# Do deep neural networks improve functional annotation of alpha-cyanobacteria?

## Juliana SILVA BERNARDES
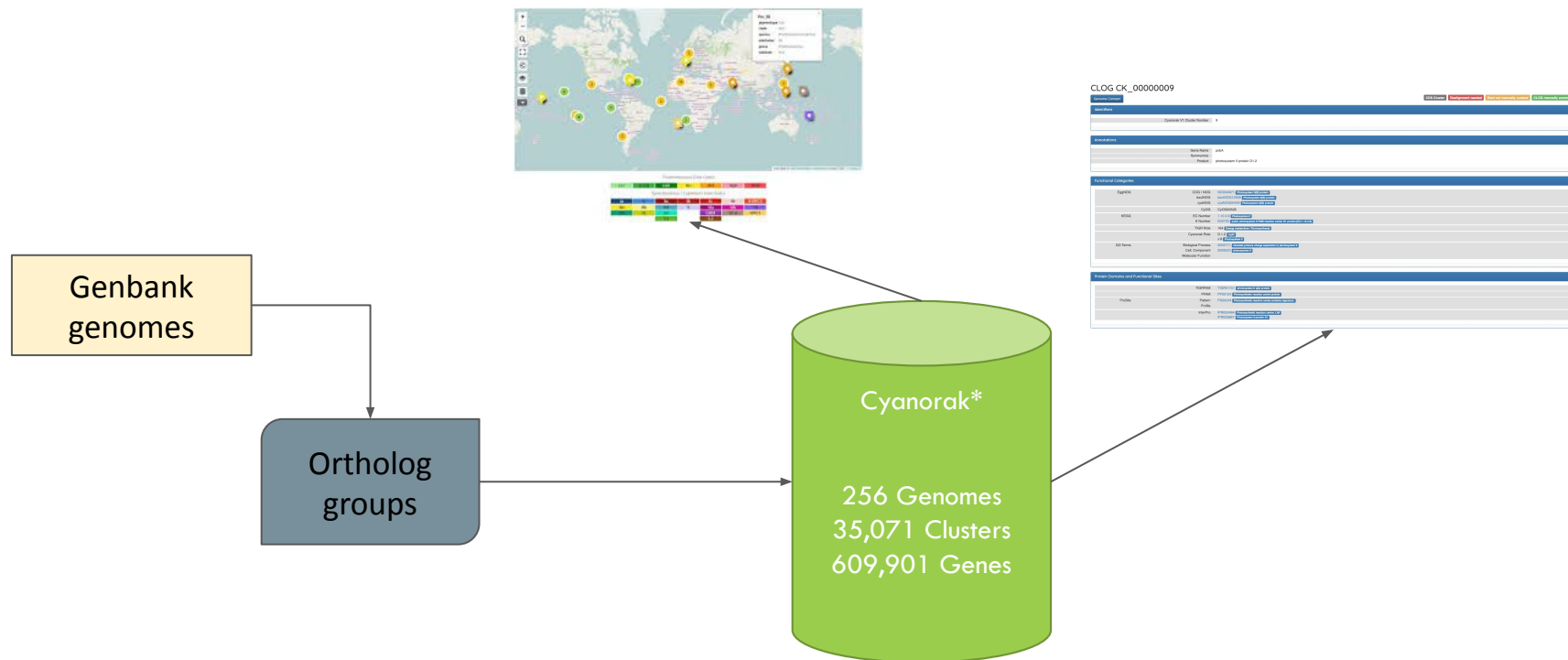
https://scholar.google.com/citations?user=a-ZYwhEAAAAJ&hl=en

