# BACBENCH: EVALUATING GENOMIC LANGUAGE MODELS FOR BACTERIA

## **Anonymous authors**

000

001

003 004

010 011

012

013

014

015

016

018

019

021

024

025

026

027

028

029

031

032

034

037

038

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

Bacteria underpin key processes in health, ecology, and biotechnology, yet machine learning in bacterial genomics lacks systematic, large-scale evaluation resources. Current resources are typically limited to single-species datasets, where the small number of available genomes leaves species-specific models underpowered, underscoring the need for approaches that can generalize across the bacterial tree of life. To address this gap, we present BacBench, the first comprehensive benchmark for bacterial genomics. BacBench consists of 11 datasets across 6 tasks, including a newly generated dataset for operon identification derived from long-read RNA sequencing. BacBench covers gene-, system-, and genome-scale prediction tasks, spanning 67k genomes, 17.6k species and 255M proteins. We analyze the performance of state-of-the-art DNA LMs, protein LMs and bacterial LMs and find that while each approach excels at different scales—the existing models fail to accurately predict the bacterial phenotype at a whole-genome level, hampering the translation to high-impact applications such as antibiotic-resistance and bioproduction. Therefore, highlighting the need to develop methods that reason over the context of the entire genomes, exploiting genomic synteny and transfer across species. We outline the key requirements for such models and release a standardized library for preprocessing, embedding, and evaluation, fostering the development of methods that accurately represent bacterial genomes, and enabling reproducible comparison of diverse approaches under a unified framework. By providing the first comprehensive benchmark dedicated to bacterial genomics, BacBench lays the ground-work for developing machine learning models that truly exploit shared evolutionary patterns and generalize across the bacterial tree of life.

## 1 Introduction

Bacteria drive indispensable processes in medicine, ecology, and biotechnology (de Steenhuijsen Piters et al., 2015; Luo et al., 2024). They produce industrial enzymes and antibiotics (Ariaeenejad et al., 2024; Santos-Júnior et al., 2024), recycle nutrients, and are being engineered for carbon capture and waste remediation (Xu & Jiang, 2024). Unlocking this potential hinges on interpreting bacterial genomes at scale. A machine learning (ML) system that can embed and reason over entire bacterial genomes could predict clinically relevant traits, surface novel enzymes for biomanufacturing, and reveal how genetic variation translates into functional capabilities across species.

Traditionally, ML approaches in bacterial genomics have been species-specific. For example, genome-wide association studies (GWAS) have been successful in identifying genotype-phenotype associations within species, such as antimicrobial resistance or virulence traits (Lees et al., 2016; Power et al., 2017; San et al., 2020). However, species-specific datasets typically include only a few genomes, leaving models statistically underpowered and prone to overfitting given the vast mutation space of bacterial genomes.

Meaningful progress therefore requires models that can share information across species. Such transfer is biologically plausible: every bacterium carries a small core of universal single-copy proteins, and many additional gene families are conserved far beyond the species level (Wang et al., 2022; Lang et al., 2013; Coleman et al., 2021). Leveraging these shared signals allows models to capture the full extent of bacterial diversity, leading to predictions that are both robust and generalizable. Training on genomes from many species, laboratories, and environmental contexts

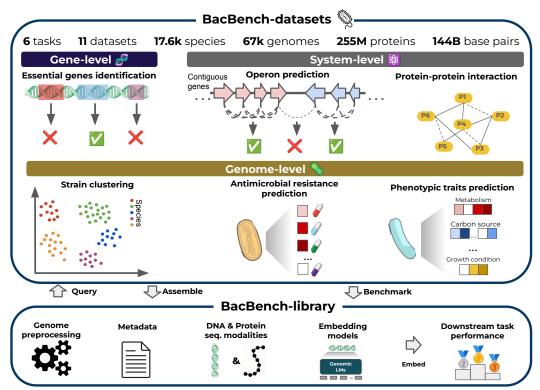


Figure 1: **BacBench overview.** We collected a diverse set of tasks at gene-, system- and genome-scales, spanning 17.6k bacterial species. BacBench-library provides support to preprocess and embed the datasets with various models. Finally, we performed systematic benchmarking for each task using diverse genomic LMs.

further guards against dataset-specific artifacts, making the resulting models less prone to sampling bias and more reliable when deployed on previously unseen strains or sequencing pipelines.

Recent breakthroughs in genomic sequence modeling show that the resulting representations can generalize across species and capture evolutionary signals (Nguyen et al., 2024; Dalla-Torre et al., 2024; Zhou et al., 2023; Lin et al., 2022; Elnaggar et al., 2021; Lin et al., 2023; Hayes et al., 2025). However, within the bacterial domain specifically, evaluation has been fragmented—either on narrow, single-task applications like antimicrobial resistance prediction (Wiatrak et al., 2024) or as part of cross-kingdom evaluations that fail to address bacteria-specific challenges (Nguyen et al., 2024). Consequently, the field lacks a dedicated, multi-scale benchmark for bacterial genomics, where genomic and metabolic mechanisms differ substantially from eukaryotes. Thus, leading to the development of models which can accurately model whole bacterial genomes.

Here, we introduce **BacBench**, the first multi-scale, multi-task and multi-species benchmark designed to evaluate ML for bacterial genomics (Fig. 1). BacBench has been collated from a diverse set of 11 datasets organized into 6 tasks spanning multiple biological scales: gene , system and genome . We consider essential gene prediction task at the gene-level, operon and protein-protein interaction prediction tasks at the system-level; and strain clustering, antibiotic resistance, and phenotypic traits prediction tasks at the genome-level. Collected from a diverse set of public resources and newly generated data, these datasets encompass 67k genomes spanning more than 17.6k bacterial species and 255M proteins. Using BacBench, we conduct comprehensive evaluation of distinct approaches to modeling bacterial genomes including DNA LMs, protein LMs (pLMs) and bacterial LMs (bLMs). We find that different modeling approaches excel at different tasks and scales, yet all models achieve low performance on phenotype prediction at a whole-genome level. Our results demonstrate (i) the need to develop ML approaches which can accurately model entire bacterial genomes, (ii) the benefits of pretraining on bacteria-specific corpora, rather than cross-kingdom ones, thus learning genomic mechanims that are unique to bacteria, and (iii) the importance of selecting the right model for the task at hand.

BacBench datasets and an accompanying toolkit for preprocessing and evaluation ensure straightforward reuse and extension<sup>1</sup>. By generating and integrating these datasets, we introduce the first benchmark suite for bacterial genomics, aiming to catalyze the development of ML methods that can transfer knowledge across species, unlocking new discoveries in bacterial genomics.

#### 2 RELATED WORK

Overall, existing resources for evaluating ML in bacterial genomics are limited in several key ways: they often lack functional labels, restrict evaluation to single species, or evaluate only one biological scale. BacBench addresses these limitations by providing a multi-scale and multi-task benchmark across the bacterial tree of life. It is specifically designed to reflect the core challenges in the field, including data sparsity within individual species, the need for cross-species generalization, and the importance of assessing model performance from the level of individual genes to entire genomes.

Bacterial genomics datasets. Large-scale sequence repositories such as MGnify (Mitchell et al., 2023), IMG/M (Markowitz et al., 2012), and AllTheBacteria (Blackwell et al., 2024) catalogue millions of bacterial assemblies. These resources are indispensable for comparative genomics, yet they provide limited metadata labels, making them ill-suited for training or benchmarking ML models. Simultaneously, task-specific collections provide information on essential genes and phenotypic traits (Zhang et al., 2004; Madin et al., 2020; Brbić et al., 2016; Weimann et al., 2016; Consortium, 2022)—supply richer annotations but usually cover only a handful of species and a single prediction setting. Diverse Genomic Embedding Benchmark (DGEB) (West-Roberts et al., 2024) offers tasks drawn from all domains of life, yet its bacterial coverage is mostly limited to single-species gene or short-segment datasets and lacks genome-scale evaluations such as broad phenotypic traits inference or strain-level clustering, leaving it unable to assess whether models generalize across thousands of species at multiple scales. BacBench complements these efforts by integrating six heterogeneous tasks, ranging from gene essentiality to genome-wide phenotype prediction - into a unified framework that explicitly tests generalization across 17.6 k bacterial species and three biological scales (Fig. 1).

Single-species bacterial genomics models. Despite the existence of large-scale bacterial genomics datasets, the ML applications in bacterial genomics have been mostly confined to single-species and single-task problems. For instance, genome-wide association studies (GWAS) have been successful in identifying genotype—phenotype associations within species, such as antimicrobial resistance or virulence traits (Lees et al., 2016; Power et al., 2017; San et al., 2020). More recent approaches such as unitig-based and deep learning models improved genotype—phenotype mapping by spanning the full pangenome (Lees et al., 2018; 2020) and predicting the effect of mutation based on the genomic context (Wiatrak et al., 2024). While the single-species models often perform well in their domains, they do not extend to other taxa and new genomic variants, and the huge genomic feature space relative to the number of labelled isolates leaves them prone to overfitting.

**Genomic LMs. DNA LMs** learn DNA sequence representations and have been shown to accurately represent long sequences (Zhou et al., 2023; Dalla-Torre et al., 2024; Jiang et al., 2023; Mourad, 2025; Nguyen et al., 2024; 2023), but are usually evaluated on human regulatory tasks, where transcriptional mechanisms and epigenomic regulation differ substantially from bacteria (Casadesús & Low, 2006). Moreover, even the DNA LMs with very large context window cannot span entire medium-sized bacterial genomes (Brixi et al., 2025). **Protein LMs** (pLMs) learn representations that correlate with structure, stability, and function (Elnaggar et al., 2021; Lin et al., 2022; 2023; Hayes et al., 2025). As bacterial genomes consist largely of coding sequence and possess simpler regulatory architectures than eukaryotes, pLMs can capture a substantial fraction of relevant biology. pLMs model proteins in isolation, however, characterizing bacterial genomes requires modeling the contextual interactions between the proteins present in the genome. Finally, we differentiate a third group of genomic LMs - Bacterial LMs (bLMs) which are recently proposed genomic LMs that are purposefully built to model bacterial genomes. These include gLM2 (Cornman et al., 2024), a mixed-modality genomic LM that represents coding regions as amino acids and intergenic regions as nucleotides to model contiguous genome context, and Bacformer, a genome-level contextual protein LM that treats each bacterial genome as an ordered sequence of proteins, refining each protein vector in the presence of

<sup>1</sup>https://anonymous.4open.science/r/BacBench-B6EF

Table 1: Summary of benchmarked models. "Max ctx." = maximum context length supported at inference; "dim" = dimensionality of the output of the last hidden layer. DNA LMs, pLMs and bLMs are separated by a horizontal line.

Model	Input	Objective	Tokenisation	Params	dim	Training corpus	Max ctx.
Mistral-DNA (Mourad, 2025)	DNA	Autoregressive	Byte-pair	138M	768	Bacteria	512
DNABERT-2 (Zhou et al., 2023)	DNA	Masked	Byte-pair	117M	768	Multi-kingdom	512
Nucleotide Transformer (Dalla-Torre et al., 2024)	DNA	Masked	k-mer	250M	768	Multi-kingdom	2,048
ProkBERT (Ligeti et al., 2024)	DNA	Masked	k-mer	27M	384	Bacteria	4,096
Evo (Nguyen et al., 2024)	DNA	Autoregressive	Single nucleotide	6.5B	4,096	Multi-kingdom	8,192
ESM-2 (Lin et al., 2022)	Single protein seq.	Masked	Single amino acid	35M	480	Multi-kingdom	1,024
ESM-C (ESM Team, 2024)	Single protein seq.	Masked	Single amino acid	300M	960	Multi-kingdom	1,024
ProtBERT (Elnaggar et al., 2021)	Single protein seq.	Masked	Single amino acid	420M	1,024	Multi-kingdom	1,024
gLM2 (Cornman et al., 2024)	Mixed modality (DNA & protein seq.)	Masked	Single nucleotide/amino acid	650M	1,280	Bacteria	4,096
Bacformer (Wiatrak et al., 2025)	Multiple protein seq.	Masked	Single protein	27M	480	Bacteria	6,000

all other proteins from the same genome, thus, encoding organism-level context. Both methods model the DNA or protein in the context of a bacterial genome and are pretrained on extensive bacterial corpora.

## 3 BACTERIAL GENOME REPRESENTATIONS & BASELINES

We selected a diverse set of five DNA LMs, three pLMs and two bLMs to represent bacterial genomes and evaluate their performance across distinct tasks and scales. These genomic LMs take as input DNA, proteins, or both (gLM2) and can therefore generalize across the bacterial tree of life. Moreover, the suite spans modalities (DNA-only, single-protein, mixed DNA-protein, genome-level protein), objectives (masked vs. autoregressive), and training corpora (bacteria-specific vs. cross-kingdom), enabling a controlled comparison of how context length, modality, and pretraining data shape genome embeddings (Table 1).

**Bacterial genomes representations.** For the gene- and system-level tasks with **DNA LMs** we embed the coding sequence plus upstream promoter of the gene; we split sequences longer than the model limit L into overlapping windows of length L and average their embeddings across all windows. For genome-level tasks, we tile each genome into chunks with overlap, embed each chunk, and average the results to extract a genome embedding (Appendix B). For the **pLMs** we embed each protein present in the bacterial genome independently and average its residue embeddings. To generate genome-scale representations, we calculate the mean of all protein vectors (Appendix B). Finally, for **bLMs** we use contiguous, genome-aware inputs: for **gLM2**, we feed mixed-modality genomic segments that encode coding regions as amino acids and intergenic regions as nucleotides, preserving local genome context across adjacent genes; for **Bacformer**, we represent each genome as an ordered sequence of proteins by obtaining per-protein tokens from its base pLM (ESM-2 35M), and pass these through a transformer to learn contextualized, genome-level protein embeddings.

#### 4 PREDICTION TASKS AND EVALUATION RESULTS

In BacBench, we consider six tasks across three scales: (i) gene , (ii) system , and (iii) genome . At gene-level, we consider the task of gene essentiality. At system-level, we assess operon identification and protein-protein interaction prediction tasks. Finally, at genome-level we evaluate methods on strain clustering, antibiotic resistance and phenotypic traits prediction tasks. We briefly describe each task and include further experimental details in the Appendices A & B.

#### 4.1 GENE ESSENTIALITY PREDICTION

Identifying essential genes is crucial for (i) defining the minimal set of cellular functions, and (ii) prioritizing drug targets. To distinguish essential from non-essential genes, the methods need to generalize to phylogenetically diverse bacteria beyond the species in the training set. For each model we report the performance by (i) fitting a linear classifier on top of the frozen gene embeddings, (ii) finetuning the model to predict a binary label (i.e. essentiality of input gene; Appendix B).

**Data.** We compile the dataset from the Database of Essential Genes (DEG) (Zhang et al., 2004), a hand-curated resource which aggregates studies published for a broad range of bacteria. After quality-control filtering, the corpus comprises 51 distinct genomes spanning 37 species (Appendix

A). This amounts to 22, 486 essential and 146, 922 non-essential genes. To prevent train-test leakage, we split by genus—placing all genomes from a genus in one split—and evaluate on held-out genera, enforcing generalization to phylogenetically distant strains.

**Metrics.** Gene essentiality prediction is a binary classification task. We evaluate performance using AUROC and AUPRC and report macro-average metrics across test genomes.

Results. pLMs and bLMs substantially outperform DNA-based models in terms of both AUROC (Fig. 2) and AUPRC using both linear probing and finetuning (Appendix C), suggesting that essential genes share conserved protein motifs that the protein-based embeddings capture. bLMs perform the best, with gLM2 achieving the best results overall, and Bacformer significantly outperforming its backbone ESM-2 model, showing the benefits of incorporating genome context. When considering DNA LMs only, Evo achieves the best performance, demonstrating the benefits of scaling model and context size.

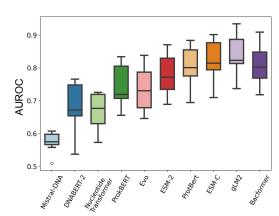


Figure 2: AUROC across genomes on essential gene prediction task using a linear model. The box spans the interquartile range with a line marking the median value.

# 4.2 OPERON IDENTIFICATION 8

Operons are multi-gene transcriptional units that underpin coordinated gene expression. Accurate operon maps are pivotal for (i) constructing gene regulatory networks (Fortino et al., 2014) and (ii) refining genome-scale metabolic models for strain engineering (Orth et al., 2010). In this task, the methods must predict whether the two neighbouring genes belong to the same operon, effectively predicting operon boundaries. Because available annotations are scarce we evaluate the task in a zero-shot setting.

**Data.** Due to the lack of experimentally validated operon annotations, we generated the operon labels by performing long-read RNA sequencing on a set of 5 diverse strains, amounting to 3, 310 unique operons (Appendix A).

**Metrics.** In BacBench, operon identification is a zero-shot binary classification task. We use

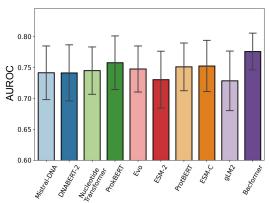


Figure 3: Zero-shot AUROC on operon identification task, the error bars represent standard error across strains.

the cosine similarity value between the two genes combined with the information on the genes' strand as a score indicating whether the genes belong to the same operon (Appendix A). We leverage AUROC and AUPRC for measuring performance. Finally, we report results across distinct strains.

**Results.** All methods except Bacformer attain similar performance (Fig. 3) in terms of both AUROC and AUPRC (Appendix C). We attribute the high performance of the Bacformer due to the (i) whole-genomic context of the model and (ii) extensive bacterial training corpus. ProkBERT performs 2nd best, showing that (i) scaling the model does not necessarily lead to improved results, (ii) the importance of a relevant pretraining corpus (ProkBERT was pretrained only on prokaryotes). Notably, both DNA and pLMs perform similarly on the task, indicating that the choice of modality is less important than incorporating genome-level context and domain-matched pretraining. Finally, the low overall performance across models highlights the need for task-specific methods and generating datasets which would allow for finetuning.

# 4.3 PROTEIN-PROTEIN INTERACTION PREDICTION

Mapping protein–protein interactions (PPI) is central to (i) reconstructing gene networks (Snider et al., 2015) and (ii) prioritizing drugtarget combinations (Wilson et al., 2022). Compared to the operon identification task where interacting genes lie next to each other, in the PPI task proteins can be separated by millions of base pairs. In this task each input is a pair of proteins, and the goal is to predict whether two proteins interact. Because interaction data are available at the protein-level, we restrict the benchmark to models which only require protein sequence data, specifically, pLMs and Bacformer. We evaluate two regimes: (i) zeroshot: the cosine similarity between the frozen embeddings of the two proteins serves as the interaction score, (ii) finetuned: the averaged

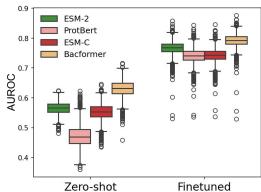


Figure 4: AUROC across genomes on PPI task in the zero-shot and finetuned setting. The box spans the interquartile range with a line marking the median value.

embeddings of the two proteins are passed through a linear classifier trained to output a binary class.

**Data.** We downloaded and processed all 10,533 bacterial strains available in STRING DB (Szklarczyk et al., 2023) together with associated PPI scores. For every strain we extract the combined interaction score for each protein pair; the median strain contains almost 640,000 scored pairs. We binarize the labels, resulting in roughly 10% of all interactions being positive (Appendix A).

**Metrics.** PPI is a binary classification task. We report performance with AUROC and AUPRC, macro-averaged over test genomes. Both the positive and negative labels are provided by STRING.

**Results.** In both zero-shot and finetuned setups, Bacformer achieves substantially higher scores than other methods. We attribute this to the rotary positional embeddings (Su et al., 2024), which increase the cosine similarity score between neighbouring genes. This is in line with experimental data, which shows that the neighbouring genes in the bacterial genomes tend to interact with each other (Dandekar et al., 1998). We also notice how the difference in performance between the methods decreases following finetuning, demonstrating how the models can adapt to the task during training. Finally, the performance of ESM-2 compared to ESM-C and ProtBERT shows that scaling up the training data and model size does not benefit PPI prediction in bacteria.

# 4.4 STRAIN CLUSTERING \*\*

Rapid clustering of whole genomes is valuable for placing newly sequenced metagenome assembled genomes (MAGs) into the bacterial tree of life and quality control. Therefore, we propose a metagenomic strain clustering task where the goal is to recover taxonomy using genome only. In this task, we feed every MAG to the model without using any species tokens or other metadata—so evaluation is fully zero-shot. We then evaluate whether models' genome embeddings preserve the taxonomy. A good embedding should cluster genomes from the same species close together and, at broader levels, conserve members of e.g. the same genus or family. We perform the clustering by computing k nearest neighbors across different k and running Leiden clustering (Traag et al., 2019) at various resolutions (Appendix B).

**Data.** In this BacBench task, we draw 6,071 strains from MGnify (Mitchell et al., 2023), spanning 25 species distributed across 10 genera and 7 families chosen to give a balanced phylogenetic spread. We process DNA and protein inputs in the same way for every model; and average the vector from the final hidden layer over the whole genome to yield one fixed-length embedding per genome.

**Metrics.** We quantify clustering performance with the adjusted Rand index (ARI), normalized mutual information (NMI) and average silhouette width (ASW). Higher is better for all metrics. To obtain cluster assignments for each model we run Leiden clustering across resolutions, retaining the resolution that maximized the mean of the three metrics over the species, genus and family ranks. The ASW is unaffected by taxonomic rank, so we report it once in the *Combined* section.

325

326 327 328

338

339

340

341

342

343

344

345

346 347

348 349

350

351

352

353

354

355

356

357

358

359 360

361

362

363

364

366 367

368

369

370 371

372

373

374

375

376

377

Table 2: Clustering performance (higher is better) across taxonomic ranks and metrics. Best value for each metric is bolded. DNA LMs, pLMs and bLMs are separated by a horizontal line.

Method	Spe	Species Genus		Family		Combined			
	ARI	NMI	ARI	NMI	ARI	NMI	ARI	NMI	ASW
Mistral-DNA	89.57	94.85	69.04	86.60	48.52	78.56	69.04	86.67	39.88
DNABERT-2 Nucleotide Transformer	97.10 98.14	98.33 98.89	64.98 64.06	85.81 85.28	43.85 $43.14$	76.35 $75.84$	68.64 68.45	86.83 86.67	63.10 $39.24$
ProkBERT Evo	<b>98.75</b> 55.56	<b>99.38</b> 76.98	63.55 $49.28$	84.88 $72.42$	$42.76 \\ 35.53$	$75.46 \\ 66.58$	68.35 $46.79$	86.58 $72.00$	<b>65.03</b> 25.33
ESM-2	50.94	71.67	56.82	72.20	46.84	69.06	51.53	70.97	16.75
ESM-C ProtBERT	$72.65 \\ 68.53$	$87.70 \\ 86.47$	79.76 <b>85.98</b>	89.87 <b>92.40</b>	60.39 <b>65.78</b>	83.43 <b>86.42</b>	70.93 <b>73.43</b>	87.00 <b>88.43</b>	$29.55 \\ 39.12$
gLM2 Bacformer	72.33 79.39	84.03 90.62	66.60 77.92	82.27 90.20	49.03 54.74	75.93 80.71	62.65 70.68	80.74 87.17	40.16 30.12

**Results.** At the species level, DNA LMs attain the highest scores, indicating that species boundaries can be recovered from sequence alone and consistent with the extra signal carried in non-coding regions. Moving up to genus and family levels, the advantage shifts: pLMs and Bacformer overtake the DNA LMs, suggesting that protein embeddings retain deeper evolutionary relationships more faithfully. Aggregated across all ranks, ProtBERT shows the highest ARI and NMI values, whereas the DNA models have the overall highest ASW. Interestingly, gLM2 which is a mixed modality model combining DNA and protein sequences performs worse than other models. Upon investigation, we believe this is due to the large variance in its embeddings and suggest that a mixed modality method with improved regularization could capture both fine-grained strain identity and higher-order phylogeny in a single embedding space.

