# "Enhancing Interpretability of Global Plankton Communities Modeling through Multi-Omics and AI Techniques"

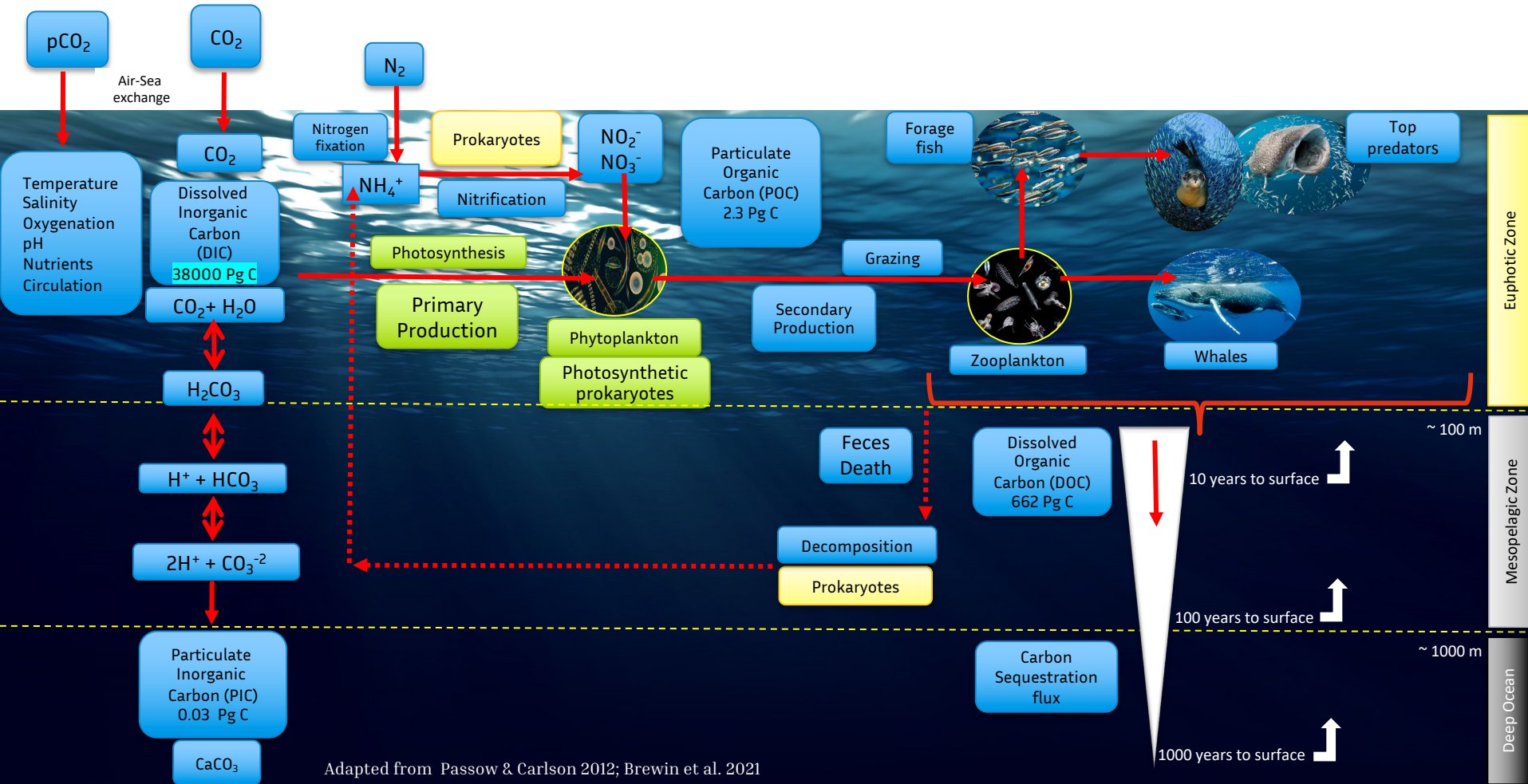Luis Valenzuela

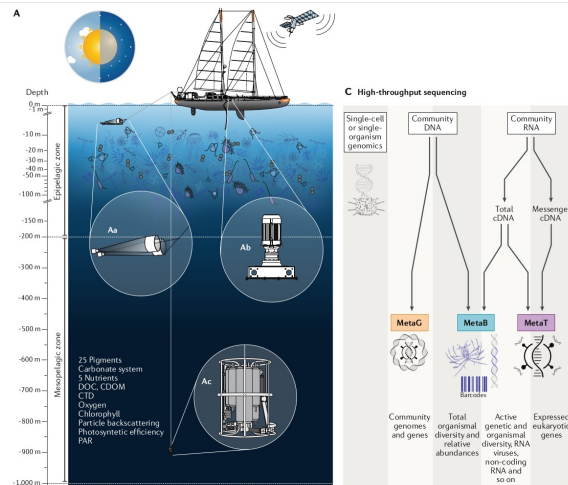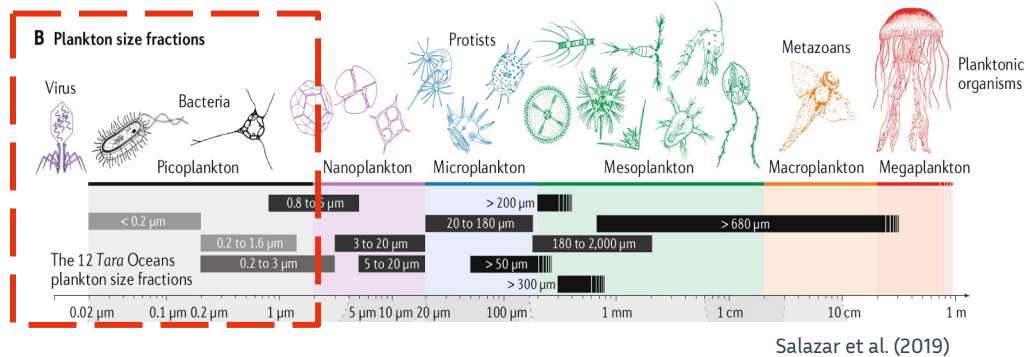Workshop: Towards a modern analysis of omics data of the Ocean