https://www.lcqb.upmc.fr/julianab/

jusilvabernardes@sb-roscoff.fr;
juliana.silva_bernardes@sorbonne-universite.fr

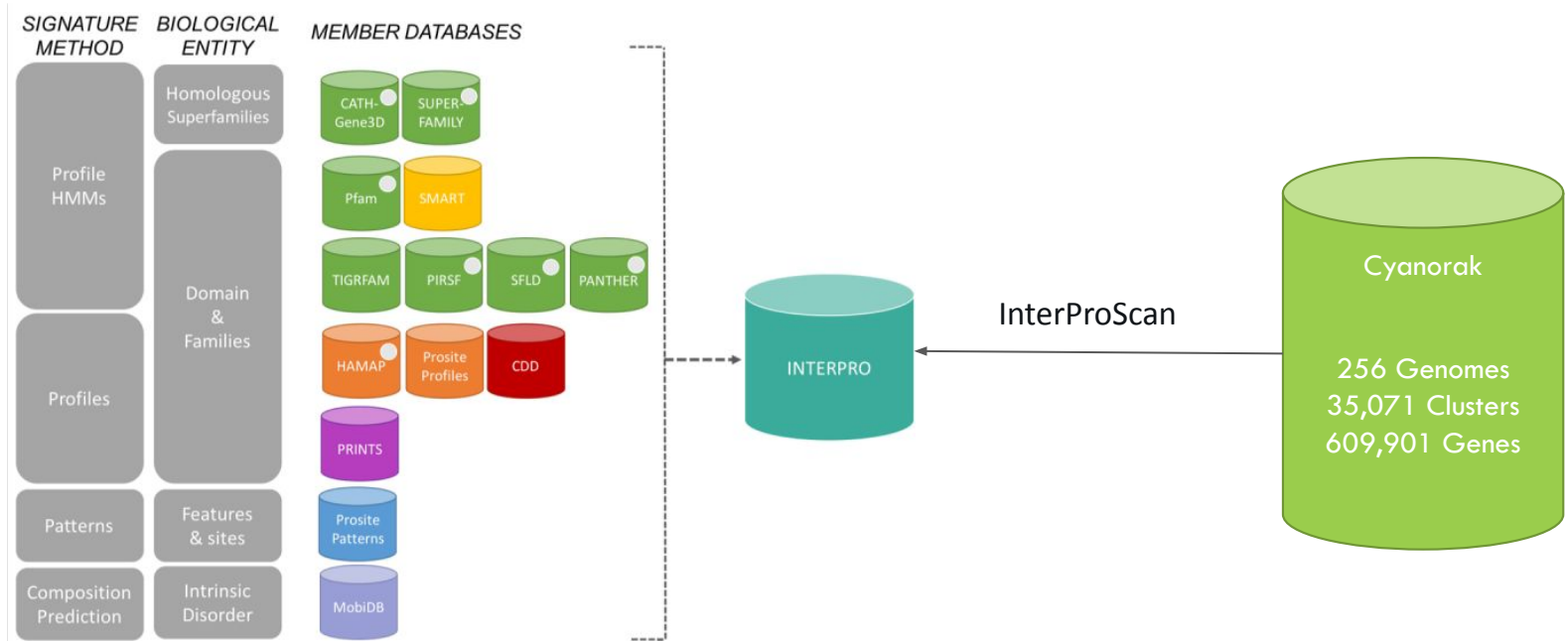CNRS · SORBONNE UNIVERSITÉ
Station Biologique
de Roscoff
1872