#### 4.5 Antibiotic resistance prediction



Predicting antibiotic resistance is a task with (i) immediate clinical value for guiding antimicrobial drug therapy, (ii) monitoring the spread of resistant lineages and (iii) prioritizing compounds in drug-discovery pipelines. In this BacBench task, we define two subtasks: (1) given a bacterial genome, predict whether the strain is resistant or susceptible to a specific drug, and (2) given a bacterial genome, estimate its minimum inhibitory concentration (MIC). The first subtask is a binary classification problem (resistant vs susceptible), while the second subtask is a regression problem. Due to the computational complexity of computing genome-level embeddings on over 25k genomes (Appendix B), we (i) do not include Evo in the analysis due to its size (6.5B parameters), which makes it computationally infeasible to embed the entire corpus in most academic environments (see Runtime analysis; Appendix B), (ii) perform evaluation by stacking a linear layer on top of the frozen genome representation, and fine-tune a separate linear classifier for each model.

Data. We assemble a cross-species panel from the NIH Antimicrobial Susceptibility Test browser (National Center for Biotechnology Information, 2025), covering 25, 032 strains drawn from 38 bacterial species. After quality control and removal of sparsely sampled drugs, the dataset retains 36 antibiotics for binary label and 56 for regression prediction that span diverse classes of antibiotics (Appendix A). For the binary task we use the resistant/susceptible calls and discard any ambiguous entries (Appendix A). For the regression task we extract the raw MIC values, apply a log1ptransformation to dampen heavy tails and train a separate linear model for each drug-model pair.

Metrics. We report performance in the classification setting with AUROC and AUPRC across drugs. For the MIC regression we compute the *Pearson* correlation coefficient and the coefficient of determination  $(R^2)$  averaged across antibiotics. We report mean scores across antibiotics and include full per-antibiotic tables in the Appendix C.

**Results.** On both binary and regression setup, bLMs and pLMs tend to outperform DNA LMs. This may be explained by the fact that resistance is usually acquired through the mutation in the coding sequence, with studies showing that 90% of characterized resistance-conferring variants reside in coding regions (Sandgren et al., 2009; Farhat et al., 2019). Finally, the bLM-Bacformer achieves the best results implying the importance of epistatic effects on antibiotic resistance (Trindade et al., 2009) which the model considers by modeling the interactions between all of the proteins present in the genome. Notably, the methods record highly variable performance across antibiotics, underscoring the difficulty of building a single model that generalizes across the resistance mechanisms.

**430** 

Table 3: Performance (%, higher is better) on the antibiotic-resistance prediction tasks across drugs. Values are mean  $\pm$  standard deviation across 3 runs; best for each metric is bolded. DNA LMs, pLMs and bLMs are separated by a horizontal line.

Method	Bin	ary	Regression		
	AUPRC	AUROC	R <sup>2</sup>	Pearson	
Mistral-DNA DNABERT-2 Nucleotide Transformer ProkBERT	$52.35 \pm 23.55$ $59.70 \pm 22.90$	$79.88 \pm 7.26$ $84.22 \pm 6.47$	$\begin{array}{c} 19.71 \pm 17.37 \\ 23.61 \pm 18.33 \\ 27.74 \pm 17.80 \\ 28.00 \pm 17.57 \end{array}$	$45.89 \pm 19.27$ $51.19 \pm 16.76$	
ESM-2 ESM-C ProtBERT	$63.41 \pm 23.43$	$85.68 \pm 6.81$	$31.18 \pm 17.39$ $33.15 \pm 17.62$ $28.23 \pm 18.46$	$56.79 \pm 15.46$	
gLM2 Bacformer			$26.15 \pm 18.13$ $33.84 \pm 18.60$		

# 4.6 Phenotypic traits prediction

Accurately predicting phenotype from a genomic sequence enables (i) inference of the biological or ecological function of bacteria (Feldbauer et al., 2015) and (ii) engineering organisms with the exact metabolic traits needed for efficient waste remediation (Rafeeq et al., 2023), accelerating sustainable industrial bioproduction (Lawson et al., 2021). We evaluate whether the models' genome embeddings are predictive of diverse phenotypic traits. Here, given a genome embedding, the task is to predict a trait. Similarly as in the antibiotic resistance prediction task above, we perform linear probing evaluation, training a separate linear classifier per phenotype and exclude the Evo model due to the computational costs (Appendix B).

**Data.** To create the benchmark we collated large trait inventories (Madin et al., 2020; Brbić et al., 2016; Weimann et al., 2016), apply stringent quality filters and discard traits represented by only a handful of isolates (Appendix A). The final corpus covers 139 discrete phenotypes spanning 15, 477 bacterial species, making it the broadest dataset of its kind and challenging the models to generalize well beyond their training clades. We group traits into carbon utilisation, biochemical activity, growth conditions and cellular morphology (Appendix A), which we use later for stratified analysis. As similar genomes often share the same phenotype, we split the data for each phenotype by genus—placing all genomes from a genus in one split—and evaluate on held-out genera, enforcing generalization to phylogenetically distant strains.

**Metrics.** For BacBench, we restrict the phenotypic traits to categorical traits due to their sufficient number of available samples, and evaluate performance using the macro-averaged AUROC and AUPRC over phenotype categories. We report mean scores for every phenotype group and include full per-trait tables in the Appendix C.

**Results.** pLMs and bLMs outperform DNA LMs across phenotype groups and metrics (Fig. 5 & Appendix C). The relative difference is especially large for biochemical activity and carbon utilization traits, likely because these phenotypes hinge on enzyme active-site composition and pathway membership—information that is explicit in amino-acid space and can be better captured

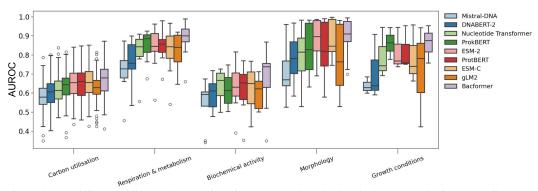


Figure 5: AUROC across diverse phenotypic traits groups and methods. The box spans the inter-quartile range with a line marking the median value.

by pLMs (Teukam et al., 2024; Lawson et al., 2021). The Bacformer bLM attains the highest scores, indicating that whole-genome context helps explain phenotypes arising from coordinated protein function and epistatic interactions (Trindade et al., 2009). In contrast, Mistral-DNA, the only autoregressive model in the task, lags well behind the other models. Overall performance remains moderate—especially for carbon utilization, biochemical activity and growth condition traits—highlighting considerable room for improvement. Progress will likely require more complex contextualized models that incorporate environmental metadata, together with larger, more balanced phenotype datasets.

# 5 DISCUSSION

Summary. BacBench addresses the lack of a comprehensive, cross-species evaluation resource for bacterial genomics by providing unified datasets and benchmarks, and by evaluating existing genomic language models across species, tasks, and biological scales. It introduces a newly generated dataset for operon identification—a key problem for refining genome-scale metabolic models—and curates five additional tasks (gene essentiality, protein-protein interaction, strain clustering, antibiotic resistance, and 139 phenotypic traits) into a single framework covering 67k genomes from 17.6k species. Our experiments show that (i) existing genomic LMs (DNA LMs, pLMs, and bLMs) capture core taxonomic structure and functional relatedness, providing a strong baseline representation of bacterial genomes, but fall short at accurate genome-to-phenotype prediction, as evidenced by antibiotic resistance and phenotype tasks; (ii) models purpose-built for bacteria (gLM2, Bacformer) or trained on bacteria-specific corpora (ProkBERT) tend to outperform broad, cross-kingdom counterparts across most tasks, underscoring the value of domain-matched pretraining and inductive biases; (iii) different modeling approaches excel at different problems—DNA LMs capture fine-grained taxonomic signals and do well on operon identification, pLMs better preserve deeper phylogeny and functional similarity for strain clustering, and bLMs like Bacformer perform best on tasks driven by multi-gene interactions (phenotypes, antibiotic resistance). All datasets are accompanied with extensive documentation; the embedding and evaluation library is provided at https://anonymous.4open.science/r/BacBench-B6EF.

**Towards accurate genome-to-phenotype bacterial prediction.** Our results suggest a practical path forward towards building model for genome-to-phenotype mapping in bacteria: (*i*) the model should represent entire genomes end-to-end to capture long-range, cross-protein dependencies (currently only feasible in Bacformer); (*ii*) pretrain on substantially larger, bacteria-focused corpora (we estimate >4M unique strains remain untapped (Mitchell et al., 2023; Markowitz et al., 2012; Blackwell et al., 2024)); (*iii*) integrate DNA and protein modalities—and where available RNA—to couple regulatory and coding signals in a single embedding while remaining computationally efficient; (*iv*) allow inclusion of structured priors (e.g., resistance gene catalogs, operon maps, HGT markers) to improve data efficiency on scarce phenotype labels; (*v*) expand high-quality phenotype supervision (including knock-outs and standardized trait panels) to close the supervision gap limiting genome-to-phenotype learning.

Limitations & Future Work. By releasing BacBench we provide a foundation for more expressive, cross-species models of bacterial genomics, yet the present benchmark covers only part of the functional landscape. Other tasksa re not yet included, and certain phenotypes and antibiotic classes remain sparse (Appendix A), underscoring the need for generating new data. We benchmarked a representative set of publicly available models capable of cross-species generalization, and expect the model suite to expand in future iterations. Finally, the current tasks do not include modalities such as transcriptomics and metabolomics due to data sparsity and inconsistent metadata. As community datasets mature, incorporating multi-omics will enable more faithful evaluation of causal, context-dependent genome-to-phenotype mappings.

We anticipate that subsequent iterations will broaden task and model coverage, ultimately enabling contextual, genome-scale representations for bacterial genomes. Community contributions of datasets, models, and evaluation routines are encouraged so that BacBench evolves into a continually updated standard for bacterial ML.

# 6 REPRODUCIBILITY STATEMENT

We release an *anonymous* codebase for preprocessing, embedding, and evaluation at https://anonymous.4open.science/r/BacBench-B6EF, together with helper utilities for bacterial genomics to lower the barrier for adding new models and tasks. All datasets in BacBench are fully documented and accompanied by anonymized MLCommons Croissant metadata files in the supplementary materials, enabling unambiguous data loading and lineage tracking. The Appendix details quality filtering, preprocessing pipelines, train/validation/test splits, model and training configurations, and hyperparameters; it also specifies random seeds, hardware, and optimization settings. Further details on the exact scripts to reproduce the results can be found in the anonymous repository linked above. These artifacts together provide the necessary pointers—code, data descriptors, and experimental specifications—to reproduce our results and to extend BacBench as an evolving benchmark for the ML community.

# REFERENCES

- Shohreh Ariaeenejad, Javad Gharechahi, Mehdi Foroozandeh Shahraki, Fereshteh Fallah Atanaki, Jian-Lin Han, Xue-Zhi Ding, Falk Hildebrand, Mohammad Bahram, Kaveh Kavousi, and Ghasem Hosseini Salekdeh. Precision enzyme discovery through targeted mining of metagenomic data. *Natural Products and Bioprospecting*, 14(1):7, 2024.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint* arXiv:1607.06450, 2016.
- Grace Blackwell, Karel M. Brinda, John A. Lees, Rebecca A. Gladstone, Philip Brownridge, et al. Allthebacteria: a uniformly assembled and searchable compendium of all public bacterial genomes. *bioRxiv*, pp. 2024.03.08.584059, 2024. doi: 10.1101/2024.03.08.584059.
- Maria Brbić, Matija Piškorec, Vedrana Vidulin, Anita Kriško, Tomislav Šmuc, and Fran Supek. The landscape of microbial phenotypic traits and associated genes. *Nucleic acids research*, pp. gkw964, 2016.
- Garyk Brixi, Matthew G Durrant, Jerome Ku, Michael Poli, Greg Brockman, Daniel Chang, Gabriel A Gonzalez, Samuel H King, David B Li, Aditi T Merchant, et al. Genome modeling and design across all domains of life with evo 2. *BioRxiv*, pp. 2025–02, 2025.
- Josep Casadesús and David Low. Epigenetic gene regulation in bacteria. *Molecular Microbiology*, 60(4):820–826, 2006. doi: 10.1111/j.1365-2958.2006.05136.x.
- Gareth A Coleman, Adrián A Davín, Tara A Mahendrarajah, Lénárd L Szánthó, Anja Spang, Philip Hugenholtz, Gergely J Szöllősi, and Tom A Williams. A rooted phylogeny resolves early bacterial evolution. *Science*, 372(6542):eabe0511, 2021.
- CRyPTIC Consortium. A data compendium associating the genomes of 12,289 mycobacterium tuberculosis isolates with quantitative resistance phenotypes to 13 antibiotics. *PLoS biology*, 20(8): e3001721, 2022.
- Andre Cornman, Jacob West-Roberts, Antonio Pedro Camargo, Simon Roux, Martin Beracochea, Milot Mirdita, Sergey Ovchinnikov, and Yunha Hwang. The omg dataset: An open metagenomic corpus for mixed-modality genomic language modeling. *bioRxiv*, pp. 2024–08, 2024.
- Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan Sirelkhatim, et al. Nucleotide transformer: building and evaluating robust foundation models for human genomics. *Nature Methods*, pp. 1–11, 2024.
- Thomas Dandekar, Berend Snel, Martijn Huynen, and Peer Bork. Conservation of gene order: a fingerprint of proteins that physically interact. *Trends in biochemical sciences*, 23(9):324–328, 1998.

- Petr Danecek, James K Bonfield, Jennifer Liddle, John Marshall, Valeriu Ohan, Martin O Pollard, Andrew Whitwham, Thomas Keane, Shane A McCarthy, Robert M Davies, and Heng Li. Twelve years of SAMtools and BCFtools. *GigaScience*, 10(2), 02 2021. ISSN 2047-217X. doi: 10. 1093/gigascience/giab008. URL https://doi.org/10.1093/gigascience/giab008.
  - Wouter AA de Steenhuijsen Piters, Elisabeth AM Sanders, and Debby Bogaert. The role of the local microbial ecosystem in respiratory health and disease. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1675):20140294, 2015.
  - Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, et al. Prottrans: Toward understanding the language of life through self-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(10):7112–7127, 2021.
  - ESM Team. Esm cambrian: Revealing the mysteries of proteins with unsupervised learning. https://evolutionaryscale.ai/blog/esm-cambrian, 2024. EvolutionaryScale website, published 4 December 2024, accessed 17 April 2025.
  - Maha R Farhat, Luca Freschi, Roger Calderon, Thomas Ioerger, Matthew Snyder, Conor J Meehan, Bouke de Jong, Leen Rigouts, Alex Sloutsky, Devinder Kaur, et al. Gwas for quantitative resistance phenotypes in mycobacterium tuberculosis reveals resistance genes and regulatory regions. *Nature communications*, 10(1):2128, 2019.
  - Roman Feldbauer, Frederik Schulz, Matthias Horn, and Thomas Rattei. Prediction of microbial phenotypes based on comparative genomics. *BMC bioinformatics*, 16:1–8, 2015.
  - Vittorio Fortino, Olli-Pekka Smolander, Petri Auvinen, Roberto Tagliaferri, and Dario Greco. Transcriptome dynamics-based operon prediction in prokaryotes. BMC bioinformatics, 15:1–14, 2014.
  - Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, et al. Simulating 500 million years of evolution with a language model. *Science*, pp. eads0018, 2025.
  - Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL https://arxiv.org/abs/2310.06825.
  - Diederik P Kingma. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
  - Jenna Morgan Lang, Aaron E Darling, and Jonathan A Eisen. Phylogeny of bacterial and archaeal genomes using conserved genes: supertrees and supermatrices. *PloS one*, 8(4):e62510, 2013.
  - Christopher E Lawson, Jose Manuel Martí, Tijana Radivojevic, Sai Vamshi R Jonnalagadda, Reinhard Gentz, Nathan J Hillson, Sean Peisert, Joonhoon Kim, Blake A Simmons, Christopher J Petzold, et al. Machine learning for metabolic engineering: A review. *Metabolic Engineering*, 63:34–60, 2021.
  - John A. Lees, Minna Vehkala, Niko Välimäki, Simon R. Harris, and Claire \*et al.\* Chewapreecha. Sequence element enrichment analysis to determine the genetic basis of bacterial phenotypes. *Nature Communications*, 7:12797, 2016. doi: 10.1038/ncomms12797.
  - John A Lees, Marco Galardini, Stephen D Bentley, Jeffrey N Weiser, and Jukka Corander. Pyseer: a comprehensive tool for microbial pangenome-wide association studies. *Bioinformatics*, 34(24): 4310–4312, 2018.
  - John A Lees, T Tien Mai, Marco Galardini, Nicole E Wheeler, Samuel T Horsfield, Julian Parkhill, and Jukka Corander. Improved prediction of bacterial genotype-phenotype associations using interpretable pangenome-spanning regressions. *MBio*, 11(4):10–1128, 2020.

- Heng Li. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–3100, 2018
  - Balázs Ligeti, István Szepesi-Nagy, Babett Bodnár, Noémi Ligeti-Nagy, and János Juhász. Prokbert family: genomic language models for microbiome applications. *Frontiers in Microbiology*, 14, 2024. ISSN 1664-302X. doi: 10.3389/fmicb.2023.1331233. URL https://www.frontiersin.org/articles/10.3389/fmicb.2023.1331233.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Sal Candido, et al. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv*, 2022.
- Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
- Zhaowei Luo, Zhanghua Qi, Jie Luo, and Tingtao Chen. Potential applications of engineered bacteria in disease diagnosis and treatment. *Microbiome Research Reports*, 4(1):N–A, 2024.
- Joshua S Madin, Daniel A Nielsen, Maria Brbic, Ross Corkrey, David Danko, Kyle Edwards, Martin KM Engqvist, Noah Fierer, Jemma L Geoghegan, Michael Gillings, et al. A synthesis of bacterial and archaeal phenotypic trait data. *Scientific data*, 7(1):170, 2020.
- Victor M Markowitz, I-Min A Chen, Krishna Palaniappan, Ken Chu, Ernest Szeto, Yuri Grechkin, Anna Ratner, Biju Jacob, Jinghua Huang, Peter Williams, et al. Img: the integrated microbial genomes database and comparative analysis system. *Nucleic acids research*, 40(D1):D115–D122, 2012.
- Alex L. Mitchell, Antonio Almeida, Martin Beracochea, Matthew Boland, Jack Burgin, Guy Cochrane, et al. Mgnify: the microbiome analysis resource in 2023. *Nucleic Acids Research*, 51(D1):D723–D730, 2023. doi: 10.1093/nar/gkac1088.
- Raphael Mourad. Mistral-DNA. https://github.com/raphaelmourad/Mistral-DNA, 2025. GitHub repository, accessed 17 April 2025.
- National Center for Biotechnology Information. Antibiotic susceptibility test (ast) browser, 2025. URL https://www.ncbi.nlm.nih.gov/pathogens/ast. Accessed 24 Apr 2025.
- Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Michael Wornow, Callum Birch-Sykes, Stefano Massaroli, Aman Patel, Clayton Rabideau, Yoshua Bengio, et al. Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. *Advances in neural information processing systems*, 36:43177–43201, 2023.
- Eric Nguyen, Michael Poli, Matthew G Durrant, Brian Kang, Dhruva Katrekar, David B Li, Liam J Bartie, Armin W Thomas, Samuel H King, Garyk Brixi, et al. Sequence modeling and design from molecular to genome scale with evo. *Science*, 386(6723):eado9336, 2024.
- Jeffrey D Orth, Ines Thiele, and Bernhard Ø Palsson. What is flux balance analysis? *Nature biotechnology*, 28(3):245–248, 2010.
- A Paszke. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703*, 2019.
- Robert A. Power, Julian Parkhill, and Tulio de Oliveira. Microbial genome-wide association studies: lessons from human gwas. *Nature Reviews Genetics*, 18:41–50, 2017. doi: 10.1038/nrg.2016.132.
- Hamza Rafeeq, Nadia Afsheen, Sadia Rafique, Arooj Arshad, Maham Intisar, Asim Hussain, Muhammad Bilal, and Hafiz MN Iqbal. Genetically engineered microorganisms for environmental remediation. *Chemosphere*, 310:136751, 2023.
- James Emmanuel San, Shakuntala Baichoo, Aquillah Kanzi, Yumna Moosa, Richard Lessells, Vagner Fonseca, John Mogaka, Robert Power, and Tulio de Oliveira. Current affairs of microbial genomewide association studies: approaches, bottlenecks and analytical pitfalls. *Frontiers in microbiology*, 10:3119, 2020.

- Andreas Sandgren, Michael Strong, Preetika Muthukrishnan, Brian K Weiner, George M Church, and Megan B Murray. Tuberculosis drug resistance mutation database. *PLoS medicine*, 6(2):e1000002, 2009.
  - Célio Dias Santos-Júnior, Marcelo DT Torres, Yiqian Duan, Álvaro Rodríguez Del Río, Thomas SB Schmidt, Hui Chong, Anthony Fullam, Michael Kuhn, Chengkai Zhu, Amy Houseman, et al. Discovery of antimicrobial peptides in the global microbiome with machine learning. *Cell*, 187 (14):3761–3778, 2024.
  - Jamie Snider, Max Kotlyar, Punit Saraon, Zhong Yao, Igor Jurisica, and Igor Stagljar. Fundamentals of protein interaction network mapping. *Molecular systems biology*, 11(12):848, 2015.
  - Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
  - Damian Szklarczyk, Rebecca Kirsch, Mikaela Koutrouli, Katerina Nastou, Farrokh Mehryary, Radja Hachilif, Annika L Gable, Tao Fang, Nadezhda T Doncheva, Sampo Pyysalo, et al. The string database in 2023: protein–protein association networks and functional enrichment analyses for any sequenced genome of interest. *Nucleic acids research*, 51(D1):D638–D646, 2023.
  - Yves Gaetan Nana Teukam, Loïc Kwate Dassi, Matteo Manica, Daniel Probst, Philippe Schwaller, and Teodoro Laino. Language models can identify enzymatic binding sites in protein sequences. *Computational and Structural Biotechnology Journal*, 23:1929–1937, 2024.
  - Vincent A Traag, Ludo Waltman, and Nees Jan Van Eck. From louvain to leiden: guaranteeing well-connected communities. *Scientific reports*, 9(1):1–12, 2019.
  - Sandra Trindade, Ana Sousa, Karina Bivar Xavier, Francisco Dionisio, Miguel Godinho Ferreira, and Isabel Gordo. Positive epistasis drives the acquisition of multidrug resistance. *PLoS genetics*, 5(7): e1000578, 2009.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
  - Saidi Wang, Minerva Ventolero, Haiyan Hu, and Xiaoman Li. A revisit to universal single-copy genes in bacterial genomes. *Scientific Reports*, 12(1):14550, 2022. doi: 10.1038/s41598-022-18762-z.
  - Aaron Weimann, Kyra Mooren, Jeremy Frank, Phillip B Pope, Andreas Bremges, and Alice C McHardy. From genomes to phenotypes: Traitar, the microbial trait analyzer. *MSystems*, 1(6): 10–1128, 2016.
  - Jacob West-Roberts, Joshua Kravitz, Nishant Jha, Andre Cornman, and Yunha Hwang. Diverse genomic embedding benchmark for functional evaluation across the tree of life. *bioRxiv*, pp. 2024–07, 2024.
  - Maciej Wiatrak, Aaron Weimann, Adam Dinan, Maria Brbić, and A. Floto, R. Sequence-based modelling of bacterial genomes enables accurate antibiotic resistance prediction. *bioRxiv*, pp. 2024.01.03.574022, 2024. doi: 10.1101/2024.01.03.574022.
  - Maciej Wiatrak, Ramon Viñas Torné, Maria Ntemourtsidou, Adam Dinan, David C Abelson, Divya Arora, Maria Brbić, Aaron Weimann, and R Andres Floto. A contextualised protein language model reveals the functional syntax of bacterial evolution. *bioRxiv*, pp. 2025–07, 2025.
  - Jennifer L Wilson, Ethan Steinberg, Rebecca Racz, Russ B Altman, Nigam Shah, and Kevin Grimes. A network paradigm predicts drug synergistic effects using downstream protein–protein interactions. *CPT: Pharmacometrics & Systems Pharmacology*, 11(11):1527–1538, 2022.
  - F Alexander Wolf, Philipp Angerer, and Fabian J Theis. Scanpy: large-scale single-cell gene expression data analysis. *Genome biology*, 19:1–5, 2018.
  - Xihui Xu and Jiandong Jiang. Engineering microbiomes for enhanced bioremediation. *PLoS biology*, 22(12):e3002951, 2024.

