

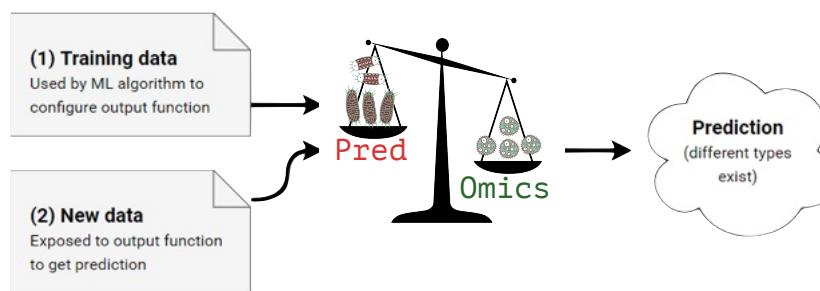
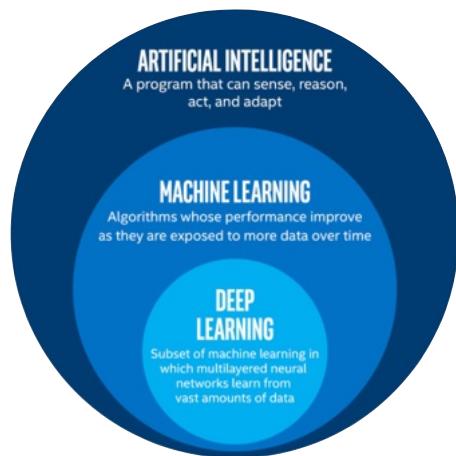


Harnessing AI in Precision Health: Interpretable AI for Microbiome Applications

Edi Prifti. PhD, DR

Institut de Recherche pour le Développement

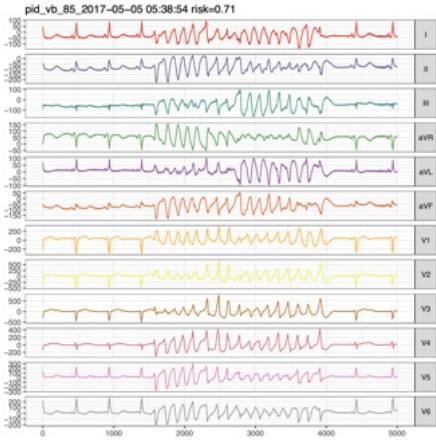
UMI 209 UMMISCO



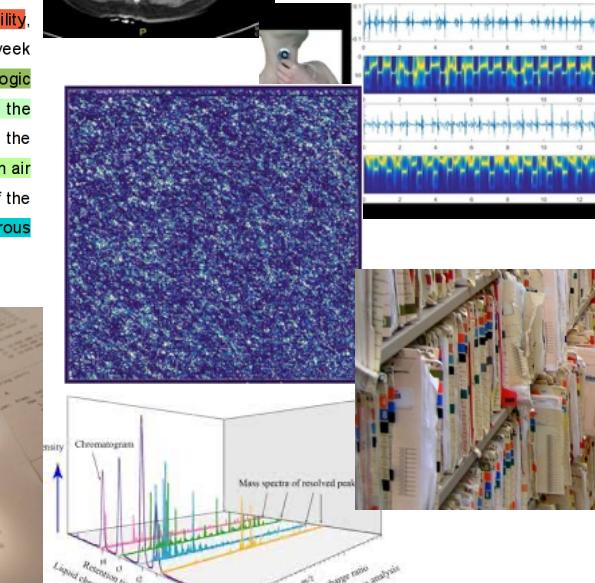
Application fields and the data

Health

biosignals (ECG, EEG), high throughput sequencing (genomics, metagenomics, barcoding, RNASeq), CT/IRM 3D, text

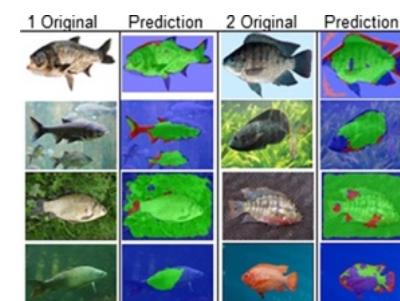
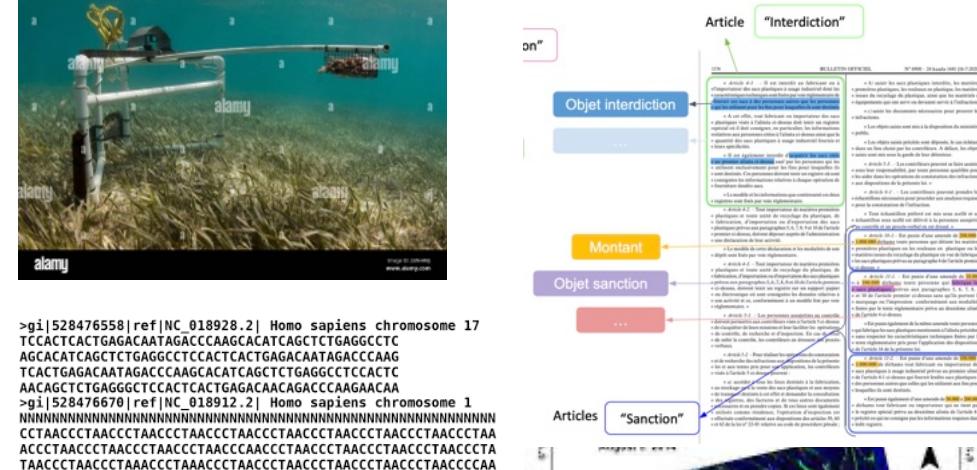


A 12-year old girl with known hyperagglutinability, presented to the emergency department with a 2-week history of headaches and facial weakness. Neurologic examination indicated sensorineural hearing loss on the right side with Weber's test lateralizing to the left, and the Rinne's test demonstrating bone conduction greater than air conduction on the right. Magnetic resonance imaging of the head revealed severe structural defects of the right petrous temporal bone. No indication of cerebral infarction.

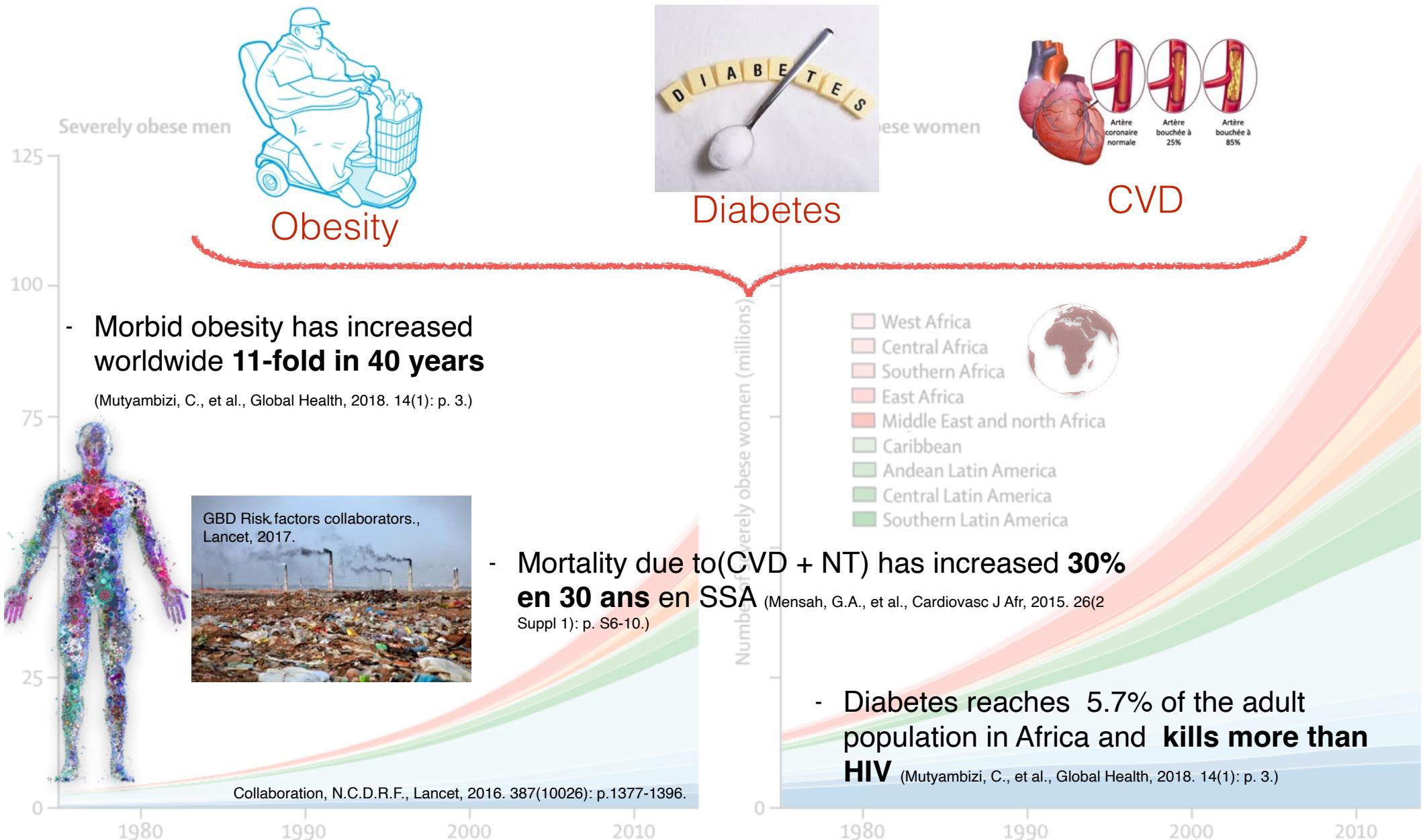


Environnement

Images, scans 2D, 3D, metabarcoding, satellite plant distribution, texte, video, populations, ecosystems

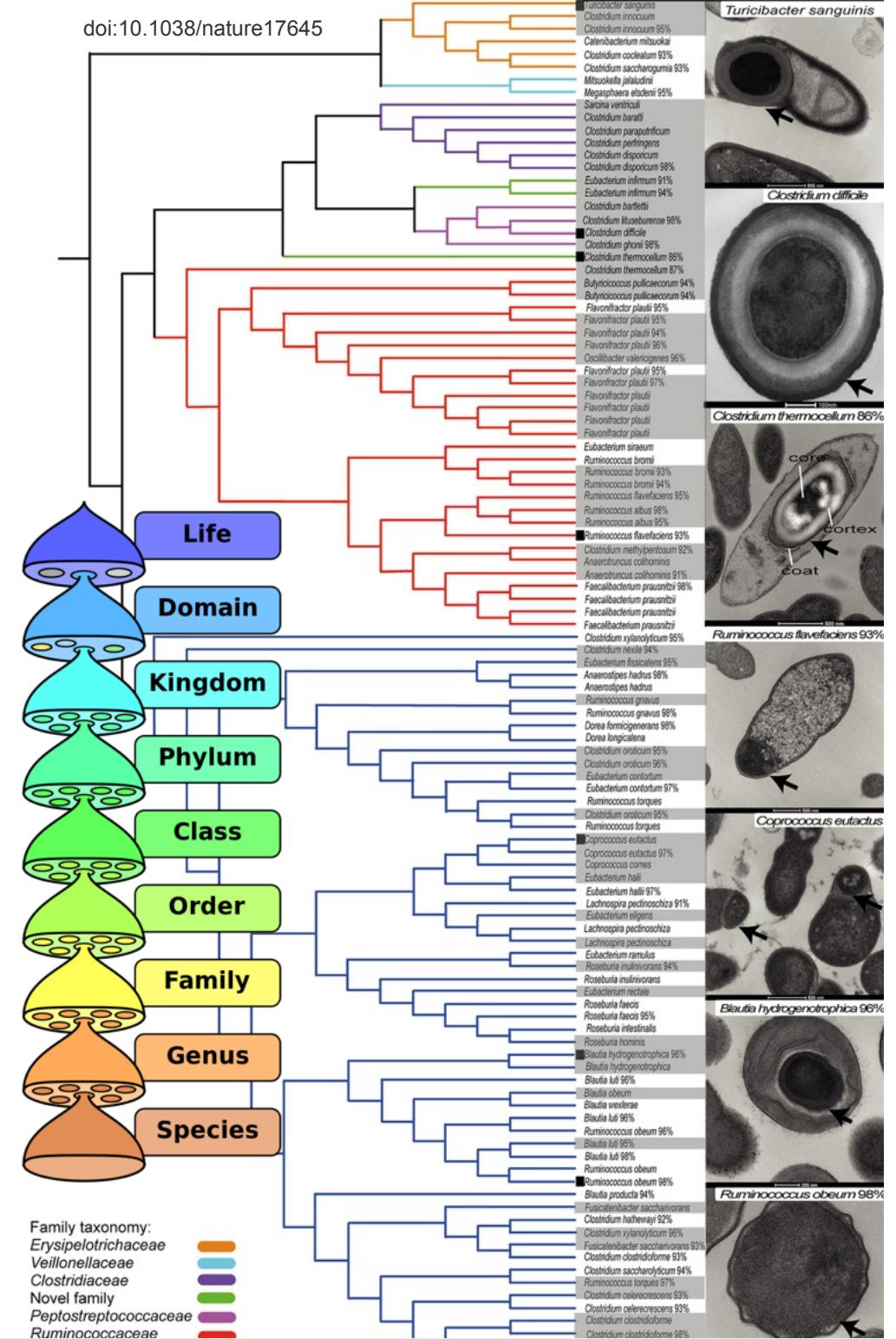


Non-Communicable Diseases (NCDs) are exponentially more prevalent. Modelling the **role of the environment**.



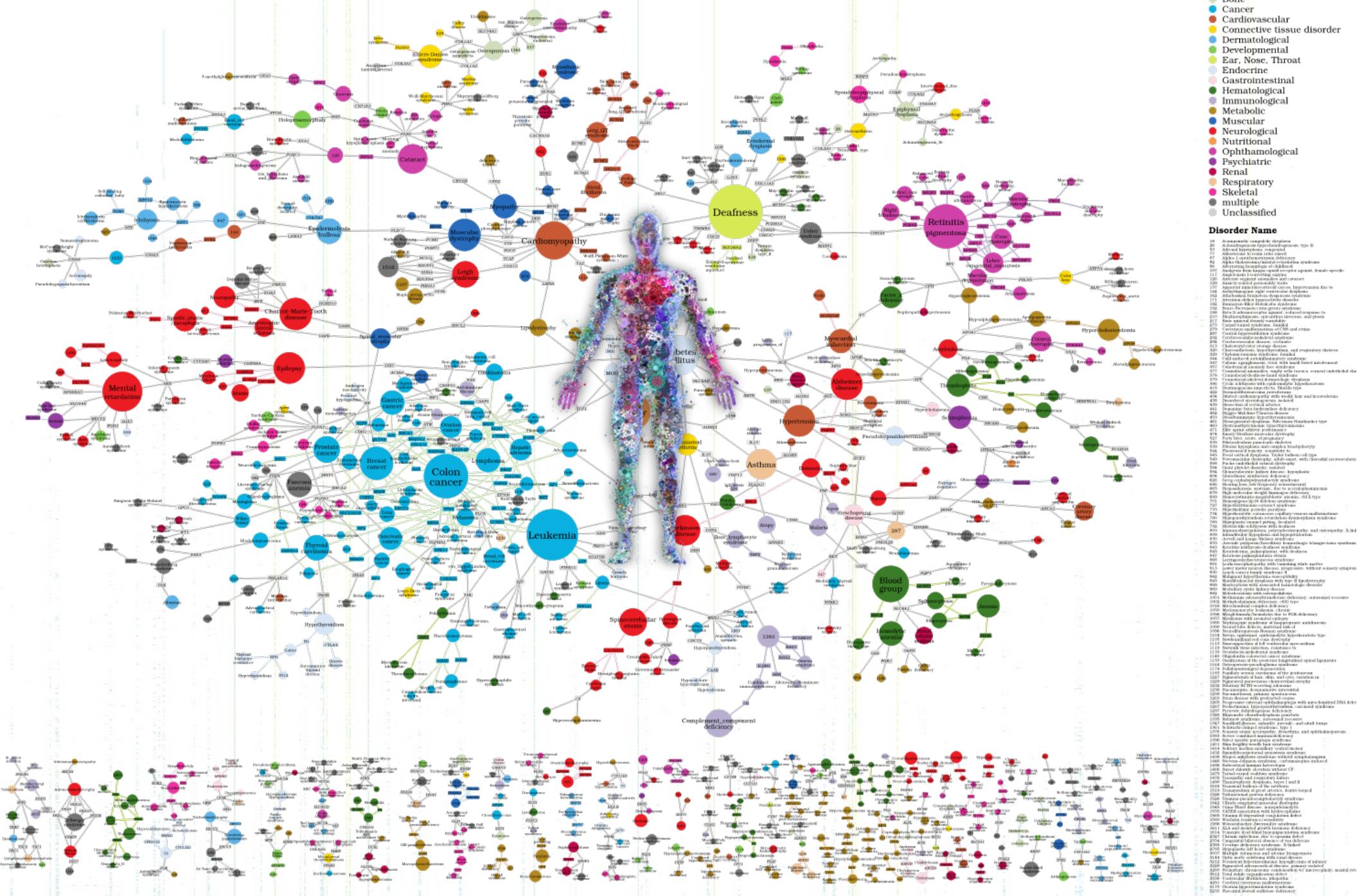


- 2Kg of microbes; **1/2 of the cells is bacterial**; 400 times more bacterial DNA.
 - ~**5000** species, ~**300** in average in a given sample
 - Important in **metabolism, immune system education, barrier** with the outside world, etc.



