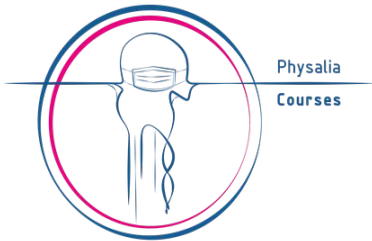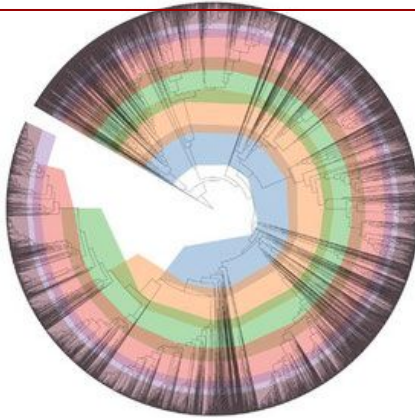# ENVIRONMENTAL METAGENOMICS
Physalia course, online, 11-15 November 2024

## MAG functional annotation

Nikolay Oskolkov, Lund University, NBIS SciLifeLab
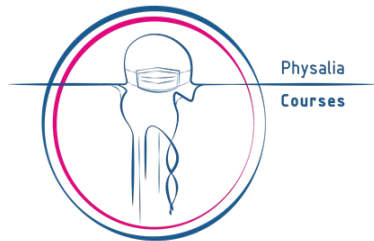Luis Pedro Coelho, Queensland University of Technology

Physalia
Courses

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

**Functional annotation**

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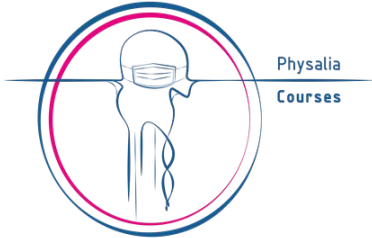- *RGI*
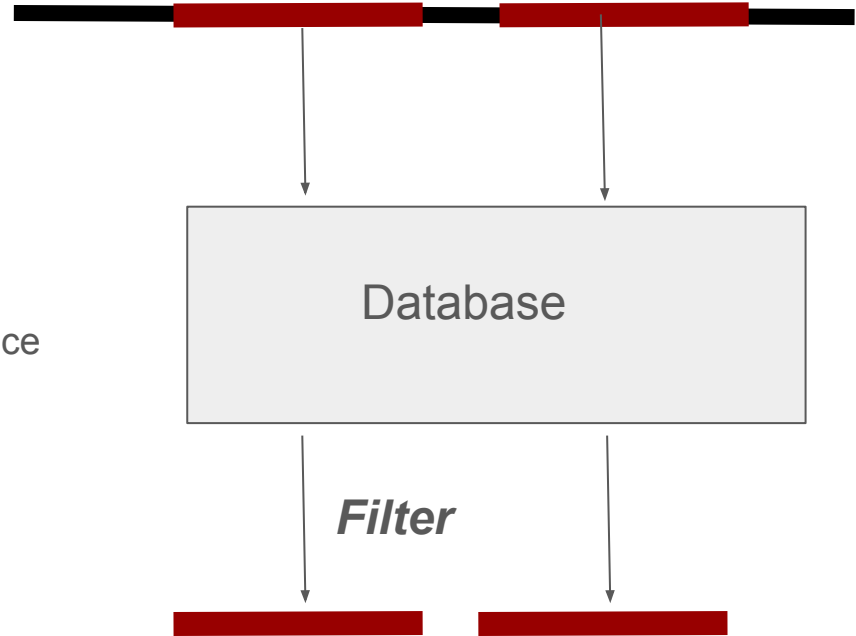- *Other gene annotations databases*
- *…*

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

Physalia
Courses

# 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*

NOG **EggNOG** 6.0.0

Search protein or OG: P53 sapiens, COG1234...

🏠 Home

🔍 Sequence search

🔍 Advance search

☰ Phylogenetic profile

💾 Downloads

ℹ Docs

🔻 eggNOG-mapper v2
(Batch Functional Annotation)

Showing matching OGs

Orthologous Groups

Filter by taxonomy...

**LCOG0787**
(root)

alanine racemase [EC:5.1.1.1],Alr-MurF fusion protein [EC:5.1.1.1 6.3.2.10],amino-acid racemase [EC:5.1.1.10]

12540 pr

Pfam domain   Ala_racemase_N (96.14%),   Ala_racemase_C (93.13%),   Mur_ligase_M (6.48%)

Smart domain   Ala_racemase_C (92.03%),   SIGNAL (1.34%),   TRANS (0.33%)

GO slim   GO:0006520 (65.04%),   GO:0071554 (1.71%),   GO:0006399 (0.67%)

KEGG pathway   map01100 (21.92%),   ko01100 (21.92%),   ko00470 (21.90%)

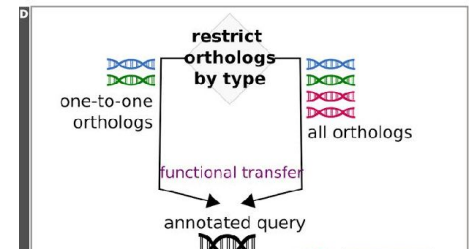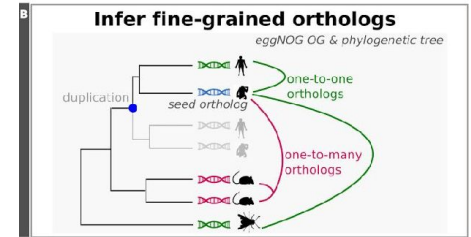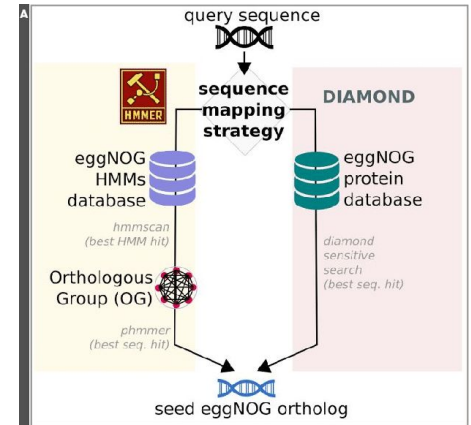KEGG module   M00652 (0.13%),   M00876 (0.02%)

