

 Open access • Posted Content • DOI:10.1101/524280

Distributed representations of protein domains and genomes and their compositionality — [Source link](#)

[Adrian Viehweger](#), [Sebastian Krautwurst](#), [Brigitte Koenig](#), [Manja Marz](#)

Institutions: [University of Jena](#)

Published on: 18 Jan 2019 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

Topics: [ENCODE](#)

Related papers:

- [An encoding of genome content for machine learning](#)
- [Efficient algorithms for gene cluster detection in prokaryotic genomes](#)
- [Exploring neighborhoods in large metagenome assembly graphs reveals hidden sequence diversity](#)
- [A graph-based algorithm for mining multi-level patterns in genomic data](#)
- [Approximate, simultaneous comparison of microbial genome architectures via syntenic anchoring of quiver representations.](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/distributed-representations-of-protein-domains-and-genomes-18bkiwqfje>

1 An encoding of genome content for machine learning

2 A. Viehweger ^{1,2,4}, S. Krautwurst ¹, D. H. Parks ³, B. König ², M. Marz ¹

3 (1) Department of Bioinformatics Faculty of Mathematics and Computer Science
4 Friedrich Schiller University Jena
5 Leutragraben 1
6 07743 Jena

7 (2) Institute of Medical Microbiology
8 University Hospital Leipzig
9 Liebigstraße 21
10 04103 Leipzig

11 (3) Australian Centre for Ecogenomics
12 Level 5, Molecular Biosciences Bldg
13 University of Queensland
14 ST LUCIA QLD 4072
15 Brisbane, Australia

16 (4) Corresponding author: adrian.viehweger@uni-jena.de

17 Abstract

18 An ever-growing number of metagenomes can be used for biomining and the study of microbial func-
19 tions. The use of learning algorithms in this context has been hindered, because they often need
20 input in the form of low-dimensional, dense vectors of numbers. We propose such a representation for
21 genomes called **nanotext** that scales to very large data sets.

22 The underlying model is learned from a corpus of nearly 150 thousand genomes spanning 750 million
23 protein domains. We treat the protein domains in a genome like words in a document, assuming that
24 protein domains in a similar context have similar “meaning”. This meaning can be distributed by a
25 neural net over a vector of numbers.

26 The resulting vectors efficiently encode function, preserve known phylogeny, capture subtle func-
27 tional relationships and are robust against genome incompleteness. The “functional” distance be-
28 tween two vectors complements nucleotide-based distance, so that genomes can be identified as similar
29 even though their nucleotide identity is low. **nanotext** can thus encode (meta)genomes for direct
30 use in downstream machine learning tasks. We show this by predicting plausible culture media for
31 metagenome assembled genomes (MAGs) from the *Tara Oceans Expedition* using their genome content
32 only. **nanotext** is freely released under a BSD licence (<https://github.com/phiweger/nanotext>).

33 Introduction

34 The content of its genome constrains the functions an organism can perform. Yet, the definition of
35 function and its representation remain elusive¹. A genome can be abstracted as a sequence of protein
36 domains, or by analogy as a document containing words. When words are basic units of “meaning”,
37 protein domains are basic units of function. Embedding protein domains in a vector space captures even
38 subtle aspects of this function. The embedding can be extended to entire genomes. This results in topic
39 vectors which *distribute* a genome’s functions across latent features². The topic of a document might
40 be “half sport, half politics”. By analogy, the topic of a genome represents its metabolic constraints.
41 Topic or genome vectors encode evolutionary and ecological relationships. They can be used as direct
42 input to learning algorithms. This enables large scale metagenomic applications such as biomining and
43 genotype-phenotype mapping.

44 In metagenomics, the bottleneck of discovery has shifted from data generation to analysis. Many
45 current sequencing efforts generate huge amounts of data. It is not uncommon to reconstruct thousands
46 of unknown genomes in a single study³⁻⁵. This wealth of data holds tremendous potential. From the
47 discovery of new enzymes and metabolites to models that predict disease. Related patterns can be
48 detected by powerful learning algorithms such as neural nets⁶. Learning is most effective when the
49 signal in the data is stable across samples. For example, the functions a microbial community encodes
50 tend to be more stable than the taxa it contains^{7,8}. To “fit” metagenome-derived functions into learning
51 algorithms, two questions need to be answered: How is function defined? And how is it represented?

52 We define function as a sequence of protein domains. Most proteins have two or more domains and
53 the nature of their interactions determines a protein’s function(s)⁹. While amino acids build proteins,
54 protein domains represent basic units of “meaning”: They evolve independently¹⁰ and their structure
55 is often more conserved than their sequence¹¹.

56 Several methods exist to represent function. Most commonly, they use protein domain counts or a
57 gene presence-absence matrix¹²⁻¹⁴. But none of those methods retains context information. Not only
58 can context help to identify protein domains¹⁵; in microbial genomes, context is vital. Genes in gene
59 clusters are often co-located in *polycistronic* open reading frames (ORFs)^{16,17}. Adjacent ORFs can even
60 transcribe as a single mRNA in a form of co-regulation¹⁸. Most protein domain representations have
61 another disadvantage. They are high-dimensional and sparse, which hinders downstream learning. For
62 example, the *one-hot-encoding* of a Pfam protein domain¹⁹ has about 17 thousand dimensions, with all
63 elements zero except one.

64 A representation that preserves context and uses dense vectors are *word embeddings*^{20,21}. They assign

65 words with shared context in a document to similar regions in vector space. This assumes that
66 shared context implies shared “meaning”. Word embeddings capture precise syntactic and semantic
67 relationships such as synonyms²². Word embeddings assign each distinct word (*vocabulary*) in a large
68 text collection (*corpus*) a vector of real numbers. The most popular training algorithm is **Word2Vec**²².
69 It extends to entire documents to create topic vectors^{2,23}. Word vectors are a popular input for learning
70 algorithms because they ease training. Here, the embedding acts as a pre-trained language model. The
71 learning algorithm can leverage this knowledge and focus on the task of interest. Embeddings have
72 been trained on biological objects such as genes²⁴ and proteins²⁵. Yet, because they model the primary
73 sequence, abundant sequence variation will act as noise and hinder training.

74 Our aim was to encode genomes using protein domains and to use them as input for learning algo-
75 rithms. The encoding should have several properties: It should represent genome content, not external
76 genome labels such as taxonomy. It should preserve known evolutionary and ecological distances be-
77 tween genomes. The representation should not be sensitive to small sequence variations or incomplete
78 genomes. Formally, it should be a low-dimensional, dense vector of real numbers. Finally, downstream
79 learning tasks using the representation should predict with high accuracy.

80 Results

81 A topic model of microbial genomes

82 The encoding model needs a large collection of protein domains (*corpus*) to train on. We collected about
83 750 million Pfam domains from 150 thousand genomes in the *Genome Taxonomy Database* (GTDB,
84 release r89)²⁶. Ten thousand domains are distinct (*vocabulary*). Their frequency is comparable to
85 English: Several protein domains are abundant and present in all genomes, like the ABC transporter
86 (PF00005). Other domains are rare and present in only a handful of genomes. The majority of domains
87 follows an approximately log-normal distribution (Fig. S1, A). The corpus was then used as input to
88 the Word2Vec algorithm to train an embedding. This weighted, context-sensitive encoding distributes
89 a genome’s content over a vector of numbers. One can then infer new vectors for genomes not seen
90 during training, such as metagenome-assembled genomes (MAGs).

91 Many hyperparameters influence training. We performed grid search over 96 combinations of param-
92 eters (Fig. S1, A and supplementary table 1). Most critical are parameters that adjust the sampling
93 weights of the domains, i.e. the *unigram distribution* (see methods, Eq. 4). Changes in those parame-
94 ters result in a spectrum of “core” and “accessory” models. Core models focus on domains shared by
95 many genomes. Accessory models focus on domains that are niche- or strain-specific or undersampled.

