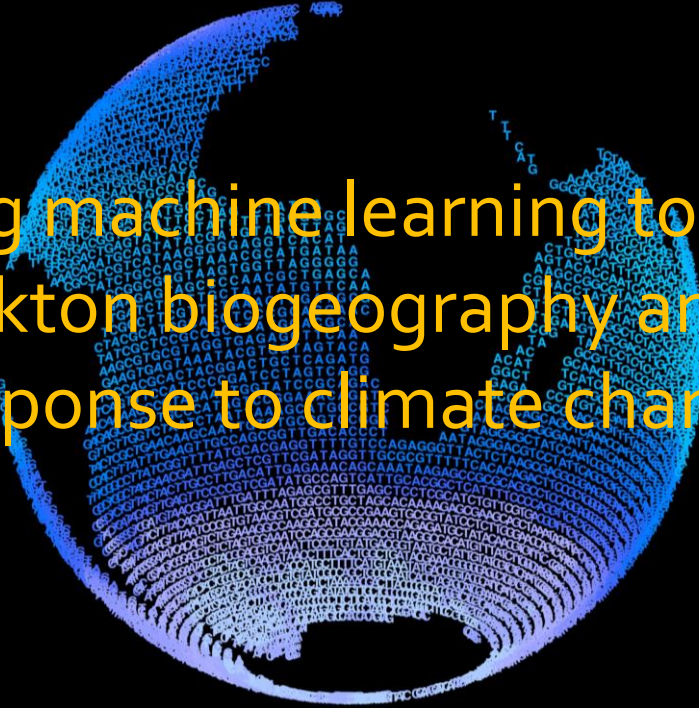




CEODEOS/AtlantEco Workshop  
Valparaiso  
15-19 May 2023



# Using machine learning to infer plankton biogeography and its response to climate change



Paul Frémont  
pfrémont@gatech.edu  
www.paulfremont.com



# Introduction: Ecological niche

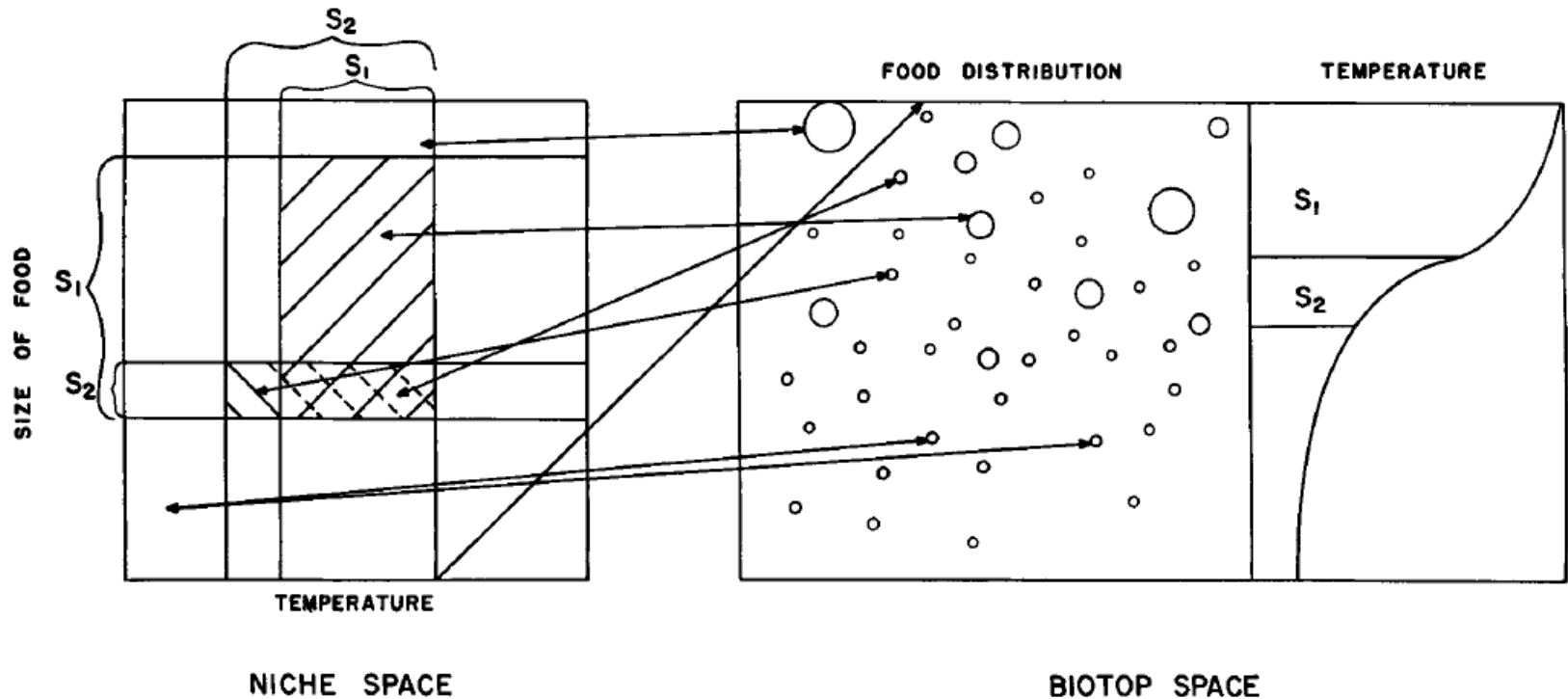


FIGURE 1. Two fundamental niches defined by a pair of variables in a two-dimensional niche space. Only one species is supposed to be able to persist in the intersection subset region. The lines joining equivalent points in the niche space and biotope space indicate the relationship of the two spaces. The distribution of the two species involved is shown on the right hand panel with a temperature depth curve of the kind usual in a lake in summer.

# Introduction: Ecological niche

## Fundamental niche:

all the *theoretical* environmental conditions where a species can live

## Realized niche:

all the environmental conditions where a species is *observed* to live

	Fundamental Niche	Realized Niche
Where the organism lives	No	Yes
Size	Large	Small
Competition for resources, predators are present	No	Yes
Other terminology	Precompetitive niche	Postcompetitive niche

⇒ In real life, when sampling an organism, we look at the *realized niche*

# Introduction: How can we calculate a model of *realized* ecological niche?

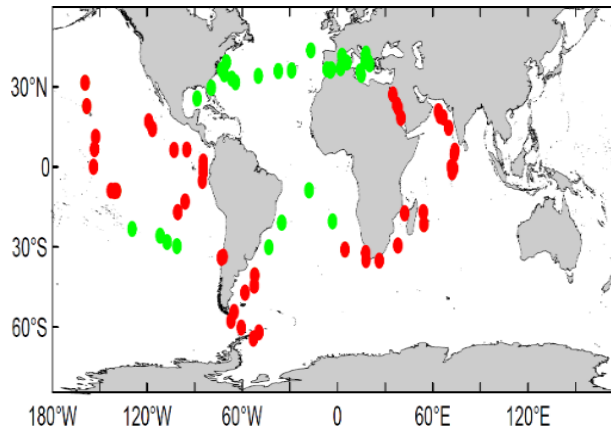
How does supervised learning works?

A model is fitted on a learning set:

**Y: Presence/absence** of a species in each *Tara* Oceans stations

**X: Associated physicochemical parameters (predictors)**

⇒ The model defines the environmental parameters in which the species lives *i.e.* it reflects the combination of parameters in which a species was found

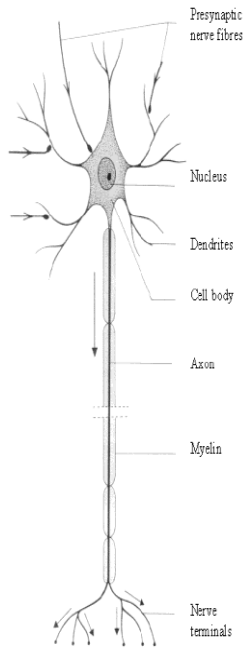


● Presence  
● Absence

$$Y = f(X)$$

We are looking for  $f$

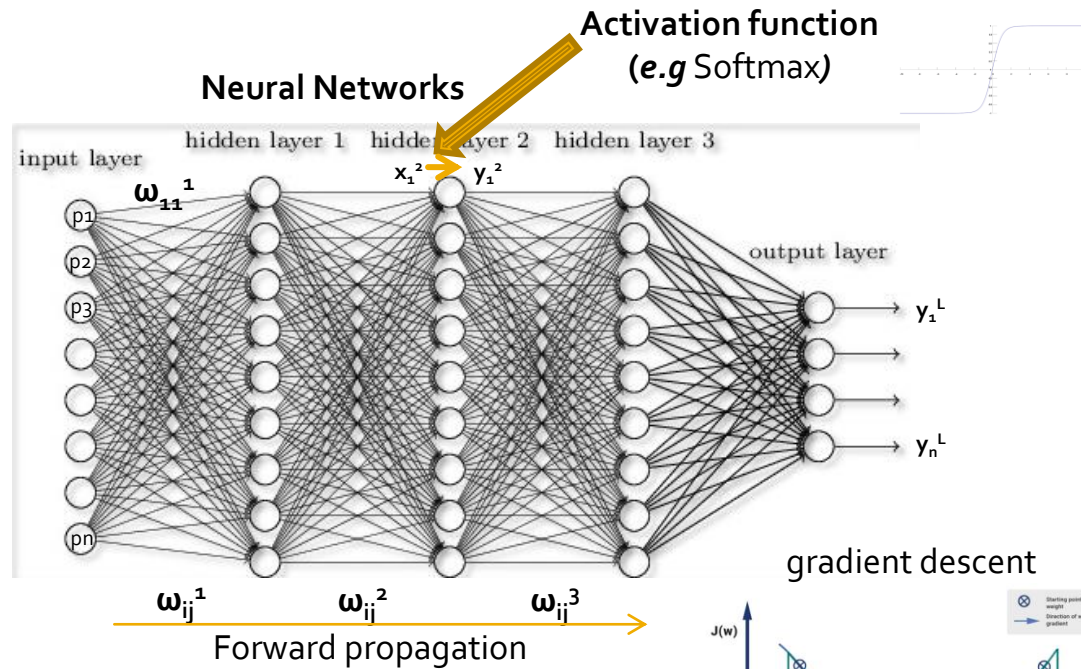
# Neural networks



**Error function** *E.g.*  
squared residuals, cross-entropy

$$\frac{\partial E}{\partial w_{ij}^l} = y_i^l \frac{\partial E}{\partial x_j^{l+1}}$$

**Back propagation**  
(calculate derivative of error with respect to weights)

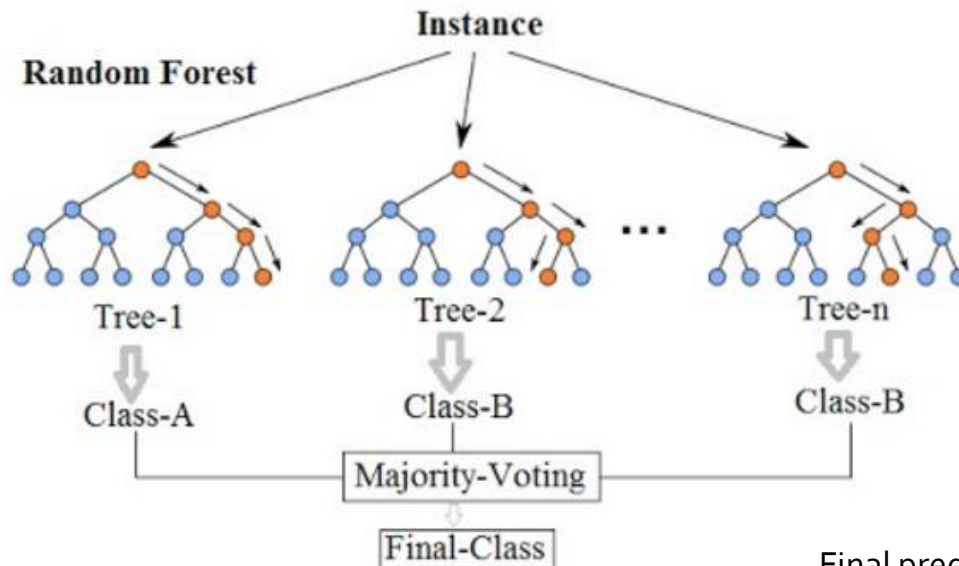


**gradient descent**



# Random forest

## Random Forest Simplified

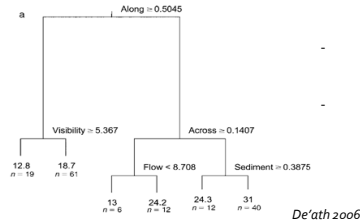


Set of B trees built based on bagging: each decision tree is built based on a subset of the set X, Y  
+  
Feature bagging: each tree is built on a random subset of the predictors

