

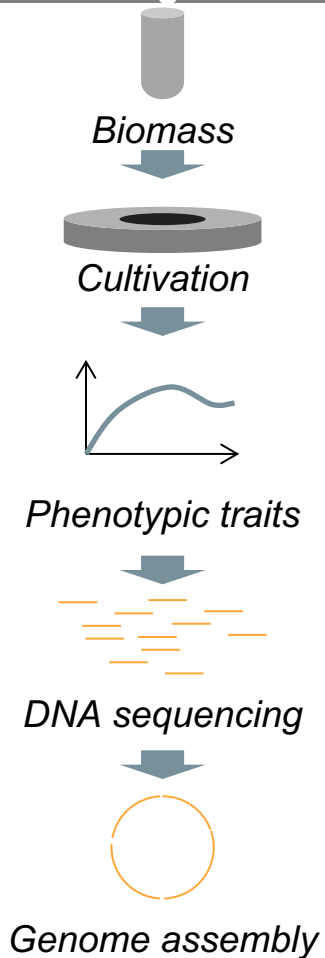
# Computational modeling and learning methods for metagenomics of microbial and viral communities

Thomas Rattei  
Department of Microbiology and Ecosystem Science  
University of Vienna

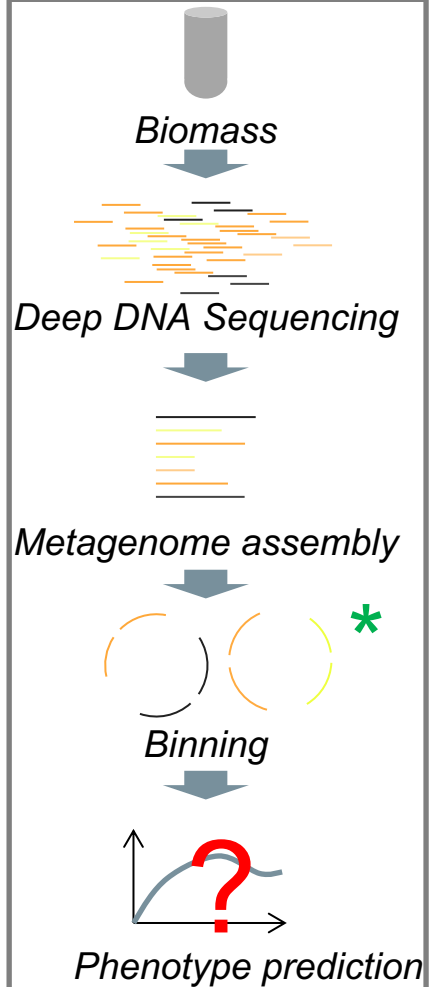
# Microbial communities



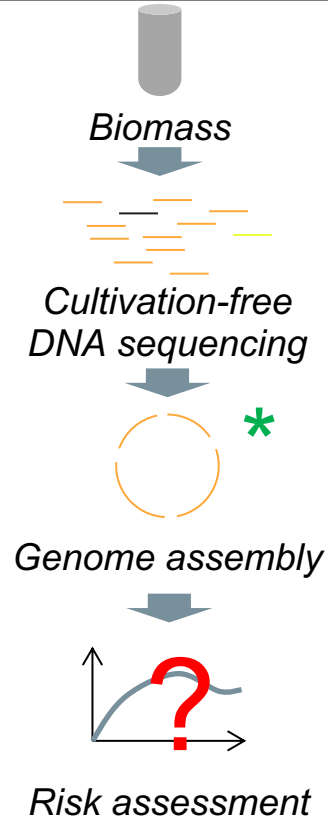
## Isolate genomes



## Metagenomes

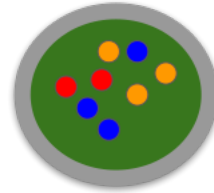


## Targeted genomics



\* Genome standards for Single-Cell genomes (MISAG) and genomes from metagenomes (MIMAG) of bacteria and archaea. Bowers et al., Nat Biotechnol 2017

# PhenDB



Prediction of bacterial phenotypes



Valerie  
Eichinger



Roman  
Feldbauer



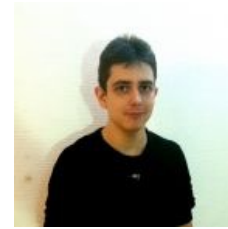
Javier  
Geijo



Patrick  
Hyden



Lukas  
Lüftinger



Florian  
Piewald



Peter  
Peneder

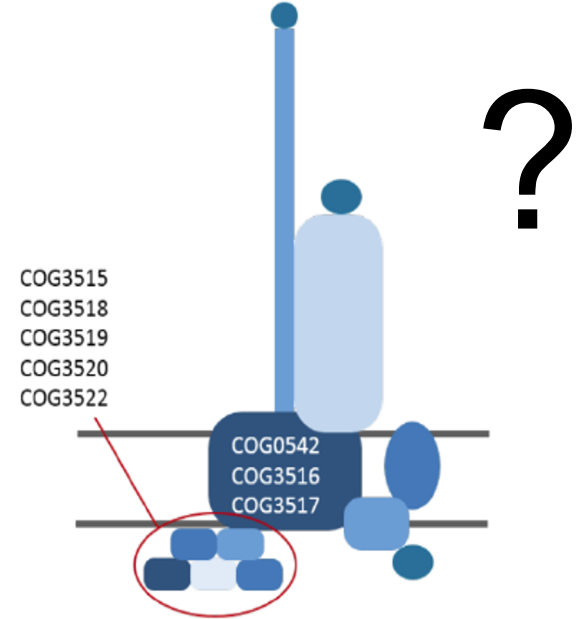
# Rationale for PhenDB

Multiple coverage binning

We can recover **nearly complete** genomes.

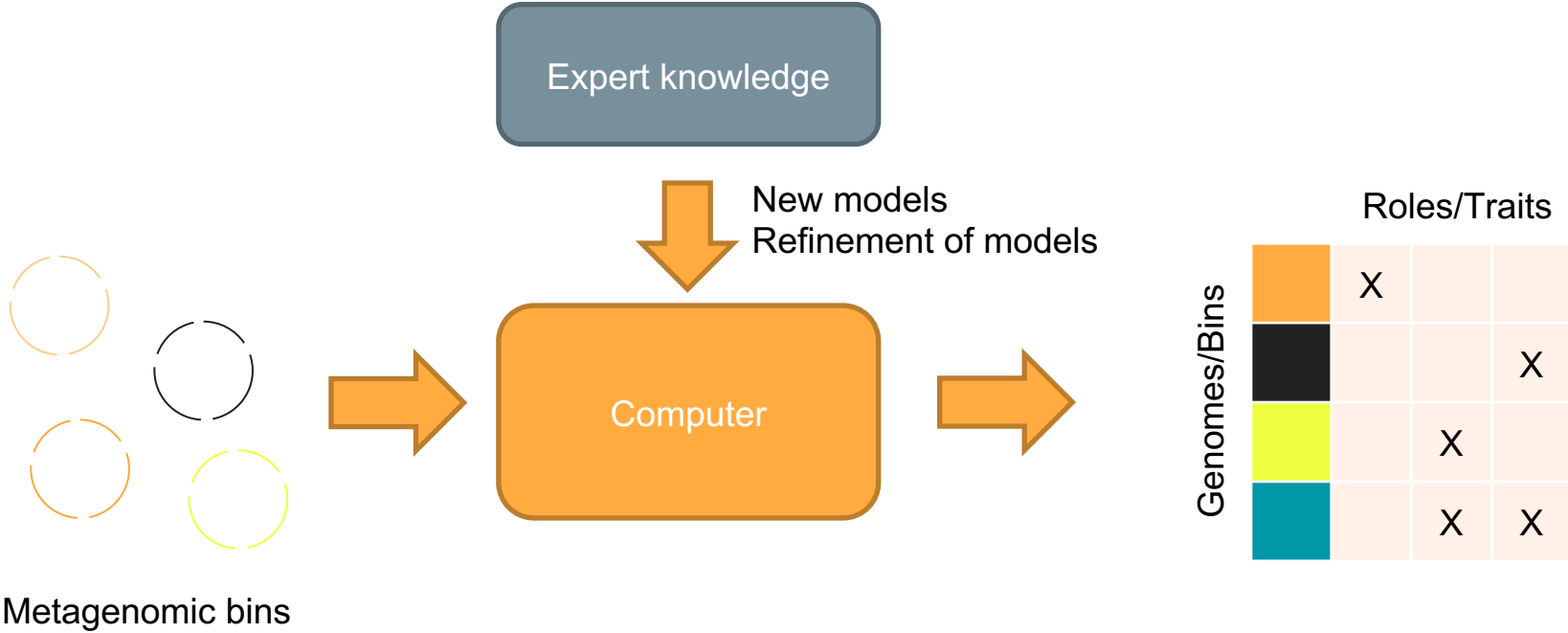
The genome is a **major source of information**.