Ren Zhang, Hong-Yu Ou, and Chun-Ting Zhang. Deg: a database of essential genes. *Nucleic acids research*, 32(suppl\_1):D271–D272, 2004.

Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2: Efficient foundation model and benchmark for multi-species genome. *arXiv preprint arXiv:2306.15006*, 2023.

# A ADDITIONAL DATASETS & TASKS DETAILS

In this section we outline the dataset details, including the overall and per dataset statistics as well as the preprocessing details. For each dataset we performed additional quality checks which we detail below. All of the genomes across datasets contain genome ID, which can be used to identify the genome source. Finally, we provide metadata on genes, proteins and genomes together with the datasets when available, together with extensive documentation.

Supplementary Table 1: Dataset summary for the six benchmark tasks used in this study.

Task	# Species	# Genomes	# Proteins	# Base pairs
Gene essentiality prediction	37	51	169 k	279 MB
Operon identification	5	5	22 k	25 MB
Protein-protein interaction prediction	6956	10 533	36 M	N/A
Strain clustering	25	6710	14 M	16 GB
Phenotypic traits prediction	15 477	24 462	100 M	111 GB
Antibiotic resistance prediction	38	26 302	105 M	112 GB

#### A.1 GENE ESSENTIALITY PREDICTION

We downloaded the gene essentiality annotations for bacteria across genomes from the Database of Essential Genes (DEG, http://origin.tubic.org/deg) (Zhang et al., 2004). Using the genome RefSeq ID provided in the database, we downloaded the associated genomes in both DNA and protein sequence modalities from the NCBI GenBank (https://www.ncbi.nlm.nih.gov/genbank/). Across 66 genomes from DEG, there were multiple genomes with more than 98% overlap when it comes to annotations. We therefore removed these genomes, as including it could lead to inflated evaluation metrics, leaving us with 51 genomes across 37 distinct species. For each genome we provide start and end for each gene together with essentiality annotations (Yes=essential, No=non-essential), verifying the gene locations are correct. We also provide the strand of the gene to allow for the extraction of the region upstream of the gene.

**Split.** We performed a random data split into training, validation, and test sets in a 60 / 20 / 20 % ratio. Additionally, to prevent train-test leakage, we split by genus—placing all genomes from a genus in one split—and evaluate on held-out genera, enforcing generalization to phylogenetically distant strains. AUROC

#### A.2 OPERON IDENTIFICATION

Due to the lack or reliable whole-genome operon annotations, we performed long-read RNA sequencing to annotate operons across five distinct strains, processing three independent biological replicates for each strain.

**Bacterial strains and culture conditions.** Five bacterial strains were used in this experiment: Staphylococcus aureus RN450 (S. aureus RN450), Mycobacterium abscessus ATCC 19977 (M. ab),  $\Delta$ leuD  $\Delta$ panCD Mycobacterium tuberculosis H37Rv 102J23 ( $\Delta$ leuD  $\Delta$ panCD M. tb), Pseudomonas aeruginosa PAO1, and Escherichia coli DH5 $\alpha$ . Each strain was cultured in triplicate under nutrientrich conditions until mid-exponential phase ( $OD_{600} = 0.4$ –0.6) was reached.

**RNA isolation.** Cells were harvested by centrifugation and processed for total RNA extraction using the MasterPure Complete DNA and RNA Purification Kit (Lucigen) with strain-specific modifications: For *S. aureus* RN450, cell pellets were pre-treated with lysostaphin (Tris buffer, pH 8.0) at 37 °C

for 30 min to aid lysis. For the mycobacterial strains (M. ab and  $\Delta$ leuD  $\Delta$ panCD M. tb), cells were mechanically disrupted by bead beating in lysis buffer, followed by extraction with the standard kit protocol. Isolated RNA was treated twice with TURBO DNase (Invitrogen) to remove residual genomic DNA and purified using RNA Clean & Concentrator columns (Zymo Research). RNA integrity was assessed on an Agilent TapeStation RNA ScreenTape system, and concentrations were measured with a Qubit fluorometer (Invitrogen).

**Library preparation and sequencing.** For each replicate, 1,000 ng of total RNA underwent rRNA depletion with riboPOOLs (siTOOLs Biotech). The depleted RNA was polyadenylated with poly(A) polymerase (PAP) in the presence of a manganese catalyst, adding 50–90 adenosines per molecule. cDNA libraries were prepared with the Nanopore cDNA-PCR kit, pooled and sequenced on a PromethION device equipped with R10 flow cells (Oxford Nanopore Technologies).

**Long-read RNA sequencing data preprocessing.** For each sequenced strain, the following genome assemblies and gene annotations from NCBI RefSeq [ref] were used: GCF\_000013425.1 (*Staphylococcus aureus RN450*), GCF\_000069185.1 (*Mycobacterium abscessus ATCC 19977*), GCF\_000195955.2 ( $\Delta$ leuD  $\Delta$ panCD *Mycobacterium tuberculosis H37Rv 102J23*), GCF\_000006765.1 - (*Pseudomonas aeruginosa PAO1*), and GCF\_002899475.1 (*Escherichia coli DH5* $\alpha$ ).

ONT reads from each replicate were polished with Pychopper (v2.7.10), polyA tails longer than 10 bases and sequencing adapters were trimmed using cutadapt and mapped against the genome assemblies using Minimap2 (v2.29)(Li, 2018) and Samtools (v1.22) (Danecek et al., 2021). Candidate operons were identified from read alignments spanning at least two genes on the same strand and then extended by combining overlapping candidates at most 50 base pairs apart. Operons were then collated from the triplicates for each strain and used as our operon annotations.

**Split.** We evaluate the operon identification in a zero-shot manner, therefore, we do not split the data into train, validation and test splits and use the entire dataset as a test set.

## A.3 PROTEIN-PROTEIN INTERACTION PREDICTION

We downloaded all the data from the STRING DB download site (https://string-db.org/cgi/download). Using the species metadata file we selected only bacterial organisms and downloaded the protein sequences for them together with protein-protein interaction scores for protein pairs. After running the download scripts we ended up with 10,533 unique strains across 6,956 species. We used the *combined* interaction score which combines information from various sources to get a final score. STRING DB provides only protein sequences and no DNA, and the interaction scores are computed mainly at the protein-level, therefore, for this dataset we only provide protein sequences and omit DNA. To binarize the interaction scores, we set the threshold at  $0.6 \geq 0.6$ =interaction, 0.6=no interaction). This threshold was chosen through conducting small-scale experiments and looking at the average performance of AUROC and AUPRC on the validation set across genomes, choosing the threshold which attains the best average performance across the two.

**Split.** We performed a random data split into training, validation, and test sets in a 70 / 10 / 20 % ratio. The larger proportion of the train set compared to the gene essentiality task is motivated by the larger size of the overall dataset, with the 10% validation set still allowing for meaningful evaluation.

# A.4 STRAIN CLUSTERING

To extract the metagenome assembled genomes (MAGs) for strain clustering, we use MGnify (Mitchell et al., 2023), which is a large-scale bacterial genomics database containing a diverse set of MAGs across numerous environments. The main reason for choosing MGnify over other potential resources is its large size combined with a uniform processing pipeline, providing comparable genomes. We wanted to evaluate whether various methods capture phylogenetic similarities across different taxonomic levels, therefore, we looked at strains which span different species, genera and families. These nested ranks provide three increasingly coarse resolutions to test whether an embedding preserves evolutionary signal. We use the taxonomic annotations provided by MGnify. The total number of unique genomes in the dataset is 6,071.

 To extract the genomes of interest, we extracted the most common species from MGnify from the corpus of 300k bacterial genomes and selected 25 species which are distributed across distinct genera and families. For a meaningful evaluation each genus (or family) must contain *at least two* species, otherwise the genus- or family-level clustering metrics become degenerate (every strain would be trivially assigned to a unique cluster at that rank). We download the chosen assemblies and use the accompanying annotations to translate gene DNA to protein sequences, while retaining the original DNA for the DNA-modality experiments.

**Split.** The strain clustering task is a fully unsupervised task, therefore, we do not split the data into train, validation and test splits and use the entire dataset as a test set.

#### A.5 ANTIBIOTIC RESISTANCE PREDICTION

We leveraged the antibiotic sensitivity readings from the NCBI AST browser (https://www.ncbi.nlm.nih.gov/pathogens/ast), which contains hundreds of thousands of antibiotic-susceptibility test (AST) records for a diverse set of antibiotics. Using the genome assembly identifiers provided by the NCBI AST browser, we downloaded the DNA and protein sequences for each genome from the NCBI GenBank (https://www.ncbi.nlm.nih.gov/genbank/) and matched them with the antibiotic resistance readings. This left us with 26,052 unique genomes with matched antibiotic resistance labels. We then processed the antibiotic resistance readings into (i) binary and (ii) regression labels decribed below.

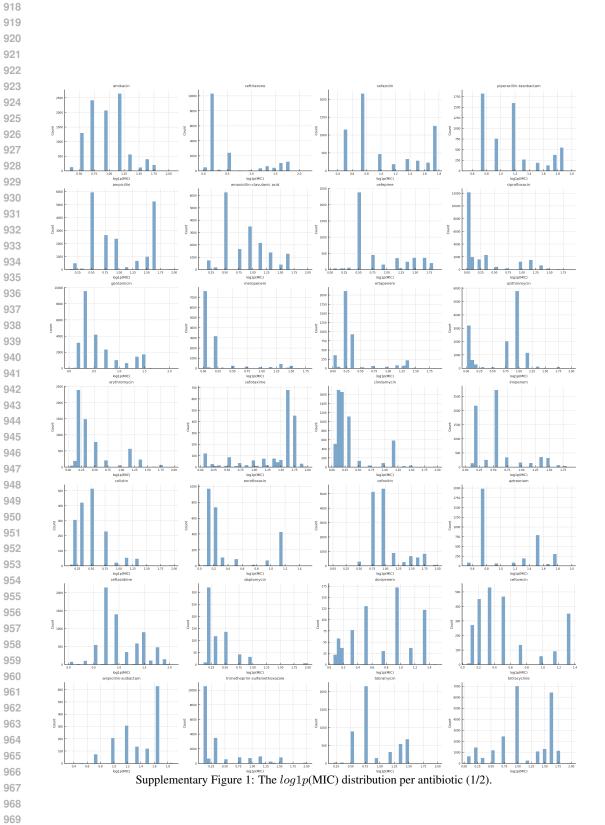
**Binary.** For binary prediction, antibiotic sensitivity labels from NCBI AST browser of either *sensitive* (S) or *resistant* (R) were used for training and testing. If the antibiotic sensitivity test had no S/R label, they were not included. We remove antibiotics which 1) have less than 500 available unique genomes in total, and 2) have less than 50 unique genomes per class. This is motivated by the need to ensure that every classifier is trained on a sufficiently large and reasonably balanced data set; with fewer than 500 genomes overall, or fewer than 50 genomes in either class, the resulting model would suffer from poor statistical power and unreliable performance estimates. This resulted in 37 unique antibiotics. The Supp. Table 2 shows the number of available genomes per drug, including the number of susceptible and resistant genomes. The number of available readings varies strongly between antibiotics, which partly explains the high variance between the per drug performance.

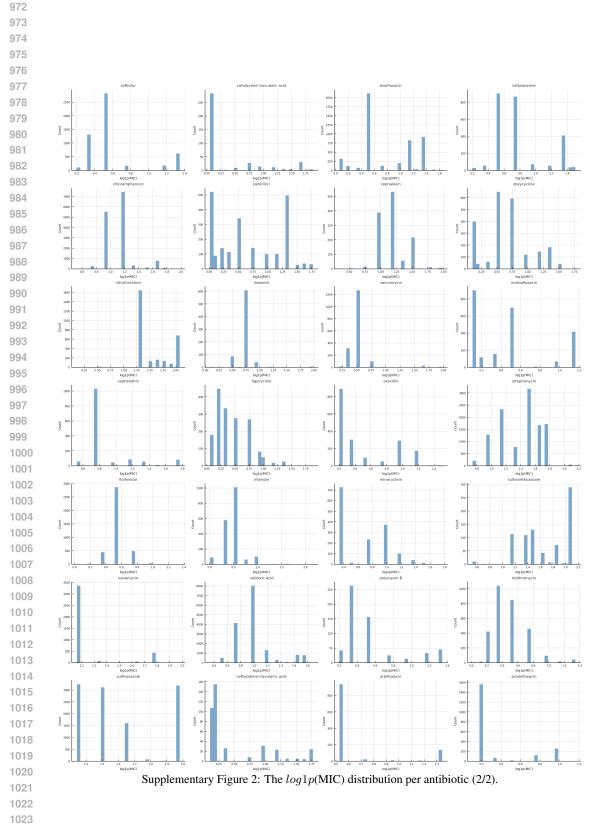
**Regression.** For regression MIC prediction, NCBI AST browser quotes minimum inhibitory concentrations (MICs) as  $< x, \le x, = x, \ge x$  and > x. These strings were translated to an actual number (y) by y = x if NCBI quoted MIC= $x, \le x$  or  $\ge x$ ;  $y = 2 \times x$  if MIC quoted as MIC > x, and  $y = 0.5 \times x$  if MIC was < x. We filter out antibiotics with less than 500 available readings to ensure that the models are trained on a sufficiently large sample. This resulted in 56 antibiotics. To dampen the long tails, we normalized the MIC with a log1p transformation. The final MIC distributions per antibiotic can be found in Supp. Fig 1 & 2.

**Split.** Due to low number of samples for many antibiotics and the variability between genomes, which may skew the results when using a single split, we trained and evaluated all models and antibiotics with k-fold split. Specifically, for each antibiotic we recommend: (1) splitting the available data into 5 equal splits using stratified split for the binary case and random split for regression. (2) In each split, further divide the larger train set into train and val, where validation makes up 20% of the train split. (3) Training the model on the train set and use the best performing model on the validation to evaluate the model on the test set.

Supplementary Table 2: Number of resistant, susceptible, and total labelled genomes for every antibiotic in the binary prediction setting.

Drug	# Resistant	#Susceptible	# Total
amikacin	7353	588	7941
ampicillin	10052	6211	16263
amoxicillin–clavulanic acid	12458	1683	14141
azithromycin	10497	2509	13006
cefazolin	1666	1897	3563
cefepime	3237	1188	4425
cefotaxime	566	446	1012
cefoxitin	12463	3902	16365
ceftazidime	2362	717	3079
ceftriaxone	13385	3469	16854
ciprofloxacin	17353	4280	21633
clindamycin	4327	337	4664
ertapenem	7708	373	8081
erythromycin	5605	1860	7465
fosfomycin	1186	396	1582
gentamicin	18242	2951	21193
imipenem	8485	409	8894
kanamycin	3768	410	4178
levofloxacin	3302	1106	4408
meropenem	11320	811	12131
nalidixic acid	1262	780	2042
nitrofurantoin	10535	5316	15851
oxacillin	6929	762	7691
piperacillin-tazobactam	3294	968	4262
rifampin	611	354	965
streptomycin	11965	2442	14407
sulfamethoxazole	3770	707	4477
sulfisoxazole	1569	323	1892
tetracycline	6936	3242	10178
tigecycline	662	357	1019
trimethoprim	2617	582	3199
ampicillin-sulbactam	1632	451	2083
aztreonam	2530	728	3258
ceftaroline	611	152	763
chloramphenicol	5345	1016	6361
colistin	572	290	862
daptomycin	492	131	623





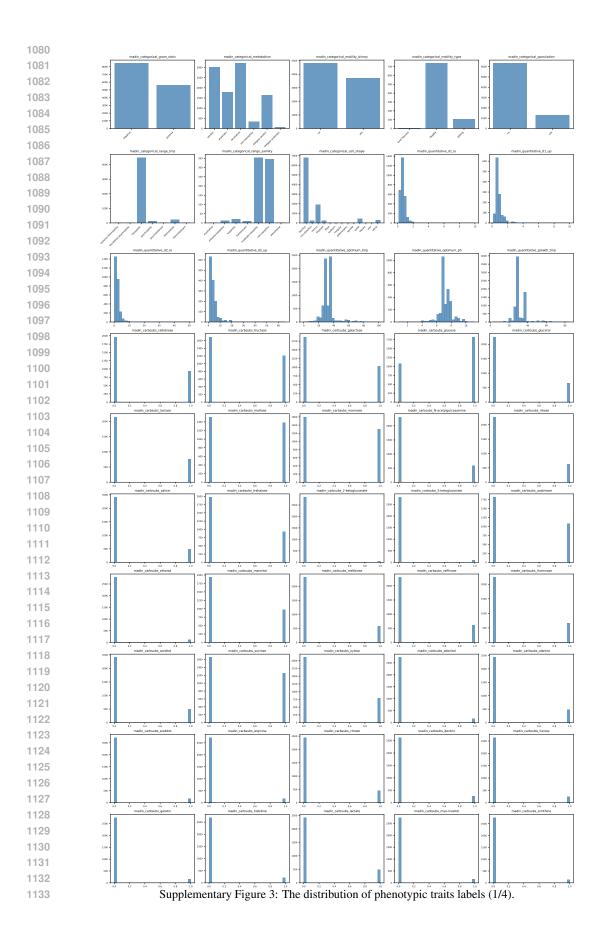
## A.6 PHENOTYPIC TRAITS PREDICTION

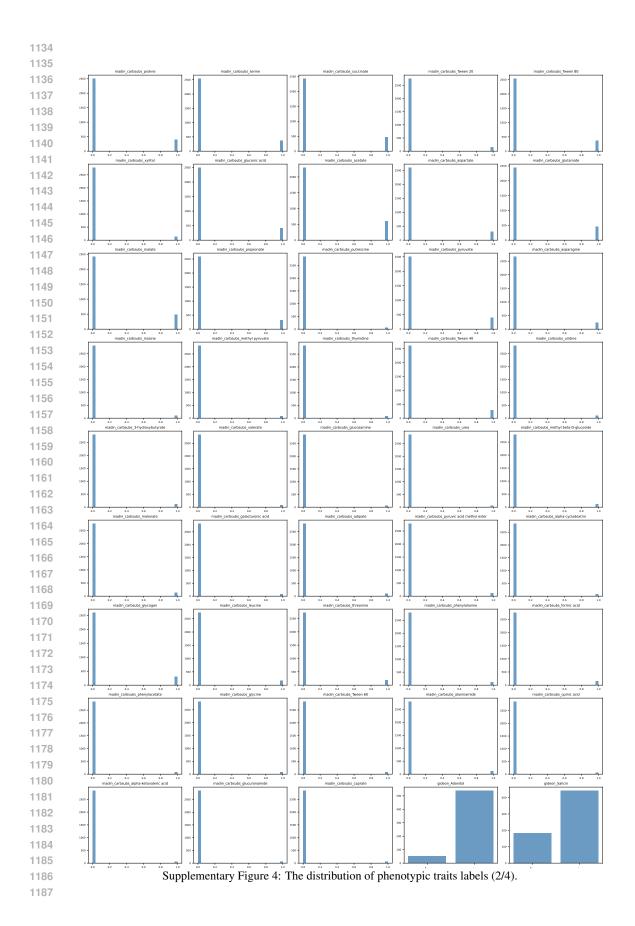
 We collected the phenotypic traits by collating two major sources (Madin et al., 2020; Weimann et al., 2016). To each phenotypic trait label we prepend its data source name so that the provenance of every label is explicit and any potential name collisions between the two catalogues are avoided. We keep duplicate traits that appear in both sources because they expand the number of labelled genomes without forcing us to merge measurements that were obtained with different experimental protocols. Using the taxonomy IDs and assembly accessions provided in the phenotypic traits datasets, we downloaded the associated genome DNA and protein sequences from the NCBI GenBank (https://www.ncbi.nlm.nih.gov/genbank/). We limit ourselves to categorical phenotypes and filter out the phenotypic traits with less than 500 genomes to ensure that the models are trained and evaluated on a sufficiently large sample. Additionally, we removed the classes with less than 50 samples. This resulted in 139 unique phenotypes across 24, 462 genomes. We group the phenotypic traits into 5 distinct groups according to the type of biological information they capture (Supp. Fig. 3). We include the distributions of each phenotypic trait label in the Supp. Fig. 3-5 which shows large variation in the number of available labels per phenotypic trait which affects the final results.

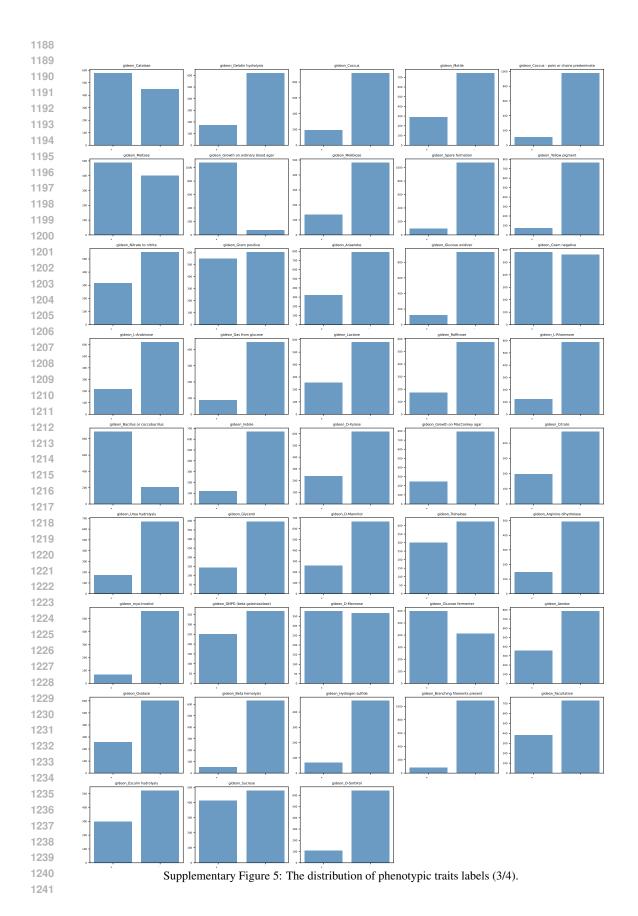
**Split.** For each phenotype, we split the data into 60/20/20 train, validation and test partitions respectively. As similar genomes often share the same phenotype, we split the data for each phenotype by genus—placing all genomes from a genus in one split—and evaluate on held-out genera, enforcing generalization to phylogenetically distant strains. To obtain stable estimates, as many traits are rare, we aggregate the per-phenotype results across 5 independent runs with different data splits.

