# Rapid Identification of AMR Determinants from Metagenomic Samples

AMRtime Progress Report

Finlay Maguire

June 22, 2018

Faculty of Computer Science, Dalhousie University

## Table of contents

# Overview

## Comprehensive Antibiotic Resistance Database

- https://card.mcmaster.ca/ (Jia et al., 2016) as of June 2018:

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
    - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
    - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
    - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
  - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
  - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation
  - rRNA gene variants e.g. Mycobacterium aminoglycoside resistance

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
    - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
    - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation
    - rRNA gene variants e.g. Mycobacterium aminoglycoside resistance
    - Efflux pump e.g. AcrAB-TolC, MexAB-OprM mutations

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
  - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
  - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation
  - rRNA gene variants e.g. Mycobacterium aminoglycoside resistance
  - Efflux pump e.g. AcrAB-TolC, MexAB-OprM mutations
  - Gene cluster e.g. Van glycopeptide resistance clusters

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
    - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
    - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation
    - rRNA gene variants e.g. Mycobacterium aminoglycoside resistance
    - Efflux pump e.g. AcrAB-TolC, MexAB-OprM mutations
    - Gene cluster e.g. Van glycopeptide resistance clusters
- Resistance Gene Identifier (RGI): contigs, predicted genes and merged metagenomic reads

## Comprehensive Antibiotic Resistance Database

- `https://card.mcmaster.ca/` (Jia et al., 2016) as of June 2018:
- Built around Antibiotic Resistance Ontology (ARO): 3996 terms
- 2536 AMR Detection Models with manually curated criteria:
    - Homology e.g. NDM beta-lactamases, aminoglycoside acetyltransferase
    - Protein Variant e.g. GyrA fluoroquinolone mutation, FolP sulfonamide mutation
    - rRNA gene variants e.g. Mycobacterium aminoglycoside resistance
    - Efflux pump e.g. AcrAB-TolC, MexAB-OprM mutations
    - Gene cluster e.g. Van glycopeptide resistance clusters
- Resistance Gene Identifier (RGI): contigs, predicted genes and merged metagenomic reads
- CARDPredicted prevalence dataset

modified from https://www.gatc-biotech.com/en/expertise/genomics/metagenome-analysis.html

Key difficulties:

- Variation in abundance and diversity

Key difficulties:

- Variation in abundance and diversity
- Short fragmentary data

# Metagenomic Analysis



modified from `https://www.gatc-biotech.com/en/expertise/genomics/metagenome-analysis.html`

Key difficulties:

- Variation in abundance and diversity
- Short fragmentary data
- Large amounts of data

# Metagenomic Analysis



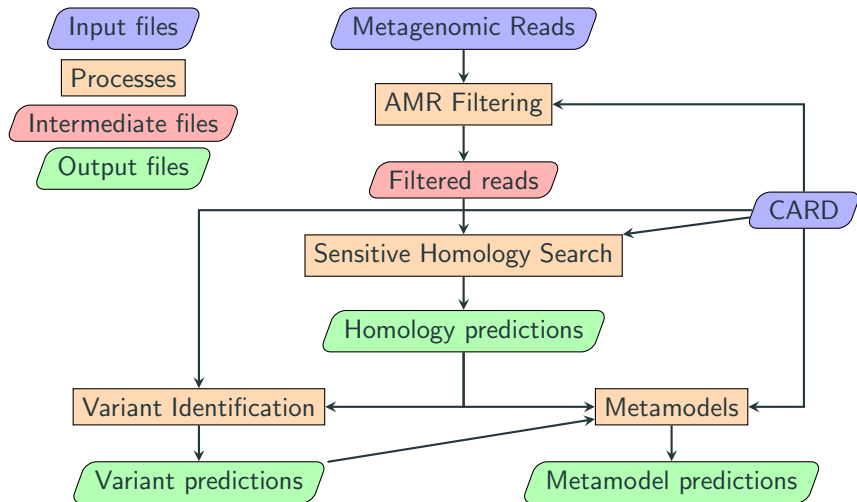modified from https://www.gatc-biotech.com/en/expertise/genomics/metagenome-analysis.html

Key difficulties:

- Variation in abundance and diversity
- Short fragmentary data
- Large amounts of data
- Compositionality

# Metagenomic Analysis



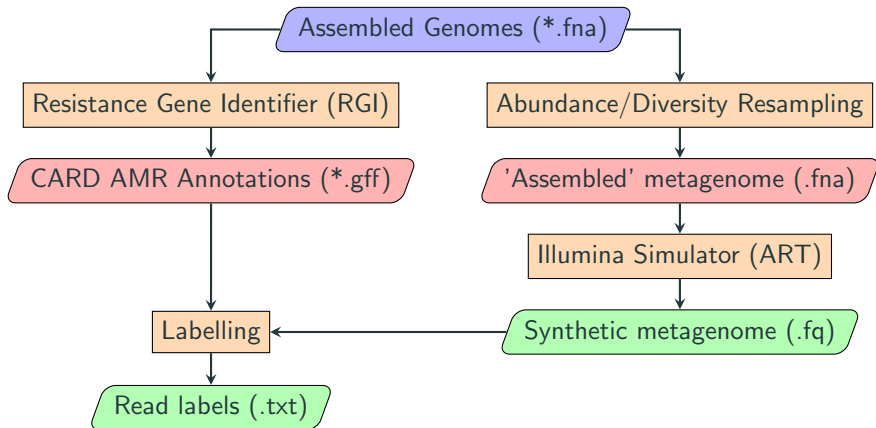modified from https://www.gatc-biotech.com/en/expertise/genomics/metagenome-analysis.html

Key difficulties:

- Variation in abundance and diversity
- Short fragmentary data
- Large amounts of data
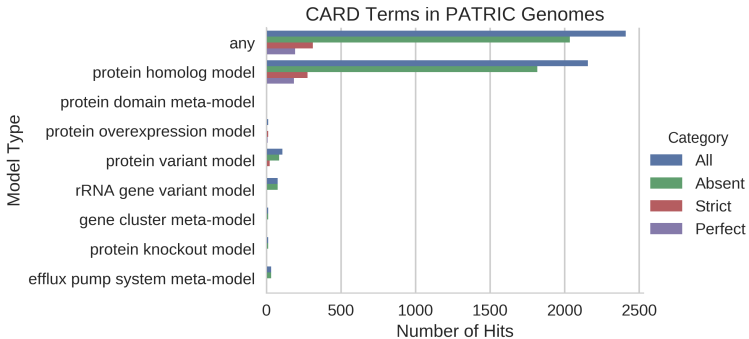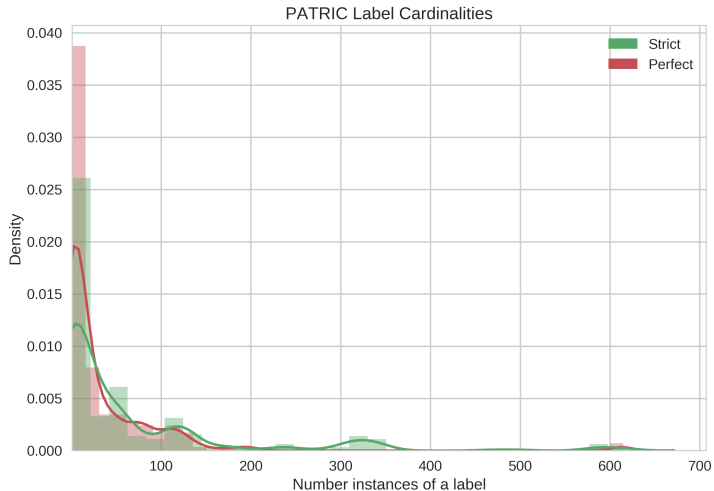- Compositionality
- Spare and imbalanced labels