Lot of genomes are coming from metagenomes. **Speed up data analysis**.

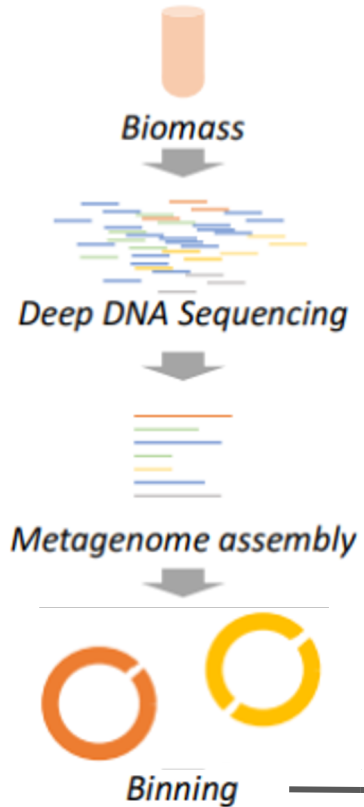



**Is there any method able to infer complex traits from nearly complete genomes?**

# Role and trait prediction in microbial communities



# Trait prediction for nearly complete genomes



| Method  | Features                               | Predictions  | Publication                               |
|---|--|--|---|
| PICA<br> | Protein family occurrence (COG/eggNOG) | Microbial roles, metabolic traits, protein modules, resistance | Feldbauer et al., BMC Bioinformatics 2015 |

Computational method based on PICA.

Learns genotype-phenotype association models from completely sequenced genomes of microbes with known traits.

These models can be used to predict the traits in novel genomes.

**The whole genome architecture can be used to predict the trait.**





# Current PICA models

## Human pathogens

Penicillin resistance of *S. aureus*

## Symbionts

Obligate intracellular lifestyle

Functional T3 Secretion System

Functional T4 Secretion System

Functional T6 Secretion System

## Metabolism

Ammonia-oxidizing bacteria

Nitrite-oxidizing bacteria

Nitrogen fixation

Hydrogen production

Alkane degradation

Arsenic detoxification

Methane utilization as carbon source

## Lifestyle

Autotrophic metabolism

Phototroph bacteria

## Habitat compatibility

Halophile

Psychrophile

Thermophile

Mesophile

Aerobic respiration

Anaerobic respiration

Facultative anaerobe

## Physiology

Motility

Spore formation

Gram negative bacteria

## Metabolism

Plant pathogenicity based on AvrE virulence factor

Secondary bile acids production

Benzoate degradation via hydroxylation

Phytate hydrolysis

Glycine and taurine are de-conjugation from bile salts

Butyrate-production

Chitine degradation

Aromatic-ring-hydroxylation

Carbon monoxide assimilation

Trimethylamine production via choline

Bibenzofuran (DF) dibenzo-p-dioxin (DD) degradation

Steroid hormone metabolism

Naphthalene degradation

Fatty acid degradation

Phosphonate hydrolysis

Recycling of organic phosphorus

Capsular polysaccharide biosynthesis

Thiosulfate oxidation

Plant pathogenicity based on Thaxtomins

Ureolytic activity

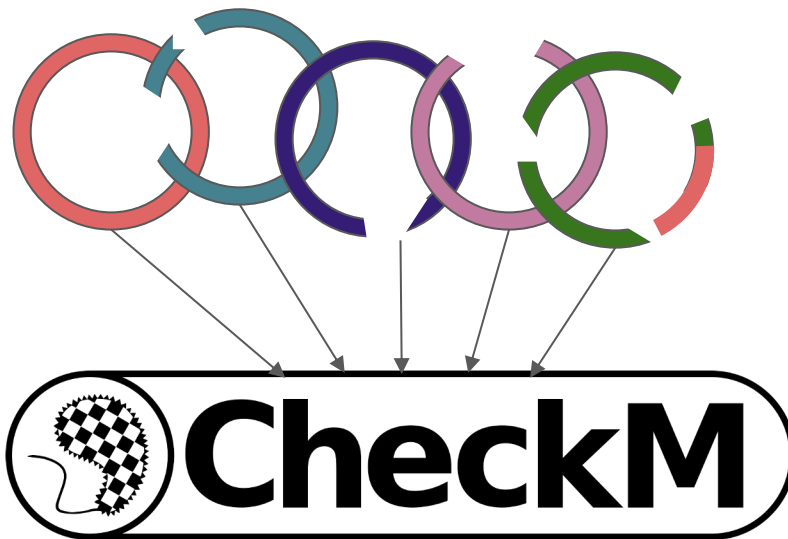
Reduction of various  $\alpha,\beta$ -unsaturated and nitro compounds

# First step: quality assessment

Metagenomic bins

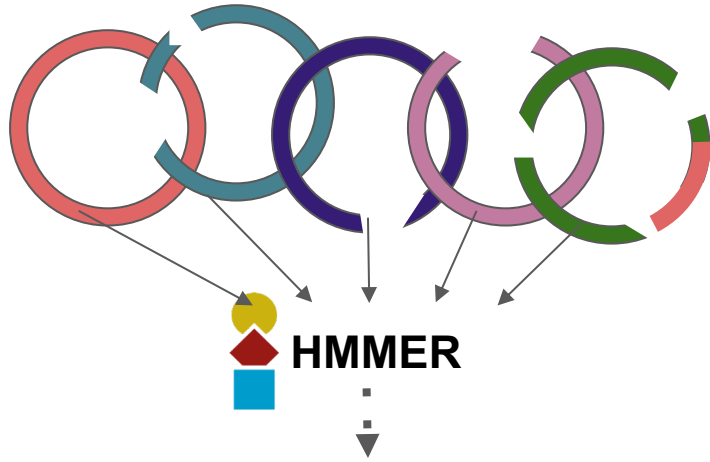
processing ....  
(roughly 10 min per bin)  
requires 20+ GB RAM

Completeness,  
contamination

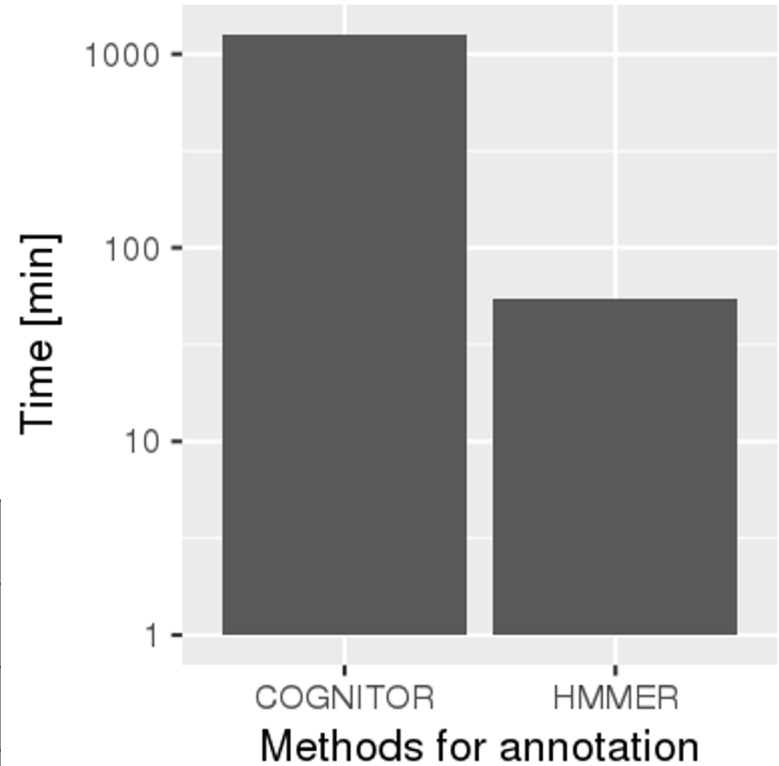


| bin_ID    | 1    | 2    | 3    | 4    | 5    |
|-----------|------|------|------|------|------|
| comp. [%] | 99.9 | 94.5 | 93.6 | 91.2 | 91.3 |
| cont. [%] | 0.0  | 0.0  | 0.0  | 0.0  | 12.2 |

