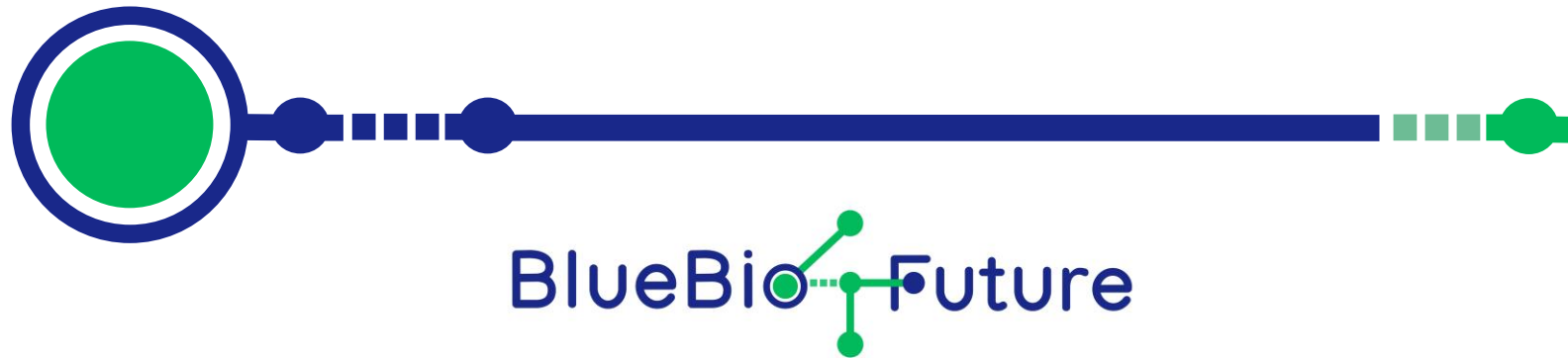# Introduction to Bioinformatics in Natural Products Discovery: an historical overview
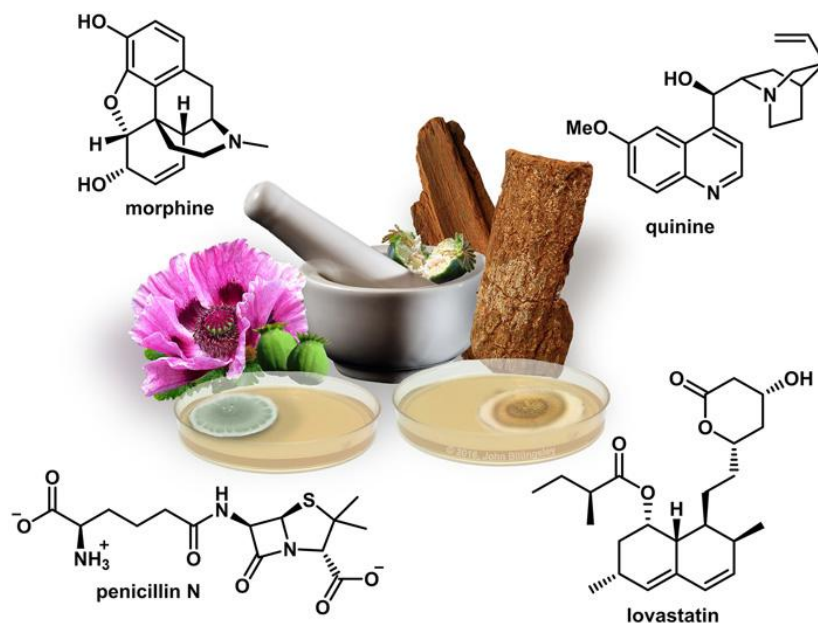
**Adriana Rego**

**Short Course in Bioinformatics – Decoding (Meta)Genomes for Natural Products Discovery**
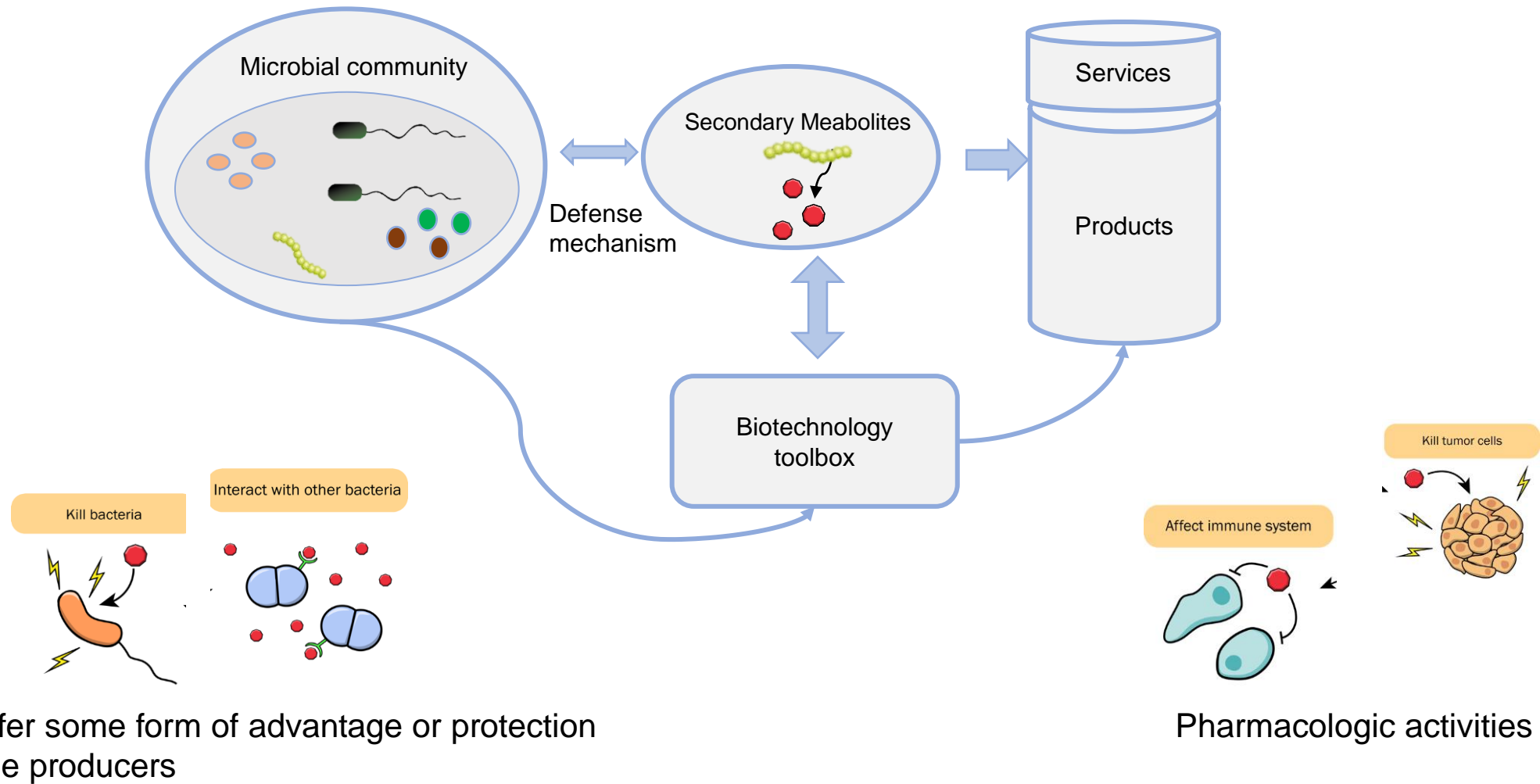
**17th – 19th March 2025**

# Microorganims as a prolific source of natural products: an historical overview

Natural products, also called secondary or specialized metabolites could be defined broadly as any **molecules found in nature**.

