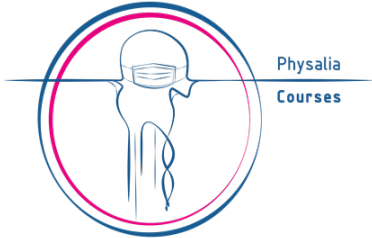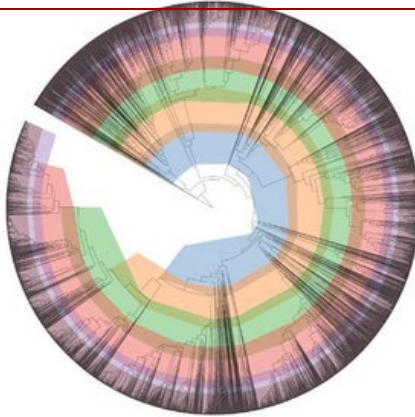**ENVIRONMENTAL METAGENOMICS**
Physalia course, online, 11-15 November 2024

**MAG functional annotation**

Nikolay Oskolkov, Lund University, NBIS SciLifeLab
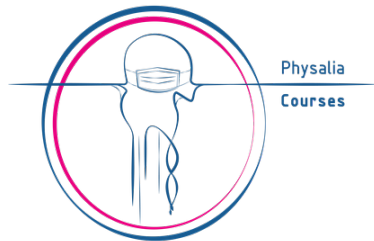Samuel Aroney, Queensland University of Technology