The human disease network

Goh K-I, Cusick ME, Valle D, Childs B, Vidal M, Barabási A-L (2007) Proc Natl Acad Sci USA 104:8685-8690



Supporting Information Figure 13 | Bipartite graph representations of the disorders. A disorder (circle) and a gene (rectangle) are connected if the gene is implicated in the disorder. The size of the circle represents the number of distinct genes associated with the disorder. Dashed disorders (disorders having no links to other disorders) are not shown. Also, only genes connecting disorders are shown.

The gut microbiome is very **complex to model**



Severine Affeldt

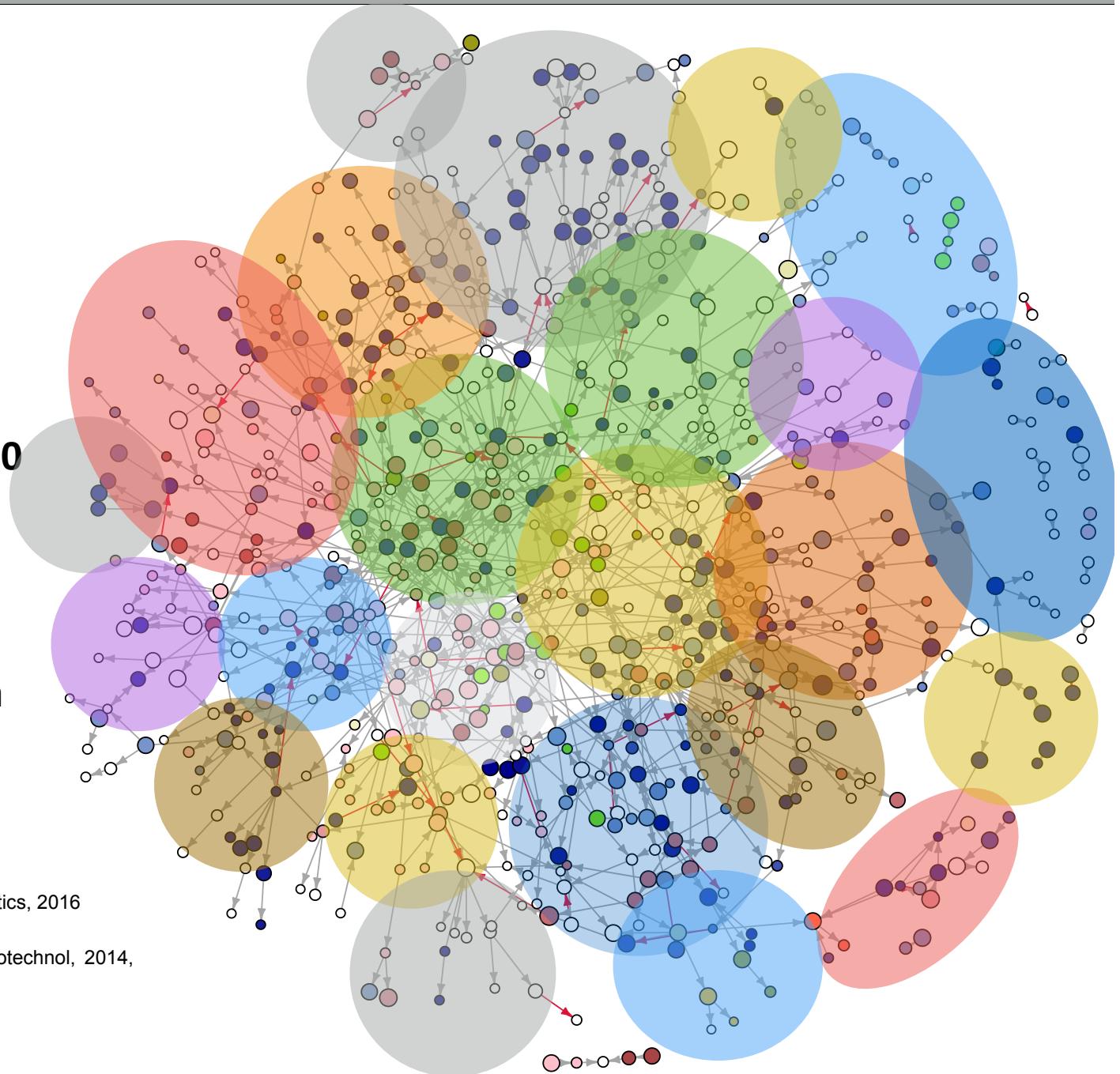


Nataliya Sokolovska



Jean-Daniel Zucker

- $p \sim 5000$ (variables) $> n \sim 100$ (observations)
- exponential complexity
- highly sparse data
- limited annotations
- **scalenet**: local reconstruction (eigenvectors / laplacian)



S Affeldt, N Sokolovska, [E Prifti](#), JD Zucker. BMC bioinformatics, 2016

S Affeldt, N Sokolovska, [E Prifti](#), JD Zucker. IJCNN, 2017

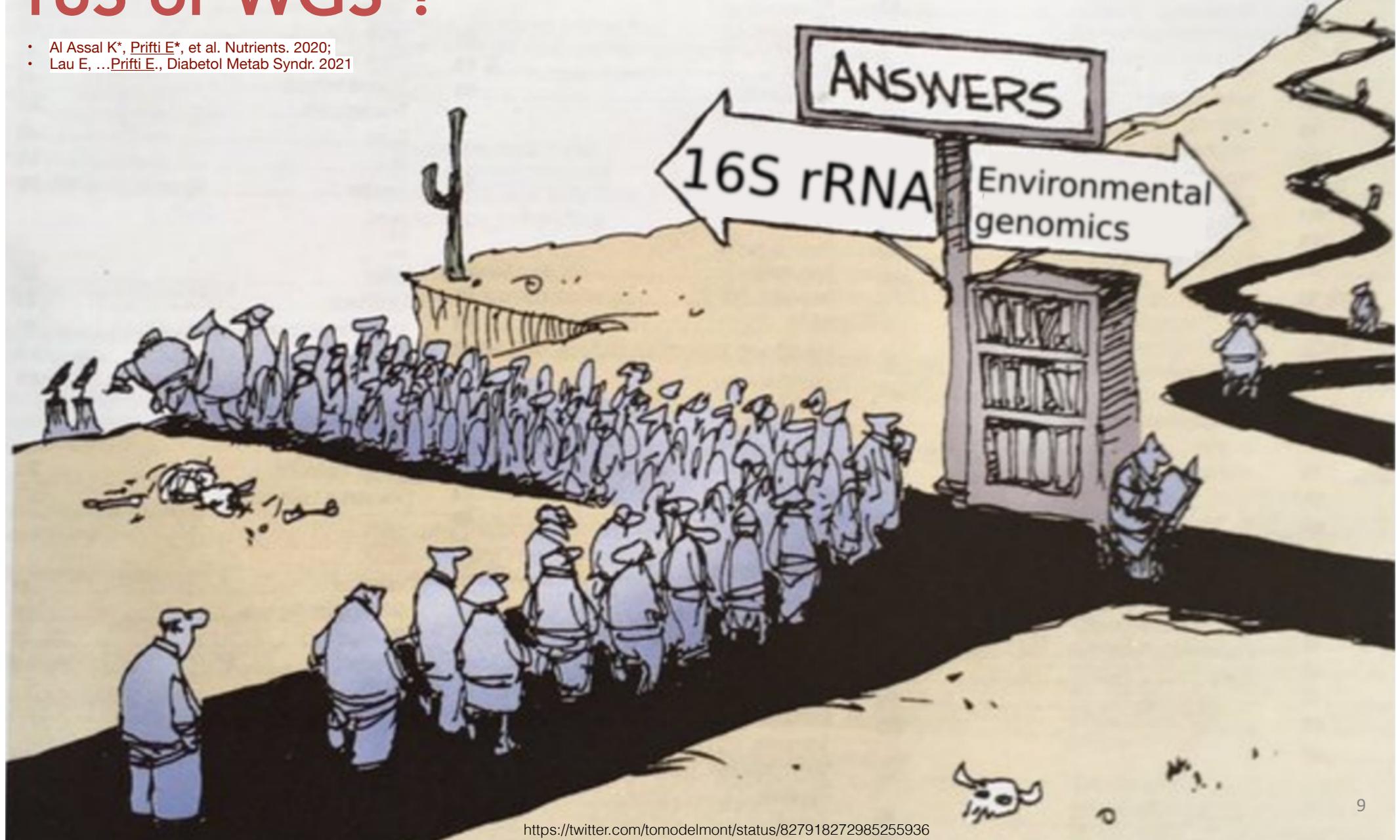
Nielsen, H. B.*., Almeida, M.*., ...[Prifti](#), E., et al. Nature Biotechnol, 2014,
32(8), 822-8

Cougoul et al 2019



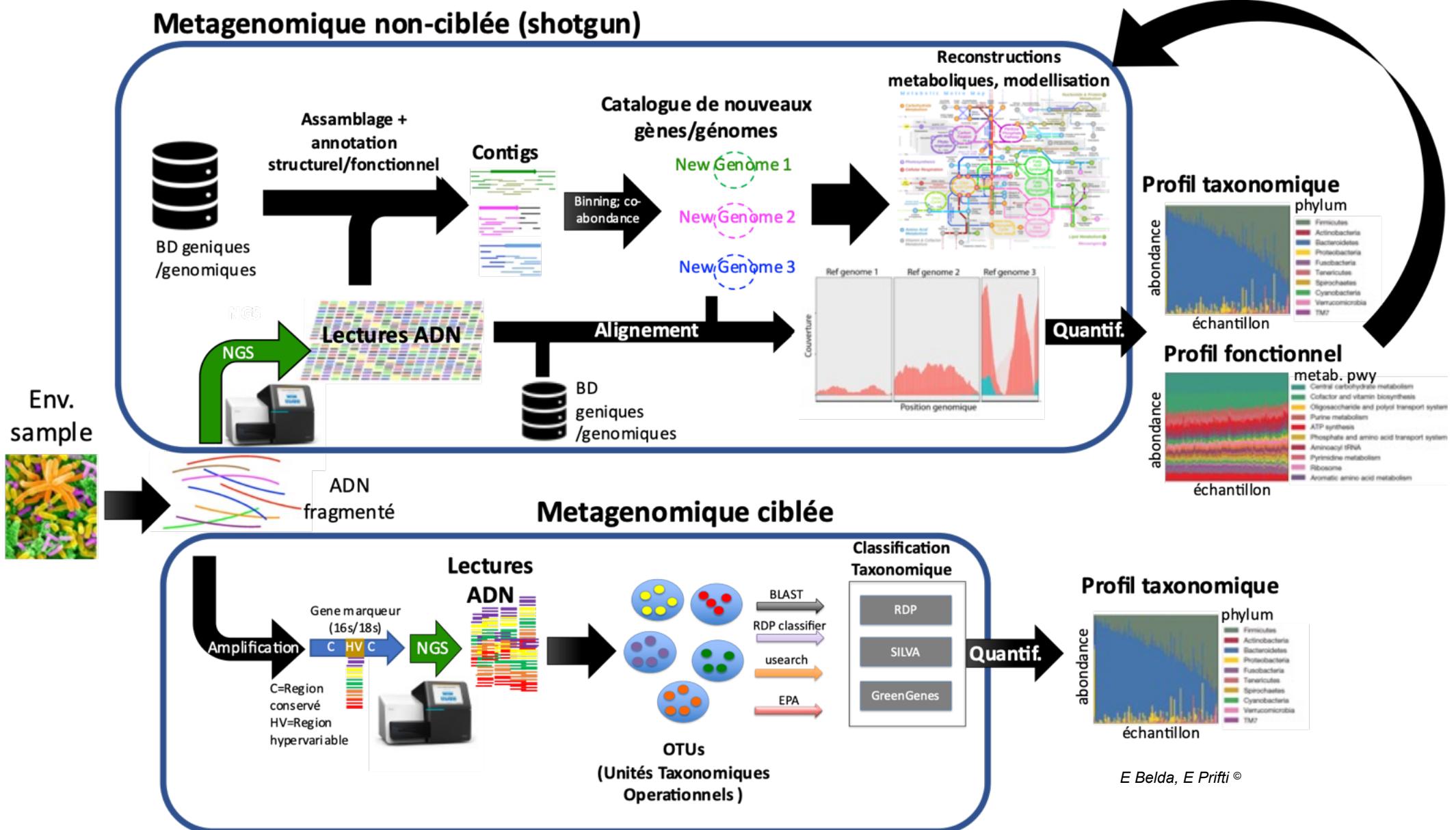
16S or WGS ?