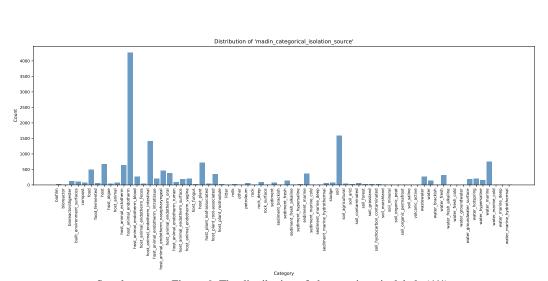
Supplementary Table 3: Phenotype groups and their associated phenotypes. The name before the first "-" symbolizes the dataset source. *madin* phenotypes were extracted from Madin et al. (2020) and *gideon* from Weimann et al. (2016).

Phenotype group	Phenotype
<b>Biochemical Activity</b>	gideon_Gelatin hydrolysis, gideon_Indole, gideon_Urea hydrolysis, gideon_Methyl red, gideon_VP (Voges Proskauer), gideon_Gal (beta-galactosidase), gideon_Beta hemolysis, gideon_Hydrogen sulfide, gideon_Esculin hydrolysis
Carbon Utilization	madin_carbsubs_cellobiose, madin_carbsubs_glucose, madin_carbsubs_glycerol, madin_carbsubs_lactose, madin_carbsubs_maltose, madin_carbsubs_mannitol, madin_carbsubs_sucrose, madin_carbsubs_xylose, gideon_D-Arabitol, gideon_D-Mannose, gideon_Sucrose, gideon_D-Sorbitol
<b>Growth Conditions</b>	madin_categorical_range_tmp, madin_categorical_range_pH, madin_quantitative_optimum_tmp, madin_quantitative_optimum_pH, madin_quantitative_optimum_O2, madin_quantitative_growth_rate, gideon_Growth on ordinary blood agar, gideon_Growth on MacConkey agar
Morphology	madin_categorical_gram_stain, madin_categorical_sporulation, madin_categorical_cell_shape, madin_quantitative_average_cell_size, gideon_Motility, gideon_Spores, gideon_Shape: bacillus or coccobacillus, gideon_Branching filaments present
Respiration Metabolism	madin_categorical_metabolism, gideon_Nitrate reduction, gideon_Alanine aminopeptidase, gideon_O/F glucose oxidizer, gideon_Gas from glucose, gideon_Glucose fermenter, gideon_Aerobe, gideon_Oxidase, gideon_Facultative









Supplementary Figure 6: The distribution of phenotypic traits labels (4/4).

## B ADDITIONAL MODELING & EVALUATION DETAILS

We outline the experimental details used for modeling & evaluation for all model types; DNA Language Models (DNA LMs), protein Language Models (pLMs) and bacterial Language Models (bLMs). Implementation and further details can be found at https://anonymous.4open.science/r/BacBench-B6EF. All of the modeling code has been implemented in PyTorch (Paszke, 2019).

**DNA LMs.** For every model we load the public checkpoint (Supp. Table 4) and keep the hidden states of the *last* encoder layer. If a DNA string has length of G tokens, the encoder produces  $\mathbf{X} \in \mathbb{R}^{G \times D}$ ; we collapse it with a simple mean-pool to obtain a vector

$$\mathbf{z} = \frac{1}{G} \sum_{g=1}^{G} \mathbf{X}_{g,\cdot} \in \mathbb{R}^{D}.$$

**Gene- & system-level.** When the input gene (plus upstream promoter) is longer than the model context C, we slide a window of length C with stride s and embed each window. Averaging the M window vectors  $\{\mathbf{z}^{(1)},\ldots,\mathbf{z}^{(M)}\}$  gives the final gene representation  $\bar{\mathbf{z}}_{\text{gene}} = \frac{1}{M} \sum_{m} \mathbf{z}^{(m)}$ .

**Genome-level.** We treat the whole genome in the same way: split into C-bp windows, embed each, and average; if several contigs are present we first average per-contig and then across contigs.

**pLMs.** A protein of K amino acids yields  $\mathbf{X} \in \mathbb{R}^{K \times D}$  and  $\mathbf{z} = \frac{1}{K} \sum_k \mathbf{X}_k$ . For genome-level tasks we aggregate the M protein vectors in that genome via  $\bar{\mathbf{z}}_{\text{prot}} = \frac{1}{M} \sum_{i=1}^{M} \mathbf{z}_i$ . Gene- and operon-level tasks use only the proteins involved.

**bLMs. gLM2** ingests *mixed-modality* genomic scaffolds in which protein-coding regions are translated to amino-acid tokens and intergenic regions remain as nucleotide tokens. We tokenize each scaffold into a single sequence up to the model context C (4,096 tokens); longer scaffolds are processed with a sliding window of length C and stride s, mirroring the DNA LM setup. For each window, we retain the last-layer hidden states and mean-pool to obtain a window vector  $\mathbf{z}^{(m)} \in \mathbb{R}^D$ ; genome-level embeddings average across windows (and across contigs when present):

$$\bar{\mathbf{z}}_{\text{genome}} = \frac{1}{M} \sum_{m=1}^{M} \mathbf{z}^{(m)}.$$

For gene essentiality and operon tasks, inputs include contextual flanking sequence to supply regulatory/positional cues, but the gene embedding is computed by averaging only the hidden states of tokens that fall within the gene's coordinates; empirically this has shown to outperform averaging the entire slice.

**Bacformer** model takes as input *local* protein vectors  $\mathbf{z}_1, \dots, \mathbf{z}_M$  obtained exactly as in the pLM setting. They are then **ordered by their genomic coordinates** (chromosome followed by plasmids) so that the model can "see" genome organisation. Rotary positional embeddings (Su et al., 2024) are added to these vectors and the ordered sequence is fed to a *genome-level* Transformer encoder (Vaswani et al., 2017) with L layers,

$$\mathbf{H}^{(0)} = [\mathbf{z}_1; \dots; \mathbf{z}_M], \qquad \mathbf{H}^{(\ell+1)} = \text{Transformer}^{(\ell)} (\mathbf{H}^{(\ell)}), \ \ell = 0, \dots, L-1.$$

The output  $\mathbf{H}^{(L)} = [\tilde{\mathbf{z}}_1; \dots; \tilde{\mathbf{z}}_M]$  contains **contextualised protein embeddings**  $\tilde{\mathbf{z}}_i \in \mathbb{R}^D$  that encode both the protein sequence and its genomic neighbourhood. During pre-training Bacformer learns to predict which proteins co-evolved, so these embeddings capture functional coupling across the genome.

Empirically, averaging token embeddings performed slightly better than using the special <code>[CLS]</code> token, so mean pooling was used for every model throughout the paper.

### B.1 GENERAL TRAINING DETAILS.

For all tasks and settings which required finetuning except gene essentiality prediction, we kept the frozen backbone encoder model frozen, and only finetuned the neural network layer(s) stacked on

Supplementary Table 4: Summary of benchmarked models. "Max ctx." = maximum context length supported at inference; for DNA models measured in base pairs, for pLMs in amino acids, and for Bacformer in number of proteins present in the genome. "dim" = dimensionsionality of the output of the last hidden layer. \*We use either bacformer-masked-complete-genomes or bacformer-masked-MAG depending on the genome type used as input. The DNA LMs, pLMs and bLMs are separated by a horizontal line.

Model	Input	Variant / Checkpoint	Objective	Tokenisation	Params	dim	Training corpus	Max ctx.
Mistral-DNA (Mourad, 2025)	DNA	Mistral-DNA-v1-138M-bacteria	Autoregressive	Byte-pair	138M	768	Bacteria	512
DNABERT-2 (Zhou et al., 2023)	DNA	DNABERT-2-117M	Masked	Byte-pair	117M	768	Multi-kingdom	512
Nucleotide Transformer (Dalla-Torre et al., 2024)	DNA	nucleotide-transformer-v2-250m-multi-species	Masked	k-mer	250M	768	Multi-kingdom	2,048
ProkBERT (Ligeti et al., 2024)	DNA	neuralbioinfo/prokbert-mini-long	Masked	k-mer	27M	384	Bacteria	4,096
Evo (Nguyen et al., 2024)	DNA	evo-1-8k-base (1.1_fix)	Autoregressive	Single nucleotide	6.5B	4,096	Multi-kingdom	8,192
ESM-2 (Lin et al., 2022)	Single protein seq.	esm2_t12_35M_UR50D	Masked	Single amino acid	35M	480	Multi-kingdom	1,024
ESM-C (ESM Team, 2024)	Single protein seq.	esmc_300m	Masked	Single amino acid	300M	960	Multi-kingdom	1,024
ProtBERT (Elnaggar et al., 2021)	Single protein seq.	prot_bert	Masked	Single amino acid	420M	1,024	Multi-kingdom	1,024
gLM2 (Cornman et al., 2024)	Mixed modality (DNA & protein seq.)	tattabio/gLM2_650M	Masked	Single nucleotide/amino acid	650M	1,280	Bacteria	4,096
Bacformer (Wiatrak et al., 2025)	Multiple protein seq.	macwiatrak/bacformer-masked-complete-genomes*	Masked	Single protein	27M	480	Bacteria	6,000

top of the model encoder. This was motivated by the computational cost required to embed all 67k genomes with all models. Further details on runtime and computational cost can be found below. We used Adam optimizer (Kingma, 2014) in all finetuning setups. All of the checkpoints used have been downloaded directly from HuggingFace and are specified in Supp. Table 4. For each task and model combination, we tuned the learning rate keeping other parameters unchanged. Further details on hyperparameters used for each task and setup can be found in task-specific sections below.

Supplementary Table 5: Model-specific context parameters used in this study. "Max ctx." is the maximum input length; "DNA-seq overlap" is the stride between consecutive windows when sliding across a genome; "Promoter length" is the upstream sequence length concatenated for promoter prediction.\* For Bacformer the maximum input size of a protein is 1,024 amino acids and maximum number of proteins in the genome is set to 6,000.

Model	Max ctx.	DNA-seq overlap	Promoter length
Mistral-DNA	512	16	128
DNABERT-2	512	16	128
Nucleotide Transformer	2,048	32	128
ProkBERT	4,096	64	128
Evo	8,192	32	128
ProtBERT	1,024	N/A	N/A
ESM-2	1,024	N/A	N/A
ESM-C	1,024	N/A	N/A
gLM2	4,096	64	128
Bacformer	1,024 / 6,000*	N/A	N/A

# B.2 TASK DETAILS

We outline the experimental details for each task, outlining the hyperparameters used and evaluation setup. All of the experiments have been performed on a single NVIDIA A100 with 32 CPU cores.

Gene essentiality prediction. We stacked a single linear layer preceded by a dropout of 0.2 and layer normalization (Ba et al., 2016) on top of the gene embeddings and trained it with binary cross-entropy loss to predict gene essentiality (1=essential, 0=non-essential). We trained the model for maximum of 100 epochs with early stopping patience of 10, monitoring the macro AUROC across genomes on the validation set. For each model we tuned the learning rate specified in Supp. Table 6. The weight decay for the Adam optimizer has been set to 0.02 for all models.

**Evo.** Evo natively does not provide straight-forward access to the output of the last hidden layer of the model. Therefore, we experimented with two ways of extracting the gene embeddings from Evo. Given a gene sequence G of size N, 1) we used the script provided as part of the Evo implementation (https://github.com/evo-design/evo) to score the log-likelihoods of the nucleotides in a sequence, resulting in a vector  $z \in \mathbb{R}^N$ , and 2) modified the Evo model to return the output of the last hidden state, resulting in a matrix  $X \in \mathbb{R}^{N \times D}$ , where D is model dimensionality which here equals 4,096. We then similarly as with other models took the average of all sequence tokens resulting in a vector  $x \in \mathbb{R}^{\mathbb{D}}$ . We include this Evo implementation in the BacBench code repository (https://anonymous.4open.science/r/BacBench-B6EF). The option 1) yielded much better results on the validation set, therefore, we used it for final benchmarking.

Supplementary Table 6: Learning rates used for essential genes linear layer.

Method	Learning rate
Mistral-DNA	0.005
DNABERT-2	0.005
Nucleotide Transformer	0.005
ProkBERT	0.01
Evo	0.001
ProtBERT	0.005
ESM-2	0.005
ESM-C	0.005
gLM2	0.001
Bacformer	0.005

**Operon identification.** The operon–identification task is evaluated *zero-shot*; no fine-tuning is performed. We formulated operon prediction as a binary boundary classification problem, evaluating whether two contiguous genes form an operon. To address this, we developed a method that incorporates three features: (1) gene embeddings and (2) the strand of each gene.

First, we compute cosine similarity between adjacent genes using embeddings from the last hidden layer of pretrained models. We also record each gene's transcriptional strand from the genome assembly.

We then define a simple score that combines similarity and strand co-orientation:

$$s_i = c_i I_{\text{strand}}, \qquad c_i = \frac{1}{2} (1 + \cos(\hat{h}_i, \hat{h}_{i+1})),$$

where  $I_{\text{strand}} = 1$  if both genes are on the same strand and 0 otherwise, and  $\hat{h}_i$ ,  $\hat{h}_{i+1}$  are the  $\ell_2$ -normalised embeddings of genes i and i+1. The score  $s_i \in [0,1]$  serves as the operon-membership score; a pair is classified as belonging to the same operon when  $s_i$  exceeds a threshold. Mapping cosine similarity to [0,1] stabilizes the scale across models, while the strand indicator provides a hard veto since genes on opposite strands do not belong to the same operon.

Performance is computed per strain.

**Evo.** Evo natively does not provide straight-forward access to the output of the last hidden layer of the model. Therefore, we experimented with two ways of extracting the gene embeddings from Evo. Given a gene sequence G of size N, 1) we used the script provided as part of the Evo implementation (https://github.com/evo-design/evo) to score the log-likelihoods of the nucleotides in a sequence, resulting in a vector  $z \in \mathbb{R}^N$ , and 2) modified the Evo model to return the output of the last hidden state, resulting in a matrix  $X \in \mathbb{R}^{N \times D}$ , where D is model dimensionality which here equals 4,096. We then similarly as with other models took the average of all sequence tokens resulting in a vector  $x \in \mathbb{R}^{\mathbb{D}}$ . We include this Evo implementation in the BacBench code repository (https://anonymous.4open.science/r/BacBench-B6EF). The option 2) yielded significantly better results on operon identification, therefore, we used it for final benchmarking.

**Protein-protein interaction prediction.** To predict whether the two proteins interact, we fed the two protein embeddings into a linear model, which is trained to predict a binary label (1=interaction, 0=no interaction). The linear model is a single-layer neural network preceded by a dropout of 0.2 and layer normalization (Ba et al., 2016). The protein embeddings are fed into the linear classifier, after which the pairs of interacting proteins are averaged and passed through a final binary classification layer preceded by a dropout of 0.2. The model is trained to minimize the binary cross-entropy loss. We experimented with different learning rates for all the models and set the final learning rate to 0.001, which has shown to perform the best for all the models. We trained the model for the maximum epochs of 10 and no early stopping patience. We set the maximum gradient norm to 2.0 and monitor the validation loss across genomes. The weight decay for the Adam optimizer is set to 0.01.

**Strain clustering.** To compute the strain-clustering metrics we run Leiden clustering (Traag et al., 2019) over a grid of parameters. We vary the *resolution* in [0.1, 0.25, 0.5, 1.0] and the

number of neighbours in [5, 10, 15], evaluating every pairwise combination. Lower resolutions are omitted to avoid collapsing many genomes into a single (or just a few) giant clusters. After computing the clustering metrics for every parameter pair, we keep for each method the combination that maximises the mean performance across species-, genus- and family-level labels. The Leiden clustering is performed using the scanpy package (Wolf et al., 2018).

Antibiotic resistance prediction. To predict the antibiotic resistance, we train a linear model for each drug and method combination. The linear model is a single-layer neural network preceded by a dropout of 0.2 and layer normalization (Ba et al., 2016). We train the models separately for the (i) binary and (ii) regression MIC prediction case. We optimize the former for the binary cross entropy loss and the latter for the mean squared error loss. We train all models for the maximum of 100 epochs with early stopping patience of 5, monitoring the validation AUPRC in the binary setup and validation  $R^2$  in the regression setup. We experimented with various learning rates for each model, setting the final learning rate to 0.005 which attained the best results across folds and seeds for all models. We set the weight decay in the Adam optimizer to 0.01.

We have also experimented with training a multi-task linear model, which simultaneously predicts antibiotic resistance of a genome to multiple drugs, however, it performed worse then a separate linear model for each drug.

**Phenotypic traits prediction.** To predict a phenotypic trait from a genome-level embedding, we train an linear model for each phenotype and method combination. The linear model is a single-layer neural network preceded by a dropout of 0.2 and layer normalization (Ba et al., 2016). As all labels are categorical, we optimize it to minimize the cross-entropy loss. We set the maximum number of epochs to 2,000 and early stopping patience to 50, monitoring the validation loss. The learning rate for all models was set to 0.01. We use the cross-entropy loss. To account for the class imbalance, we weigh each class according to:

$$w_c = \frac{N}{K n_c},\tag{1}$$

where  $n_c$  is the number of training samples in class c, K is the total number of classes, and N is the total number of training samples.

**Runtime analysis.** A typical bacterial genome contains on the order of 3,000-5,000 genes and  $\sim 4-6$  Mbp of DNA, with average protein lengths of  $\sim 300$  amino acids; embedding an entire genome therefore entails thousands of forward passes for protein LMs and millions of tokens for DNA LMs, making raw throughput a practical bottleneck for population-scale studies. We measured the wall-clock time required to embed the 11 genomes used in the operon identification task on a single NVIDIA A100 and extrapolated to BacBench's full collection of 67,000 genomes (Table 7). The fastest models are **ESM-2** and **Bacformer** (64–67 s for 11 genomes;  $\approx 1.2-1.25$ k GPU-hours for 67k genomes). DNA LMs add a modest cost (e.g., 93 s for Mistral-DNA; 168 s for DNABERT-2) yet remain tractable on a single GPU. In stark contrast, the 6.5B-parameter **Evo** is  $\sim 10^4 \times$  slower (5.76×10<sup>5</sup> s for 11 genomes, i.e.,  $\sim 14$  h per genome), yielding an impractical  $\sim 1.07 \times 10^7$  GPU-hours for 67k genomes on one GPU. These measurements underline that—even when accuracy is the primary goal—runtime quickly becomes the limiting factor for population-scale analyses.

Beyond embedding costs, fine-tuning on genome-level tasks requires backpropagating through *all* proteins per genome (median  $\sim$ 2,506 proteins across our datasets) or through multi-megabase DNA contexts, which dramatically amplifies memory and compute, requiring often ¿200 NVIDIA A100 GPUs for a single genome. Practically, this makes genome-scale fine-tuning out of reach for most academic labs for DNA LMs and pLMs, while remaining feasible for Bacformer-style bLMs; thus, linear-probe evaluations are a necessary, controlled proxy for model selection at scale.

#### B.3 Datasets & models licenses

To ensure that our datasets and benchmarks are reusable in the academic setting, we checked the license for each resource used in the manuscript. The Supp. Table 8 details a license for all models and datasets used in the manuscript.

Supplementary Table 7: Embedding runtime (s=seconds) on the operon identification task (11 genomes) on one NVIDIA A100 GPU and the extrapolated wall-clock time (h=hours) to process the entire collection of 67k bacterial genomes on the same hardware.

Model	Runtime (s)	Estimated time for 67 k genomes (h)
Mistral-DNA	93	1,731
DNABERT-2	168	3,127
Nucleotide Transformer	100	1,861
ProkBERT	106	1,973
Evo	575,629	10,712,540
ProtBERT	95	1,768
ESM-2	64	1,191
ESM-C	137	2,550
gLM2	182	3,387
Bacformer	67	1,247

Supplementary Table 8: External resources used in this study and their licences.

Resource	Type	Licence
ESM-2	Model	MIT
ESM-C	Model	Cambrian Open License Agreement
ProtBERT	Model	Academic Free License v. 3.0
gLM2	Model	Apache 2.0
Bacformer	Model	Apache 2.0
DNABERT-2	Model	Apache 2.0
ProkBERT	Model	MIT
Evo	Model	Apache 2.0
Nucleotide Transformer	Model	CC BY-NC-SA 4.0
Mistral-DNA	Model	Apache 2.0
MGnify	Data	CC0 1.0 Universal
NCBI AST Browser	Data	Public domain (U.S. Gov data)
NCBI GenBank/RefSeq	Data	Public domain (U.S. Gov data)
Database of Essential Genes	Data	CC BY 4.0
Operon DB	Data	CC BY 4.0
STRING DB	Data	CC BY 4.0

## C ADDITIONAL RESULTS

We show and discuss further results across all tasks. To allow for comparability in future work, we include tables with all numerical results as well as per antibiotic and phenotypic traits scores.

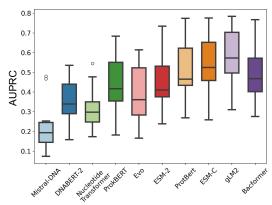
## C.1 GENE ESSENTIALITY PREDICTION

The results on AUPRC show the bLMs and pLMs outperforming DNA LMs (Supp. Fig 7). gLM2 performs the best, showing the benefits of taking as input both DNA as well as protein sequence information. Moreover, increasing the pLM size appears to further boost performance, as demonstrated by ESM-C (300M) and ProtBERT (420M) outperforming ESM-2 (35M). Bacformer outperforms its protein representation backbone ESM-2, showing the benefits of incorporating whole-genome context. Evo outperforms other DNA LMs, except ProkBERT, demonstrating the performance gain by conducting pretraining on a relevant corpus. Finally the results on AUPRC show that there is large room for improvement. We believe that increasing the number and diversity of annotated genomes would significantly boost model performance.

The Supp. Table 9 shows the exact results across AUROC and AUPRC measured across disinct genomes.

**Finetuning performance.** To further analyse model performance, we have conducted finetuning on the gene essentiality task. The results show that finetuning boosts performance (Table 10); however, the model ranking remains largely unchanged, with the pLMs and bLMs outperforming DNA LMs. We have excluded Evo from finetuning due to the computational complexity required to finetune

 Evo (Table 7). Gene essentiality is the only task for which we finetuned all models; genome-level tasks such as antibiotic-resistance prediction require end-to-end, whole-genome context and are prohibitively expensive for all models except Bacformer—underscoring the need for models that can be finetuned efficiently for whole-genome tasks.



Supplementary Figure 7: AUPRC across test genomes on the gene–essentiality prediction task. The box spans the inter-quartile range with a line marking the median value.

Supplementary Table 9: Overall performance on gene–essentiality prediction. Values are  $mean \pm standard\ deviation$  over 3 random seeds; the best score for each metric is highlighted in bold.

