

Assessing and characterising the repertoire of constitutive promoter elements in soil metagenomic libraries in *Escherichia coli*

Cauã Antunes Westmann¹, Luana de Fátima Alves^{2,3}, Rafael Silva-Rocha^{1‡} and María-Eugenia Guazzaroni^{2‡*}

¹ Department of Cellular and Molecular Biology, FMRP, University of São Paulo, Ribeirão Preto, SP Brazil

² Department of Biology, FFCLRP, University of São Paulo, Ribeirão Preto, SP Brazil

³ Department of Biochemistry, FMRP, University of São Paulo, Ribeirão Preto, SP Brazil

‡These authors contributed equally to this work

*Correspondence to:

María-Eugenia Guazzaroni
Faculdade de Filosofia, Ciências e Letras de Ribeirão Preto, Universidade de São Paulo
Av. Bandeirantes, 3.900. CEP: 14049-901.
Ribeirão Preto, São Paulo, Brazil
Tel.: +55 16 3315 3695; Fax: +55 16 3633 4886
E-mail: meguazzaroni@gmail.com

1 SUMMARY

2 Although functional metagenomics has been widely employed for the discovery of genes relevant to biotechnology
3 and biomedicine, its potential for assessing the diversity of transcriptional regulatory elements of microbial
4 communities has remained poorly explored. Here, we have developed a novel framework for prospecting,
5 characterising and estimating the accessibility of promoter sequences in metagenomic libraries by combining a bi-
6 directional reporter vector, high-throughput fluorescence assays and predictive computational methods. Using the
7 expression profiling of fluorescent clones from two independent libraries from soil samples, we directly analysed
8 the regulatory dynamics of novel promoter elements, addressing the relationship between the “metaconstitutome”
9 of a bacterial community and its environmental context. Through the construction and screening of plasmid-based
10 metagenomic libraries followed by *in silico* analyses, we were able to provide both (i) a consensus exogenous
11 promoter elements recognizable by *Escherichia coli* and (ii) an estimation of the accessible promoter sequences
12 in a metagenomic library, which was close to 1% of the whole set of available promoters. The results presented
13 here should provide new directions for the exploration through functional metagenomics of novel regulatory
14 sequences in bacteria, which could expand the Synthetic Biology toolbox for novel biotechnological and biomedical
15 applications.

16

17 INTRODUCTION

18 The study of prokaryotic transcriptional regulation is essential for understanding the molecular
19 mechanisms underlying decision-making processes in microorganisms ¹, comprising populational (e.g. colony
20 structure, quorum sensing detection), ecological (e.g. nutrient acquisition, biomass degradation) and pathogenic
21 behaviours (e.g. host recognition, biofilm formation). The activity of most bacterial promoters is usually dependent
22 on the combined action of transcription factors and sigma factors in response to multiple environmental stimuli ².
23 For instance, in *Escherichia coli*, the compilation of decades of experimental data indicate that approximately 50%
24 of its promoters are under the control of a single specific regulator, while all other genes are regulated by at least
25 two transcription factors ³. Moreover, the recent development of experimental and large-scale sequencing
26 techniques, together with powerful computational approaches have allowed both the discovery of insightful

1 information about bacterial transcriptional systems and the development of novel approaches for studying those
2 systems in higher depth ⁴⁻⁷. However, despite technical innovations, most of the studies are still centred on the
3 model organism *E. coli*, a single bacterial species among at least 30,000 other already sequenced ⁸, in an estimated
4 total of 1 trillion species ⁹.

5 With the advent of Metagenomics ¹⁰, the exploration of unculturable bacteria (approximately 99% of a
6 bacterial community¹¹ widely expanded genomic information, providing resourceful data about populational
7 structures and genetic diversity in a myriad of environmental samples ¹²⁻¹⁴. Two main approaches are commonly
8 adopted for those metagenomic studies ¹⁵: the sequence-based metagenomic approach, which relies on massive
9 sequencing of metagenomic DNA and powerful bioinformatics tools for extracting information from the
10 metagenomic sequences; and functional metagenomics ^{16,17}, which directly explores the functionality of enzymes
11 and other structural elements through a wide range of stress/substrate/product-based assays ¹⁸⁻²¹. In this context,
12 although a large number of genes/ORFs has been discovered through the previously described approaches, the
13 detection of novel bacterial regulatory elements using high-throughput technologies has been poorly explored,
14 presenting so far a single well-defined method for the discovery of substrate-inducible regulatory sequences -
15 SIGEX - ¹⁹ and a limited assay for exploration of constitutive promoters ²². This narrow range of methodologies is
16 directly related to the biased functional search towards novel genes and to a lack of both experimental and
17 computational tools for finding and validating promoter sequences in metagenomic libraries ²³.