- Al Assal K*, Prifti E*, et al. Nutrients. 2020;
- Lau E, ...Prifti E, Diabetol Metab Syndr. 2021

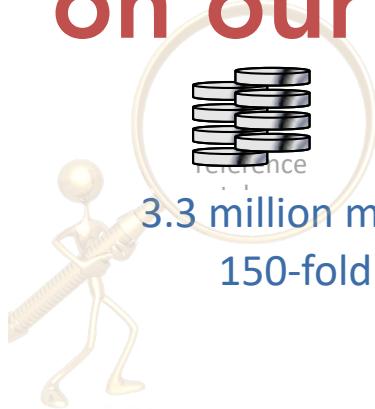


<https://twitter.com/tomodelmont/status/827918272985255936>

There are two main approaches of microbiome “quantitative” profiling



The way how we see the world depends on our reference



3.3 million microbial gene catalogue
150-fold human genome !!!
Qin et al Nature, 2010

Human Animal Environment

Terabases

Gigabases

+ 800 Mb Open ocean

150 Mb Lean mouse
80 Mb Deep sea

Megabases

2006

2007

2008

2009

2010

2011

2012

2013



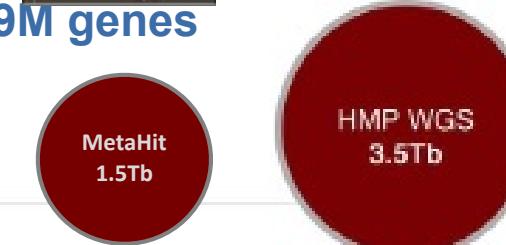
3.3M genes



3.9M genes

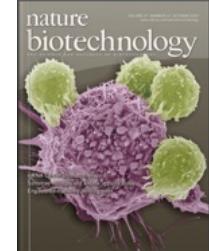


5.1M genes



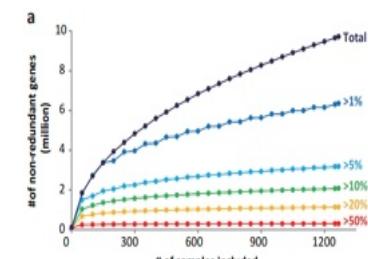
100 Gb human gut

30 Gb HMP 16S



10M genes

MetaHit HMP
and others



The wheel is turning as we go and our view of the world becomes more complex



Nicolas Pons



E. Le Chatelier



Mathieu Almeida

UHGP
>170M
proteins
205k genomes

nature
biotechnology

RESOURCE
<https://doi.org/10.1038/s41587-020-0603-3>

OPEN
A unified catalog of 204,938 reference genomes from the human gut microbiome

Alexandre Almeida^{1,2}✉, Stephen Nayfach^{3,4}, Miguel Boland¹, Francesco Strozzi⁵, Martin Beracochea⁶, Zhou Jason Shi^{6,7}, Katherine S. Pollard^{8,9,10,11}, Ekaterina Sakharova¹, Donovan H. Parks¹², Philip Hugenholtz¹², Nicola Segata¹³, Nikos C. Kyrpides^{14,3,4} and Robert D. Finn¹²✉

Comprehensive, high-quality reference genomes are required for functional characterization and taxonomic assignment of the human gut microbiota. We present the Unified Human Gastrointestinal Genome (UHGG) collection, comprising 204,938 non-redundant genomes from 4,644 gut prokaryotes. These genomes encode >170 million protein sequences, which we collated in the Unified Human Gastrointestinal Protein (UHGP) catalog. The UHGP more than doubles the number of gut proteins in comparison to those present in the Integrated Gene Catalog. More than 70% of the UHGG species lack cultured representatives, and 40% of the UHGP lack functional annotations. Intraspecies genomic variation analyses revealed a large reservoir of accessory genes and single-nucleotide variants, many of which are specific to individual human populations. The UHGG and UHGP collections will enable studies linking genotypes to phenotypes in the human gut microbiome.

<https://www.nature.com/articles/s41587-020-0603-3#article-info>

Published: 20 July 2020

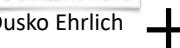


Florian Plaza
Onate



Dusko Ehrlich

Le Chatelier et al Nature, 2013
Cotillard et al Nature, 2013
Prifti et al. IHMC, 2013
Nielsen, et al. Nat Biotech. 2014
Li et al. Nat Biotech 2014
Almeida et al The ISME Journal, 2016

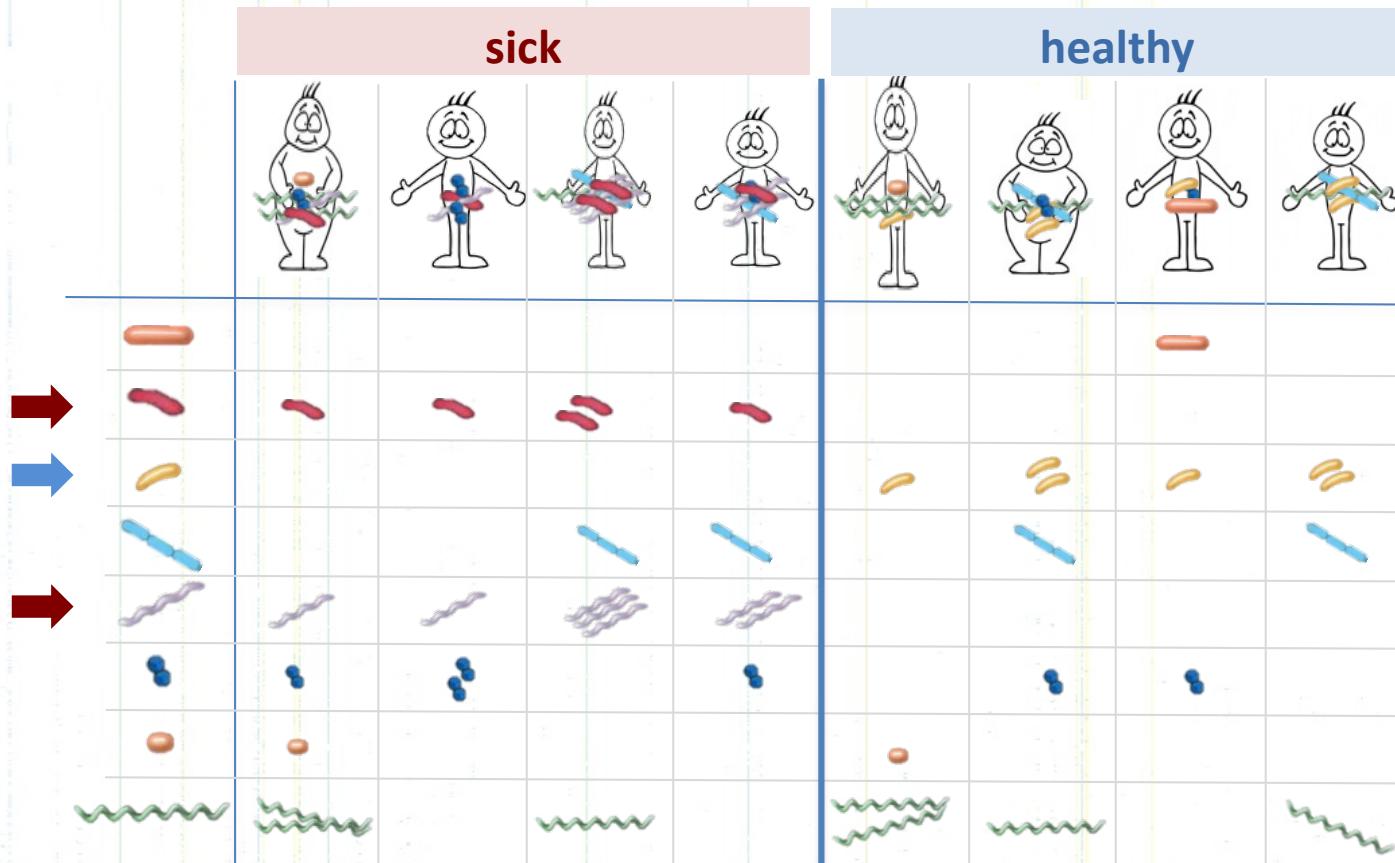


teams





The gut microbiome offers important disease classification



Stratification of the individuals & personalised medicine.

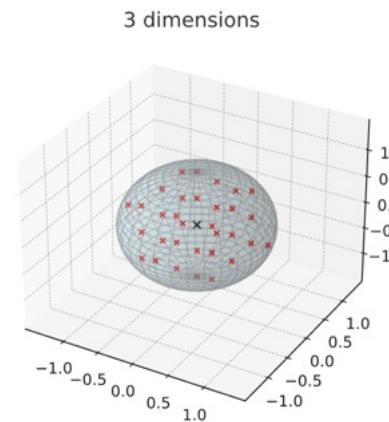
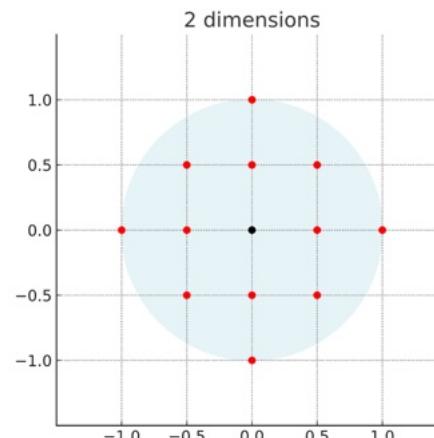
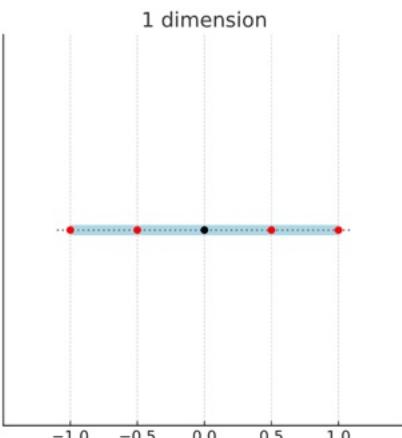
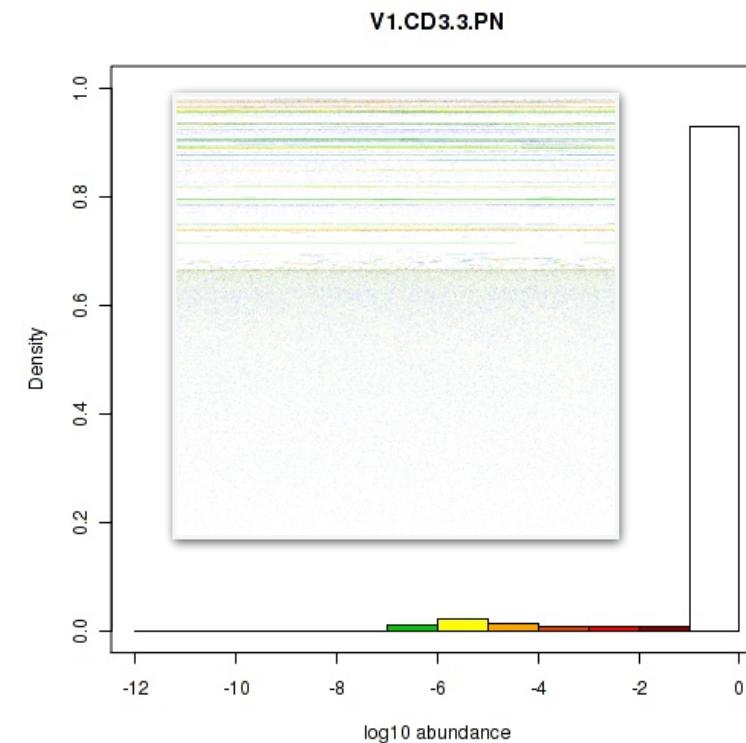
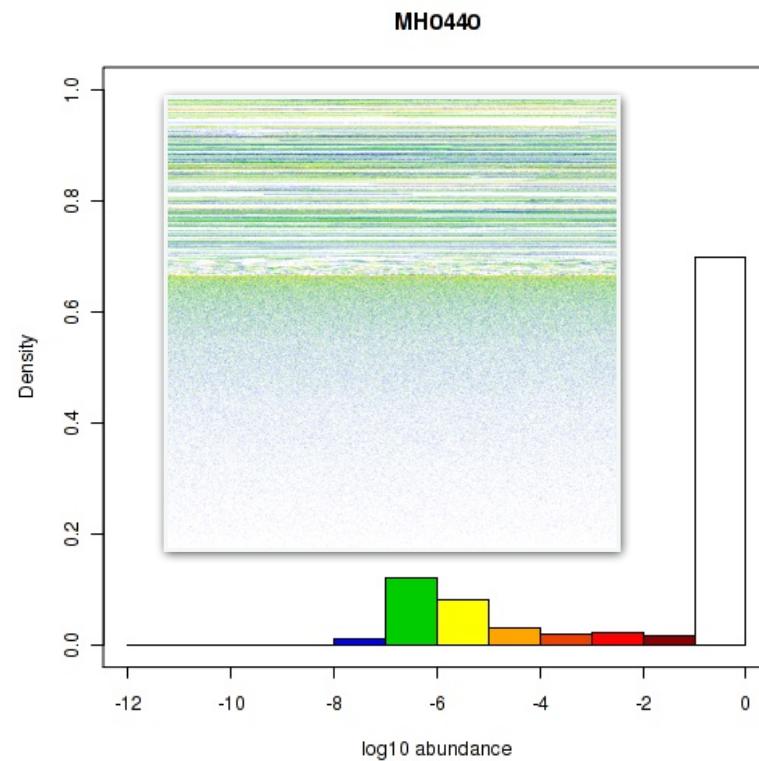
Health vs. disease biomarkers.

identify microbiota signature used as a diagnostic or prognostic tools.

Interventional medicine.

modify the gut ecosystem using prebiotics, probiotics, microbiota transplantation or microbiome targeted drugs.

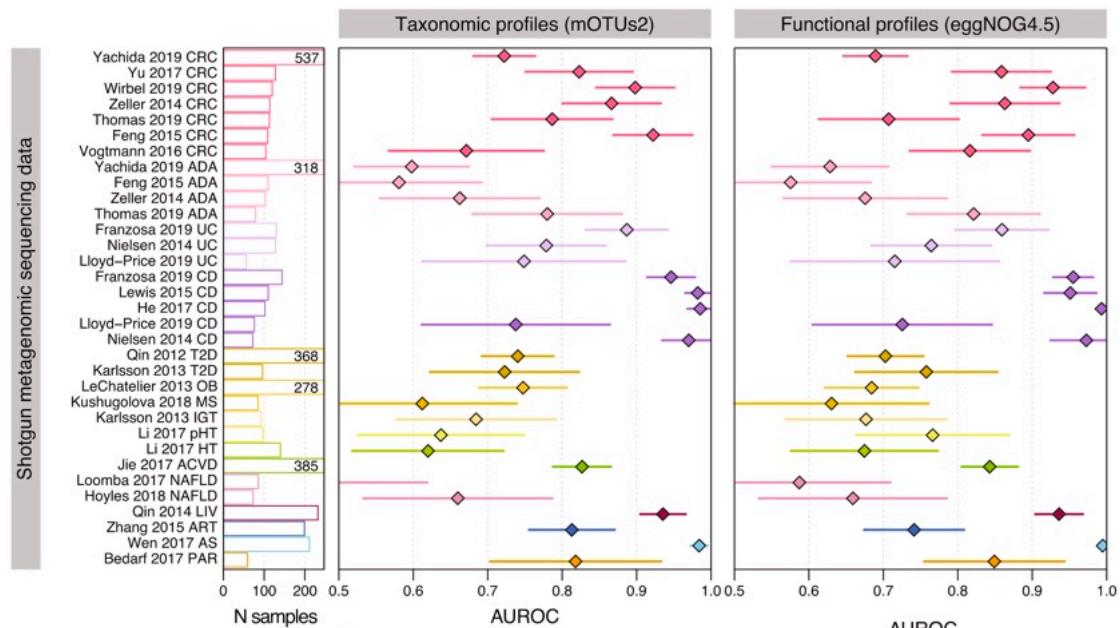
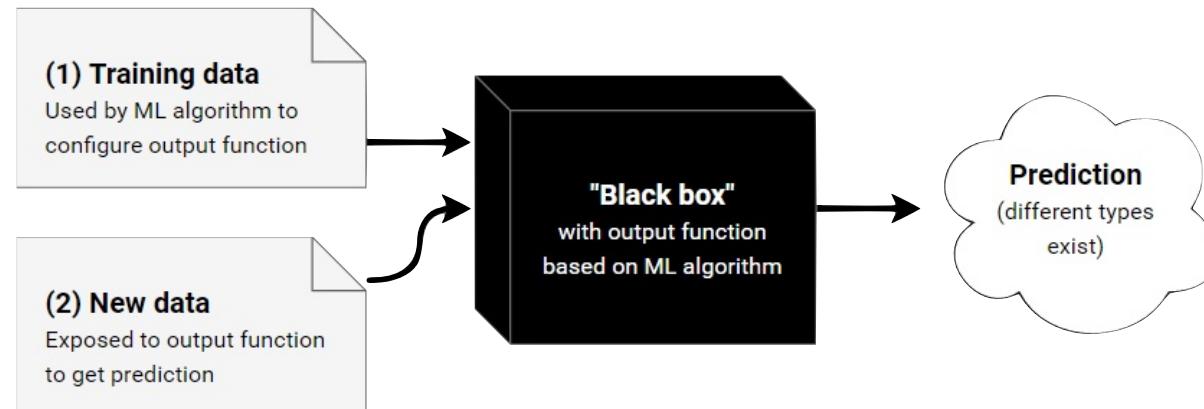
These are **not straightforward** data to work with.
They are **compositional** and full of **empty space**.