96 Model selection based on ecotype separation

97 To select between alternative embeddings, many so-called *tasks* exist, on which these models are evalu-
98 ated. One such task is for example to recognise synonyms^{25,27}. We found only one task to translate to
99 genomes – the *semantic odd man out* (SOMO) task²⁸. For a set of words (protein domains), we try to
100 identify the one that does not fit the context. For example, “cereal” is odd in the context “zebra, lion,
101 flamingo”. But this task only scores how accurate word-level vectors are. To score document-level vec-
102 tors, we propose the *ecotype* task: For a given species, the aim is to separate known ecotypes. Ecotypes
103 are subpopulations which are adapted to specific environmental niches. We obtained published ecotype
104 labels for three representative bacterial species (<https://github.com/phiweger/nanotext#data>). First,
105 the cyanobacterium *Prochlorococcus* (Fig. 1, C) can grow at a range of depths in the oceans’ surface
106 waters²⁹. There are two recognized ecotypes that reflect oceanic niche differentiation³⁰ – high-light
107 (HL) and low-light (LL). The HL and LL ecotypes can be further refined into subclades¹⁴ (Fig. 1,
108 D). Second, *Pseudomonas koreensis* (Fig. 1, B) is a group of related species inside the *Pseudomonas*
109 *fluorescens* complex. Two ecotypes colonise either bulk soil or the rhizosphere³¹. Third, *Vibrio para-*

110 *haemolyticus* (Fig. 1, A) is a human gastrointestinal pathogen that lives in warm brackish waters.
111 Subpopulations are associated with different geographic locations³².

Table 1: Accuracy of selected models on the SOMO and ecotypes task. The core model performs best on the SOMO task, while the accessory model excels in ecotype separation.

model	no.	SOMO	VppAsia	VppUS1	VppUS2	VppX
core	93	84.8	86.1	19.4	56.5	77.1
accessory	22	82.9	89.9	52.8	71.7	80.7
uniform	45	82.8	87	38.9	65.2	79.3

112 The SOMO and ecotype tasks helped select three models. Model 93 and 22 correspond to core and
113 accessory models, respectively. Model 45 weights all domains equal. Against a baseline of 16.7%
114 accuracy for random guessing, these three models obtained 84.8, 82.9 and 82.8% accuracy on the
115 SOMO task (Table 1). Generally, models with larger context sizes tend to perform better (Fig. S2).
116 On the ecotype task, the “Vpp” subpopulations of *Vibrio parahaemolyticus* were hardest to separate.
117 This makes them valuable for model selection. Accessory model 22 was best on the ecotype task
118 across all labels (Table 1). Generally, models that give more weight to rare domains outperform core
119 models (S3, supplementary table 2). But they fail when genomes are incomplete. By definition, there
120 are fewer rare domains in a genome than common ones. When the assembly is incomplete, these
121 rare domains are more likely to be missing. This distorts the encoding and results in less accurate
122 genome relationships (see below). Model selection should thus not only rely on task metrics, but on
123 desiderata and on principle, too. As our final model, we averaged the three best-performing models in
124 an *ensemble*. The ensemble model provides clear separations of ecotypes (Fig. 1, A-C). The ensemble
125 model can even distinguish several subclades of *Prochlorococcus* (Figure 1, D). This was not possible
126 in the original study using information on gene presence and absence alone¹⁴. Our results show that
127 subtle ecological relationships can be inferred from genomic content.

128 **Ecological and evolutionary relationships overlap at higher taxonomic ranks**

129 A fundamental unit of biological diversity is the *species*. Yet, for bacteria and archaea, this concept
130 is subject to ongoing debate³⁴. One ecological definition uses shared genome content to estimate how
131 similar two organisms are^{35–37}. Shared genome content can also predict phylogeny^{38,39}. Whether this
132 is true at the level of species has been questioned⁴⁰.

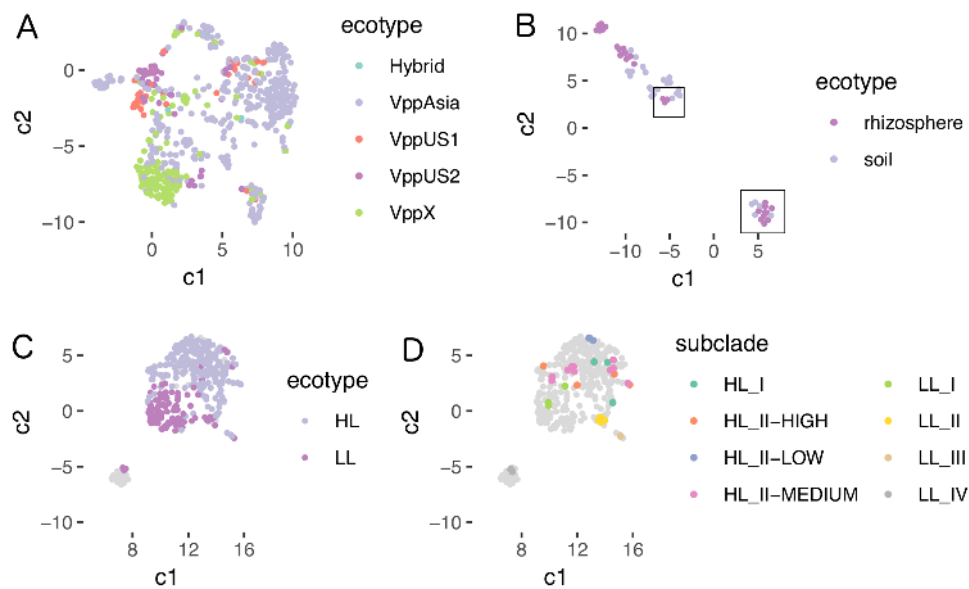


Figure 1: Clear separation of ecotypes using the ensemble model. Vectors were inferred for each genome in the task set and projected into the plane with UMAP³³. Point colors correspond to ground truth ecotype labels^{14,31,32}. **(A)** *Vibrio parahaemolyticus* ecotypes are hardest to separate. They also discriminate best between models. **(B)** *Pseudomonas koreensis* is well separated into its ground truth ecotypes. The original study marked ecotypes by location, not genome content (supplementary figure 3 in *Lopes et al*³¹). This explains why several rhizosphere samples (R15, R23, R24, R45) group with the soil ecotype and vice versa (B2, B3, B7, B11, B35) (boxes). The model distinguishes further clusters within ecotypes, suggesting even deeper niche-differentiation. **(C)** *Prochlorococcus* ecotypes are well separable into high-light (HL) and low-light (LL) ecotypes. **(D)** *Prochlorococcus* ecotypes separate further into subclades.

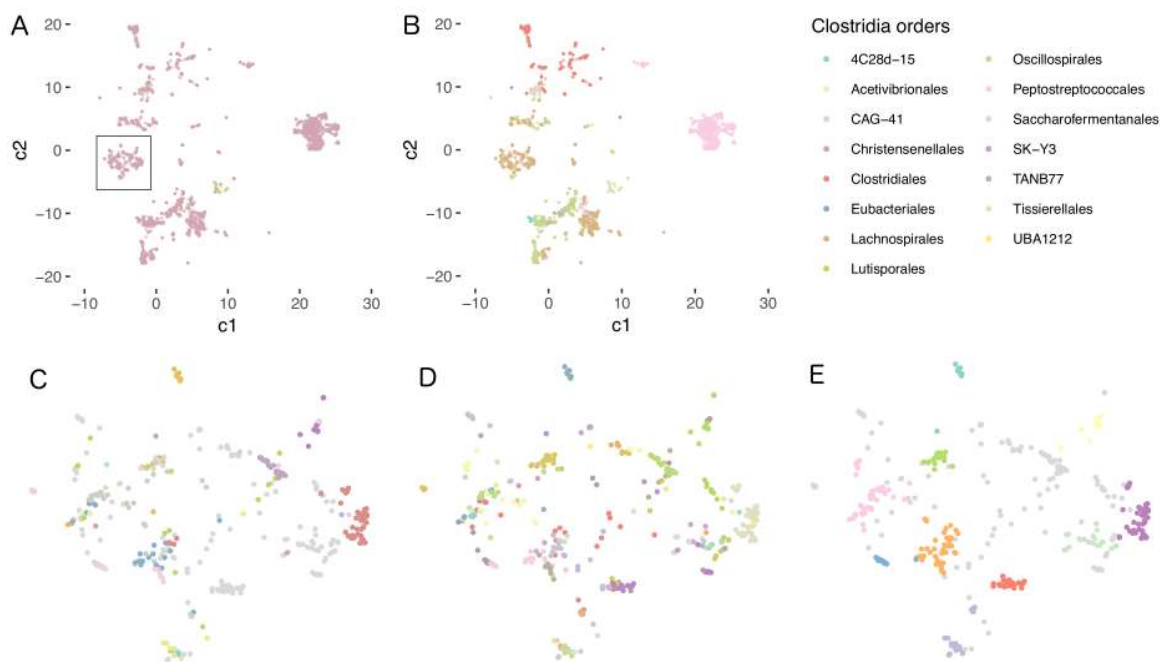


Figure 2: Clustering of genome vectors corresponds to known taxa for ranks down to genus level. Genome vectors from class *Clostridia* were colored by order. *Parks et al.* recently reclassified these orders based on relative evolutionary divergence²⁶ (see figure 5 therein). **(A)** The original NCBI taxonomy, where most of class *Clostridia* consists of order *Clostridiales* genomes. **(B)** The reclassified taxonomy released as *Genome Taxonomy Database* (GTDB, r86). It introduces new orders which correspond well to large clusters of genome vectors. **(C)** Detail from A (box); color labels correspond to genera. Many genomes have no assigned NCBI taxonomy (grey). Their clustering matches taxon labels only in part. **(D)** The corresponding GTDB taxonomy is congruent with *nanotext* clusters **(E)**.

