genomics
aotearoa

# Day 4

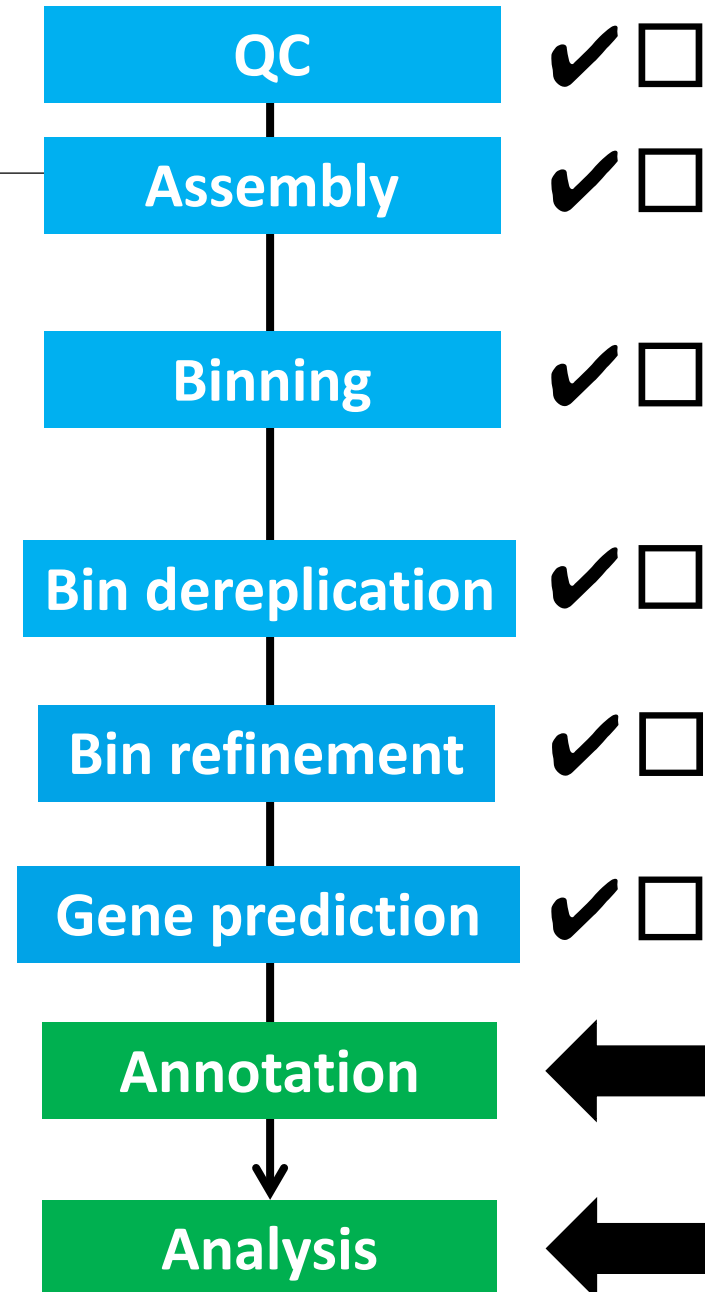Genome annotation (continued)
Report findings
Presentation of data

# Day overview

- **Goals:**

  - **Gene annotations (continued)**

  - **Presentation of group findings**

  - **Visualizing metagenomic data**

| QC | ✔ ☐ |
| --- | --- |
| **Assembly** | ✔ ☐ |
| **Binning** | ✔ ☐ |
| **Bin dereplication** | ✔ ☐ |
| **Bin refinement** | ✔ ☐ |
| **Gene prediction** | ✔ ☐ |
| **Annotation** | ⬅ |
| **Analysis** | ⬅ |