# Why α-Cyanobacteria ?

- Phototrophic **microorganisms**, appeared **3.5 billion** years ago

- **Responsible** for the apparition of **Oxygen** on Earth

- One of the **most diverse/widely distributed** prokaryotic phyla

- Colonize both **terrestrial** and **aquatic environments**

- Origin of land **plants /marine algae** over the next millions of years

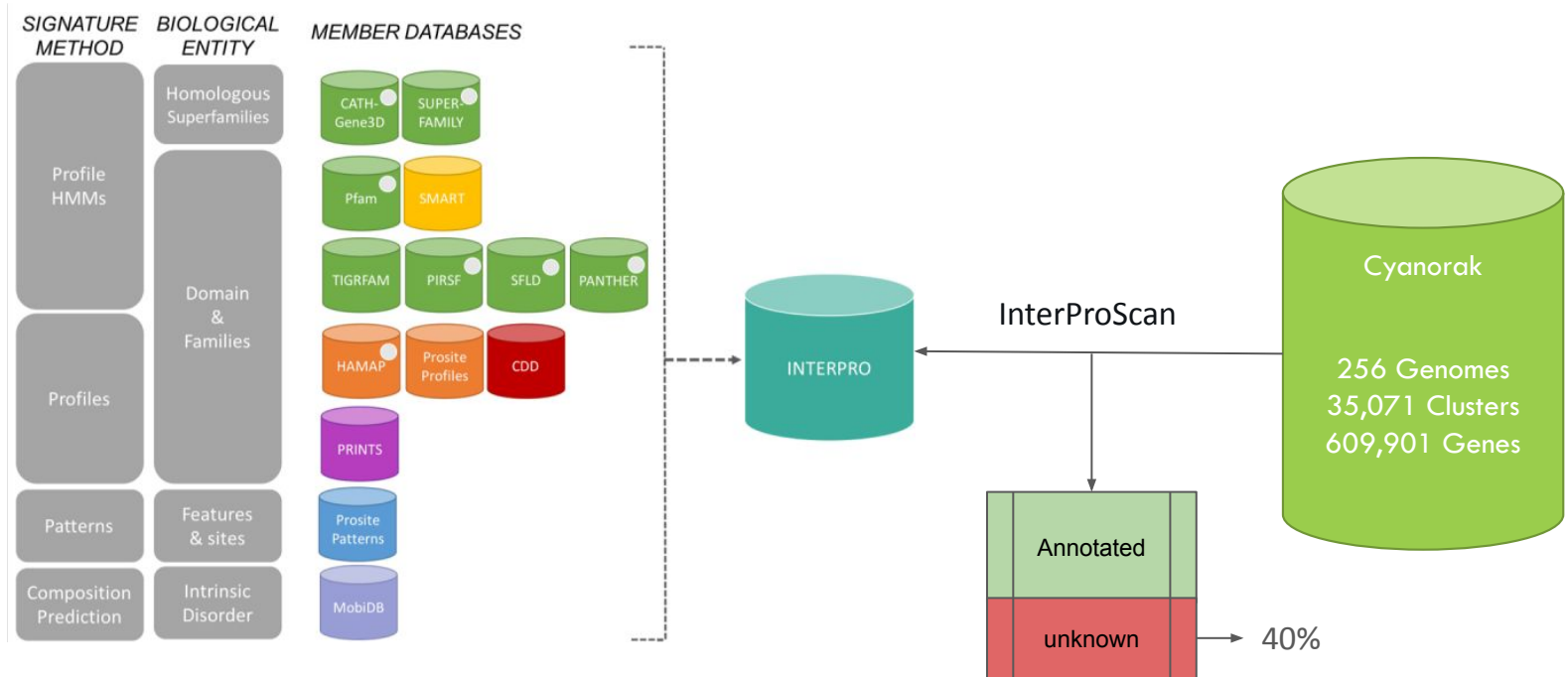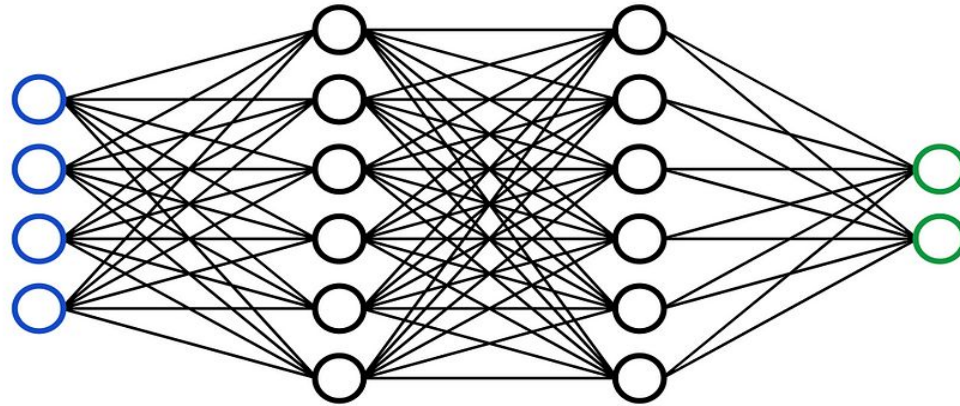- They are **most abundant photosynthetic organisms** in the ocean and large lakes.