10 000 000 dimensions ?

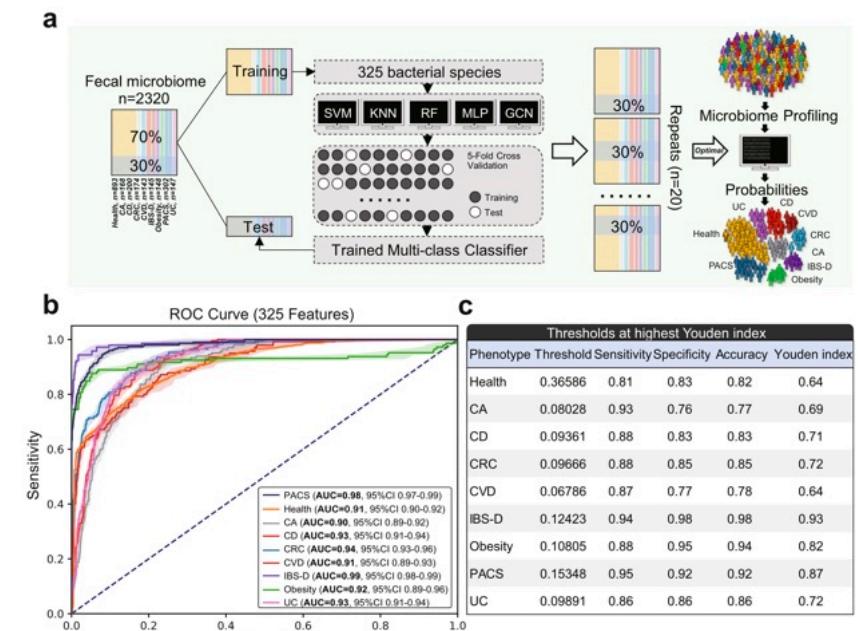
5000 dimensions ?

Using the gut microbiome and AI to find **novel biomarkers**



Microbiome meta-analysis and cross-disease comparison enabled by the SIAMCAT machine learning toolbox

Jakob Wirbel¹, Konrad Zych^{1,2}, Morgan Essex^{1,3}, Nicolai Karcher^{1,4}, Ece Kartal¹, Guillem Salazar⁵, Peer Bork^{1,6,7,8}, Shinichi Sunagawa⁵ and Georg Zeller¹



nature communications



Article

<https://doi.org/10.1038/s41467-022-34405-3>

Faecal microbiome-based machine learning for multi-class disease diagnosis

Received: 2 September 2022

Accepted: 21 October 2022

Published online: 10 November 2022

Qi Su^{1,2,3,4,6}, Qin Liu^{1,2,3,4,6}, Raphaela Iris Lau^{1,2,3}, Jingwan Zhang^{1,2,3,4},

Zhilu Xu^{1,2,3,4}, Yun Kit Yeoh¹, Thomas W. H. Leung², Whitney Tang²,

Lin Zhang^{1,2,3,4}, Jessie Q. Y. Liang^{2,3,4}, Yuk Kam Yau^{1,2,3}, Jiaying Zheng^{1,2,3},

Chengyu Liu^{1,2,3}, Mengjing Zhang^{1,2,3}, Chun Pan Cheung^{1,2,4},

Jessica Y. L. Ching^{1,2,3}, Hein M. Tun^{1,2,3}, Jun Yu^{1,2,3,4}, Francis K. L. Chan^{1,2,3,4}

Back to the drawing board !

Develop a method that is Interpretable by nature, full of constraints and thought for the microbiome data



GigaScience, 9, 2020, 1-11

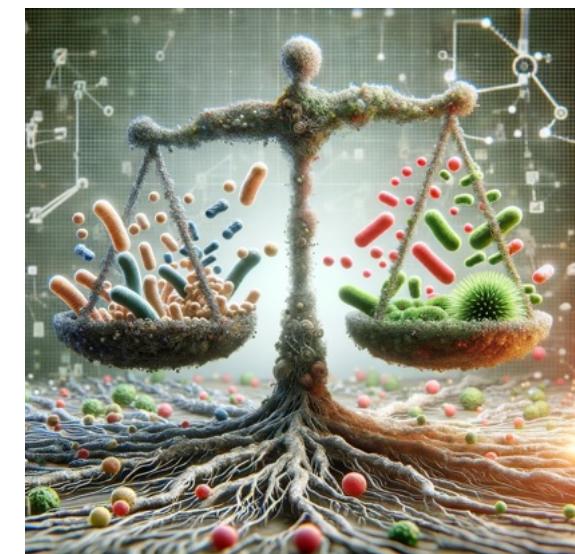
doi: 10.1093/gigascience/giaa010
Research

RESEARCH

Interpretable and accurate prediction models for metagenomics data

Edi Prifti ^{1,2,*}, Yann Chevaleyre ³, Blaise Hanczar ⁴, Eugeni Belda ²,
Antoine Danchin ⁵, Karine Clément ^{6,7} and Jean-Daniel Zucker ^{1,2,6,*}

¹IRD, Sorbonne University, UMMISCO, 32 Avenue Henri Varagnat, F-93143 Bondy, France; ²Institute of Cardiometabolism and Nutrition, ICAN, Integromics, 91 Boulevard de l'Hopital, F-75013, Paris, France;



<https://academic.oup.com/gigascience/article/9/3/giaa010/5801229>

Veillonella unclassified +# Clostridium perfringens
 < 18.0 % then class = healthy

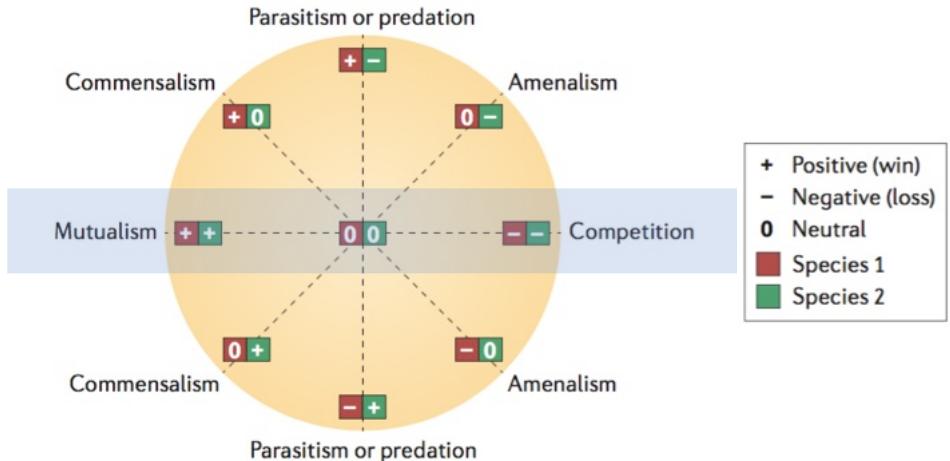
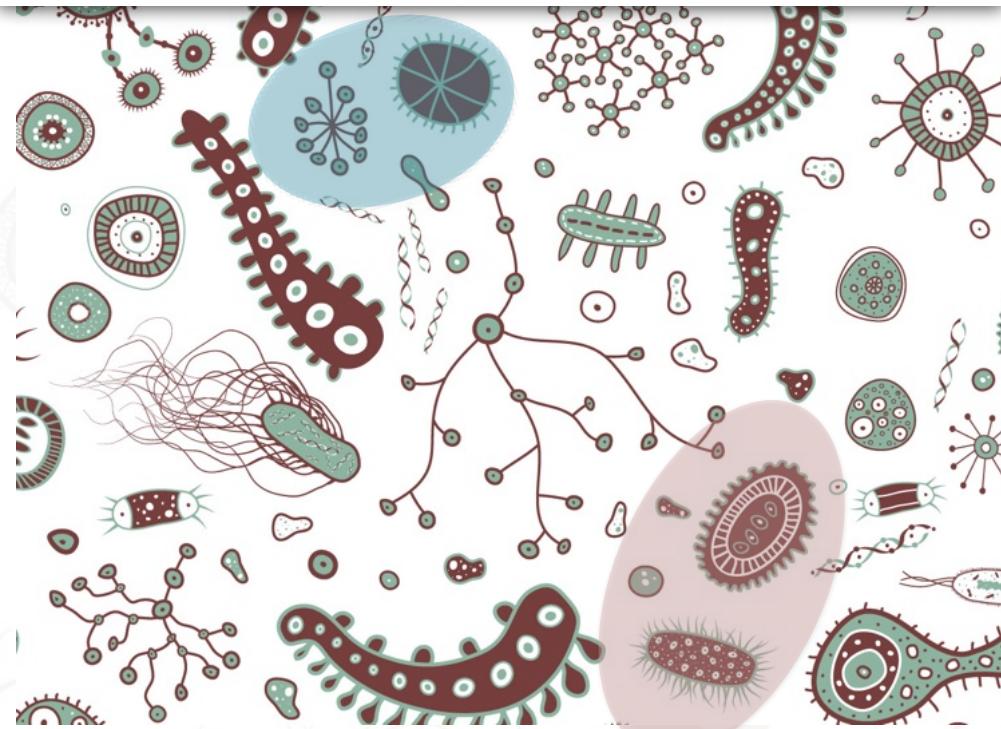


Figure 1 | Summary of ecological interactions between members of different species. The wheel display introduced by Lidicker¹ has been adapted to summarize all possible pairwise interactions. For each interaction partner, there are three possible outcomes: positive (+), negative (-) and neutral (0). For instance, in parasitism, the parasite benefits from the relationship (+), whereas the host is harmed (-); this relationship is thus represented by the symbol pair +−.

Faust, Raes. *Nature Reviews Microbiology* 10, no. 8 (August 1, 2012): 538–550.

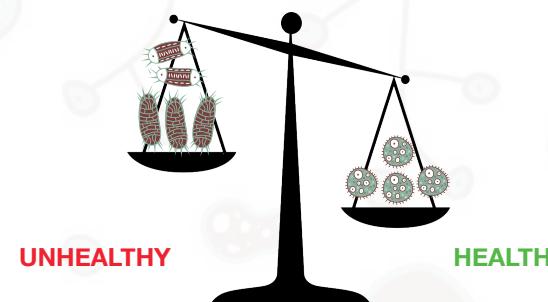
Inspired by ecological relations we proposed the **Bin/Ter/Ratio** (BTR) models, which quantify interactions by comparing abundance.



Veillonella unclassified + Clostridium perfringens
 < 18.0 % then class = healthy

BINARY

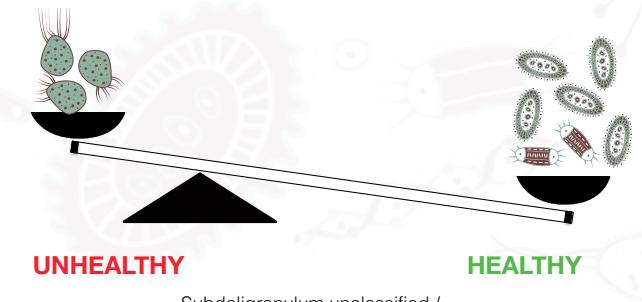
HEALTHY



Streptococcus_anginosus + Veillonella_unclassified - Alistipes_indistinctus
 < 8.3 % then class = healthy

TERNARY

UNHEALTHY



UNHEALTHY

HEALTHY

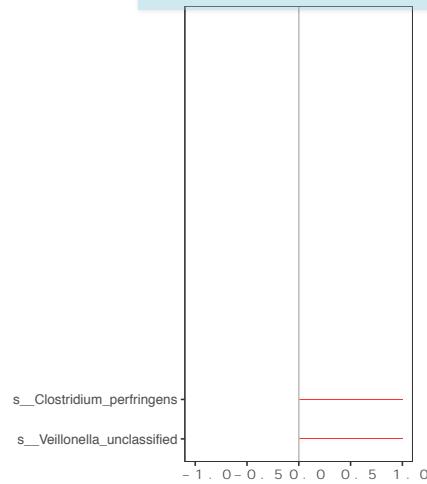
Subdoligranulum unclassified /
 (Megasphaera micronuciformis + Streptococcus anginosus)
 > 81 then class = healthy

RATIO

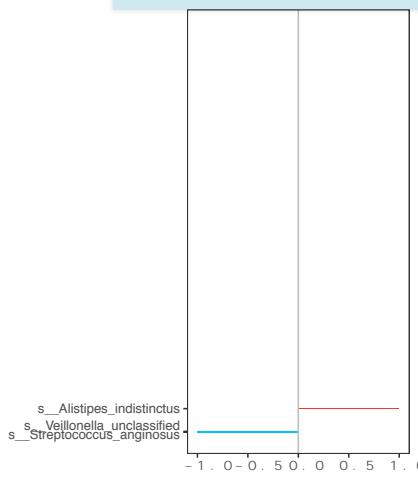
- Prifti et al. GigaScience, 2020
- Nguyen, Prifti et al. Sci Rep, 2020.

BTR models are **extremely simple** ...

