

1 **PredicTF: a tool to predict bacterial transcription factors in complex**  
2 **microbial communities.**

3  
4 Lummy Maria Oliveira Monteiro<sup>1,2,3</sup>, Joao Saraiva<sup>1</sup>, Rodolfo Brizola Toscan<sup>1</sup>, Peter F  
5 Stadler<sup>2</sup>, Rafael Silva-Rocha<sup>3</sup>, Ulisses Nunes da Rocha<sup>1\*</sup>

6  
7 <sup>1</sup> Helmholtz Center for Environmental Research (UFZ), Leipzig, Germany

8 <sup>2</sup> Bioinformatics Group, Institute of Computer Science, Universität Leipzig, Leipzig, Germany

9 <sup>3</sup> Ribeirão Preto Medical School (FMRP), University of São Paulo (USP), Ribeirão Preto, Brazil

10 \_\_\_\_\_  
11  
12 \*Correspondence: Ulisses Nunes da Rocha, [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de)

13  
14 Authors' email addresses:

15 Lummy Maria Oliveira Monteiro - [lummymaria@gmail.com](mailto:lummymaria@gmail.com)

16 João Saraiva - [joao.saraiva@ufz.de](mailto:joao.saraiva@ufz.de)

17 Rodolfo Brizola Toscan - [rodolfo.toscan@ufz.de](mailto:rodolfo.toscan@ufz.de)

18 Peter F. Stadler - [studla@bioinf.uni-leipzig.de](mailto:studla@bioinf.uni-leipzig.de)

19 Rafael Silva-Rocha - [silvarochar@usp.br](mailto:silvarochar@usp.br)

20 Ulisses Nunes da Rocha – [ulisses.rocha@ufz.de](mailto:ulisses.rocha@ufz.de)

21

22

23

24

25

26

## 27 **Abstract**

28 Transcription Factors (TFs) are proteins that control the flow of genetic information by  
29 regulating cellular gene expression. Here we describe PredicTF, a first platform  
30 supporting the prediction and classification of novel bacterial TF in complex microbial  
31 communities. We evaluated PredicTF using a two-step approach. First, we tested  
32 PredicTF's ability to predict TFs for the genome of an environmental isolate. In the  
33 second evaluation step, PredicTF was used to predict TFs in a metagenome and 11  
34 metatranscriptomes recovered from a community performing anaerobic ammonium  
35 oxidation (anammox) in a bioreactor. PredicTF is open source pipeline available at  
36 <https://github.com/mdsufz/PredicTF>.

37

38 **Keywords:** Gene regulation, Transcription factors, Deep Learning, Transcription factor  
39 database, Microbial Communities

40

## 41 **Background**

42 The functional potential of microbial communities can be determined by the  
43 genetic content of its constituent members. However, genetic content alone does not  
44 guarantee that a given function or enzymatic reaction will be performed [1]. In this  
45 scenario, Transcription Factor proteins (TFs) play a central and critical role in gene  
46 regulation. These proteins are responsible for optimizing proteins and structural RNAs  
47 and the subsequent levels of metabolites and other properties, ensuring the survival and  
48 adaptation of organisms to the most diverse types of stress and environmental changes  
49 [2]. The activity of bacterial TFs is modulated by environmental signals (e.g. changes in  
50 the oxygen condition, temperature, pH or the lack of a specific substrate) [3].

51 Additionally, for many promoters, combinations of transcription factors work together

52 to integrate different signals [2,4]. TFs can also work with other DNA-binding proteins  
53 whose primary role is to sculpt the bacterial folded chromosome [2,5]. Knowledge of  
54 the TFs profile expressed by an organism is the first step to better understand the  
55 regulatory network that controls protein expression in an organism or a community.

56 Since TFs may determine when and which genes are expressed, profiling TFs  
57 can help understand the regulation of gene expression and to build regulatory networks  
58 in complex microbial communities. Further, defining which factors control gene  
59 expression may offer insights into the mechanisms controlling ecosystem processes and  
60 even interactions between species of a microbial community. However, current TF  
61 databases are focused on single or small groups of genomes. These databases are largely  
62 manually curated based on literature evidence and pairwise sequence comparison of  
63 genomes from model organisms. Examples of these databases include RegulonDB for  
64 *Escherichia coli* K-12 [6], DBTBS for *Bacillus subtilis* [7], FlyBase for *Drosophila* [8],  
65 and FTFD for fungal species [9]. DBD [10], is a database generated from the prediction  
66 of TFs from 150 sequenced genomes from across the tree of life. Unfortunately, DBD  
67 has not been updated for more than 9 years.

68 One of the major goals in the manipulation of microbiomes for ecological and  
69 biotechnological applications is to control the outcome of their functions [11]. As TFs  
70 are key to potentially control which genes are expressed, one of the best ways to study  
71 and understand gene regulation in a microbiome may be to profile its TFs. To date, no  
72 platform supports prediction and classification of novel bacterial TF from 'omics data  
73 recovered from microbial communities.

74 Deep Learning approaches have been used to predict DNA sequence affinities  
75 [12] and to identify TF-binding sites in humans [13]. Although deep learning has been  
76 used in gene regulation, it has never been used to predict bacterial TFs. Further, the  
77 need for a user-friendly tool for prediction of TFs that could assist in gene regulation

78 analysis motivated the development of PredicTF. PredicTF is a deep learning tool used  
79 to predict and identify TFs from full protein-length sequences. Further, we constructed a  
80 robust database for bacterial transcriptional factors (BacTFDB) that was used to train  
81 our Deep Learning model.