More traditionally in organic and medicinal chemistry communities, natural products **are defined as small organic molecules** (MW < 1500 daltons) generated **from secondary metabolic pathways.**
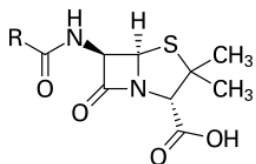


morphine

quinine

penicillin N

lovastatin

*Nat. Prod. Biosynth.* **2022**, RSC Publishing

# Microorganims as a prolific source of natural products: an historical overview



Microbial community

Defense mechanism

Secondary Meabolites

Services

Products

Biotechnology toolbox

Kill bacteria

Interact with other bacteria

Confer some form of advantage or protection to the producers

Kill tumor cells

Affect immune system
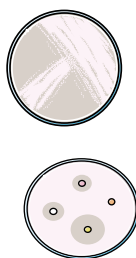
Pharmacologic activities

(Adapted from  Senft 2020)

# Microorganims as a prolific source of natural products: an historical overview

**1929**

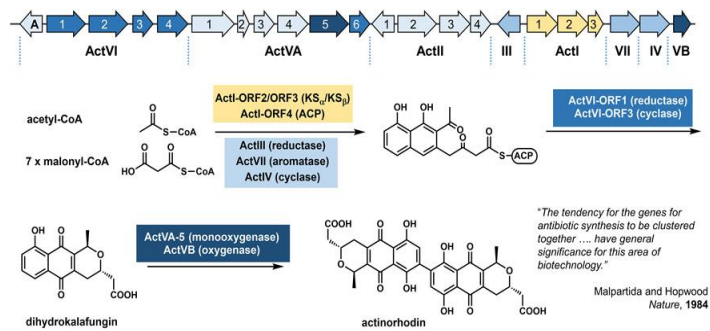| Penicillin | Easily cultivable bacteria | Re-discovery |
|---|---|---|



23 000 new NPs

Beginning of an age of natural product-driven medicine

**bioactivity-guided fractionation**

Potential of cultured microbes had been exhausted

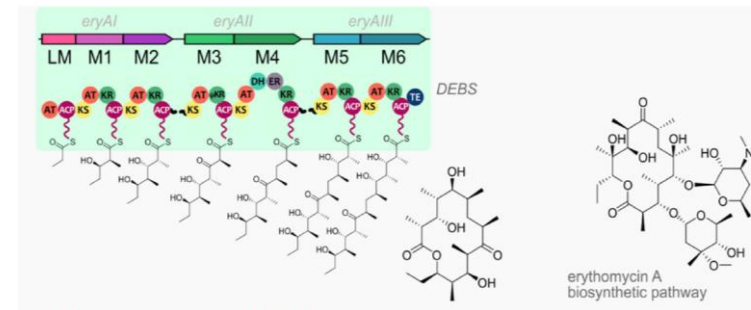# Microorganims as a prolific source of natural products: an historical overview

## 1984



*Nature* **1984,** 309, 462–464

This study established that the necessary **biosynthetic genes are clustered on a contiguous stretch of DNA.**
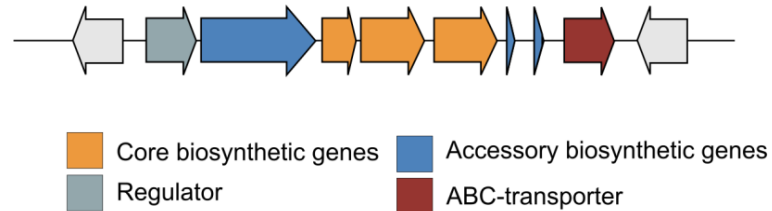
## 1990, 1991



*Nature* **1990**, *348*, 176; *Science***1991**, *252*, 671

The genes required to assemble the 14-member macrolactone core of erythromycin are **contiguous** in *Saccharopolyspora Erythraea.*

# Microorganims as a prolific source of natural products: an historical overview

**A Biosynthetic Gene Cluster (BGC)** can be defined as a physically clustered group of two or more genes in a particular genome that together encode a biosynthetic pathway for the production of a specialized metabolite (including its chemical variants).

Core biosynthetic genes    Accessory biosynthetic genes
Regulator    ABC-transporter

**1998**

'A new frontier of science is emerging that unites biology and chemistry - the **exploration of natural products from previously uncultured soil microorganisms.**
The methodology has been made possible by advances in molecular biology and eukaryotic genomics, which have laid the groundwork for cloning and functional analysis of the collective genomes of soil microflora, which we **term the metagenome of the soil**.'
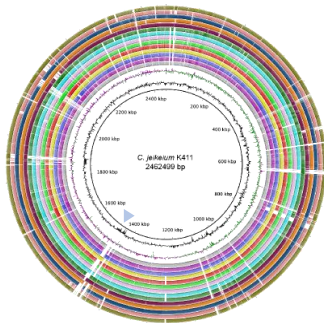
*Chemistry & Biology* **1998,** 5, 245-249



Soil sample

Separate bacteria

Isolate large DNA fragments

Cut DNA with restriction enzyme

(+ BAC)

Clone DNA in to BAC vector

Transform *E. coli*

Screen transformants

Chemistry & Biology

# Bioinformatic analysis of natural product biosynthetic capacity

## The emergence of a new core discipline - the genome mining

**2000**

Genome sequencing and mining



Recovery of complete BGCs

In the early 2000s, the sequencing of the first *Streptomyces* bacterial genomes revealed that the vast majority of small molecules produced by microbes had yet to be discovered, thus opening the door for future discovery efforts and for the emergence of a new core discipline – the **genome mining.**

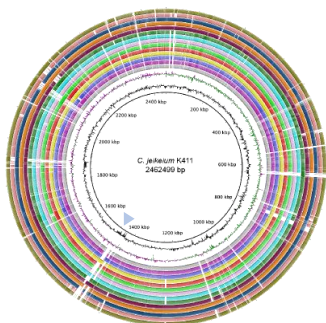Renaissance and Genome-guided discovery of new molecules

# Bioinformatic analysis of natural product biosynthetic capacity

## The emergence of a new core discipline - the genome mining

**2000**

Genome sequencing and mining



Recovery of complete BGCs

**Genome mining** describes the targeted bioinformatic analysis of (meta-)genomes to identify gene clusters involved in the biosynthesis of NPs.
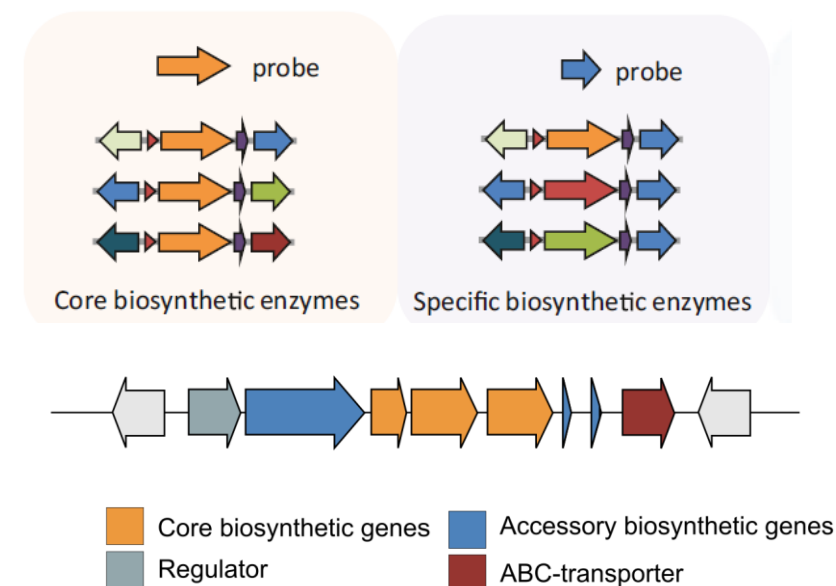
Renaissance and Genome-guided discovery of new molecules