# Climate Change, Oceanic Carbon pumps and plankton roles

pCO$_2$

CO$_2$

Air-Sea exchange

N$_2$

Nitrogen fixation

Prokaryotes

NO$_2^-$ NO$_3^-$

CO$_2$

Temperature Salinity Oxygenation pH Nutrients Circulation

Dissolved Inorganic Carbon (DIC) 38000 Pg C

NH$_4^+$

Nitrification

Particulate Organic Carbon (POC) 2.3 Pg C

Forage fish

Top predators

CO$_2$+ H$_2$O

Photosynthesis

Grazing

Euphotic Zone

Primary Production

Phytoplankton

Secondary Production

Zooplankton

Whales

H$_2$CO$_3$

Photosynthetic prokaryotes

~ 100 m

H$^+$ + HCO$_3$

Feces Death

Dissolved Organic Carbon (DOC) 662 Pg C

10 years to surface

Mesopelagic Zone

2H$^+$ + CO$_3^{-2}$

Decomposition

Prokaryotes

100 years to surface

~ 1000 m

Particulate Inorganic Carbon (PIC) 0.03 Pg C

Carbon Sequestration flux

Deep Ocean

CaCO$_3$

Adapted from  Passow & Carlson 2012; Brewin et al. 2021

1000 years to surface

Salazar et al. (2019)



Sunagawa et al. (2020)

Ocean Microbial Reference Gene Catalog v2:

**Metagenomic dataset:**

~ 57.000 million reads (90 pb ±2.6pb)

$ATGC...CTGGG_{90}$  ···  $TTCA...TTTCC_{90}$

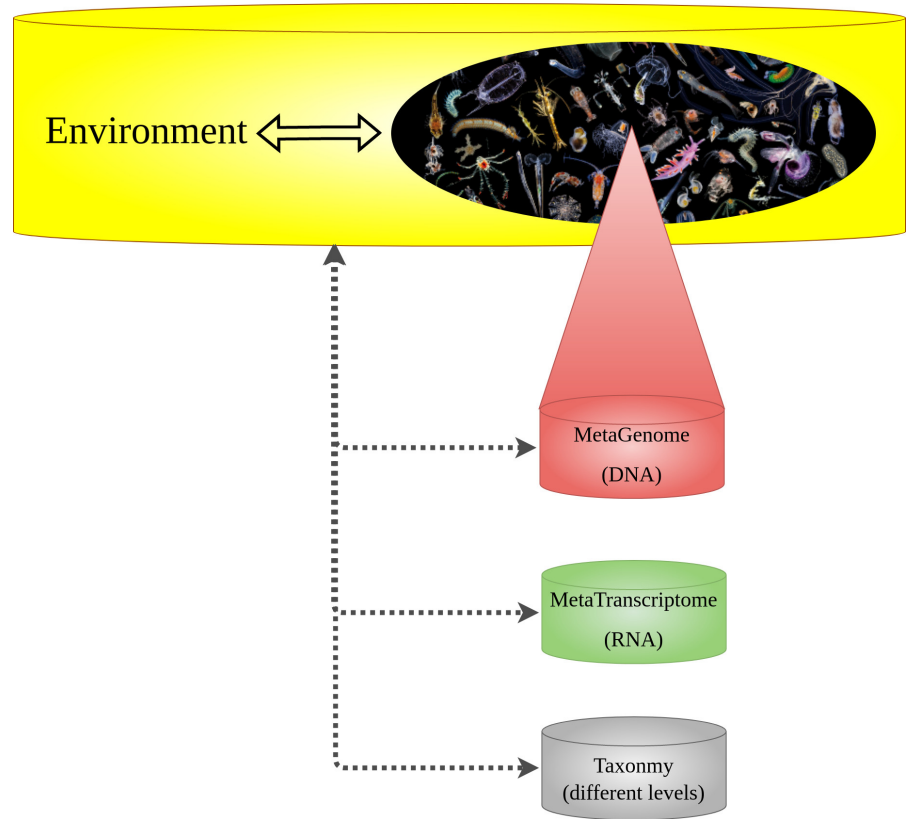$GCCC...AAAAG_{90}$  $GGGGA...AGCTA_{90}$

(meta)genome assembly

~ 200 metagenomes:

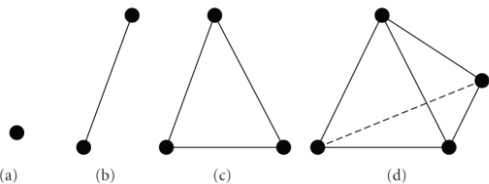~ 230000 scaffolds, N50: 1300pb (average)

~ 42 million genes

Functional annotation
Quantification
Normalization

|  | $Gene_1$ | $Gene_2$ | ... | $Gene_{42M}$ |
|---|---|---|---|---|
| $metaG_1$ |  |  |  |  |
| $metaG_2$ |  |  |  |  |
| ... |  |  |  |  |
| $metaG_{200}$ |  |  |  |  |

Salazar et al., Gene expression changes and community turnover
differentially shape the global ocean metatranscriptome, Cell, 2019.

https://www.ebi.ac.uk/biostudies/studies/S-BSST297

Environment ⟺

MetaGenome
(DNA)

MetaTranscriptome
(RNA)

Taxonmy
(different levels)

## K-Simplicial Complexes



(a) 0-simplex is a vertex, (b) 1-simplex is a line,
(c) 2-simplex is a triangle, (d) 3-simplex is a tetrahedron.

## Persistent Homology filtrations



r = 0.6     r = 0.8     r = 1.2



Compute a graphical representation
of the dataset

Learn an embedding that preserves
the structure of the graph

https://umap-learn.readthedocs.io/en/latest/

An Autoencoder is a neural architecture designed to learn an identity function in an unsupervised way to reconstruct the original input while (usually) compressing the data in the process to discover a more efficient representation.

Shapley value is the average expected marginal contribution of one feature after all possible combinations have been considered.



Coalition
$(C_{1234})$

Coalition value
$(V_{1234})$

How much contribute each one to the coalition valuue?

$C_{1234}$

$C_{234}$

$Value_{1234}$

$Value_{234}$

$V_{1234} - V_{234} =$ Marginal contribution of member 1 to C234

Shapley values are the mean marginal contribution.

$$\varphi_i = \frac{1}{\#\ Members} \sum_{\forall C\ s.t.\ i \neg \in C} \frac{Marginal\ Contribution\ of\ i\ to\ C}{\#\ Coalitions\ of\ size\ |C|}$$

Lundberg & Lee, 2018. A Unified Approach to Interpreting Model Predictions. 31st Conference on Neural Information Processing Systems.

SHAP is a game theoretic approach to explain the output of any machine learning model.