Physalia
Courses

# You got MAGs!
# You know that they are good
# Now what?

**Functional annotation**
- Genes
- *eggNOG*
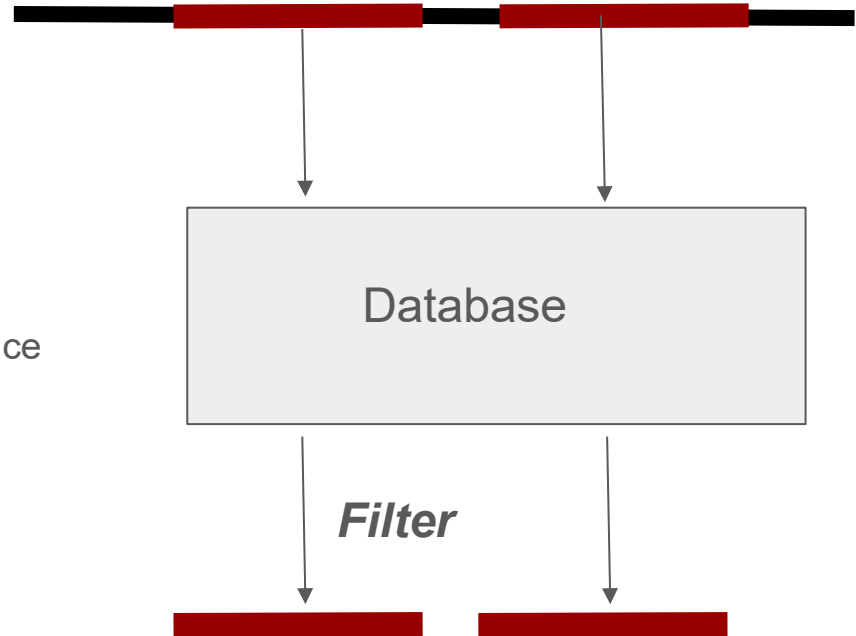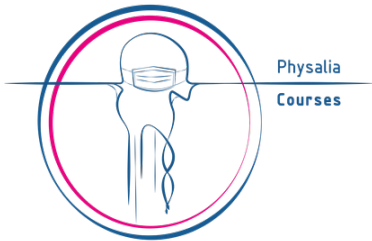- *RGI*
- *Other gene annotations databases*
- *…*

```
ACTTCTCGGGCCGACGCCTGTTTCGATCGATGATTTGGTGCGTCTTGCGCAGACCACGGCGGCGATCGTGCGCGCTGCCCTGCTGGAACTCGAAATTGCCGGGCGGATTGAACC
CGATACCCCCAATCGTCATGTTCATCCGATTTCACCTGGCTGTCGCGGCTATACTGATCGTGCTGGTGGCGTCATCGCGGGGCCAGGAACAACCCGTTTCGGCTTTTCTCGCAC
CGCACGACGTTCGATACCGCCATGAAGGGATTCACCCCCGATCAACGGGTGATCACCGCGACCCGGCGTCAGCCGGAATACGGCAAACCCGTCGGCGACTACGTGAACGCCATC
AGGCCGCGAATGGGCAAAAACATTCGATGTCGTCGAAAAGAAGTTTCAGGTCGAACGCTGGGTCTTGCTCGCCTTGTGGGGCATGGAATCGGACTTCGGGTCAGAAAAAGATCC
ATGTGAAATTCCGCGACCCCTATTTTCGCAACGAGCTGATCGTGGCGATGCGCATCATGCAAGACAATCGCATCGCCCGCGAAAAGATGGTCAGCTCTTGGGCCGGCGCGATGG
TATGCGATCGATTTTTCCGGCGACGGACGGGCGGATATCTGGGGCAACGTACCGGATGTGCTTGCGTCGACCGCCAACTACCTGCGCAAATGGAAATGGAATCCGGCGCTTCCC
CTACATGCGCAGCCGCGCAAGTTTTGCGGAATGGCAAGCGCTCGGCGTGCGCCGCGCGACGGCAAGGCGTTCCCCAACTCCGGACAGGGCATCCTCTTCTTTCCCAGCGGCGG
ATGTGCTCAAGGAATACAACAACTCCGATGCTTACGCGCTCGCGGTCGGGCACCTCGCCGACCGAATCCACGGCGGCGATCTGATCAAGCACGCCCTGGCCTAAGGACGATCGC
AGACTCGCGGCACTCGGCTACAAAGTGAACGAGTTTGAGGCCCACATCGATTTCGATTTGGTCGGACAACATTCGCGTCGAGCAAAAGAAGCTGGGGATGGTCCCCGACGGCAAT
GCCTCGGCTCTAGGTCTCCAGCAGTACGCGGACCGCCTTGAAGCTGGCAATCCAGCCTTGATACAATTTACGGTTGTCGTCTCTTTCGTGCCCTATTTCGCTGGTCTGACACGG
TTCCCCATGTTCCGGGCCAAACATGCGGCCGATTCCGCCTTGATTCTCAACAGGCCTGGCGGCTAACCCATTGGATTTTCATGAATGTGTCGTAGTCGAGTCGCCGTCCAAGGC
GAGGTTCTGGCCTCATTTGGCCATATCCGGGACCTGCCCCCCAAGGATGGCTCCGTCGATCCCGACAATGATTTCCGCATGCTCTGGGAGGTCGACGCCAGGTCGAACCAGCGC
CAAGCTGATCCTCGCCACCGACCCGGATCGCGAGGGCGAAGCAATTTCCTGGCACGTGCTCGAGGTGCTGAAGGAAAAGAAAGGCGCTCAAGGACCACAAGATCGAGCGCGTCGT
CGATGAAGCATCCACGGATGATCGACGCCGCATTGGTCGATGCCTACCTCGCCGCGCGCGCTCGACTATCTCGTCGGCTTCACCCTTTCACCGGTGCTGTGGCGAAGCTGGC
GCGCTTCGGCTTGTGTGCGATCGCGAACTCGAAATCGAGAAGTTTGTTGCGAAGGAGTATTGGTCGATTCTCGCCAGGCTCGCGACGCCGCGCAACGAAGTGTTCGAAGCGCGT
CGACATAGGTTCGGGCGCCGAAGCGGAAGCTTTCACCCGAGACCTCGAGAATGCGACCTTCAAGGTGACGTCGGTCGAGGCAAAGCCCGCACGGCGCAATCCGCCGCCGCCTT
AGCTCGGCTTTGCACCGGCGGTCGCCATGCGTCTCGCCCAGCGGCTCTACGAAGGCGTGGAAATCGACGGCGAAGCGACCGGCTTGATTACGTATATGCGTACTGACGGCATCC
ATGTTGGGCCGCAACTACGGCAAGGAGTACGTCCCTGCGTCGCCGCGCGAGTACCACAACAAATCCAAGAACGCGCAGGAGGCGCACGAAGCGGTGCGCCCGACCAGCGCGAAG
CGACCAGGCGCGGCTCTACGAGTTGATCTGGAACCGCGCGGTCGCGAACGCCGCCCAAATGGAAATCCGCCGAGCTCGAGCGCACCACGGTCGACATTGTTGCGAAGGCGGGCTCACGCAA
TCGACGGCTTTCTGACGCTCTATCAGGAAGGCCAGGACGAGACGCCGGACGATGACGAGTCGCGGCGTCTGCCCGCGATGTCGGAAGGCGAAACGCTCAGCAAGCAGGCGATCC
TTCTCGGAAGCGGCGCTGGTCAAGCGGATGGAAGAGCTCGGCATCGGCCGGCCCTCCACTTATGCCTCCGTCCTCCAAGGTGCTGCAGGATCGCGGCTATGTGCGGATCGACAAG
CGTCGCCTTCCTTGAAAGCTTCTTCGCGCGCTACGTCGAATACGACTTCACAGCCAGCCTCGAGGAAAAACTCGACGAGATTTCGGCGGGGCAATATCGACTGGCGGGCCGTGCT
ACGACATCAAGGAAGTACGCAATCGCGTCGTGCTCGATGCGCTCAACGACCTACTTGAGCCGCATATTTTCCCGAGCGCACTGACGGCAAACCGCGCCGCCAATGCCCGCAGT
TTCGGGCTTTCGTCGGTTGCTCGAATTATCCGGAATGCAATTTCACGCGGCAATTACTCCGAGCGCCGATGGCATTCAGGGCAAAAGGGTCCTGGGTGAAGATCCGGCCCACA
GCCCTATCTGCAGCTTGGCGAGCAGATCAGGCCGCCCAAGCCGAAGAAAGGCGAGAAAAAACCGGAAGTCGAAAAGCCCAAGCGCGCCGGAATTCCAAAAGGTGTGTCGCCCGA
CGCTGCCGCGTGAAATTGGATTGTCGCCGGAAGACGGCGAGCCGATAGTCGCCGGCATCGGACGCTTCGGCTCATACGTGAAGCACGGCAAGATCTACGCCAACCTGGAAGAAG
GTCACATTGATCGCCGAGAAGAAGGCGAACCCGAAGAAGGGCCGGCGGTTCGGCGCCGATCCAGGCAAGGTCTTGGGCGAACATCCCGATCAGGGCGGCCCGGTTGTCGTCAAG
CATCAACGCGACGATAACCGGCGACAGAACGCCCGAGACCATCACGCTGACTGAAGCGGTCGTGCTGCTTGACGCGCGCGCCGATCAACTGAGTTCTCAGCCGCGGCGCGCGAG
CGAAGGGGCGCAAGAAAAAGACGGCGAACCGGCAGACAAGGCCGAAAAATCAGCCAAGCCGCGACAAGAAAAGCGCCAAGAAAGGCAGAGTCAAAAAAGACTGG
AGCGCCGTTTTGGCGCAGGACTGTCTATTTTGAATAAGAAATCCAAACGAACACCCTCCCTCCCTTCCAAAGCCGACATACTCGCCTTCATCGGTAATCAGCCCGGAAAAGTCC
```