# Bioinformatic analysis of natural product biosynthetic capacity

**What makes BGCs being easily detected by Genome Mining?**

While NPs are chemically diverse, their biosynthetic machineries **are often highly conserved**. Core biosynthetic enzymes are characterized by **high amino-acid sequence similarity**, which allows screening of genomic data for the presence of **specific biosynthetic genes that encode the required enzymatic activity**.
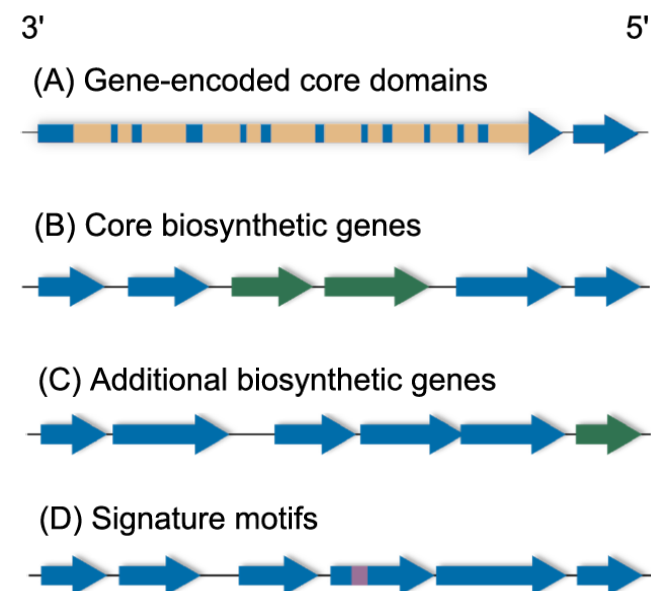
# Bioinformatic analysis of natural product biosynthetic capacity

**What makes BGCs being easily detected by Genome Mining?**

Bacterial biosynthetic genes are **clustered on a contiguous stretch of DNA (co-localized).**

**Existence of core biosynthetic genes –** highly conserved across different organisms (e.g. KS and A domains)

Similar **genetic pattern organization** – core genes, regulatory genes and transporters.

3'                                                                    5'

(A) Gene-encoded core domains

(B) Core biosynthetic genes

(C) Additional biosynthetic genes

(D) Signature motifs

# Bioinformatic analysis of natural product biosynthetic capacity

## Homology search tools/Algorithms for BGCs prediction

**Bioinformatic analysis of natural product biosynthetic capacity**
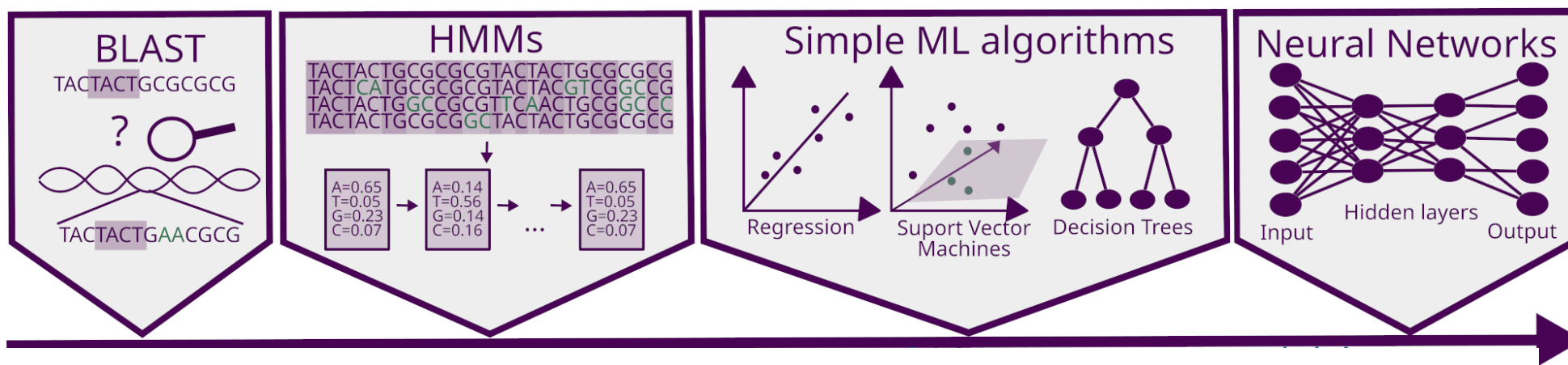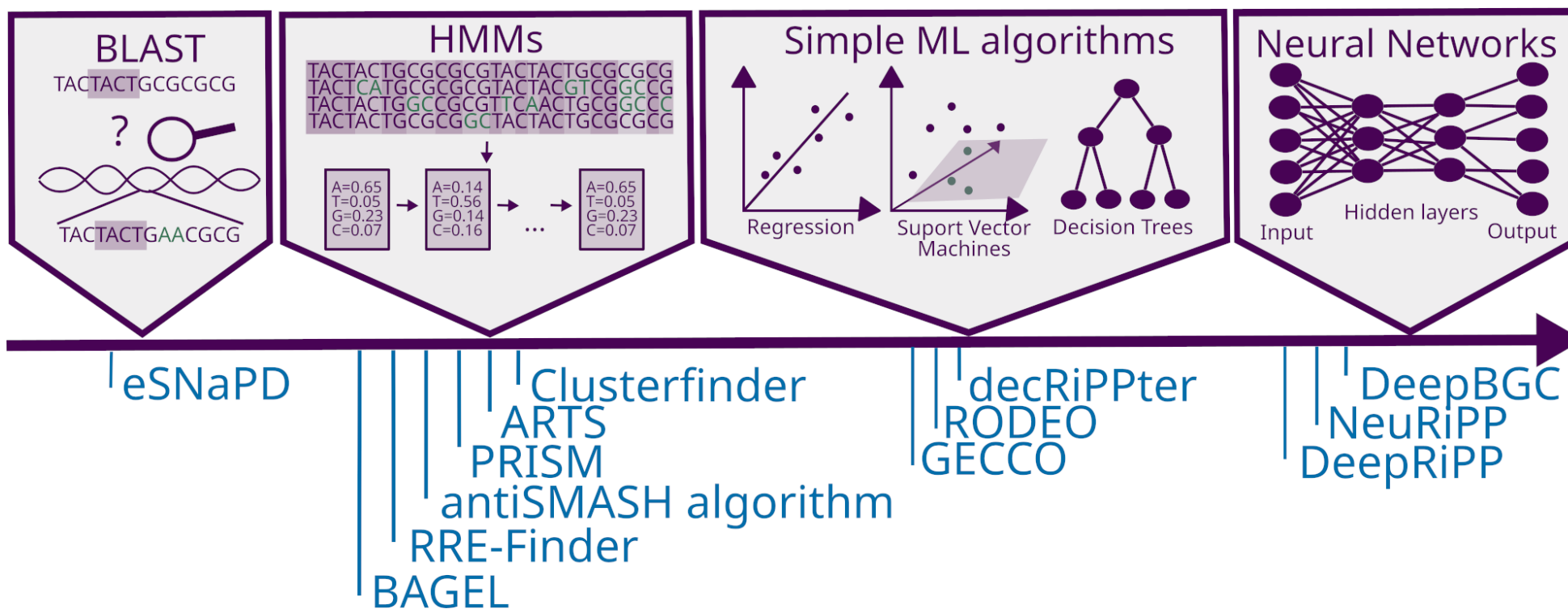
**Homology search tools/Algorithms for BGCs prediction**

*J. Org. Chem.* **2022,** *18,* 1656–1671.