Final prediction is the average of the decision trees

$$\hat{f} = \frac{1}{B} \sum_{b=1}^B f_b(x')$$

# Other types of models



- Single trees are poor classifiers and large trees are hard to interpret => idea of combining several trees = **Boosting**
- Each tree is optimized by a least square method => **Regression**

Example of a single decision tree

A boosted regression tree is the sum of weighted trees:

$$y = f(x) = \sum_m f_m(x) = \sum_m \beta_m b(x; \gamma_m)$$

$y$  is the variable to be predicted

$x$  are the predictors

$b(x; \gamma_m)$  represent a tree :  $\gamma_m$  are the split variables,  $x$  their values at each node

$\beta_m$  is the 'learning rate': weight given at each tree (ranges from 1 to 0,001)

## Boosted Regression Trees

- In the regression setting, a generalized additive models has the form:

$$E(Y|X_1, X_2, \dots, X_p) = \alpha + f_1(X_1) + f_2(X_2) + \dots + f_p(X_p). \quad (9.1)$$

Here  $f_i$ 's are unspecified smooth and nonparametric functions.

Instead of using LBE in chapter 5, we fit each function using a scatter plot smoother(e.g. a cubic smoothing spline)

## Generalized Additive Models

MAXENT,  
Support Vector  
Machine  
Regression,  
Lasso regression,  
Linear regression  
etc

# How can we say that a model performs well? Cross validation

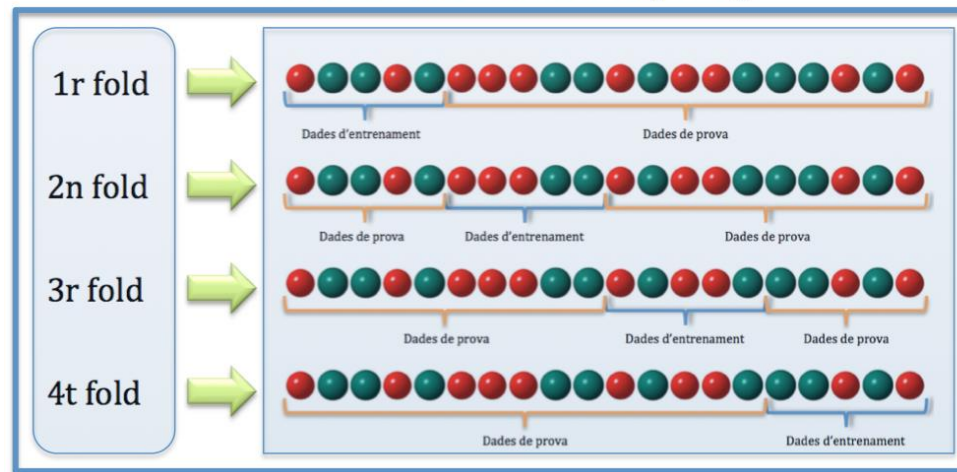
Most of the times, models **overfit data** => how can we deal with this?

⇒ Cross validation

⇒ Hyperparameter optimization and model performance assessment

For a given set of hyperparameters:

## K-cross fold validation (K=4)



Mean model performance over the k-fold (RMSE, AUC etc)

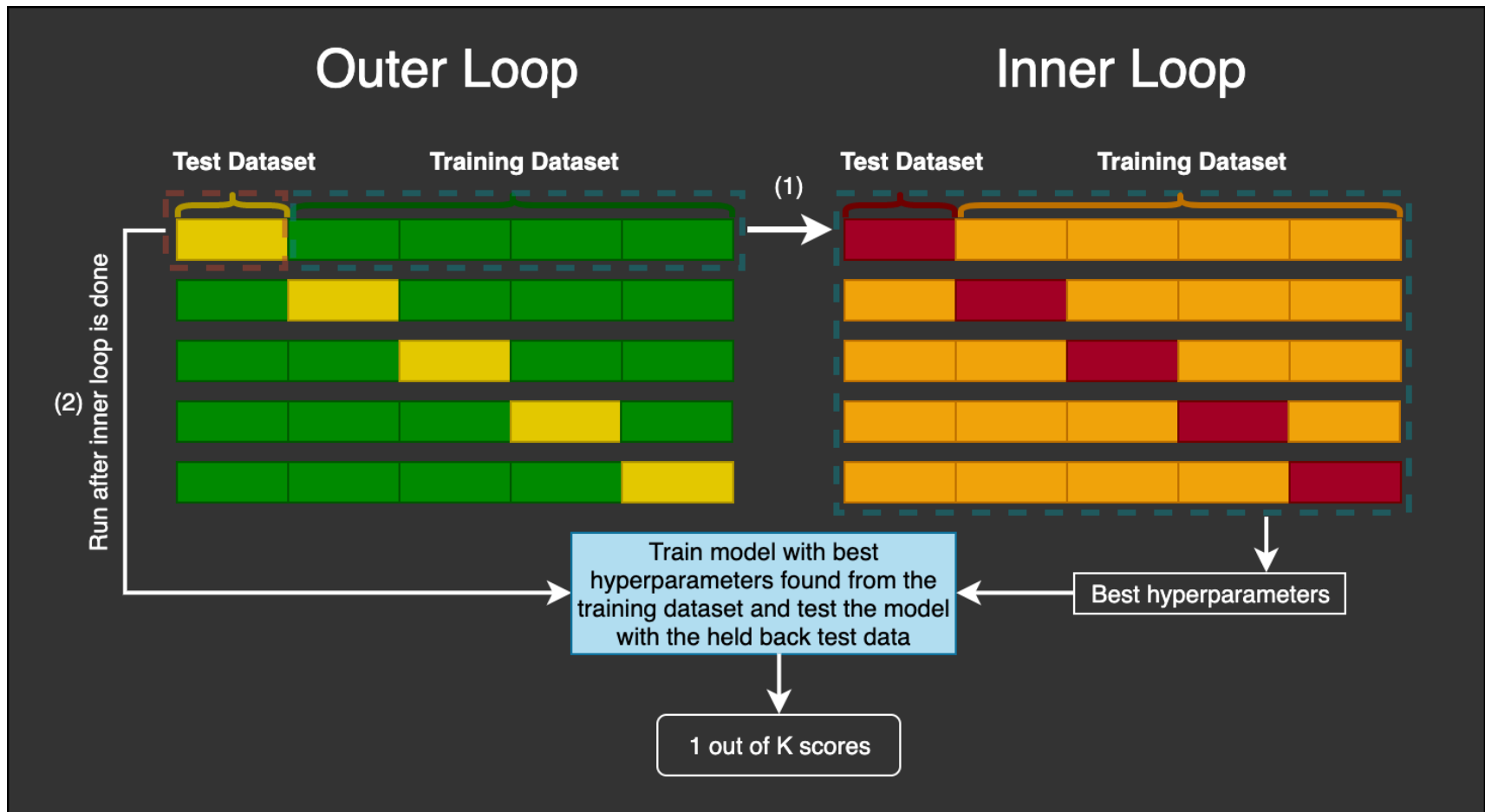
Choose set of hyperparameter that has the best performance for the final model training

Other possibilities: Monte Carlo cross validation (random subsamples), exhaustive cross validation (leave-p-out cross validation), **nested cross validation**



# Nested cross validation

Problem: Flat cross validation uses some points on which it has been trained to estimate model performance => positive bias on the model performance assessment



# Overzealous?



arXiv > cs > arXiv:1809.09446

Search...

Help |

Computer Science > Machine Learning

[Submitted on 25 Sep 2018]

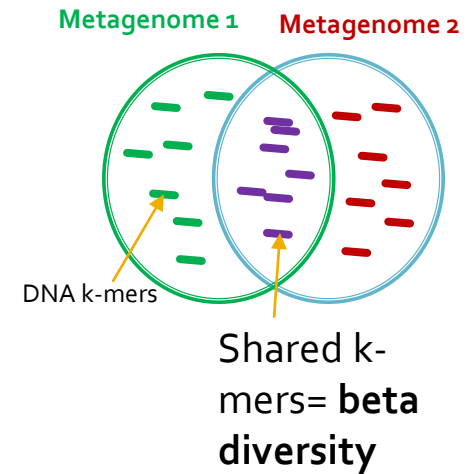
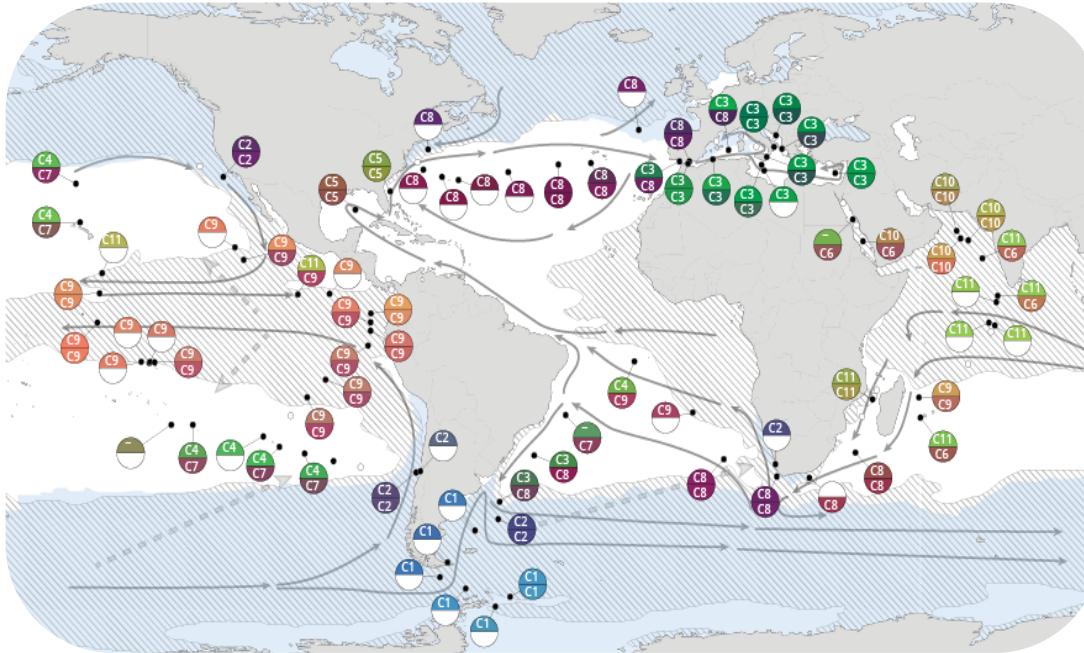
## Nested cross-validation when selecting classifiers is overzealous for most practical applications