**SHAP Kernel** approximates Shapley values through much fewer samples.



$$W_C = \frac{\# \, total \, features - 1}{\# \, coalitions \, of \, size \, |C| * \#included \, features \, in \, C \, * \#excluded \, features \, in \, C}$$

Lundberg & Lee, 2018. A Unified Approach to Interpreting Model Predictions.
31st Conference on Neural Information Processing Systems.

# Exploration of the Omics datasets

Metagenomic assays: Functional Composition

# Metatranscriptomic assays: Gene expression

**i) KEGG db: 8937 features**

**ii) eggNOG db 71662 features**

**iii) eggNOG+GC db 314715 features**

Taxonomy datasets

Domain level · Phylum level · Class level · Order level · Family level

Genus level

Polarity

Layer

$$\mathbb{R}^{9024} \rightarrow \mathbb{R}^{3}$$

Clustering

min dist: 0.1

[10 models]

Sankey plot: exploring min dist hyper-parameter:
gif: 0.0, 0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40, 0.45, 0.50

[110 models]

Polar Data:
(42, 9027)
SRF 22
DCM 10
MES 10

Non Polar Data:
(131, 9027)
SRF 61
DCM 42
MES 28

All Data:
(173, 9073)
SRF 83
DCM 52
MES 38

Cluster0    36
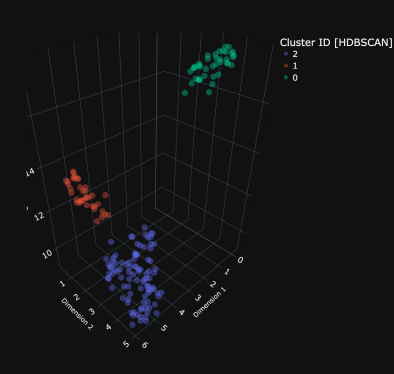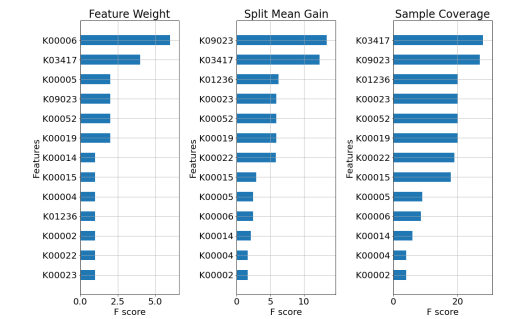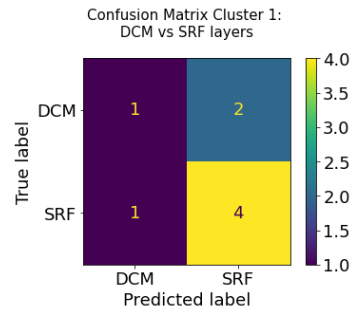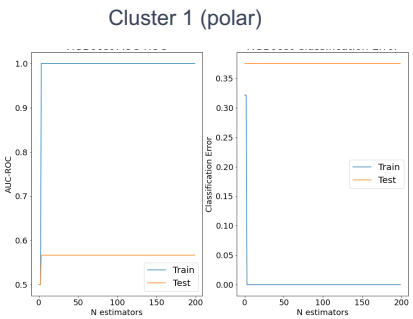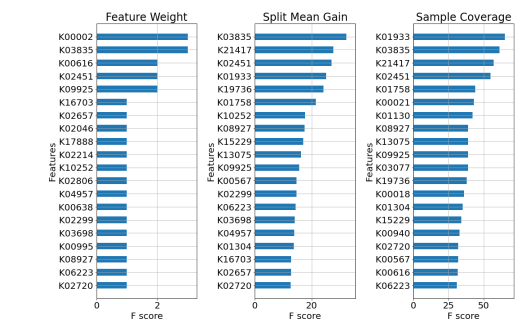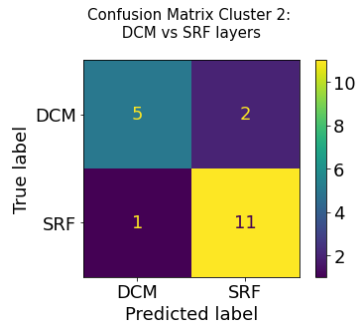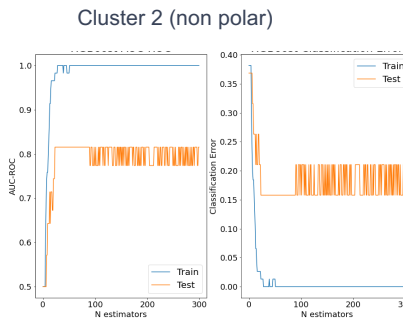Cluster1    32
Cluster2    95

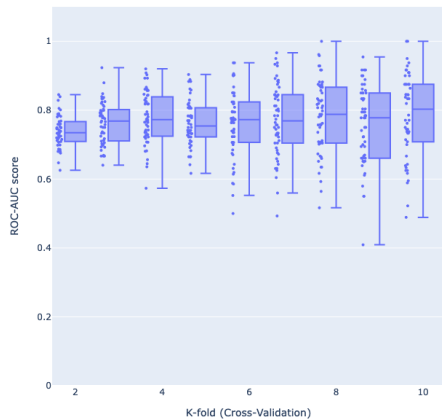Cluster 0 vs Cluster 1 (mes vs polar)

Cluster 1 vs Cluster 2 (polar vs non polar)

Cluster 0 vs Cluster 2 (mes vs non polar)

General exploration: Binary Classification [KEGG data]

No parametric test: Wilcoxon, adjusted p value via fdr.
1808/9027 features with adj. P val < 0.05 as input for XGBoost hyperparameter tuning.