# Genome annotation (continued)

# Bin taxonomic classification

- **16S rRNA commonly not recovered by *de novo* assembly**

- **Can recover 16S and 18S using EMIRGE**

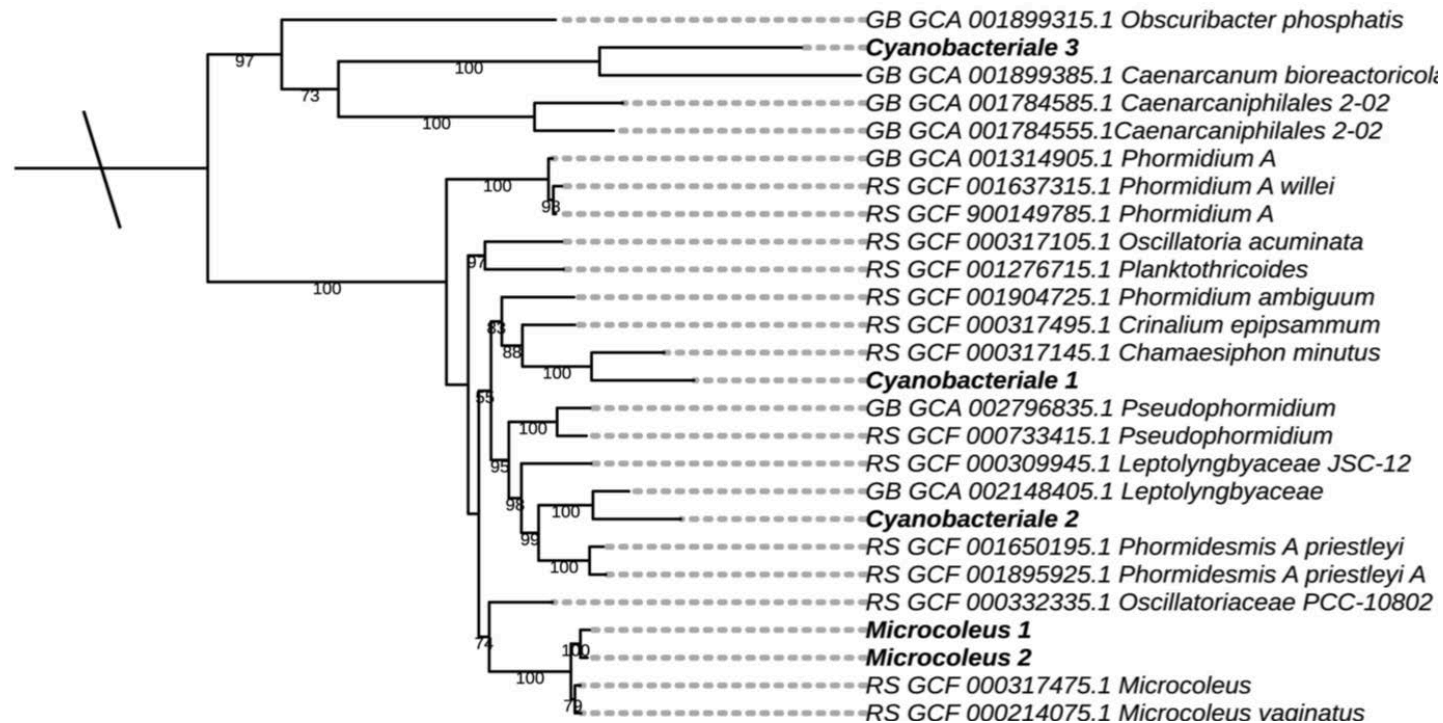- **Caveat: can be difficult to assign to genomes in complex communities with many similar taxa**

# Bin taxonomic classification

**Solution:**

- **Use one or more single copy core genes**

- **Concatenate protein sequences of multiple single copy core genes**



Tree scale: 0.1

Concatenated protein sequence tree: Phylogenetic placement of cyanobacterial genome bins (Wai-iti River, Nelson)

# Bin taxonomic classification

- **Genome Taxonomy Database Toolkit (GTDB-Tk)**

- **Use to classify bins against reference genome trees (GTDB)**



Welcome to GTDB

GENOME TAXONOMY DATABASE

145,904 genomes

Release 04-RS89 (19th June 2019)