Jacques Wainer, Gavin Cawley



# A concrete example: Tara Ocean genomic biogeography

## Genomic provinces



Richter *et al.* 2022

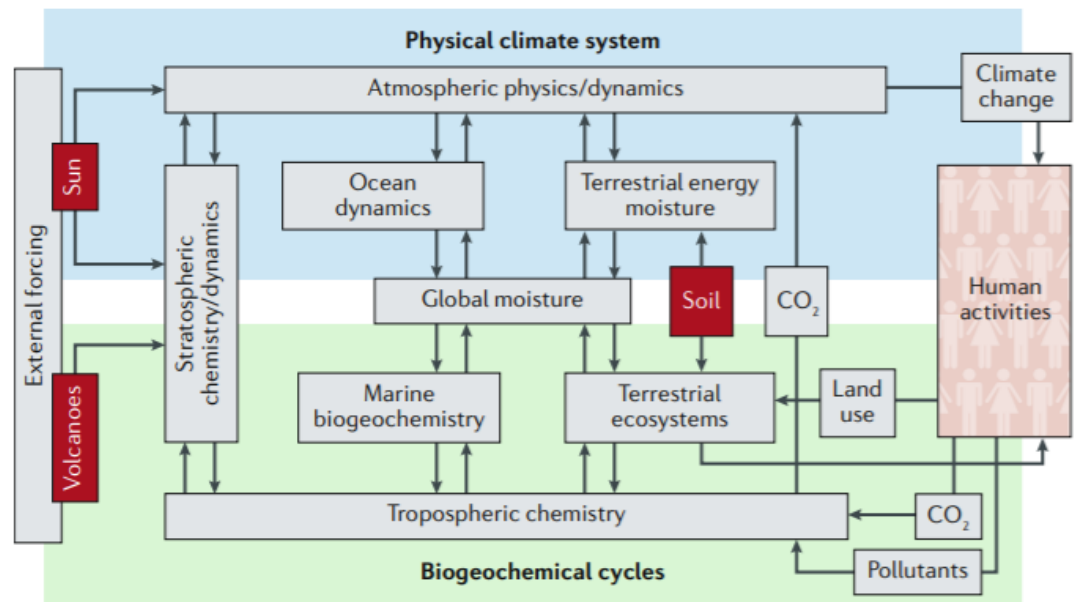
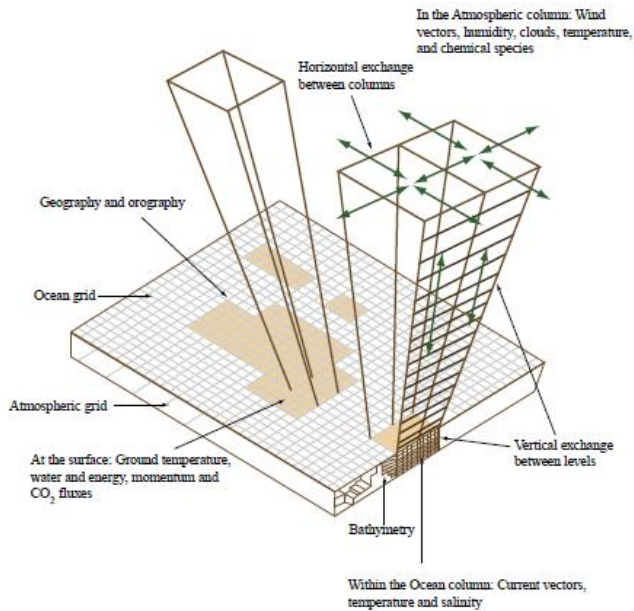
Metagenomic Biogeography of size fraction 0.8-5  $\mu\text{m}$   
definition of **genomic provinces**

⇒ **Hypothesis of associated environmental niches**

# A concrete example: Tara Ocean genomic biogeography

## Earth system science

### Earth System Science



The NASA Bretherton diagram of the Earth System.

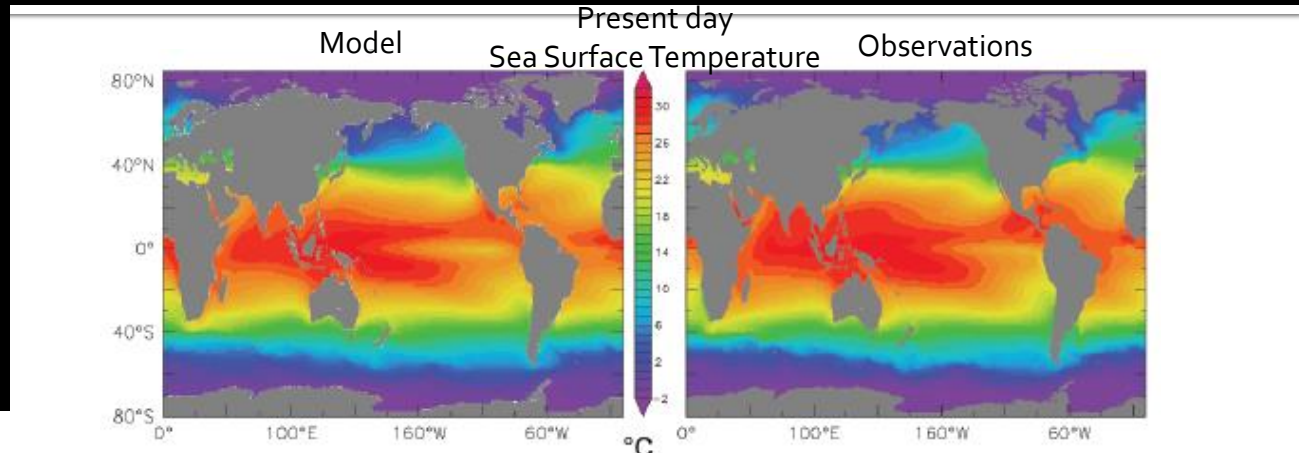
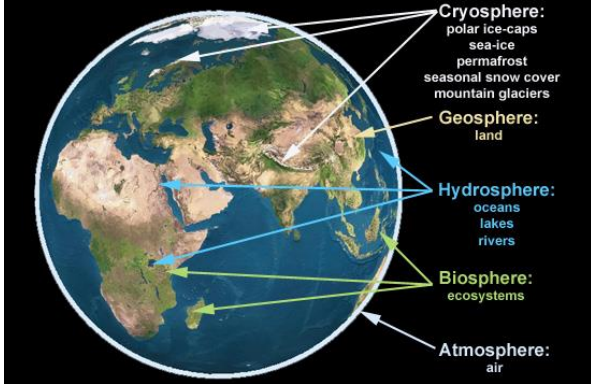
Steffen et al. 2020

"The IPSL Climate Model (IPSL-CM) is developed since 1995."

# A concrete example: Tara Ocean genomic biogeography

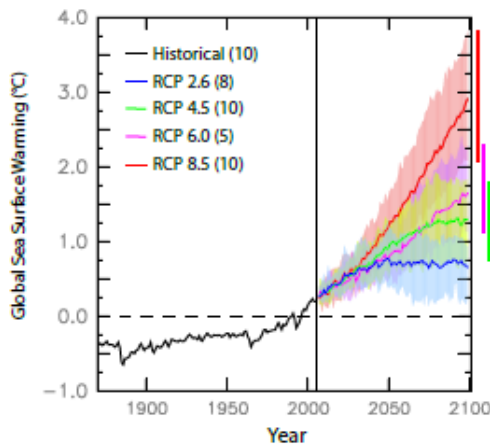
## Climate models of the ocean

The components of Earth's Climate System

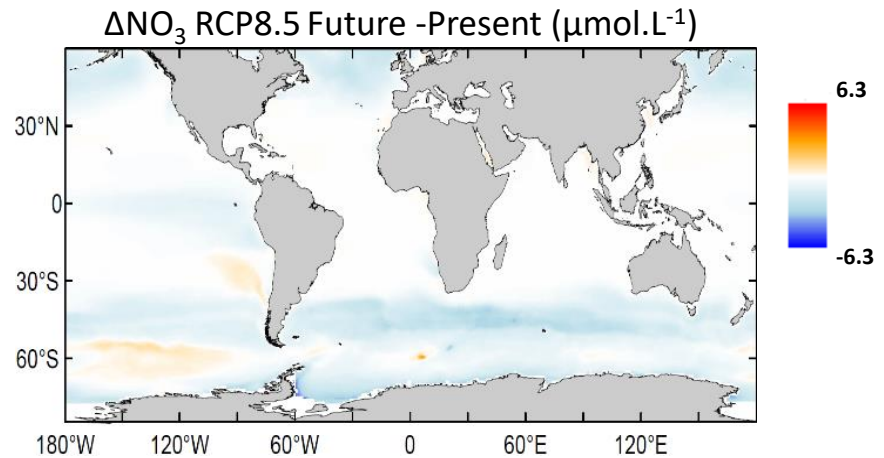


⇒ Complex climate models reproducing present day climatologies

*Bopp et al. 2013*



Projections of ocean warming

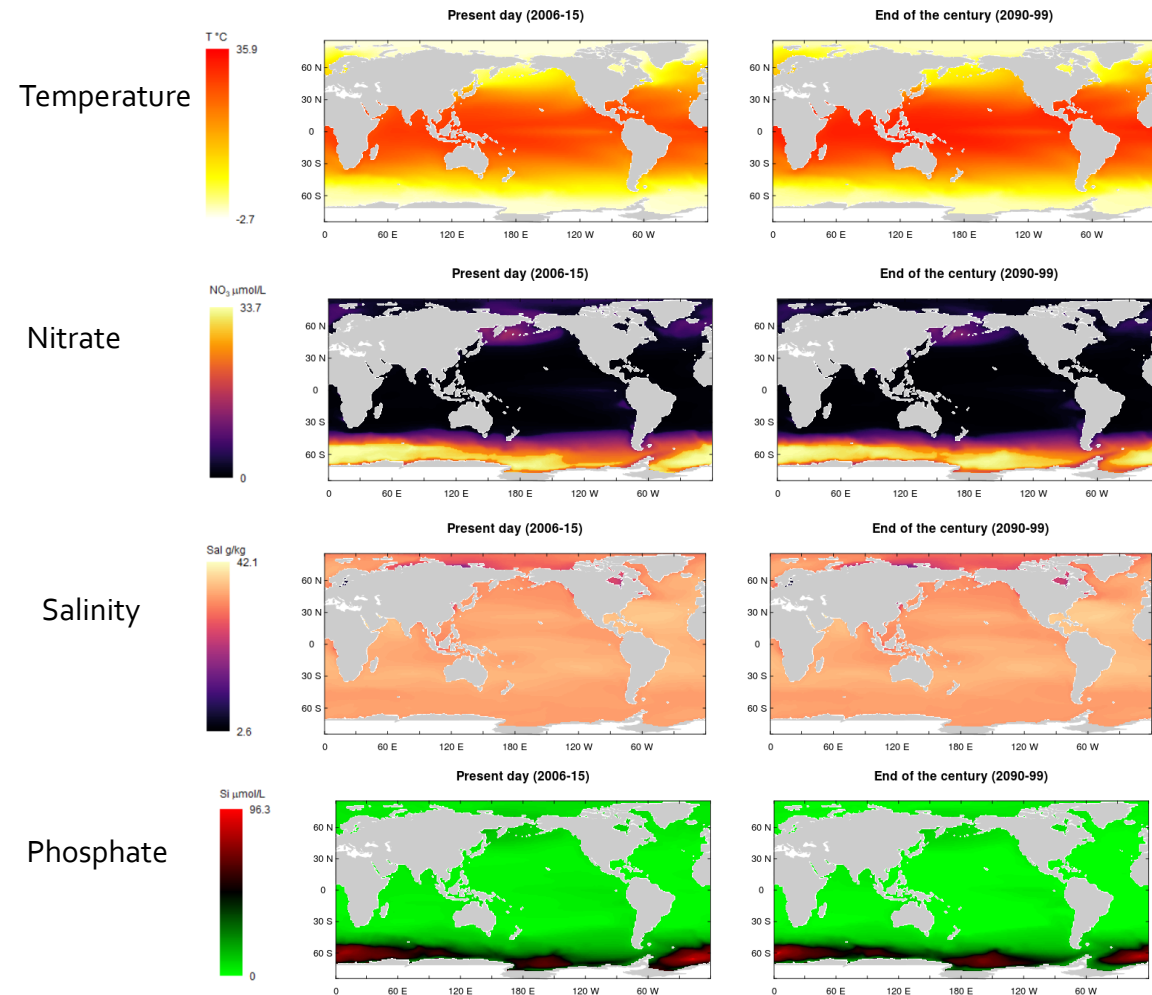


⇒ These models allow exploring the possible hypothetical futures depending on different greenhouse gas emission scenarios