82

## 83 **Results and Discussion**

84 PredictTF is a command line software for prediction of novel transcription  
85 factors from genomic and metagenomic data. We created a bacterial transcription factor  
86 database (BacTFDB) by merging and manually curating TFs present in CollectTF [14]  
87 and the Universal Protein Resource (UniProt) [15]. CollectTF provides well described  
88 and characterized, *in vivo* validated, TFs while UniProt is a comprehensive resource for  
89 protein sequence and annotation data. We used BacTFDB to train a deep learning model  
90 to predict new TFs and their families in genomes and metagenomes. Five model  
91 organisms (*Escherichia coli*, *Bacillus subtilis*, *Pseudomonas fluorescens*, *Azotobacter*  
92 *vinelandii* and *Caulobacter crescentus*) were used to test the performance and accuracy  
93 of PredicTF. We used the same approach to predict TFs from a clinical isolate (*P.*  
94 *aeruginosa* PAO1) and a metagenome sample isolated from an anaerobic ammonium  
95 oxidation community. We also determined if the predicted TFs were expressed in  
96 transcriptomes (isolate) and metatranscriptomes (microbial community), respectively  
97 (Fig. 1).

98

## 99 **Database**

100 BacTFDB is a robust and versatile bacterial TF database, it contains 11,691 TFs  
101 amino acid sequences spanning 1049 TF families and 720 different bacterial species.  
102 Fig. 2 shows the database distribution based on TF families and regulatory elements

103 (Fig. 2A) and the distribution based on bacterial species (Fig. 2B). Although BacTFDB  
104 is composed by 11.961 TFs elements from 1049 different families and 720 organism's  
105 species, Fig. 2 shows TFs families and species that accumulate more than 50 sequences.  
106 We will update BacTFDB annually by adding novel entries deposited in UniProt and  
107 CollecTF. BacTFDB was used in PredicTF's deep learning model training. This model  
108 was later used to predict new TFs and their families in genomes and metagenomes.

109

### 110 **Performance and Accuracy**

111 The performance and accuracy of PredicTF were evaluated through the  
112 prediction of TFs in five model organisms (*E. coli*, *B. subtilis*, *P. fluorescens*, *A.*  
113 *vinelandii* and *C. crescentus*). For each model organism a different PredicTF model was  
114 trained to predict TFs from full protein-length sequences (described in the  
115 implementation section).

116 The performance of PredicTF to identify TFs in the different model organisms  
117 ranged from 27% to 60% of the proteins described as TFs in the genomes of model  
118 organisms and the accuracy for experimentally validated TFs ranged from 73.91% and  
119 91.43% (Table 1). Further, PredicTF was able to identify putative annotated TFs in the  
120 genomes of *E. coli* and *B. subtilis* with accuracies 85.71% and 100%, respectively  
121 (Table 1). No novel TF was predicted in the genome of *C. crescentus*, *P. fluorescens*  
122 and *A. vinelandii* (Table 1). TFs predicted by PredicTF for each organism, sorted by TF  
123 family, are shown in Fig. 3. For all organisms tested the most predicted TF family was  
124 LysR followed by OmpR/PhoB. The degree of accuracy obtained by PredicTF suggests  
125 that the deep learning strategy used is promising for the prediction of TFs in genomic or  
126 metagenomic data of bacterial species. PredicTF performance and accuracy can be  
127 further improved by expanding the number and diversity of sequences present in

128 BacTFDB. As BacTFDB will be update yearly, we expect an improvement in TF  
129 identification of with every update.  
130  
131 **Table 1.** PredicTF performance, accuracy for experimentally validated Transcription  
132 Factors (Accuracy EV) and accuracy for putative Transcription Factors (Accuracy PU)  
133 in genomes of model organisms.

Organism	Performance <sup>a</sup>	Accuracy EV <sup>b</sup>	Accuracy PU <sup>c</sup>
<i>E. coli k12</i>	35.40%	88.51%	85.71%
<i>B. subtilis</i>	27.23%	73.91%	100%
<i>C. crescentus</i>	38.04%	83.93%	- <sup>d</sup>
<i>P. fluorescens</i>	51.19%	91.43%	-
<i>A. vinelandii</i>	60.53%	90.40%	-

134 <sup>a</sup>Performance was calculated by the ratio of the total number of TFs predicted by PredicTF (*Predicted TFs*) to the total number of  
135 proteins annotated as TFs in NCBI (*Annotated TFs*) multiplied by 100;  
136 <sup>b</sup>Accuracy EV was determined by the ratio of the total number of TFs predicted by PredicTF in agreement with NCBI annotation  
137 (*TFs predicted correctly*) to the total number of TFs predicted by PredicTF (*TFs predicted*) multiplied by 100;  
138 <sup>c</sup>Accuracy TU was determined by the total number of putative TFs predicted correctly divided by putative TFs predicted multiplied  
139 by 100; *Putative TFs predicted correctly* is the total number of putative TFs predicted correctly by PredicTF in agreement with  
140 NCBI annotation; and, *Putative TFs predicted* is the total number of putative TFs predicted by PredicTF;  
141 <sup>d</sup>Currently there are no putative annotated TFs described in the genome of *C. crescentus*, *P. fluorescens* and *A. vinelandii*  
142

## 143 Mining and Predicting TFs in Genomes and Transcriptomes from a bacterial 144 isolate using PredicTF

145 PredicTF was used to predict TFs on the genome of *P. aeruginosa* PAO1 and  
146 these TFs were mapped in transcriptomes from the same isolate [16]. PredicTF  
147 predicted a total of 199 TFs in the *P. aeruginosa* PAO1 genome shown in Additional  
148 file 1: Fig. S1A by a family's distribution graphic. These 199 TFs were mapped in the  
149 transcriptomic data of a reference of *P. aeruginosa* PAO1. Initially, the mapping was  
150 done in the transcriptome of *P. aeruginosa* PAO1 cultured in LB media. Using this  
151 strategy, we were able to map 69 of the 199 predicted TFs to the transcriptomes under  
152 the experimental conditions carried out by Hwang & Yoon, 2019 (Additional file 1: Fig.  
153 S1B) [16]. Next, the mappings were done for another three clinical mutants of *P.*  
154 *aeruginosa* PAO1 (Y82, Y71, Y89) cultured in LB media (absence of an antibiotic  
155 cocktail) (Additional file 2: Fig. S2A, S2C and S2F). The TFs family's distribution for

156 each *P. aeruginosa* PAO1 mutant cultured in presence of antibiotic cocktail is shown in  
157 the supplementary data (Additional file 2: Fig. S2B, S2D and S2F). These results  
158 demonstrate the potential of PredicTF in mapping regulatory elements in bacterial  
159 genomes and the use of this tool to map and compare TFs profiles after under different  
160 environmental conditions.