- **Uses set of 120 concatenated protein sequences (of single copy core genes)**

https://github.com/Ecogenomics/GtdbTk
https://gtdb.ecogenomic.org/

**Rank assignment based on:**

- Tree placement

- Relative Evolutionary Divergence (value between 0=root and 1=tip)

- Species assignment:

  - Average Nucleotide Identity

# Discriminate species

**Proxy for DNA-DNA hybridization**

**Pairwise genome comparisons:**

- **Average Nucleotide Identities (ANI)**

  - **gene comparisons**

- **Average Amino Acid Identities (AAI)**

  - **predicted protein comparisons**

- **Alignable Fraction (AF)**

  - **proportion of genes that align**

**Determine via: Pairwise BLAST-like search**

# Phylogenetic trait distributions

- **Interactive phylogenetic and trait based tree**

- **Annotree (http://annotree.uwaterloo.ca/app/)**

- **Trait searches by:**

  - **Taxonomic hierarchy**

  - **KEGG (KO number)**

  - **Pfam**

  - **TIGRFAM**

# Phylogenetic trait distributions

**Get KEGG KO number from the KEGG website or your annotations**



https://www.genome.jp/kegg-bin/get_htext#C17

# Phylogenetic trait distributions

**Add to ANNOTREE search box and select hierarchy**
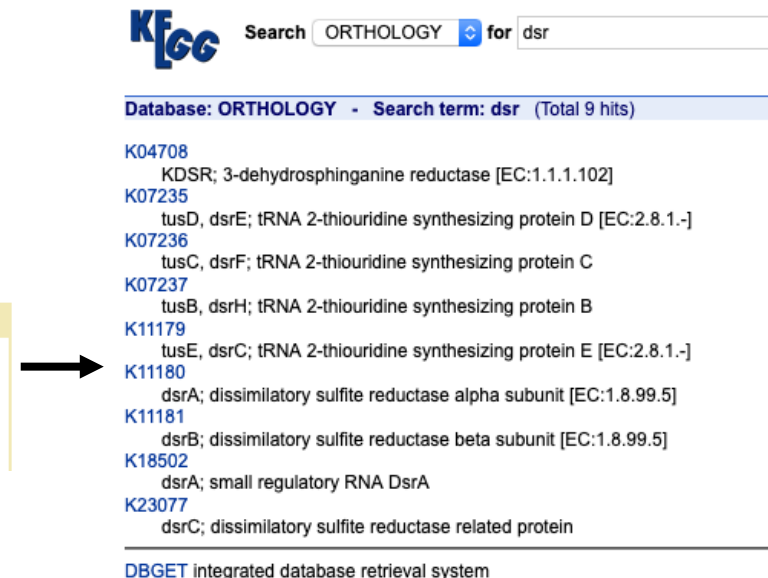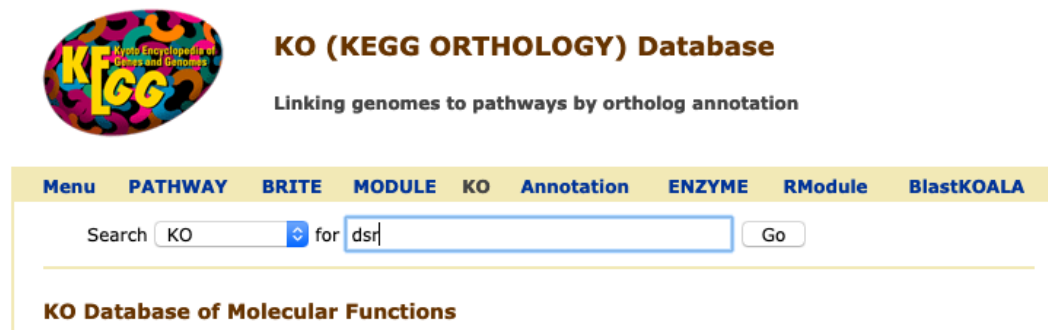


http://annotree.uwaterloo.ca/app/

# Use ANNOTREE
# and
# Prepare group presentations

# Task 1: Use ANNOTREE

- Use ANNOTREE to explore the phylogenetic distribution of functions

- Try using attribute annotations for your group task

- You can use your KEGG Orthology (KO) numbers

- **Note:** You can also get KO numbers from the KEGG website (https://www.genome.jp/kegg/ko.html) by searching for gene names

