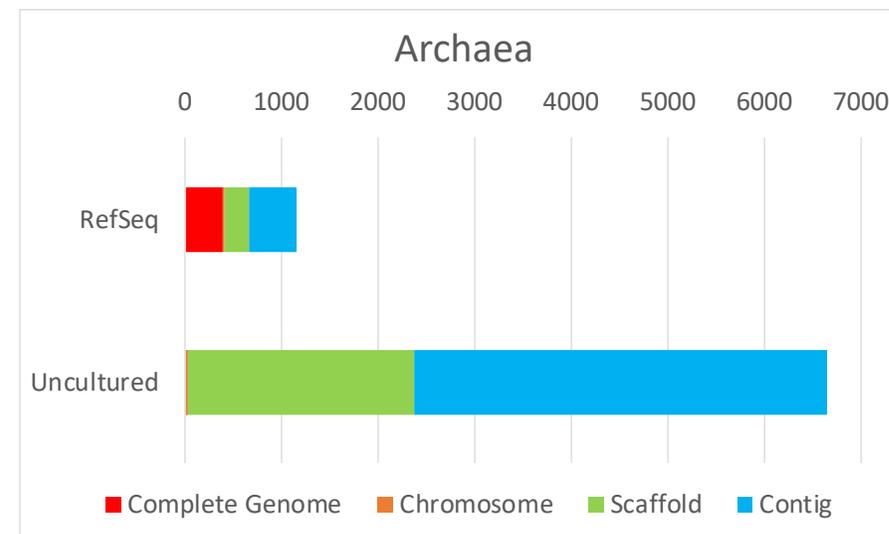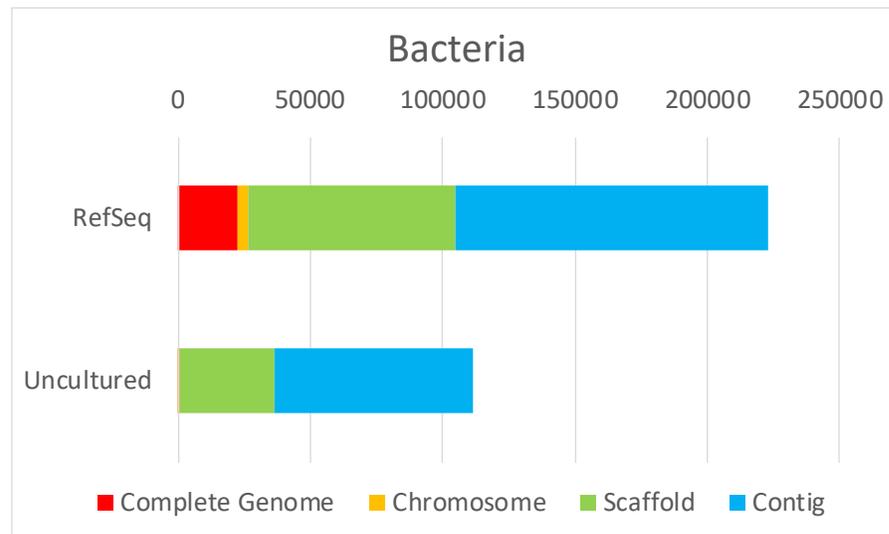# 15 MBGDのオーソログ推定機能を用いた新規ゲノム配列の機能推定

〇内山郁夫1)、三原基広2)、西出浩世1)、千葉啓和3)、高柳正彦4)、髙見英人5)

1) 基礎生物学研究所、2) ダイナコム、3)ライフサイエンス統合データベースセンター、4) ウェブブレイン、5)東大大気海洋研究所

# モチベーション：未培養菌ゲノムデータの活用

- メタゲノムからゲノム再構築などの手法で、未培養菌のゲノム解読が進み、環境中の未知微生物の生態をゲノム配列から推定することが重要になりつつある。

- MBGDのオーソログデータを用いて、新規ゲノムのアノテーションおよびゲノム機能推定に活用する。

The number of genomes released from NCBI (from Assembly Reports)