161

## 162 **Mining and Predicting TFs in a Metagenome and Metatranscriptome using** 163 **PredicTF**

164 PredicTF was used to profile TFs in one metagenome recovered from an  
165 anaerobic ammonium oxidation community [17] followed by the mapping of the  
166 predicted TFs in metatranscriptomes recovered from the same community  
167 (metatranscriptomes accession numbers can be found in Additional file 3: Table S1). A  
168 total of 792 TFs (Fig. 4A) were predicted in LAC\_MetaG\_1, an anaerobic ammonium  
169 oxidizing microbial community from an anammox membrane bioreactor [17]. These  
170 792 TFs are distributed across 27 TF families (Fig. 4A) and are related to the regulation  
171 of functions such as the oxygen limitation response and late symbiotic functions  
172 (NarL/FixJ), phosphate regulon response (OmpR/PhoB), transcriptional activator for  
173 nitrogen-regulated promoters (NtrC/DctD) and ferric uptake regulation (Fur). To  
174 determine how a traditional annotation pipeline identify potential TF we used Prokka  
175 [18]. This tool was able to identify 1815 ORFs (Additional file 4: Table S2). PredicTF  
176 can be used with no previous knowledge regarding transcription factors, it is fast and it  
177 requires low memory when compared to Blast based annotation and it indicates only  
178 results of TFs with a specific TF family annotation. On the other hand, to identify TFs  
179 using Prokka one would need specialized training to mine the general annotation.  
180 Therefore, scientists with general microbiology background may take a long time to  
181 undergo this task. Further, Prokka gives no indication to the TF families of the putative

182 annotated TFs. Time is also a drawback of using Prokka to mine TFs, we calculated we  
183 needed over 400 h to perform mine one single metagenomics library; in comparison,  
184 PredicTF needed 2 h to identify TF in the same metagenomics library.

185         Next, the 792 TFs were mapped in 11 metatranscriptomes collected in different  
186 dates from the same bioreactor where the metagenome was recovered (Additional file 5:  
187 Table S3, Fig. 4B). Clustering analysis demonstrated the presence of five different  
188 groups of TFs families based on the number of transcription factor families expressed in  
189 each library (Fig. 4B). It is interesting to note that the two most abundant clusters in the  
190 heatmap are directly related to the oxygen limitation caused by the anaerobic  
191 ammonium oxidizing cultivation. In a bioreactor where oxygen is limited, an increase in  
192 the amount of nitrogen and phosphate is expected. The presence of N and P diverts the  
193 metabolism of the microbial community towards the production of regulators (TFs) that  
194 help to maintain community stability. Clustering analyzes can be helpful to demonstrate  
195 the similarity between metatranscriptomic libraries based on the occurrence of TFs. This  
196 strategy can be useful to compare the profiles of TFs expressed in different  
197 environmental situations (comparing libraries with different metadata) creating patterns  
198 of TFs expression. Exploration of TF profiling in microbial communities (metagenomes  
199 or metatranscriptomes) will allow the exploration of regulation within complex  
200 microbial communities. Further, The recovery of metagenome assembled genomes is  
201 becoming standard in metagenomics studies [19–21]. The use of PredicTF together with  
202 the recovery of metagenome assembled genomes will allow the exploration of species-  
203 specific molecular mechanisms involved in the regulation of different ecosystem  
204 processes.

205

206 **Conclusions**



207           A better understanding of TFs in a bacterial community context open revenue  
208   for the exploration of gene regulation in ecosystems where bacteria play a key role. Our  
209   deep learning strategy was based on a novel and robust TF bacterial database  
210   (BacTFDB) with over 11 thousand TFs and their respective families. BacTFDB is a  
211   unique resource for the exploration of TFs and it provided the data to train a model  
212   within PredicTF capable of predicting novel TFs from genomes and metagenomes.  
213   PredicTF is the first pipeline designed to predict and annotate TFs in complex microbial  
214   communities. The prediction of TFs can provide information for those aiming to study  
215   and understand bacterial communities within a context of gene regulation. We also  
216   demonstrated that PredicTF can be used to predict novel TFs in metagenomes and  
217   metatranscriptomes creating the potential profile for regulatory elements in complex  
218   microbial communities.

219           PredicTF is a flexible open source pipeline able to predict and annotate TFs in  
220   genomes and metagenomes and can be found at <https://github.com/mdsufz/PredicTF>.

221

## 222   **Methods**

### 223   **BacTFDB - Bacterial Transcription Factor Data Base**

224           To create a novel Bacterial Transcription Factor Data Base (BacTFDB), we  
225   collected data from two publicly available databases. Initially, we chose to collect data  
226   from CollecTF [14], a well described and characterized database. Since CollecTF does  
227   not provide an application programming interface (API) for bulk download, we  
228   developed a Python code (version 2.7) using the Beautiful Soup 4.4.0 library to recover  
229   the data from CollecTF. With this strategy we listed 390 TF experimentally validated  
230   amino acid sequences distributed over 44 TF families. The script can be found at  
231   <https://github.com/mdsufz/PredicTF>.

232           Additionally, we retrieved TF amino acid sequences from UniProt using  
233   UniProt's API. We downloaded sequences of interest by adding a filter with the key  
234   words (Transcription factor, transcriptional factor, regulator, transcriptional repressor,  
235   transcriptional activator, transcriptional regulator). After, we filtered for Reviewed  
236   (Swiss-Prot) - Manually annotated sequences that belonged to the bacteria taxonomy.  
237   The UniProt API was accessed on 8<sup>th</sup> September-2019 and a total of 21.581 TF amino  
238   acid sequences, with applied filters, were collected. We merged the data collected from  
239   CollecTF and UniProt databases which resulted in a total of 21.971 TFs. Next, we  
240   removed redundant TF entries and TF sequences lacking a TF family since PredictTF  
241   was designed to also assign TF family. Finally, a manual inspection was performed to  
242   remove case sensitive and presence of characters associated to the database header. The  
243   first version of BacTFDB contains a total of 11.691 unique TF sequences. A summary  
244   of the information contained in BacTFDB can be found in the supplementary data  
245   (Additional file 6: Fig. S4). To evaluate PredictTF in model organisms we created 5  
246   subsets of BacTFDB. The description of these subsets can be found in the  
247   supplementary data (Additional file 7: Table S4).