KEGG ortholog   K01775 (20.37%),   K01798 (1.04%),   K25317 (0.46%)

KEGG gene symbol   alr (20.84%),   alr-murF (1.04%),   bsrV (0.46%)

KEGG gene name   alanine racemase [EC:5.1.1.1] (20.37%),   Alr-MurF fusion protein [EC:5.1.1.1 6.3.2.10] (1.04%),   amino-acid (0.46%)

OG members   Taxonomic profile   Functional profile



A

query sequence

**sequence mapping strategy**    DIAMOND

HMMER

eggNOG HMMs database

*hmmscan (best HMM hit)*

Orthologous Group (OG)

*phmmer (best seq. hit)*

eggNOG protein database

*diamond sensitive search (best seq. hit)*

seed eggNOG ortholog

B

**Infer fine-grained orthologs**

eggNOG OG & phylogenetic tree

duplication

seed ortholog

one-to-one orthologs

one-to-many orthologs

C

**Discard distant orthologs**

Taxonomic adjustment across 107 eggNOG levels

D

**restrict orthologs by type**

one-to-one orthologs

all orthologs

functional transfer

annotated query

# Gene annotation with RGI (CARD)



Similar principle: annotate to existing database of functions

# Gene annotation with RGI (CARD): Annotation foibles
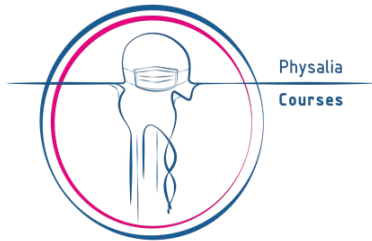
**Strict**: a match above threshold
**Loose**: a match below threshold
**Nudged**: a match below threshold, but high identity over a fragment

**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

We (and others) have started to look into it.

## A catalog of small proteins from the global microbiome

Yiqian Duan, Célio Dias Santos-Júnior, Thomas Sebastian Schmidt, Anthony Fullam, Breno L. S. de Almeida, Chengkai Zhu, Michael Kuhn, Xing-Ming Zhao ✉, Peer Bork & Luis Pedro Coelho ✉

Physalia
**Courses**

# 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 availa... 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

Physalia
Courses

Contigs    pyrodigal    smORFs    Predict

Predicted smORFs

**Diamond**    **MMseqs2**

GMSC    GMSC

diamond    mmseqs
blastp    search    Search
Hits

smORFs candidates

**Taxonomy  Habitat  Domain  Quality**    Annotate

Annotated smORFs

# Public databases of MAGs



## Global Microbial Gene Catalog v1.0

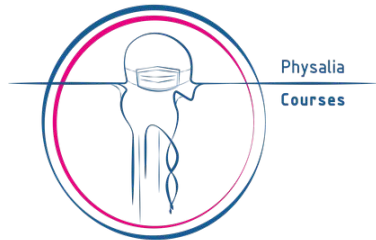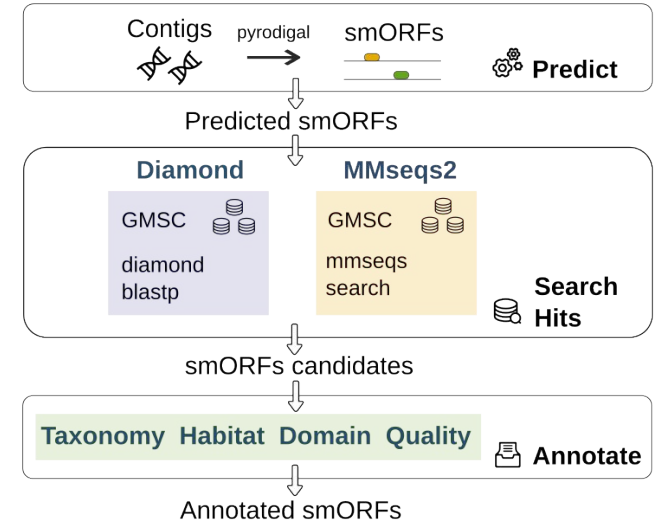The Global Microbial Gene Catalog is an integrated, consistently-processed, gene catalog of the microbial world, combining [...] billion ORFs from 13,174 m[...] database were clustered to[...] unigenes. Read more...

**GlobDB**

Query by sequence or identifier

Find a u[...] (eggNO[...]

Find homologues by sequenc[...]

Physalia Courses

### 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 genome analyses, and individual (meta)genome studies can generate thousands of microbial genomes. Although multiple databases are available to store these datasets, the integration of large scale studies sometimes has proven challengin[...] The GlobDB aims to integrate several resources that are currently not (yet) consolidated otherwise.

SPIRE

## SPIRE

The microbial world at your fingertips.

**S**earchable **P**lanetary-scale m**I**crobiome **RE**source: a one stop shop for microbial data, integrated and consistently processed across habitats and phylogeny, at global scales.

Explore Environments

Explore Taxonomy

# What can you do with these databases?

1. Download their MAGs
2. Compare your MAGs to theirs (dereplication/pangenomes)
3. Check geographic distributions &c

Physalia
Courses