SRF vs DCM [all data]
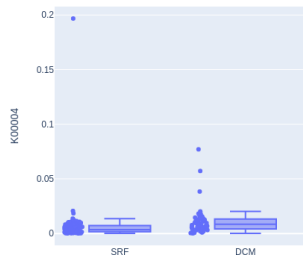(significant KOs)

SRF vs DCM [all data]
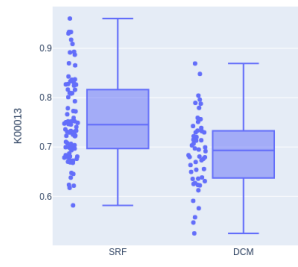(significant KOs)

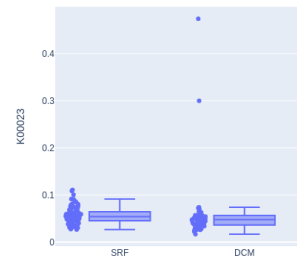Classification Performance (SRF vs DCM)
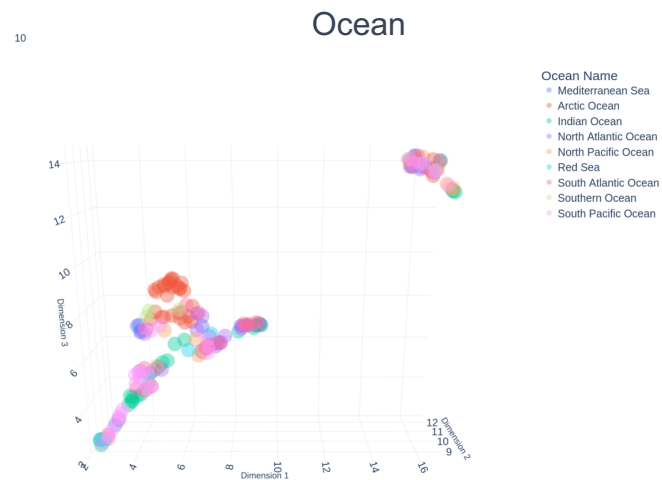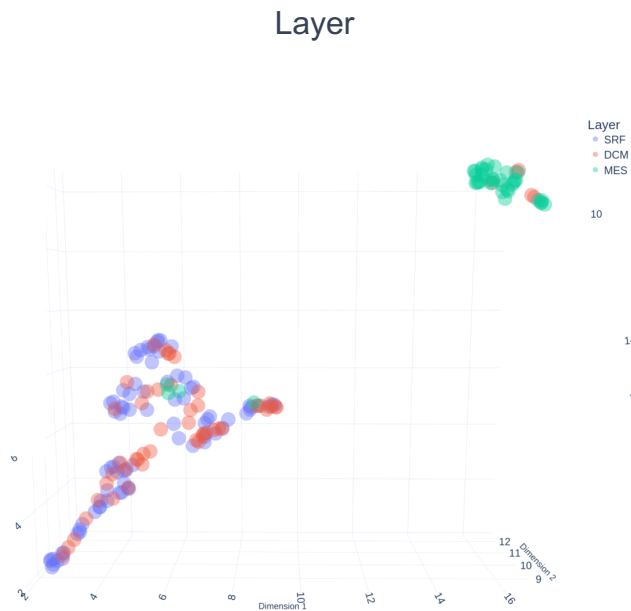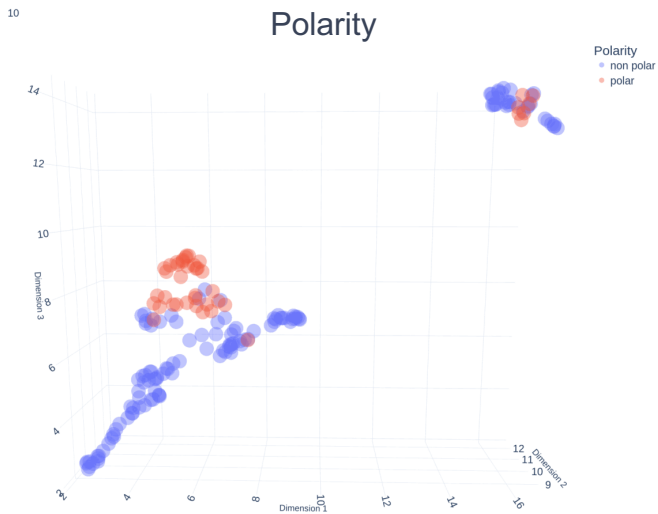
Homoserine dehydrogenase