248

### 249   **Mapping Transcription Factors using PredictTF**

250           We used a deep learning approach similar to that found in DeepARG [22].  
251   Supervised machine learning models are usually divided into characterization, training,  
252   and prediction units. Briefly, our approach uses the concept of dissimilarity-based  
253   classification [23] where sequences are represented and featured by their sequence  
254   similarity to known genes. BacTFDB was used to train and test the deep learning model  
255   (<https://github.com/mdsufz/PredictTF>) and latter validated in model organisms. Next,  
256   PredictTF was used to predict novel TFs from full protein-length sequences in genomes

257 and in one metagenome. After prediction, the data was mapped in transcriptomes and  
258 metatranscriptomes from samples where the genetic potential was determined.

259 Using PredicTF, we trained five different models – one for each model organism  
260 (Additional file 3: Table S1). For each model, the TFs affiliated with the respective  
261 model organism were removed prior to training to avoid overfitting. PredicTF-no-coli  
262 was trained to predict TFs in *E. coli*, PredicTF-no-subtilis was trained to predict TFs in  
263 *B. subtilis*, PredicTF-no-crescentus was trained to predict TFs in *C. crescentus*,  
264 PredicTF-no-fluorescens was trained to predict TFs in *P. fluorescens* and PredicTF-no-  
265 vinelandii was trained to predict TFs in *A. vinelandii*.

266

## 267 **Performance and accuracy calculation**

268 We evaluated PredicTF by calculating accuracy and performance. Performance  
269 can be deemed to be the fulfillment of a task. In PredicTF case, performance is how  
270 good TF predictions are. Using model organisms (see later in the session *Prediction of*  
271 *Transcription Factors in model organisms*), performance was calculated by quantifying  
272 the number of TFs that PredicTF was able to predict divided by number of TFs already  
273 described and annotated for our model organisms (Additional file 7: Equation 1).  
274 Accuracy indicates how correct the predictions performed by PredicTF are. Also using  
275 data of model organism, accuracy was determined by calculating the number of TFs  
276 correctly predicted divided by the total number of TFs predicted by PredicTF. We  
277 divided accuracy in two categories. In the first accuracy category, we determined  
278 accuracy against experimentally validated TFs (Additional file 7: Equation 2). In the  
279 second accuracy category, we determined accuracy against TFs without experimental  
280 validation (Additional file 7: Equation 3); i.e., putative TFs. The performance, accuracy,  
281 and accuracy for putative TFs were calculated as the ratio of predicted to annotated TFs.

282 Accuracy was quantified as the fraction of correctly predicted TFs among all  
283 predictions.

284

### 285 **Prediction of Transcription Factors in model organisms**

286 We selected bacterial species that have been widely studied as model organisms.  
287 Some bacterial species became model organisms for TF studies because they are easy to  
288 maintain and grow in a laboratory setting and to manipulate in pure culture experiments.  
289 Five complete genomes from model organisms (*E. coli*, *B. subtilis*, *P. fluorescens*, *A.*  
290 *vinelandii* and *C. crescentus*) were downloaded directly from NCBI. The strains details  
291 and accession number (RefSeq) for all selected organisms are listed in the  
292 supplementary data (Additional file 3: Table S1). By evaluating PredicTF using model  
293 organisms (Additional file 6: Table S3) we extrapolated performance and accuracy of  
294 our deep learn model. Since known TFs for each organism were removed from each the  
295 training dataset, we eliminate the possibility of mapping TFs already known and  
296 annotated for each of the different species. Performance, accuracy and accuracy for  
297 putative TFs of PredicTF for these five model organisms were calculated using  
298 Equations 1, 2 and 3.

299

### 300 **Prediction of Transcription Factors in a clinical isolate**

301 We demonstrated the use of PredicTF in a previously sequenced *P. aeruginosa*  
302 (PAO1) genome, a clinical isolate publicly available in NCBI (accession number  
303 **NC\_002516.2**). *P. aeruginosa* PAO1 was selected because its genome has been  
304 sequenced and because of the availability of transcriptomes from three clinical mutants  
305 of PAO1 (Y71, Y82, and Y89) grown in the presence and absence of an antibiotic  
306 cocktail. The transcriptomes of *P. aeruginosa* PAO1 mutants Y71, Y82, and Y89 are  
307 available in NCBI (Bioproject identifier **PRJNA479711**) [16]. These clinical *P.*

308 *aeruginosa* PAO1 mutants were isolated from the sputa of three different pneumonia  
309 patients. Transcriptomes of *P. aeruginosa* PAO1 wild type and its mutants cultured in  
310 two different conditions (LB medium and LB medium in presence of antibiotic cocktail)  
311 have been previously described [16]. We used this data to determine the TF profile in  
312 these *P. aeruginosa* PAO1 mutants grown in two different conditions.

313         PredicTF was first used to predict TFs in the *P. aeruginosa* PAO1 genome.  
314 Next, the predicted TFs were mapped to the transcriptomes of the *P. aeruginosa* PAO1  
315 mutants Y71, Y82 and Y89 (see later). Further description of the mapping of the  
316 transcriptomes to the genomes is available at <https://github.com/mdsufz/PredicTF>. The  
317 PredicTF model used in this step was trained with the full database BacTFDB. All  
318 accession numbers used in this work are listed in the supplementary data (Additional  
319 file 3: Table S1).