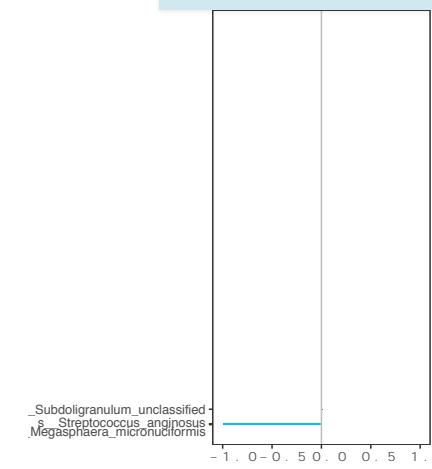
A + B > s



A + B - C > s

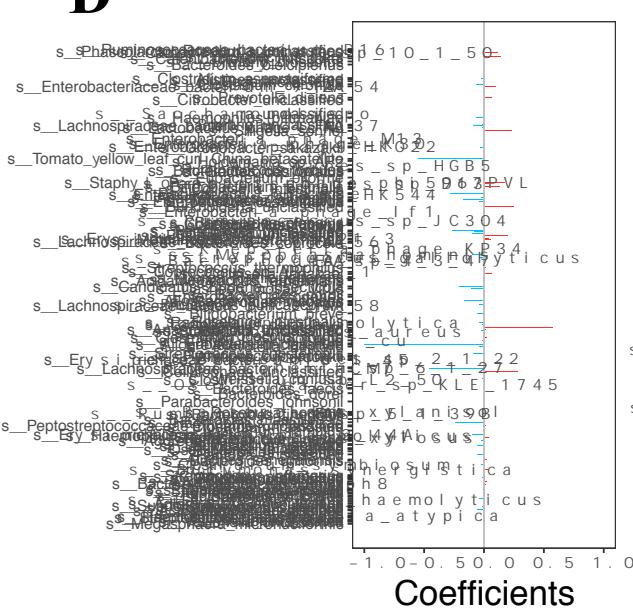


A / (B + C) > s

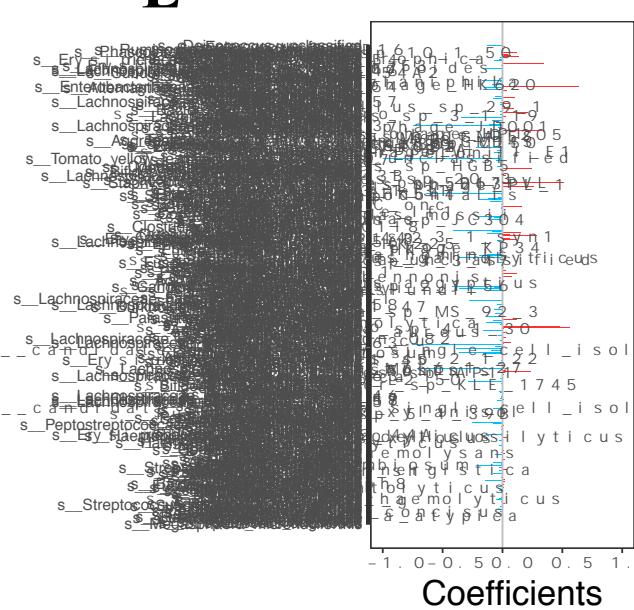


... compared with **SOTA**

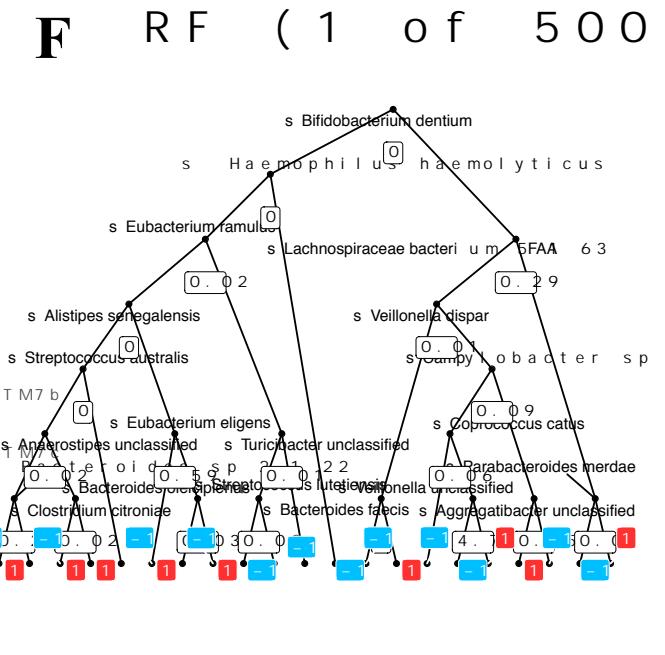
D



E



F



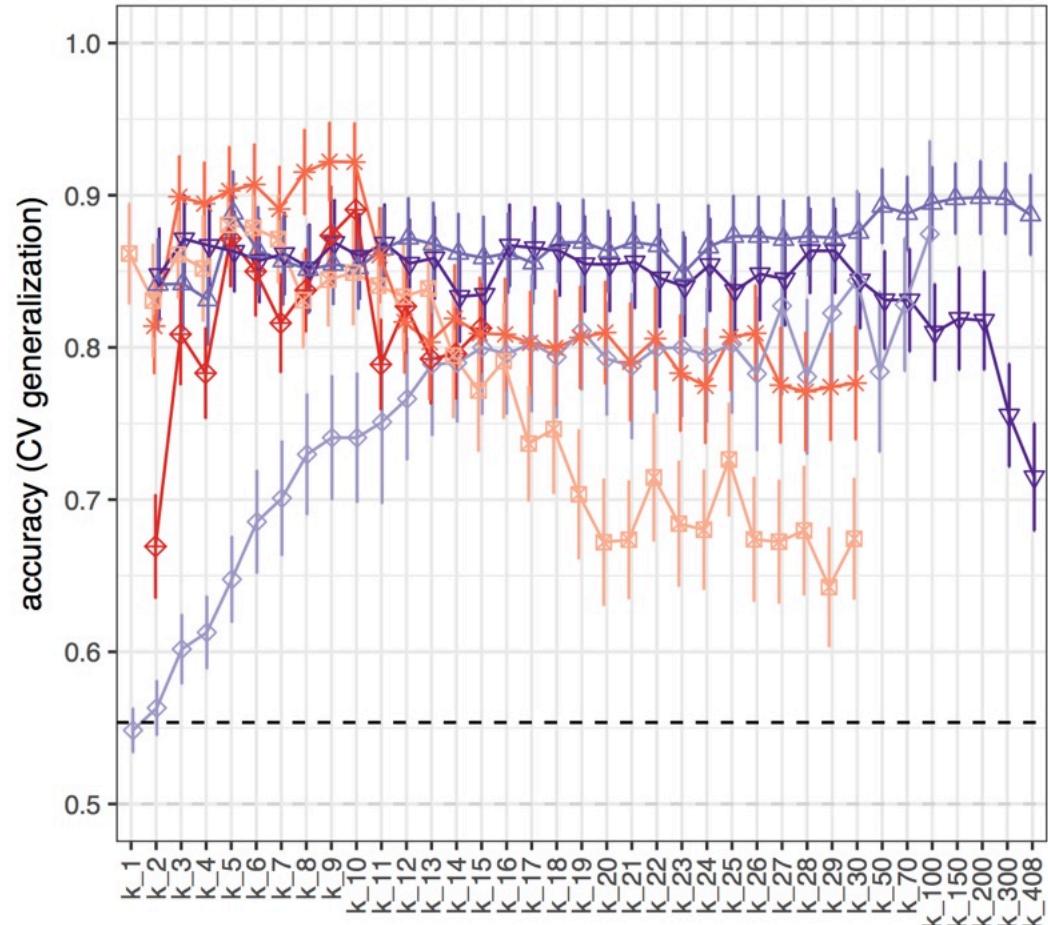
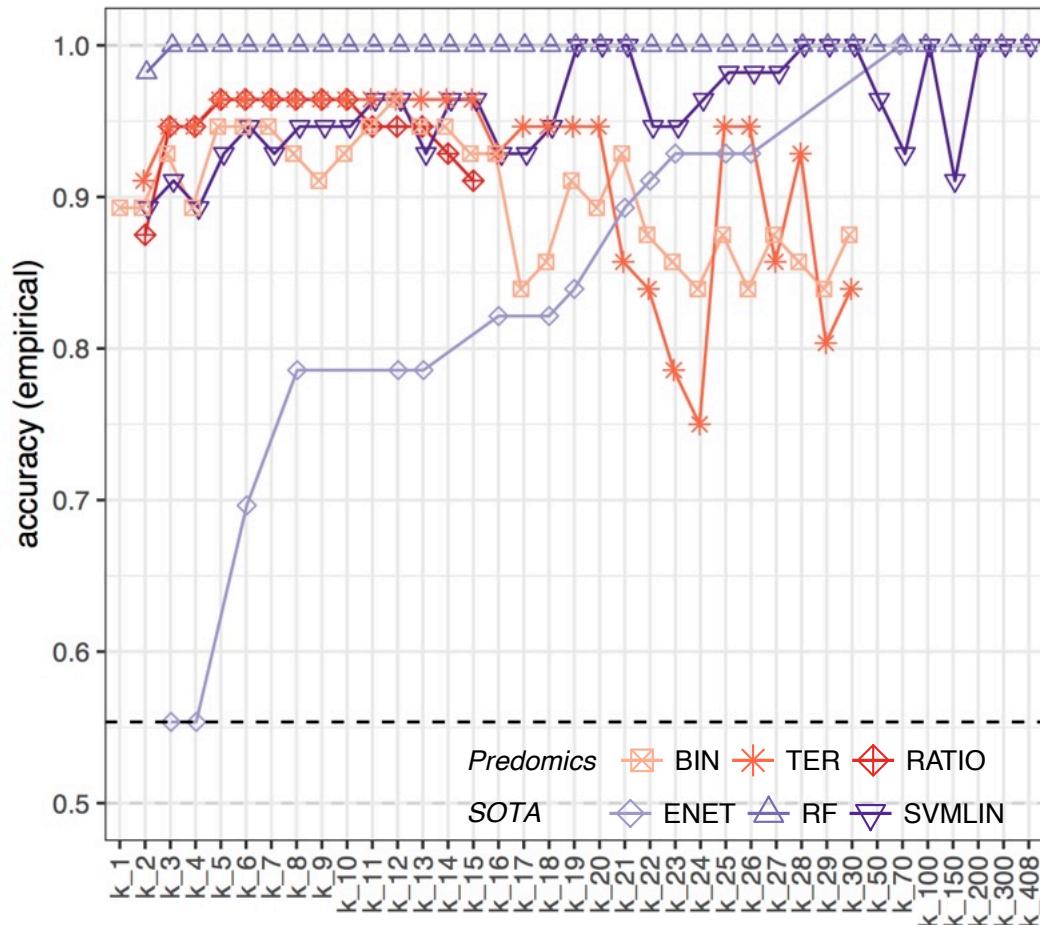
Alterations of the human gut microbiome in liver cirrhosis

Nan Qin^{1,2*}, Fengling Yang^{1*}, Ang Li^{1*}, Edi Prifti^{1*}, Yanfei Chen^{1*}, Li Shao^{1,2*}, Jing Guo¹, Emmanuelle Le Chatelier¹, Jian Yao^{1,2}, Lingliao Wu¹, Jiawei Zhou¹, Shujun Ni¹, Lin Liu¹, Nicolas Pons², Jean Michel Batté², Sean P. Kennedy³, Pierre Leonard³, Chunhua Yuan¹, Wenchao Ding², Yuanming Chen¹, Xinjun Hu¹, Beiven Zheng^{2,3}, Guirong Qian¹, Wei Xu¹, S. Dusko Ehrlich^{3,4}, Shusen Zheng^{2,5} & Lanjuan Li^{1,2}

Accessible, curated metagenomic data through ExperimentHub

Edoardo Pasolli^{*1}, Lucas Schiffer^{*2}, Audrey Renson², Valerie Obenchain³, Paolo Manghi¹, Duy Tin Truong¹, Francesco Beghini¹, Faizan Malik², Marcel Ramos², Jennifer B. Dowd^{2,4}, Curtis Huttenhower^{5,6}, Martin Morgan³, Nicola Segata^{*1}, Levi Waldron^{*1,2}

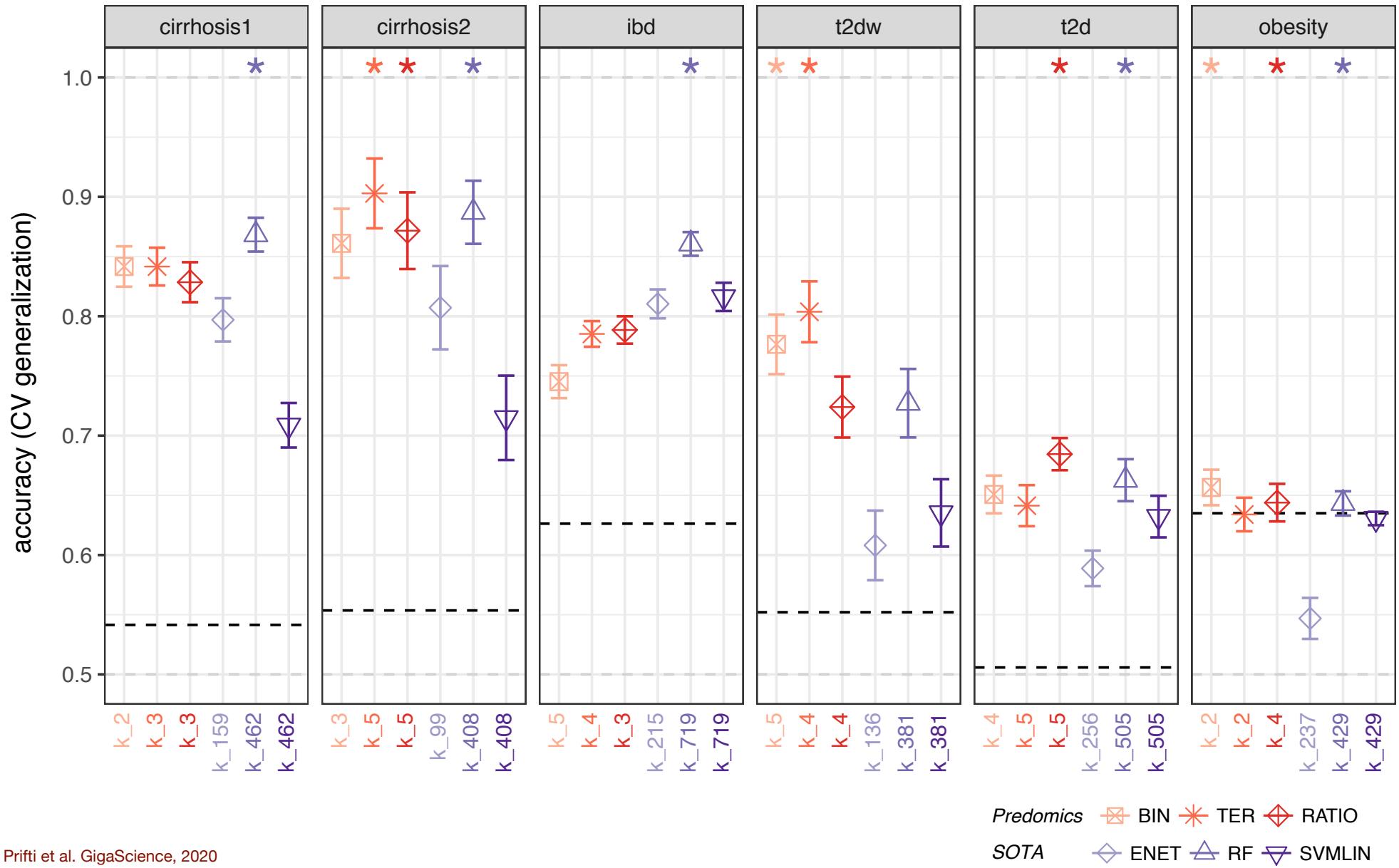
Cirrhosis stage 2 - species dataset



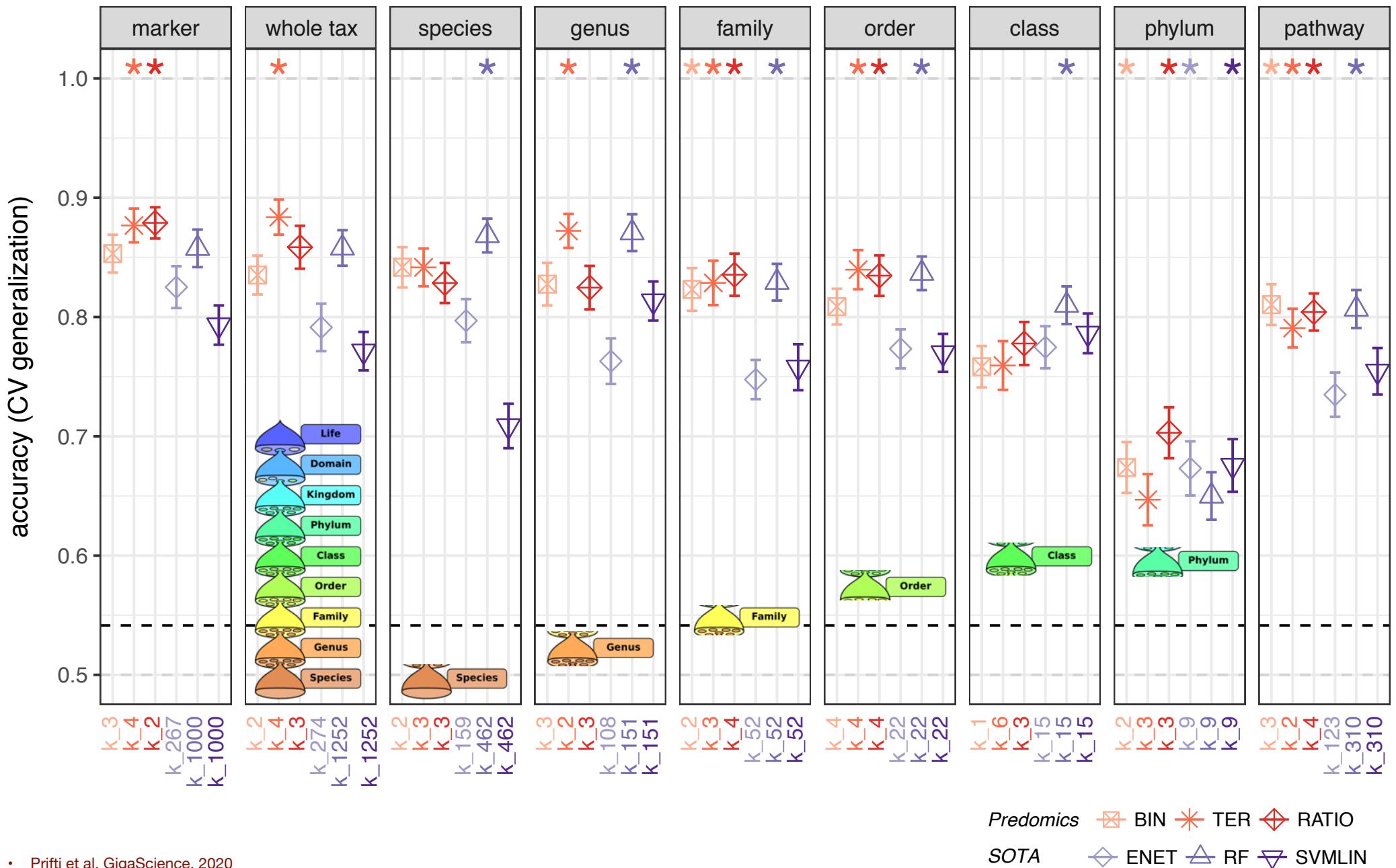
- large exploration of models
- k-selection based on training db using penalty (acc-k*penalty); penalty = 0.1/100
- cross-validation of the algorithms = 10*10 fold; statistical tests possible