Butanediol dehydrogenase

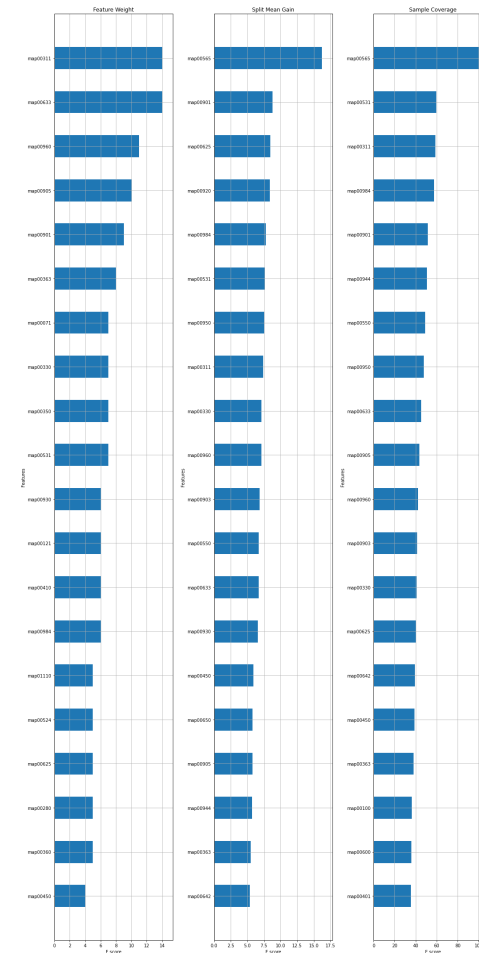histidinol dehydrogenase

Acetoacetyl-CoA reductase

# Transcription Factors regulating KEGG pathways

## Polarity

Polarity
- non polar
- polar

## Layer

Layer
- SRF
- DCM
- MES

## Ocean

Ocean Name
- Mediterranean Sea
- Arctic Ocean
- Indian Ocean
- North Atlantic Ocean
- North Pacific Ocean
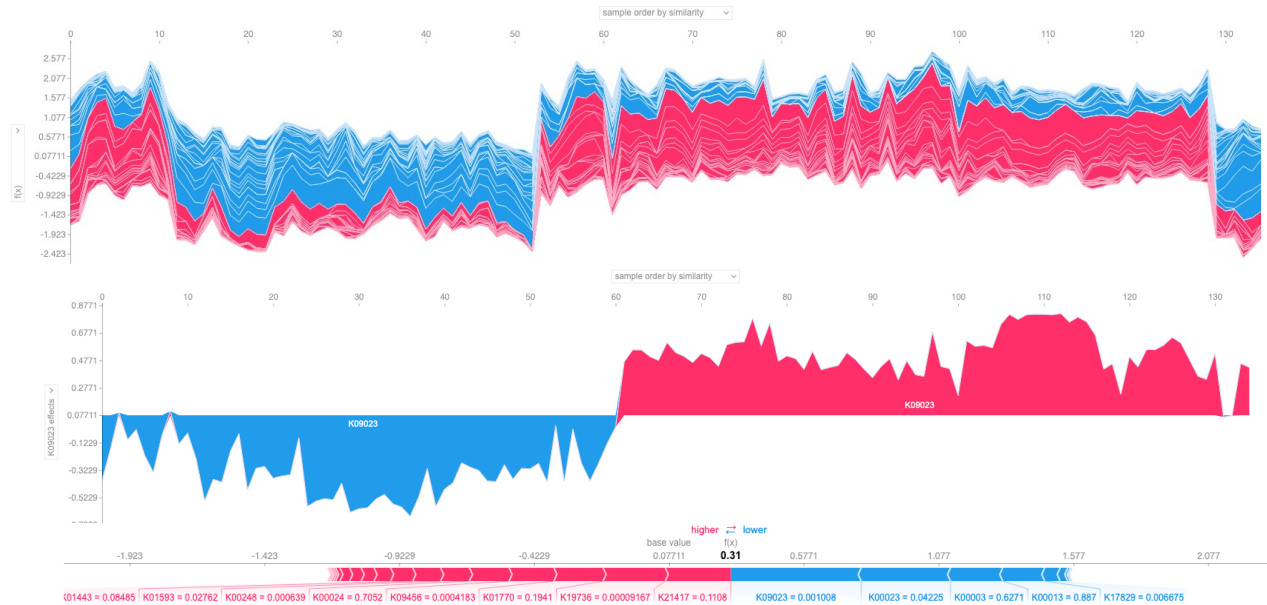- Red Sea
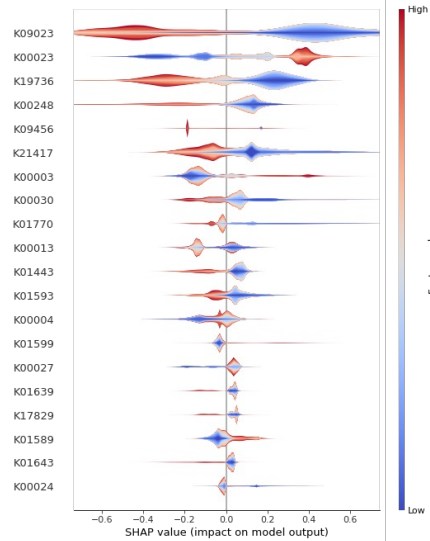- South Atlantic Ocean
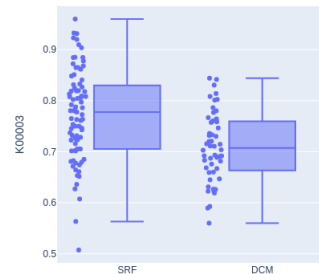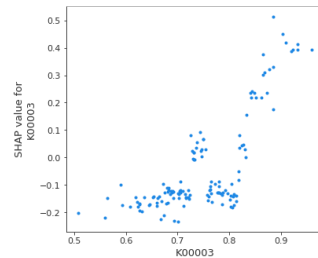- Southern Ocean
- South Pacific Ocean

# Maps matrix
173 x 163

Map00311: Penicillin and cephalosporin biosynthesis
Map00633: Nitrotoluene degradation
Map00960: Tropane, piperidine and pyridine alkaloid biosynthesis
Map00905: Brassinosteroid biosynthesis
Map00901: Indole alkaloid biosynthesis
Map00363: Bisphenol degradation
Map00071: Fatty acid degradation
Map00330: Arginine and proline metabolism
Map00350: Tyrosine metabolism
Map00930: Caprolactam degradation
Map00121: Secondary bile acid biosynthesis
Map00410: beta-Alanine metabolism
Map00984: Steroid degradation
Map01110: Biosynthesis of secondary metabolites
Map00524: Neomycin, kanamycin and gentamicin biosynthesis
Map00625: Chloroalkane and chloroalkene degradation
Map00280: Valine, leucine and isoleucine degradation
Map00360: Phenylalanine metabolism
Map00450: Selenocompound metabolism

Aminoacrylate hydrolase (pyrimidine metabolism)

TetR/AcrR family transcriptional regulator

No parametric test: Wilcoxon, adjusted p value via fdr.
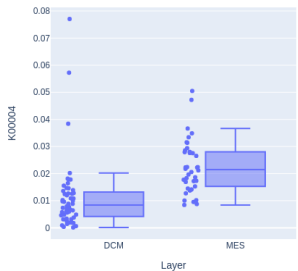6937/9027 features with adj. P val < 0.05 as input for XGBoost hyperparameter tuning.
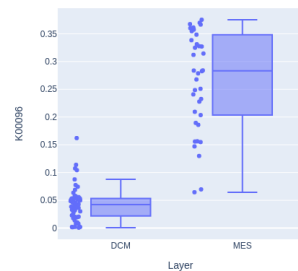


Classification Performance (DCM vs MES)