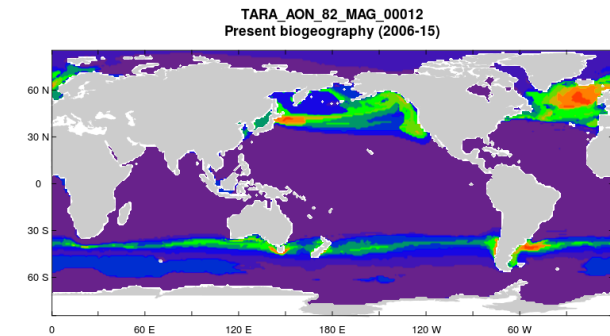
⇒ We focus on the 'business as usual' scenario RCP8.5

# What are ocean climate models good for?

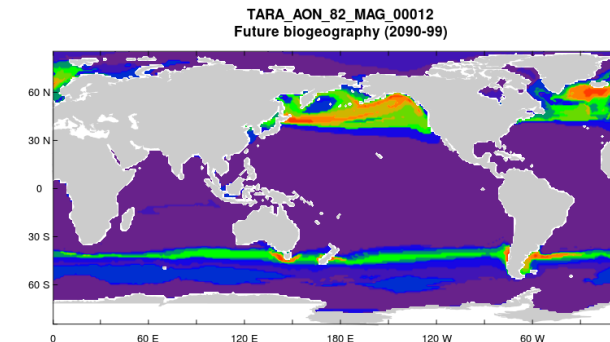
Project **SDMs** on new physicochemical datasets and evaluate climate change effects