Input files

Processes

Intermediate files

Output files

Metagenomic Reads

AMR Filtering

Filtered reads

CARD

Sensitive Homology Search

Homology predictions

Variant Identification

Metamodels

Variant predictions

Metamodel predictions

# Training Data

# Dataset Generator

CARD Terms in PATRIC Genomes
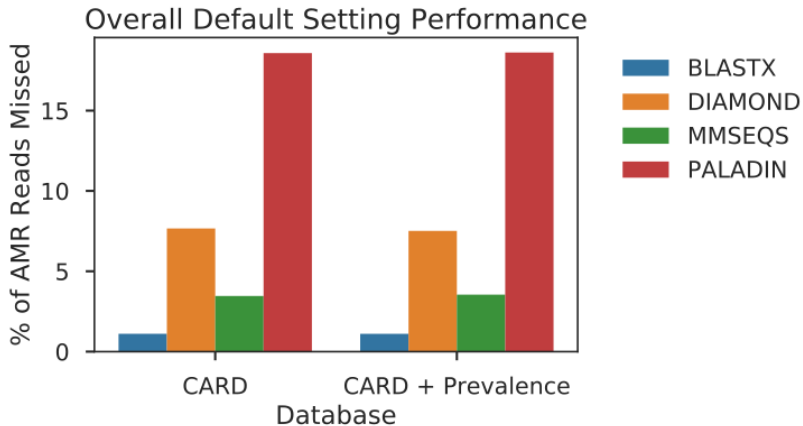
PATRIC Label Cardinalities

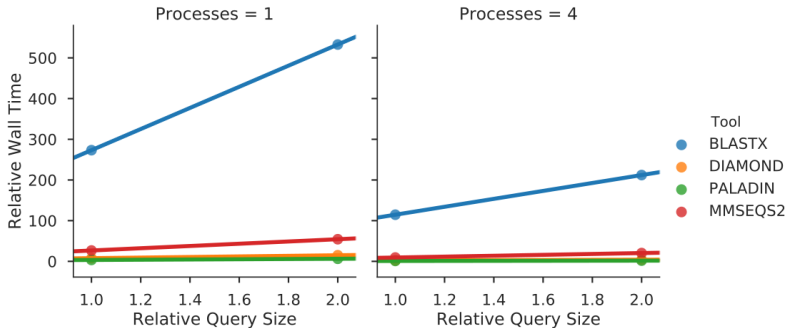# AMR sequence space is biased

# Read filtering

## Homology Filter Approaches

- BLASTX (Gish et al., 1993)
- DIAMOND (Buchfink et al., 2015)
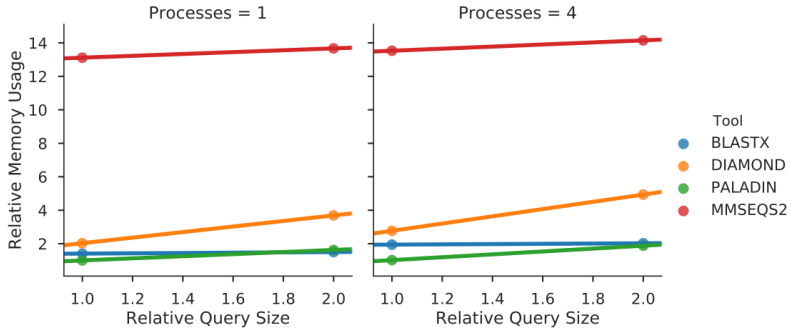- PALADIN (Westbrook et al., 2017)
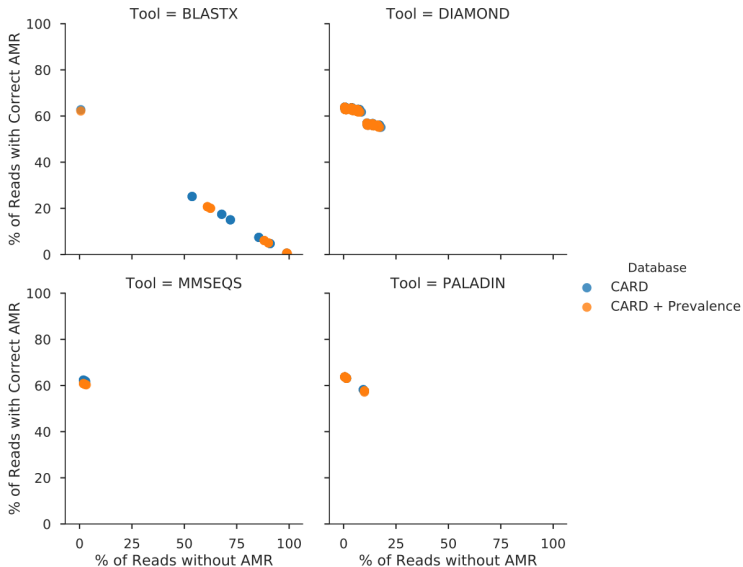- MMSeqs2 (Steinegger and Söding, 2017)