Homoserine dehydrogenase

Butanediol dehydrogenase

glycerol-1-phosphate dehydrogenase

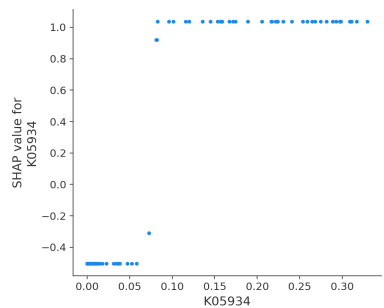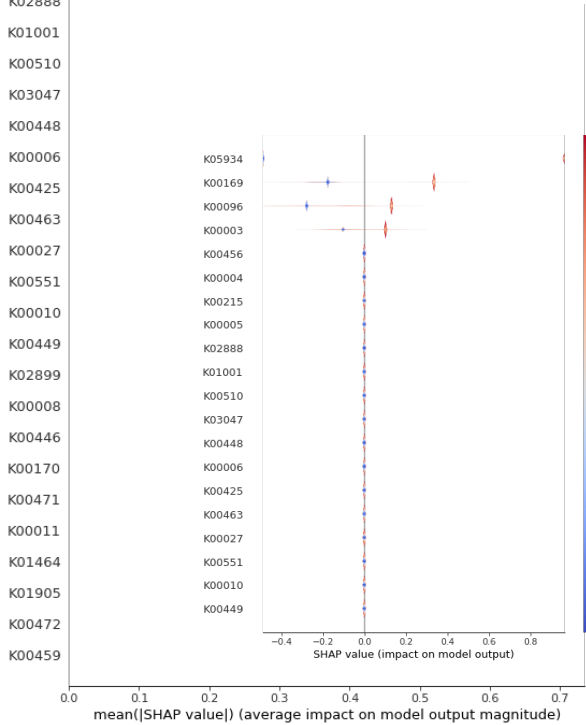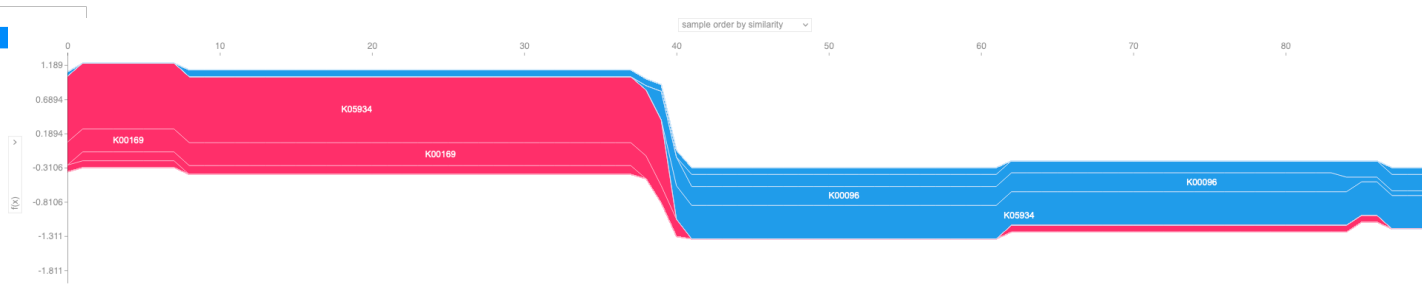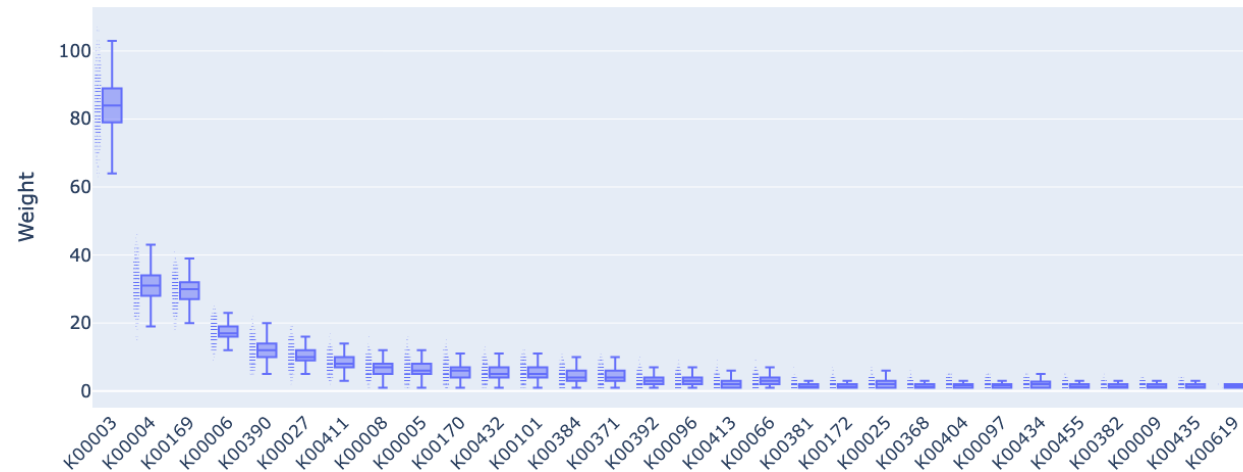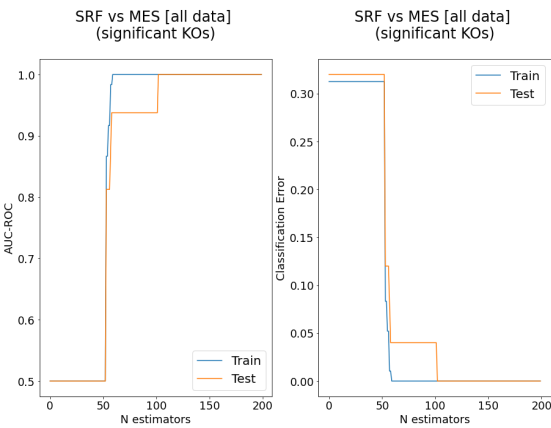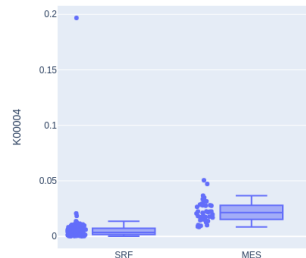glycerol dehydrogenase