Bioinformatic analysis of natural product biosynthetic capacity

Homology search tools/Algorithms for BGCs prediction

BLAST
HMMs
Simple ML algorithms
Neural Networks

eSNaPD 2014
Clusterfinder 2014
ARTS 2016
PRISM 2014
antiSMASH algorithm 2011
RRE-Finder 20220
BAGEL 2006
decRiPPter 2020
RODEO 2017
GECCO 2021
DeepBGC 2019
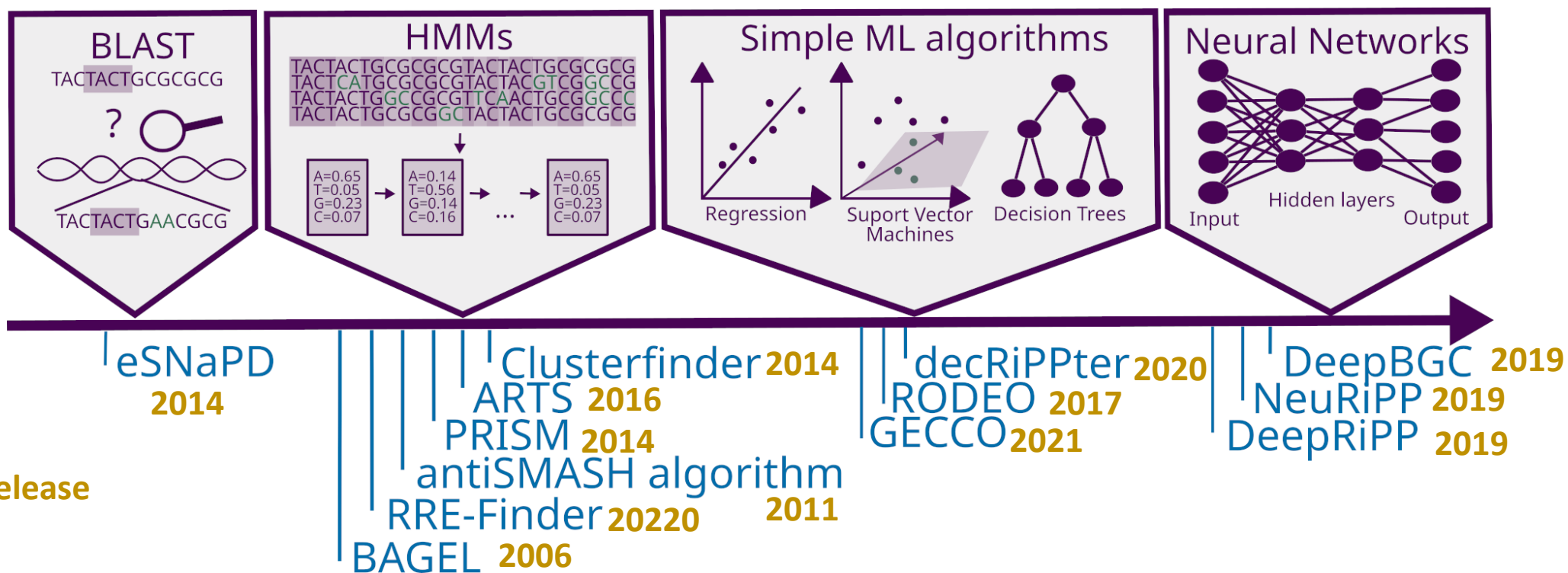NeuRiPP 2019
DeepRiPP 2019

1st version release

*J. Org. Chem.* **2022,** *18,* 1656–1671.

# Bioinformatic analysis of natural product biosynthetic capacity

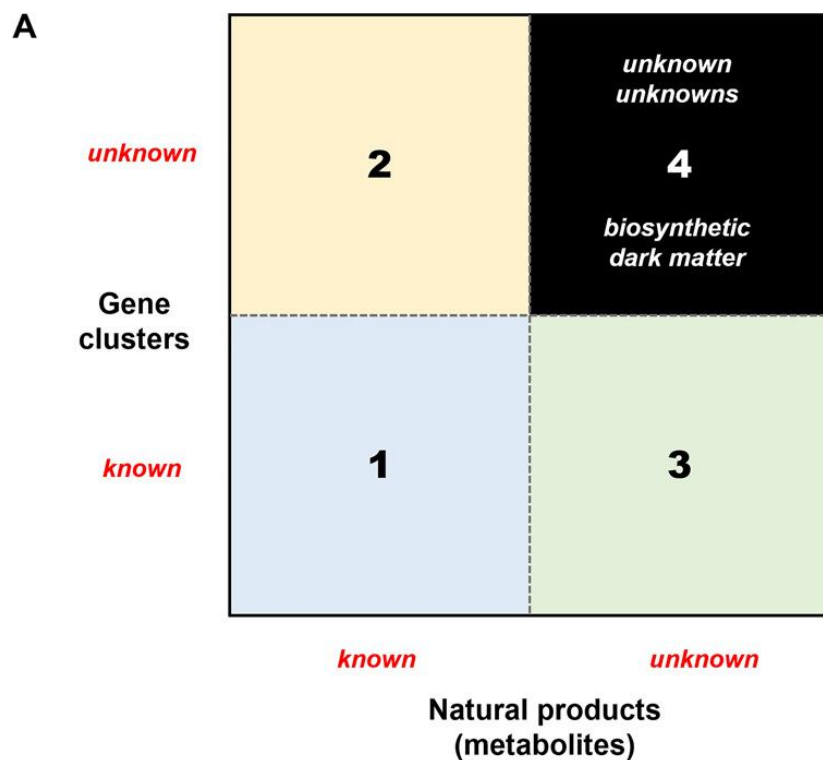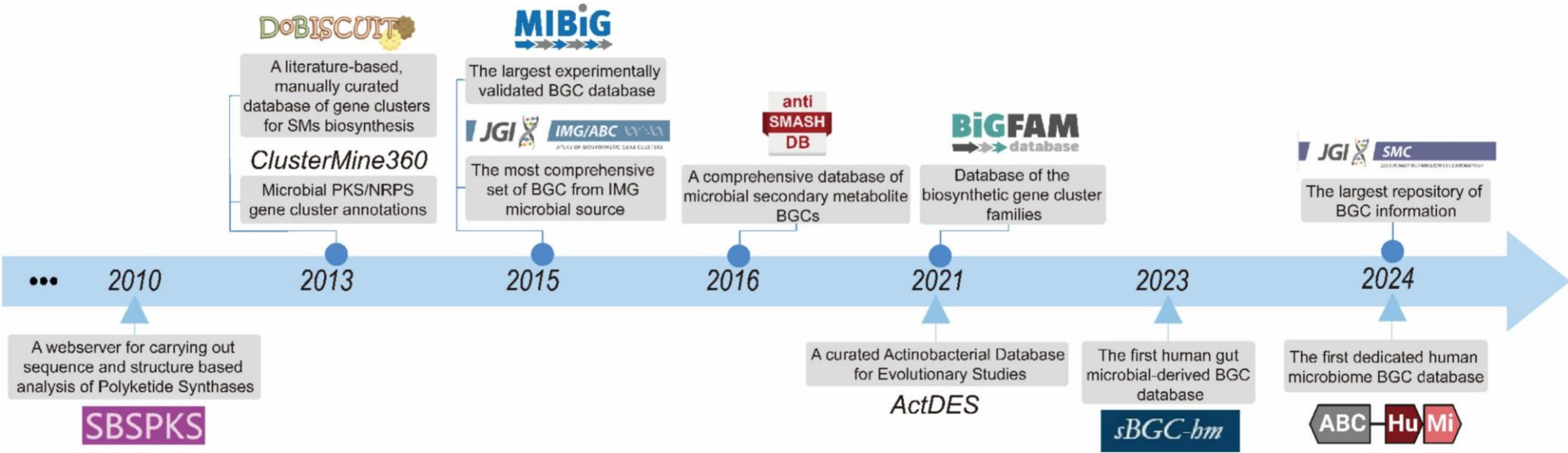## Homology search tools/Algorithms for BGCs prediction



Genome mining tools rely on existing datasets for accurate predictions.

Large databases are essential to train algorithms effectively.