Overall Default Setting Performance

# What about in terms of memory?

# Is there a cap on overall performance?

False Negative Distribution

Truncated False Negative Distribution

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
    - Protein 2790824-2789724

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
  - Protein 2790824-2789724
  - DNA 1-732

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
  - Protein 2790824-2789724
  - DNA 1-732
- OXA-2 (M95287.4):

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
    - Protein 2790824-2789724
    - DNA 1-732
- OXA-2 (M95287.4):
    - Protein 2456-3280

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
  - Protein 2790824-2789724
  - DNA 1-732
- OXA-2 (M95287.4):
  - Protein 2456-3280
  - DNA 1-828

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
  - Protein 2790824-2789724
  - DNA 1-732
- OXA-2 (M95287.4):
  - Protein 2456-3280
  - DNA 1-828
- Acinetobacter OprD conferring resistance to imipenem (CP006768.1):

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
  - Protein 2790824-2789724
  - DNA 1-732
- OXA-2 (M95287.4):
  - Protein 2456-3280
  - DNA 1-828
- Acinetobacter OprD conferring resistance to imipenem (CP006768.1):
  - Protein 3513470-3514777

## Why are these 10 always missed?

- Enterococcus faecalis liaS mutant conferring daptomycin resistance (AE016830.1):
    - Protein 2790824-2789724
    - DNA 1-732
- OXA-2 (M95287.4):
    - Protein 2456-3280
    - DNA 1-828
- Acinetobacter OprD conferring resistance to imipenem (CP006768.1):
    - Protein 3513470-3514777
    - DNA 3514887-3515414

# CARD Full Length Alignment QC

- 11 AROs protein not detected from DNA

# CARD Full Length Alignment QC

- 11 AROs protein not detected from DNA
- 2 AROs different top protein hit from DNA

# CARD Full Length Alignment QC

- 11 AROs protein not detected from DNA
- 2 AROs different top protein hit from DNA
- Warnings: 119 AROs with different top protein but $ID\% > 99$
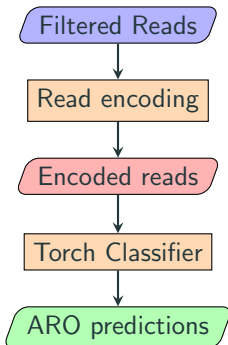
## CARD Full Length Alignment QC

- 11 AROs protein not detected from DNA
- 2 AROs different top protein hit from DNA
- Warnings: 119 AROs with different top protein but $ID\% > 99$
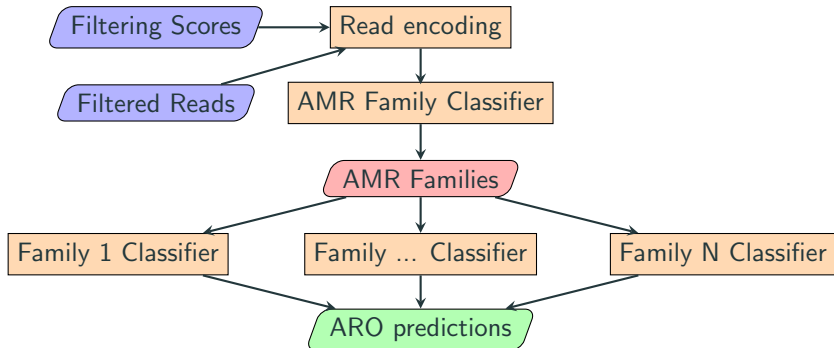- Warnings: 2 AROs with $ID\% < 99$ to correct protein