Genbank genomes

Ortholog groups

Cyanorak*

256 Genomes
35,071 Clusters
609,901 Genes

CLOG CK_00000009

*https://academic.oup.com/nar/article/49/D1/D667/5943826

- Functional annotation of α-Cyanobacteria is based on Interpro database

- Functional annotation of α-Cyanobacteria is based on Interpro database

- Deep neural networks have already reached impressive results in protein three-dimensional structure prediction.

- Here, we investigated if a such methodology could improve function annotation in α-Cyanobacteria

Masking

Language model

"Would you tell me, please, which way I ought to go from here?"

"That depends a good deal on where you want to get to," said the Cat.

"I don't much care where—" said Alice.

"Then it doesn't matter which way you go," said the Cat.

"—so long as I get *somewhere*," Alice added as an explanation.

"Oh, you're sure to do that," said the Cat, "if you only walk long enough."
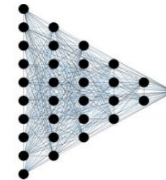
Original text

"Would you tell me, ███, which way I ███ to go from here?"

"That ███ a ███ deal on where you want to get to," said the Cat.

"I ███ much care where—" ███ Alice.

"Then it doesn't matter ███ you go," said the Cat.

"—so long as I get *somewhere*," Alice ███ as an explanation.

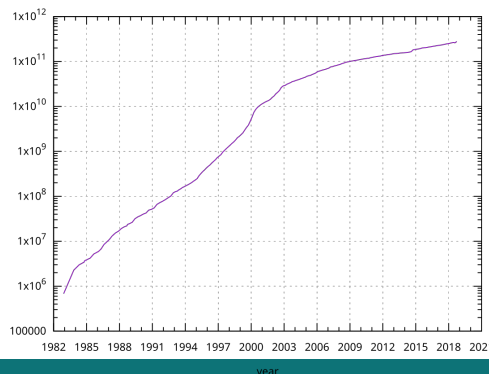"Oh, ███ to do that," said the Cat, "if ███ only ███ long enough."

Masked text

"Would you tell me, sir, which way I need to go from here?"

"That depends a good deal on where you want to get to," said the Cat.

"I don't much care where—" said Alice.

"Then it doesn't matter which way you go," said the Cat.

"—so long as I get *somewhere*," Alice added as an explanation.

"Oh, no need to do that," said the Cat, "if one only waits long enough."

Predicted text

Large corpus
(unlabeled text)
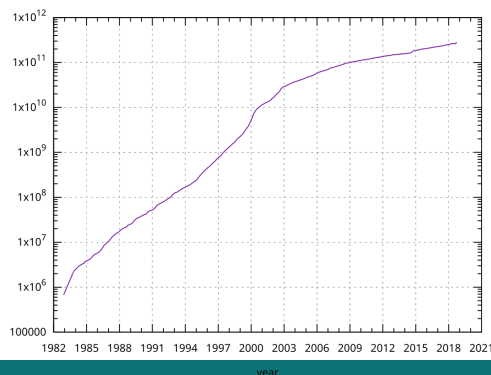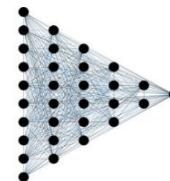
*Loss*

Masking

Language model

Original text

Masked text

Predicted text

*Loss*

Large corpus
(unlabeled text)

Protein sequences

Growth of GenBank

# Protein language models



Masking

Language model

Original text

Masked text

Predicted text

*Loss*

Large corpus
(unlabeled text)

Protein sequences

Growth of GenBank

Embedding

a
c
d
f
…
a
r

Encoder

Decoder

a
**e**
d
f
…
a
**w**

*https://www.pnas.org/doi/abs/10.1073/pnas.2016239118

PlmSearch*
Fantasia+

*https://www.nature.com/articles/s41467-024-46808-5

+https://www.nature.com/articles/s42003-025-08651-2

# Results : comparing predictions



PlmSearch | 85% | (annotated) | (unknown)
InterPro | 74%

annotated*

unknown



identique | sous-ensemble | different

PlmSearch = InterPro + PlmSearch ⊆ InterPro + InterPro ⊆ PlmSearch + PlmSearch ≠ InterPro + PlmSearch ∩ InterPro