320

### 321 **Prediction of Transcription Factors in Complex Microbial Communities**

322         To test PredicTF in a complex microbial community, we used an anaerobic  
323 ammonium oxidizing (anammox) microbial community from an anammox membrane  
324 bioreactor metagenome (LAC\_MetaG\_1) (data publicly available at NCBI bioproject  
325 via accession number **PRJNA511011**) [17]. We removed short and low-quality reads  
326 using Trim Galore - v0.0.4 dev according developer's instructions [24]. Over 50 million  
327 reads survived this step and were assembled using the *de novo* assembler metaSPADES  
328 - v3.12.0 [25]. The assembly was translated from nucleotide to amino acid sequences,  
329 considering all possible translation frames, using emboss transeq [26]. The translated  
330 assembly was then used as input for the prediction of transcription factors using  
331 PredicTF. The region from each predicted TF was extracted. These putative TFs were  
332 later used in the mapping TFs to metatranscriptomes.

333 We checked if the putative TFs predicted in the metagenomes were transcribed  
334 by checking if the metatranscriptomic libraries were mapping to those regions. The  
335 metatranscriptomic and metagenomic libraries used in this step belonged to the same  
336 bioreactor. These metatranscriptomes are publicly available at the European Nucleotide  
337 Archive under the accession numbers SRR7091385, SRR7523233, SRR7523244,  
338 SRR7523245, SRR7091400, SRR7091401, SRR7091381, SRR7091402, SRR7091406,  
339 SRR7523243, SRR7523246. These 11 metatranscriptomes were used to demonstrate the  
340 effectiveness of the pipeline and to indicate the potential of PredicTF to profile  
341 transcription factors in complex microbial communities. All accession numbers used in  
342 this work are listed in the supplementary data (Additional file 3: Table S1).

343 To have a baseline comparison with a traditional annotation pipeline, we used  
344 Prokka [18] to annotate the same anammox membrane bioreactor metagenome  
345 (LAC\_MetaG\_1). We mined the annotation by hand with specialized knowledge of  
346 scientists specialized in Transcription Factors. We did not determine the families as this  
347 work would need to be done for every single hit individually using the output of Prokka.

348

### 349 **Mapping transcription factors to transcriptomes and metatranscriptomes**

350 Each transcriptomic and metatranscriptomic library was quality controlled by  
351 removing short and low-quality reads using Trim Galore - v0.0.4 dev [24]. The 7  
352 transcriptomic libraries for the *P. aeruginosa* PAO1 wild type and mutants showed at  
353 least 26 million paired end reads after quality checking. The 11 metatranscriptomic  
354 libraries yielded over 50 million reads per library after quality check. After, the  
355 remaining transcriptomic and metatranscriptomic reads were mapped to their respective  
356 assembled genome or metagenome using Bowtie2 - v2.3.0 [27]. The number of reads  
357 mapped, and the regions covered was extracted using SAMTools - v1.9 [28] and  
358 python 2.7. The regions of the genome or metagenome assembly covered by

359 transcriptomic or metatranscriptomic reads were then cross-referenced with the  
360 regions of their respective assembly which PredicTF assigned as putative TFs creating a  
361 TF profile for each transcript and metatranscriptome. A detailed description on the  
362 mapping of RNA-seq data to their respective genome or metagenome assembly can be  
363 found at the PredicTF github (<https://github.com/mdsufz/PredicTF>).

364

### 365 **List of Abbreviations**

366 Transcription Factors (TFs)

367 Bacterial Transcription Factor Data Base (BacTFDB)

368 Transcription factor binding sites (TFBSs)

369 anaerobic ammonium oxidizing (anammox)

370

### 371 **Declaration Sections**

#### 372 **Ethics approval and consent to participate**

373 Not applicable

374

#### 375 **Consent for publication**

376 Not applicable

377

#### 378 **Availability of data and materials**

379 Project name: PredicTF

380 Project home page: <https://github.com/mdsufz/PredicTF>

381 Operating system: Linux64

382 Programming languages: Python 2.7

383 Other requirements: DIAMOND [29]; Nolearn Lasagne deep learning library

384 [30]; Sklearn machine learning routines (<https://scikit-learn.org/stable/>) [31]; Theano

385 (<http://deeplearning.net/software/theano/>) [32]. Trim Galore - v0.0.4 dev  
386 ([https://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)) [24]. MetaSPADES  
387 - v3.12.0 (<https://github.com/ablab/spades#meta>) [25]. Emboss transeq  
388 (<http://www.bioinformatics.nl/cgi-bin/emboss/transeq>) [26]. Bowtie2 - v2.3.0  
389 (<https://sourceforge.net/projects/bowtie-bio/>) [27]. SAMTools - v1.9  
390 (<http://github.com/samtools/>) [28].

391 Genomes of the model organisms used in the Tool Validation step are available  
392 at the National Center for Biotechnology Information (<https://www.ncbi.nlm.nih.gov/>)  
393 under the accession numbers NC\_000913.3, NC\_000964.3, NC\_011916.1,  
394 NC\_021149.1, and NC\_016830. The datasets supporting the Prediction of Transcription  
395 Factors in a clinical isolate of this article are available at National Center for  
396 Biotechnology Information (<https://www.ncbi.nlm.nih.gov/>) under the accession  
397 number NC\_002516.2 (genome) and study accession PRJNA479711 (transcriptomes).  
398 The datasets used for the Prediction of Transcription Factors in Complex Microbial  
399 Communities of this study are available at National Center for Biotechnology  
400 Information (<https://www.ncbi.nlm.nih.gov/>) under the study accession PRJNA511011.  
401 The respective data sets of metatranscriptomes used are available at National Center for  
402 Biotechnology Information (<https://www.ncbi.nlm.nih.gov/>) under the SRA numbers  
403 SRR7091385, SRR7523233, SRR7523244, SRR7523245, SRR7091400, SRR7091401,  
404 SRR7091381, SRR7091402, SRR7091406, SRR7523243, SRR7523246 and the Joint  
405 Genome Institute (<https://jgi.doe.gov/>) under the Gold Analysis Project identifiers  
406 Gp0267156, Gp0267150, Gp0267154, Gp0267155, Gp0267157, Gp0267158,  
407 Gp026715, Gp0267159, Gp0267152, Gp0267153, Gp0267160. All analysis, results and  
408 scripts used to generate figures are available at <https://github.com/mdsufz/PredicTF>.