# Second step: annotating genomes

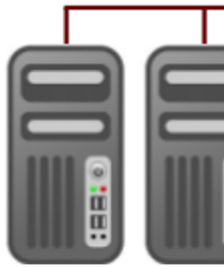


| bin_ID | 1   | 2   | 3   | 4   | 5   |
|--------|-----|-----|-----|-----|-----|
| ENOG1  | 0   | 1   | 0   | 0   | 1   |
| ENOG2  | 1   | 1   | 0   | 0   | 2   |
| ...    | ... | ... | ... | ... | ... |



Computational costs for annotation of an *E. coli* genome with EggNOG database

# Improvement by parallelization



*Bioinformatics*, 34, 2018, i254–i262  
doi: 10.1093/bioinformatics/bty275  
ISMB 2018

OXFORD

---

## DeepFam: deep learning based alignment-free method for protein family modeling and prediction

Seokjun Seo<sup>1</sup>, Minsik Oh<sup>1</sup>, Youngjune Park<sup>2</sup> and Sun Kim<sup>1,2,3,\*</sup>

<sup>1</sup>Department of Computer Science and Engineering and <sup>2</sup>Interdisciplinary Program in Bioinformatics and <sup>3</sup>Bioinformatics Institute, Seoul National University, Seoul 08826, Korea

Further speedup  
Annotation running on 30 compute  
nodes of the LiSC



Methods for annotation

Computational costs for annotation of an *E. coli* genome with EggNOG database

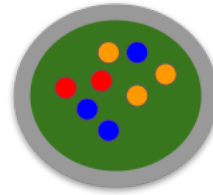
# Pipeline - results and summaries

Individual Results file for each uploaded bin

- For each model:
  - **Verdict: Prediction of PhenDB (YES/NO/NA)**
    - “N/A” if balanced accuracy below **cut-off** (default 0.7)
  - **PICA probability**
    - Depends on found COGs
  - **Balanced accuracy**
    - Depends on Completeness/Contamination and model

|    | A            | B       | C           | D                 |
|----|--------------|---------|-------------|-------------------|
| 1  | Model_name   | Verdict | Probability | Balanced_accuracy |
| 2  | AEROBE       | YES     | 0.77        | 0.87              |
| 3  | ANAEROBE     | NO      | 0.92        | 0.87              |
| 4  | AOB          | NO      | 0.91        | 0.96              |
| 5  | AUTO         | NO      | 0.84        | 0.8               |
| 6  | FACULTATIVE  | NO      | 0.84        | 0.79              |
| 7  | GRAMNEGATIVE | YES     | 0.99        | 0.99              |
| 8  | HALO         | NO      | 0.62        | 0.73              |
| 9  | METHANOTROPH | NO      | 0.77        | 0.84              |
| 10 | MOTILE       | YES     | 0.84        | 0.86              |
| 11 | PEN_180      | NO      | 0.98        | 0.96              |
| 12 | PHOTO        | NO      | 0.9         | 0.92              |
| 13 | PSYCHRO      | N/A     | N/A         | 0.68              |
| 14 | SPORE        | NO      | 0.97        | 0.92              |
| 15 | SYMBIONT     | NO      | 0.98        | 1                 |

# PhenDB

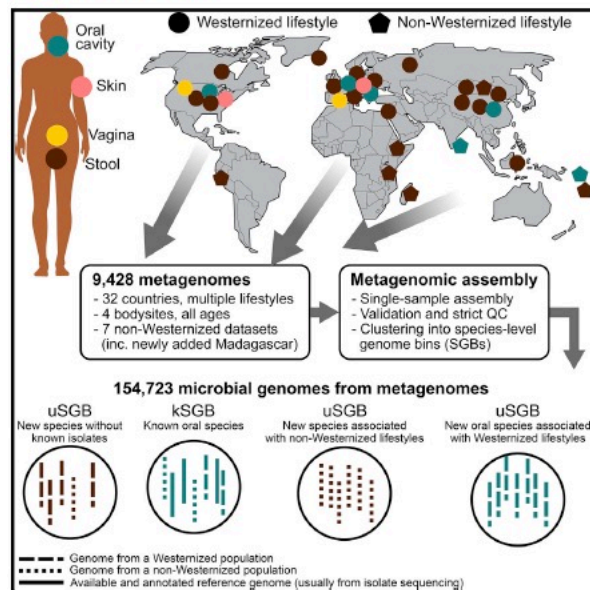


Prediction of bacterial phenotypes

<http://phendb.org>

# Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle

## Graphical Abstract



## Authors

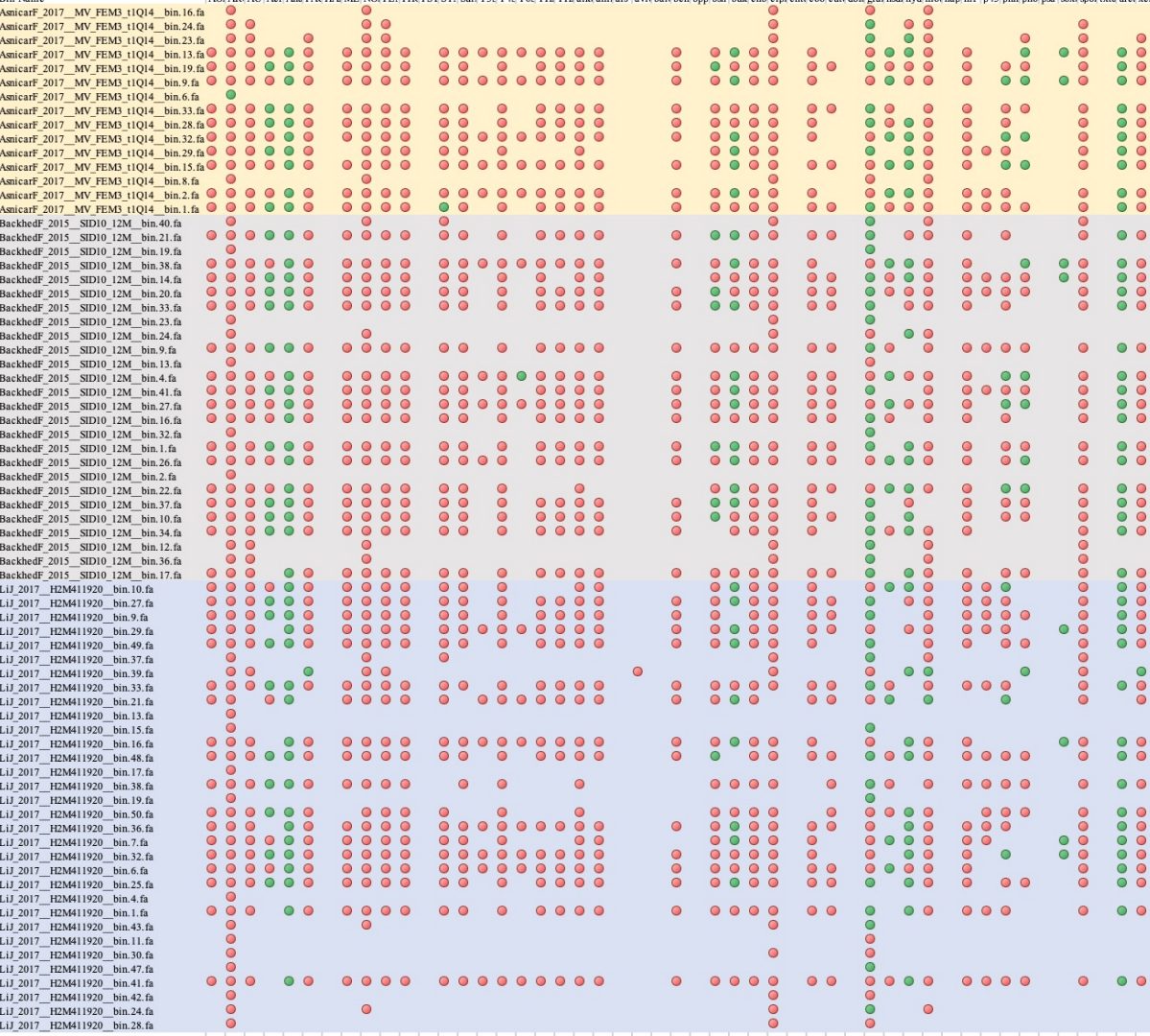
Edoardo Pasoli, Francesco Asnicar, Serena Manara, ..., Christopher Quince, Curtis Huttenhower, Nicola Segata

## Correspondence

nicola.segata@unitn.it

## In Brief

The human microbiome harbors many unidentified species. By large-scale metagenomic assembly of samples from diverse populations, we uncovered >150,000 microbial genomes that are recapitulated in 4,930 species. Many species (77%) were never described before, increase the mappability of metagenomes, and expand our understanding of global body-wide human microbiomes.