Cobalamin biosynthesis

No parametric test: Wilcoxon, adjusted p value via fdr.
7517/9027 features with adj. P val < 0.05 as input for XGBoost hyperparameter tuning.
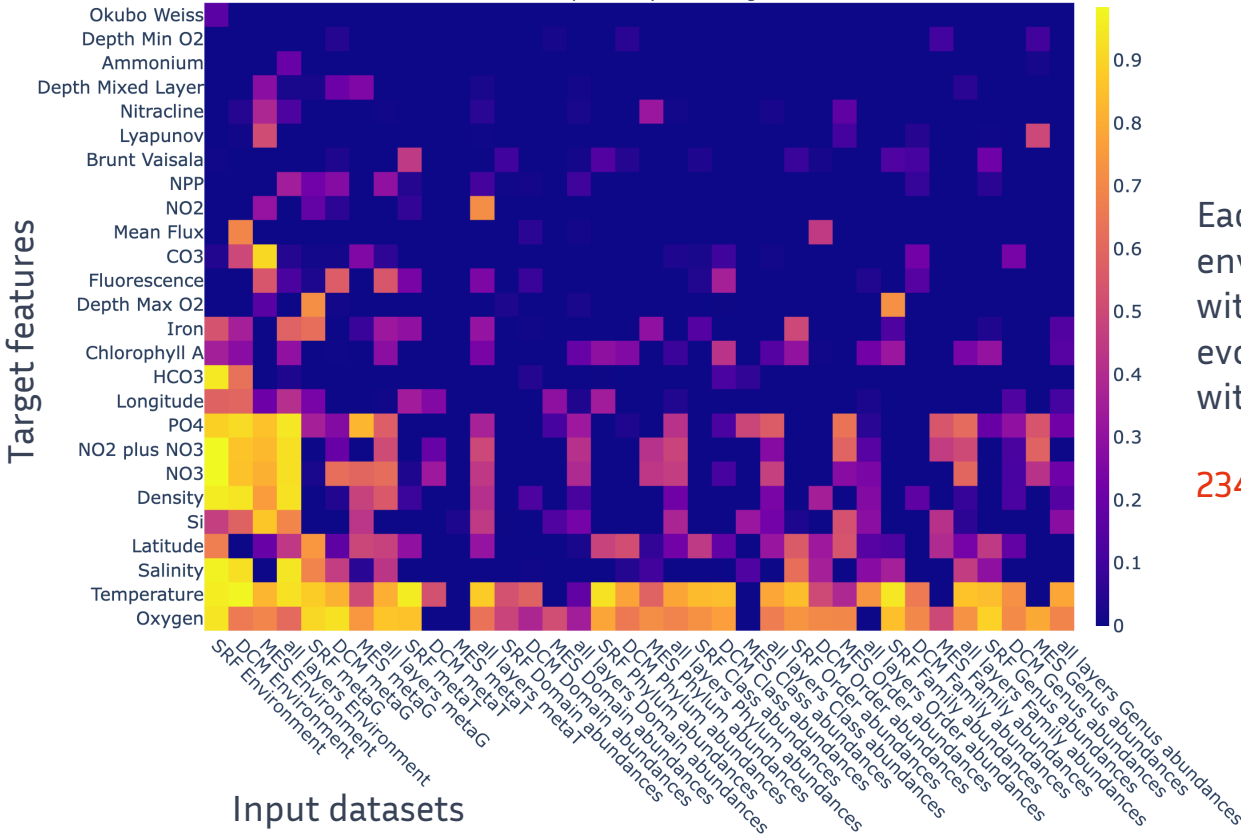
pyruvate ferredoxin oxidoreductase

Y = (X8/X1) * (X9 - 0.500)



Population Initialization

Tournaments and Selections

Reproduction (mutations + crossover)

Termination

Subtree Mutation

Hoist Mutation

Point Mutation

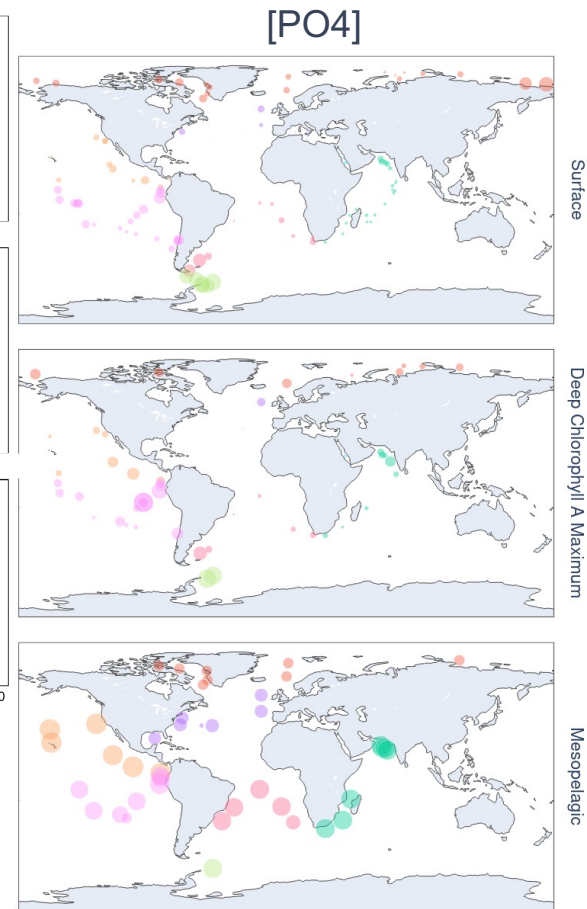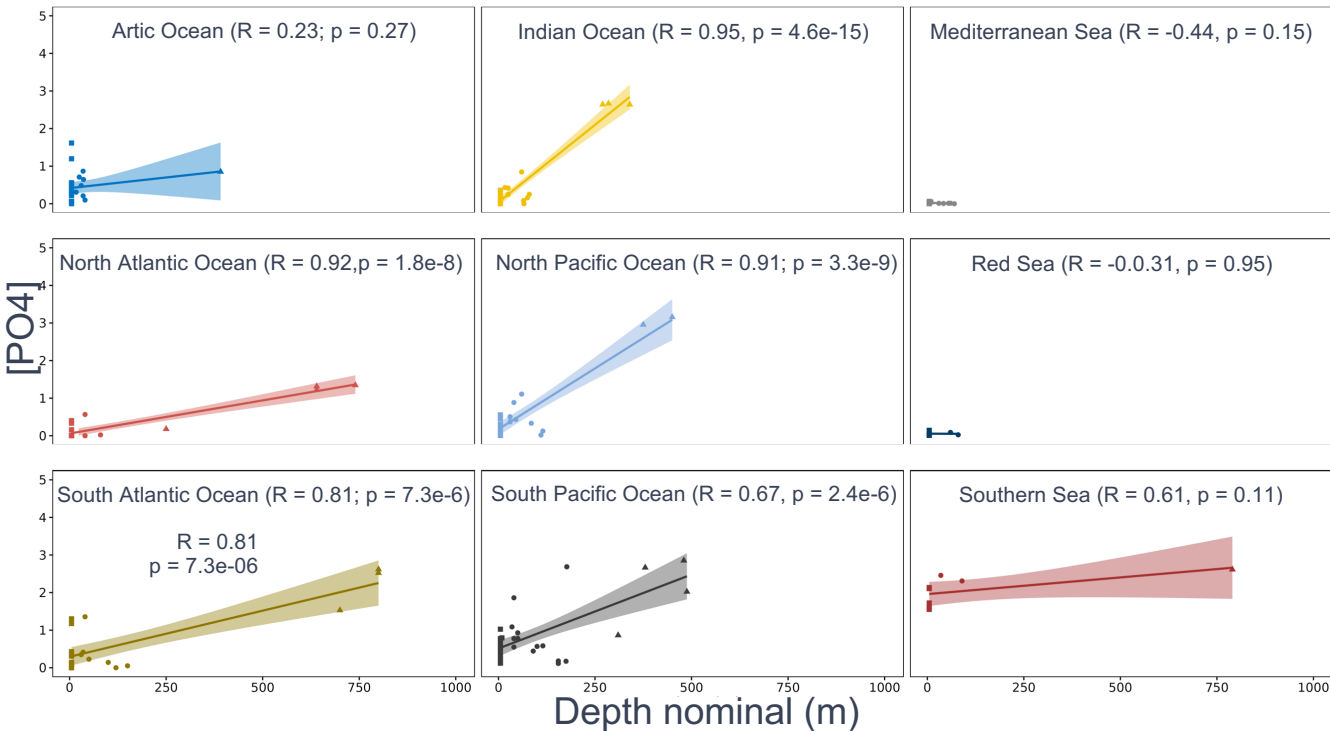R² Coefficients Heatmap from Symbolic Regressions Models