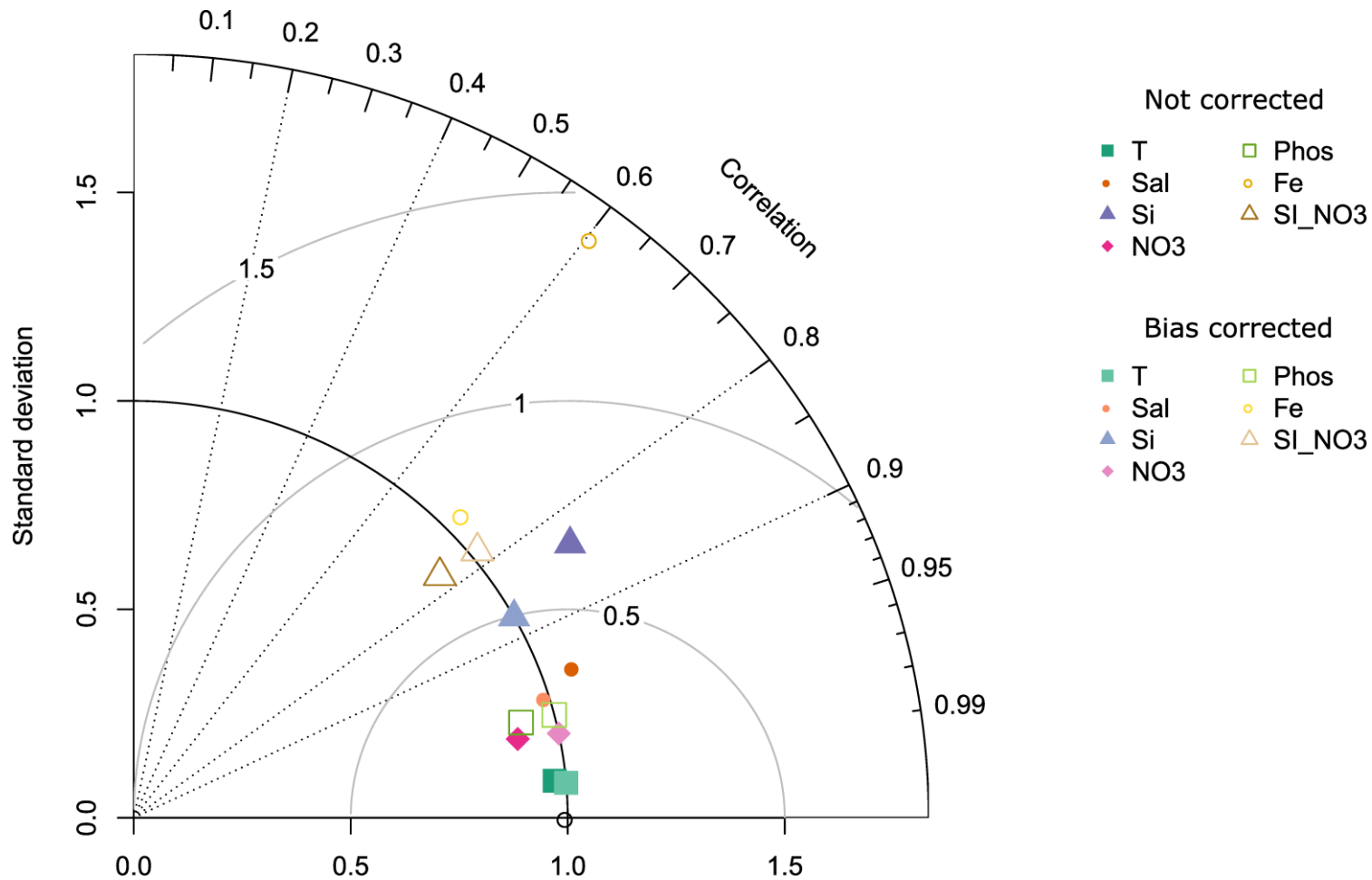
Present day



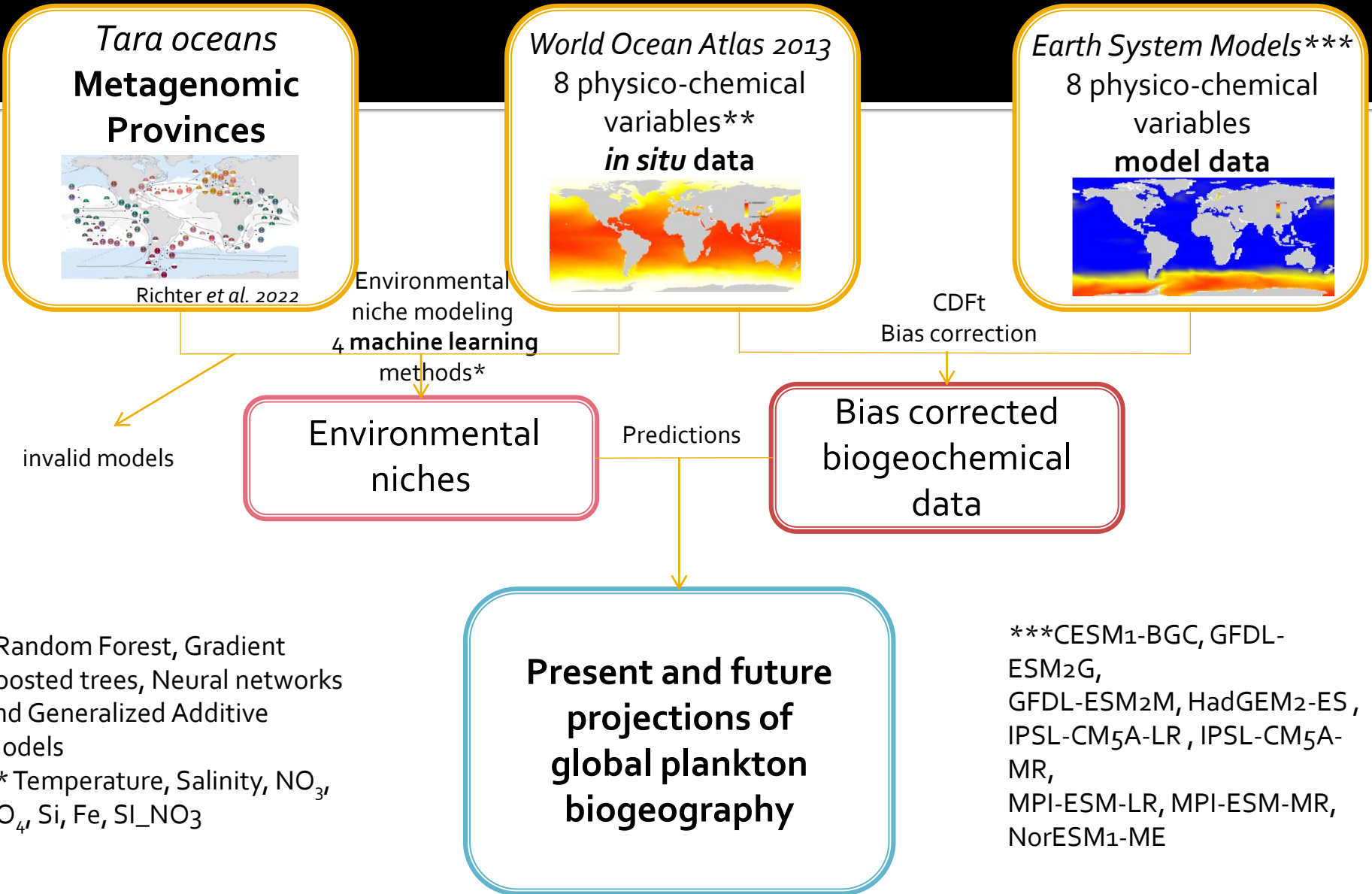
End of the century



# Do they perform well?

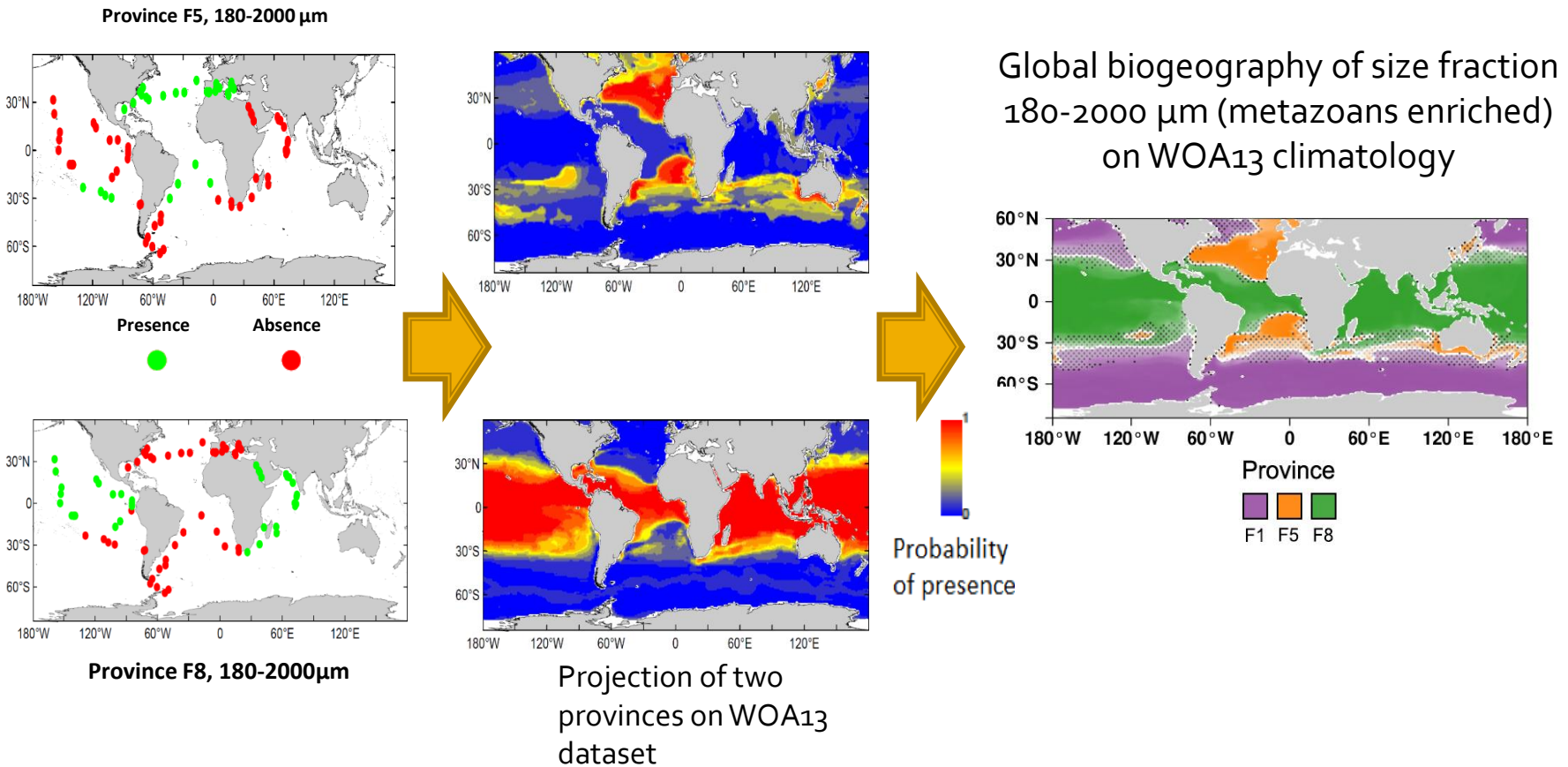


# A concrete example: Tara Ocean genomic biogeography



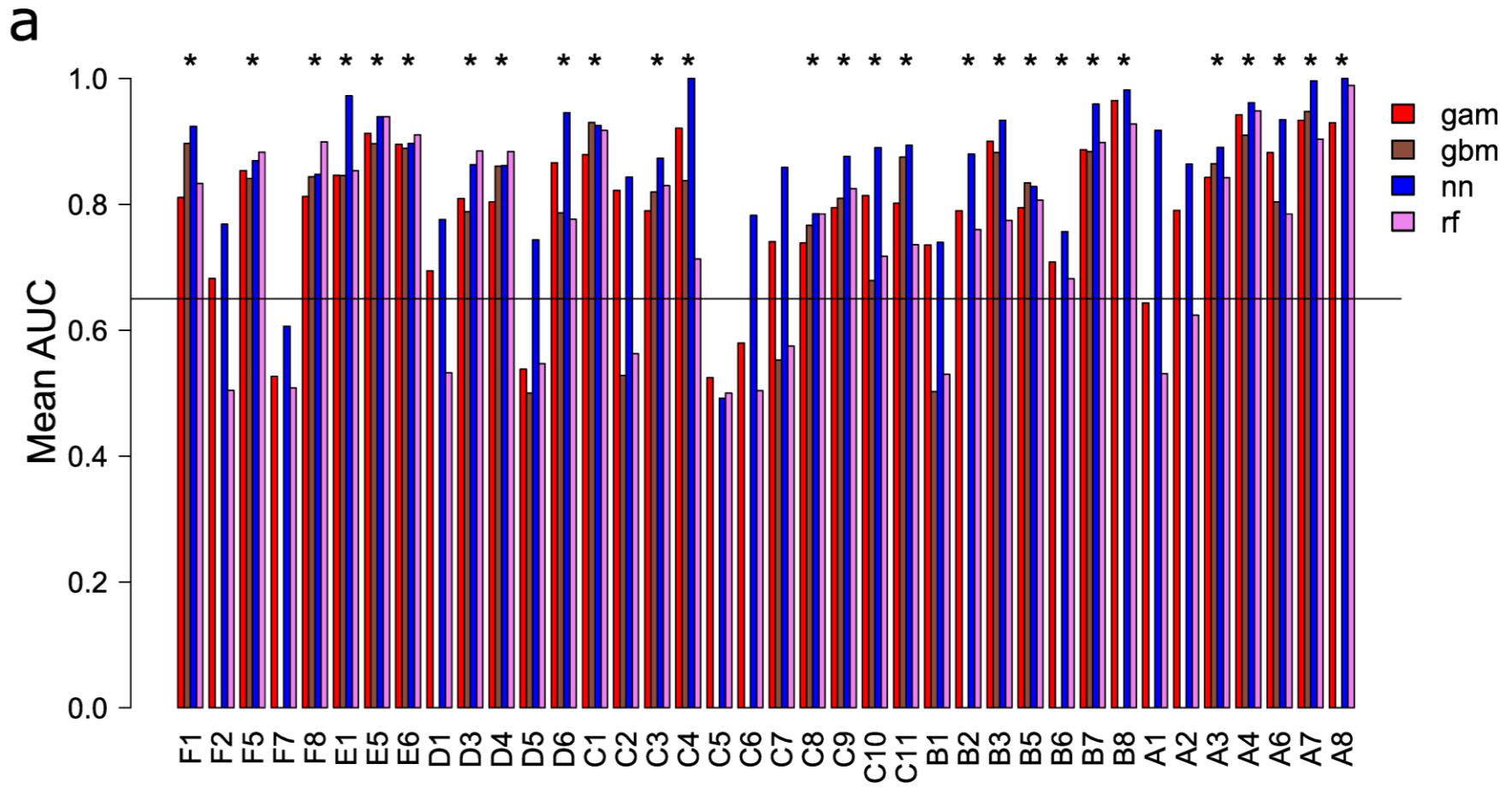


# A concrete example: Tara Ocean genomic biogeography Projections



- ⇒ Multiclass classifier divided in multiple single class classifiers
- ⇒ Simple biogeography for metazoans (separation of temperate, equatorial + tropical and polar)