# Task 2: Prep for group presentations

- **Use the white board for illustrations**

- **Things to include:**

  - The attribute you found

  - Details about the attribute

  - The organism(s) you found it in

  - A brief explanation of biological relevance

  - The tools and annotations you used

  - Anything else?

# Presentation of findings

# Task: Report findings

**Report your findings with regards to your objective**

1. What did you find?
2. How did you make the discoveries?
   - Functional prediction only
   - Functional prediction and taxonomic context
3. Each group present for 5 mins each (max!)

# Presentation of data

# Presentation of data

**Finding the answer to your question is only half of the issue**

- How do we report/present our data?

1. **Inference of gene trees**

2. **Heatmaps of genomic features**

3. **Gene synteny analysis**

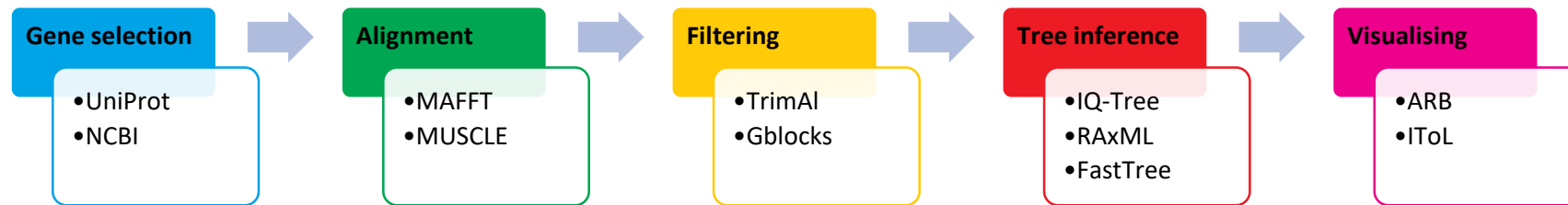4. **Creating metabolic schematics**

# Presentation of data
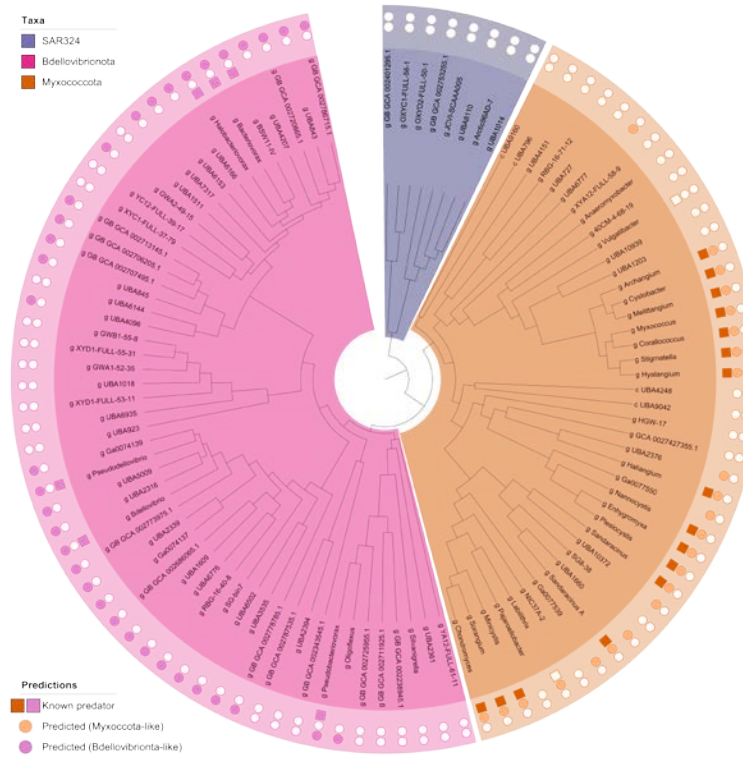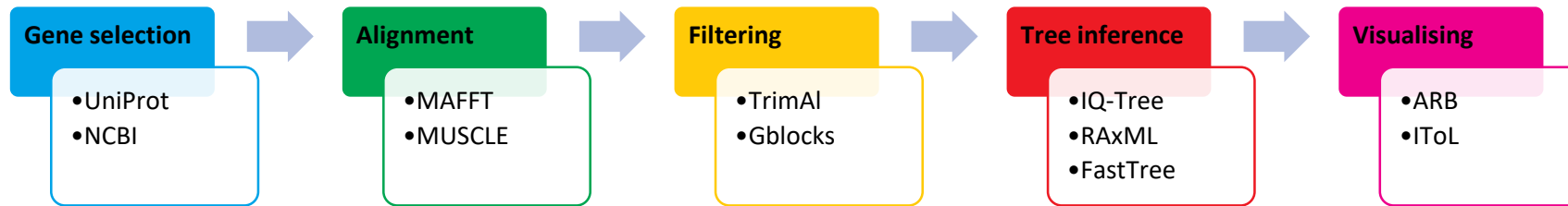
1. **Inference of gene trees**

**Gene trees are a great way to present**

- Confirmation of annotation

- Novelty of detection / horizontal gene transfer

- Rate of evolution in the feature

| Gene selection | Alignment | Filtering | Tree inference | Visualising |
|---|---|---|---|---|
| •UniProt<br>•NCBI | •MAFFT<br>•MUSCLE | •TrimAl<br>•Gblocks | •IQ-Tree<br>•RAxML<br>•FastTree | •ARB<br>•IToL |

