ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
UNIVERSITY OF CRETE

# Microbial communities through the lens of high throughput sequencing, data integration and metabolic networks analysis

Haris Zafeiropoulos

Dissertation presented in partial fulfillment of the requirements for the degree of Doctor of Science (PhD) in Biology

**Promotors:**
Prof. Emmanouil Ladoukakis
Dr Evangelos Pafilis
Dr Christoforos Nikolaou

Academic year 2021 – 2022

# Members of the examination committee
# &
# reading committee

**Prof. Emmanouil Ladoukakis**
Univeristy of Crete
Biology Department

**Dr Evangelos Pafilis**
Hellenic Centre for Marine Research
Institute of Marine Biology, Biotechnology and Aquaculture

**Dr Christoforos Nikolaou**
Biomedical Sciences Research Center "Alexander Fleming"
Institute of Bioinnovation

**Dr Jens Carlsson**
University College Dublin
School of Biology and Environmental Science/Earth Institute

**Here are some thoughts of mine for the rest of the committe!**
I am considering of **Prof Faust** that I will join her lab for a few months (for an EMBO short term fellowsip)
Also, **Prof. Elias Tsigaridas** that we have worked together on flux sampling.
Finally, in case that it is ok to have non permanent researchers in the committe, **Dr. Christina Pavloudi** with whom I have worked all these years.

# Preface

*Haris Zafeiropoulos*

# Contents

# Abstract

# Περίληψη

Και στα ελληνικά

# List of Figures and Tables

## List of Figures

## List of Tables

# List of Abbreviations and Symbols

## Abbreviations

NGS      Next Generation Sequencing
HPC      High Performance Computing
MCMC    Markov Chain Monte Carlo
MMCS    Multiphase Monte Carlo Sampling
PREGO    PRocess Environment OrGanism
PEMA     Pipeline for Environmental DNA Metabarcoding Analysis
DARN     Dark mAtteR iNvestigator

# Chapter 1

# Introduction

## 1.1 Microbial communities: structure & function

Microbes, i.e. Bacteria, Archaea and small Eukaryotes such as protoza, are omnipresent and impact global ecosystem functions [1] through their abundance [2], versatility [3] and interactions [4].

### 1.1.1 The role of microbial communities in biogeochemical cycles

### 1.1.2 Microbial interactions: unravelling the microbiome

## 1.2 Bioinformatics challenges in HTS approaches

## 1.3 Data integration & data mining in the era of omics

## 1.4 Sampling the flux space of a metabolic model: challenges & potential

## 1.5 The hypersaline Tristomo swamp: a case study of an extreme environment

## 1.6 Systems biology from a computational resources point-of-view

## 1.7 Aims and objectives

The aim of this PhD was double; on the one hand, to enhance the analysis of microbiome data by building algorithms and software to address some of the on-going computational challenges on the field. On the other, to exploit these methods to identify taxa, functions, especially related to sulfur cycle, and microbial interactions that support life in

microbial community assemblages in hypersaline sediments. All parts of this work are computational.

In **Chapter 2**, challenges derived from the analysis of HTS amplicon data are examined. A bioinformatics pipeline, called pema, for the analysis of several marker genes was developed, combinining several new technologies that allow large scale analysis of hundreds of samples. In addition, a software tool called darn, was built to investigate the unassigned sequences in amplicon data of the COI marker gene.

In **Chapter 3**, data integration, data mining and text-mining methods were exploited to build a knowledge-base, called prego, including millions of associations between:

1. microbial taxa and the environments they have been found in

2. microbial taxa and biological processes they occur

3. environmental types and the biological processes that take place there

In **Chapter 4**, the challenges of flux sampling in metabolic models of high dimensions was presented along with a Multiphase Monte Carlo Sampling (MMCS) algorithm we developed.

In **Chapter 5**, sediment samples from a hypersaline swamp in Tristomo, Karpathos Greece were analysed using both amplicon and shotgun metagenomics. The taxonomic and the functional profiles of the microbial communities present there were investigated. Key microbial interactions for the assemblages were infered. All the methods developed and presented in the previous chapters were used to enhance the analysis of this microbiome.

In **Chapter 6**, the history of the IMBBC-HCMR HPC facility was presented indicating the vast needs of computing resources in modern analyses in general and in microbial studies more specifically.

Finally, in the **Conclusions** chapter, general discussion and conclusions that have derived from this research were presented.

# Chapter 2

# Software development to establish quality HTS-oriented bioinformatics methods for microbial diversity assessment

## 2.1 Environmental DNA metabarcoding: challenges and caveats

## 2.2 PEMA: a flexible Pipeline for Environmental DNA Metabarcoding Analysis of the 16S/18S ribosomal RNA, ITS, and COI marker genes

Publication relative to this chapter: [5].