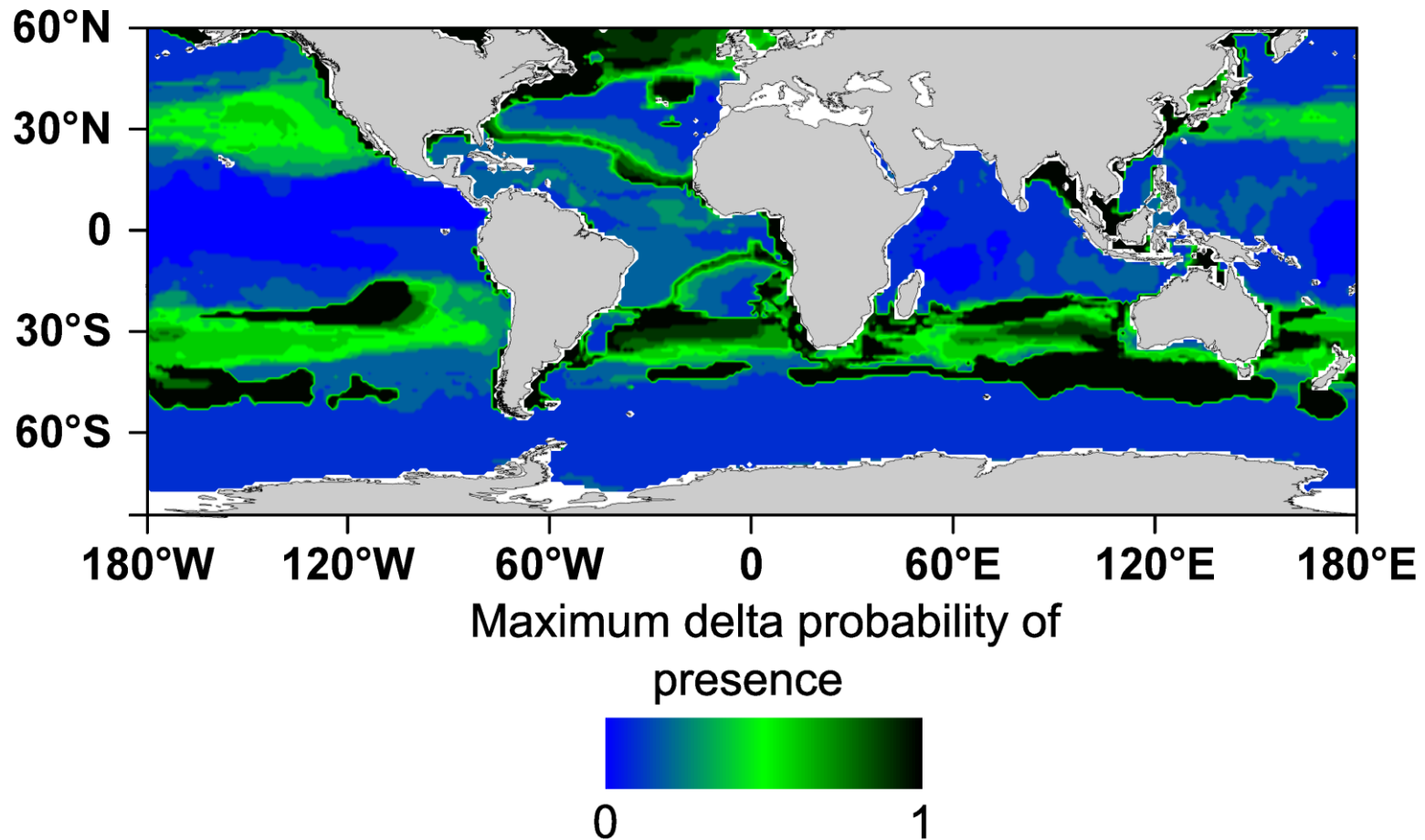
# Model performances



27 valid models

Neural network performs better

# Do they agree with each other?



# Genomic biogeographies of plankton



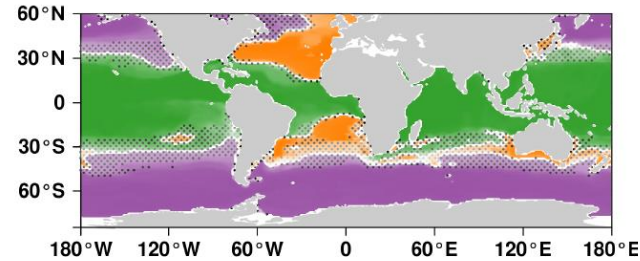
180-2000  $\mu\text{m}$

Small metazoans

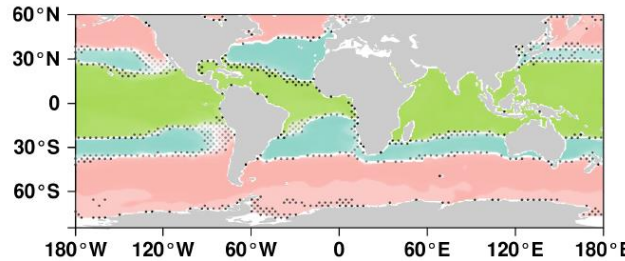
20-180  $\mu\text{m}$

Small metazoans

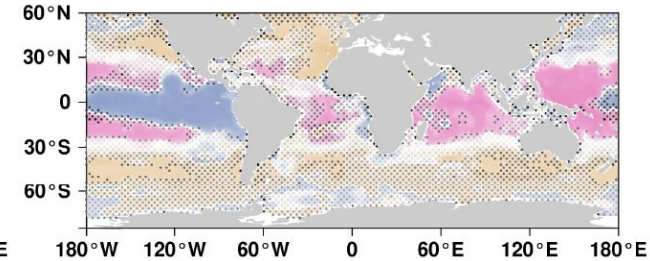
5-20  $\mu\text{m}$



Province



Province



Province

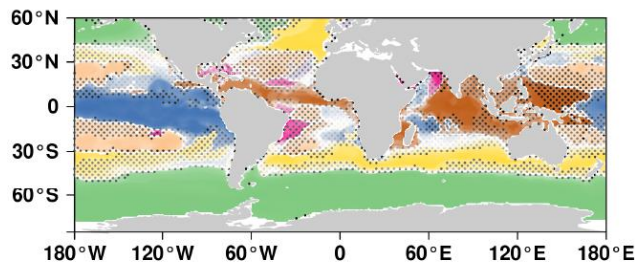


⇒ Simple latitudinal biogeography of large plankton organisms

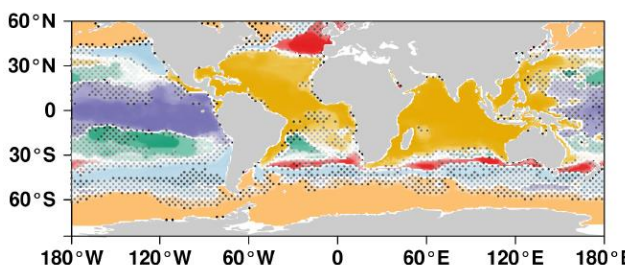
0.8-5  $\mu\text{m}$

0.22-3  $\mu\text{m}$

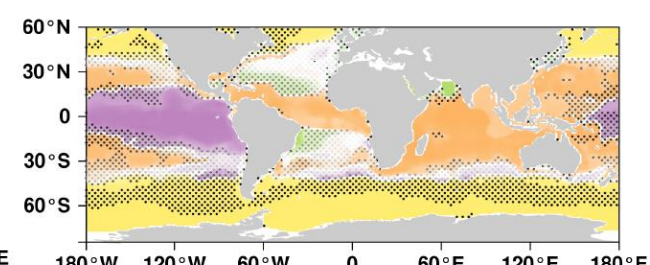
0-0.2  $\mu\text{m}$



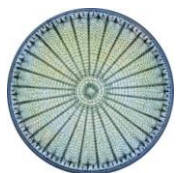
Province



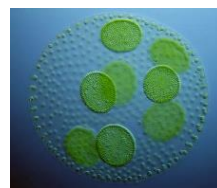
Province



Province



Unicellular algae



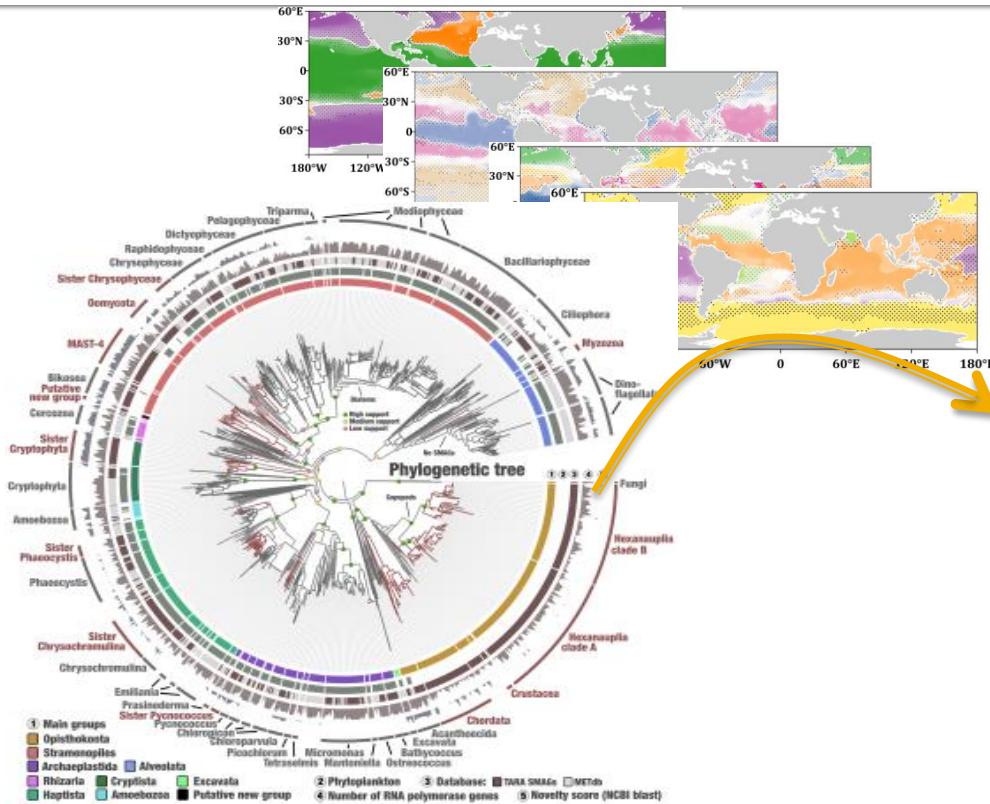
Bacteria



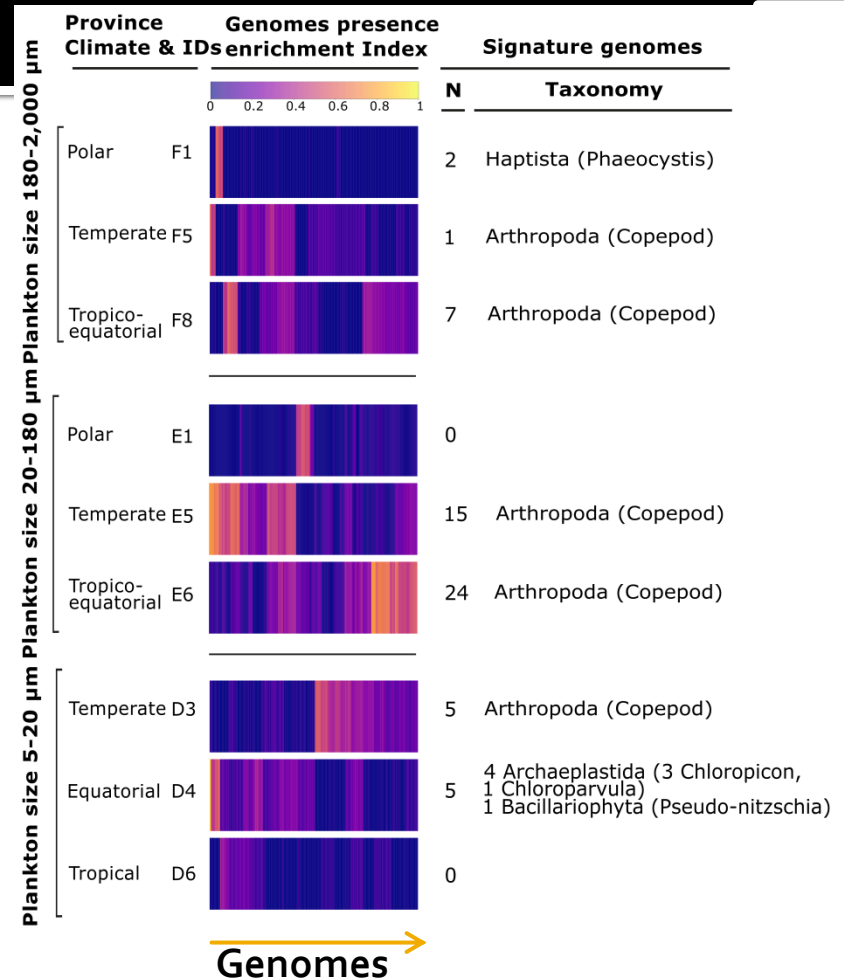
Viruses

⇒ More complex and patchy biogeography of small plankton organisms

# Signature genomes



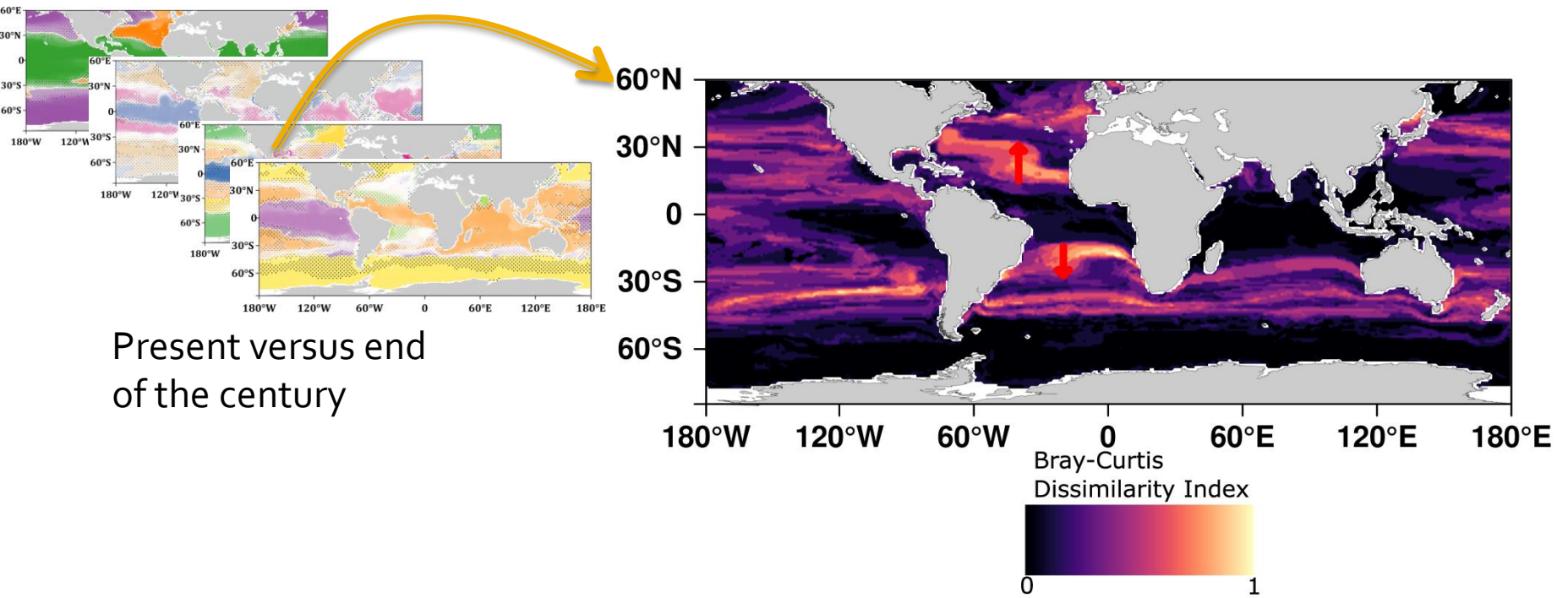
Delmont et al. 2021 *BioRxiv*



⇒ Signature genomes of climato-genomic provinces: highlights species and genomes that structure plankton biogeography at the global scale

⇒ *Climato-genomic* provinces structure plankton biogeography at a higher level than individual genomes