*Nat. Prod. Biosynth.* **2022**, RSC Publishing

# Bioinformatic analysis of natural product biosynthetic capacity
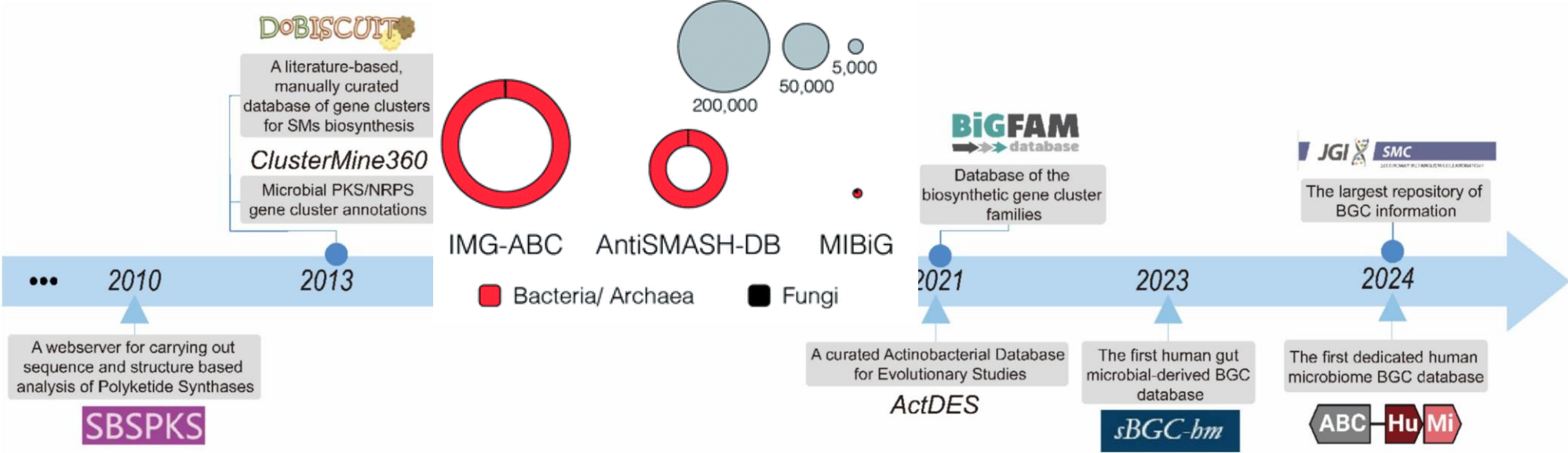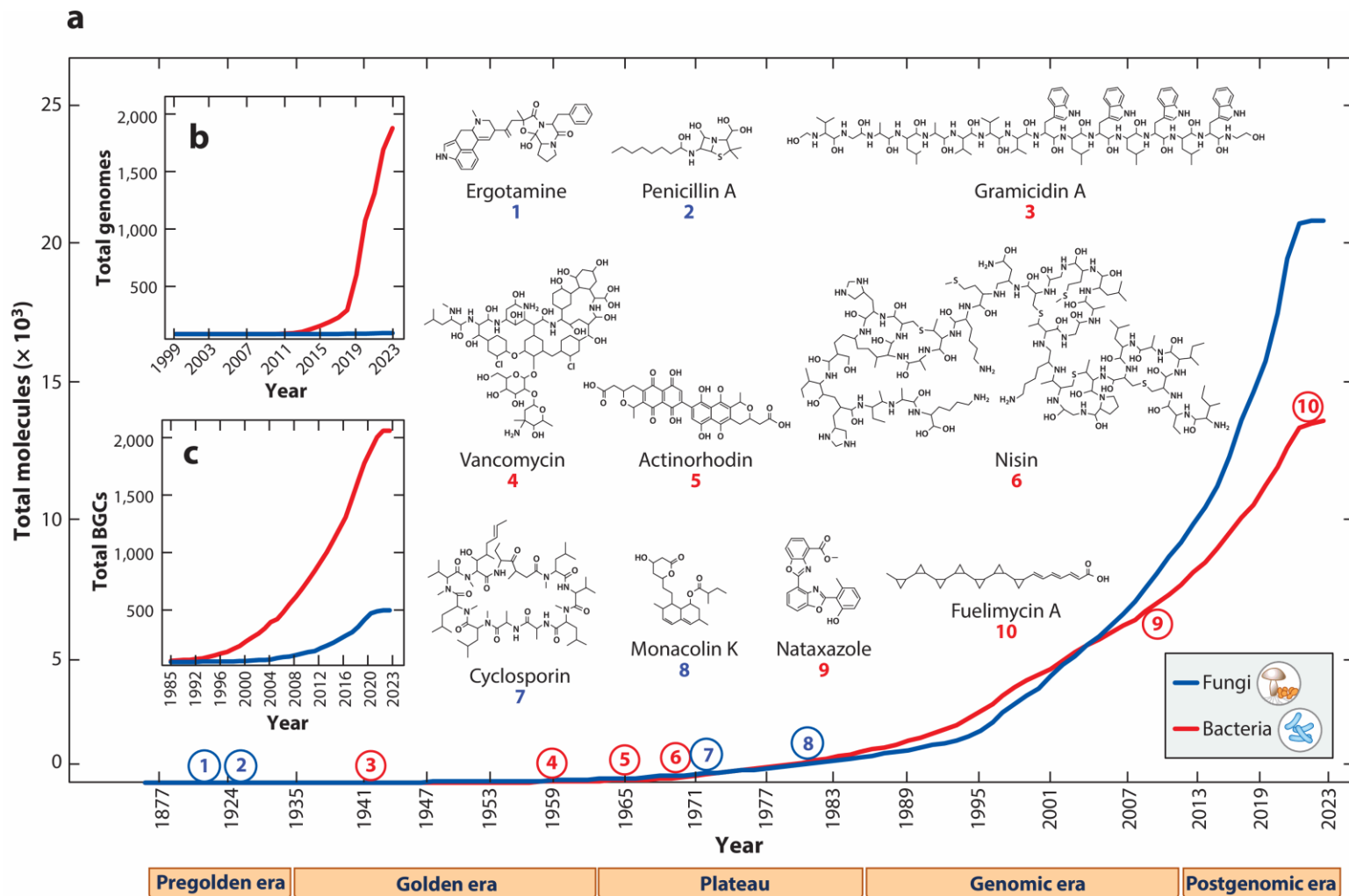
## Biosynthetic gene clusters (BGCs) databases

# Bioinformatic analysis of natural product biosynthetic capacity

## Biosynthetic gene clusters (BGCs) databases

# Mapping the biosynthetic gene clusters diversity

## Mapping the biosynthetic gene clusters diversity

**BiG-SCAPE** - Biosynthetic Gene Similarity Clustering and Prospecting Engine

**BiG-SLICE** - Biosynthetic Gene clusters - Super Linear Clustering Engine
Is a powerful tool for analyzing BGCs at scale, enabling clustering of millions of BGCs based on sequence similarity.

**BIG-FAM database**
Is an online repository for "homologous" groups of biosynthetic gene clusters (BGCs) putatively encoding the production of similar specialized metabolites. BiG-FAM facilitates querying putative BGCs to rapidly **find their position on the diversity map and gain a better understanding of their novelty or (probable) function**s, based on relationships with other known and predicted BGCs from publicly available data.