133 We quantified the overlap between phylogeny and genome content across all ranks. For this, we
134 compared GTDB and NCBI taxa labels to **nanotext** clusters. The GTDB labels derive from an
135 evolutionary model for 120 ubiquitous single-copy marker proteins²⁶. In it, *Parks et al.* extensively
136 reclassified many taxa based on relative evolutionary divergence. **nanotext** is a weighted model of
137 genome content. Its genome vectors form nested clusters which correspond to known taxa. Clusters
138 are more congruent with taxa labels from the GTDB taxonomy, as compared to NCBI (Fig. 2, A and
139 C). Compared with GTDB taxa, coarser clusters correspond to orders (Fig. 2, B), while the most
140 granular clusters correspond to genera (Fig. 2, D and E). At the level of GTDB species, there is little
141 correspondence. **nanotext** clusters contain homogeneous taxa down to the genus level for bacteria
142 and to the family level for archaea (Fig. 3, A). Archaea make up less than 5% of the corpus. Our
143 model has thus too little data to separate archaeal genomes as well as it can separate bacterial ones.
144 It is also too little data to train a separate “archaea-only” model. Our results suggest that metabolic
145 potential and taxonomy are coherent only at higher taxonomic ranks. This was observed recently in
146 an independent study using other methods⁴⁰.

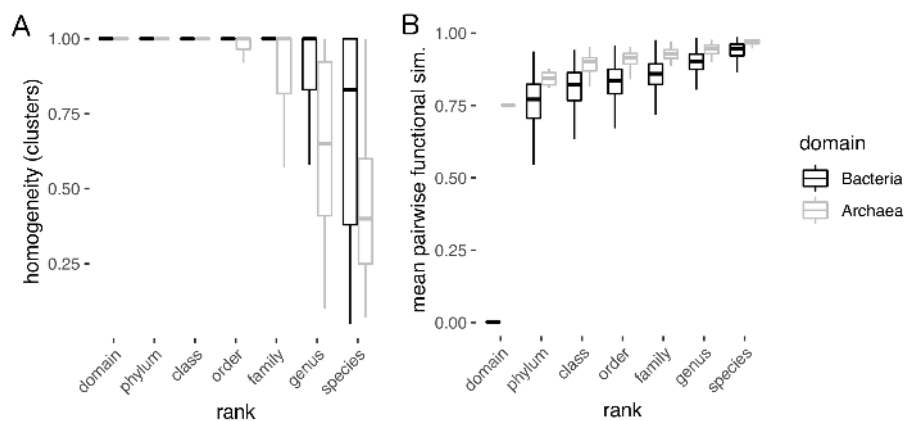


Figure 3: Metabolic potential and taxonomy are coherent at higher taxonomic ranks. **(A)** In the most granular **nanotext** clusters, taxa labels are homogenous at higher ranks. Homogeneity is defined as the fraction of the most abundant taxon within a cluster. At the species level, clusters contain heterogeneous labels. This indicates that species fill specific functional guilds. The larger heterogeneity of archaea is an artifact of the low number of archaeal genomes (< 5%) in the GTDB. Given more data, the model could separate these genomes better. **(B)** The mean pairwise cosine similarity was calculated across all rank labels in the GTDB, i.e. each estimate is calculated “intra-rank” for a given rank label such as “class *Clostridia*”. We downsampled each rank randomly to $n = 200$ genomes. As expected, lower ranks have a high self-similarity compared to higher ranks, approaching zero at the domain level. Again, this effect differs for archaea, due to far fewer genomes in the corpus.

147 We then “flipped” the analysis to quantify self-similarity within ranks. Here, self-similarity is the
148 average pairwise cosine similarity between pairs of genomes at each rank. As expected, self-similarity

149 increases for lower taxonomic ranks (Fig. 3, B). The average relatedness of any two bacteria approxi-
150 mates null. For archaea, the suboptimal separation of genomes results in a higher self-similarity across
151 all ranks. In summary, **nanotext** clusters are coherent with known taxa at higher ranks. At lower
152 ranks, our model captures ecological features. This is supported by the ability of the model to separate
153 ecotypes (Fig. 1).

154 **Genome vectors are stable even for highly incomplete genomes**

155 One main use case for genome vectors is to encode MAGs for downstream machine learning. Because
156 most MAGs are incomplete⁴¹, we assessed the sensitivity of our model to incomplete data. All near-
157 perfect GTDB genomes were selected for the truncation experiment ($n = 89$, 100% complete, 0%
158 contaminated). The GTDB taxa are considered ground truth. Genomes were then truncated *in silico*
159 in 10% increments. Each time, the taxonomy was inferred using **sourmash lca**⁴² and the **nanotext**
160 core model. **sourmash lca** assigns taxa based on shared k-mers between query and reference genomes.
161 This approach was first popularized by **Kraken**⁴³. The **nanotext** core model was chosen because its
162 inference is robust towards incomplete genomes (Fig. S5). It puts little weight on rare domains, which
163 are more likely to be missing from incomplete genomes. For up to 70% missing genome content, both
164 methods returned correct taxa (Fig. 4). **nanotext** alone can correctly assign about 5-10% of taxa. Its
165 accuracy drops at about 70% truncation. **sourmash lca** can classify taxa up to about 90% truncation.
166 Note that both methods are very different. **sourmash lca** can “lock in” on a genome if it contains a
167 set of k-mers that is present in its reference database. **nanotext** is able to learn higher-level features
168 of ranks. This makes it more applicable to searches where few close reference genomes are available.
169 In summary, the **nanotext** genome encoding is robust towards even very incomplete genomes.

170 **Functional similarity complements nucleotide similarity**

171 **nanotext** is able to infer relationships for incomplete and distant genomes. This is possible because it
172 complements nucleotide-based distance in two ways. First, **nanotext** uses protein domains for training.
173 These are annotated using *Hidden Markov Models* (HMMs). An HMM in turn distills groups of multi-
174 ple sequence alignments (MSAs). This aggregates nucleotide variation. Second, **nanotext** embeds the
175 protein domains instead of relying on their names or ID. It places similar domains in similar positions
176 in vector space. This creates the notion of “synonymous” protein domains and aggregates groups of
177 HMMs. **nanotext** quantifies the relationship between any pair of genomes with “functional similarity”.
178 This metric is defined as the cosine similarity between genome vectors. Functional similarity com-

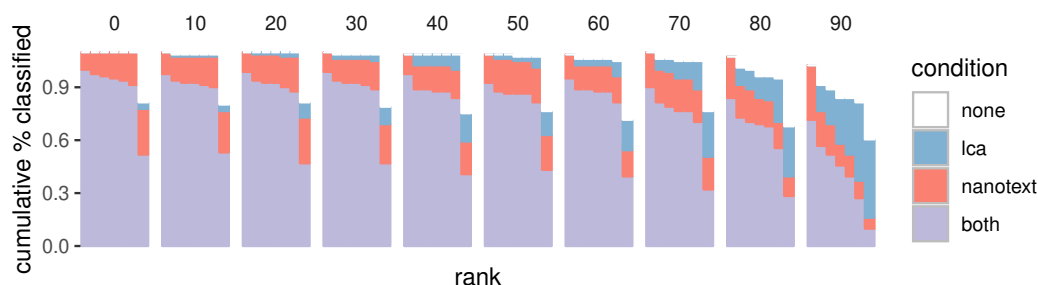


Figure 4: Genome vectors from core model are robust towards incomplete genomes. Near-perfect genomes from the GTDB ($n = 89$) were truncated *in silico* in 10% increments (facets). Each time, both **nanotext** and **sourmash lca** inferred taxa labels across all ranks (x-axis). Either only **nanotext** (red) or **sourmash lca** (blue) were correct, both (violet) or none (white). This resulted in a score for cumulative accuracy (y-axis). Note that about 20% of genomes in the test set had no species label (left most facet, right most bar). **nanotext** is more correct in its assignment of taxa up to about 70% incompleteness. Its accuracy declines thereafter. **sourmash lca** assigns less taxa in general, but does so for up to 90% truncation.