# Take home

- Fast prediction of complex traits for large genome-centered metagenomes.
- PhenDB can not predict completely new mechanisms
- We improve models and expand the web interface
- Only available for bacteria.
- Can we quantitatively compare traits?





Hans-Jörg Hellinger



Sabrina Jutz



Jean Mainguy

FWF

Der Wissenschaftsfonds.



Nicole Webster

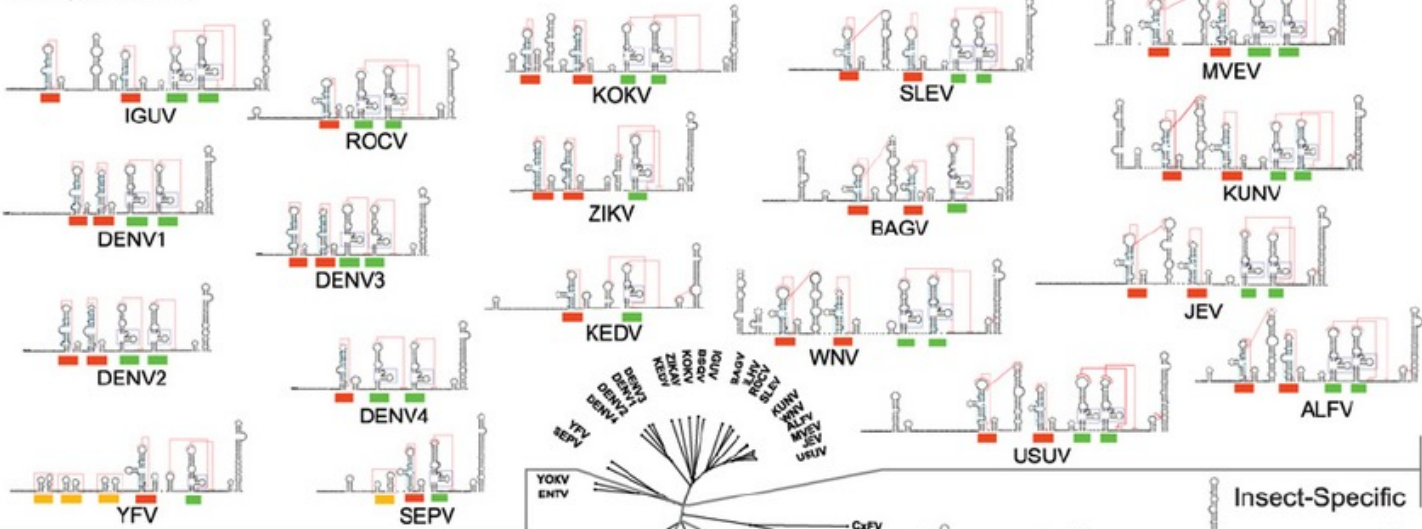
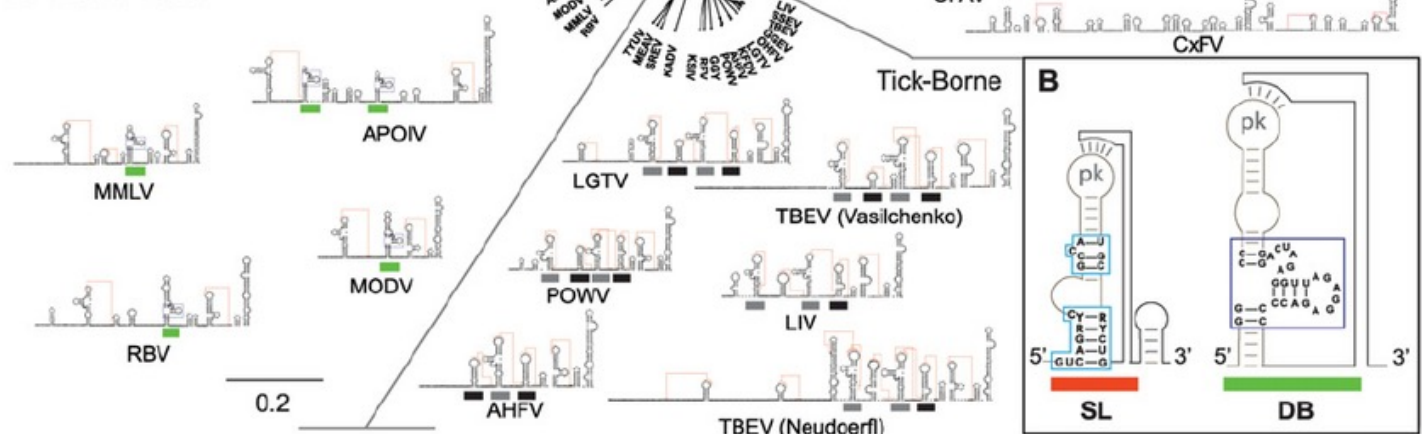
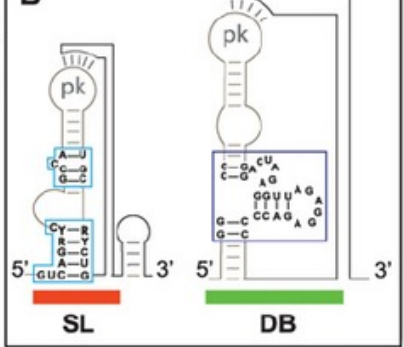