• Prifti et al. GigaScience, 2020

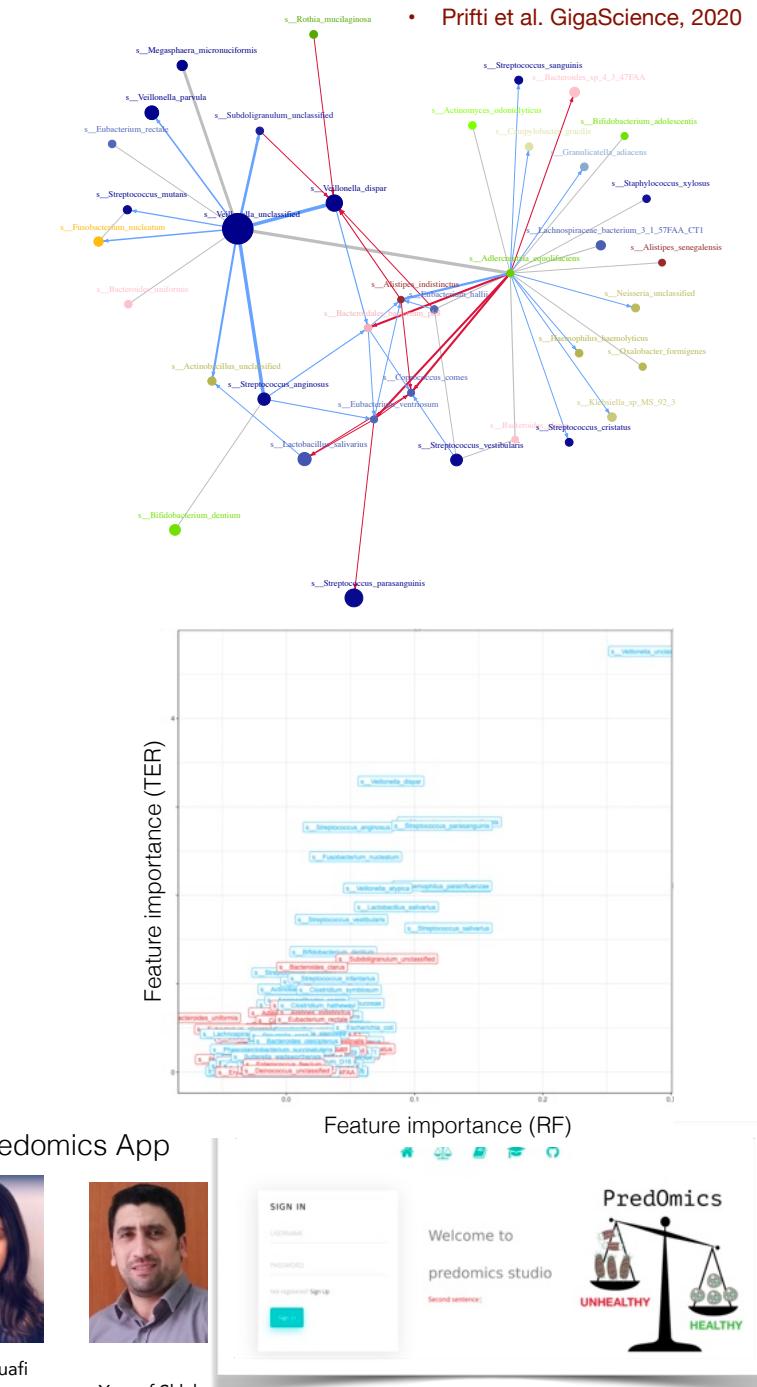
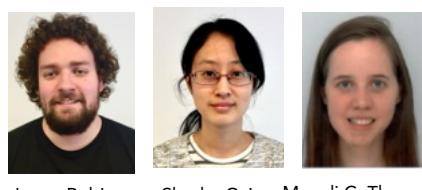
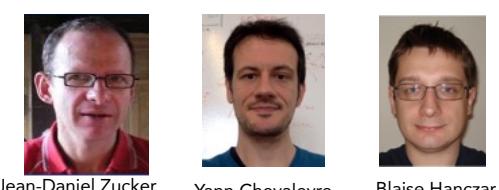
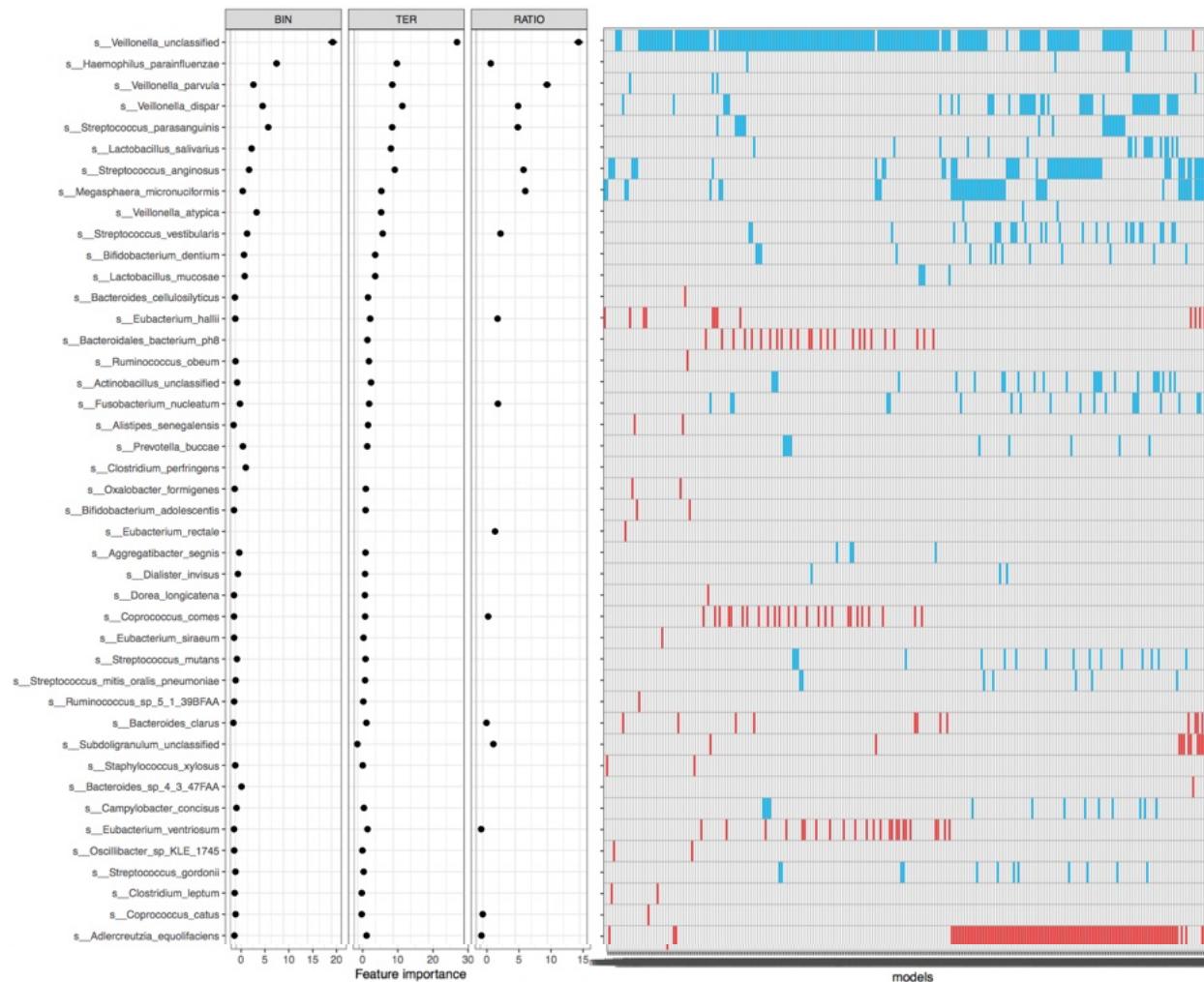
BTR models display **similar performance** with SOTA across datasets



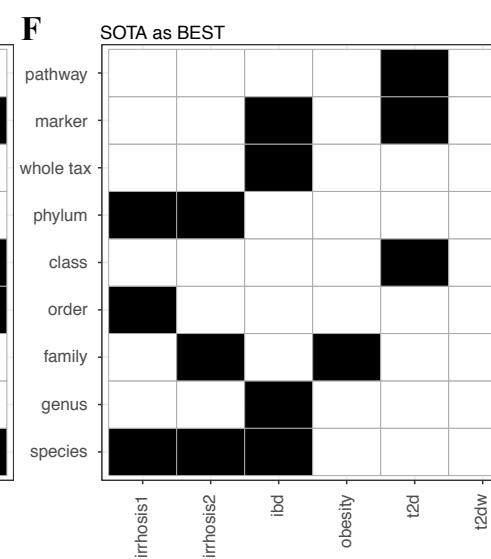
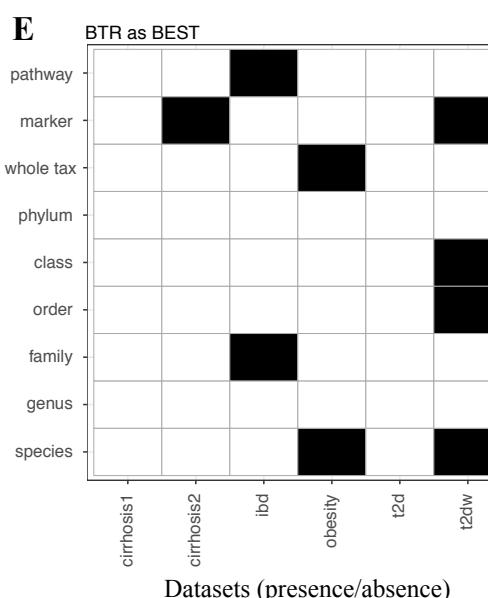
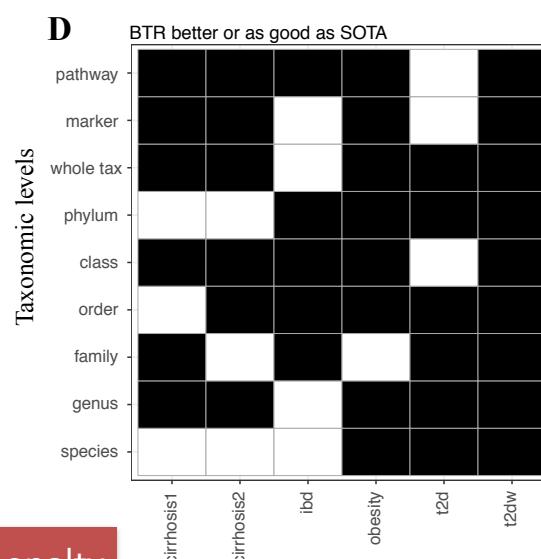
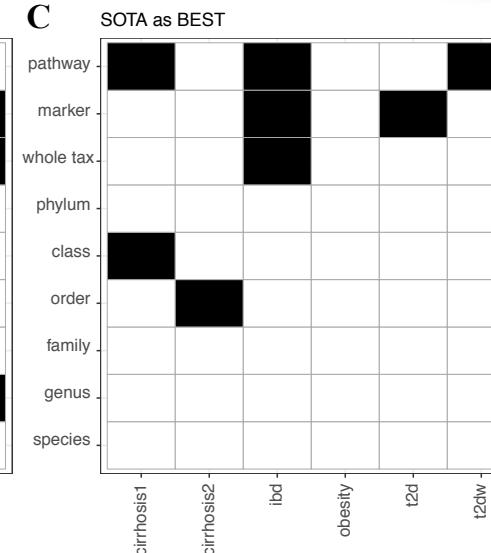
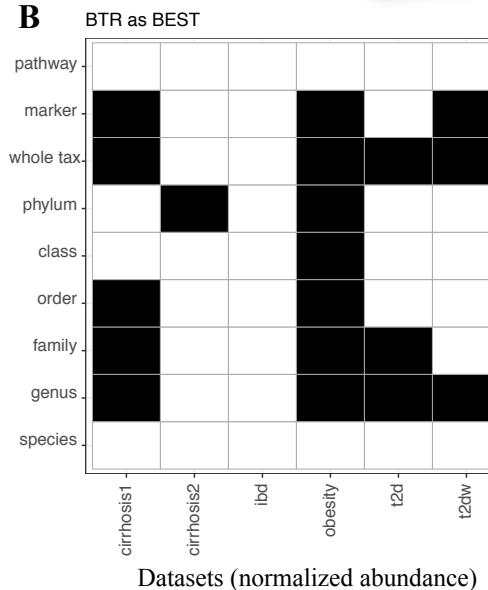
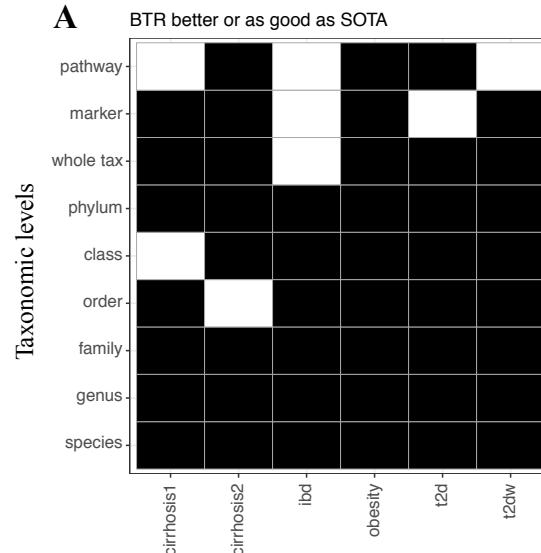
Predictive performance **decreases** with taxonomic specificity (Cirrhosis db)



The analysis of the Family of best models (FBM) **provides** the most predictive features

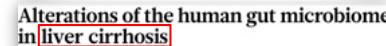
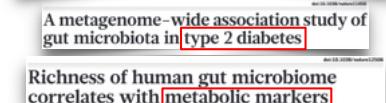
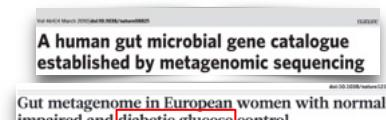


Predictive performance **decreases** with taxonomic specificity (Cirrhosis db)