179 plements nucleotide-based distance metrics like *Jaccard* similarity⁴⁴ and average nucleotide identity
180 (ANI)⁴⁵.

181 We illustrate this with a use case, where we aim to identify MAGs of an understudied organism with few
182 reference genomes. We search the *Tara Oceans* MAG collection by *Delmont et al.*⁴ for members of the
183 phylum *Gemmatimonadota*. Little is known about this phylum: Most cultured members live in soil⁴⁶.
184 Recently, marine isolates were discovered that photosynthesize⁴⁷ using a non-canonical mechanism⁴⁸
185 acquired by horizontal gene transfer (HGT)⁴⁹. We first retrieved all genomes labelled *Gemmatimonadota*
186 from the GTDB ($n = 64$). We then overlaid the corresponding vectors with inferred vectors for the
187 *Tara* data (core model, Fig. 5, A). A MAG was returned as putative *Gemmatimonadota*, if it had a
188 cosine similarity $\cos(\theta) \geq 0.7$ to any GTDB record of the same phylum, a conservative threshold.

189 Any similarity search based on exact sequences of nucleotides or amino acids requires close genomes
190 in the reference database. **nanotext** does not share this restriction. Because it relates genomes
191 by content, they can have very different sequence compositions. If their metabolic potential
192 is similar, functional similarity will remain high. We found five genomes that pass our thresh-
193 old (TARA_ANE_MAG_00005, TARA_RED_MAG_00040, TARA_ION_MAG_00042, TARA_RED_MAG_00069,
194 TARA_RED_MAG_00065 – for additional metadata, see supplementary table 3 in *Delmont et al.*⁴
195 and <http://merenlab.org/data/tara-oceans-mags/>). None of these MAGs was assigned a taxon in
196 the original, marker-gene based study⁴. We confirmed these assignments using k-mer based ANI
197 estimation⁴², *in silico* DNA-DNA hybridization⁴⁵ as well as with a phylogenetic approach²⁶ (Fig. 5,

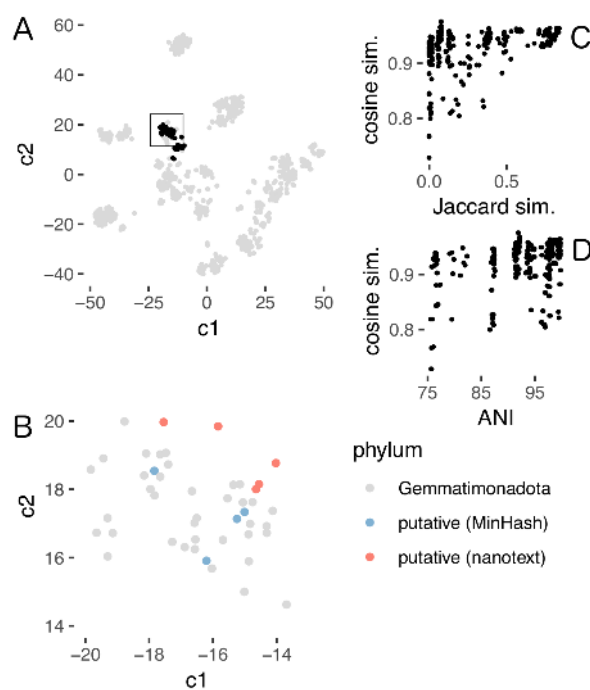


Figure 5: Functional similarity between genomes generalizes nucleotide similarity. **(A)** Members of the understudied phylum *Gemmatimonadota* were identified from a large *Tara Oceans* MAG collection⁴. We overlaid vectors for all *Gemmatimonadota* in the GTDB (black, $n = 64$) with vectors inferred from *Tara* MAGs (grey, $n = 957$). The GTDB genomes form two distinct clusters. **(B)** Detail of A (black points in box now in grey). Putative *Gemmatimonadota* in the *Tara* collection were identified using functional and nucleotide similarity. *nanotext* and *sourmash* both identified the same five MAGs (red). Additionally, the nearest neighbors identified by *sourmash* based on nucleotide similarity using the MinHash algorithm are shown in blue. Note how these are further from the identified MAGs (red) than are the closest GTDB genomes in vector space (grey). *nanotext* can relate these closest genomes with similar metabolic potential even when their nucleotide similarity is low. **(C, D)** Pairwise nucleotide distance between reference and putative *Gemmatimonadota* shows large variance. Yet functional similarity remains stable even past common thresholds of nucleotide-based relatedness (ANI < 0.8, Jaccard similarity < 0.5).

198 C and D; supplementary table 3). In summary, functional similarity can relate genomes that have
199 very different nucleotide compositions.

200 **Accurate prediction of plausible culture media from genome content alone**

201 To test our genome encoding on a machine learning task, we trained a model to predict culture
202 media for MAGs based on their genome content alone. Many unculturable bacteria are thought to
203 be culturable⁵⁰. But it is infeasible to test thousands of medium recipes. Yet, many media share a
204 significant number of ingredients. We therefore trained embeddings of ingredients⁵¹ for each medium
205 in the catalogue of the *German collection of microorganisms and cell cultures* (DSMZ)^{52,53}. To train
206 the genotype-phenotype mapping, we had to link the GTDB to the DSMZ BacDive database⁵³. We
207 then used a fully-connected neural net to learn this mapping.

208 Given a genome vector, we were able to predict plausible culture media with high accuracy (Fig. 6,
209 A). The prediction result is a medium vector. Related media can be retrieved via nearest neighbor
210 search. A common-sense baseline is to always predict the most common label of the data set (medium
211 no. 514), which would result in an accuracy of 17%. We classify a prediction as correct, if the target
212 medium is within the top (1, 10) closest media by cosine similarity. This is a common evaluation
213 scheme in multi-class image labelling tasks⁶. On the test set, our model obtains a top-1-accuracy of
214 63.5% and a top-10-accuracy of 82.5% (Fig. 6, A). On *Tara* MAGs, we obtained a top-1-accuracy of
215 50% and a top-10-accuracy of 73.2% (Fig. 6, B). The lower accuracy on the *Tara* data is likely due to
216 genomes without a close representative in the media training data.

217 A good model should generalise to unseen genome-media pairs. We assessed this on a *Prochlorococcus*
218 MAG (TARA_ION_MAG_00012). There were no records for *Prochlorococcus* in BacDive at the time of
219 writing. Yet, there exist established culture protocols such as “Artificial based AMP1 Medium”⁵⁶. We
220 labelled the AMP1 ingredients with the same protocol used for the training data⁵². We then inferred
221 a medium vector by summing over the ingredient vectors. Of the 10 nearest media in our embedding
222 space, medium no. 737 (“Defined Propionate Minima Medium”, DPM) has a cosine similarity of 0.979
223 to the target AMP1. Half of the AMP1 ingredients are found in DPM, including vital trace elements.
224 Other ingredients that are not shared belong to buffers and could be substituted. In summary, genome
225 vectors can be used to train non-trivial learning tasks. We showed this by predicting plausible media
226 for MAGs as a starting point for culture experiments.

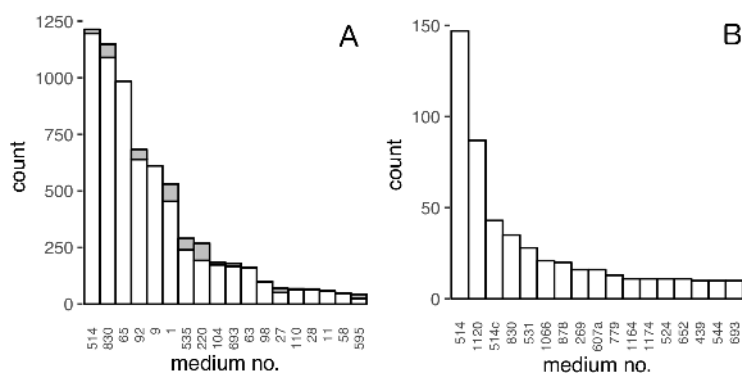


Figure 6: Genome embeddings allow accurate predictions when used as input to learning algorithms. **(A)** Prediction of culture media from genome content. We linked GTDB genomes to DSMZ media at the level of genus. We then trained a fully-connected neural net to learn the genotype-phenotype mapping. Media (x-axis) can be predicted with high accuracy (y-axis, white bars count correct cases, grey bars count wrong cases). A classification was correct, if the the target medium was among the 10 nearest neighbors of the predicted medium vector. Only the 20 most common media in the database are displayed. **(B)** Predicted top media for *Tara* MAGs. The most common media (excluding their variants) in the prediction set are no. 514 (“Bacto Marine Broth”), no. 1120 (“PYSE Medium”) – e.g. used to study *Colwellia maris* isolated from seawater⁵⁴, no. 830 (“R2A Medium”) – developed to study bacteria which inhabit potable water⁵⁵, no. 1066 (“Marinobacter Lutaoensis Medium”), no. 878 (“Thermus 162 Medium”), no. 269 (“Acidiphilium Medium”) and no. 607 (“M13 Verrucomicrobium Medium”) – which includes artificial seawater as an ingredient. All these media correspond to marine isolates. They are plausible starting media for the query MAGs.