*Congrats!*
*You got big FASTA files!*

# Functional annotation

1. Predict genes
   a. Prodigal
   b. Pyrodigal
2. Predict function of genes
   a. Eggnog-mapper (generic)
   b. RGI (specialized to antibiotic resistance genes)
   c. …
3. Use informed human judgement!!

Database

*Filter*

# Gene annotation with eggnog-mapper

Basic concept: *eggNOG orthologous group*

# Gene annotation with RGI (CARD)



Similar principle: annotate to existing database of functions

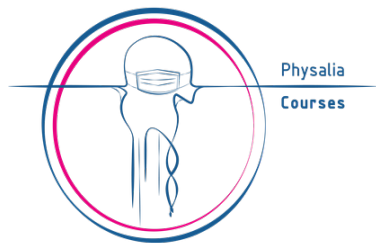# Gene annotation with RGI (CARD): Annotation foibles

**Strict**: a match above threshold
**Loose**: a match below threshold
**Nudged**: a match below threshold, but high identity over a fragment

**Beware the false positives!**

- Efflux pumps!
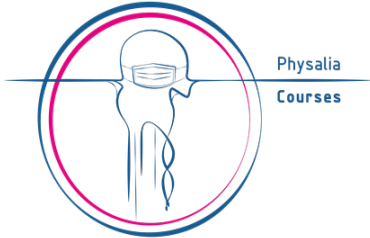- vanR (the regulator) without the regulatees
  - probably not ARG
- …

Physalia
Courses

# There are other specialized databases

- CAZy is very popular
- https://www.cazy.org/

# What about small genes?

Traditionally ignored

Some groups have started to look into them (and [others]).