18 Unravelling novel bacterial promoters is essential for understanding the regulatory diversity of
19 microorganisms, addressing important questions, such as the abundance of both constitutive and inducible
20 elements in a metagenomic library, the bottlenecks regarding host choices (i.e. the constraints limiting the diversity
21 of exogenous regulatory sequences that can be recognized by different hosts) and the correlation between
22 promoter strength, transcriptional noise and the functional role of the regulated gene/operon ²³⁻²⁶. Furthermore,
23 prospecting and characterising novel regulatory sequences is crucial for expanding the current Synthetic Biology
24 toolbox and generating novel biotechnological applications. For instance, there is a high demand for novel
25 constitutive and inducible promoters responding to process-specific parameters imposed by a wide variety of
26 processes, such as industrial applications, heterologous protein expression and biosensors generation ^{19,23,27-29}.

1 In this context, the most common strategy for prospecting regulatory sequences is the usage of
2 unidirectional promoter trap-vectors, which consist in transcriptional fusions between DNA fragments and a
3 reporter gene. This method has been widely employed for assessing regulatory sequences in genomic DNA³⁰⁻³³,
4 however its application in metagenomic DNA fragments has remained poorly explored¹⁹. The main constraint
5 regarding the use of unidirectional systems is that bacterial genomes present a large variation in the percentage
6 of their leading-strand genes, ranging from ~45% to ~90%^{34,35}. Thus, a bi-directional promoter reporter system
7 would be preferable, by increasing the probability of finding promoter sequences. In the present work, we have
8 developed a novel strategy for in-depth prospection, characterisation, and quantification of accessible promoter
9 elements from soil metagenomic samples in *E. coli* as a standard host.

10 Although both constitutive and inducible promoters were potentially detectable by this method, we have
11 focused exclusively on the study of the former, as a proof of concept, by avoiding substrate-based induction assays
12 as previously reported¹⁸⁻²¹. We have collected soil samples from two differentially biomass-enriched sites of a
13 Secondary Atlantic Forest in South-eastern Brazil and generated metagenomic libraries in a bi-directional probe
14 vector for primary screenings. We have characterised the expression behaviours of a large set of GFP_Iva
15 expressing clones from both libraries and narrowed down our selection to 10 clones for an in-depth analysis
16 regarding potential ORFs and endogenous promoters. By cross-validating *in silico* analyses and experimental data
17 of predicted regulatory sequences, we have located and profiled the expression of 33 endogenous promoters
18 within the selected clones (see Supplementary Table S2 online), providing resourceful information concerning the
19 architecture and transcriptional dynamics of promoters from metagenomic fragments. Thus, in order to contribute
20 to this set of accessible genetic features, we have used our gathered data to provide for the first time a direct
21 estimation of the whole set of accessible constitutive promoters in a soil metagenomic library hosted in *E. coli*,
22 which we have called the “metaconstitutome” of an environmental sample.