227 Discussion

228 In this work we showed how genome embeddings capture many functional aspects of the underlying
229 organism. The main assumption of our approach is that the function of a genome can be abstracted as
230 a sequence of protein domains. Much like words determine the topic of a document, protein domains
231 act as atomic units of “meaning” in a genome. This view of function is very reductive and more
232 comprehensive definitions exist¹. For example, we do not consider functional RNAs⁵⁷. Yet, our results
233 suggest that **nanotext** performs well on a variety of tasks. This success is likely due to a focus on
234 bacteria and archaea, where most functions are protein-mediated⁵⁸.

235 We omit other organisms such as eukaryotes and viruses from the corpus. But these can be included
236 without changes to the method. A major bottleneck in corpus expansion is annotation. Currently
237 most approaches are based on *Hidden Markov Models* (HMMs)⁵⁹. These do not scale well to thousands
238 of genomes. It would be interesting to replace protein domain HMMs with homology-based protein
239 clusters. Such clusters could derive from large protein collections such as *UniRef*⁶⁰ or the *Soil Reference*
240 *Catalog* (SRC)⁶¹. A vocabulary compiled from protein clusters would include many proteins of unknown
241 function. Yet, they could still be used in predictive tasks because they are distinct tokens. But the
242 vocabulary size needs restriction. The **nanotext** embedding was trained with a corpus-to-vocabulary
243 ratio of about $10^5 : 1$. This is in line with current corpora in *Natural Language Processing* (NLP)
244 such as **CommonCrawl** (<http://commoncrawl.org/>). Once a model has been trained, search can scale to
245 billions of vectors⁶².

246 For training the embeddings we used the **Word2Vec** algorithm. It is a special case of exponential
247 family embeddings²¹. Other embedding methods could be better suited. For example, a *market basket*
248 embedding might be more appropriate to embed media ingredients. The embeddings could be further
249 enriched by subword information⁶³. This character-level training on nucleotide sequences would allow
250 inference on *out-of-vocabulary* words. One main limitation of all latent embeddings is that they are
251 not interpretable. Only relative distances within a training run are meaningful. It is thus difficult to
252 reverse engineer our model to identify important domains.

253 Because **nanotext** encodes genome content, it captures even subtle ecological features. At higher
254 taxonomic ranks, these features are congruent with evolution-based methods^{38,39}. At lower ranks,
255 these features complement phylogenetic ones. It has long been recognized that members of a species
256 can be functionally very distinct⁸. For example, *E. coli* can be commensal, enterohemorrhagic or
257 probiotic to the host. There is enormous genomic variation within this species. This diversity is missed
258 in most approaches that focus on marker-genes within the small core genome⁶⁴. **nanotext** encodes

259 the entire genome, an advantage in metagenomics where genomes tend to be incomplete. **nanotext**
260 also uses functional similarity to relate genomes of similar metabolic potential. This complements
261 methods that are based on nucleotide composition. Relationships between genomes can be discovered
262 even though their nucleotide-similarity is below common thresholds. **nanotext** is thus well suited to
263 encode incomplete, distantly related MAGs to either explore or base learning tasks on.

264 We also showed how genome vectors can be used in as input to machine learning tasks. We illustrated
265 this by predicting plausible culture media for MAGs with good generalization to unseen data. One
266 problem was the lack of full genome sequences in current phenotype databases. We had to create a
267 mapping between the GTDB genome collection²⁶ and the DSMZ BacDive database⁵³. This compro-
268 mise likely reduces the predictive power of the learned model. Several strain collections have started
269 to whole-genome sequence their inventory. We thus expect a more accurate model in the future. More
270 generally, learning algorithms can become much more efficient when working from embeddings. The
271 algorithms can focus on the learning task and need not learn a “language model” in parallel. Less
272 training data is needed as a consequence. Existing data can be leveraged for machine learning using
273 **nanotext** in two more ways. First, one can use 16S amplicons and “look up” the corresponding func-
274 tional vector. This approach is similar to PICRUST⁶⁵. Second, for large collections of metagenomes,
275 HMM-based domain annotation becomes a bottleneck. A comprehensive reference database⁶⁶ offers
276 a shortcut. One can then use fast nucleotide-based methods⁴⁴ to search for similar genomes. Once
277 found, their corresponding genome vector can be retrieved by simple lookup.

278 In conclusion, we trained a genome encoding that captures function in low-dimensional vectors. It
279 solves the “curse of high dimensionality” of previous sparse encodings. Genome vectors are insensitive
280 to missing genome content and nucleotide variation. Pre-trained **nanotext** embeddings can be used
281 for taxonomic classification, biomining and phenotype prediction.

282 **Methods**

283 **Model training**

284 We used all roughly 150 thousand genomes from the *Genome Taxonomy Database* (GTDB, release
285 r89)²⁶ for model training. Because the associated taxonomic assignments for release r89 were not yet
286 available at the time of writing, we used the metadata from the previous r86 release. The taxonomy
287 is largely consistent with the expected r89 one (personal communication). Genomes were annotated
288 using `pfam_scan.pl` (v1.6, <https://bit.ly/2CHXIVI>) which resulted in a corpus of 750 million domains.
289 Each line in the corpus is the sequence of Pfam protein domains on a contig. Strand information is not
290 preserved. We did not filter any protein domains. The vocabulary has 10,879 domains, about 60% of
291 the domains in Pfam (v32)¹⁹.

292 We obtained word and document vectors using the `Word2Vec` algorithm^{22,23}. We trained over a grid of
293 hyperparameter combinations (see below) with a linearly decreasing learning rate (0.025 to 0.0001) over
294 10 epochs using the *distributed bag of words* (PV-DBOW) training option as implemented in `Gensim`
295 (v3.4.0, <https://radimrehurek.com/gensim/>). The result were 100-dimensional vectors for each domain
296 and genome. The similarity of any two vectors can be calculated using cosine similarity (Eq. 1), with
297 a range from -1 (no similarity) to 1 (identical). Cosine distance is defined over the range (0, 2) (Eq.
298 2). For inference on unseen genomes, we concatenated the annotation results of all contigs. To reach
299 convergence, inference used 1000 epochs. This resulted in stable vector estimates with a pairwise cosine
300 distance less than 0.01 for repeated inferences of the same genome.

$$\text{cosine similarity } s = \cos(\theta) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} \quad (1)$$

$$\text{cosine distance } d = 1 - s \quad (2)$$

301 **Hyperparameter optimization**

302 The `Word2Vec` algorithm has many tunable parameters, and they matter⁶⁷. We performed grid search
303 over a range of plausible parameters and trained 96 different models. The most relevant parameters all
304 influence the weight given to individual words based on their frequency (Fig. S1, A). We found that
305 the following parameters influenced training outcome most: the subsampling parameter, the negative

306 sampling exponent, the number of negative samples per frame and context window size (Supplementary
307 table 1).

308 The *subsampling* parameter s determines how likely a word w_i with frequency $z(w_i)$ in the corpus is
309 kept during training, expressed as a probability $P(w_i)$ (Eq. 3). Based on this formula, we determined
310 how many of the most frequent domains were affected from subsampling given s (Figure S1, B). We
311 set s to (3e-5, 1.4e-4) which translates into the most frequent (100, 1000) domains being affected from
312 downsampling.

$$P(w_i) = \left(\sqrt{\frac{z(w_i)}{s}} + 1 \right) \cdot \frac{s}{z(w_i)} \quad (3)$$

313 Negative sampling is used during training to teach the neural net which words are not part of the
314 context of a given target word. For efficiency, not all word weights outside the target context are
315 adjusted, but only a small number (between 2 and 20). Negative samples are drawn from a so-called
316 *unigram distribution* (Eq. 4). It scales the count of each domain $f(w_i)$ by an exponent v which was
317 empirically set to 0.75 in the original publications^{22,68}. While this parameter is not usually optimized,
318 it can improve model performance on certain problems⁶⁷. Intuitively, v provides a handle to focus
319 on less frequent protein domains during training. $v \geq 0$ attends to protein domains shared by many
320 genomes (“core”) while $v \leq 0$ focuses on domains particular to only a handful of genomes (“accessory”).
321 For $v = 0$ all protein domains from the “gene pool” are treated equal (Fig. S1, C).

$$P(w_i) = \frac{f(w_i)^v}{\sum_{j=0}^n (f(w_j)^v)} \quad (4)$$

322 For the remaining parameters the following were used: Window (context) size – (2, 5, 10), number of
323 negative samples (2, 5). The remaining settings were left set to their default values.