409

410 **Competing of interests**



411 Not applicable

412

413 **Funding**

414 LMOM were supported by FAPESP PhD (award # 2016/19179-9) and FAPESP

415 Research Internship Scholarship Abroad (award # 2018/21133-2). RSR was supported

416 by FAPESP (award # 2019/15675-0). JS and UNR were supported by the Helmholtz

417 Young Investigator grant VH-NG-1248 Micro ‘Big Data’.

418

419 **Authors’ contributions**

420 LMOM, PFS, RSR, and UNR developed the concept of PredicTF. LMOM, JS, UNR

421 developed the PredicTF workflow. LMOM, JS, and UNR performed the benchmarks.

422 LMOM provided information and data for the creation BacTFDB dataset. RBT and

423 UNR performed the metagenome and metatranscriptome analysis. LMOM and UNR

424 wrote the manuscript. All authors read and approved the manuscript.

425

426 **Acknowledgements**

427

428

429

430

431

432

433

434

435

436

437 **Figure legends:**

438

439 **Fig. 1**

440 **PredicTF workflow and testing.** We collected publicly available data on TFs from two  
441 different databases: CollecTF and UNIPROT. After removing redundancies and  
442 filtering TFs well characterized, this data (BacTFDB) was used to train a deep learning  
443 model to predict new TFs and their families. Five model organisms (*Escherichia coli*,  
444 *Bacillus subtilis*, *Pseudomonas fluorescens*, *Azotobacter vinelandii* and *Caulobacter*  
445 *crescentus*) were used to test the accuracy of PredicTF. Later, we used the same  
446 approach to predict TFs from an isolate (*P. aeruginosa*) and mapped TFs predicted in  
447 transcriptomics data (*P. aeruginosa* and mutants in two experimental conditions).  
448 Finally, we used our tool to predict TF for complex communities (metagenome) and  
449 mapped these TFs in their respective meta-transcriptomes.

450

451 **Fig. 2**

452 **Database composition: Transcription Factor Database (BacTFDB) distribution. A)**  
453 Database distribution based on the TFs and **B)** Regulatory Elements families and  
454 Organisms species. In these graphics only families with up to 50 sequences and only  
455 organisms that contributed with more than 50 sequences are shown.

456

457 **Fig. 3**

458 **Prediction of TFs by PredicTF for genomes of model organisms.** Prediction of TFs  
459 or 5 model organisms sorted by family. **A)** *Escherichia coli* **B)** *Bacillus subtilis* **C)**  
460 *Caulobacter crescentus* **D)** *Pseudomonas fluorescens* **E)** *Azotobacter vinelandii*

461

462

463 **Fig. 4**

464 **Recovery of novel Transcription Factors in one metagenome and eleven**

465 **metatranscriptomes. A)** PredicTF predicted 792 TFs were predicted in one anaerobic

466 ammonium oxidizing microbial communities from anammox membrane bioreactor

467 (LAC\_MetaG\_1) and were grouped by family. **B)** Using 792 TFs predicted in one

468 metagenome, we mapped these TFs for 11 metatranscriptomes of reference from the

469 same bioreactor where the metagenome was recovered.

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489 **Additional files**

490 **Additional file 1: Fig. S1**

491 Transcription factor (TF) families predicted for *Pseudomonas aeruginosa* PAO1  
492 genome (accession number NC\_002516.2) [18] using PredicTF and their mapping to *P.*  
493 *aeruginosa* PAO1 growing in LB medium. A) A total of 199 TFs distributed in 25 TF  
494 families were predicted in the *P. aeruginosa* PAO1 genome. B) These 199 TFs were  
495 mapped in the transcriptomic data of a reference of *P. aeruginosa* PAO1 (Bioproject  
496 identifier PRJNA479711) [18]. Initially, the mapping was done in the transcriptome of  
497 *P. aeruginosa* PAO1 cultured in LB media. Using this strategy, we were able to map 69  
498 of the 199 predicted TFs to the transcriptome.

499

500 **Additional file 2: Fig. S2**

501 Transcription Factor (TF) family profiles in three *Pseudomonas aeruginosa* PAO1  
502 mutants. After the prediction of Transcription Factors (TFs) using *P. aeruginosa* PAO1  
503 genome, we mapped transcriptomes from three *P. aeruginosa* PAO1 mutants (Y82,  
504 Y71, Y89) cultured in LB media (A, C and F). After, the mapping was done for each *P.*  
505 *aeruginosa* PAO1 mutant cultured in presence of antibiotic cocktail (B, D and E). *P.*  
506 *aeruginosa* PAO1 mutant Y82 (A, B); *P. aeruginosa* PAO1 mutant Y71 (C, D); *P.*  
507 *aeruginosa* PAO1 mutant Y89 (E, F).

508

509 **Additional file 3: Table S1**

510 Accession number for 5 model organisms, *Pseudomonas aeruginosa* PAO1 genome and  
511 transcriptomes and Complex Microbial Communities used to validate and test PredicTF.

512

513

514

515 **Additional file 4: Table S2**

516 Transcription factors from the metagenome of an anaerobic ammonium oxidizing  
517 microbial community from an anammox membrane bioreactor (LAC\_MetaG\_1) mined  
518 and hand curated from a general annotation generated using Prokka.

519

520 **Additional file 5: Table S3**

521 Number of Transcription Factors (TFs) per TF family mapped to each of the 11  
522 metatranscriptomes of reference from the same bioreactor where the metagenome  
523 (accession number PRJNA511011, NCBI) used to predict the putative TFs was  
524 collected. The different metatranscriptomes are represented by their European  
525 Nucleotide Archive accession numbers.