23 **RESULTS AND DISCUSSION**

24 **Generating metagenomic libraries and screening for fluorescent clones**

25 We have constructed and assessed two metagenomic libraries hosted in *E. coli* DH10B strain for the
26 analysis of bacterial promoters in environmental samples (Figure 1). The libraries were generated from soil

1 microbial communities of two sites bearing differential tree litter composition (*Anadenanthera spp.* and *Phytolacca*
2 *dioica*) within a Secondary Atlantic Forest zone at the University of Sao Paulo, Ribeirão Preto, Brazil. Both
3 metagenomic DNA were cloned into the pMR1 bi-directional reporter vector, which has GFP_{Iva} and mCherry
4 reporter genes in opposite directions³⁶. Each metagenomic library presented about 250 Mb of environmental DNA
5 distributed into approximately 60.000 clones harbouring insert fragments size ranging from 1.5 Kb to 7 Kb, with an
6 average size of 4.1 Kb (Table 1). We have chosen fragments of 1.5-7 Kb in order to validate our strategy on
7 standard-sized functional metagenomic libraries based on plasmid vectors^{18,19,37-39}. In total, 1,100 fluorescent
8 clones, resulting in a rate of approximately one fluorescent clone every one hundred fifty clones (USP1) or every
9 ninety clones screened (USP3), were manually selected under blue light exposition. Then, these fluorescent clones
10 were directly recovered from LB agar plates supplemented with chloramphenicol. The direct screening was
11 preferred over the use of metagenomic clone pools from stocks as it reduces the chances of both biased clone
12 enrichment (e.g. clones with higher growth rates, usually clones bearing small inserts or without insert) and dilution
13 of positive clones with impaired growth (e.g. clones with high expression of GFP and/or other exogenous genes),
14 avoiding thus clonal amplification.

15 **Evaluating the expression dynamics of fluorescent clones**

16 In order to analyse the expression patterns of the isolated clones, we evaluated the intrinsic dynamics of
17 GFP_{Iva} and mCherry by randomly selecting 20 clones expressing each reporter (as schematically represented in
18 Figure 1). As represented in Figures 2A-B, we found that clones expressing mCherry were not suitable for standard
19 microplate 8 hour assays, as the fluorescence intensity values differed dramatically between 8 and 24 hours after
20 the beginning of the experiment. The slow kinetics of mCherry expression has already been reported as a
21 consequence of a two-step oxidation process for protein maturation when compared to the one-step maturation
22 process found in GFP reporters⁴⁰. On the other hand, the clones expressing GFP_{Iva} presented the enhanced
23 intrinsic properties for microplate assays, supported by the observation of very similar fluorescence intensities
24 between the two time points tested. Furthermore, the GFP_{Iva} has an LVA-degradation tag attached to its C-
25 terminal, which reduces GFP accumulation and increases protein turnover, generating a more precise fluorescence
26 output on analysis of expression patterns⁴¹.

1 Thus, 260 clones expressing GFP_{lva} (160 clones from the USP1 library and 100 from USP3) were
2 selected for further analysis of expression patterns on microplate reader assays with biological and technical
3 triplicates. The dynamic profiles for each clone were converted into heat maps and hierarchically clustered by a
4 Euclidean Distance algorithm into a dendrogram, concisely representing the expression patterns of each
5 metagenomic library. In order to assess the diversity of promoter strengths among the generated metagenomics
6 libraries, three previously characterized constitutive promoters (see Experimental Procedures for further
7 information) positioned upstream a GFP_{lva} reporter were used as standards for strong, medium and weak
8 expression profiles (referred here as p100, p106 and p114, respectively). Considering both metagenomics libraries,
9 we have found a total of 30 strong promoters showing a strength similar to the p100 control, 40 medium strength
10 promoters similar to the p106 control, 60 weak promoters similar to the p114 control and a wide range of promoters
11 with particular expression patterns which did not cluster with any of the previously mentioned positive controls
12 (Figure 2C and Supplementary Fig. S1 online). Since the exploration of distinct expression behaviours is essential
13 for expanding the current set of commercial promoters, the diversity of expression profiles highlighted in this study
14 has supported the current framework as a promising strategy for finding novel promoters for downstream
15 applications.

16 Furthermore, concerning the hierarchical organization of the expression profiles, the dendrogram of the
17 USP3 library (Figure 2C) suggests the presence of at least four well-defined expression clusters comprising: (i)
18 high, (ii) medium, (iii) low and (iv) very low expression profiles. A very similar pattern was identified in the
19 expression dendrogram independently generated for the USP1 metagenomic library (see Supplementary Fig. S1
20 online), suggesting those clusters might be depicting broader trends of organizational expression patterns in
21 nature. Independent studies on microbial communities from aquatic environments have described similar patterns
22 by evaluating gene expression through metatranscriptomic analysis⁴²⁻⁴⁵, indicating that our observations are not
23 restricted to the assessed soil samples. However, further studies with a systematic application of the
24 methodologies described here over a broader range of environmental samples would be required for evaluating
25 these profiles.