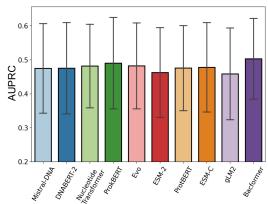
Method	AUROC	AUPRC
Mistral-DNA	$57.52 \pm 2.86$	$22.69 \pm 14.21$
DNABERT-2	$68.21 \pm 7.17$	$35.76 \pm 12.46$
Nucleotide Transformer	$67.03 \pm 5.56$	$31.78 \pm 11.88$
ProkBERT	$74.79 \pm 6.22$	$44.79 \pm 15.25$
Evo	$73.71 \pm 7.10$	$39.08 \pm 15.21$
ESM-2	$77.99 \pm 5.82$	$46.23 \pm 15.87$
ESM-C	$82.25 \pm 6.04$	$55.08 \pm 15.83$
ProtBERT	$80.72 \pm 5.85$	$52.00 \pm 15.90$
gLM2	$\textbf{83.77} \pm \textbf{6.17}$	$58.61 \pm 14.96$
Bacformer	$80.72 \pm 5.87$	$50.33 \pm 15.43$

Supplementary Table 10: Overall finetuning performance on gene–essentiality prediction. Values are  $mean \pm standard\ deviation$  over 3 random seeds; the best score for each metric is highlighted in bold.

Method	AUROC	AUPRC
Mistral-DNA	$62.18 \pm 7.87$	$28.81 \pm 13.51$
DNABERT-2	$74.88 \pm 16.75$	$56.12 \pm 29.71$
Nucleotide Transformer	$73.06 \pm 14.99$	$48.56 \pm 28.65$
ProkBERT	$69.88 \pm 8.60$	$44.98 \pm 14.36$
ESM-2	$85.36 \pm 9.12$	$64.42 \pm 17.79$
ESM-C	$89.96 \pm 7.53$	$71.85 \pm 14.49$
ProtBERT	$\textbf{91.31} \pm \textbf{73.02}$	$73.02 \pm 16.56$
gLM2	$90.12 \pm 7.87$	$\textbf{74.03} \pm \textbf{21.70}$
Bacformer	$90.31 \pm 7.40$	$72.13 \pm 15.16$

#### C.2 OPERON IDENTIFICATION

The AUPRC follows the same trend as the AUROC described in the main manuscript: Bacformer attains the best scores, ProkBERT ranks second, and the remaining DNA LMs and pLMs form a tight cluster with broadly comparable performance—reinforcing that genome-level context and bacteria-specific pretraining matter more than modality alone; absolute AUPRC values are lower (as expected under class imbalance) but track the same model ordering.



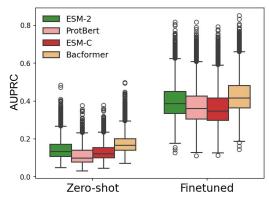
Supplementary Figure 8: AUPRC across test genomes on gene essentiality prediction task. The box spans the inter-quartile range with a line marking the median value.

Supplementary Table 11: Overall performance on operon identification. Values are  $mean \pm standard$  deviation; the best score for each metric is highlighted in bold.

Method	AUROC	AUPRC
Mistral-DNA	$74.16 \pm 9.70$	$47.48 \pm 29.47$
DNABERT-2	$74.13 \pm 10.14$	$47.53 \pm 30.07$
Nucleotide Transformer	$74.50 \pm 8.57$	$48.19 \pm 27.47$
ProkBERT	$75.77 \pm 9.69$	$49.03 \pm 30.07$
Evo	$74.77 \pm 8.33$	$48.21 \pm 28.27$
ESM-2	$73.02 \pm 10.30$	$46.27 \pm 29.55$
ESM-C	$75.25 \pm 9.26$	$47.79 \pm 29.39$
ProtBERT	$75.11 \pm 8.65$	$47.58 \pm 27.98$
gLM2	$72.85 \pm 10.77$	$45.86 \pm 30.18$
Bacformer	$\textbf{77.59} \pm \textbf{6.64}$	$\textbf{50.31} \pm \textbf{26.61}$

#### C.3 PROTEIN-PROTEIN INTERACTION PREDICTION

The protein-protein interaction (PPI) prediction results on AUPRC (Supp. Fig. 9) show that contextual pLM, Bacformer, consistently outperforms other methods. We credit it to its usage of the genomic context. Moreover, the performance does not increase by scaling the model size, as shown by ESM-C and ProtBERT underperforming ESM-2. The overall results (Supp. Table 12) show relatively low performance considering the trainins set size ( $\[ in ] \]$ 7k genomes) even in the finetuned setting. We believe this is due to the 1) complexity of the task, 2) noisy source data, as STRING DB where the interactions have been extracted from collates information from a variety of sources and is not limited to experimentally validated interactions, highlighting the importance of building high quality PPI datasets.



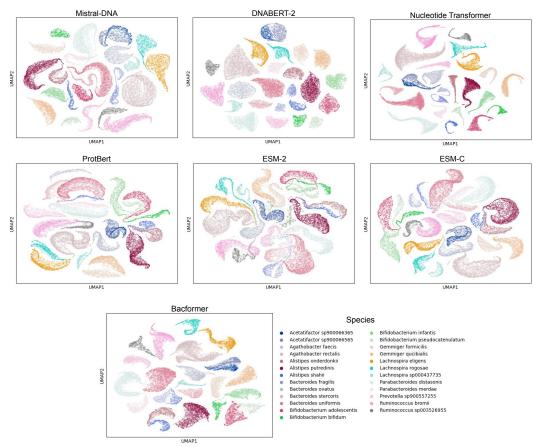
Supplementary Figure 9: AUPRC across test genomes on the protein–protein-interaction task in both zero-shot and fine-tuned settings. The box spans the inter-quartile range with a line marking the median.

Supplementary Table 12: Protein–protein-interaction prediction on the held-out genomes. Values are  $mean \pm standard\ deviation$  over five seeds; the best score for each metric is shown in bold.

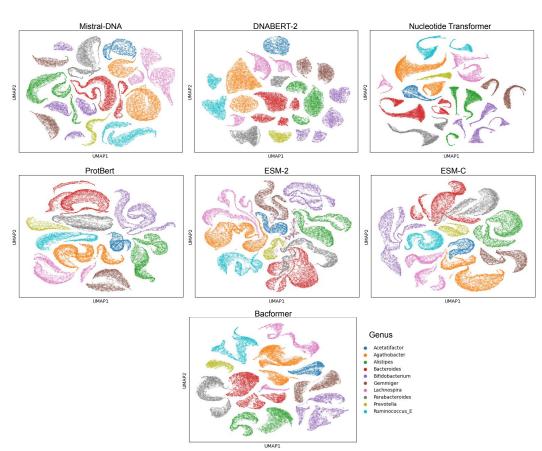
Method	Zero	-shot	Fine	tuned
	AUROC	AUPRC	AUROC	AUPRC
ESM-2	$56.46 \pm 2.10$	$14.94 \pm 6.04$	$76.62 \pm 2.35$	$40.68 \pm 10.71$
ProtBERT	$47.20 \pm 3.88$	$11.45 \pm 4.89$	$74.05 \pm 2.53$	$38.15 \pm 11.16$
ESM-C	$55.17 \pm 2.66$	$13.33 \pm 4.61$	$74.21 \pm 2.37$	$37.30 \pm 10.98$
Bacformer	$\textbf{63.09} \pm \textbf{2.73}$	$\textbf{18.20} \pm \textbf{6.28}$	$\textbf{79.09} \pm \textbf{2.25}$	$\textbf{43.47} \pm \textbf{10.61}$

#### C.4 STRAIN CLUSTERING

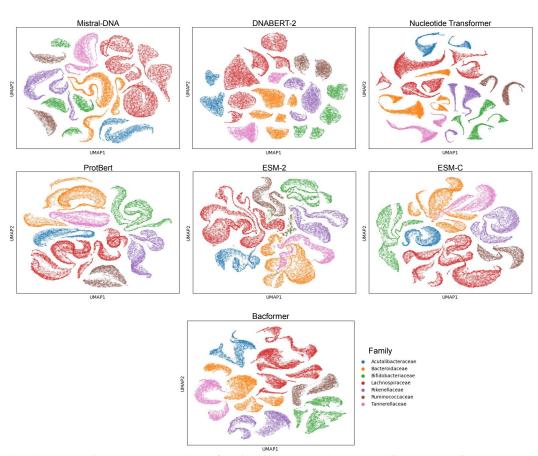
In addition to the strain clustering metrics included in the main manuscript, we plotted 2-dimensional UMAP results colored by species (Supp. Fig. 10), genus (Supp. Fig. 11) and family (Supp. Fig. 12) for a subset of models for further investigation. The UMAPs show that the strain representations differ between the models. All models cluster strains into separate species clusters, however, not all of them retain the phylogenetic similarities between species in the same genus or family. We also notice that the DNA LMs tend to output less overlapping species clusters, which boosts its performance at species level, but leads to lower results at higher taxonomic levels (genus and family).



Supplementary Figure 10: UMAP plots of strains (i.e. genome-level embeddings) across diverse methods colored by species.



Supplementary Figure 11: UMAP plots of strains (i.e. genome-level embeddings) across diverse colored by genus.



Supplementary Figure 12: UMAP plots of strains (i.e. genome-level embeddings) across diverse methods colored by family.

#### C.5 ANTIBIOTIC RESISTANCE PREDICTION

On the following pages we present the results per antibiotic across metrics. This includes binary prediction results (susceptible/resistant, Supp. Table 13) and MIC regression prediction results (Supp. Table 14). Each antibiotic was run across 5-folds with 3 random seeds to avoid variance stemming from random initialization and data-split bias. In the binary prediction setting the bLM-Bacformer outperforms other methods, achieving the best AUROC on 26 drugs, AUPRC on 27 drugs and F1 on 24 drugs out of 37 in total. Thus, showcasing the benefits of considering the interactions between proteins present in the genome. In the regression setting, Bacformer attains the best result on 44 drugs on  $R^2$  and 45 on *Pearson* correlation coefficient out of total of 56 antibiotics. Finally, we see large variance across as well as within drugs. The former can be partly explained by the variable number of labels available per antibiotic, while the latter shows that, even within a single antibiotic, model performance can fluctuate markedly across folds and random seeds—highlighting sensitivity to sample composition and pointing to the need for larger, more balanced datasets or stronger regularisation to obtain stabler estimates.

Supplementary Table 13: Per-antibiotic binary antibiotic resistance prediction. Values are mean  $\pm$  standard deviation across 3 seeds. Bold indicates the highest mean

3.51 0.70 0.70 4.63 0.74 2.71 12.81 1.26 1.26 8.95 10.96 1.52 3.68 0.62 5.93 7.76 1.12 4.82 3.95 3.95 1.99 6.66 6.66 1.70 1.29 0.42 11.09 11.80 11.80 11.80 11.73 for each metric.

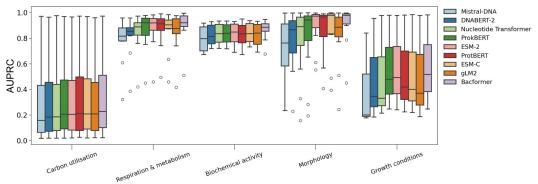
Supplementary Table 14: Per-antibiotic MIC regression prediction results across 3 seeds. Values are mean  $\pm$  std. Bold indicates the best mean for each metric.

Antibiotic				· ·							******							
	Mistral-DNA Pearson	-DNA R <sup>2</sup>	DNABEKI-2 Pearson	EKI-2 R <sup>2</sup>	Nucleonde 1 Pearson	son R <sup>2</sup>	Pearson	EKI R <sup>2</sup>	FSM-2 Pearson	F2 R2	FSM-C Pearson	ي س	Pearson	DEKI R <sup>2</sup>	Pearson	gLM2	Bactormer Pearson	mer R <sup>2</sup>
amikacin	67.35 ± 0.91	45.08 ± 1.48	68.37 ± 1.03	46.63 ± 1.44	70.01 ± 1.02	48.82 ± 1.47	69.78 ± 0.74	4 2	69.95 ± 1.08	48.63 ± 1.71	71.13 ± 0.99	50.50 ± 1.47	69.36 ± 0.85	48.01 ± 1.21	69.42 ± 1.67	48.00 ± 2.42	71.32 ± 1.59	50.46 ± 2.08
amoxiciilin-ciavutanic acid	46.27 ± 1.43 46.43 ± 0.78	20.05 ± 1.03	48.39 ± 2.33	25.24 ± 2.06 25.35 + 1.03	52.77 ± 0.51		51.78 + 0.63					31.94 ± 0.41	59.47 ± 0.95	29.15 ± 1.12	51.14 ± 1.11	25.24 ± 1.01	57.06 ± 1.09	32.09 ± 0.70 31.79 + 1.36
ampicillin-sulbactam	40.77 ± 3.28	16.57 ± 2.74	41.23 ± 3.29	16.77 ± 2.37	1 +			_				21.33 ± 2.96	$42.40 \pm 1.66$	17.46 ± 1.20	47.21 ± 3.75	22.04 ± 3.45	47.78 ± 7.02	22.09 + 7.02
azithromycin	$79.15 \pm 0.38$	$62.59 \pm 0.66$	$79.57 \pm 0.06$	$63.18 \pm 0.10$	+	_		$63.42 \pm 0.09$		_	$80.43 \pm 0.58$	$64.50 \pm 1.02$	$79.97 \pm 0.14$	$63.82 \pm 0.24$	$79.73 \pm 0.14$	$63.40 \pm 0.28$	$80.27 \pm 0.33$	$64.33 \pm 0.70$
aztreonam	$72.00 \pm 1.08$	$50.51 \pm 2.22$	$72.36 \pm 1.57$	$51.22 \pm 2.83$	+	_		_		_	$78.07 \pm 2.12$	$60.23 \pm 3.90$	$77.07 \pm 2.08$	$58.61 \pm 3.02$	$74.41 \pm 0.82$	54.77 ± 1.17	$\textbf{79.38} \pm \textbf{1.96}$	$62.61 \pm 3.12$
cefazolin	$65.50 \pm 0.57$	$41.05 \pm 1.29$	$66.18 \pm 1.10$	$43.27 \pm 1.77$	+	$47.84 \pm 1.98$		$46.19 \pm 2.16$		_	$75.37 \pm 1.85$	$56.05 \pm 2.89$	$67.57 \pm 1.06$	$45.26 \pm 1.74$	$70.19 \pm 0.19$	$48.77 \pm 0.96$	$73.18 \pm 0.83$	$53.21 \pm 1.30$
cefepime	48.40 ± 2.99	23.05 ± 2.62	49.95 ± 2.10	24.28 ± 1.43	53.21 ± 2.92	28.21 ± 3.25		28.14 ± 1.85		33.82 ± 3.44	59.39 ± 1.85	34.51 ± 2.59	56.44 ± 3.72	31.56 ± 4.15	52.34 ± 0.59	27.11 ± 0.72	$60.53 \pm 1.95$	$36.34 \pm 2.49$
cerotaxime	15.28 ± 5.95	1.04 ± 0.84	35.60 ± 12.14	10.01 ± 6.23	Н-	25.09 ± 9.05	49.80 ± 7.62	_	57.55 ± 1.90	31.6/ ± 2.40	28.41 ± 2.50	33.13 ± 3.77	54.12 ± 5.05	27.74 ± 5.29	46.18 ± 4.42	20.12 ± 5.28	63.09 # 1.48	38.04 ± 3.43
cefovecin	9.42 ± 6.07	20 82 ± 1 66	60.37 ± 0.74	36.04 ± 0.86	41.05 ± 2.51 63.07 ± 0.88	30 67 ± 1.15	55.65 ± 19.51	38 31 + 1 18	49.22 ± 0.72	30 14 + 1 13	51.34 ± 1.83	42 60 ± 1.05	59.08 ± 0.29	30.75 ± 0.22	51.62 ± 9.08	37.35 ± 0.50	52.00 ± 1.88	43.21 ± 1.46
cefoorime	32 75 + 4 16	10.27 + 2.00	32 72 + 3 35	10.00 + 1.46	4 +	15 11 + 2 73	42 37 + 0 25	_	47.28 + 0.40		48 70 + 2 08	22 66 + 1 90	45 07 + 1 70	18 94 + 0 33	30.86 + 5.52	8 71 + 3 37	52.81 + 1.99	26 55 + 1 40
ceffazidime	44.76 ± 0.70	19.70 ± 0.71	46.81 ± 0.74	21.63 ± 0.69	1 +	29.66 ± 2.24	54.68 ± 1.50	_	59.62 ± 1.22		60.87 ± 0.62	36.49 ± 0.84	50.81 ± 2.20	25.48 ± 1.95	49.58 ± 1.50	24.17 ± 1.82	61.88 ± 1.62	36.95 ± 2.55
ceftiofur	$4.56 \pm 3.02$	$-0.02 \pm 0.30$	$38.41 \pm 1.93$	$9.77 \pm 0.87$		_	$38.19 \pm 2.72$	12.09 ± 2.73	$47.37 \pm 0.13$	_	$48.21 \pm 0.04$	$21.82 \pm 0.64$	$41.84 \pm 1.81$	$16.04 \pm 2.22$	$32.45 \pm 3.89$	$9.84 \pm 2.07$	$45.78 \pm 2.60$	$19.01 \pm 2.91$
ceftriaxone	$56.54 \pm 0.78$	$31.78 \pm 1.01$	$57.93 \pm 1.22$	$33.37 \pm 1.43$	$63.17 \pm 1.83$				$64.81 \pm 1.96$	_	$69.00\pm1.59$	$\textbf{47.11} \pm \textbf{2.05}$	$63.53 \pm 1.58$	$39.97 \pm 1.81$	$62.63 \pm 1.41$	$38.71 \pm 1.61$	$67.66 \pm 1.71$	$45.43 \pm 2.16$
cephalexin	$6.48 \pm 2.58$	$0.05 \pm 0.19$	$7.01 \pm 4.62$	$0.24 \pm 0.54$	$19.53 \pm 6.34$	$3.33 \pm 1.91$		$2.25 \pm 1.92$	$20.39 \pm 6.07$	_	$35.04 \pm 5.97$	$9.38 \pm 2.38$	$18.65 \pm 1.38$	$2.94 \pm 0.97$	_	$2.38 \pm 3.24$	$35.36 \pm 12.54$	$10.53 \pm 7.83$
cephalothin	58.16 ± 8.99	30.88 ± 10.42	58.18 ± 8.67	30.84 ± 9.15		_		_	61.98 ± 8.07	_	60.87 ± 8.15	35.29 ± 10.44	$63.70 \pm 10.05$	$36.95 \pm 14.57$		35.71 ± 10.05	$66.14 \pm 10.27$	$41.47 \pm 13.32$
chloramphenicol	11.59 ± 0.68	1.27 ± 0.12	16.99 ± 1.50	2.82 ± 0.44	22.13 ± 2.63	4.62 ± 0.89	18.82 ± 2.22	_	29.32 ± 2.20	7.60 ± 0.94	26.42 ± 5.41	6.74 ± 2.38	20.58 ± 4.21	4.07 ± 1.53		3.64 ± 1.38	$34.47 \pm 1.39$	$10.99 \pm 0.56$
cipronoxacin	55.09 ± 0.62	30.40±2.62	57.20 ± 1.51	36.30 ± 4.38	58.82 ± 0.82	39.12 ± 1.30		33.72 ± 0.76	59.10 ± 0.89	_	59.15 ± 1.11	34.00 ± 1.00	28.33 ± 1.48	34.02 ± 1.33	20:44 ± 0:48	30.75 ± 3.01	70.27 ± 0.07	37.71 ± 0.45
cundaniyem	22.28 + 6.30	4 73 + 2 59	26.95 + 3.60	7.03 + 1.72	4 +				32.42 + 3.01			17.47 + 4.24	45.48 + 5.02	19.00 + 4.77		8 86 + 1 53	50.27 + 7.21	22 54 + 5 30
dantomycin	39.30 + 20.02	13 52 + 13 63	40.87 + 18.10	13.90 + 11.44					43.66 + 15.77			11 59 + 30 18	43 99 + 15 22			26 53 + 5 14	53 00 + 0 03	28.48 + 9.28
doripenem	44.81 ± 7.03	18.72 ± 6.94	44.48 ± 7.11	$19.01 \pm 6.14$	50.00 ± 4.35			20.23 ± 2.64		22.04 ± 14.10	$60.01 \pm 6.35$		47.19 ± 7.18		48.11 ± 5.74	22.84 ± 5.62	51.85 ± 3.34	25.74 ± 2.71
doxycycline	$48.62 \pm 2.95$	$23.17 \pm 3.47$	$50.81 \pm 2.43$	$25.53 \pm 2.70$	$52.66 \pm 1.16$	_		_					$53.76 \pm 1.33$		_	$26.02 \pm 2.57$	$58.22 \pm 1.21$	$33.35 \pm 1.46$
enrofloxacin	$23.72 \pm 2.70$	$5.16 \pm 0.99$	$26.65 \pm 1.42$	$6.82 \pm 0.52$	$47.42 \pm 2.82$	_	$51.24 \pm 2.65$	_	$58.55 \pm 1.40$			_	$47.04 \pm 3.52$		_	$14.29 \pm 1.70$	$60.94 \pm 3.94$	$35.63 \pm 4.48$
ertapenem	64.83 ± 1.13	40.73 ± 2.02	68.0 ± 0.99	42.60 ± 1.09	67.35 ± 0.75	_		_				_	$69.96 \pm 1.36$		_	39.23 ± 2.09	$71.99 \pm 1.53$	51.63 ± 2.18
erythromycm	42.05 ± 1.24	16.50 ± 0.80	50.18 ± 1.28	23.97 ± 1.66	53.06 ± 0.29	26.90 ± 0.36	54.88 ± 2.60	28.60 ± 1.17					54.07 ± 2.06			25.13 ± 0.98	59.05 ± 1.36	$33.53 \pm 0.57$
norrement	42.43 ± 1.80 55.70 ± 0.61	30.47 ± 0.75	58 11 + 0 70	33.46 ± 0.94	42.77 ± 8.03 50 10 ± 1.03		43.80 ± 9.39 58 66 ± 0.67					35 36 ± 1 56	58.61 + 0.60	34 24 + 0 68	_	34 30 ± 1 36	45.50 ± 9.11	36.01 ± 0.02
iminenem	70.39 + 1.64	49.28 + 2.42	70.87 + 1.50	49.89 + 2.37	71.28 + 1.70		71.27 + 1.81					50.89 + 1.85	72.22 + 1.30	51.62 + 1.96	_	50.42 + 1.38	73.39 + 1.88	53.62 + 2.66
kanamycin	$0.75 \pm 5.27$	$-0.19 \pm 0.24$	$23.16 \pm 2.69$	$4.62 \pm 0.85$	$36.53 \pm 2.17$	_	$36.88 \pm 2.80$	_				$11.61 \pm 2.80$	$33.68 \pm 2.35$	$10.99 \pm 1.38$		$8.57 \pm 1.94$	$36.26 \pm 3.37$	$12.59 \pm 2.59$
levofloxacin	$28.22 \pm 10.03$	$7.61 \pm 4.71$	$42.69 \pm 1.33$	$17.54 \pm 0.91$	$48.47 \pm 0.58$	_		_			$58.67 \pm 2.39$	$32.76 \pm 2.13$	$49.48 \pm 2.47$	$23.69 \pm 2.09$	_	24.93 ± 1.49	$58.59 \pm 1.18$	$33.35 \pm 1.64$
linezolid	$8.05 \pm 15.49$	$\textbf{0.22} \pm \textbf{2.85}$	$7.73 \pm 14.70$	$-5.96 \pm 10.47$	$7.37 \pm 1.10$	_	$16.07 \pm 17.49$	$-7.48 \pm 20.22$		$-5.13 \pm 12.04$	5.68 ± 10.41	$-14.01 \pm 26.52$	$3.85 \pm 25.12$	$-21.28 \pm 35.29$	_	$-15.08 \pm 25.90$	$2.86 \pm 9.77$	$-34.65 \pm 30.70$
marbofloxacin	41.18 ± 6.12	15.70 ± 6.36	42.54 ± 4.39	16.75 ± 5.53	50.89 ± 2.76	_		23.46 ± 11.50		35.53 ± 5.81		42.74 ± 4.01	53.88 ± 4.92	28.06 ± 5.22		21.83 ± 9.38	62.21 ± 5.45	37.09 ± 8.55
meropenem	73.23 ± 2.30	25.57 ± 2.49	50 32 ± 0.95	33.72 ± 0.88	/6.01 ± 1.35	36.28 + 2.54	61 22 + 3.01	37.25 ± 2.04	69.82 ± 0.98	20.73 ± 1.20	70.19 ± 1.05	48 27 + 2 01	70.12 ± 1.23	27.38 ± 1.73	60.95 ± 0.80	36.28 ± 1.07	/6.95 ± 1.48	28./8 ± 1.98 43.54 ± 5.42
nalidixic acid	42.25 ± 2.47	17.63 ± 1.76	43.20 ± 2.09	18.45 ± 1.44	45.64 ± 2.06	20.62 ± 1.96	43.62 ± 2.13	19.01 ± 1.85				24.44 ± 1.54	$45.39 \pm 1.63$	$20.08 \pm 1.74$	44.65 ± 2.80	19.76 ± 2.41	$54.10 \pm 1.14$	$28.32 \pm 1.20$
nitrofurantoin	$91.48 \pm 1.91$	$62.92 \pm 2.32$	$94.36 \pm 0.93$	$88.58 \pm 1.68$	+	$89.60 \pm 2.22$	$95.05 \pm 1.36$	$90.30 \pm 2.51$	$94.93 \pm 0.94$	_	$94.87 \pm 0.81$	$89.97 \pm 1.50$	$95.10 \pm 1.00$	$90.42 \pm 1.88$	$94.76 \pm 1.08$	$89.78 \pm 2.02$	$95.18 \pm 1.05$	$90.55 \pm 1.99$
orbifloxacin	6.61 ± 4.56	-4.38 ± 4.73	10.49 ± 5.31	-2.52 ± 3.35	ш.	-5.91 ± 6.96	3.74 ± 1.34	-5.61 ± 5.73	14.22 ± 8.99	-4.05 ± 4.23	24.15 ± 9.09	-2.85 ± 13.33	8.81 ± 7.84	$-5.34 \pm 5.05$	17.65 ± 9.42	-0.58 ± 3.66	$26.94 \pm 13.87$	$4.06 \pm 3.81$
oxacıllın	39.98 ± 5.97	15.98 ± 6.21	24.68 ± 5.45 43.53 ± 2.53	18 16 + 3 47	26.74 ± 2.49	30.98 ± 2.18	65.44 ± 0.42	38.44 ± 0.91	64.06 ± 1.88	40.62 ± 2.06	66.14 ± 3.02 55.16 ± 5.06	45.00 ± 5.14	62.72 ± 5.66 52.00 ± 2.74	38.87 ± 4.20	28.80 ± 4.05	33.19 ± 3.74	70.36 ± 2.44	48.92 ± 2.99
piperacillin-tazobactam	56.97 ± 0.92	32.01 ± 1.01	58.82 ± 0.85	34.18 ± 0.79	1 +	44.81 ± 2.91	66.31 ± 1.20	_	68.88 ± 2.18	46.45 ± 2.72	71.07 ± 1.08	49.57 ± 1.36	67.21 ± 0.99	44.58 ± 0.79	64.52 ± 0.58	41.01 ± 0.57	$72.30 \pm 2.04$	51.81 ± 2.56
polymyxin B	$23.37 \pm 1.83$	$0.06 \pm 8.91$	$29.67 \pm 9.17$	$6.46 \pm 8.12$	$43.30 \pm 5.78$	$18.56 \pm 4.72$	$46.37 \pm 3.09$	_	$46.17 \pm 6.72$	_	$44.44 \pm 11.11$	$19.31 \pm 10.42$	$\textbf{46.82} \pm \textbf{5.68}$	$21.02 \pm 6.42$	$31.96 \pm 13.30$	$8.46 \pm 9.67$	$41.83 \pm 12.62$	$17.56 \pm 9.45$
pradofloxacin	$13.71 \pm 4.66$	$1.44 \pm 0.97$	$14.28 \pm 0.85$	$1.47 \pm 0.25$	+1	$7.07 \pm 3.19$	$35.05 \pm 11.51$	_	$48.88 \pm 4.54$	_	$57.13 \pm 3.18$	$30.88 \pm 3.04$	$35.33 \pm 2.14$	$12.10 \pm 1.33$	$33.54 \pm 11.82$	8.98 ± 4.66	$51.13 \pm 3.26$	$25.49 \pm 2.82$
rifampin	65.70 ± 10.55	38.90 ± 12.54	65.43 ± 11.25	42.08 ± 13.97	64.59 ± 7.60	39.63 ± 7.41	67.21 ± 8.07	44.45 ± 9.87	68.89 ± 8.92	46.59 ± 11.82	66.77 ± 7.60	43.27 ± 10.46	62.28 ± 5.46	$37.65 \pm 6.19$	67.05 ± 8.26	44.57 ± 11.06	68.45 ± 6.02	46.60 ± 8.24
streptomycin sulfamethoxazole	34.32 ± 2.22	-0.74 + 1.05	41.54 ± 2.45	-2 67 + 2 93	46.45 ± 0.99	21.34 ± 0.77	20.07 ± 2.00	3.08 + 2.02	28.34 + 9.00		24.80 ± 2.49	5.01 + 2.30	10 19 + 2 35	10.90 ± 3.93	24 86 + 5 28	4 20 ± 1.57	57.69 ± 2.42 23.63 ± 3.14	34.43 ± 4.43 255 + 283
sulfisoxazole	19.83 ± 0.36	3.56 ± 0.64	45.59 ± 2.22	19.43 ± 1.51	55.93 ± 3.70	30.50 ± 4.43	53.53 ± 1.14	27.34 ± 1.24	62.66 ± 0.37	_	65.77 ± 0.46	$42.18 \pm 0.92$	57.25 ± 1.24	31.15 ± 1.75	46.49 ± 1.57	20.29 ± 1.30	60.71 ± 1.40	$35.84 \pm 1.63$
telithromycin	$41.16 \pm 2.74$	$12.00 \pm 8.06$	$43.11 \pm 2.49$	$16.37 \pm 2.49$	$42.63 \pm 3.04$	$17.63 \pm 2.78$	$41.81 \pm 2.89$	$16.78 \pm 3.28$	$45.30 \pm 3.50$	_	$45.50 \pm 2.78$	$19.87 \pm 2.55$	$42.61 \pm 3.40$	$17.77 \pm 3.03$	_	$11.00 \pm 1.42$	$45.80 \pm 0.55$	$20.27 \pm 0.62$
tetracycline	$25.72 \pm 0.94$	$6.50 \pm 0.38$	$33.08 \pm 1.60$	$10.85 \pm 1.00$	$36.88 \pm 2.22$	$13.39 \pm 1.55$	$39.37 \pm 1.38$	$15.15 \pm 0.91$	$44.68 \pm 1.24$	_	$49.46 \pm 1.65$	$23.74 \pm 1.49$	$40.02 \pm 7.56$	$16.03 \pm 5.47$	_	20.04 ± 1.46	$52.35 \pm 0.84$	$26.20 \pm 1.24$
tigecycline	37.97 ± 3.94	13.42 ± 2.80	48.96 ± 4.84	23.52 ± 4.10	51.95 ± 4.02	26.61 ± 4.52	51.92 ± 3.85	26.81 ± 3.90	55.00 ± 3.40	29.33 ± 3.28	54.76 ± 4.72	29.39 ± 4.57	53.01 ± 4.70	27.83 ± 5.41	52.17 ± 2.08	26.74 ± 2.32	59.21 ± 3.78	34.44 ± 4.14
trimethonrim-sulfamethoxazole	58 87 + 0 77	34 10 ± 0 77	52.80 ± 2.02 61 77 + 1 48	37.84 + 1.71	63 94 + 1 22		53.13 H 4.06 64.00 + 1.32	40.80 ± 4.33			68.82 + 0.78	46.89 ± 0.98	57.10 ± 4.72 65.03 + 1.25	41.87 ± 1.53		39.83 + 0.44	68.07 ± 0.69	35.10 ± 2.45 45 93 ± 0 96
vancomycin	$14.32 \pm 5.27$	$-0.57 \pm 4.11$	$30.27 \pm 7.49$	$8.39 \pm 4.89$	$49.51 \pm 4.69$	$23.33 \pm 5.10$	$43.46 \pm 10.95$	$17.68 \pm 10.86$		$13.18 \pm 6.24$	$\textbf{51.69} \pm \textbf{2.31}$	$\textbf{25.36} \pm \textbf{2.71}$	$40.21\pm10.64$	$12.70 \pm 10.03$	$33.83 \pm 9.94$	$10.08 \pm 8.51$	$45.52 \pm 8.99$	$18.05 \pm 10.09$