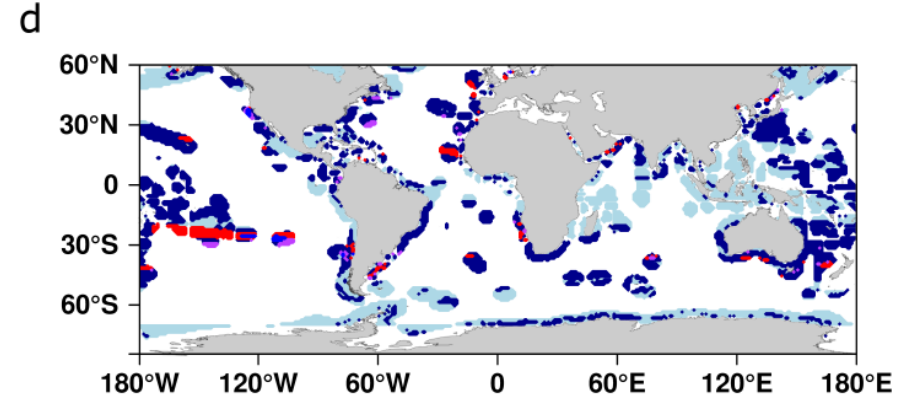
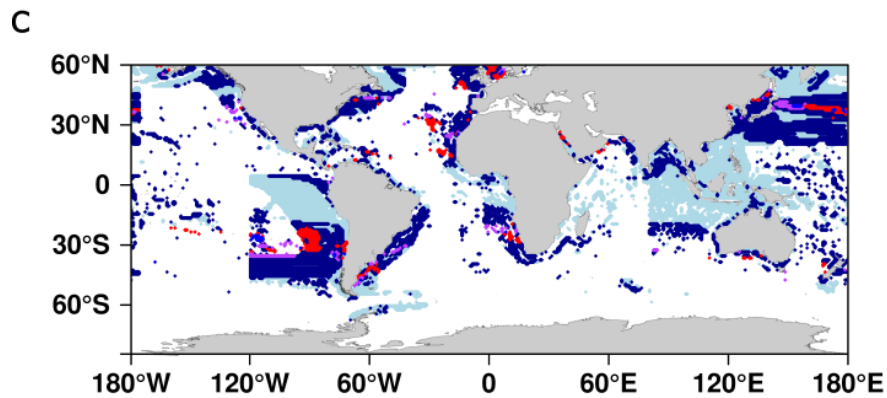
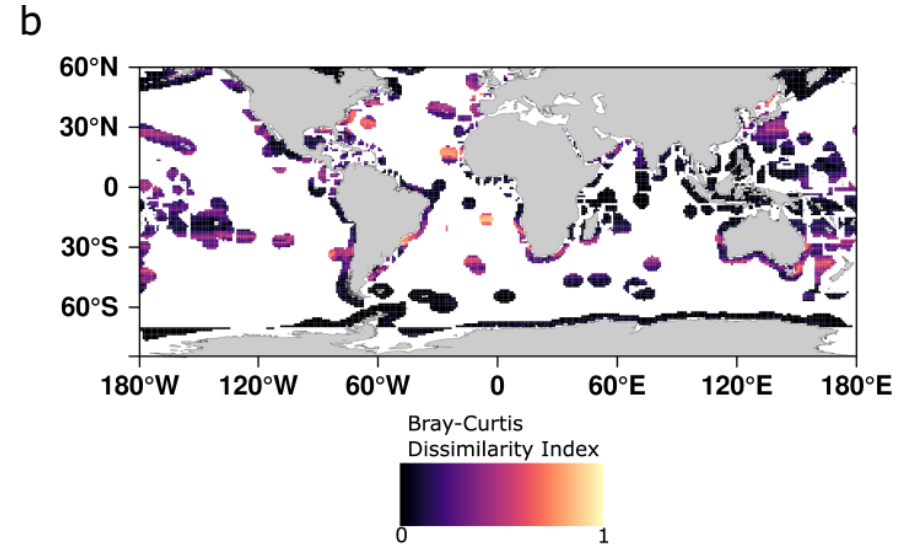
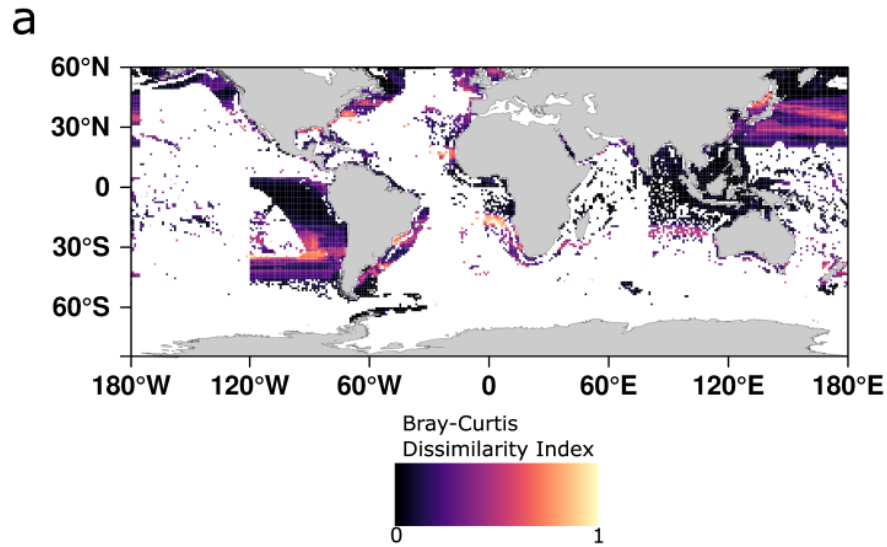
# Effects of climate change



⇒ Most important reorganization in temperate regions ( $25^{\circ}$  to  $60^{\circ}$ ): mean dissimilarity of 0.29 (north) and 0.24 (south)

⇒ **45 % to 57 %** of considered ocean area with an important change

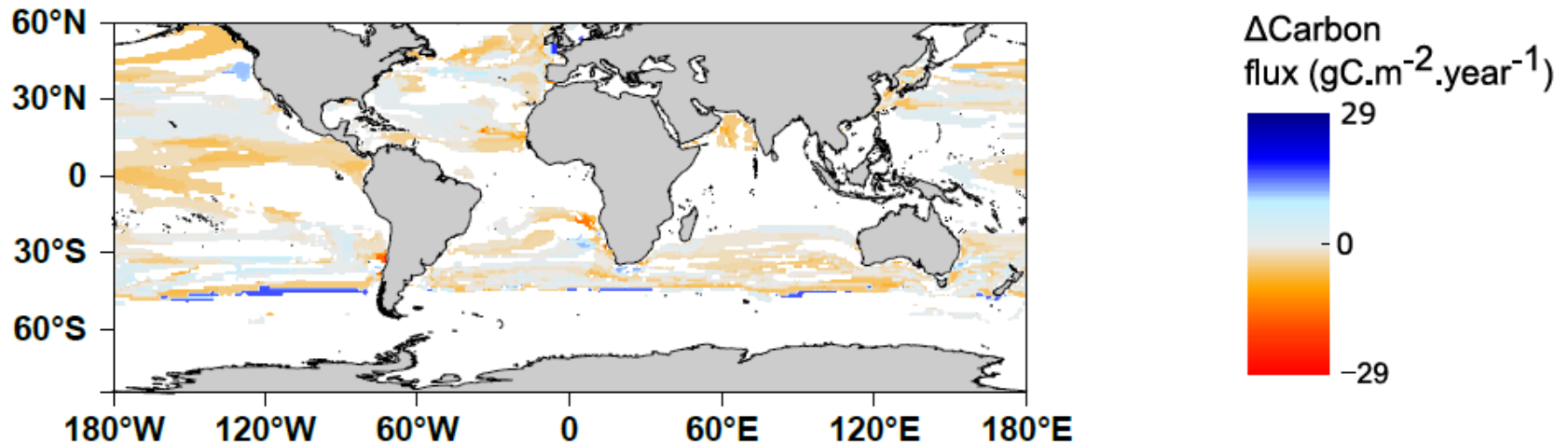
# Socio-economic impact



□ No significant change    □ Same assemblage    □ Change of assemblage    □ 2006 specific assemblage    □ 2009 specific assemblage    □ 2006 & 2009 specific

# Feedback on climate change

d



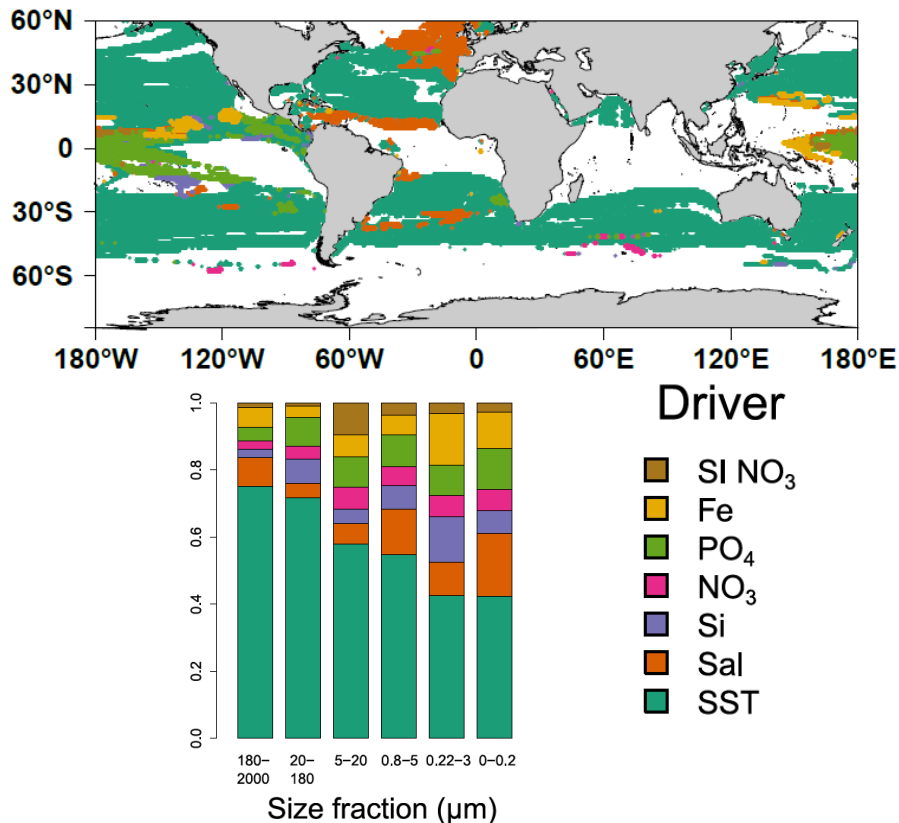
⇒ Based on the provinces' assemblages model: **decrease of 4% on average of POC export fluxes** (based on three extrapolated models of POC export: Eppley et al. 1979, Laws et al. 2000 and Henson et al. 2012)

⇒ Feedback on climate change: **reinforcement**



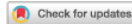
# Drivers

## Physico-Chemical Drivers of reorganizations around 2090



⇒ **Temperature changes** explains only **50% of changes** in dominant communities followed by phosphate (11.0%) and salinity (10.3%)

⇒ **Nutrients changes** are more important in driving small organisms' changes (consistent with their trophic modes)



## Restructuring of plankton genomic biogeography in the surface ocean under climate change

Paul Frémont<sup>1,2</sup>, Marion Gehlen<sup>3</sup>, Mathieu Vrac<sup>3</sup>, Jade Leconte<sup>1,2</sup>, Tom O. Delmont<sup>1,2</sup>, Patrick Wincker<sup>1,2</sup>, Daniele Iudicone<sup>4</sup> and Olivier Jaillon<sup>1,2</sup>

The impact of climate change on diversity, functioning and biogeography of marine plankton remains a major unresolved issue. Here environmental niches are evidenced for plankton communities at the genomic scale for six size fractions from viruses to meso-zooplankton. The spatial extrapolation of these niches portrays ocean partitionings south of 60° N into *climato-genomic* provinces characterized by signature genomes. By 2090, under the RCP8.5 future climate scenario, provinces are reorganized over half of the ocean area considered, and almost all provinces are displaced poleward. Particularly, tropical provinces expand at the expense of temperate ones. Sea surface temperature is identified as the main driver of changes (50%), followed by phosphate (11%) and salinity (10%). Compositional shifts among key planktonic groups suggest impacts on the nitrogen and carbon cycles. Provinces are linked to estimates of carbon export fluxes which are projected to decrease, on average, by 4% in response to biogeographical restructuring.