Patrick Laffy



Simon Roux

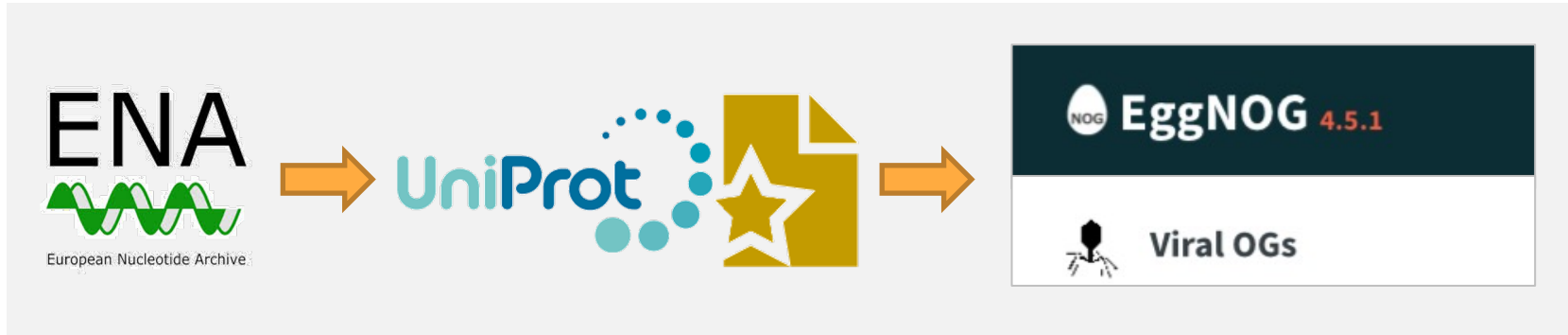


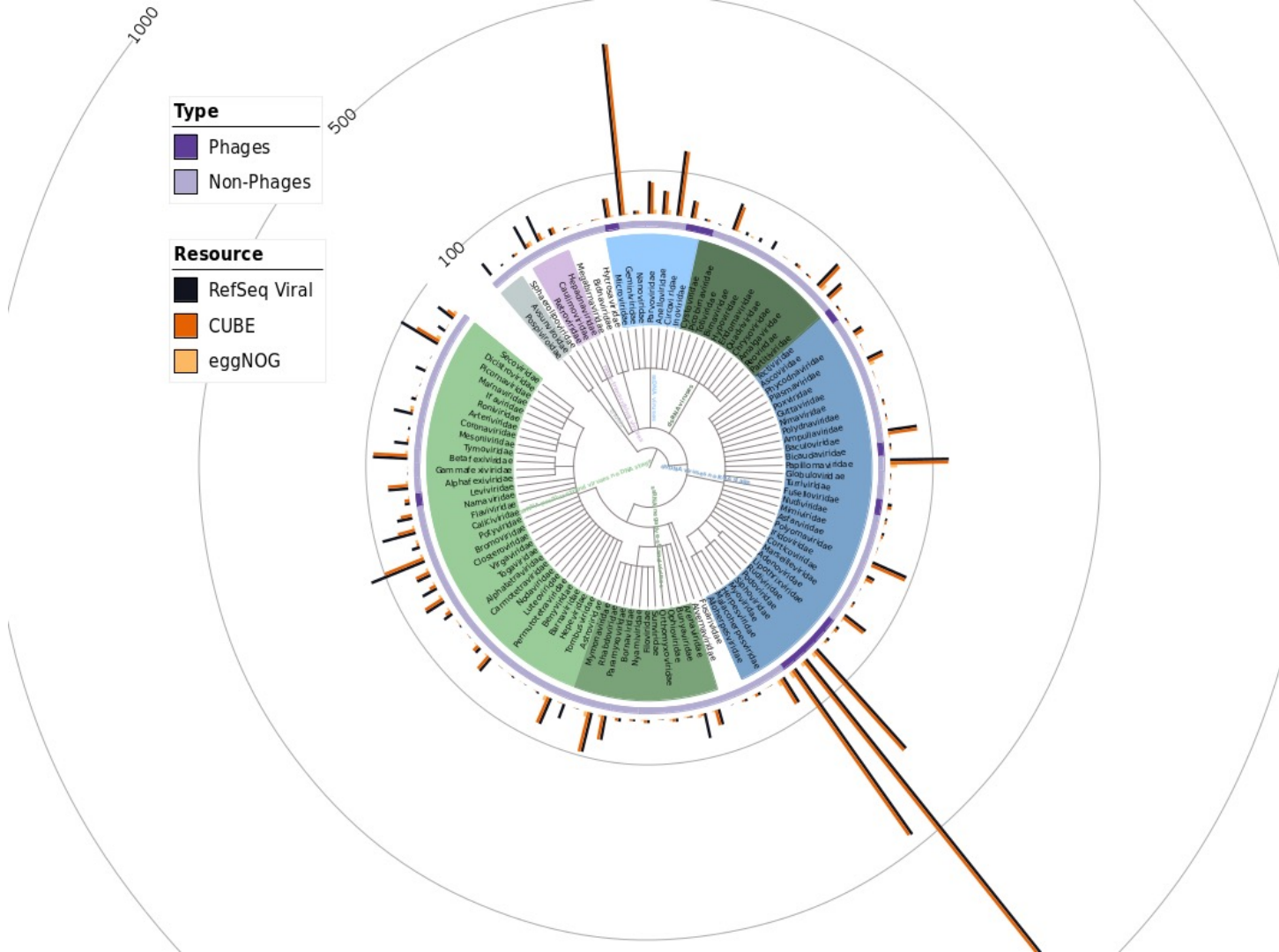
**A****Mosquito-Borne****No-Known Vector****B**

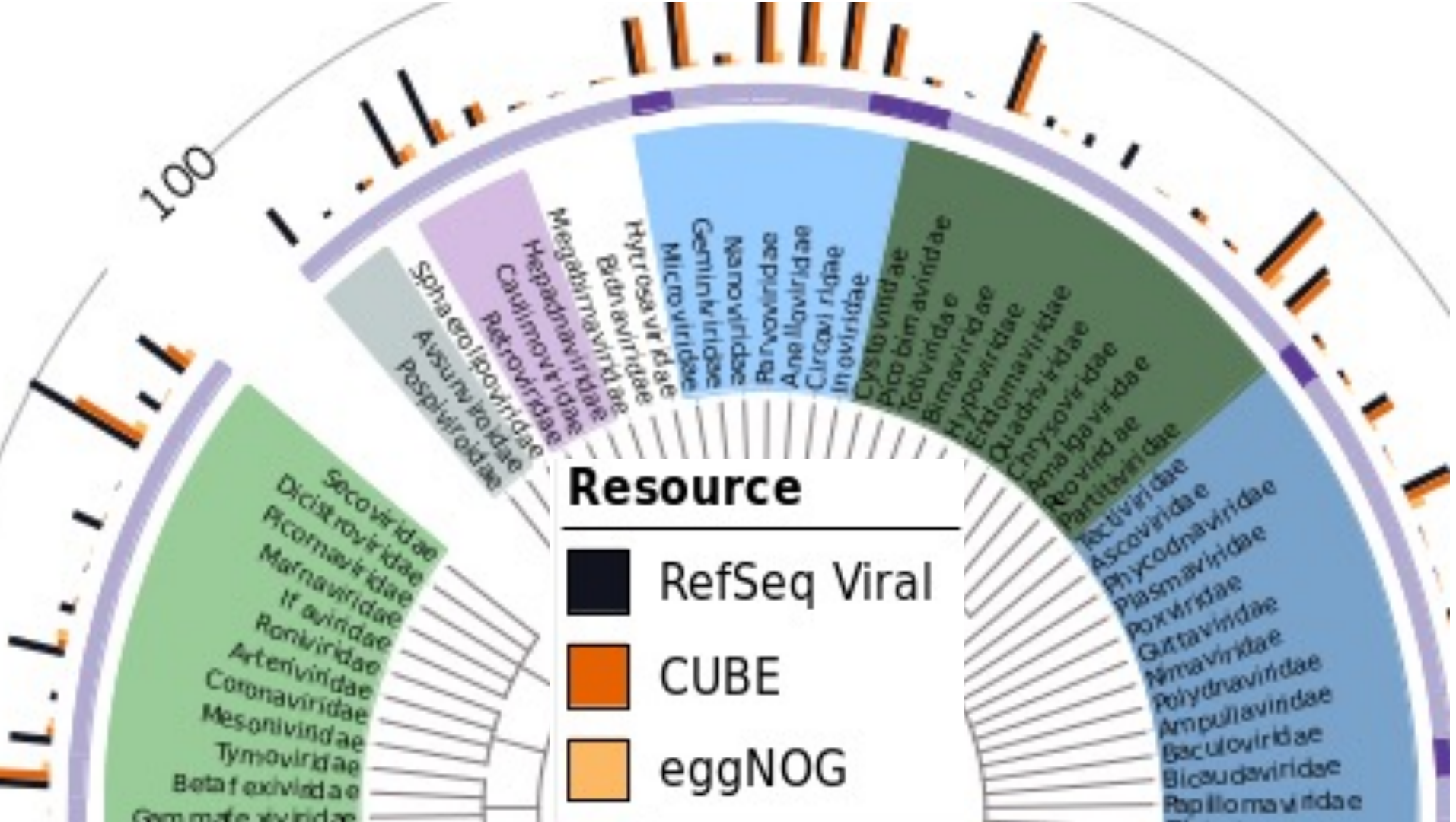
# Catalogs of protein families in viruses



