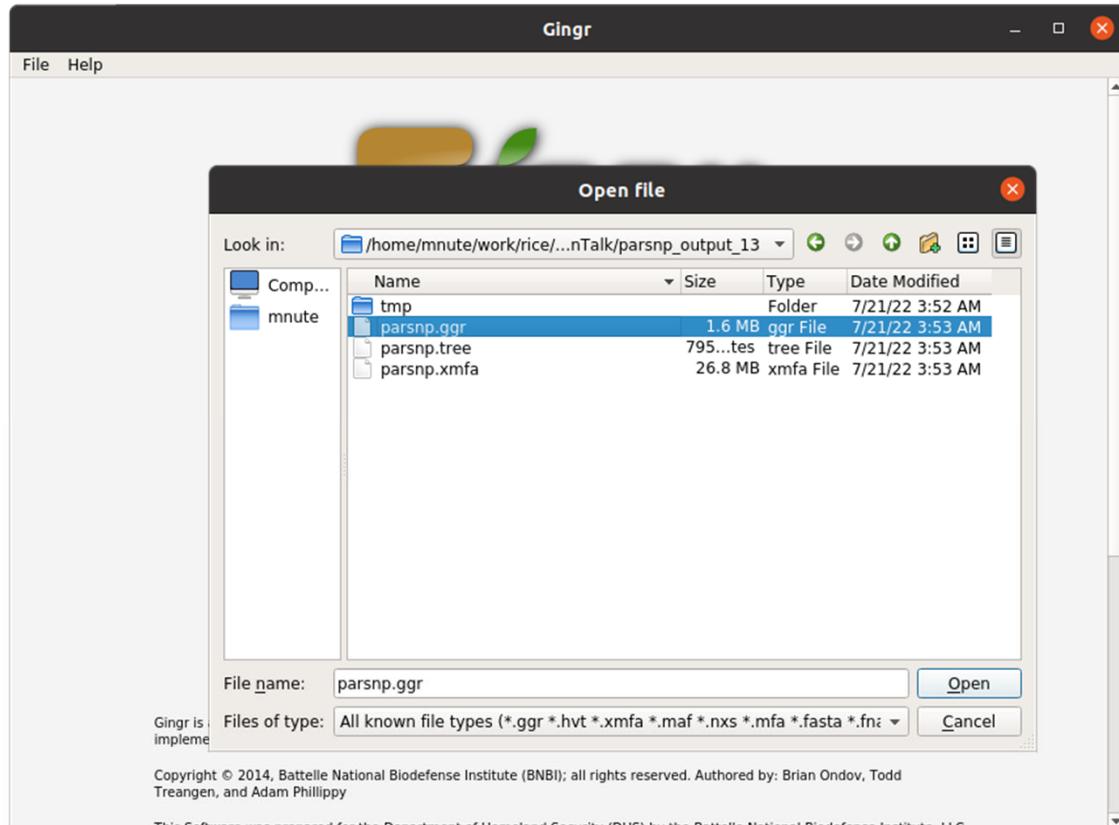
---

- Whole-genome alignment will give a detailed comparison specifically of the *core* genome
  - Maybe also auxiliary genes (*pan*-genome)
- Visualization can get you up close and personal with the data
  - (This statement applies to almost everything, not just genomes)
- Strains can differ from one another in weird ways.
  - Selective mutation at points of interest
  - Gene gain/loss depending on environment
  - Genome-wide phylogenetic signal vs. Locus-specific signal
  - Etc...?