Each dataset was used to predict one environmental variable at a time starting with a **population size of 20000**, and evolving the models during **20 generations** with a crossover probability of 0.65.
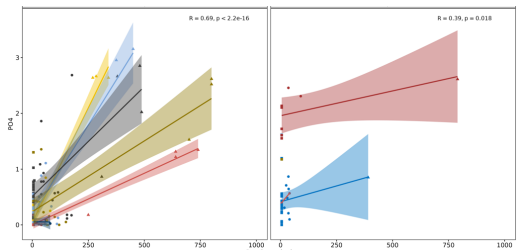
234 models

$$Target_i = SymReg_{GeneticProgramming}(dataset_j)$$

Exploration: Phosphate concentration vs Depth nominal (m)

# PO4 prediction from metaG dataset

$$PO4 = K00315 - K03040 + K07483 - K09701 + \log(\log(K06982 * K19229))$$

$(R^2: 0.66)$

Input: 'K09023', 'K00023', 'K19736', 'K00248','K09456','K21417','K00003'
Output: 1 environmental feature.

Symbolic Regressions: population_size=20000, generations=20, p_crossover=0.65



Latitude

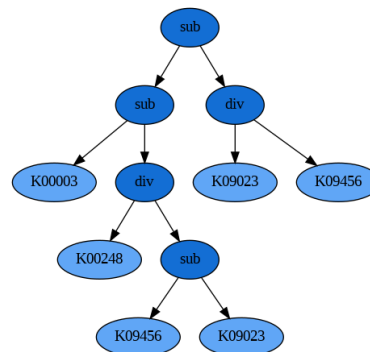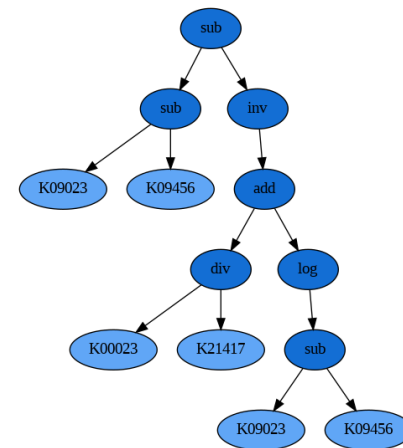log((K09456**0.5/(-K00003*K00248/K09456 + K21417))**0.5)

(R$^2$: 0.6600)

Temperature

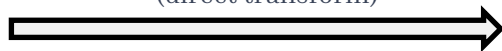log((K09456**0.5/(-K00003*K00248/K09456 + K21417))**0.5)

(R$^2$: 0.6688)

Salinity

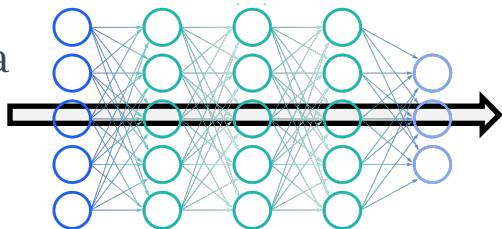log((K09456**0.5/(-K00003*K00248/K09456 + K21417))**0.5)

(R$^2$: 0.5883)

Meta-Genomic data

$R^{9024}$

(direct transform)

$\longrightarrow$

$R^3$ 3D embedding

Environmental data

$R^{29}$



$\longrightarrow$

$R^3$ 3D embedding

reconstruction

3D embedding

reconstruction

$R^3$

(reverse transform)

$\longrightarrow$

$R^{9024}$ Meta-Genomic data

reconstruction

Transcriptome Autoencoder

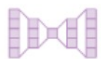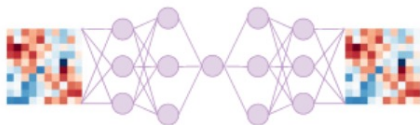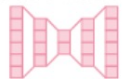Transcriptome- to-Transcriptome
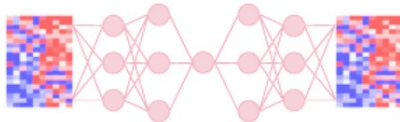
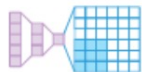Microbiome Autoencoder

Microbiome- to-Microbiome

Image Autoencoder

Image- to-Image

Transcriptome-to-Image Autoencoder

Transcriptome- to-Image
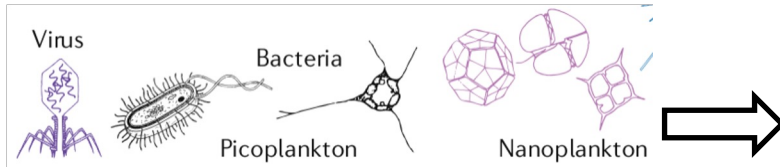
Image-to-Transcriptome Autoencoder

Image-to-Transcriptome

Challenges for implementations of generative models
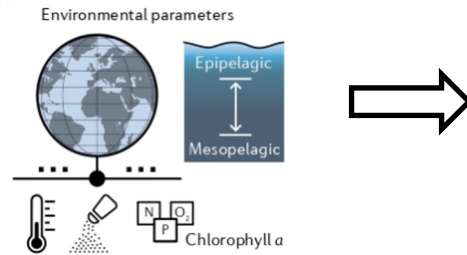
Data augmentation estrategy:

Variational autoencoders

Optimal Transport

Edited and adapted from Yang, K. D., Belyaeva, A., Venkatachalapathy, S. (2021). Multi-domain translation between single-cell imaging and sequencing data using autoencoders. Nature Communications, 12, 31. doi: 10.1038/s41467-020-20249-2.

Thank you!