26

27 ***In silico* analysis of DNA metagenomic fragments from selected clones**

1 From the 260 assessed samples, we have selected 10 clones displaying particular profiles (see
2 Supplementary Fig. S2 online) depicting the diversity of expression behaviours found in both libraries. The inserts
3 from selected clones were sequenced and analysed for both potential ORFs and RpoD-related promoter regions
4 (-10 and -35 conserved regions). In the case of the identification of putative genes, twenty-nine ORFs with
5 significant *E-values* (<0,001) were found (Table 2 and Supplementary Table S1 online) unevenly distributed
6 between both DNA strands, in line with a lack of strong directional trends regarding bacterial genome organization
7 ⁴⁶. The ORFs were also classified within a range of functional classes (delineated by MultiFun⁴⁷ and potential
8 bacterial phyla (see Supplementary Fig. S3 online). For this, we carried out the analysis of the microorganisms
9 associated with the closest similar protein of the identified ORFs (Table 2). The most abundant ORFs were related
10 to unknown functions (31%) and metabolism (31%), followed by stress adaptation cell processes (17%) (see
11 Supplementary Table S1 online), while the most abundant phyla related to the recovered ORFs were
12 Proteobacteria (35%), followed by Bacteroidetes (22%) and Chloroflexi (14%) (see Supplementary Fig. S3 online).
13 The relative abundance of the guanine-cytosine content of each insert was also assessed (Table 2), resulting in a
14 median of 54%, varying from 43% to 61%, indicating their diverse phylogenetic affiliation. These results are in
15 agreement with previous G-C content diversity analyses of soil samples which ranged from 50% to 61% ⁴⁸⁻⁵⁰. Even
16 with a limited sample size when compared to NGS-based metagenomic studies, the abundance of gene functions
17 and bacterial groups predicted in this work was similar to the ones found in previous studies in soil microbial
18 communities ⁵¹⁻⁵³. Considering the above, these results suggest that different bacterial groups could be the sources
19 of accessible promoters in *E. coli*, that is, regulatory sequences recognizable by the molecular transcriptional
20 machinery of *E. coli* that allowed the expression of the reporter genes.

21 The *in silico* promoter prediction has also provided relevant information concerning the potential number
22 of regulatory regions on each selected fragment. The BPROM software ⁵⁴ has been extensively employed in other
23 promoter prediction studies and is based on the analysis of the -35 and -10 consensus sequence of RpoD
24 promoters. The main sigma subunit, sigma-70 encoded by *rpoD*, plays a major role in transcription of growth-
25 related genes, the so-called housekeeping genes ⁵⁵⁻⁵⁷. From the *in silico* analysis, a total of 140 promoters were
26 predicted among the 10 selected clones, suggesting an average of 5 RpoD-related promoters/Kb. This led us
27 reasoning that most of the expression profiles previously described (Figure 2C and Supplementary Figure S1

1 online) were representing the dynamics of the merged promoters present in the metagenomic fragment.
2 Considering that, we delineate a strategy to experimentally assess the number and location of accessible
3 promoters from our selected clones, contrasting experimental results with *in silico* data.

4 **Experimental identification, characterisation, and cross-validation of promoter regions**