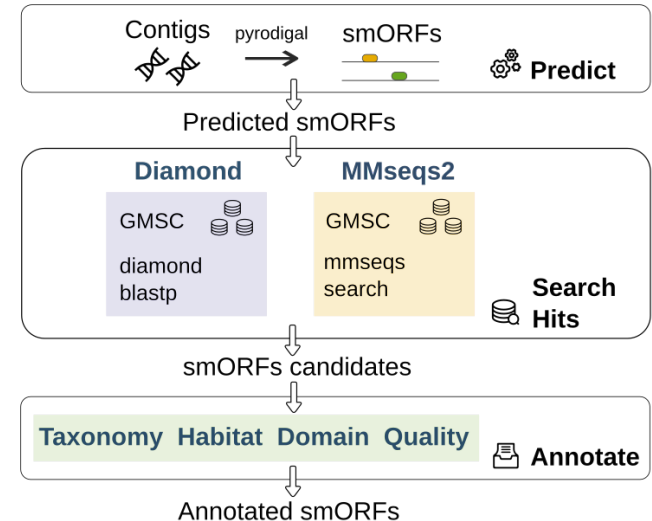| Home | Browse | Downloads | Help |
|------|--------|-----------|------|

## GLOBAL MICROBIAL smORFs CATALOGUE v1.0

The global microbial smORF catalogue (GMSC) is an integrated, consistently-processed, smORFs catalogue of the microbial world, combining publicly avail genomes. A total of non-redundant ~965 million 100AA ORFs were predicted from 63,410 metagenomes across global habitats from the SPIRE database and the ProGenomes2 database. The smORFs were clustered at 90% amino acid identity resulting in ~288 million 90AA smORFs families.

- The annotation of GMSC contains:
  - taxonomy classification
  - habitat assignment
  - quality assessment
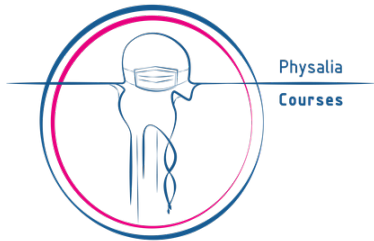  - conserved domain annotation
  - cellular localization prediction

Search from identifier or find homologues by sequence (GMSC-mapper

# Public databases of MAGs

Large-scale genome recovery not submitted to NCBI (so not in GTDB):

- SPIRE (50k new species)

- mOTUs (50k new species)

- SMAG (13k new species)

- GEM (4k new species)

- OceanDNA (2k new species)
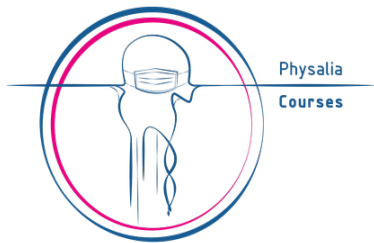
# Public databases of MAGs – GlobDB

Globdb includes all of GTDB and SPIRE, mOTUs, etc.

- 306,260 species (~2x in GTDB)

- Standardised naming, QC,
  gene calling, annotation

- Complete taxonomy (including new clades)

- SingleM metapackage (taxonomic profiling)

- 82 million protein clusters + ProtT5 protein language model embeddings

**GlobDB**

**Welcome to the GlobDB genomes database**

This website hosts the GlobDB, a dereplicated set of species representative microbial genomes. The genomic era offers great opportunities for microbial

Physalia
Courses

# What can you do with these databases?

1. Download their MAGs
2. Compare your MAGs to theirs
   - Dereplication
   - Pangenome
   - Taxonomy (novel GlobDB clades coming soon to Sandpiper!)
3. Check geographic distributions

Physalia
Courses