#### C.6 PHENOTYPIC TRAITS PREDICTION

 We measure the AUPRC across phenotype groups (Supp. Fig. 13) as well as overall performance across metrics and all phenotypic traits (Supp. Table 15), and groups (Supp. Table 16). Finally, we include the results for each individual phenotype (Supp. Table 17 & 18). Analyzing performance across phenotype groups and metrics (Supp. Fig. 13 & Supp. Table 16), bLM-Bacformer achieves the best or combined best result across all groups and metrics, showing the benefits of incorporating genomic interactions into phenotype prediction models. Across metrics and phenotypes, Bacformer achieves the best result on 80 phenotypes on AUROC, 78 on AUPRC and 86 on F1 out of a total of 139 phenotypes. Therefore, showing that there is a considerable variation between phenotypes and one should choose a model specific for a phenotype. We believe this is due to the variable number of labels available for a phenotype as well as the inherent differences in the phenotypes themselves, which make some phenotypes easier to predict using a computational approach than others.

Predicting phenotypes is often a very challenging task which includes understanding the effect of mutations and multi-level genomic interactions. However, accurately predicting phenotypes could allow us to engineer genomes for a desired purpose, such as sustainable bioproduction, thus having potentially massive positive impact. We believe that next-generation models should consider the genomic context and incorporate the prior knowledge, such as genome-scale metabolic models to make the most out of available data for a given phenotype.



Supplementary Figure 13: AUPRC across diverse phenotypic traits groups and methods. The box spans the inter-quartile range with a line marking the median value. The results were macro-averaged across classes for each phenotype.

Supplementary Table 15: Overall phenotypic traits prediction performance across all phenotypes. Values are mean  $\pm$  standard deviation across 5 seeds. The results were macro-averaged across classes for each phenotype.

Method	AUROC	AUPRC	<b>F</b> 1
Mistral-DNA	$60.34 \pm 9.81$	$39.72 \pm 33.09$	$47.25 \pm 6.67$
DNABERT-2	$63.83 \pm 11.06$	$42.10 \pm 34.24$	$49.42 \pm 8.99$
Nucleotide Transformer	$65.98 \pm 11.45$	$43.22 \pm 34.85$	$52.04 \pm 11.07$
ProkBERT	$67.64 \pm 11.87$	$44.80 \pm 34.72$	$52.03 \pm 11.56$
ESM-2	$68.71 \pm 12.36$	$45.61 \pm 35.45$	$53.10 \pm 11.90$
ESM-C	$69.07 \pm 11.98$	$45.54 \pm 35.38$	$52.28 \pm 11.39$
ProtBERT	$68.47 \pm 11.97$	$45.47 \pm 34.89$	$53.71 \pm 11.65$
gLM2	$65.82 \pm 11.41$	$44.00 \pm 34.40$	$51.11 \pm 10.86$
Bacformer	$\textbf{71.34} \pm \textbf{12.91}$	$\textbf{47.80} \pm \textbf{35.99}$	$56.14 \pm 13.19$

Supplementary Table 16: Phenotype-group prediction results;  $mean \pm standard\ deviation$  over five random seeds. The highest score in each column is shown in bold. Results are macro-averaged across classes for each phenotype.

Method	В	iochemical activi	ty	-	arbon utilisation	n	(	Frowth condition	s		Morphology		Resp	iration & metabo	olism
	AUROC	AUPRC	F1	AUROC	AUPRC	F1	AUROC	AUPRC	F1	AUROC	AUPRC	F1	AUROC	AUPRC	F1
Mistral-DNA	54.67 ± 12.73	$79.60 \pm 10.06$	$44.48 \pm 2.40$	58.17 ± 7.25	$27.95 \pm 27.92$	46.67 ± 3.25	63.73 ± 4.34	$40.58 \pm 37.70$	$36.80 \pm 17.74$	$70.07 \pm 11.96$	$72.05 \pm 23.63$	$52.81 \pm 9.92$	$71.79 \pm 11.54$	$77.30 \pm 18.47$	50.88 ± 14.54
DNABERT-2	$59.88 \pm 9.16$	$82.39 \pm 8.55$	$46.19 \pm 6.05$	$60.46 \pm 7.60$	$29.24 \pm 28.29$	$47.63 \pm 3.90$	$71.22 \pm 17.07$	$49.41 \pm 40.69$	$45.89 \pm 27.95$	$79.15 \pm 12.39$	$79.75 \pm 22.57$	$56.83 \pm 14.23$	$76.88 \pm 11.06$	$80.97 \pm 16.72$	$60.52 \pm 15.97$
Nucleotide Transformer	$63.74 \pm 8.59$	$83.40 \pm 8.47$	$50.62 \pm 4.88$	$62.07 \pm 7.58$	$30.06 \pm 28.72$	$49.20 \pm 4.03$	$79.33 \pm 13.37$	$50.82 \pm 41.12$	$51.73 \pm 33.88$	$79.63 \pm 14.10$	$78.94 \pm 26.27$	$61.55 \pm 20.66$	$81.74 \pm 10.26$	$83.90 \pm 16.08$	$66.91 \pm 14.99$
ProkBERT	$61.84 \pm 8.89$	$83.41 \pm 8.29$	$49.52 \pm 5.28$	$63.50 \pm 7.72$	$31.50 \pm 28.50$	$49.34 \pm 4.49$	$86.10 \pm 8.33$	$56.82 \pm 37.42$	$49.07 \pm 29.41$	$82.76 \pm 12.95$	$81.29 \pm 26.40$	$61.87 \pm 23.45$	$84.91 \pm 8.16$	$85.55 \pm 16.14$	$66.74 \pm 14.85$
ESM-2	$62.86 \pm 12.75$	$83.30 \pm 9.29$	$51.30 \pm 9.54$	$64.60 \pm 8.18$	$31.77 \pm 29.34$	$49.38 \pm 4.58$	$81.70 \pm 11.20$	$56.97 \pm 37.48$	$52.70 \pm 29.63$	$87.97 \pm 10.47$	$89.36 \pm 14.86$	$71.41 \pm 18.24$	$83.46 \pm 11.48$	$86.63 \pm 15.62$	$66.11 \pm 16.31$
ESM-C	$62.71 \pm 11.90$	$83.90 \pm 8.16$	$50.82 \pm 9.24$	$65.37 \pm 8.16$	$32.20 \pm 29.54$	$49.27 \pm 4.24$	$78.71 \pm 15.27$	$53.44 \pm 40.05$	$55.72 \pm 35.81$	$88.05 \pm 9.67$	$87.26 \pm 20.11$	$67.68 \pm 20.43$	$81.79 \pm 12.13$	$85.05 \pm 16.95$	$61.05 \pm 14.90$
ProtBERT	$61.93 \pm 13.25$	$82.30 \pm 10.18$	$50.78 \pm 7.31$	$64.86 \pm 8.29$	$32.22 \pm 29.14$	$50.30 \pm 5.05$	$82.07 \pm 11.45$	$54.06 \pm 39.32$	$54.87 \pm 33.97$	$83.27 \pm 13.36$	$86.10 \pm 18.70$	$68.60 \pm 18.70$	$84.71 \pm 8.40$	$86.48 \pm 14.80$	$68.69 \pm 14.17$
gLM2	$60.39 \pm 8.39$	$82.71 \pm 9.45$	$47.76 \pm 4.82$	$62.78 \pm 7.22$	$30.97 \pm 28.27$	$48.65 \pm 4.06$	$71.55 \pm 26.57$	$51.10 \pm 41.36$	$45.49 \pm 27.39$	$77.44 \pm 16.90$	$82.79 \pm 21.88$	$62.18 \pm 20.77$	$82.26 \pm 9.24$	$84.22 \pm 16.68$	$64.43 \pm 16.91$
Bacformer	$67.42 \pm 16.45$	$86.32 \pm 9.03$	$55.09 \pm 11.26$	$66.94 \pm 8.94$	$33.94 \pm 30.24$	$51.99 \pm 6.44$	$86.18 \pm 9.85$	$58.25 \pm 37.22$	$59.53 \pm 33.08$	$89.53 \pm 10.01$	$91.13 \pm 15.07$	$70.20 \pm 19.26$	$88.18 \pm 9.31$	$89.19 \pm 14.12$	$77.15 \pm 15.49$

Supplementary Table 17: Per-phenotype performance (1/2). Values are mean  $\pm$  standard deviation across 5 random seeds. Bold highlights the best mean per metric. The results were macro-averaged across classes for each phenotype.