# Sensitive Homology Search

# Revised classifier structure

- Raw sequence

## Encodings

- Raw sequence
- Filtering homology search family similarity/dissimilarity

## Encodings

- Raw sequence
- Filtering homology search family similarity/dissimilarity
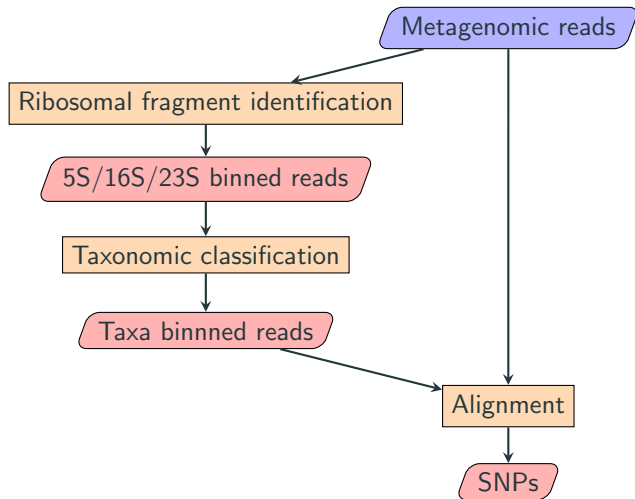- Manual feature extraction (GC/TNF/compositional)

- Raw sequence
- Filtering homology search family similarity/dissimilarity
- Manual feature extraction (GC/TNF/compositional)
- One-hot K-mer representation

- Raw sequence
- Filtering homology search family similarity/dissimilarity
- Manual feature extraction (GC/TNF/compositional)
- One-hot K-mer representation
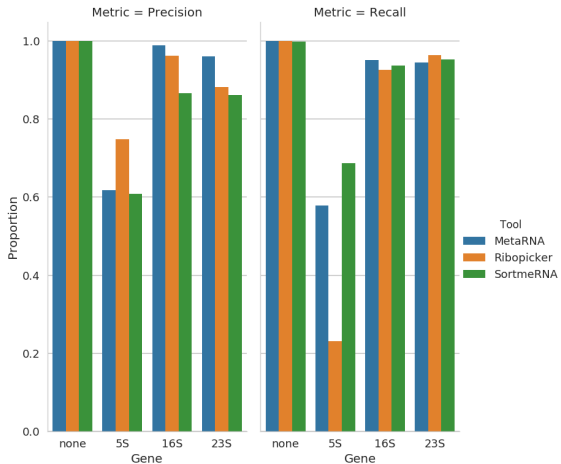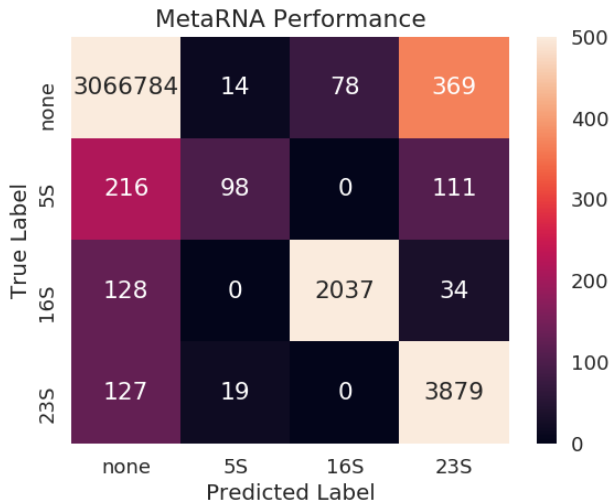- K-mer embeddings (DNA2vec/BioVec)

# Variant Models

## Identifying Ribosomal Reads

- MetaRNA (Huang et al., 2009)
- Ribopicker (Schmieder et al., 2011)
- SortmeRNA (Kopylova et al., 2012)
- 77 models
- Reads simulated from the underlying 30 species reference genomes
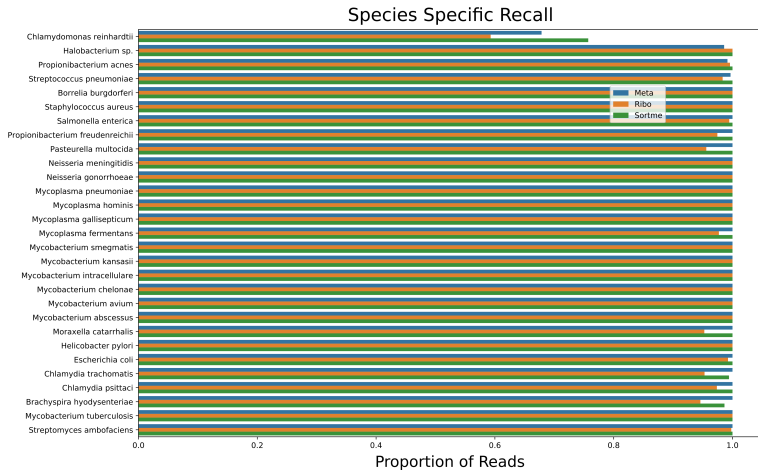
# Identifying Ribosomal Reads