# Mapping the biosynthetic gene clusters diversity

**BiG-SCAPE** - Biosynthetic Gene Similarity Clustering and Prospecting Engine

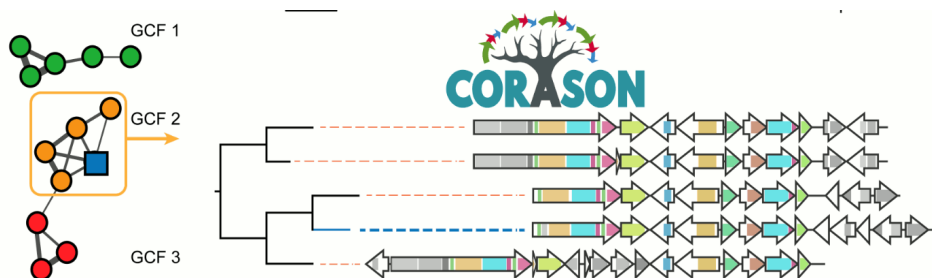| | |
|---|---|
| Introduction to antiSMASH and BiG-SCAPE workflows | Catarina Loureiro |
| | |
| Hands- on session antiSMASH and BiG-SCAPE | Catarina Loureiro |

# Mapping the biosynthetic gene clusters diversity

**CORASON** - CORe Analysis of Syntenic Orthologs to prioritize Natural Product-Biosynthetic Gene Cluster

**cblaster** - a remote search tool for rapid identification and visualization of homologous gene clusters

**clinker** - a pipeline for easily generating publication-quality gene cluster comparison figures.



11h00  30'+45'

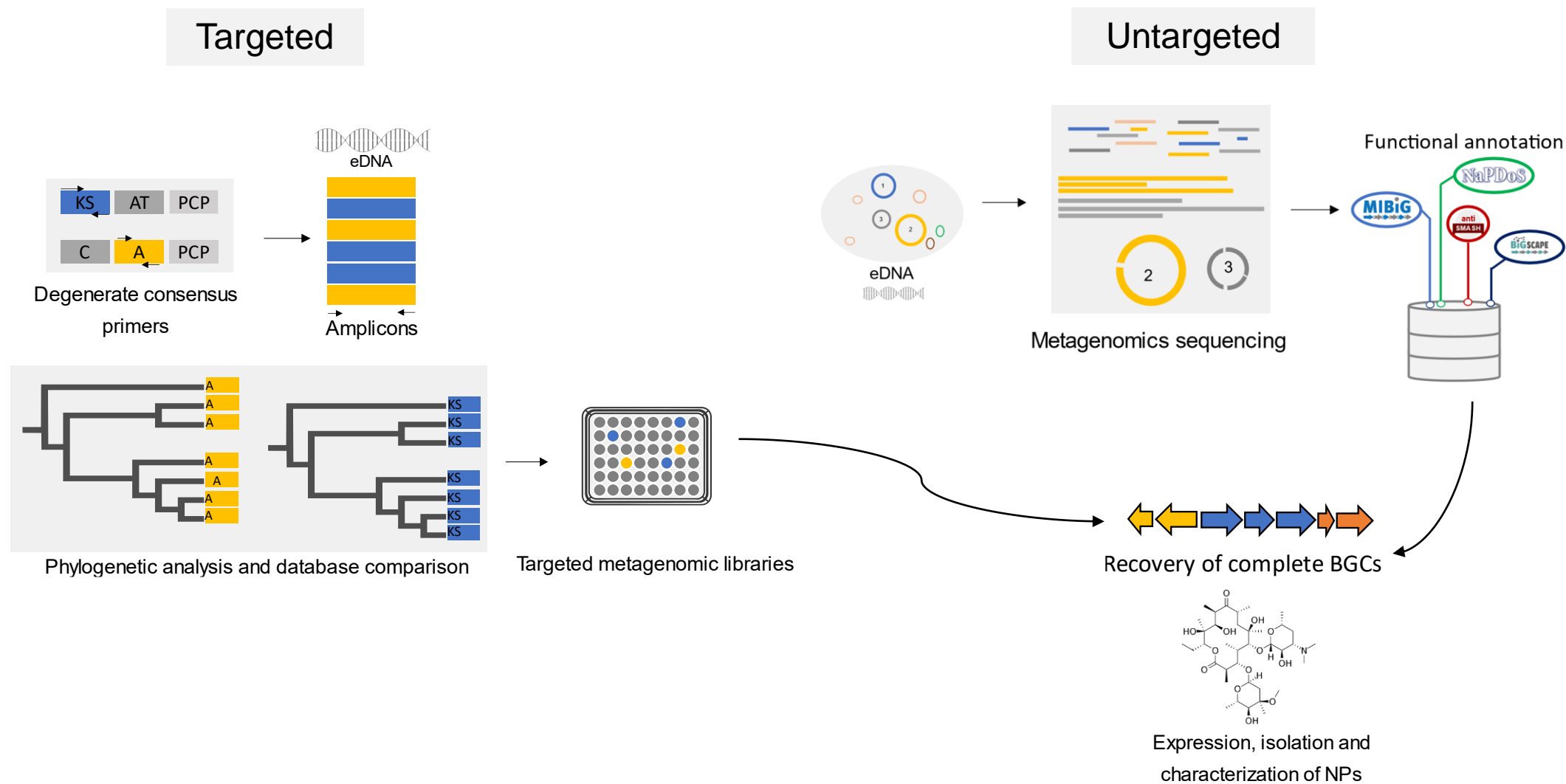Genomic context of target biosynthetic genes

(short) Hands-on session genomic context

Adriana Rego

Adriana Rego and Catarina Loureiro

# A new era of metagenomics-driven NPs discovery

Targeted

Untargeted

eDNA

KS AT PCP
C A PCP
Degenerate consensus primers

Amplicons

eDNA

Metagenomics sequencing

Functional annotation

MIBiG NaPDoS antiSMASH BiGSCAPE

A A A A A A A
KS KS KS KS KS KS KS
Phylogenetic analysis and database comparison

Targeted metagenomic libraries

Recovery of complete BGCs

Expression, isolation and characterization of NPs

# A new era of metagenomics-driven NPs discovery

## Amplicon-based or targeted approaches



**A**

PKS Gene

PKSEnzyme

| KS | AT | PCP | KS | AT | PCP | KS | AT | PCP |

domain / module

NRPS Gene

NRPSEnzyme

| C | A | PCP | C | A | PCP | C | A | PCP |

domain / module

**B**

| KS | AT | PCP |
| C | A | PCP |

Degenerate consensus primers

eDNA

Amplicons

Phylogenetic analysis and database comparison

- Known compound
- New derivative
- Unknown clade
- New class of compounds

**a** Rare and **common** BGCs in the soil metagenome

Untargeted sequencing

de novo assembly

Targeted amplification of biosynthetic domains

Amplicon sequencing

A B C D

**Main applications:**
Poorly assembled genomes
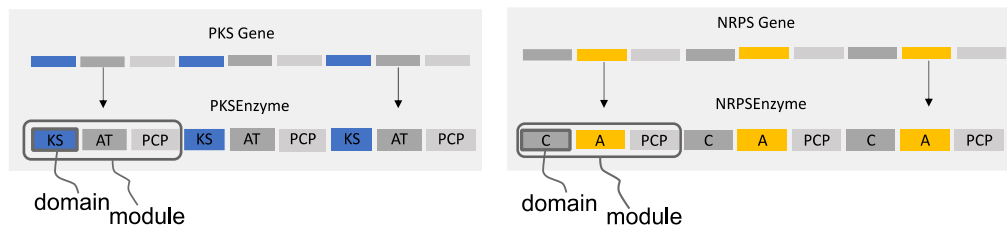Amplicons
Metagenomes
**Advantage** - quick estimate biosynthetic potential
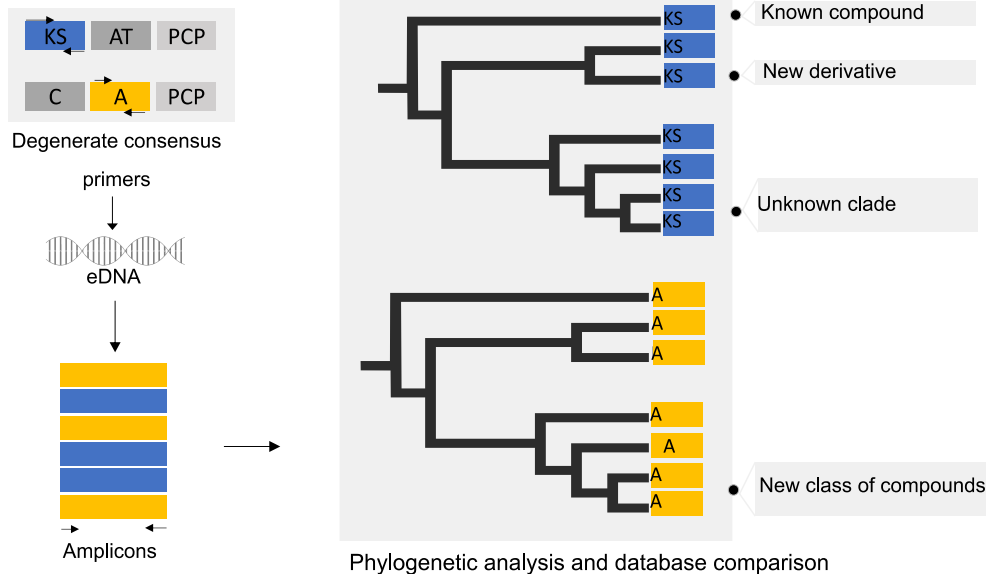New clades = new functionality