# All-virus orthologous groups

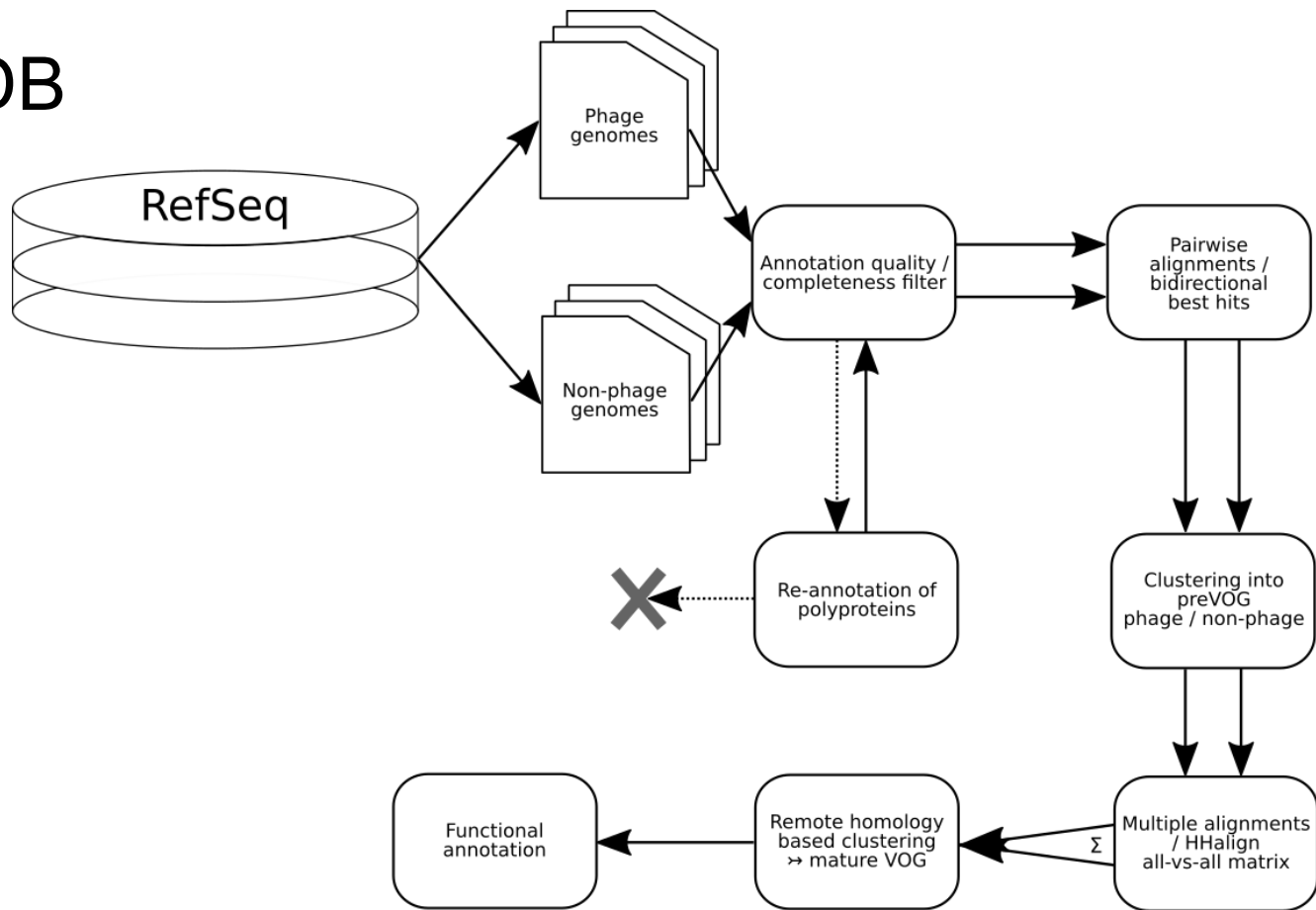








# VOGDB



```
SAKNVHVPYIQASSGFEMWKNNSGRPLQETAPFGCKIAVNP LRAVDCSYIGNIFISIDI  
PNAAFIRTSADPLVSTVKCEVSECTYSADFGGMATLQYVSDREGQCPVHSHSSTATLQ  
ESTVHVLEKGAVTVHFSTAS PQANFIVSLCGKKTTCNAECKPPADHIVSTPHKNDQEF  
QAAISKTSWSWLFALFGGASSLLIIGLMIFACSMMLTSTRR"
```

mat peptide

```
7647..8438  
/locus_tag="SINVgp3"  
/product="capsid (c) protein"  
/protein_id="NP_740673.1"
```

mat peptide

```
8439..8630  
/locus_tag="SINVgp3"  
/product="e-3 structural protein"  
/protein_id="NP_740674.1"
```

mat peptide

```
8631..9899  
/locus_tag="SINVgp3"  
/product="e-2 structural protein"  
/protein_id="NP_740675.1"
```

mat peptide

```
9100..11064  
/locus_tag="SINVgp3"  
/product="6k structural protein"  
/protein_id="NP_740676.1"
```

mat peptide

```
10065..11381  
/locus_tag="SINVgp3"  
/product="e-1 structural protein"  
/protein_id="NP_740677.1"
```

gene

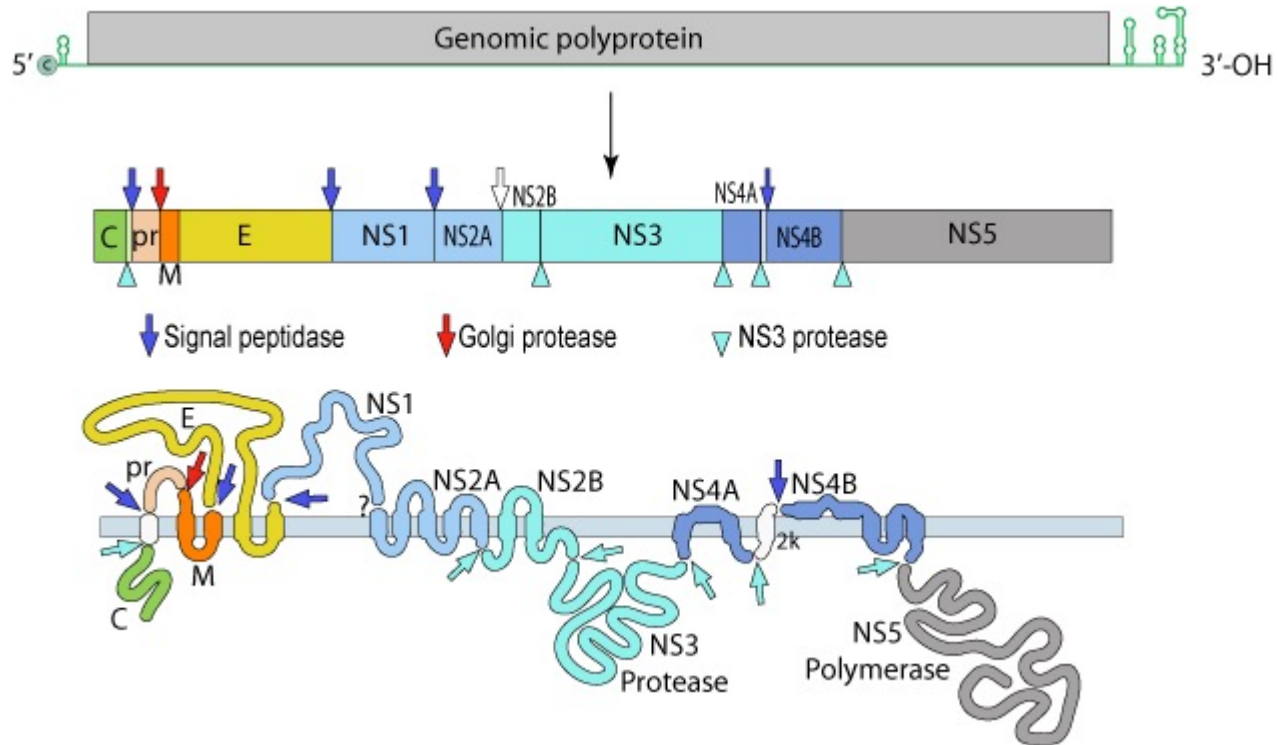
```
7647..10111  
/locus_tag="SINVgp4"  
/db_xref="GeneID:13165406"
```

CDS

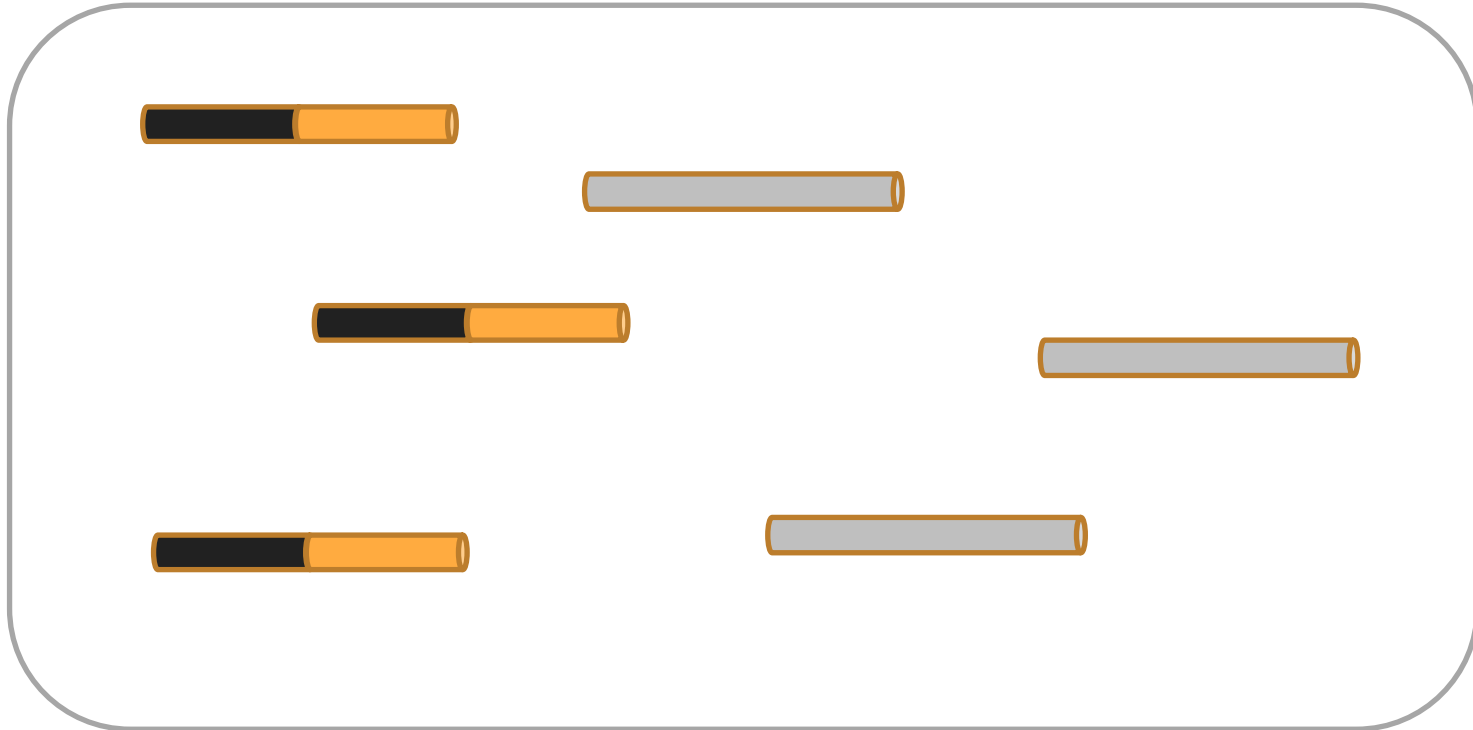
```
join(7647..10028,10028..10111)  
/locus_tag="SINVgp4"  
/note="Truncated version of structural polyprotein that  
will be produced when frameshifting occurs at nt 10028;  
The amino acid sequence of the protein is identical to that of
```