324 **Model postprocessing and nearest neighbor search**

325 After training of word vectors, the average similarity between random pairs of vectors is not zero.
326 Instead, it approximates a mean vector⁶⁹. We subtracted this mean vector from all model vectors
327 and then unit-normalized them. For fast nearest neighbor search, we used the **Faiss** library (v1.5,
328 <https://github.com/facebookresearch/faiss>)⁶².

329 **Task design for model selection**

330 For the SOMO task, we selected one thousand contigs selected at random from the corpus. For each
331 frame of five consecutive domains, one random sample from the vocabulary was added. The mean
332 of the embedding vectors of this set was then calculated. The “odd” domain has the largest cosine
333 distance to this mean. A frame is correct, when the odd domain is the random vocabulary sample.
334 This results in a baseline accuracy of 16.7 % for random guessing. In the ecotype task, each point
335 received the label of its nearest neighbor. If both shared a label, this counted as correct. The average
336 over all points for each ecotype (e.g. soil) resulted in the final score for that ecotype. An accuracy of
337 1 indicates complete separation of clusters.

338 **Annotation of Tara genomes**

339 To annotate protein domains for a collection of 957 MAGs⁴ we first identified open reading frames
340 using Prodigal (v2.6.3)⁷⁰. We then used `pfam_scan.pl` wrapping HMMER (v3.2.1)⁵⁹ to search against
341 the Pfam database.

342 **Estimation of nucleotide distance and taxonomy**

343 To estimate average nucleotide identity (ANI) between pairs of genomes we used the MinHash
344 algorithm^{44,71} as implemented in `sourmash` (v2.0.0a11, <https://github.com/dib-lab/sourmash>)⁴². The
345 GTDB was used as search index except for the truncation experiments. There, the test genomes
346 were removed before indexing and search. To generate MinHash signatures from genomes, we chose a
347 sketch size of `--scaled 1000` and a k-mer size of 31. For taxonomy inference we used the `sourmash`
348 `lca` subcommand. We validated all distance estimates using *in silico* DNA-DNA hybridization as
349 implemented in `FastANI` (v1.1)⁴⁵.

350 **Clustering and dimension reduction**

351 To visualize clusters of genome vectors in high-dimensional space, we projected them into the plane us-
352 ing either the UMAP algorithm³³ implemented in `umap-learn` (v0.3.7, <https://github.com/lmcinnes/umap>)
353 for large vector sets or t-SNE⁷² as implemented in `scikit-learn` (v0.20.0)⁷³, both with default
354 settings. For clustering we first reduced the vectors to 10 dimensions using UMAP. We then performed
355 clustering using HDBSCAN⁷⁴ as implemented in `hdbscan` (v0.8.19, <https://github.com/scikit-learn->

356 contrib/hdbscan) using default parameters and selecting clusters using the conservative leaf
357 method.

358 **Training of media embeddings**

359 To quantify media similarity, we created a media embedding. The current media collection of the
360 DSMZ lists over 1500 media. To reduce the effective number of media, we treated a medium recipe
361 as a sequence of ingredients. We could then use `Word2Vec` to create a latent representation of
362 ingredients⁵¹. The DSMZ media are not easily parsable and contain many non-unique ingredient
363 tags such as “beef extract” and the synonymous “meat extract”. We therefore used preprocessed data
364 from the KOMODO database of known media⁵². To download all 3,637 recipes, we used a custom crawl-
365 ing script (`scrape_komodo.py`). Note that some current additions to the DSMZ media list do not
366 figure in the KOMODO database. From each recipe we extracted a list of ingredients⁵². We excluded wa-
367 ter (`SEED-cpd00001###`) and agar (`SEED-cpd13334###`) because these ingredients are non-informative.
368 We trained with a window size of 5 and a learning rate as described above over 100 epochs using
369 negative sampling of 15 words per window. To make sure that pairs of media ingredients could occur
370 in the same window, we augmented the data set by shuffling each ingredient list 100 times. The result
371 was a 10-dimensional vector for each media ingredient. To represent a medium, we summed across
372 its ingredient vectors. The similarity of any two DSMZ media could then be compared using cosine
373 similarity. For example, the closest media to medium no. 1 are medium no. 306 (0.99) and no. 617
374 (0.99), one adding *yeast extract* and the other *NaCl* to medium no. 1; an ID-based representation
375 would treat these media as distinct, although they are near identical. Indeed, medium no. 617 and
376 953 have identical ingredients, which is reflected by a cosine similarity of 1. Any predicted medium
377 can suggest n similar media via nearest neighbor search. We visualized the media vector space using
378 t-SNE (Figure S6). Media vectors cluster and enable learning algorithms to discriminate between
379 media classes.

380 **Linking GTDB genome assemblies to BacDive culture media**

381 To predict a medium from a genome we needed to create a training set that matches the two. The
382 DSMZ BacDive database holds taxonomic and phenotypic information including culture media for
383 currently over 60 thousand strains⁵³. However, these strains do not directly correspond to genomes
384 in the GTDB collection²⁶. To link these two, we had to pair records using taxonomy at the rank of
385 genera.

386 **Culture medium prediction**

387 For the medium prediction task, we used a multi-layer, fully connected neural net. We selected the
388 training data as follows: For each genus used to link the two databases, we first sampled records from
389 **BacDive** at the genus level. Because this data is highly skewed towards medically relevant genera such
390 as *Mycobacterium*, we randomly selected a maximum of 100 records per genus. As target \mathbf{y} , we used
391 the vector of the most common medium for each genus. For the same genus we then randomly sampled
392 a genome vector from **nanotext** (\mathbf{x}).

393 We repeated this process 10 times. Data augmentation is a common practice when training neural nets.
394 It enables the training of more complex models, which then generalize better. Using data augmentation,
395 we can circumvent the need to collect more data by varying the input slightly⁶. We used a total of
396 73,916 genome-media pairs for training, optimized hyperparameters on a validation set of 3891 (5%)
397 and tested the final model on a holdout set of size 8,646 (10%). The neural net architecture consisted of
398 three fully connected layers with (512, 128, 64) nodes. Before applying the non-linear transformation
399 (rectified linear units, ReLU), we normalized the batches of size 128. After each layer we applied
400 Dropout (0.5, 0.3, 0.1). The output layer has 10 nodes to represent a culture medium vector with
401 10 latent elements, activated by a linear transformation. We optimized a cosine similarity loss using
402 **Adam** with a learning rate of 10^{-2} over the course of 10 epochs. We did not rescale the variables before
403 training. We implemented the model using the deep learning library **Keras** (<https://keras.io>).

404 **Code availability**

405 All relevant resources to reproduce the major results in this article have been deposited in a dedi-
406 cated **nanotext** repository (<https://github.com/phiweger/nanotext>). This includes source code for
407 the **nanotext** library and command line tool, preprocessing and training workflows as well as an in-
408 troductory tutorial. The trained models, corpus, taxonomy metadata and test data were deposited on
409 the *Open Science Framework* (OSF, <https://osf.io/pjf7m/>).