adipe	AUROC	Mistral-DNA AUPRC	Œ	AUROC	DNABERT-2 AUPRC	FI	Nucleo AUROC	ideotide Transformer AUPRC	E N	ProkBERT AUROC AUPRC	ORT C FI	AUROC	ESM-2 AUPRC	FI	AUROC	ESM-C AUPRC	E	AUROC AU	ProdBERT AUPRC F1	AUROC	gLM2 OC AUPRC	2 IC FI	AUROC	Bacformer	F1
Adonitol	78.65±8.07		47.84 ± 1.49		96.42 ± 4.56	47.77 ± 1.40		97.56 ± 2.88 51	-	8.40	51.76 ±	177.70 ±	∓ 1996	47.84 ± 1.49	78.03 ± 9.47	14	I-	± 6.84	97.24 ± 3.01 47.77 :	1.40   77.22 ±	95.94	47.84	± 29'08	97.22 ±	£1519
Acrobe				85.67±3.39 80			88.49±1.52 9		74 ± 4.43 92.43	+ 246 94.96 + 1.72	82.49 ±	4.69 93.05 ± 1.45	95.67 ± 0.59	87.14 ± 2.79	90.71 ± 2.42	9424±1.75 65.	2 8	28	11.14 72.43 1	5 5	92±258 94.70±	±197 79.99±	±4.48 96.51±0	0.82 98.23 ± 0.31	92.05 ± 1.21
dihydrolase	35.07 ± 9.42	69.15 ± 4.22	42.87 ± 1.67		77.45 ± 7.76	42.87 ± 1.67			-	19.51	42.16 ±	39.06	72.64	42.87 ± 1.67	42.30 ± 17.13		29		67.28 ± 5.88 42.87	1.67 52.26 ±	78.49	42.87	34.86	£2979	1 66.14
obacillas							_		1.9	E 11.47	51.87 ±	75.66±	55.22 ±	$^{\rm H}$			±7.36	4.65	± 26.37 62.78 ±	8	24.28 ±	48.21	78.62 ±	7	$^{\rm H}$
	$63.40 \pm 15.20$		47.08±0.64	62.26 ± 16.78 9.					47.08 ± 0.64 53.06	± 5.76	47.08 ±	\$4.80 H	# 19%8	#	54.08 ± 6.86	90.87±3.83 47			± 5.24	8	89.60	47.08	\$2.05 ±	24 88.81 ± 5.	5 47.08±0.6
ng filaments present									98	1.63	\$3.06 ±	98.23±	99.30 ±	#		+031 56	~	17.41	+355	8.8	99.78 ±	48.37	+ 69'96	99.75 ±	48.37±
n,Catalase		77.79 ± 13.31	_			-			1 4 30	+6.13	84.23 H	94.35 H	± 2906	87.54 ± 3.48	90.54 ± 5.98	9.12	8	2.40	₩.	61.6	5.12 83.68±		94.73 H	<u> </u>	90.22±
		76.67 ± 9.44	_			_	~ .	GS1±823 56.97	10.62	+ 12.62	H :	7241 ±	84.18±	+1 -	72.19 ± 10.49	83	± 14.02	± 10.85	11.73	14.8	8.46 85.35 ±	6.48 53.78 ± 14.29	98.99	62 82.79 ± 12	41.
	68.59±4.91	9121 ± 2.80	_		93.78 ± 2.55			+ 241	17.13	± 10.39	62.52	91.16 ±	# 05 TS	75.01 ± 17.61	+1 -	8 5	18 18 10 11 11 11 11 11 11 11 11 11 11 11 11	H 15.74 98	F6259 #51#	18.71 S3.07 ±	41.4	8.74 45.85	93.94	2.10 98.66 ± 0.	81.98 ± 4.3
its or chains predominate		98.98±0.85	_			48.77 ± 1.18		99.26±0.52 51	200	-		±4.80 86.92 ± 16.60	99.57±0	67.14 ± 18.28		7:	77±1.18	23 ± 6.30 97.60	±2.71 48.35	25	38.25	159 48.77	1.18 90.94±	28 99.71±0.	52.42 ± 3.2
		8139±1.66	_		86.62 ± 3.52 4		75.88±4.69 8		9.71		49.10	77.25 ±	~	56.58 ± 15.49	H	+2.48 49.	± 16.45	± 5.57	± 2.87 63.15	8	± 96'88 80'90	62.9	77.72 #	,	55.76 ± 14.0
		48.12 ± 6.55	_					± 7.16	# .	12.28	46.49	_		+	# .	113.77 56	30±331 66.58	12.40	±835 5858:	3 1	1519	5.73	4.28 73.86±	6.47 76.60 ± 12	9 63.37 ± 12.2
	03.14 ± 12.32	92.08±3.01	_		00/ ± 9000	-		# ! H	7	± 16.30 89.39	4900	H # 70	H 57 TA	н.	н.	107	_	18.81	1000	8	18.72 88.40±	4000	H	#01.00 10.00	
	$62.12 \pm 2.18$	$83.82 \pm 1.21$	_			-		333	± 0.83	±6.16 89.51 ::	67.75	98.99	87.35 ±	н	H	Š	98	± 5.16 s	± 2.82	88.38	100 8593	91.49	± 6994	89.85	66.33±
probysis and a second	59.34 ± 12.27		_			_	-	-	H	±3.62 71.78 :-	47.22 ±	± 29.65	69.005 ± 1	55.48 ± 16.58	H	75.15 ± 17.42 60.51	± 19.25	± 13.64 SM	±8.55 53.54±	17.86 63.50±	73.48 ±	488	76.008 ±	18.25 78.69 ± 18.4	_
	71.07 ± 12.49		43.12 ± 2.72		83.53 ± 14.34 6	_	86.07±7.98	± 6.43	90 9	157	++ -	24.39 83.51 ± 8.39	93.07 ± 4.47	67.21 ± 23.75	77.57 ± 13.56	±5.19 62	± 22.18	86	±605 81.17±	10.33 73.03 ±	99.99	12.70 52.26 ±	35 89.95±	4.08 93.09±4.	84.52 ± 4.7
	45.55 ± 12.51			53.39 ± 14.32 &		_	~ .	1330	± 9.12	± 6.90 89.21	\$1.10 ±	56.26 ±	88.14 ±	Ξ.	#	171 48	#183	+ 17.75	± 642 5933	1.92 64.60±10.16	88.42 ±	'n.	± 66.693 ±	# 96°68	•
	59.28 ± 7.85	77.91 ± 5.26	_			_		297	570	±13.82 78.84±	80.24	± 90'09	78.78 ±	41.80 +	H	7	+153	± 11.62	600	2.86 50.75±	15.99±	11.95 42.95	1.31 72.97 ±	S.20 87.75 ± 3.8	_
be .	67.50±3.09	00.96 ± 6.11	_			_	80.33 ± 10.15 72	±16.46	1137	± 10.00 75.05 ±	0830		77.21 ±	+0989		8	_	659	12.56	8.41 80.43 ± 7.18	4.02±	_	7.79 89.98 ±	90 86.31±7.2	83.65 ± 4.7.
oxidizer		95.84 ±3.32	_		95.79±3.33	46.11 ± 2.76		52	20.27	±3.81 97.83	54.38	# 19'S8	97.25±	н	#	±2.86 49	+3.85	± 13.91		_	₹97.65	46.06	0000 ⊕	08.10 ±	_
						_	~	±15.36	± 15.29	± 13.95 8	43.72 ±	71.71 ±	87.31 ±	50.81 ± 5.52	H	88.99±9.87 61.	.74 ± 15.38   66.91	± 16.02	± 9.57 46.73:	_	\$5.01±		74.23 ±	# 88.93 ±	\$6.39±9.20
	67.08 ± 13.42		50.41 ± 12.62		85.84 ± 4.62 6	_	89.11±8.97 8	+ 9.71	81.49 ± 8.53 97.09	+ 246	193 86.66±9	# SS 03 +	+ 61.86	Ξ.	98.70 ± 1.90	1 + 2.39	_	+ 0.42	+ 0.58	2.96 88.52±	12.16		_	99.16 ±	97.75±
		_				_		± 6.03	+	-	93.17 ±	-		93.26 ± 6.92	#1	Z 70+	388	+239	8:	88	4.29 96.69±	92.29	3.52 97.61 ±	98.27 ±	76.39 ±
	62.83 ± 6.12	84.09±6.04	_		95.45±2.59	-		1.48	+ 3.73	13.77 97.88 ±	74.04 ±	H 121 H	97.73 ±	44	Π.	8	5	· .	1.52	3	10076	0.88 70.89±	10.57 95.31 ±	H 15 86 31	82.14 ±
ury blood agar	68.45 ± 12.66	19.93 ± 24.85	_			-		12925	± 9.92	± 16.24 47.97 ±	\$6.51 ±	76.84 ±	49.20 ±	# .	66.33 ± 24.10	135.50 64	~	± 32.78 4	±3831	19.56 78.07±9	74 36.90±	'n.	75.73 ±		
n_Hydrogen salfide	34.06±9.54	8132±420	_	± 16.33		-		+ 5.73	2.96±5.56 59.32	15.68 90.54	\$000		93.01 ±	51.20 ± 8.92	Η.	92.92±3.09 46	_	± 6.94	+141	8	42 91.17 ±	2.81 47.72	75.28 +	94.75±	+ 68'89
n_Indole		87.96±0.99	_					00.93±3.07 50	43 ± 9.13 66.92	+ 17.26	4.93 45.39 ±0	Ξ.	9031±	46.74 ±		5.14	46±0.60 62.03	± 2009	10.60	8	16±13.40 90.61±	282 46.54	+	90.38 ±	46.42 ±
n.LAmbin ose		84.21 ± 11.77	-			_		+1	74 ± 9.35 79.7	±4.59 88.26±	8.30 65.82±6	95 79.00±4.77	87.32 ± 1	56.32 ± 12		±10.66 59	±13.48	500 H	+6.55	11.38 80.32 ±	2.27 89.52±	530 62.72±	82.80 ±	•	
mose	69.02 ± 7.46	91.83±2.26	_	_ `		-	-	+0.42 6	Ξ.	± 10.33 91.08 ±	3.66 46.24 ±2	20 76.16 ± 9.27	H 28 H	+1	79.36 ± 1.75	95.02 ± 1.47 44	_	+ 5.73	±1.70 55.96	7.14 68.98±	17.61 92.67±	1.59 59.57±	12.06 79.61 ± 1	1.42 93.93±4.	59.57±
	49.82 ± 11.14		_			_		388		±7.74 82.04 ±	±9.18 68.67 ±7	+1	82.13 ±	62.63 ± 6.81	74.99 ± 5.59		_	1.24.89		5		7.47 52.75 ±	20.94 82.05±	+ 0076	74.48
	57.91±1.00		47.22±6.44	_		\$4.22 ± 14.94		9.12	-	20.26 ±	\$6.15 ±	74.67±	72.94 ± 7.65	+1	+1	69.61±1.98 48.	_	1.89 ± 7.93 71.22 =	ъ.	61.95	523 5	9.43 37.46	+1	3.76 82.49 ± 8.20	-
36	53.02 ± 10.99		40.56±4.24			_		+ 3.54	43.008 ± 3.55 49.00	± 6.05 75.76 ±	44.91 H	23 52.09 ± 16.81	79.72 ±	$40.81 \pm 3.83$		± 2.66	±3.15	18.6 =	4.75	47.86	7.5	3.27 50.47±	6.12 64.98 ± 1	.08 86.77±5.	н.
	66.98±9.54	77.79±6.46	43.13 ± 5.84			_		± 2.55	56.00 ± 14.77 64.5	+522 7720+		+1	88.24 H	70.20±11.65	#1	22	+23.70	9.50±6.52 81.25	± 5.43 70.29 ±	11.87 63.98±	-	4.27 47.11 =	_	71 91.10 ± 7.	_
	65.98 ± 7.21	78.20±6.12	_		84.73 ± 4.76	_	77.91±7.68 8	9			58.99 H	75.60±	85.66±	53.70 ± 20.78	80.16 ± 11.98	90 ±0.0±	_	5 ± 9.63	603	10.70 75.32 ±	2	_	81.54 H	89.02±	70.42 ±
beta galactosidase)	59.67 ± 8.24	67.14 ± 8.62	_			_		w.		1 8.38 75.81 ±	57.43 ±	81.60±	81.15±	Ξ.	76.65 ± 9.60	6	_	12.20	18.47	12.46	50	e 1	# SeSI ±	5.65 88.21±2.3	# \$6'09
	S7. H 47.47	87.80 ± 7.50	_	_	87.28 ± 10.78	-			0.5	D 56 47.74	H :	~	# 07.46 # 1		85.55 ± 15.50	8	± 10.95	44 -	56.0	8 7 7		3.38 78.07±	#64.78	H 05'56	81.04 ± 7.89
9	28.43 ± 15.93	87.67±2.53	-			_		80.25 ± 4.49 30	1004 ± 7.58 74.47 ± 000	1471 9034	H 0000	н-		57.54 ± 10.13	75.86 ± 11.52	2 -	00 T T T T T T T T T T T T T T T T T T	8852 ± 14.54 89.59	±3290 39393±	12.11 70.24 ± 15.00	1500 9009±	542 5285	# ST 90	978 9483 759	# CY2
Same Commenter	25 E T T T T T T T T T T T T T T T T T T	100 ± 077	30.00 H 10.00	40.99 ± 24.89 08				9.0	3 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6	T 00 30 1733 H	1971 #F23 H	1.00 CO.14 H 13.05	9000	15 H T T T T T T T T T T T T T T T T T T	39.36 ± 18.03	0001 T 1000 44	1179	6.4	1000	_	- 6	170 4679	17.00	190	
	56 37 ± 6.40	69577 E 3009	_		24.75 ± 10.60	00 T T T CO	70.70 + 7.70		1 16 m 21	+ 4.74	1000	201 23 30 ± 4.40	T 69'04	H .	Н 1	H 3.34	. 8	777	1 1 1 1 1 1	1000	2000	1970 0771	1 00 01 50 1	5 to 55 to 1 5 5	
	54 66 ± 8.83	06.00 ± 00.00		1 0 10					10.53	100 22 2011	1 1 1 1 1	74.71	. [-	00.2 + 23.89		188+	-	1000	10989	417 60 77 + 13.70	70.46+		+ 92 22 91	76.44 +	
ohsis	67 30 + 9 31	91.28 + 0.89		-		-			757 71	+ 92.26 - 05.11 +	15.95	70.27+		51.50 + 5.70		101+	26+521 68		+3.86 5103	_	93.25+		+ 374.68+	13.08 94.34+1.7	
	61.93 + 23.77	93.80+5.14	_	58.44 + 23.82 9		-		129	566	85 01	-	66 69.04 + 25.7	+ 95 50	47.78 + 0.42		98.18+1.06 47			+431 5082	526 61 95 + 16 66	94.32	_	72.38+	+ 55 50	56.44 + 8.49
	74.39 ± 6.10		_			-		± 2.49 4	6.51 ± 1.43 75.7	± 2.18 95.94 ±	222 S4.50 ±7	.70 69.13 ± 8.02	92.67 ± 5.39	46.77 ± 1.44	72.10 ± 8.69	93.41±3.65 46	75±1.33 69	17 ± 8.41 93.68	± 3.98 46.94 :	1.51 75.62 ±	2.83 95.79±	226 46.94 ±	1.51 69.43 ± 6	94.34 H 24 H	3 46.94 ± 1.5
	51.01 ± 6.03		49.41 ± 0.06	_	2.56±0.88	_	64.84 ± 6.91	1.96 ± 2.88 49.	41 ± 0.06 60.4	±3.39 4.19 ±	1.03 49.41 ±0	JOS 56.51 ± 12.69	9 4.74 ± 3.08	$49.41 \pm 0.06$	52.26 ± 10.65	2.96±1.38 49	41 ± 0.06 57.83	± 10.23 3	±130 49.41:	: 0.06 52.71 ±	7.89 4.75±	2.11 49.41 ±	:0.06 60.10 ± 0	+00+	49.41 ± 0.00
te	$65.81 \pm 5.31$		49.00 ± 0.05			_	_	± 1.65	9,00 ± 0,00 7,00 ± 0,00 7,00 1,00 1,00 1,00 1,00 1,00 1,00	~	2.14 49.00 ±0	JOS 67.14 ± 6.7	$10.21 \pm 4.27$	48.96 ± 0.07	73.28 ± 4.30		00 ± 0.05   69.81	+7.30	±3.03 49.00:	± 10.08 59.04 ±	4.15 6.65±3	221 49.00±	.0.05 68.38±9	$52  11.35 \pm 4.44$	1 49.00 ± 0.00
Leaths ubs.,5-keto gluconate	$46.85 \pm 0.59$	$4.59 \pm 2.76$	49.14 ± 0.17			_		98	Ή	_	0.17 49.14±0	.17 43.41±1.70	3.26 ± 0	$49.14 \pm 0.17$	50.00 ± 3.48		+1	± 0.25	± 0.67 49.13:	9	7	1.42 49.14	.0.17 41.16±	52 3.46±0.89	49.14 ± 0.15
carbs ubs. N acetylgulc esamine	62.78±4.36	28.16±7.31	_		28.06 ± 4.85	44.29±0.54		8	7.26 ± 1.45 66.3		5.58 45.01 ± 1	39 70.79 ± 2.87	38.08 ±	49.16 ± 3.49	Η.	38.01±123 47	H :	57 ± 0.95 38.41	±2.73 47.67	457 0430#	4.89 32.64 ±	664 45.53	1.45	88 40.44±3.0	
a carbs ubs. Tween 20	57.51 ± 3.75	7.59±1.55	_			_	61.13 ± 5.07	1280	9000	+ 1.47	3.75 48.51 ±0	288 + 3.17	10.70	48.51 ± 0.06	H -	19.07 ± 4.10 48	51 ± 0.06		±3.36 48.51	: 0.06 65.13±	230 14.14 ±	4.10 48.51	0.06 77.59 ±	20.20 ± 2	48.51 ± 0.00
	SECH 0148	CI C# 8771	47.42 ± 0.04		12.76 ± 1.74	<b>5</b>	60.10±3.77	788	# COD #	1539	343 47.42 ±0	54 - 00/US + 4/43	T I STORY	47.78 ± 0.89	H :	7079 ± 9707	#2 H COD # 24	_	18.77 47.42	# #6.20 # #0.00	1000H	492 4742	10.04	1830 H 8	48.60 ± 1.7
	49.02 ± 11.61	3.35 ± 1.52	49.29 ± 0.20	_		ą.	54.38±9.27	1.1	8.29 ± 0.20 61.66	± 5.61 4.85 ±		72.06±	989	49.29 ± 0.20	73.99 ± 3.29	6.87±1.71 49	29 ± 0.20 74	70±2.86 6.80	+1.42 49.29	2	2.76 4.14 ±0	193 49.29	0.20 72.46 ±	33 II.36±6.	49.29 ± 0.20
	53.72±2.95	14.20±1.97	47.03±0.44			# :	62.02 ± 1.27	4.	7.03 ± 0.44 67.7	± 1.67 20.81 ±	+1 -	# 9F-99	23.61 ±	48.54 ± 1.27	69.32 ± 2.87	25.76 ± 3.99 48	#1.73	±6.17 2	±636 47.03:	81	257 20	2.78 47.03	0.44 68.42 ±	70 24.59 ± 2	49.73±5.L
	32 D ± 0.25	2132±3007	# 30 H 0.15	00.10 ± 1.33	31.30 ± 3.03	_		т.	870 97 ± 1967	12.95	# 5 5 5	823 6231±8.01	470 H 470	46.52 ± 1.90	67.65 ± 1.67	77 107 17 17	25 H 5.84	•	187	150.00	10715	392 44.39	T 17 80 CI 10	70 H 200 H 20	28.80
	00.72 H 3.30	0 0 t + 3 1 0	49.09 H 0.10				100 H 400	00 H 00 H	100 H		1 00 00 10 10 10 10 10 10 10 10 10 10 10	H	14.17 H 7.32	49.00 H 0.10	73.16 H 4.00	15.03 H 5.78 49	100	000 /00 H CO	1 2 2 4 40.08	0.10 /0.16 H	H 60 C 1	10.50	H 400 100 100 100 100 100 100 100 100 100	7 H 04 7 1	10 H 07 H 07 H
a Carros unos an omnol	29.30 H 0.30	800 H 128	10 O O O						1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	6471 647 H	1 2 2	000 H 2016 H 000	H - 2	H (704		54 CINCHONS	100 H 000 H	0071	170 P	100.00	- '	707	0.04 0.04 0.04 0.04	CH0071	40 H (704
	28.61 H 9.95		49.20 H 0.00			49.20 H 0000	32 H 1 22	00 H 2.34 45	20 H 0000	1 01 90 1 01 90	H 407 W 107 H	01-32 H 33-4	100 H	45.73 ± 0.43	H 4	3,10 H L33	20 H 0000	12 H 3/61 0.003	+ 264 4647	1 63 63 63 63 1	30.00 +240H	100 97	H 10 10 10 10 10 10 10 10 10 10 10 10 10	7 1 22 17 20	50.76 H 0.00
of columnia.	37.30 H 4.83		45.02 H 0.45	•		_		00 T T T T T T T T T T T T T T T T T T	0.02 E 0.45 00.21	114.00	H 74'04 10'4	20 CH 20 H 27 CH	H 00 10	49.49 ± 0.45	Н 1	22.00 ± 0.51 +5	72 T 0.43	10071	1 7 00 do 40 40 40 40 40 40 40 40 40 40 40 40 40	8 6	H+7/7 900	7070	T 1000 701	00 THOO HO	10 He He He He
-	10' T 67' D		00.00 - 00.00						1 9 9	100	100 to 10	98 ( + 98 ( )	1 12 1	40.40 + 0.00		00 20 TO TO TO	_	8874	10.10	4 3		1000	0.00	20 AVE + 24	W 0 + 0 7 0 F
	58 57 ± 0.82				48 00 + 2 50 4		Sem+171 4	675+116 43	97 29 95 11 +90 9	1 2 63 6	377 44.75 ±7.75	80 + 21 89	+12.95	47.60 + 8.46	971 00 19	CO 17 + 7 1 2	1071	4		1 30 09	T 457 48 6400	0707 313	0.76	1 + 30 19 4	31 + 58 19
	96 1 + 58 55						25.14.1.77	7	+ 0.74	1 1 200	1.08 48.71 ± 0	31 + 89 F5	780+080	48.71 + 0.74	1 +	704+057	_	0 017+	1 2 65 48 10	8	136 1031+	14871	+1925 17.0	17 8 70 + 17	48.21 + 0.2
	89 5 + 29 85						209+2925	8	49+0.26 61.5	+406 968+	16.24	26 64.14 + 4.2	921+996	48.45 + 0.26	56.46 + 10.85	851+350 48	- 0	11568	+142 4848	3	2.58 9.78+	150 48.49	0.26 57.04+	57 8 14 + 2.7	48.49 + 0.24
carbsubs assertaine	53.33 + 8.71	8.75+1.83	48.11 + 0.19	1 12 1 10 10	13.22 + 1.14	48.11 + 0.19	1 199+8109	187 48	0.59 61.0 + 11	+ 151 13.73 +	100 48.11+0	10 63.30 + 2.3	12.50 + 1.61	48.11 + 0.19	62.38 + 2.00	1097+094 48	. 8	+	1,000 48.11	0.00	242 1240+	156 4811	+8119	1027+1	48.05+0.25
								100	1			1	1000					- POV T /		-		-			1

Supplementary Table 18: Per-phenotype performance (2/2). Values are mean  $\pm$  standard deviation across 5 random seeds. Bold highlights the best mean per metric. The results were macro-averaged across classes for each phenotype.