# Presentation of data

## 1. Inference of gene trees

# Presentation of data

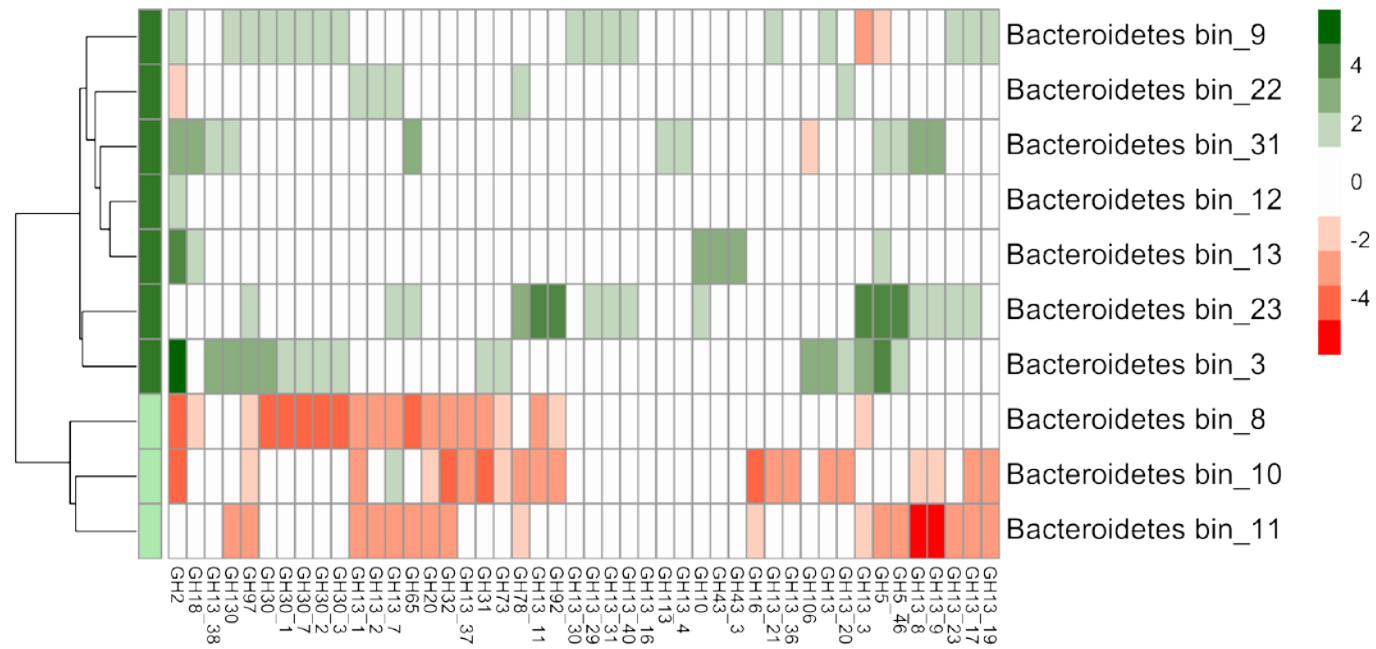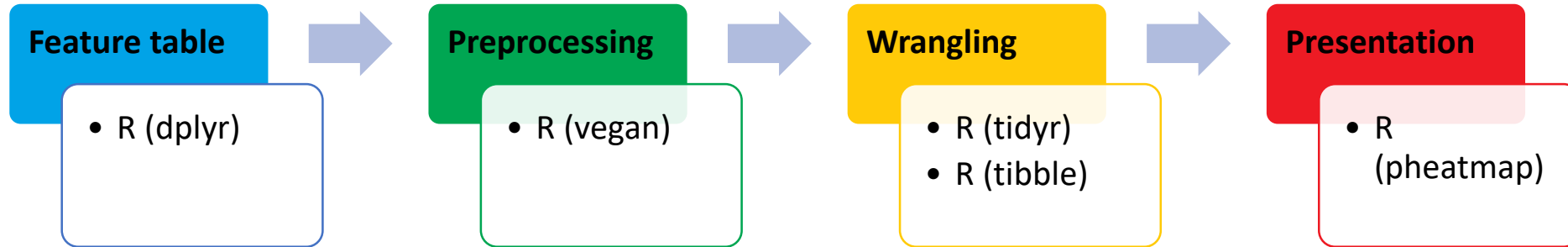**2. Heatmaps of genomic features**

**Simple figure to display complex data tables**

- *M* genomes x *N* features in one place

- Presence/absence or relative abundance (multi-copy)