410 Supplement

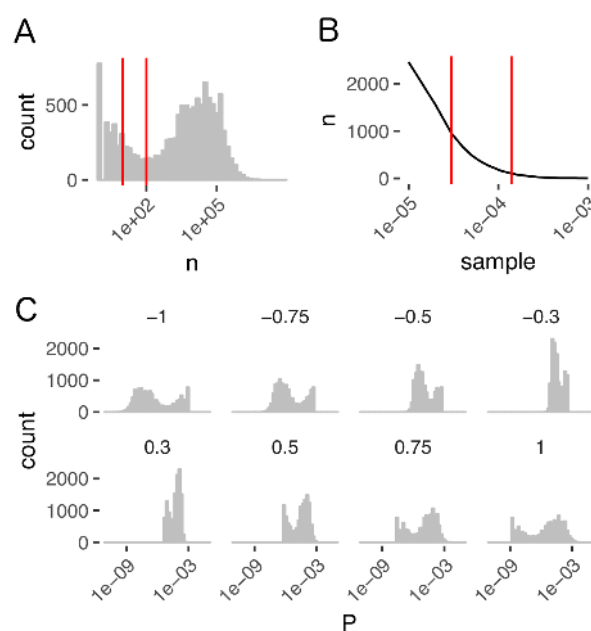


Figure S1: Protein frequency distribution informs hyperparameter selection. **(A)** Count distribution of protein domains. Few domains are present in nearly all genomes while the majority follows a log-normal distribution. There are also many domains present in few copies. Red vertical lines indicate the selected thresholds for minimum count for inclusion in the corpus (left 10, right 100). **(B)** Number of the most frequent domains that are subject to subsampling as a function of different settings for the `Word2Vec` subsampling parameter (Eq. 3). Red vertical lines indicate two values chosen for inclusion in the hyperparameter grid search (left 100, right 1000). **(C)** Influence of the negative exponent v on the *unigram distribution* (Eq. 4). For $v = 1$ the unigram and frequency distributions are identical. For $v = -1$, they are inverted, i.e. the rarest domains are sampled with the highest probability during training. At $v = 0$, the unigram distribution becomes uniform, i.e. all domains are sampled equal. We chose $(-0.3, 0, 0.3$ and 0.75 – the canonical default in most applications) for grid search during training.

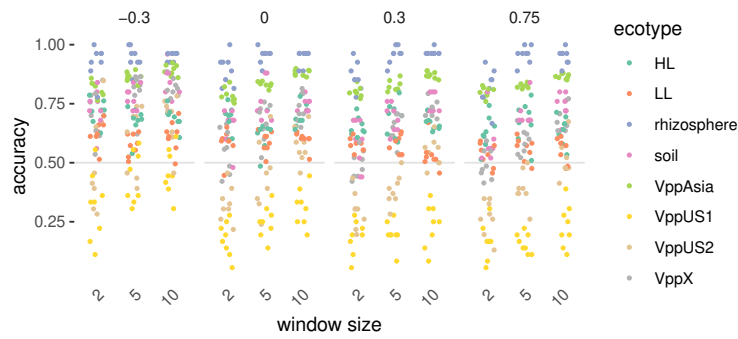


Figure S2: Effect of window size parameter on ecotype task accuracy. Across all negative exponent values for v (facets) larger window (context) size performs better for all ecotype labels.

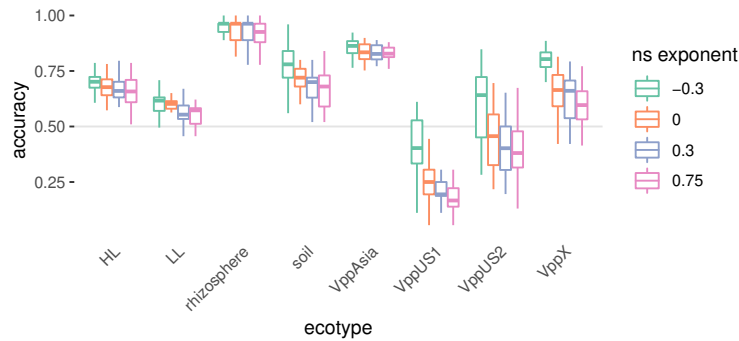


Figure S3: Effect of negative sampling exponent v on ecotype task accuracy. Lower exponents lead to more accurate models on the ecotype task across all ecotype labels.

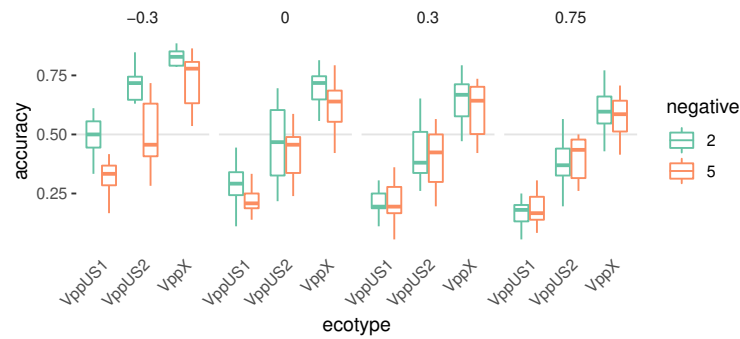


Figure S4: Effect of number of negative samples on ecotype task accuracy. For the hardest to separate and thus most informative ecotypes (x-axis) a smaller count increases accuracy for smaller values of the negative sampling exponent v (facets). For $v \geq 0.3$, the effect is negligible.

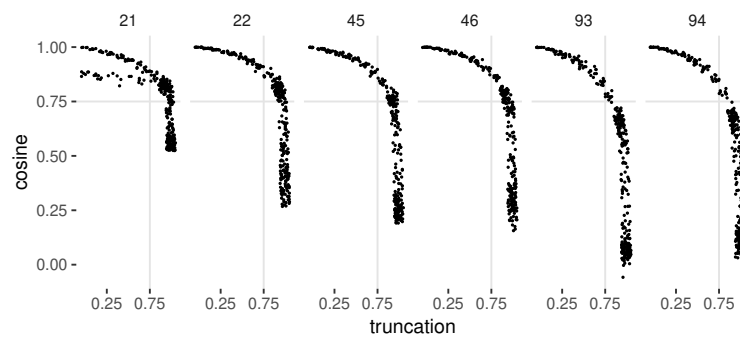


Figure S5: Simulation of genome incompleteness for a single genome (uncharacterized isolate from biogas, unpublished). The genome was truncated in 1% increments (10 repetitions). Each time a genome vector was inferred with a representative model from the hyperparameter grid search (facets). Cosine similarity was calculated against the complete genome. For all training parameters, see supplementary table 1. For accessory models (21, 22), genome vector inference was not stable under truncation. This instability is not observed for models with larger values of the negative sample exponent parameter (45, 46, 93 and 94).

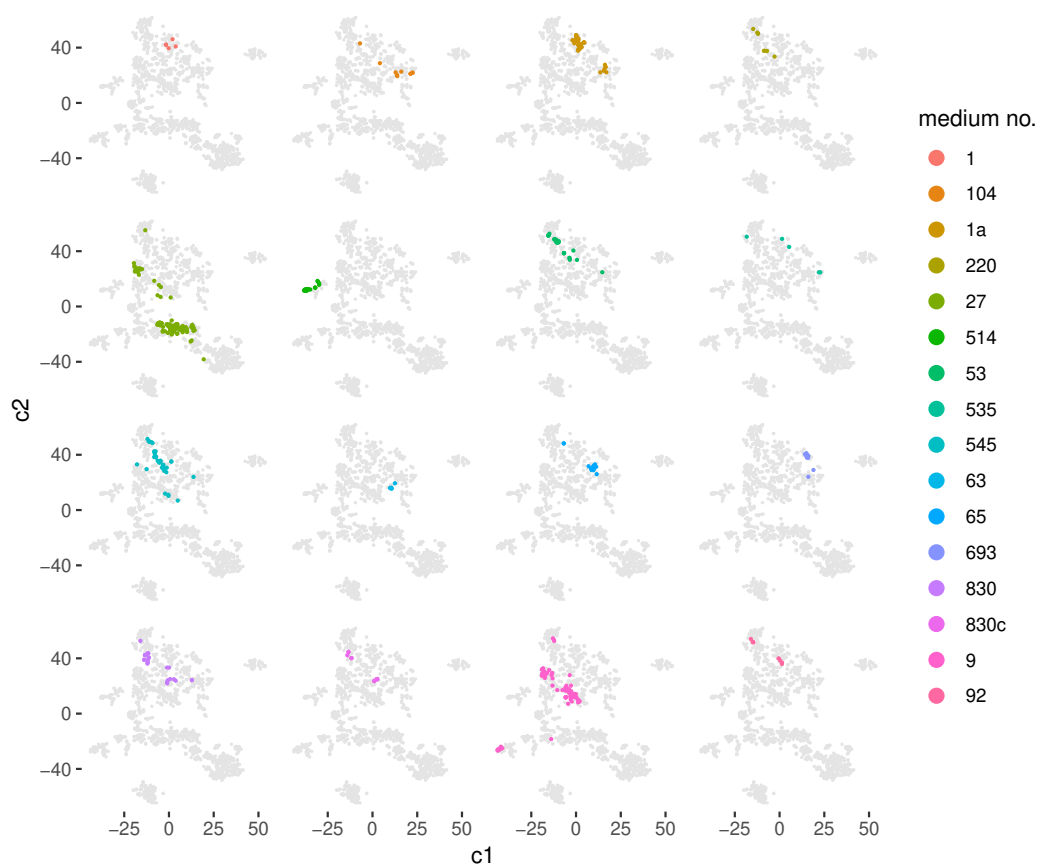


Figure S6: Visualisation of medium embedding space. We used t-SNE to project medium vectors into the plane (grey points). All media with more than 0.95 cosine similarity to any of the top 16 most common DSMZ media were colored. We observe clear clusters of similar media. These clusters can be used by learning algorithms to discriminate media classes. Note how near-identical media such as no. 830 and no. 830c are embedded in near identical vector space.