# A new era of metagenomics-driven NPs discovery

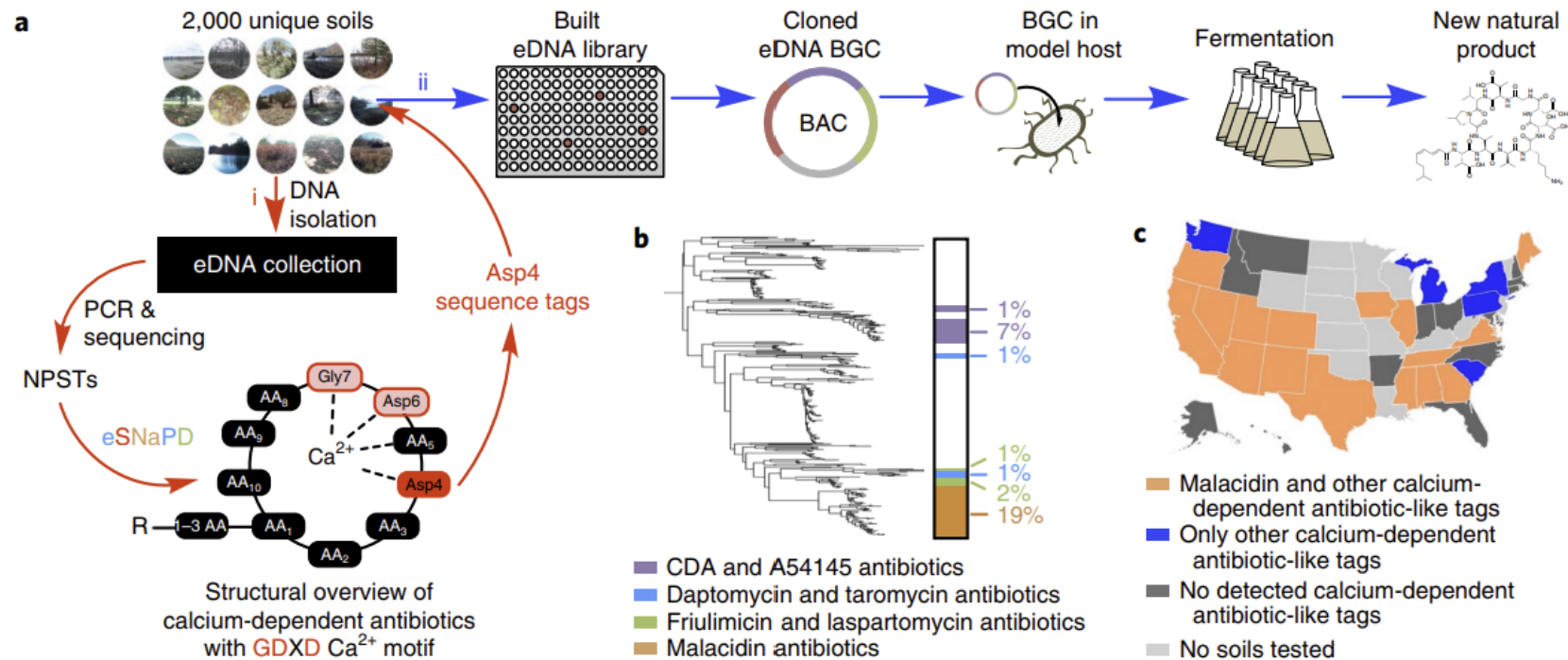## Amplicon-based or targeted approaches



**A**

PKS Gene / PKSEnzyme — KS AT PCP KS AT PCP KS AT PCP — domain, module

NRPS Gene / NRPSEnzyme — C A PCP C A PCP C A PCP — domain, module

**B**

KS AT PCP / C A PCP — Degenerate consensus primers → eDNA → Amplicons

Phylogenetic analysis and database comparison

- Known compound
- New derivative
- Unknown clade
- New class of compounds

**Hans Singh**
(University of Hawaiʻi at Mānoa, USA)

| Time | Duration | Topic |
|------|----------|-------|
| 9h00 | 60'+30' | Phylogenetic approaches to natural product discovery (short) Hands-on session NaPDos2 |

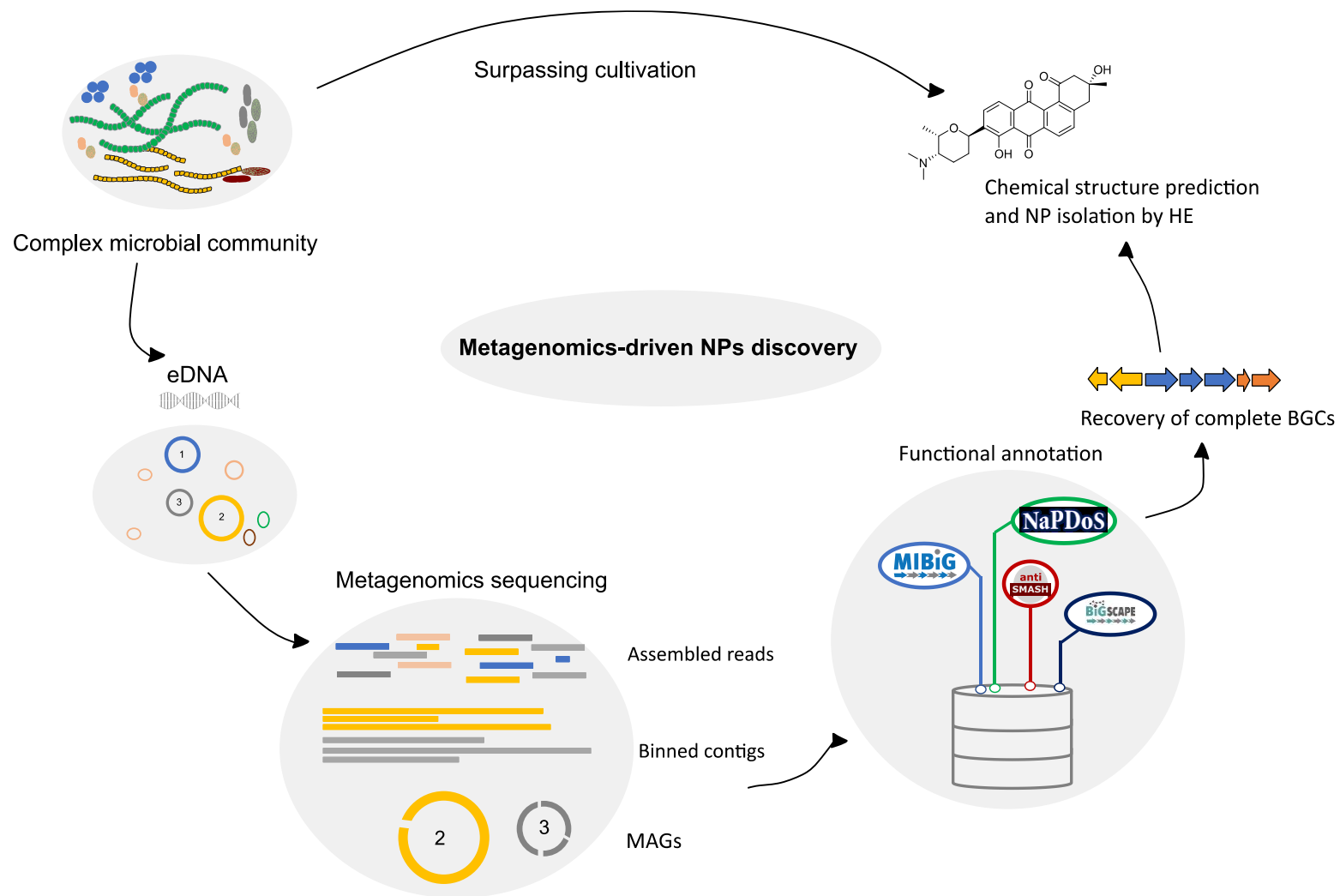# Bioinformatic analysis of natural product biosynthetic capacity