526

527 **Additional file 6: Fig. S4**

528 Bacterial Transcription Factor Data Base (BacTFDB) were created from from two  
529 publicly available databases. We collect 390 TFs from CollecTF and 21.581 from  
530 UniProt (accessed 8-Sep-2019) accumulating 21.581 Transcription Factor (TF) amino  
531 acid sequences. We merged the data from CollecTF and UniProt databases resulting in a  
532 total of 21.971 TFs amino acid. We removed redundant TF entries and since PredicTF  
533 was designed to also assign TF family, TF sequences lacking a TF family were  
534 removed. Finally, a manual inspection was performed to remove misleading of spelling,  
535 case sensitive and presence of characters associate to the database header. The final  
536 database (BacTFDB) contains a total of 11.691 TF unique sequences.

537

538 **Additional file 7: Table S4**

539 Description of the bacterial transcriptional factors database (BacTFDB) subsets used to  
540 train models to predict Transcription Factors in model organisms

541 **Additional file 8**

542 Equations used to calculate PredicTF's accuracy and performance.

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

## 567 **References**

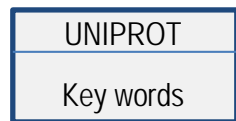
- 568 1. Liu J, Meng Z, Liu X, Zhang XH. Microbial assembly, interaction, functioning,  
569 activity and diversification: a review derived from community compositional data. *Mar*  
570 *Life Sci Technol.* 2019;1:112–28.
- 571 2. Browning DF, Busby SJW. The regulation of bacterial transcription initiation. *Nat*  
572 *Rev Microbiol.* 2004;2:57–65. Nature Publishing Group.
- 573 3. Browning DF, Butala M, Busby SJW. Bacterial Transcription Factors: Regulation by  
574 Pick “N” Mix. *J. Mol. Biol. Academic Press.* 2019;4067–77.
- 575 4. Browning DF, Busby SJW. Local and global regulation of transcription initiation in  
576 bacteria. *Nat Rev Microbiol.* 2016;14:638–50. Nature Publishing Group.
- 577 5. Browning DF, Grainger DC, Busby SJ. Effects of nucleoid-associated proteins on  
578 bacterial chromosome structure and gene expression. *Curr Opin Microbiol.*  
579 2010;13:773–80.
- 580 6. Gama-Castro S, Salgado H, Santos-Zavaleta A, Ledezma-Tejeida D, Muñiz-Rascado  
581 L, García-Sotelo JS, et al. RegulonDB version 9.0: high-level integration of gene  
582 regulation, coexpression, motif clustering and beyond. *Nucleic Acids Res.*  
583 2016;44:D133–43.
- 584 7. Siervo N, Makita Y, De hoon M, Nakai K. DBTBS: A database of transcriptional  
585 regulation in *Bacillus subtilis* containing upstream intergenic conservation information.  
586 *Nucleic Acids Res.* 2008;36.
- 587 8. Gelbart, W. M., Crosby, M., Matthews, B., Rindone, W. P., Chillemi, J., Twombly,  
588 S. R., et al. FlyBase: a *Drosophila* database. Flybase Consortium. *Nucleic Acids Res.*  
589 1998;26:85–8.
- 590 9. Park J, Park J, Jang S, Kim S, Kong S, Choi J, et al. FTFD: An informatics pipeline  
591 supporting phylogenomic analysis of fungal transcription factors. *Bioinformatics.*  
592 2008;24:1024–5.
- 593 10. Kummerfeld SK. DBD: a transcription factor prediction database. *Nucleic Acids*  
594 *Res.* 2006;34:D74–81.
- 595 11. Widder S, Allen RJ, Pfeiffer T, Curtis TP, Wiuf C, Sloan WT, et al. Challenges in  
596 microbial ecology: Building predictive understanding of community function and  
597 dynamics. *ISME J. Springer Nature.* 2016;2557–68. Nature Publishing Group.
- 598 12. Alipanahi B, Delong A, Weirauch MT, Frey BJ. Predicting the sequence  
599 specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol.*  
600 2015;33:831–8. Nature Publishing Group.
- 601 13. Pan X, Shen H Bin. RNA-protein binding motifs mining with a new hybrid deep  
602 learning based cross-domain knowledge integration approach. *BMC Bioinformatics.*  
603 BioMed Central Ltd. 2017;18.
- 604 14. Kiliç S, White ER, Sagitova DM, Cornish JP, Erill I. CollecTF: A database of

- 605 experimentally validated transcription factor-binding sites in Bacteria. *Nucleic Acids*  
606 *Res.* 2014;42.
- 607 15. UniProt Consortium T. UniProt: the universal protein knowledgebase. *Nucleic*  
608 *Acids Res.* 2018;46:2699–2699.
- 609 16. Hwang W, Yoon SS. Virulence Characteristics and an Action Mode of Antibiotic  
610 Resistance in Multidrug-Resistant *Pseudomonas aeruginosa*. *Sci Rep.* 2019;9.
- 611 17. Keren R, Lawrence JE, Zhuang W, Jenkins D, Banfield JF, Alvarez-Cohen L, et al.  
612 Increased replication of dissimilatory nitrate-reducing bacteria leads to decreased  
613 anammox bioreactor performance. *Microbiome.* 2020;8. BioMed Central Ltd.
- 614 18. Seemann T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics.*  
615 2014;30:2068–9. Oxford University Press.
- 616 19. Parks DH, Rinke C, Chuvochina M, Chaumeil PA, Woodcroft BJ, Evans PN, et al.  
617 Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the  
618 tree of life. *Nat Microbiol.* 2017;2:1533–42. Nature Publishing Group.
- 619 20. Pasolli E, Asnicar F, Manara S, Zolfo M, Karcher N, Armanini F, et al. Extensive  
620 Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from  
621 Metagenomes Spanning Age, Geography, and Lifestyle. *Cell.* 2019;176:649-662.e20.
- 622 21. Tully BJ, Graham ED, Heidelberg JF. The reconstruction of 2,631 draft  
623 metagenome-assembled genomes from the global oceans. *Sci Data.* 2018;5.
- 624 22. Arango-Argoty G, Garner E, Pruden A, Heath LS, Vikesland P, Zhang L.  
625 DeepARG: A deep learning approach for predicting antibiotic resistance genes from  
626 metagenomic data. *Microbiome.* 2018;6. BioMed Central Ltd.
- 627 23. Sørensen L, Loog M, Lo P, Ashraf H, Dirksen A, Duin RPW, et al. Image  
628 dissimilarity-based quantification of lung disease from CT. *Lect Notes Comput Sci.*  
629 2010;37–44.
- 630 24. Krueger F. Babraham Bioinformatics - Trim Galore!. Version 0.5.0. 2018. Available  
631 from: [https://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)
- 632 25. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. MetaSPAdes: A new versatile  
633 metagenomic assembler. *Genome Res.* 2017;27:824–34.
- 634 26. Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, et al. The EMBL-  
635 EBI search and sequence analysis tools APIs in 2019. *Nucleic Acids Res.*  
636 2019;47:W636–41.
- 637 27. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat*  
638 *Methods.* 2012;9:357–9.
- 639 28. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
640 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25:2078–9.
- 641 29. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using  
642 DIAMOND. *Nat. Methods.* 2014;59–60. Nature Publishing Group.