411 References

- 412 **1** Stadler, P. F. *et al.* *Theory Biosci.* 128, 165–170 (2009)
- 413 **2** Blei, D. M. *Commun. ACM* 55, 77–84 (2012)
- 414 **3** Parks, D. H. *et al.* *Nat Microbiol* 2, 1533–1542 (2017)
- 415 **4** Delmont, T. O. *et al.* *Nat Microbiol* 3, 804–813 (2018)
- 416 **5** Stewart, R. D. *et al.* *bioRxiv* 489443 (2018)
- 417 **6** Krizhevsky, A. *et al.* in *Advances in neural information processing systems 25* (eds. Pereira, F.,
418 Burges, C. J. C., Bottou, L. & Weinberger, K. Q.) 1097–1105 (Curran Associates, Inc., 2012)
- 419 **7** Doolittle, W. F. *et al.* *Biol. Philos.* 32, 5–24 (2017)
- 420 **8** Franzosa, E. A. *et al.* *Nat Microbiol* 4, 293–305 (2019)
- 421 **9** Vogel, C. *et al.* *Curr. Opin. Struct. Biol.* 14, 208–216 (2004)
- 422 **10** Marsh, J. A. *et al.* *Genome Biol.* 11, 126 (2010)
- 423 **11** Illergård, K. *et al.* *Proteins* 77, 499–508 (2009)
- 424 **12** Mendler, K. *et al.* *bioRxiv* 463455 (2018)
- 425 **13** Zhu, C. *et al.* *Nucleic Acids Res.* 46, D535–D541 (2018)
- 426 **14** Delmont, T. O. *et al.* *PeerJ* 6, e4320 (2018)
- 427 **15** Ochoa, A. *et al.* *Bioinformatics* 33, 2471–2478 (2017)
- 428 **16** Gordon, S. P. *et al.* *PLoS One* 10, e0132628 (2015)
- 429 **17** Blin, K. *et al.* *Nucleic Acids Res.* 45, W36–W41 (2017)
- 430 **18** Burkhardt, D. H. *et al.* *Elife* 6, (2017)
- 431 **19** Finn, R. D. *et al.* *Nucleic Acids Res.* 44, D279–85 (2016)
- 432 **20** Hinton, G. E. *et al.* in (eds. Rumelhart, D. E., McClelland, J. L. & PDP Research Group, C.)
433 77–109 (MIT Press, 1986)
- 434 **21** Rudolph, M. R. *et al.* (2016)

- 435 **22** Mikolov, T. *et al.* in *Advances in neural information processing systems 26* (eds. Burges, C. J. C.,
436 Bottou, L., Welling, M., Ghahramani, Z. & Weinberger, K. Q.) 3111–3119 (Curran Associates, Inc.,
437 2013)
- 438 **23** Le, Q. V. *et al.* (2014)
- 439 **24** Asgari, E. *et al.* *PLoS One* 10, e0141287 (2015)
- 440 **25** Yang, K. K. *et al.* *Bioinformatics* 34, 2642–2648 (2018)
- 441 **26** Parks, D. H. *et al.* *Nat. Biotechnol.* 36, 996–1004 (2018)
- 442 **27** Bakarov, A. (2018)
- 443 **28** Conneau, A. *et al.* (2018)
- 444 **29** Moore, L. R. *et al.* *Nature* 393, 464–467 (1998)
- 445 **30** Rocap, G. *et al.* *Nature* 424, 1042–1047 (2003)
- 446 **31** Lopes, L. D. *et al.* *Environ. Microbiol.* 20, 4401–4414 (2018)
- 447 **32** Cui, Y. *et al.* *bioRxiv* (2018)
- 448 **33** McInnes, L. *et al.* (2018)
- 449 **34** Cohan, F. M. *Annu. Rev. Microbiol.* 56, 457–487 (2002)
- 450 **35** Cohan, F. M. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 361, 1985–1996 (2006)
- 451 **36** Koepfel, A. *et al.* *Proc. Natl. Acad. Sci. U. S. A.* 105, 2504–2509 (2008)
- 452 **37** Wright, E. S. *et al.* *BMC Genomics* 19, 724 (2018)
- 453 **38** Snel, B. *et al.* *Nat. Genet.* 21, 108–110 (1999)
- 454 **39** Snel, B. *et al.* *Annu. Rev. Microbiol.* 59, 191–209 (2005)
- 455 **40** Royalty, T. *et al.* *bioRxiv* 520973 (2019)
- 456 **41** Sczyrba, A. *et al.* *Nat. Methods* 14, 1063–1071 (2017)
- 457 **42** Brown, C. T. *et al.* *The Journal of Open Source Software* (2016)
- 458 **43** Wood, D. E. *et al.* *Genome Biol.* 15, R46 (2014)

- 459 **44** Ondov, B. D. *et al. Genome Biol.* 17, 132 (2016)
- 460 **45** Jain, C. *et al. Nat. Commun.* 9, 5114 (2018)
- 461 **46** DeBruyn, J. M. *et al. Appl. Environ. Microbiol.* 77, 6295–6300 (2011)
- 462 **47** Zeng, Y. *et al. Proc. Natl. Acad. Sci. U. S. A.* 111, 7795–7800 (2014)
- 463 **48** Dachev, M. *et al. PLoS Biol.* 15, e2003943 (2017)
- 464 **49** Fischer, W. W. *et al. Annu. Rev. Earth Planet. Sci.* 44, 647–683 (2016)
- 465 **50** Browne, H. P. *et al. Nature* 533, 543–546 (2016)
- 466 **51** Tansey, W. *et al.* (2016)
- 467 **52** Oberhardt, M. A. *et al. Nat. Commun.* 6, 8493 (2015)
- 468 **53** Reimer, L. C. *et al. Nucleic Acids Res.* (2018)
- 469 **54** Wannicke, N. *et al. FEMS Microbiol. Ecol.* 91, (2015)
- 470 **55** Sandle, T. *PDA J. Pharm. Sci. Technol.* 58, 231–237 (2004)
- 471 **56** Moore, L. R. *et al. Limnol. Oceanogr. Methods* 5, 353–362 (2007)
- 472 **57** Waters, L. S. *et al. Cell* 136, 615–628 (2009)
- 473 **58** Land, M. *et al. Funct. Integr. Genomics* 15, 141–161 (2015)
- 474 **59** Eddy, S. R. *PLoS Comput. Biol.* 7, e1002195 (2011)
- 475 **60** Suzek, B. E. *et al. Bioinformatics* 31, 926–932 (2015)
- 476 **61** Steinegger, M. *et al. bioRxiv* 386110 (2018)
- 477 **62** Johnson, J. *et al.* (2017)
- 478 **63** Bojanowski, P. *et al.* (2016)
- 479 **64** McInerney, J. O. *et al. Nat Microbiol* 2, 17040 (2017)
- 480 **65** Langille, M. G. I. *et al. Nat. Biotechnol.* 31, 814–821 (2013)
- 481 **66** Forster, S. C. *et al. Nat. Biotechnol.* 37, 186–192 (2019)

482 **67** Caselles-Dupré, H. *et al.* (2018)

483 **68** Mikolov, T. *et al.* (2013)

484 **69** Mu, J. *et al.* (2017)

485 **70** Hyatt, D. *et al.* *BMC Bioinformatics* 11, 119 (2010)

486 **71** Broder, A. Z. in *Compression and complexity of sequences 1997. Proceedings* 21–29 (IEEE, 1997)

487 **72** Maaten, L. van der. *J. Mach. Learn. Res.* 15, 3221–3245 (2014)

488 **73** Pedregosa, F. *et al.* *J. Mach. Learn. Res.* 12, 2825–2830 (2011)

489 **74** McInnes, L. *et al.* (2017)