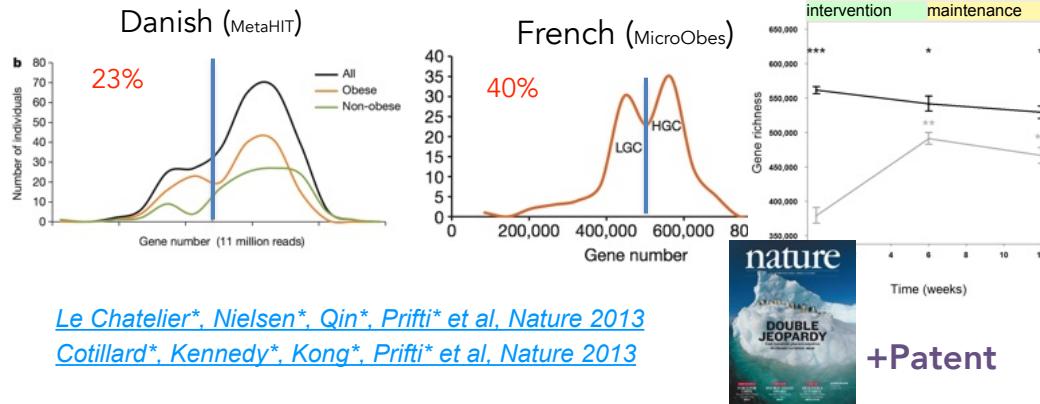
Accessible, curated metagenomic data through ExperimentHub

Edoardo Pasolli^{*1}, Lucas Schiffer^{*2}, Audrey Renson², Valerie Obenchain³, Paolo Manghi¹, Duy Tin Truong¹, Francesco Beghini¹, Faizan Malik², Marcel Ramos², Jennifer B. Dowd^{2,4}, Curtis Huttenhower^{5,6}, Martin Morgan³, Nicola Segata^{*1,2}, Levi Waldron²

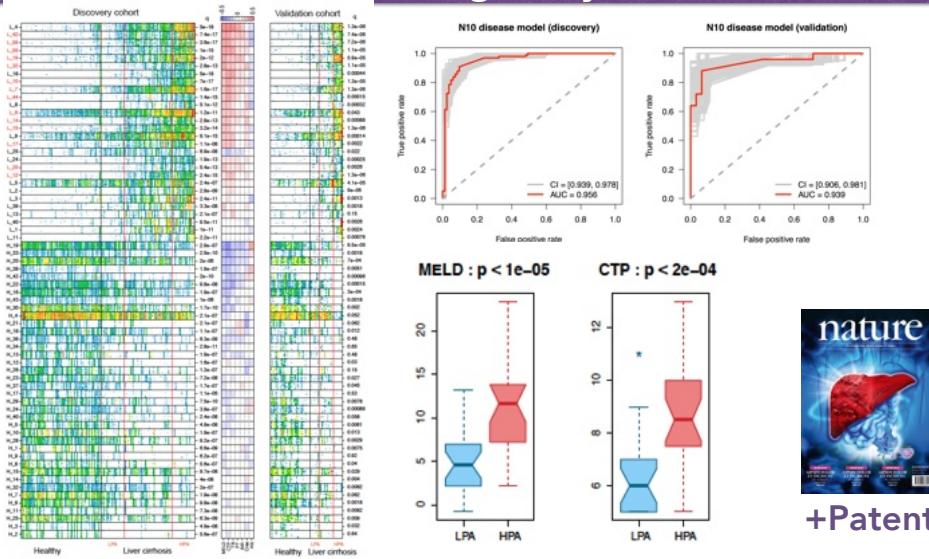


We have applied these models with **success** in numerous human cardiometabolic diseases

Discovery of the bimodal distribution of gut microbiome richness and its association with the success of dietary interventions

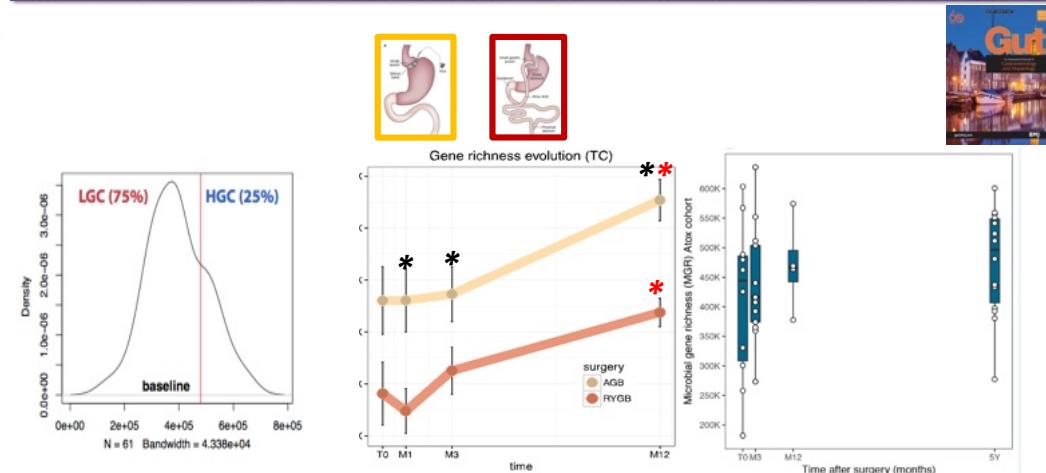


Liver cirrhosis patients have an oral originated microbiome, which is altered and can predict disease gravity.

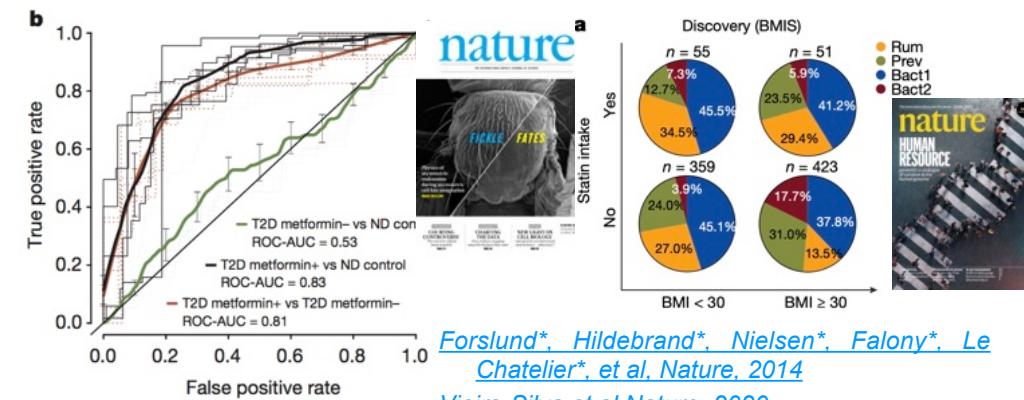


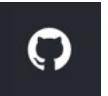
Qin*, Yang*, Li*, Prifti* et al. *Nature*, 2014

Bariatric surgery improves gut microbial richness, although it can not entirely restore it as other phenotypes.

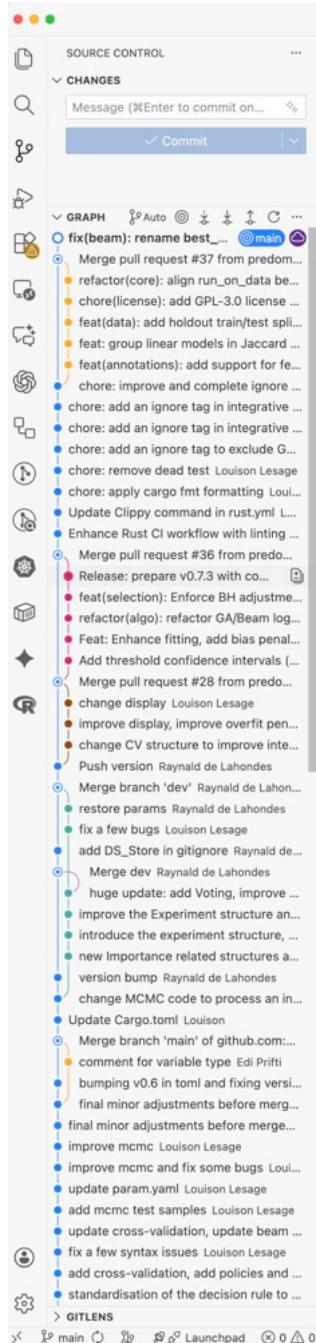


Medications confounding effects should be taken into account compte (metformin in T2D).





An **open source project** which is actively improving (multiclass, GPU, R and python integration, Shiny app and more) ...



predomics

1 follower <http://predomics.integromics.fr> @ediprisci

Overview Repositories 4 Projects Packages People 1

Popular repositories

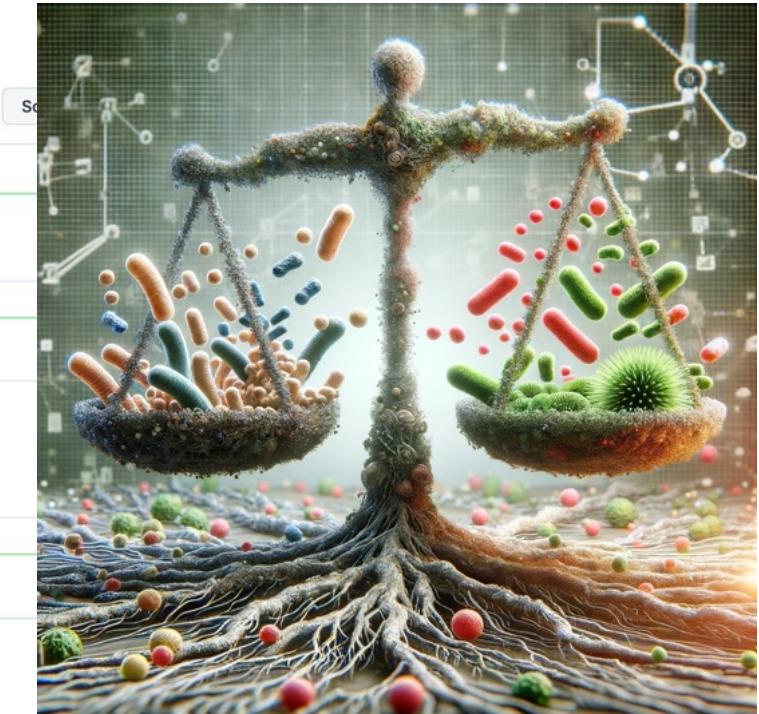
- predomicspkg** Public
HTML ⭐ 8 🏷 2
- predomics.github.io** Public
CSS
- gpredomics** Public
GPU based acceleration of predomics
Rust 🏷 1
- gpredomicsR** Public
The R package that includes the #gpredomics package

Repositories

Find a repository... Type Language Sort

- gpredomics** Public
GPU based acceleration of predomics Label: Public nics
Rust ⭐ 0 🏷 1 ⚡ 13 📈 0 Updated yesterday
- predomicspkg** Public
HTML ⭐ 8 GPL-3.0 🏷 2 ⚡ 4 📈 0 Updated 2 days ago
- gpredomicsR** Public
The R package that includes the #gpredomics package
⭐ 0 🏷 0 ⚡ 0 📈 0 Updated 3 days ago
- predomics.github.io** Public
CSS ⭐ 0 🏷 0 ⚡ 0 📈 0 Updated on Nov 2, 2023

<https://github.com/predomics/>



The gpredomics parameter space

```
general:
  seed: 42
  algo: ga
  cv: false
  thread_number: 8
  gpu: true
  language: ter,ratio
  data_type: raw,prev
  epsilon: 1e-5
  fit: auc
  k_penalty: 0.0001
  fr_penalty: 0.0
  nb_best_model_to_test: 100
  #log_base: ""
  log_level: info
  display_level: 2
  display_colorful: true
  keep_trace: false
  save_exp: exp.mp

data:
  X: "../../../../data/X/nsatDisc.tsv"
  y: "../../../../data/y/yDisc_OS12p.tsv"
  Xtest: "../../../../data/X/nsatValid.tsv"
  ytest: "../../../../data/y/yValid_OS12p.tsv"
  features_maximal_number_per_class: 0
  feature_minimal_prevalence_pct: 10
  feature_minimal_feature_value: 1e-4
  feature_selection_method: wilcoxon
  feature_maximal_pvalue: 0.05
  feature_minimal_log_abs_bayes_factor: 2
  classes:
    - "OS12-"
    - "OS12+"
    - "unknown"

cv:
  inner_folds: 5
  overfit_penalty: 5
  outer_folds: 5
  fit_on_valid: true
  cv_best_models_ci_alpha: 0.05

importance:
  compute_importance: false
  n_permutations_oob: 100
  scaled_importance: true
  importance_aggregation: mean

voting:
  vote: true
  use_fbm: true
  min_perf: 0.5
  min_diversity: 10
  method: Majority

# used in parent selection, child conception (cross over) and mutation, all of which is single thread
# ga for genetic algorithm, beam for beam algorithm and mcmc for MCMC-based algorithm
# should cross-validation be enabled?
# the number of thread used in feature selection and fit computation
# should Gpredomics use GPU ? (ga and beam only)
# possible values are ter,bin,ratio,pow2, see README.md for detail. A comma separated list (no spaces) is accepted, which means the initial population
# possible values are raw,prev,log, see README.md for detail. Same as above, comma separated list is fine.
# this is only usefull for data_type prevalence (where it is a threshold) or log (where it replaces values below)
# possible values are auc,specificity,sensitivity (classification), see README.md for details
# this penalty is deduced from fit function multiplied by k, the number of variables used in the model
# used only when fit is specificity or sensitivity, deduce (1 - symmetrical metrics) x fr_penalty to fit
# nb of models to test in the last generation (default to 10, 0 means all models)
# uncomment to print log and results in log_file_name
# possible values are trace, debug, info, warning or error
# precision in variable display (0=anonymized features, 1=feature line index, 2=feature names (default))
# should the terminal results be coloured to make them easier to read?
# keep this setting to false when using gpredomics as a binary
# uncomment to save experiment in timestamp-save_exp, which can be reloaded with --load timestamp-save_exp. Extension should be .json, .mp/.msq

# the features of the train data set
# the class description of the train data set (0=class 0, 1=class 1 (the class to be predicted), 2=unknown status)
# the features of the test data set
# the class description of the test data set
# 0: all significant features ; else first X significant features (per class!) sorted according to their pvalue/log_abs_bayes_factor
# per class, e.g. features are retained if any of the class reaches this level
# features which mean is below that value are discarded
# possible values are wilcoxon, student and bayesian_fisher. wilcoxon is recommended in most cases.
# features with differences less significant (p value above that threshold) than this will be removed
# features with a fewer log absolute bayes factor will be removed (bayesian method only)

# number of folds used to penalize overfit if overfit_penalty > 0
# this penalty is deduced from fit function (fit = fit on k-1 - abs(delta with last fold) * overfit_penalty)
# number of folds used for cross-validation (launch algorithm on each k-1 folds then merge Families of Best Models).
# if true, FBM is based on validation fold fit favouring generalisation, else on k-1 folds ; DO NOT fit on valid without external validation !
# alpha for the Family of Best Models confidence interval based on the best fit on validation fold. Smaller alpha, larger best_model range.
```