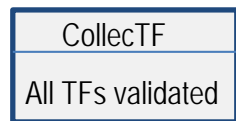


- 643 30. van Merriënboer B, Bahdanau D, Dumoulin V, Serdyuk D, Warde-Farley D,  
644 Chorowski J, et al. Blocks and Fuel: Frameworks for deep learning. arxiv.org. 2015.
- 645 31. Pedregosa F, Michel V, Grisel O, Blondel M, Prettenhofer P, Weiss R, et al. Scikit-  
646 learn: Machine Learning in Python. J. Mach. Learn. Res. 2011; 2825-2830.
- 647 32. Al-Rfou R, Alain G, Almahairi A, Angermueller C, Bahdanau D, et al. Theano: A  
648 Python framework for fast computation of mathematical expressions. arXiv. 2016;1605.
- 649  
650



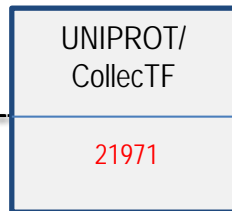


21581



390

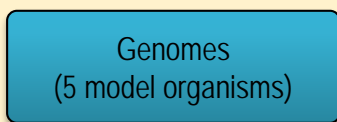
Merge



Remove  
redundancy  
Remove TFs  
without FAMILY

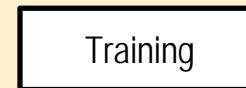
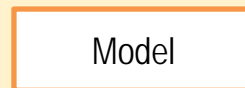
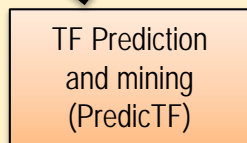


Manual  
inspection  
+  
Handle the  
headers\* to use  
Deep Learning



- Validation for model organisms
  - Statistical analysis
  - Performance

*Escherichia coli*  
*Bacillus subtilis*  
*Pseudomonas fluorescens*  
*Azotobacter vinelandii*  
*Caulobacter crescentus*



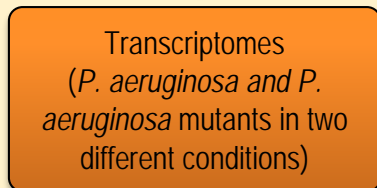
- PredicTF for Genomes and Transcriptomes from isolates for an Isolate

199 TFs predicted



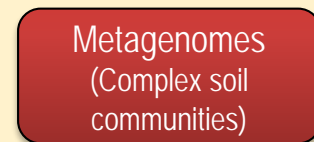
mapping

70 TFs mapped



- PredicTF for Metagenomes and Metatranscriptomes

792 TFs predicted

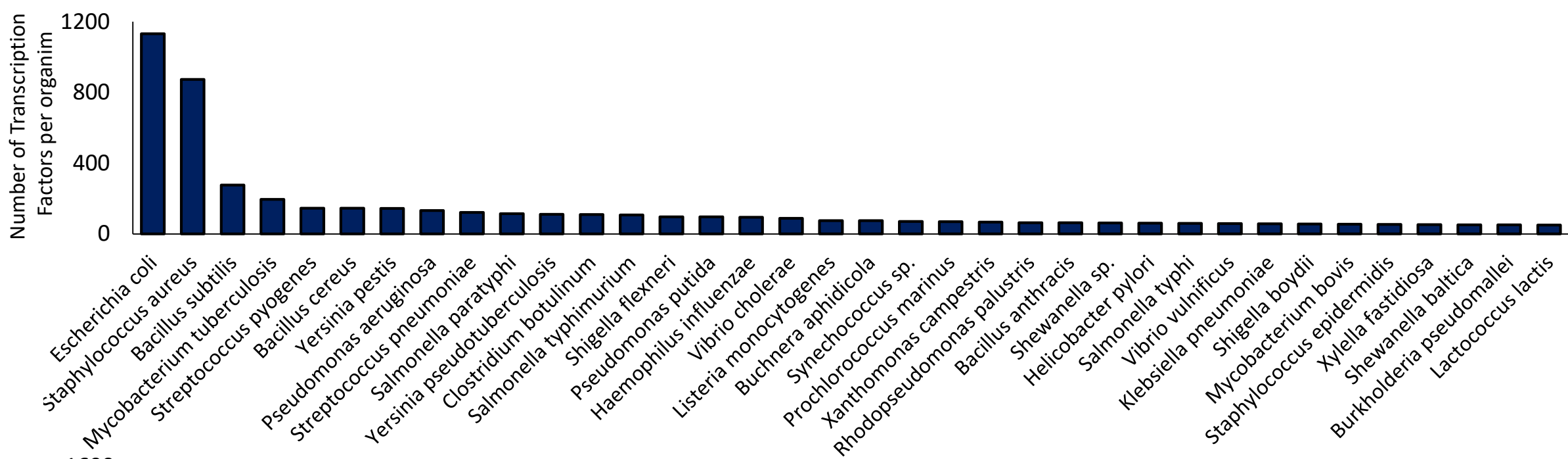


mapping



Mapping: 11 metatranscriptomes of reference from the same bioreactor where the metagenomes were collected

A



B