FIGURE 2.1: The PEMA workflow: figure from publication

## 2.3 The Dark mAtteR iNvestigator (DARN) tool: getting to know the known unknowns in COI amplicon data

Publication relative to this chapter: [6]



FIGURE 2.2: DARN methodology: figure in the publication

## 2.4 A workflow for marine Genomic Observatories data analysis

# Chapter 3

# Software development to build a knowledge-base at the systems biology level

**3.1  Metadata: a key issue for robust metanalyses**

**3.2  Ontologies & databases: the corner stone of mordern biology**

## 3.3 PREGO: a literature- and data-mining resource to associate microorganisms, biological processes, and environment types

Publication relative to this chapter: under submission



FIGURE 3.1: PREGO methodology: figure in the publication under submission

# Chapter 4

# Software development to establish metabolic flux sampling approaches at the community level

## 4.1 Genome-scale metabolic model analysis

## 4.2 A New MCMC Algorithm for Sampling the Flux Space of Metabolic Networks

Publication relative to this chapter: [7]



FIGURE 4.1: Our MMCS algorithm and its first phases. Figure published on SoCG21

## 4.3 Flux sampling at the community level

# Chapter 5

# Microbial interactions inference in communities of a hypersaline swamp elucidate mechanisms governing taxonomic & functional profiles

Publication relative to this chapter: ongoing work, to be submitted before phd defense, probably not accepted by then though.

## 5.1 Amplicon & shotgun metagenomic analysis

darn and PEMA will be used at this point, among other software

## 5.2 Inferring microbial interactions

PREGO and dingo will be used to this end

# Chapter 6

# 0s and 1s in marine molecular research

Publication relative to this chapter: [8]

## 6.1 Computing resources: a prerequisite & a limitation in modern microbial ecology

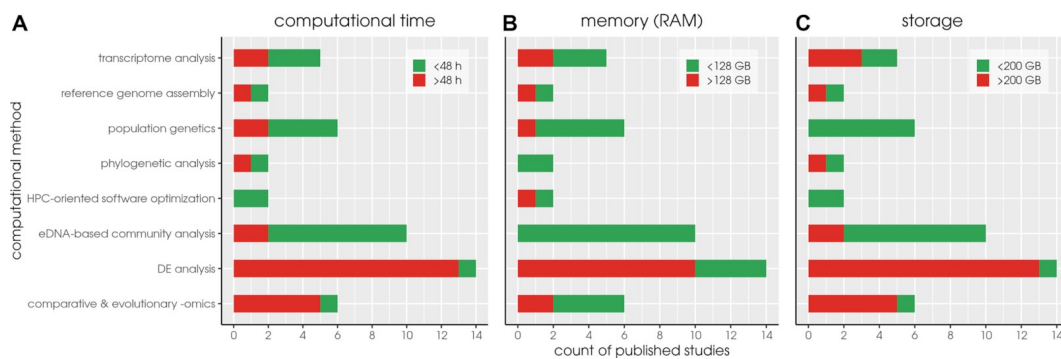## 6.2 High Performance Computing and Cloudification: scaling up bioinformatics analysis

FIGURE 6.1: Computing requirements of the published studies performed on the IMBBC HPC facility over the last decade. Figure from publication.

# Chapter 7

# Conclusions

# Appendices

# Bibliography

[1] P. G. Falkowski, T. Fenchel, and E. F. Delong, "The microbial engines that drive earth's biogeochemical cycles," *science*, vol. 320, no. 5879, pp. 1034–1039, 2008.

[2] Y. M. Bar-On, R. Phillips, and R. Milo, "The biomass distribution on earth," *Proceedings of the National Academy of Sciences*, vol. 115, no. 25, pp. 6506–6511, 2018.

[3] H. A. Rees, A. C. Komor, W.-H. Yeh, J. Caetano-Lopes, M. Warman, A. S. Edge, and D. R. Liu, "Improving the dna specificity and applicability of base editing through protein engineering and protein delivery," *Nature communications*, vol. 8, no. 1, pp. 1–10, 2017.

[4] L. Röttjers and K. Faust, "From hairballs to hypotheses–biological insights from microbial networks," *FEMS microbiology reviews*, vol. 42, no. 6, pp. 761–780, 2018.

[5] H. Zafeiropoulos, H. Q. Viet, K. Vasileiadou, A. Potirakis, C. Arvanitidis, P. Topalis, C. Pavloudi, and E. Pafilis, "Pema: a flexible pipeline for environmental dna metabarcoding analysis of the 16s/18s ribosomal rna, its, and coi marker genes," *GigaScience*, vol. 9, no. 3, p. giaa022, 2020.

[6] H. Zafeiropoulos, L. Gargan, S. Hintikka, C. Pavloudi, and J. Carlsson, "The dark matter investigator (darn) tool: getting to know the known unknowns in coi amplicon data," *Metabarcoding and Metagenomics*, vol. 5, p. e69657, 2021.

[7] A. Chalkis, V. Fisikopoulos, E. Tsigaridas, and H. Zafeiropoulos, "Geometric Algorithms for Sampling the Flux Space of Metabolic Networks," in *37th International Symposium on Computational Geometry (SoCG 2021)* (K. Buchin and E. Colin de Verdière, eds.), vol. 189 of *Leibniz International Proceedings in Informatics (LIPIcs)*, (Dagstuhl, Germany), pp. 21:1–21:16, Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2021.

[8] H. Zafeiropoulos, A. Gioti, S. Ninidakis, A. Potirakis, S. Paragkamian, N. Angelova, A. Antoniou, T. Danis, E. Kaitetzidou, P. Kasapidis, *et al.*, "0s and 1s in marine molecular research: a regional hpc perspective," *GigaScience*, vol. 10, no. 8, p. giab053, 2021.