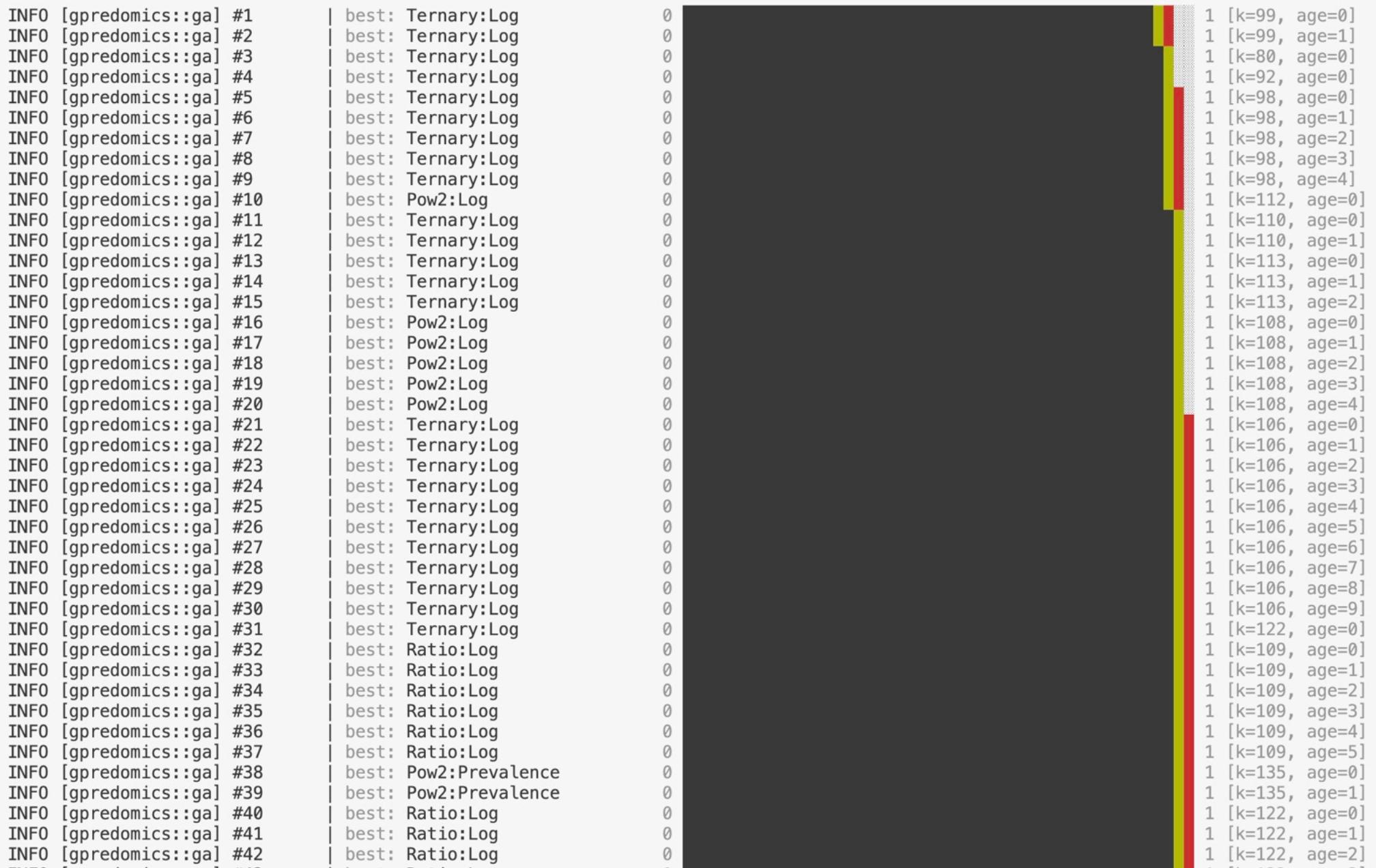
INFO [gpredomics] Training using Genetic Algorithm

INFO [gpredomics::data] Selecting features...

INFO [gpredomics::data] 447 features selected

INFO [gpredomics::ga] Population size: 4991, kmin 1, kmax 199

INFO [gpredomics::ga] Legend: [# diversity filter] [σ resampling] [█: AUC] [██: penalized fit]



Model #1 Ternary:Prevalence [k=18] [gen:138] [fit:0.868] AUC 0.914/0.533 | accuracy 0.857/0.512 | sensitivity 0.800/0.312 | specificity 0.886/0.708
 Class R: $(\text{msp_0134}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0850}^0 + \text{msp_0874}^0 + \text{msp_1060c}^0 + \text{omsp_2357}^0) \geq 2$

Model #2 Ternary:Prevalence [k=20] [gen:189] [fit:0.865] AUC 0.889/0.522 | accuracy 0.797/0.465 | sensitivity 0.889/0.469 | specificity 0.750/0.462
 Class R: $(\text{msp_0208}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0592}^0 + \text{msp_0657}^0 + \text{msp_0785}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1698}^0 + \text{omsp_1410}^0) - (\text{msp_0850}^0 + \text{msp_1060c}^0 + \text{msp_1556}^0 + \text{omsp_0212}^0 + \text{omsp_2357}^0) \geq 2$

Model #3 Ternary:Prevalence [k=19] [gen:109] [fit:0.864] AUC 0.910/0.497 | accuracy 0.789/0.504 | sensitivity 0.911/0.547 | specificity 0.727/0.462
 Class R: $(\text{msp_0134}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1698}^0 + \text{msp_1729}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0850}^0 + \text{msp_1060c}^0 + \text{msp_1556}^0 + \text{omsp_0212}^0 + \text{omsp_2357}^0) \geq 0$

Model #4 Ternary:Prevalence [k=17] [gen:190] [fit:0.861] AUC 0.913/0.540 | accuracy 0.850/0.543 | sensitivity 0.778/0.312 | specificity 0.886/0.769
 Class R: $(\text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0850}^0 + \text{msp_0874}^0 + \text{msp_1060c}^0 + \text{omsp_2357}^0) \geq 2$

Model #5 Ternary:Prevalence [k=20] [gen:190] [fit:0.861] AUC 0.912/0.493 | accuracy 0.842/0.527 | sensitivity 0.778/0.391 | specificity 0.875/0.662
 Class R: $(\text{msp_0134}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0850}^0 + \text{msp_0874}^0 + \text{msp_1060c}^0 + \text{msp_1381}^0 + \text{msp_1556}^0 + \text{omsp_2357}^0) \geq 1$

Model #6 Ternary:Prevalence [k=18] [gen:197] [fit:0.859] AUC 0.905/0.509 | accuracy 0.827/0.512 | sensitivity 0.822/0.422 | specificity 0.830/0.600
 Class R: $(\text{msp_0134}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1461}^0 + \text{msp_1729}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0850}^0 + \text{msp_0874}^0 + \text{msp_1060c}^0) \geq 2$

Model #7 Ternary:Prevalence [k=21] [gen:197] [fit:0.858] AUC 0.898/0.536 | accuracy 0.805/0.504 | sensitivity 0.867/0.484 | specificity 0.773/0.523
 Class R: $(\text{msp_0208}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0740}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1296c}^0 + \text{msp_1356}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{msp_1729}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0793}^0 + \text{msp_0850}^0 + \text{msp_1060c}^0 + \text{msp_1556}^0 + \text{omsp_0212}^0) \geq 2$

Model #8 Ternary:Prevalence [k=19] [gen:177] [fit:0.858] AUC 0.912/0.514 | accuracy 0.835/0.512 | sensitivity 0.867/0.422 | specificity 0.818/0.600
 Class R: $(\text{msp_0134}^0 + \text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{msp_1729}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0850}^0 + \text{msp_0874}^0 + \text{msp_1060c}^0) \geq 2$

Model #9 Ternary:Prevalence [k=21] [gen:101] [fit:0.856] AUC 0.901/0.521 | accuracy 0.812/0.504 | sensitivity 0.911/0.500 | specificity 0.761/0.508
 Class R: $(\text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0740}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1296c}^0 + \text{msp_1356}^0 + \text{msp_1461}^0 + \text{msp_1698}^0 + \text{msp_1729}^0 + \text{omsp_1410}^0) - (\text{msp_0173}^0 + \text{msp_0212}^0 + \text{msp_0793}^0 + \text{msp_0850}^0 + \text{msp_1060c}^0 + \text{msp_1556}^0 + \text{omsp_0212}^0) \geq 1$

Model #10 Ternary:Prevalence [k=15] [gen:138] [fit:0.854] AUC 0.883/0.485 | accuracy 0.782/0.481 | sensitivity 0.933/0.531 | specificity 0.705/0.431
 Class R: $(\text{msp_0254}^0 + \text{msp_0294}^0 + \text{msp_0414}^0 + \text{msp_0785}^0 + \text{msp_0942}^0 + \text{msp_1048}^0 + \text{msp_1146}^0 + \text{msp_1244}^0 + \text{msp_1356}^0 + \text{msp_1399}^0 + \text{omsp_1410}^0) - (\text{msp_0850}^0 + \text{msp_1060c}^0 + \text{omsp_0212}^0 + \text{omsp_2357}^0) \geq 1$

VOTING ANALYSIS

Majority jury [14 experts] | AUC 0.923/0.503 | accuracy 0.859/0.504 | sensitivity 0.932/0.484 | specificity 0.821/0.524 | rejection rate 0.038/0.031

DETAILED METRICS

CONFUSION MATRIX (TRAIN)

	Pred 1	Pred 0	Abstain
Real 1	41	3	1
Real 0	15	69	4

CONFUSION MATRIX (TEST)

	Pred 1	Pred 0	Abstain
Real 1	30	32	2
Real 0	30	33	2

PREDICTIONS BY SAMPLE (TEST)

Sample	Real	Predictions	Result	Consistency
--------	------	-------------	--------	-------------

ERRORS (12 shown of 62, sorted by inconsistency)						
SAMN37671286	0	01110110011100	→ 1	x		57.1%
SAMN37671581	1	01100110101000	→ 0	x		57.1%
SAMN37671687	0	01101010101101	→ 1	x		57.1%
SAMN37671700	0	001001000111111	→ 1	x		57.1%
SAMN37671479	0	11011111000011	→ 1	x		64.3%
SAMN37671506	0	011001100111111	→ 1	x		64.3%
SAMN37671669	1	000000000101111	→ 0	x		64.3%
SAMN37671683	1	001000101000110	→ 0	x		64.3%
SAMN37671688	0	011001101011111	→ 1	x		64.3%
SAMN37671705	0	01111110101100	→ 1	x		64.3%
SAMN37671715	1	01110000101000	→ 0	x		64.3%
SAMN37671732	0	011001100111111	→ 1	x		64.3%

ABSTENTIONS (4 shown of 4, sorted by inconsistency)						
SAMN37671287	1	001001000111111	→ 2	~		50.0%
SAMN37671511	0	001010001010111	→ 2	~		50.0%
SAMN37671682	1	010100101010111	→ 2	~		50.0%
SAMN37671694	0	010100101010111	→ 2	~		50.0%

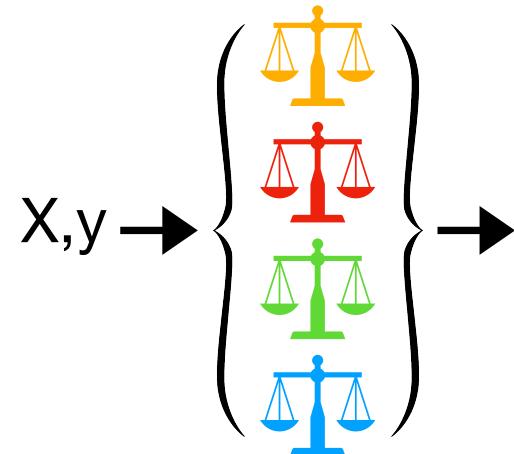
CORRECT (4 shown of 63, sorted by inconsistency)						
SAMN37671280	1	00111111100100	→ 1	✓		57.1%
SAMN37671512	0	01100000101101	→ 0	✓		57.1%
SAMN37671519	1	10101111001001	→ 1	✓		57.1%
SAMN37671522	0	001000001011111	→ 0	✓		57.1%

... 109 additional samples not shown

What else ?



Multi-class
classification
OVO/OVA



Fabien
Kambu

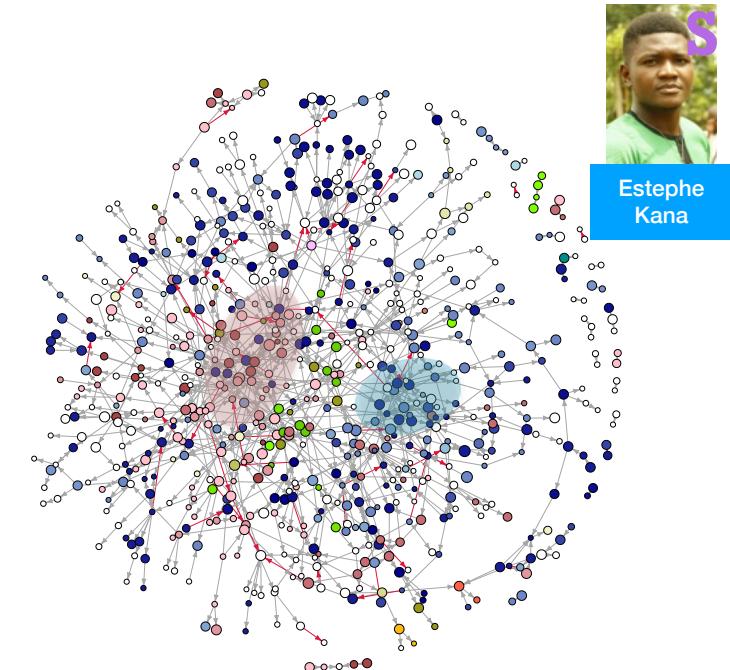
GPU/rust acceleration
Model robustness
Industrial applications



Louison
Lesage

Raynald de
Lahondès

Ecosystem awareness
of predictive signatures



Estephe
Kana

A shiny web application

Welcome to Predomics Studio

Predicting and giving sense to your data
You are currently logged in as Gray

[Log in](#) [Sign up](#)

Login (or email) :

Password :

[Get started](#) [Continue as guest](#)

PredOmics

UNHEALTHY

HEALTHY

Best Predomics Models Summary

Here you can see a quick summary of your best Predomics models by type of Predomics learner, click on the corresponding button to see the best algorithm of a type.

The best model in your analysis is found with learner terBeam and language terinter . It predicts a class of -1 with auc_ of 0.8836 and a sparsity of 5 It is described by the equation :

$(+ SP_142 + SP_291) - (+ SP_391 + SP_214 + SP_397) > -0.00088$

TerBeam - Ter

TerDA - Bin

The best model in your analysis with learner terBeam uses language terinter . It predicts a class of -1 with auc_ of 0.8836 and a sparsity of 5 It is described by the equation :

$(+ SP_142 + SP_291) - (+ SP_391 + SP_214 + SP_397) > -0.00088$



Gaspar
Roy

Take-away messages ...

- **See the Data, don't trust blindly the models and p-values**
- The gut microbiome's role in our health is still under-appreciated mostly because of challenges with standards, quantification, reference and catalogs...
- AI is powerful in numerous fields, but can be **challenging in biomedical data**, which are highly resolute but **with a limited number of observations**.
- **Interpretability and clinical validation** are key for downstream applications.
- **Ad hoc specialised tools and methods** such as predomics... need to be developed and **adapted to the fields of study**. Highly constrained models such as BTR @predomics are both **extremely simple and accurate and Interpretable**
- There is always room for improvement. **Don't let a good idea die, follow up... it usually pays off**



Do you want to **put your skills, time and effort into exciting and useful projects that can save lives?** This is what we offer you by joining a **dynamic, agile and passionate team.**

We are hiring
3 postdocs +
1 PhD in

Position 1: (2 year postdoc; DL + ECG)

- Deep learning engineer: transformer models **embedding** research, **robustness** and **interpretability** for ECG (electrocardiograms) clustering in the context of sudden death risk.
- Apply (**CV + letter + 3 references**) to talents@ummisco.fr

Position 2: (3 year postdoc; DL + ECG)

- Data scientist in AI and ECG. Analysing one of the world's largest and diverse ECG datasets.
- Apply (**CV + letter + 3 references**) to talents@ummisco.fr

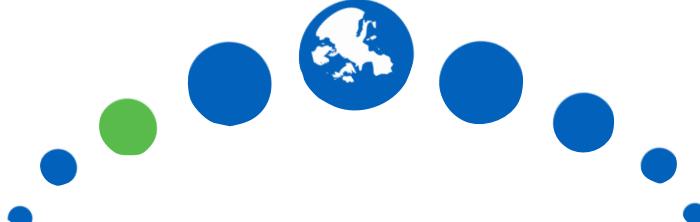
Position 3: (3 year postdoc; DL + ECG)

- Data manager and biostatistics. Analysing one of the world's largest and diverse ECG datasets.
- Apply (**CV + letter + 3 references**) to talents@ummisco.fr

Position 4: (3 year; PhD; DL + ECG)

- PhD position: multimodal transformer based models **for ECG** research.
- Apply (**CV + letter + 3 references**) to talents@ummisco.fr

+ Multiple internships



Thanks

Merci

Faleminderit