Planktonic communities are composed of complex and heterogeneous assemblages of small animals, single-celled eukaryotes (protists), bacteria, archaea and viruses. They contribute to the regulation of the Earth system through primary production via photosynthesis<sup>1</sup> and carbon export to the deep ocean<sup>2,3</sup> and are at the base of the food webs that sustain the whole trophic chain in the oceans<sup>4</sup>.

The composition of planktonic communities varies over time at a given site with daily<sup>5</sup> to seasonal fluctuations<sup>6</sup> following environmental variability<sup>7</sup> (for example, temperature, nutrients).

Time series observations have highlighted recent changes in the planktonic ecosystem such as changes in community composition<sup>24</sup> or poleward shifts of some species<sup>25</sup>, mainly attributed to temperature increase caused by anthropogenic greenhouse gas emissions. These changes are expected to intensify with ongoing global warming<sup>26</sup> and could lead to major reorganization of plankton community composition<sup>22</sup> with a potential decline in diversity<sup>27–29</sup>. Another major consequence of global reorganization of the seascape on biological systems would be a decrease of primary production at mid-latitudes and an increase at higher latitudes<sup>26</sup>.

Codes for SDMs: <https://github.com/institut-de-genomique/NCLIM-20102618B>

## Genomic evidence for global ocean plankton biogeography shaped by large-scale current systems

Daniel J Richter<sup>1,2†</sup>, Romain Watteaux<sup>3,4†</sup>, Thomas Vannier<sup>5,6,7†</sup>, Jade Leconte<sup>5,6</sup>, Paul Frémont<sup>5,6</sup>, Gabriel Reygondeau<sup>8,9</sup>, Nicolas Maillet<sup>10</sup>, Nicolas Henry<sup>1,6</sup>, Gaëtan Benoit<sup>11</sup>, Ophélie Da Silva<sup>6,12</sup>, Tom O Delmont<sup>5,6</sup>, Antonio Fernández-Guerra<sup>13,14,15</sup>, Samir Suweis<sup>16</sup>, Romain Narcis<sup>17</sup>, Cédric Berney<sup>1,6</sup>, Damien Eveillard<sup>6,18</sup>, Frederick Gavory<sup>5</sup>, Lionel Guidi<sup>6,12</sup>, Karine Labadie<sup>19</sup>, Eric Mahieu<sup>19</sup>, Julie Poulain<sup>5,6</sup>, Sarah Romac<sup>1,6</sup>, Simon Roux<sup>20</sup>, Céline Dimier<sup>1,21</sup>, Stefanie Kandels<sup>22,23</sup>, Marc Picheral<sup>6,12</sup>, Sarah Seaton<sup>6,12</sup>, Tara Oceans Coordinators, Stéphane Pesant<sup>24,25</sup>, Jean-Marc Aury<sup>5</sup>, Jennifer R Brum<sup>20,26</sup>, Claire Lemaitre<sup>11</sup>, Eric Pelletier<sup>5,6</sup>, Peer Bork<sup>22,27,28</sup>, Shinichi Sunagawa<sup>22,29</sup>, Fabien Lombard<sup>6,12,30</sup>, Lee Karp-Boss<sup>31</sup>, Chris Bowler<sup>6,21,21</sup>, Matthew B Sullivan<sup>20,32,33,34</sup>, Eric Karsenti<sup>6,21,23</sup>, Mahendra Mariadassou<sup>17</sup>, Ian Probert<sup>1,6</sup>, Pierre Peterlongo<sup>11</sup>, Patrick Wincker<sup>5,6</sup>, Colomán de Vargas<sup>1,6\*</sup>, Maurizio Ribera d'Alcalá<sup>3\*</sup>, Daniele Iudicone<sup>3\*</sup>, Olivier Jaillon<sup>5,6\*</sup>

## Cell Genomics

CellPress  
OPEN ACCESS

### Article

## Functional repertoire convergence of distantly related eukaryotic plankton lineages abundant in the sunlit ocean

Tom O. Delmont<sup>1,2,9,\*</sup>, Morgan Gaia<sup>1,2</sup>, Damien D. Hingsinger<sup>1,2</sup>, Paul Frémont<sup>1,2</sup>, Chiara Vanni<sup>3</sup>, Antonio Fernandez-Guerra<sup>4</sup>, A. Murat Eren<sup>5</sup>, Artem Kourlaiev<sup>1,2</sup>, Leo d'Agata<sup>1,2</sup>, Quentin Clayssen<sup>1,2</sup>, Emilie Villar<sup>1</sup>, Karine Labadie<sup>1,2</sup>, Corinne Cruaud<sup>1,2</sup>, Julie Poulain<sup>1,2</sup>, Corinne Da Silva<sup>1,2</sup>, Marc Wessner<sup>1,2</sup>, Benjamin Noel<sup>1,2</sup>, Jean-Marc Aury<sup>1,2</sup>, Tara Oceans Coordinators, Colomán de Vargas<sup>2,6</sup>, Chris Bowler<sup>2,7</sup>, Eric Karsenti<sup>2,6,8</sup>, Eric Pelletier<sup>1,2</sup>, Patrick Wincker<sup>1,2</sup> and Olivier Jaillon<sup>1,2</sup>

# Plankton biogeography website

TO  
EN  
DB  
Tara  
Oceans  
Ecological  
Niches  
Database



Tara Oceans niches | [IGOB: Interactive Generator Of Biogeography](#)

Visualize theoretical niches (convex hulls)

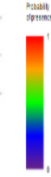
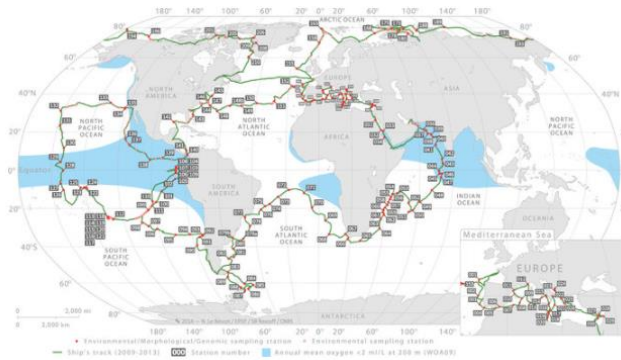
An ecological niche is defined as the envelopp of physico-chemical parameters in which a species lives. The niche concept was popularized by the zoologist G. Evelyn Hutchinson in 1957. Here you can visualize the ecological niches of 374 Metagenome Assembled Genomes (MAGs) built from the Tara Oceans dataset. Select below the SMAG of interest and visualize the effect of climate change (high emission scenario RCP8.5) on its niche by clicking on 'Severe climate change'. Colors on the maps allows you to visualize the probability of presence at each location of all oceans of the selected MAG

Select kingdom: All  
Select phylum: All  
Select class: All  
Select order: All  
Select family: All  
Select genre: All

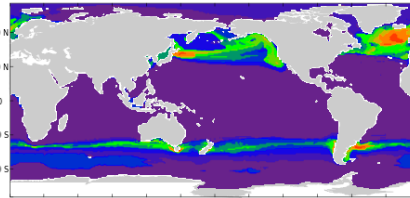
Select SMAG of interest

No SMAG Selected

Select a genome and visualize its biogeography and projected impacts of climate change

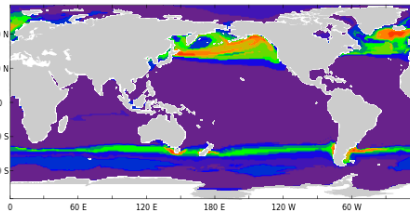


TARA AON\_82\_MAG\_00012  
Present biogeography (2006-15)

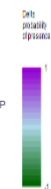
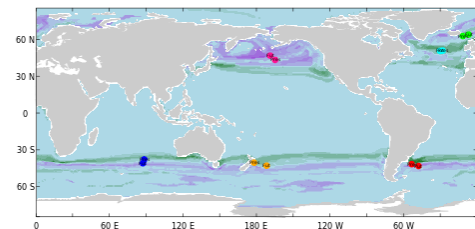


this map displays the niche of AON\_82\_MAG\_00012\_Ostreococcus in projected end of the century conditions under high emission scenario RCP

TARA AON\_82\_MAG\_00012  
Future biogeography (2090-99)



TARA AON\_82\_MAG\_00012  
Delta (Future - Present)



[http://end.mio.osupytheas.fr/Ecological\\_Niche\\_database/](http://end.mio.osupytheas.fr/Ecological_Niche_database/)

# Good practices in Species Distribution Modeling

- Spatial species distribution models are **powerful tools** to project species distributions or organism or whole community distribution
- **Evaluate model performances** and **optimize hyperparameters** using **cross validation**
- **Don't model species for which there is not enough data** (presence points): equilibrium in presence/absence data is preferable though rarely achievable
- Be careful in the **choice of predictors** and **extrapolate in predictability zones** **OR/AND acknowledge uncertainties/biases:**
  - choose predictors based on knowledge
  - use z-scores (for better performance of neural networks)
  - spatio-temporal sampling and extrapolation
- Other types of SDM exists pseudo-absence models, presence only models, models informed by biotic interactions, **quantitative models (very uncertain, strong errors)**

# SDM using machine learning: drawbacks and limits

- 0/1 SDMs are **Statistical** species distribution models => they remain correlative models, they are not causal (and don't account for currents, adaptation, acclimation...)
- Mechanistic species distribution models exists => predict the fundamental niche of a species based on its physiology and project it on the seascape (competition can be added too)
- Purely statistical models can't predict beyond extremes and training datasets values (can physics do it?)

RESEARCH ARTICLE | BIOLOGICAL SCIENCES | 8

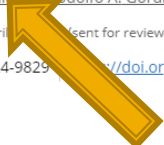


## Present and future global distributions of the marine Cyanobacteria *Prochlorococcus* and *Synechococcus*

Pedro Flombaum, José L. Gallegos, Rodolfo A. Gordillo, , and Adam C. Martiny  [Authors Info & Affiliations](#)

Contributed by David M. Karl, April 2013; accepted for review January 22, 2013)

May 23, 2013 | 110 (24) 9824-9829 | <https://doi.org/10.1073/pnas.1307701110>



**Bad practice!** extrapolation beyond training, no acknowledgement

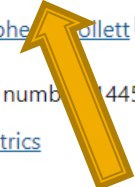
Article | [Open Access](#) | [Published: 15 March 2023](#)

## Understanding opposing predictions of *Prochlorococcus* in a changing climate

Vincent Bian, Merrick Cai & Christopher J. Collett 

*Nature Communications* **14**, Article number 4445 (2023) | [Cite this article](#)

996 Accesses | 10 Altmetric | [Metrics](#)



**Warns on the dangers/limits of (quantitative) statistical models**

# Acknowledgement

## My PhD supervisors

Olivier Jaillon (Genoscope), Marion Gehlen



## Stazione Zoologica

Daniele Iudicone

Lucia Campese

Bruno Hay Mele



## LSCE

Marion Gehlen, Mathieu Vrac



## Genoscope (LAGE)

Patrick Wincker, Eric Pelletier, Tom Delmont, Amin Madoui, Adrien Thurotte, Quentin Carradec, Betina Porcel, Julie LeHoang, Romual Laso-Jadart, Jade Leconte, Marie Burel, Janaina Rigonato, Nina Guérin, Margaux Crédeville, Lucas Pavlovic