**111** (4,88%) (A) + **4** (0,18%) (B) **1310** (57,61%) + **639** (28,10%) (C) + **210** (9,23%) (D)

PlmSearch 1458 | 2274 | 973 InterPro

identique | sous-ensemble | different

| PlmSearch = InterPro | + | PlmSearch ⊆ InterPro | + | InterPro ⊆ PlmSearch | + | PlmSearch ≠ InterPro | + | PlmSearch ∩ InterPro |

**111** (4,88%) (A) | **4** (0,18%) **1310** (57,61%) (B) | **639** (28,10%) (C) | **210** (9,23%) (D)



Tableau récapitulatif des tests statistiques de Mann-Whitney
(ns : non significatif ; *** : < au risque α = 1‰)

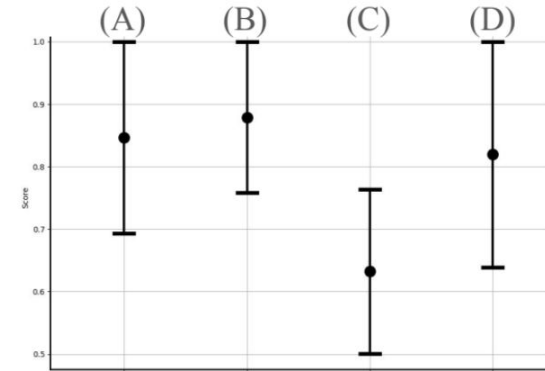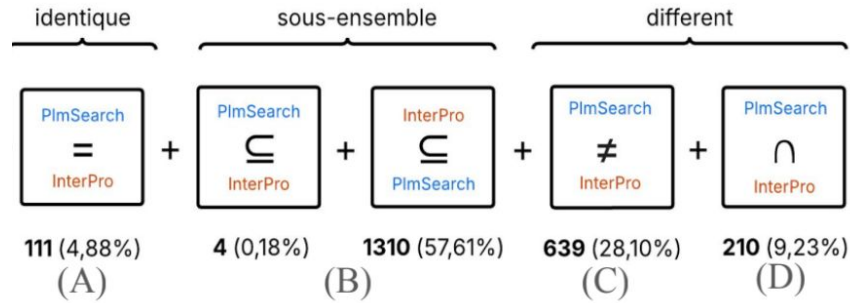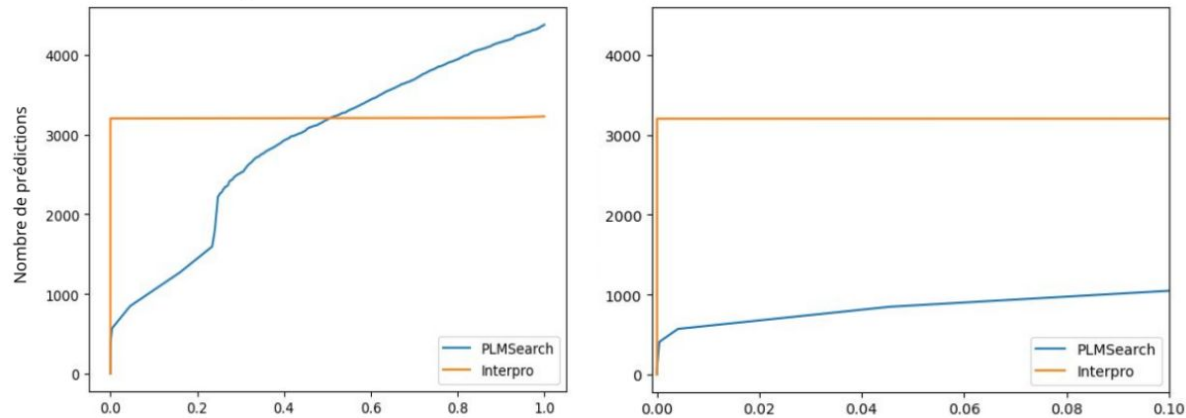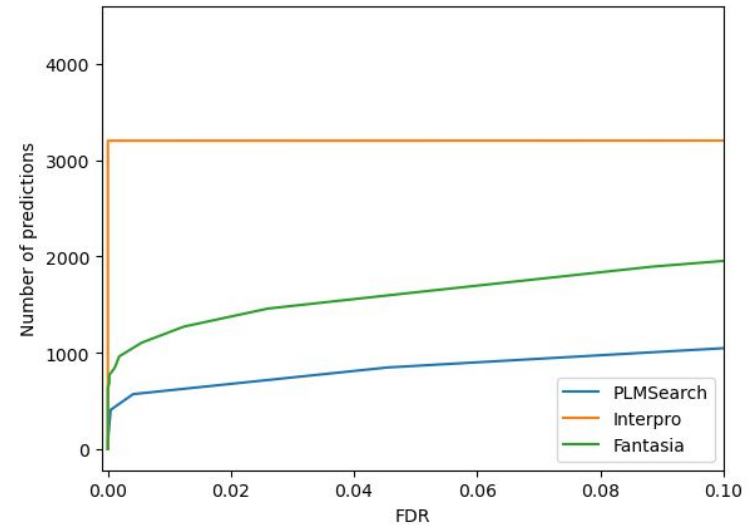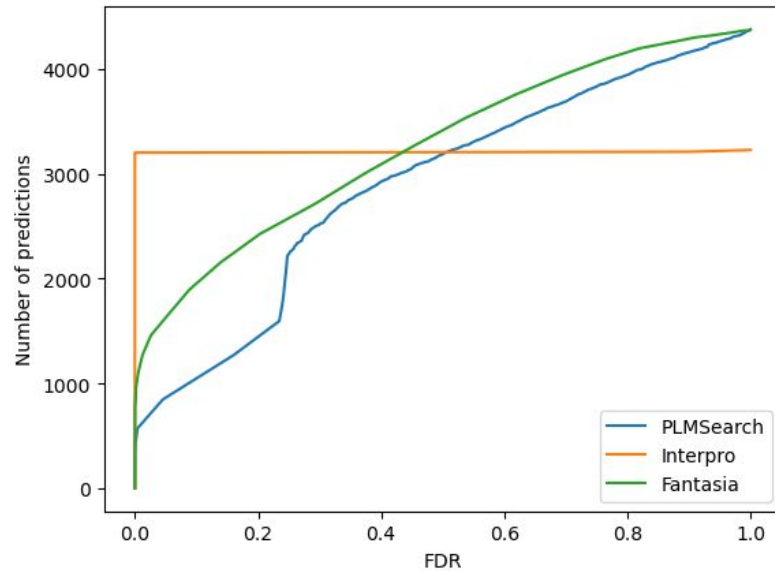| TEST | P-VALUE | SIGNIFICATIVITÉ |
|------|---------|-----------------|
| A & B | 0.556 | ns |
| A & C | $2.16 \times 10^{-27}$ | *** |
| A & D | $4.28 \times 10^{-4}$ | *** |
| B & C | $4.95 \times 10^{-140}$ | *** |
| B & D | $1.78 \times 10^{-11}$ | *** |
| C & D | $3.69 \times 10^{-27}$ | *** |

- We shuffled the amino acid order of each sequence with original dataset (D) to create an artificial database (R)

- FDR for a given score S is :

$$\frac{\text{number of predictions in R}}{\text{number of predictions in R + number of predictions in D}}$$

- We also tested Fantasia and measured FDR

- PLMSearch seems to complete/enrich the InterPro annotations.

- However, when it finds completely different annotations, the scores are significantly lower, indicating low accuracy.

- For a given FDR PLMSearch annotated less proteins than Interpro, the difference is even more evident for an FDR bounded to 10%.

- Higher FDR were also observed in Fantasia results.

# Acknowledges

- Emile Hembert - L2 student - Sorbonne Université

- Dorian Le Roux - L2 student - Sorbonne Université

- Fabio RJ Vieira - IR - Sorbonne Université

- Laurence Garczarek - CNRS-DR

- Frédéric Partensky, CNRS-DR

Thanks for your attention !