# Polypeptides and their peptides

# A typical flavivirus genome

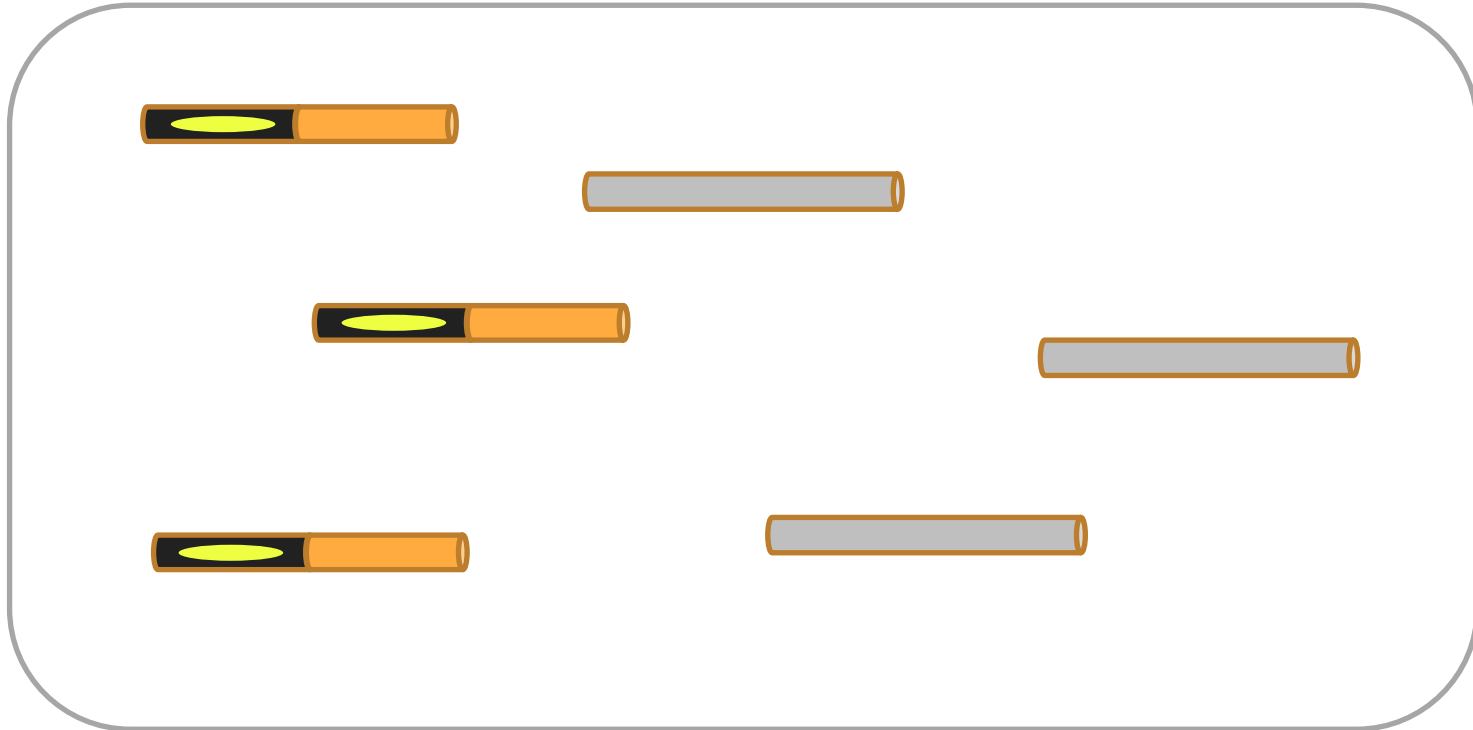


# Automatic polyprotein annotation: concept



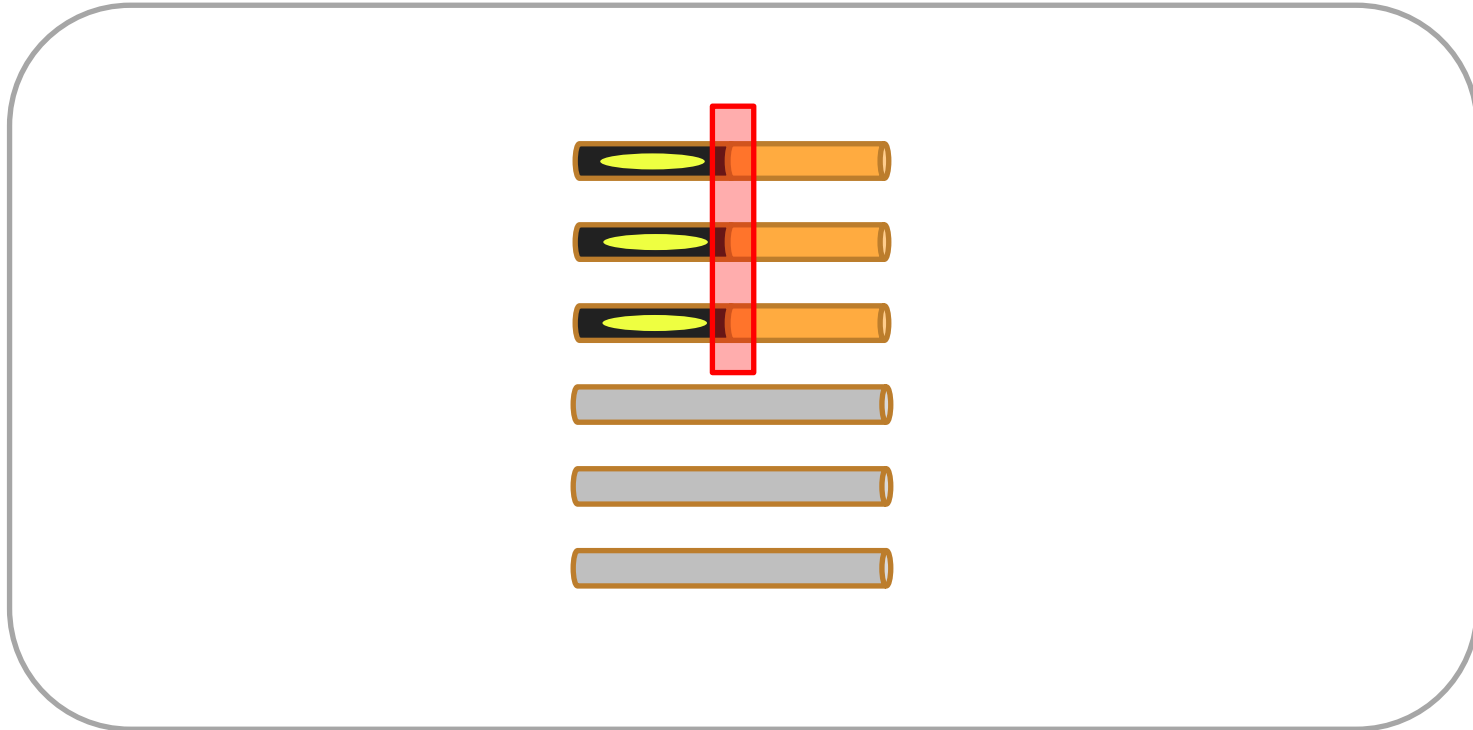
1. Extraction of homologous polyproteins from a genus or a homologous cluster