-De	AUROC MA	Mistral-DNA AUPRC	FI AUROC		DNABERT-2 AUPRC F1	AUROC	Nucleo	dide Transformer AUPRC F1	AUROC	ProkBERT AUPRC	E	AUROC	ESM-2 AUPRC		AUROC AU	ESM-C AUPRC F	F1 AUROC	ProtBERT NC AUPRC	E	AUROC	gLM2 AUPRC	E	AUROC	Bactomer	FI
arbeabs,caprate	ı				ı		11.25 10.79 ± 9.42		18   76.79 ± 13.71	L	49.53 ± 0.18		14.20 ± 8.97 45	l–		ı	49.53 ± 0.18   84.68 ± 7.18		ľ	5   78.00 ± 13.09	11.39 ± 8.08	49.53 ± 0.18   8		ı	49.53 ± 0.18
8		42.87 ± 2.60 41.3	41.28±2.38 61.98±3.47		42.19 ± 3.58 40.40 ± 0.89	0.89 61.67 ± 2.17	9	4.66 53.39 ± 8.56	61.85±	$42.32 \pm 8.10$	$46.20 \pm 5.28$	75		48.65 ± 4.62 62	253	40.64 ± 3.97 43.35	43.35 ± 3.52 63.12 ±	5.60 43.44 ± 7.88	47.19	9	$42.21 \pm 7.75$		65.19 ± 3.25	43.23 ± 2.36 5	51.73 ± 10.61
arbeabs citrate							4.41 24.93 ± 3.17		± 19:19	$32.09 \pm 3.96$				_	۳,		± 0.38 72.12 ±	3.83 35.53±3	\$2.01						$51.01 \pm 6.63$
							1.77 10.46±	1.53 47.75±0.34	_	$12.83 \pm 1.63$	•	2	•	150		47.8	± 0.31 S4.12 ±		47.8	-			58.35±3.78		78±0.30
	61.06±3.96 S		06±0.17 SS.11±15.45			0.17 66.31 ± 5.40	0, 1	4906±	-	13.89 ± 14.99		_	10.84 ± 3.16 45	_	73.45 ± 4.49 9.84	9.84 ±1.92 49.06	75.78	1234	49.06	_	6.11 ± 0.64	_	71.87±6.50	+ 5.57	50.47 ± 2.59
9							5.80 15.75±5.00	±97.0±	-		Τ.			12 T T T T T T T T T T T T T T T T T T T	28 ± 298		74.87	H 443 10.73 H	48.78	_	7.01 ± 2.82	61	18.97 ± 2.43	13.40	20 ± 1.61
	35.77 ± 1.57 42		40.88±7.74		49.59±2.50 46.99±		48.49	* 77.77	200	SI.70±1.46	47.86±9.29			100		S4.51±1.77 S5.28	0	3337±	47.01	_	30.73 ± 2.08	52.68 ± 2.85	65.87±2.11		SL65 ± 9.04
							100	#7.71	S7.24 ±	1250±5.57	47.71 ± 0.29			-		-	20.00	# 15.47 14.54 # 8.71		-	11.75 ± 4.00	9	00.08±4.30	•	7.71 ± 0.29
	35.52 ± 4.32 35		_		±66.66	_	3.05 39.40±4.05	42.09 ±	241 0121 ± 221	44.61 ± 3.10	46.10 ± 3.71	8!	1.3.13	228	٠.	~ ·	1.75 ± 3.62 66.22 ±	2.88 47.04 ± 3.58	35.43	_	42.91 ± 3.89	42.03±2.55	68.15±1.36	17.80	0.54 ± 4.76
romic acid			_		49.41 +	88	11.57 2.62 ±	1.00 49.41±00	SS 1 = 1.83	3.18 ± 0.95	49.41 ± 0.08	48.75 ± 12.17	+1.57	# 0008	2015	4	± 0000 ± 49.58 ±	Η.		_	3.18 ± 1.09	970	6.89±18.86	н	41 ± 0.08
			_		+89'84	_	0.60 7.24±	+1	16 SSS++5.55	7.00±0.59	48.69 ± 0.16	51.45±6.22	•	1000	-	24 ± 0.70 48.69	± 0.16 58.53 ±	2.57 7.75±1	1.77 48.69±0.10	5939±438	7.07 ± 0.82	+ 0.16	59.42±8.79		8.69 ± 0.16
_			_		18.41 ± 4.89 46.32 ± 0.97	69	3.16 20.64 ±	\$ 46.32 ±	97 62.53±6.26	$22.93 \pm 6.92$	$46.32 \pm 0.97$	61.08±3.84	•	± 0.97 ←	-	7	8	$3.13  23.41 \pm 5.70$	•	€3.23±5.65	$23.15 \pm 7.43$	+ 0.97	62.14 ± 1.01	755 ± 6.87 ±	31 ± 1.00
mine	57.91 ±4.33 3		2		+9.40 +	_	6	49.40 H	0.16 S8.52 ± 8.40	$3.66 \pm 0.97$	49.40 ± 0.16	61.62 ± 1.08	7	± 0.16	7.94 ± 8.00 3.88	•	± 0.16 57.92 ±	5.69 <b>4.91</b> ±1	+9.40	_	$3.91 \pm 0.87$	+ 0.16	54.93 ± 5.82		9.40 ± 0.16
arbeades glacose			S,	76±7.58 67.77±	Ì	_	5.73 67.56 ± 5.49	38.94 ±	$10.45 + 51.91 \pm 8.63$	$65.74 \pm 6.61$	$43.27 \pm 7.69$	59.15±3.74 7	$2.53 \pm 3.03$ 54	54.78 ± 2.12   58	8.44 ± 3.22 71.43	±120 49.16±	± 6.42   55.78 ±	6.18 69.74 ± 4.57	57 46.88±5.93	S2.07±7.54	65.61 ± 5.02	9.35	62.28 ± 1.50	3.86 ± 1.05 5	5.38 ± 4.42
mide			×	89±4.00 4.05±	4.05 ± 2.14 49.56 ± 0.04	Ξ	10.28 2.70±1.24		M 53.00 ± 7.24	$3.28 \pm 2.26$	49.56 ± 0.04	49.67 ± 17.08	2.13±1.24 48	9.56 ± 0.04 70	96 ± 6.30 4.72	L72 ± 2.17 49.56	19.56 ± 0.04   50.92 ±	3.94 3.97±2.77	77 49.56 ± 0.00	57.63±9.18	$3.46 \pm 1.22$	49.56 ± 0.04	69.06 ± 1.44	3.49±0.69 4	19.56 ± 0.04
	57.83 ± 2.42 18		-		19.59 ± 1.55 46.27 ± 0.31	-	21	+8.76 ±	S2 69.27 ± 3.06	24.58 ± 0.90	46.75 ± 0.51	66.30±1.26	0.78 ± 1.79 46	16.65 ± 0.92 65	23 ± 7.23 21.5;		46.27 ± 0.31   66.74 ±	$3.42   20.84 \pm 0.92$	92 4621 ±0.26	69.54 ± 3.41	24.44 ± 1.00	6.3	66.39 ± 4.20	150 ±4.11	7.06 ± 1.07
						_	617 2695+3 Q	4521+	71 4031 + 743	3007 + 487	45.45 + 3.14	58 00 + 7 m2 1	_	353+070 62	S0 + 485 32.4	242+385 4525	+100 4176	+ 12 11 109		4072+674	30.78 + 3.85	1 36	50 47 + 7 18		372+071
	C 381 T 98 FT	OF C90 TOCC	10 76 + 0 10 40 70 4		7 90 08	010 6104 430	1.00	0.81 40.76±0	10 10 14 1 1 2 2	730 ± 50.0	90 0 1 9 10	30 (1 + 07 97	1133		20 T PO T 98	90 09	1010 6101	130 7 361 139		44804718	190 + 990	010	40 84 ± 6.47	115 + 000	0.00 + 0.10
			į		1 100	-	12.20		20107 31	14 67 + 106	47.05 + 0.16	C) 80 - 4 77	9	910	100 1 212 163		91.0	276 14 21 4 164		51611737	13.41 + 3.21	910	50 FE T 1 26		C 05 + 0 16
arceans g by ogen		C77 H0071	5 8			-		H 70774	20.10 H 2.7	90 H (9'41	47.03 H 0.16	0770 H 4770		9 8	110		2 2 2	H 154 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		104 H101	13.41 H 3.21	9 8 8	20.00 H 47.00		00 H 00
arrestics histidine			_		10.82 ± 1.58 +8.01 ±	000 39.75 ± 3.55	_	3	20'SC # 20'SC	13.30 ± 3.28	48.01 ± 0.09	39.30±3.90	Ξ.	301 ± 000	CDI 77.7 TOS	10.57 ± 1.94 +8.01		±1571 65+	180	888	1207 ± 3.14	8	11.7 = 57.79	e e	730 ± 0.12
			_			-		49.24	-	454 ± 136	49.24 ± 0.12	62.84 ± 5.40	•	± 0.12	_	•	19.24 ± 0.12   59.67 ±	7.03 6.60 ± 3.47	48,75	_	5.33 ± 0.80	7	55.59±8.76	Ĩ,	19.24 ± 0.12
	67.30 ±3.44 2S		_			-	23	+08'60 +	7.70 68.68 ± 4.59	$32.63 \pm 4.21$	$45.60 \pm 0.37$	2	•	+ 0.37	±634	•	7.88±4.30   66.86±	1.94 25.18 ± 3.56	45.78	64.94 ± 2.26	$29.40 \pm 1.21$	-	±6.79	±6.77	51.96 ± 5.91
arbsubs_Jactose	e,		_			-	3.63 37.65 ± 9.32	₹90%	Ť	$42.95 \pm 4.25$	48.11 ± 3.78	1 2 70 4	_		±3.56 4	2005	-	$2.01  44.74 \pm 6.73$	47.41	64.78±3.09	$42.28 \pm 4.59$	265	± 2.07	47.52 ± 5.56	0.79 ± 6.66
	47.12 ± 3.82 7.		ģ		Ī	-		48.49 ±	_	$9.05 \pm 3.10$	$48.49 \pm 0.03$	59.81±2.11	7	18.49 ± 0.03   60.94	± 242	48.49	± 0003 = 59.67 ±	1.50 9.54 ± 1.29	48.49	\$ 57.76 ± 5.74	10.91 ± 6.77	+ 0003	+ 2.49	$10.10 \pm 1.95$ 4	$8.49 \pm 0.03$
			6		26.94±1.88 45.30±0.36	0.36 68.46 ± 4.12	4.12 28.90 ± 4.45	4.45 50.52 ± 4.62	62 69.27 ± 4.59	31.78 ± 8.91	$48.94 \pm 2.81$	72.60±5.35 3	7	5.30±0.36 72	$2.99 \pm 3.65$ $33.40$	3.40±5.06 45.62	$5.62 \pm 0.91$ 72.07 $\pm$	3.99 31.34 ± 6.49	45 45 55 ± 0.3	I 69.38 ± 3.71	$31.48 \pm 6.88$	45.30 ± 0.36	73.26±5.44	483 ± 7.36 4	5.73 ± 0.75
			8		+250 48.63+0.35	-	5.83 7.51+13	+8984	0.35 64.71 + 5.91	11.71 + 4.85	48.63 + 0.35	69.07 + 4.32	0.57 + 1.74 48	89 + 0.35 68	28+4.67 10.7	0.71+1.15 48.63	+ 80.99 6.48 6.08 +	4.55 8.62 + 1.67	67 48.61 ±0.34	67.06 + 2.24	$12.43 \pm 1.95$	48.63 + 0.35	67.62 + 5.16	+268	SE 0 + 6.35
arbeite maltoer		47.85 + 0.77 45	3			-		+ 90 97		55 F + F2 55	48.71+3.11	6121+134		222	39 + 48		_	777 4871+	8 4	3	92 5 + 12 95		109+9219	+433	787 + 176
						-		101 4884 01	3	50 CO T 3 84	\$1 \$0 T 7 78	3	1 2 10	30 11 7	1361		_	91 C T P3 C3 19 C		_	47.05 ± 5.99	1 1	71 71 72 46	1361	20.76 + 1.13
	# # # # # # # # # # # # # # # # # # #		_			_			-	SULP HOSE	51.69 H 4.46		•	_			_	107		_	47.30 H 3.88	_	_	i i	70 H 1-12
			_		7	_	1.80 35.56±3.59	±066	19 61.85 ± 1.76	30.77 ± 4.40	41.17 ± 12.30	60.52±1.61		989	53.98 ± 2.02 59.89		01	8	57.50	2	36.85±5.89	2	_	2	0000 ± 5.68
arbeates, melaboose		27.00 ± 6.20 44.	_			0.08   62.47 ± 5.08	3.08 52.17±	48.75 ±	438 08.23 ± 3.00	47.54 ± 0.67	3009 ± 6.07	5		0.10	# SAS #	48.49	- 687			-	4004 ± 9.18	5.13	-	97	н
hacoeide			_				6.62 6.12±	H	-	$10.21 \pm 5.55$	$49.04 \pm 0.40$	66.51±0.78	7	± 0.40	+ 10.40	•		7.74 6.04 ± 1	1.70 <b>49.04</b> ± <b>0.4</b> 0	198 = 1599	10.66 ± 6.08	± 0.40	3.24	•	9.04 ± 0.40
arbsubs methyl pyruvate	56.03 ±4.29 3	3.02 ± 1.50 49.2	_			_	3.53 8.76±	Ŧ	0.32 64.92 ± 16.00	_	$49.39 \pm 0.32$	67.61 ± 9.27	7	8739 ± 0.32   61.38	± 12.90	•	3	7	49.37	Ť	$6.60 \pm 3.61$	± 0.32	434	0.25	19.39 ± 0.32
			48.20±0.56 51.61±9.17			_	5.32 10.01 ±	3.32 48.20 ± 0.	56 66.59±6.33	$16.10 \pm 3.79$	$48.20 \pm 0.56$	54.32 ± 7.51	9.30 ± 5.15 48	8.20 ± 0.56 53.	$77 \pm 13.71$ 9.64	364 ± 430 48.20	± 0.56 66.83 ±	10.94 17.62 ± 3	3.53 49.03 ± 1.68	8 65.83±8.00	14.23 ± 4.46		59.61 ± 12.50		48.20 ± 0.56
		6.13 + 3.42 48.3	48.88 + 0.29 50.18 +	+12.41 6.80 + 3.65	+ 3.65 48.88 + 0.29	-	+968 565	3.74 48.88 + 0.	0.29 71.02 ± 5.71	8.08 ± 3.12	48.88 + 0.29	71.28 ± 5.19	1043 + 3.20 48	8.88 + 0.29 73	3.00 + 4.99 9.83	0.83 + 2.43 48.88	+ 0.29 74.62 +	4.74 11.58 + 3	3.83 48.88 + 0.29	70.58 + 6.25	8.17 ± 3.08	48.88 + 0.29	75.43 + 3.46	11.03 ± 3.34 4	SSS + 0.29
49.00						85 51 + 87 59 67 0	15.58 5.00 + 4.00	46 24 +	-	12.09 + 14.75	4023+040	. 7	100+	40.42	+ 16.06	4	9 24 + 0.47 64 (7) +	14	40.71	- 1-	11 57 + 13 04	9		1911 + 080	0.24 ± 0.42
			_			_	5 4	10101	27.69 + 3.35	7 64 + 1 70	4010-013	66.33 ± 1.50	•		303 +			10.65 6.35 ± 1.11	91.00	_	716 + 0.60	1		137 + 001	
			46.36 ± 0.21			000 T 100 T	TOTAL CO.	6 11 46 36 ± 0.	20 1 T C1 177	30 C T 170C	45.16 ± 0.13	66.44 + 3.46		900			46.36 ± 0.21	100 325E	46.34	-	21.42 ± 262	100	00 CT 1 30 07	417 + 277	0.00 ± 0.00
						_	1007	1.46 47.77 + 0.	04.14.10	Н -	40.00 ± 0.20	00-77 H 600	107	_			21111070 62421	H +	200	_	16.66 + 3.74	9 0	20 00 T T TO 00	100	10 10 10 10 10 10 10 10 10 10 10 10 10 1
	1 100 H 2011		ŝξ			_	170 PACE 170 PACE 1	DH (7774 047	200	100 100 100 1	40.67 - 0.17	2020 H 2020			100 77 17 1000		H (# 00 H )	200 316±3	90.01	_	10.23 H 3.74	3	0 T T T T T T T T T T T T T T T T T T T	28	00 H 01 H
arcounce, prancocure						-	H077 000	0 H 70 00 00 0	200 H 1970	SWI H 57.7	49.55 H 0.17	40.76 ± 0.02	•	_	•		H 100 CT   100 CT   100 CT		49.49		01 E E E E	100	100 H 100 H		100 H 000
			40.43 ± 0.18 04.22 ±			-	4	# 9CTPC	200	2022 ± 1.77	21.23 ± 0.80	/1.55±5.45	+ '	8			# 777/ CI7 #			0070 = 170	24.0/ ± 1.10	47.	17.541 # 5.77		STS # II
d methyl ester						-	0	*	8	81.8 + 8.71	49.02 ± 0.16	CC-01 = 10.90	7	97.10	٦	•	9707 ± 0.16   57.51 ±	Https	49.07	_	11.23 ± 9.57	9.16	71.49 ± 5.40	2.18	45.02 ± 0.16
DI.			_		•	_	-	*	NS = 12.00	4.01 ± 5.04	49.59 ± 0.08	25.52 ± 16.07	7	9000	00.15 ± 18.80 2.46	•	9730 ± 0008   30°31 ±	10.86 2.47 ± 0.86	49.59	-	4.45 ± 5.27	8000 H	31.30±17.78	81.4	49°29 ± 0.008
			43.87±0.38 65.04±7.59		34.35 ± 10.88 43.84 ± 0.6	_	3	4589±	8	3331 ± 8.54	30.23 ± 5.57	5	7	7.55 ± 5.10   65	75 ± 7.62 55.38	_	3		48.14	_	SOUN # 7.55	8	971	98	16.83 ± 2.94
280	67.25 ± 5.98 42		_			_	4.48 42.45±0.33		9	41.60 ± 4.79	49.44 ± 9.04	5 :		77 T T T T T T T T T T T T T T T T T T	70±05/	•	13 30 ± 0.72 71.57 ±	1.59 42.18±	4800	-	57.62 ± 15.79	98	11.13	15.89 ± 8.22	3.73 ± 6.27
			8		29.13 ± 2.06 45.97 ± 0.06	-	1.40 52.17±	1.43 45.91 ± 2.	00 04:01 ± 4:00	32.49 ± 4.46	45.42±1.89	00.75±5.18	7	5 1	929		#24 ± 1.11   65.06 ±	4.30 50.55 ± 5		_	33.17±6.38	± 1.52	67.590±4.540	887	12 ± 6.63
		25.00±2.50 45.	-			_	8.52 Z9.18 ±	11.28 51.00±5.	79 69:68 ± 2.62	29:59 ± 7.38	$45.76 \pm 1.12$	70.42±2.82 3	4	+ 0.71	3.07 ± 4.39 34.03		45.74 ± 1.002   71.77 ±		46.53 ±	6353±825	23:59 ± 9:00	н	71.53 ± 3.15	± 6.49	45.95 ± 0.98
							_	2.03 47.83 ± 1.	SS 6621±263	$20.60 \pm 3.51$	46.98 ± 0.44	65.48±1.62	_	+	$21 \pm 2.15  20.79$	•	46.98 ± 0.44   66.07 ±	1.47 21.05±	4	66.02±2.53	$20.73 \pm 2.18$	-		22	16.98 ± 0.44
		18.23 ± 0.73 45.2	_		20.44 ± 1.29 45.29 ± 0.43	-	×	2.71 45.29±0.	43 57.83±4.01	$23.80 \pm 1.27$	$45.25 \pm 0.39$	59.62 ± 3.63 2	± 203 +	5.29 ± 0.43 55	$05 \pm 7.73$ 19.64		15.25±0.39 S8.49±	$3.76  22.08 \pm 0.38$	38 45.29 ± 0.43	\$ 58.01±3.25	23.77 ± 0.90	± 0.43	_	8	15.23 ± 0.45
3			8			-	3.47 29.69 ±	2.40 49.39±30	MS 69.03 ± 2.12	34.26 ± 7.19	48.50±3.70	69.38±4.09 2	5.02	5.45 ± 0.15 70	$03 \pm 3.56  30.38$		45.45 ± 0.15 72.44 ±	$^{\rm H}$	45.75	0.46±2.03	$32.80 \pm 4.06$	± 0.82	72.43±1.33	2.53	$47.41 \pm 3.34$
		\$2.39±7.55 45.1	45.14 ± 6.73 60.52 ±		··	-	2.84 56.34±	7.85 58.24±3.	45 60.94±3.59	54.53 ± 5.34	$48.51 \pm 12.23$	66.26±4.07 6	±7.05 5	7.08 ± 7.78 63	46±7.06 55.71	+11.34	48.71±9.70 65.90±	1.99 59.91±7		_	53.73 ± 8.54	+1		12.54 5	$6.71 \pm 15.99$
arbeates threenine	54.40 ± 2.07 9.		_		8.50±2.46 48.29±0.23	_	1.92 9.58 ±	2.90 48.29±0.	23 S8.47 ± 1.10	$11.10 \pm 2.13$	48.29 ± 0.23	56.76±2.34	*	8.29 ± 0.23 57.	00° ± 3°€ 300°		48.29 ± 0.23   55.02 ±	$6.81 - 10.18 \pm 3.64$	.64 48.29 ± 0.2	SS54±2.42	$11.28 \pm 3.06$	48.29 ± 0.23	57.79±1.92		$19.44 \pm 2.13$
			_			-	5.56 3.97±	1.70 49.46±0.	D9 53.62 ± 12.1	2.78 ± 1.18	49.46 ± 0.09	54.95±2.55	± 0.34	9.46 ± 0.09 42	$61 \pm 4.86$ 1.98	•	19:46 ± 0:09   50:07 ±	$323  250 \pm 021$	Ī	53.82 ± 11.94	$2.88 \pm 1.08$	- 0000 H	49.41 ± 7.21	E 0.57	9.46 ± 0.09
arbeabs_trehalose		3634±338 41.7	_			-	423 38.10±	6.12 44.55 ± 2.	40 61.06 ± 2.76	$41.33 \pm 5.35$	49.75 ± 5.24	63.65±2.27 4	339 ± 423 45	523±1.94 63	$21 \pm 3.34 + 43.07$	±6.30 47.00	± 6.62   62.10 ±	2.90 41.12±	41 48.92 ±5.50	5 60.50±2.87	$41.08 \pm 5.77$	49.49 ± 0.44	62.89 ± 3.71	3.30±4.70 §	$1.63 \pm 1.37$
			48		Ī	_	9.86 2.39±1	0.57 49.46±0.	12 36.66±9.77	$1.81 \pm 0.13$	$49.46 \pm 0.12$	47.46±11.56	2.15±0.17 49	0.46 ± 0.12 47.	$58 \pm 9.26$ 2.17	±0.19 49.46	± 0.12 52.96 ±	8.87 2.98±1	1.32 49.46 ± 0.15	44.22±1.38	$2.12 \pm 0.33$	49.46 ± 0.12	51.99±9.37	2.26±0.32 4	46 ± 0.12
arbsubs an dine					Ī	-	1.67 5.02 ±	1.05 49.21±0.	13 58.68±2.96	$4.24 \pm 0.98$	$49.21 \pm 0.13$	63.29 ± 2.21	4.98±0.70 48	9.21 ± 0.13 56	94 ± 4.66 4.52 ::	±1.57 49.21	± 0.13 61.27 ±	7.91 5.06±1	1.63 49.21 ± 0.13	\$ 61.19±2.35	$4.88 \pm 0.33$	49.21 ± 0.13	56.70±6.48	4.37±1.23 4	21 ± 0.13
arbsabs yalerate	59,41 ± 2.81 3		\$		5.28 ± 2.58 49.35 ± 0.31	-	2.94 5.11±	236 49.35±0.	31 58.51 ± 3.52	3.97 ± 1.41	$49.35 \pm 0.31$	60.27±5.90	4.43 ± 2.11 45	9.35 ± 0.31 62	90±7.86 4.89	.80 ± 2.37 49.35	± 0.31 60.39 ±	6.62 S.10 ± 2	250 49.35 ± 0.31	56.59 ± 3.11	4.27 ± 1.90	49.35 ± 0.31	63.68 ± 9.04	9.55±5.59 4	35 ± 0.31
		4.76±0.27 48.8	48.81±0.12 50.24±	24±7.82 5.50±	± 0.56 48.81 ± 0.12	-	6.86 5.02±1	0.73 48.81 ± 0.	12 54.22 ± 3.70	$6.38 \pm 1.16$	48.81 ± 0.12	51.09±8.30	5.47 ± 0.45 48	8.81 ± 0.12 49	04±4.20 4.77	±0.07 48.81	± 0.12 50.83 ±	4.28 6.20±1	14 48.81 ± 0.15	\$0.79±2.57	5.93 ± 1.05	48.81 ± 0.12	51.02 ± 4.85	7.77 ± 1.88 4	8.81 ± 0.12
			8			-	2.26 41.15±	2.52 47.80 ± 11.	42 64.99 ± 2.3	41.14 ± 1.05	50.16 ± 10.52	68.52±2.74 4	2.71 ± 5.41 45	5.80 ± 3.85 67	54 ± 1.92 44.14	±0.77 54.36±	10.80 66.75 ±	4.12 44.08±	.80 49.77±7.72	62.15 ± 5.21	39.38 ± 5.84	41.65 ± 0.82	70.98±1.66	6.54 ± 2.86	7.29 ± 9.36
strin			8		93.15+3.14 91.68+2.64	-	0.20 96.72+	1 + 52 50 19 0	28 08 20 + 0.04	96.97 + 1.29	03.68+0.50	08.70+0.22	96 580 + 29 2	10 1 1 0 1 0 8	83+0.09 97.9	+0.87 97.12	7.12+0.57 99.02+	0.19 98.12+0	90 + 95 26 19	08.12+0.01	96.47 + 1.51	03.46+0.40	08.02 + 0.27	786+087	7.25 + 0.48
atezorical_metabolism		32.14 ± 2.82 24.9				1.15 78.81 ± 2.89	2.89 41.98 ±	3.42 35.66 ± 6.	54 8133 ± 1.64	45.04 ± 5.20	41.21 ± 3.47	83.01 ± 1.89 4	391	8.75 ± 3.69   8.3	28 ± 3.42 45.85	+4.95 39.30	+327 8641+	1.13 49.38 ±	11 4607 ±27	7950+432	41.65 ± 6.47	36.22 ± 6.56	86.68 ± 1.62	0.99 ± 1.73	7.58±1.38
aterorical, motility, binary		37 ± 15.70 63	63.57 ± 4.18 64.77 ± 3.61		59.78 + 12.23 58.77 + 5.81	-	9.84 71.49 ±	3.32 64.51 ± 4.	45 82 39 ± 5 24	81.20 ± 6.01	76.59 ± 3.37	8573±6.11 8	2.24 ± 7.14 77	7.58 ± 2.21 8.3	77 ± 4.53 83.3	+1.64 72.39	+ 3.69 78.93 +	3.06 72.34 ± 1	\$21 68.58 ± 6.70	8241±534	81.05 ± 3.71	75.79 + 2.25	89.80 ± 3.04	8.87 ± 3.73 8	1.67 ± 2.29
ateacairal anomalation		75 77 11 + 05.				_	8	8 17 80 45 + 2	31 - 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	95 9 + LU 68	80 48 + 1 80	97.11+1.06	0 61 + 4 68 07	76 100 107	11 10 291 7 43	92.70 15.74	+ 2.90 07.10+	1.06 88.42+5	87.77 + 10.8	8 9627+233	88 40 + 8 00	01 26 + 273	07 63 11 17	071 701	11+408

# D Broader Impact & Limitations

Broader impact BacBench provides the first public, multi-task testbed that spans gene-, system-and genome-scale prediction problems over 67k genomes from 17.6k species. By standardising data splits, evaluation code and baseline embeddings, it lowers the entry barrier for machine-learning researchers who lack domain-specific pipelines yet want to work on microbial genomics. In the near term, more reliable essential-gene or antibiotic-resistance predictors could shorten drug-development cycles and inform stewardship policies, while better phenotype-from-genome models will accelerate the search for chassis strains that sequester carbon, degrade waste or synthesise valuable biochemicals. Because the benchmark emphasises cross-species generalization, methods that succeed on BacBench are naturally suited to poorly studied or newly sequenced taxa, helping global health laboratories track emerging pathogens even when only draft assemblies are available. Finally, releasing all data under permissive licences and exposing a HuggingFace hub invites continual community contributions, which should foster an open, comparative culture similar to computer-vision or NLP benchmarks and, in turn, drive rapid, reproducible advances in microbial bio-AI research.

Limitations Despite its breadth, BacBench still samples an uneven slice of bacterial diversity: phenotypic-trait labels cluster heavily around medically important genera, and some antibiotic classes remain sparsely represented, which could bias models towards well-studied lineages and mechanisms. Tasks that matter for ecology and biotechnology—horizontal-gene-transfer detection, host–phage interaction, metabolic-flux prediction or transcriptome conditioning—are absent, so performance on BacBench should not be interpreted as general mastery of bacterial genomics. Moreover, the benchmark inherits experimental noise from upstream databases: STRING DB interaction scores mix heterogeneous evidence; operon annotations are incomplete; and phenotype labels amalgamate disparate growth protocols, introducing label uncertainty that caps achievable accuracy. Finally, computing embeddings for every update is resource-intensive, which may hinder participation from groups without access to multi-GPU servers, although smaller surrogate splits are planned for future releases.