- Fixed layout, or clustering by patterns

# Presentation of data

## 2. Heatmaps of genomic features

# Presentation of data
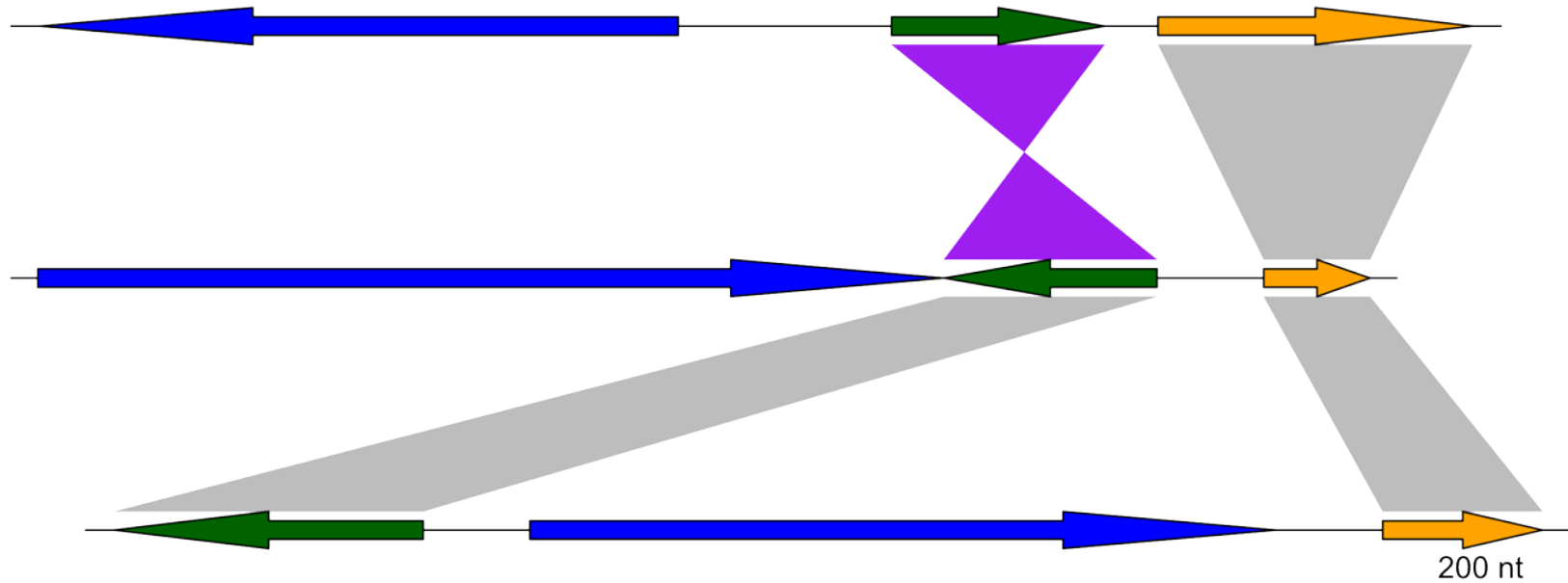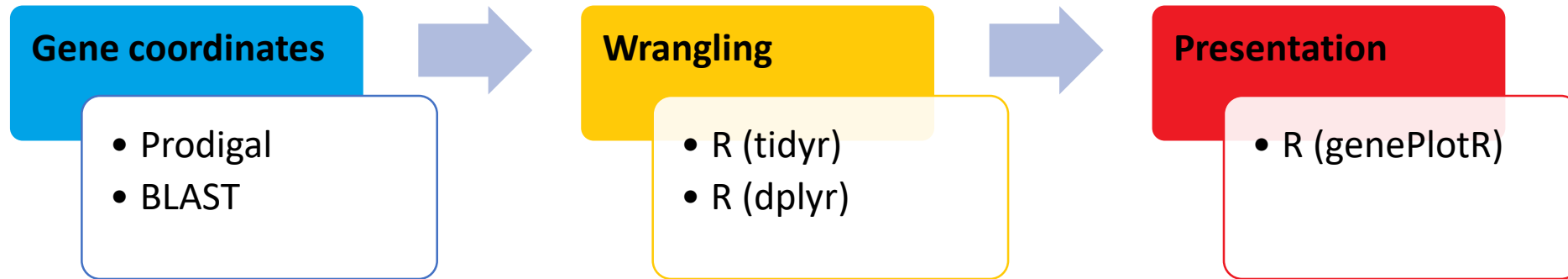
3. **Gene synteny analysis**

**A more informative view of gene content**

- Shows local gene context

- More detailed than reporting gene table

- Sometimes absence of genes from operon is biologically informative

# Presentation of data

3. **Gene synteny analysis**

# Presentation of data

**4. Creating metabolic schematics**

Summarise the entire core metabolism of an organism into a single figure

**What's the magic tool for producing these?**

- Illustrator, Inkscape, GIMP (and a lot of time)

- Use tools for picking colour schemes

  - ColorBrewer2 (http://colorbrewer2.org)

  - IWantHue (https://medialab.github.io/iwanthue/)

# Task: Presentation of data

- **Using Pheatmap to build CAZY plot**

- **Using genoPlot R to build comparative a gene map**

# Wrap up and Q/A