# Automatic polyprotein annotation: concept



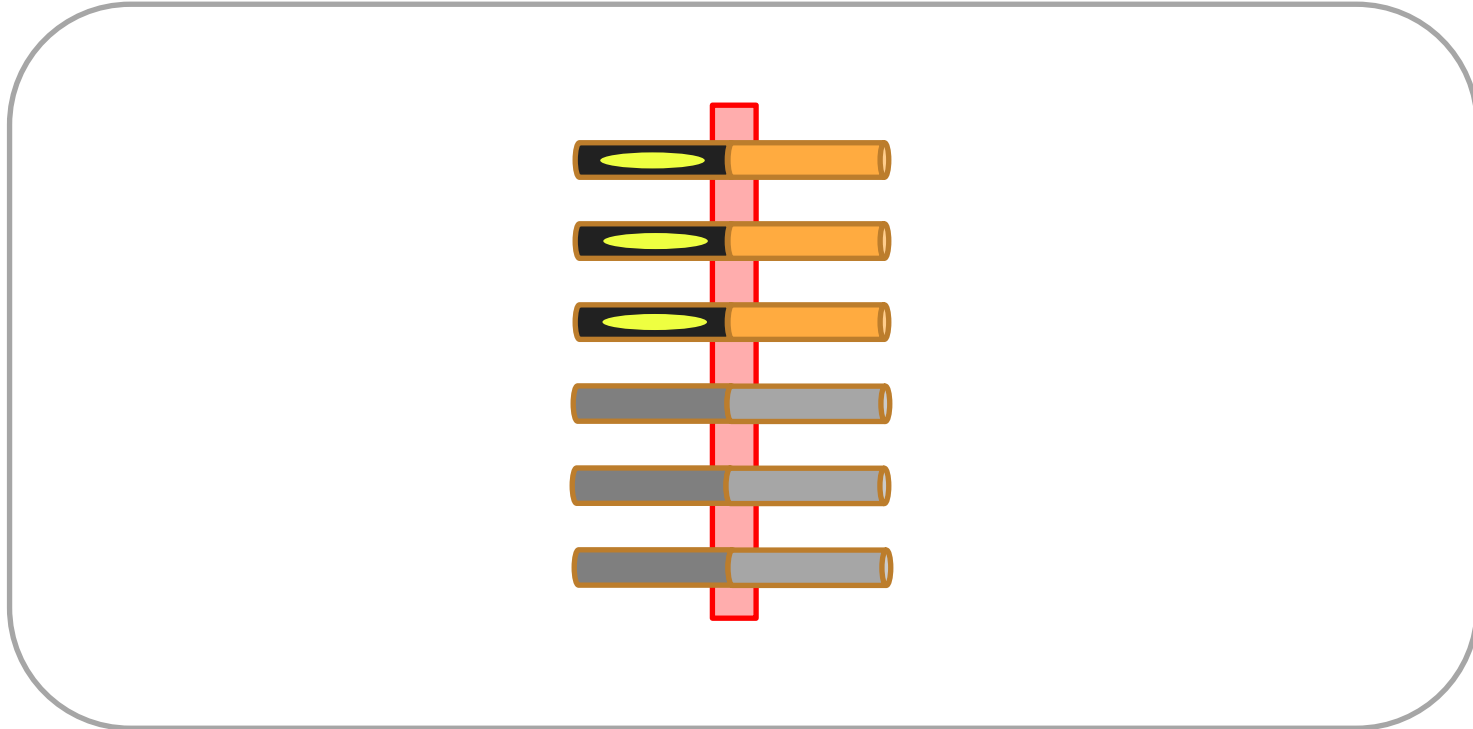
2. Validation of peptide annotation with domains

# Automatic polyprotein annotation: concept



3. Multiple alignment, validation of conservation of sequence and cleavage sites

# Automatic polyprotein annotation: concept

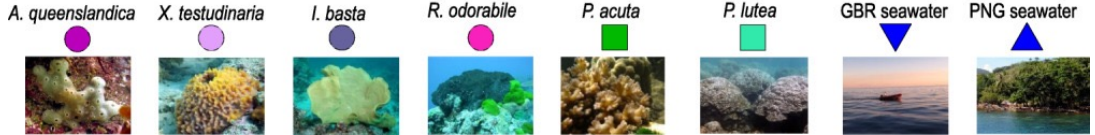


4. Propagation of cleavage sites to unannotated proteins

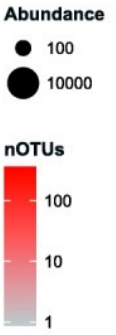
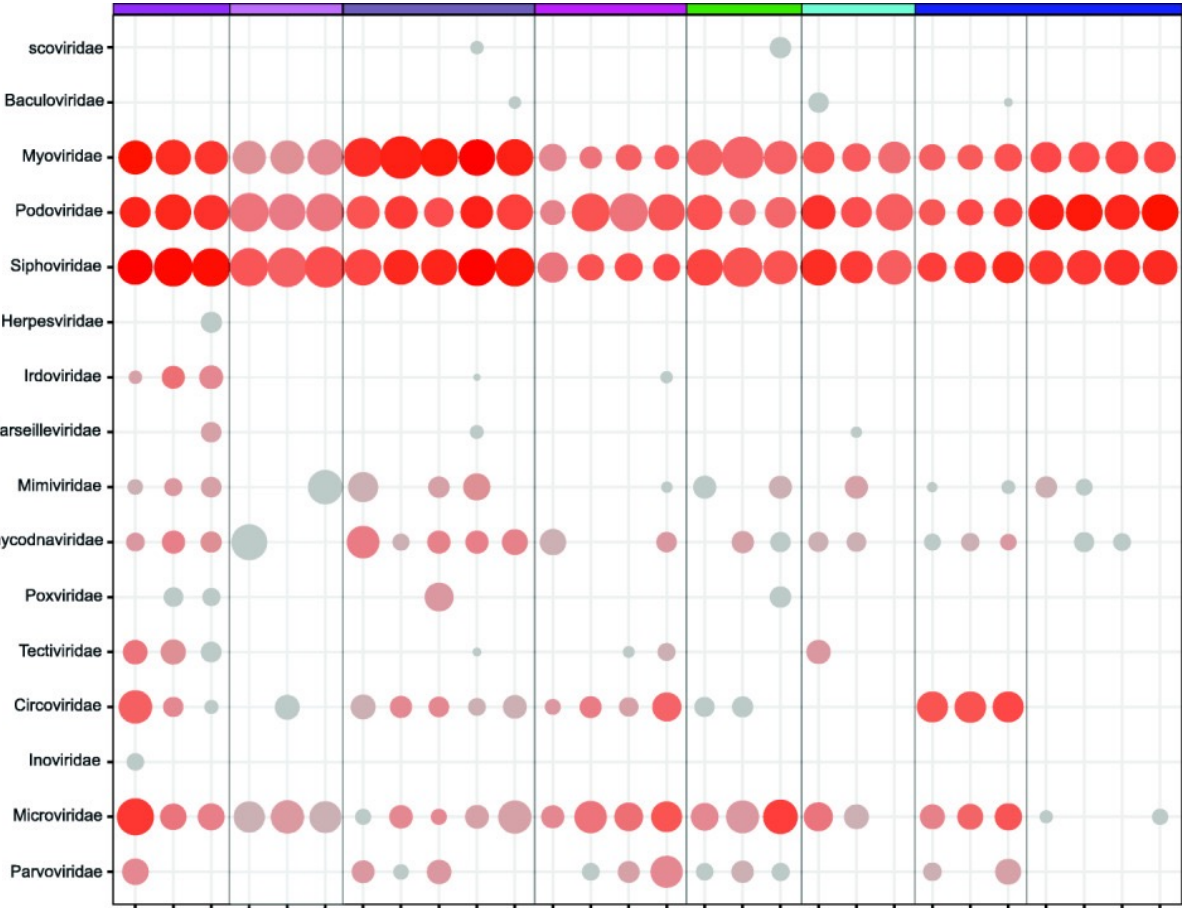


<http://vogdb.org>





Composition of viruses across holobiont species and reef environments



# Take home

- Grouping of all viral proteins into orthologous groups and families.
- Virus genomes require annotation check and re-annotation
- Many families are present in prokaryotic and eukaryotic viruses
- Applications e.g. in viral ecology, virus evolution, study of host interactions...
- Lots of features in database and web interface to come