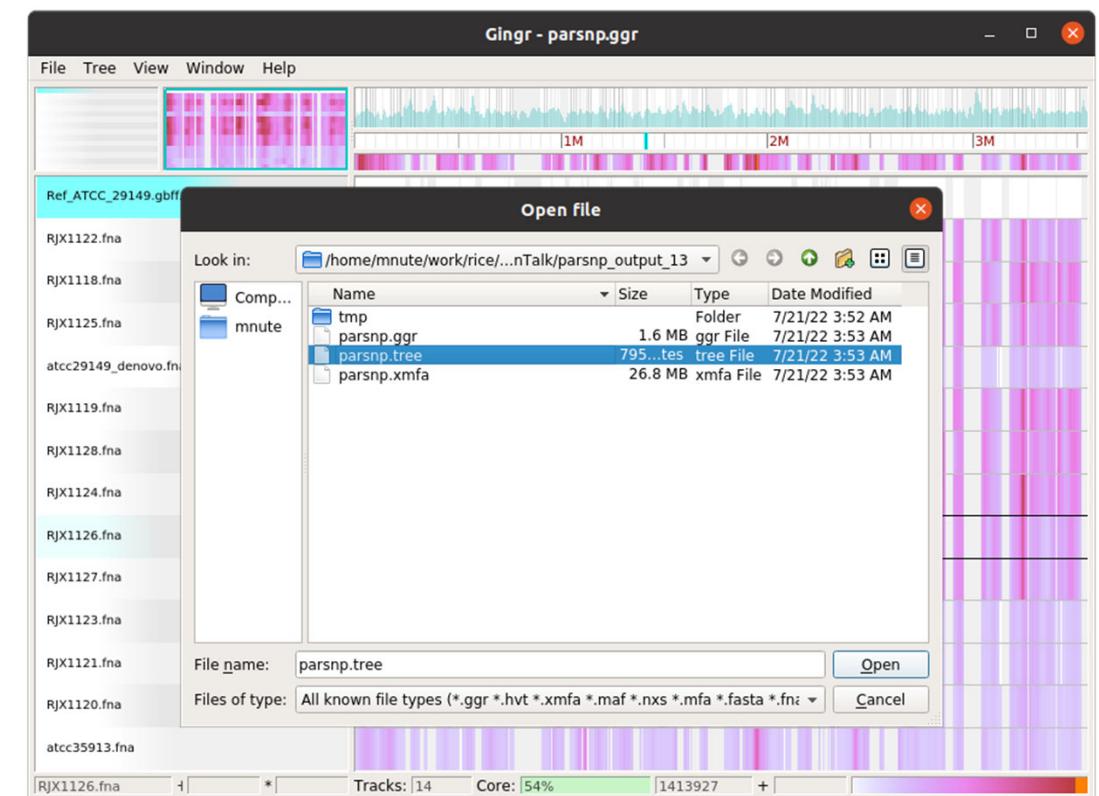
## **Special Thanks To:**

- The Treangen Lab (Rice)
  - Todd Treangen
  - Bryce Killie
  - Kristen Curry
  - Nick Sapoval
  - Yunxi Liu
  - Yilei Fu
  - Advait Balaji
- The Savidge Lab (Baylor College of Medicine)
  - Qinglong Wu
  - Charlie Seto
- Taylor Reiter (for the *R. gnavus* idea)

# Appendix: Quick How-to with Gingr (1 of 2)

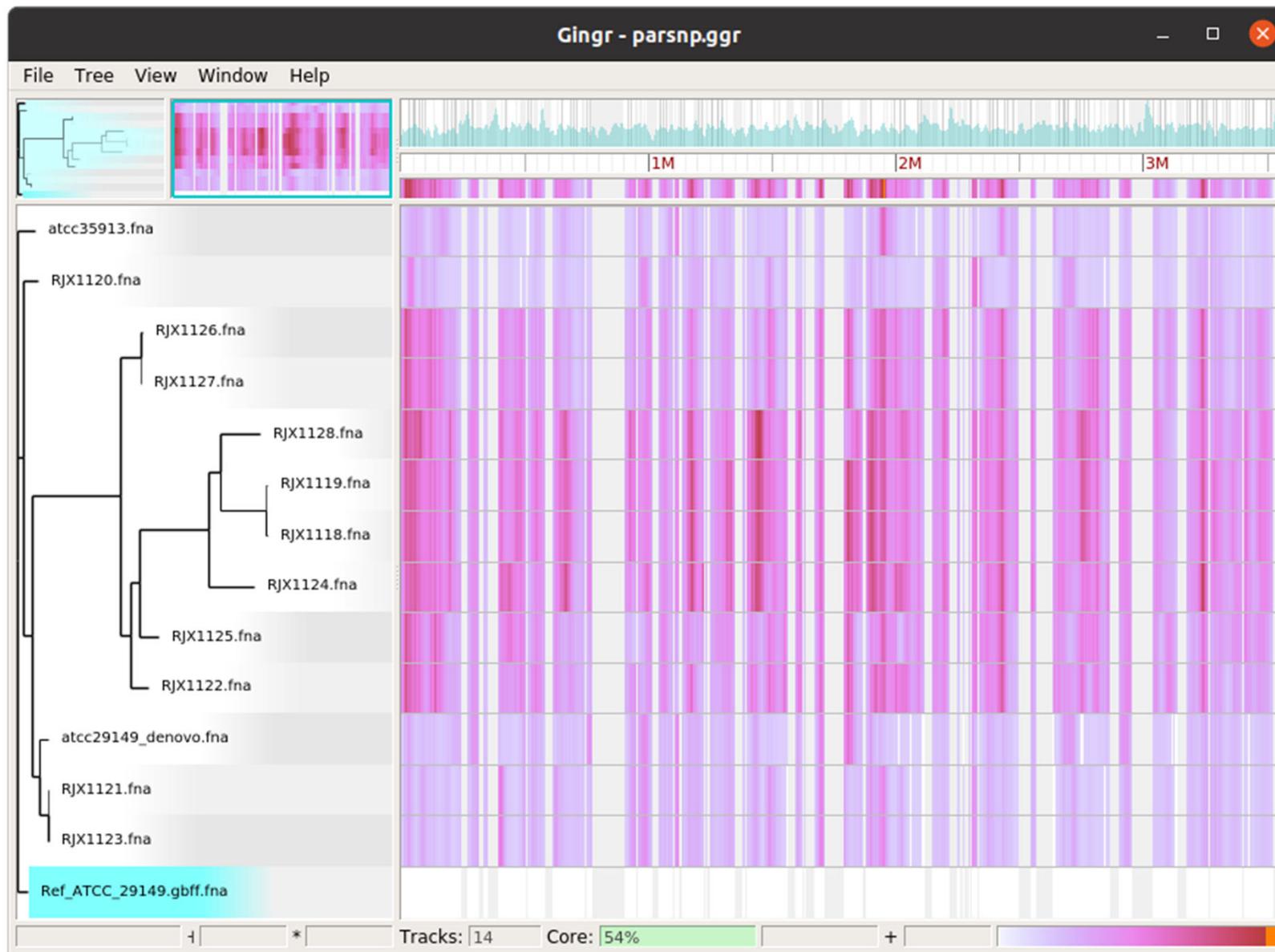


1.) Open the \*.ggr file created in the parsnp output folder.



2.) Once it is open, go back to the “Open” dialogue and open the \*.tree file in the same folder.

# Appendix: Quick How-to with Gingr (2 of 2)



3.) This will give you the standard Gingr view. Other options to re-root the tree or to switch to Synteny view are available under the “Tree” and “View” menus.