## Amplicon-based or targeted approaches



*Nat Microbiol* **2018**, **3**, 415–422

# A new era of metagenomics-driven NPs discovery


©Helena Klein

**Article**

## Biosynthetic potential of the global ocean microbiome

*Nature* **2022** 607, 111–118

**Article**

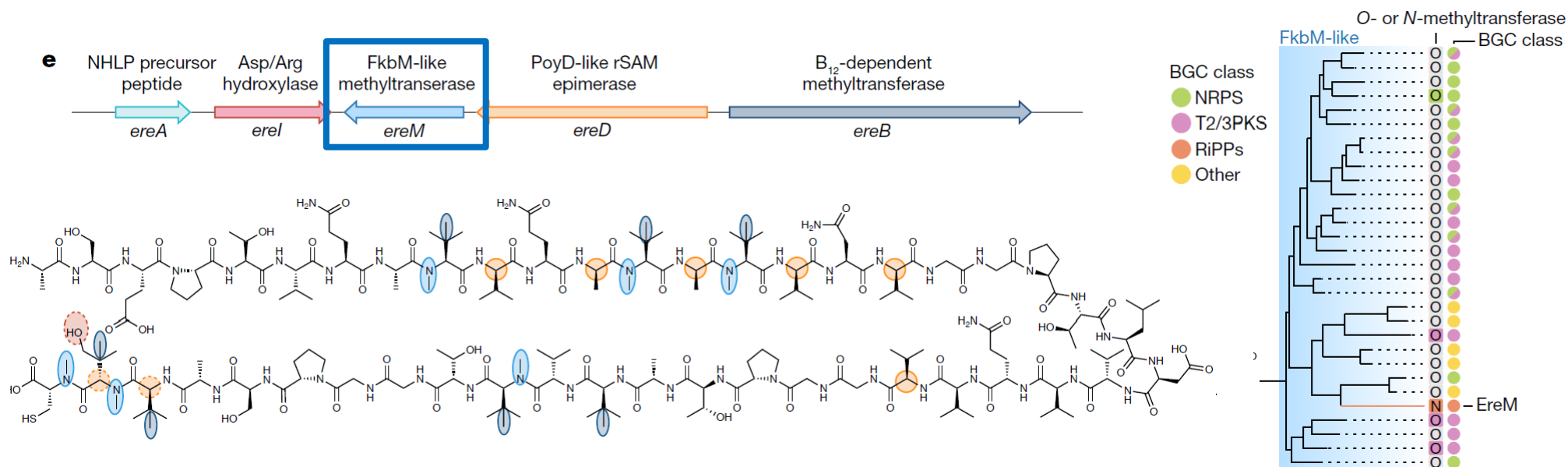## Global marine microbial diversity and its potential in bioprospecting

*Nature* **2024** 633, 371–379

# A new era of metagenomics-driven NPs discovery

**Are existent computational approaches sufficiently powerful to predict new enzymology and natural products ?**

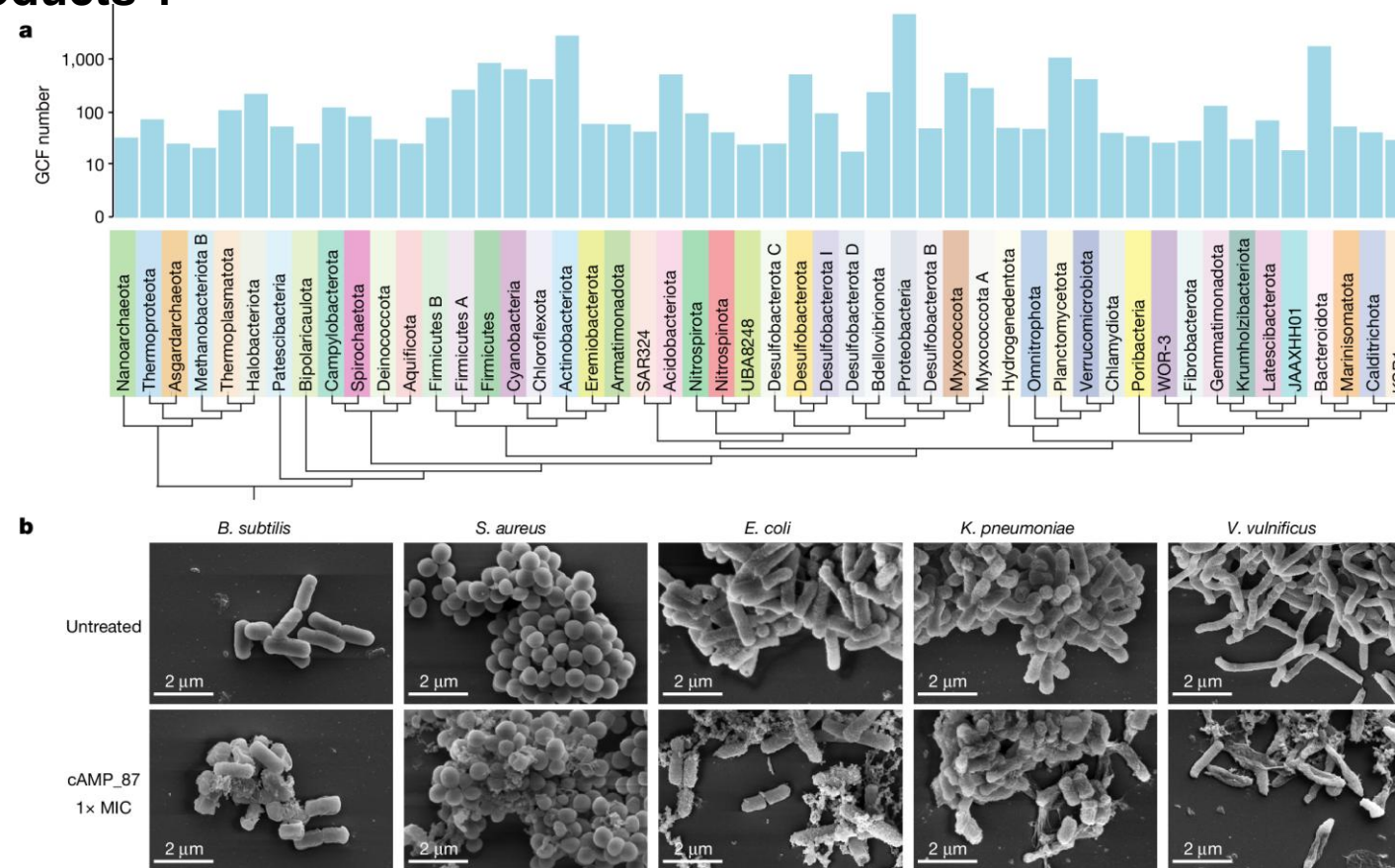unique trans amide-N-methylation



*Nature* **2022** 607, 111–118

# A new era of metagenomics-driven NPs discovery

**Are existent computational approaches sufficiently powerful to predict new enzymology and natural products ?**

**Outlook**

We still only grasp a small diversity of the existent **biosynthetic potential** and chemical diversity, including the understanding of their **ecological functions and applications**.

Genome mining is a **continuously evolving field,** although current sequence homology- and ML-based computational approaches are capable of identifying new biosynthetic pathways and enzymes, that share even low levels of similarity with known BGCs.

Future bioprospection is **environmental and metagenomics-driven**.