# 微生物比較ゲノムデータベースMBGD



6318 genomes
2547 species
1019 genus

5861 Bacteria
254 Archaea
203 Eukaryota

ドメイン単位のオーソロググルーピング

DomClust

DomRefine

オーソロググループ

オーソロググループ

オーソロググループ

オーソロググループ

キーワード検索
（遺伝子／オーソロググループ／生物種）

オーソログテーブル
サマリービュー

オーソログテーブル
ゲノム

遺伝子

オーソログ比較マップ

マルチプル
アライメント

# プロファイル検索によるオーソログ推定

MBGD ortholog groups

↓

Multiple sequence alignment

↓
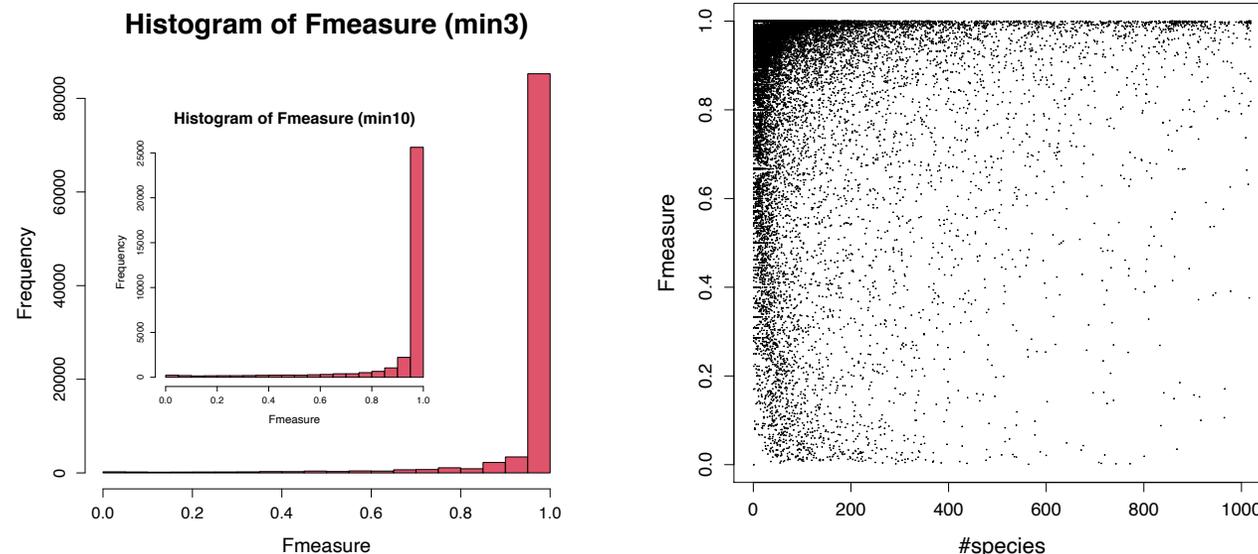
User Query → MMseqs/HMMER profiles

↓

Ortholog assignment

|  | # clusters |
|---|---|
| all | 755852 |
| min3 (members ≥ 3) | 97661 |
| min10 (members ≥10) | 38446 |

## Ortholog group assignment test (5000 queries)



|  | SearchTime(s) |
|---|---|
| hmmer-all | 20298.6 |
| mmseqs-all | 4632.0 |
| mmseqs-min3 | 1758.0 |
| mmseqs-min10 | 1105.6 |

■ correct  ■ partially correct  ■ wrong

## Result of profile quality test

Histogram of Fmeasure (min3)

Histogram of Fmeasure (min10)

# Genomapleを用いた機能モジュール充足判定

## MBGD-Genomaple pipeline



MAPLE (Takami et al. 2016):
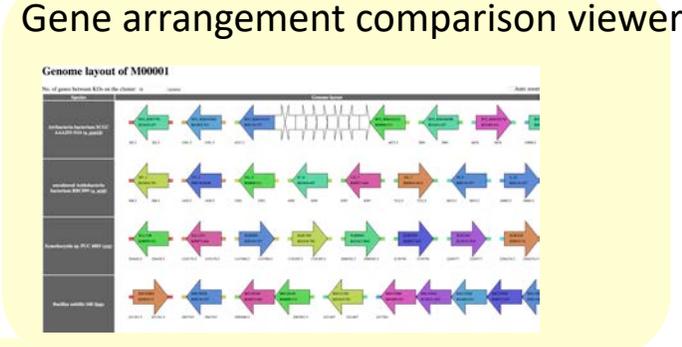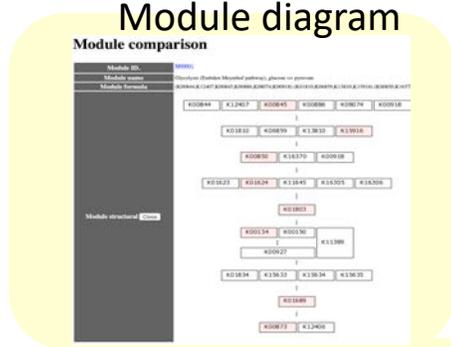ゲノム中に同定されたKOに基づいて KEGG Module の充足率を計算するツール

# MBGD-Genomaple解析のための新規MyMBGDインターフェイス

MyMBGD: upload user genomes

User Genomes

MBGD-Genomaple pipeline

Choose genomes to compare

user genome

public genome

Public genome selection

Module diagram

Gene arrangement comparison viewer

Module comparison viewer

Detailed module comparison

# 欠損遺伝子に対する代替候補遺伝子の検索