MetaRNA Performance

Species Specific Recall

# Identifying Taxonomy

# Some are relatively easy



Correct_index:390_species:Streptomyces ambofaciens

# Others are a mess



Misspredict_index:750_species:Chlamydomonas reinhardtii

Misspredict_index:259 species:Mycobacterium chelonae

Probably a Mycobacterium?

Ambiguous_index:212_species:Escherichia coli

# Ambiguity in classification

## Next Steps

- Mapping reads to reference to assess presence or absence of mutation related SNP

- Comparison of whole pipeline with just direct mapping to database of ribosomal sequences and SNP calling approaches.

- Tuning of sensitivity for number of potential SNPs required to make a prediction of AMR.

# Summary

- AMRtime still not a 'fait accompli'

## Conclusions

- AMRtime still not a 'fait accompli'
- Filtering analysis possibly needs redone for fixed CARD

## Conclusions

- AMRtime still not a 'fait accompli'
- Filtering analysis possibly needs redone for fixed CARD
- False positive analysis pending for best settings

## Conclusions

- AMRtime still not a 'fait accompli'
- Filtering analysis possibly needs redone for fixed CARD
- False positive analysis pending for best settings
- Framework and code developed for sensitive homology classification but optimisation and evaluation work still required

## Conclusions

- AMRtime still not a 'fait accompli'
- Filtering analysis possibly needs redone for fixed CARD
- False positive analysis pending for best settings
- Framework and code developed for sensitive homology classification but optimisation and evaluation work still required
- Not shown but preliminary family level classification shows 100x improvements over previous ARO attempts

## Conclusions

- AMRtime still not a 'fait accompli'
- Filtering analysis possibly needs redone for fixed CARD
- False positive analysis pending for best settings
- Framework and code developed for sensitive homology classification but optimisation and evaluation work still required
- Not shown but preliminary family level classification shows 100x improvements over previous ARO attempts
- Ribosomal Variant Model work progressing well with full pipeline metrics available soon.

# Acknowledgements

## Acknowledgements

# References

Buchfink, B., Xie, C., and Huson, D. H. (2015). Fast and sensitive protein alignment using diamond. *Nature methods*, 12(1):59.

Gish, W. et al. (1993). Identification of protein coding regions by database similarity search. *Nature genetics*, 3(3):266.

Huang, Y., Gilna, P., and Li, W. (2009). Identification of ribosomal rna genes in metagenomic fragments. *Bioinformatics*, 25(10):1338–1340.

Jia, B., Raphenya, A. R., Alcock, B., Waglechner, N., Guo, P., Tsang, K. K., Lago, B. A., Dave, B. M., Pereira, S., Sharma, A. N., et al. (2016). Card 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. *Nucleic acids research*, page gkw1004.

Kopylova, E., Noé, L., and Touzet, H. (2012). Sortmerna: fast and accurate filtering of ribosomal rnas in metatranscriptomic data. *Bioinformatics*, 28(24):3211–3217.

Schmieder, R., Lim, Y. W., and Edwards, R. (2011). Identification and removal of ribosomal rna sequences from metatranscriptomes. *Bioinformatics*, 28(3):433–435.

Steinegger, M. and Söding, J. (2017). Mmseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nature biotechnology*, 35(11):1026.

Westbrook, A., Ramsdell, J., Schuelke, T., Normington, L., Bergeron, R. D., Thomas, W. K., and MacManes, M. D. (2017). Paladin: protein alignment for functional profiling whole metagenome shotgun data. *Bioinformatics*, 33(10):1473–1478.