5 In order to explore the potential set of accessible promoter regions from our metagenomic libraries, we
6 developed a small DNA insert library generation approach (Figure 1). Firstly, the plasmids from the previously 10
7 selected clones (original clones) were pooled together for insert amplification in a single PCR reaction. The
8 resulting amplicons were fragmented by Sau3AI digestion and DNA fragments ranging from 0.2 Kb to 0.5 Kb were
9 selected for subsequent cloning into the pMR1 vector. The generation of this sub-fragment library allowed the
10 screening for both red and green fluorescent colonies as they would represent the accessible set of promoters
11 among the metagenomic DNA fragments studied. It is important to highlight that as the cloning process was not
12 directed, small fragments bearing promoter regions had a 50% chance of getting cloned in any direction, thus
13 clones expressing mCherry were also isolated for subsequent sequencing. A total of 100 clones coming from the
14 small DNA insert library (80 expressing GFP_{Iva} and 20 expressing mCherry) were sequenced and then align
15 against the original metagenomic fragments. As a result, we have identified at least 33 promoter regions within the
16 initial set of the selected metagenomic clones (Figure 3, Supplementary Fig. S4 and Supplementary Table S2
17 online). These findings showed that the *in silico* prediction of 140 RpoD-related promoters was overestimated in
18 comparison with the experimental results. The above can be explained since prediction algorithms usually
19 misrepresent nature by underestimating or overestimating results due to a lack of information regarding diversity
20 and variability of natural *cis*-regulatory sequences⁵⁸⁻⁶⁰.

21 Additionally, the current experimental approach allowed us not only to identify novel promoter regions but
22 also to determine promoter directionality. The evaluation of promoter localization within the 10 selected clones
23 revealed that from the 33 experimentally selected small fragments, 7 (21%) were considered intragenic promoters
24 while the remaining 79% (26 promoters) were considered primary promoters, defined as the furthest upstream
25 promoter in a gene/operon⁶¹. This small-scale analysis slightly diverges from architectural features found in *E. coli*
26 K-12 genome in which the promoter dataset was dominated by primary promoters (66.3%), with a lower number
27 of secondary promoters (19.6%), defined as intergenic and downstream of primary promoters⁶¹, internal promoters

1 that are intragenic (9.8%), and antisense (4.2%) promoters^{61,62}. This observation might reflect the diversity of
2 genomic architectures in metagenomic libraries and highlight the current underestimation of bacterial intragenic
3 promoters, which doubled the number in comparison to *E. coli*.

4 Based on the alignment results, we selected a defined set of small fragment clones related to each original
5 sequence for dynamic expression profiling on a microplate reader. The results showed that for each set of small-
6 fragments belonging to a DNA metagenomic clone, there was at least one with an expression pattern
7 corresponding to the original clone previously observed (Figure 3 and Supplementary Fig. S4). Similarly, we
8 identified other clones bearing small-inserts with individual profiles different to the primarily observed, representing
9 alternative promoter regions in the original sequence that were not mapped in the initial approach (Figure 3). The
10 diversity of the promoter expression profiles found in a single original metagenomic clone has a multifactorial
11 nature, ruled by different processes. Firstly, it should be considered the inherent relationship between the
12 regulatory dynamics and the functional role of the regulated gene²⁶. Secondly, the transcriptional bias imposed by
13 the *E. coli* molecular machinery, which would recognize orthologous sequences, but not necessarily reproduce the
14 original behaviours found in natural hosts^{23,39,63,64}. Finally, another point to be considered is that the increase in
15 expression levels can be the result of the artificial juxtaposition of the promoter to the fluorescent reporter ribosome
16 binding site, as a consequence of the cloning process.

17 Regarding *in silico* cross-validation, from the 33 experimentally validated promoters, 23 RpoD-related
18 promoters (70%) were supported by the algorithmic analysis as they were aligned to their respective original
19 sequences (Figure 3). On the other hand, the remaining 10 sequences (30%) were considered as promoters
20 exclusively identified by experimental approaches. We hypothesized that these sequences could be either
21 recognized by other sigma factors than sigma70 or presented unusual consensus sequences for -10 and -35 boxes
22 which has bypassed the algorithmic analysis. However, experimental validation in *E. coli* strains lacking diverse
23 sigma factors genes should be necessary for a more accurate conclusion.

24 Finally, sequences of the above experimentally validated promoters were characterised accordingly to
25 previous studies reported in the literature. For this, we adopted an *in silico* classification proposed by Shimada et
26 al⁶⁵ (2014), in which constitutive promoters present a high-level conservation of the consensus sequence for the
27 major sigma factor RpoD, that is, the elements TTGACA (-35) and TATAAT (-10) separated by approximately 17

1 bp (Figure 4A and B). Constitutive promoters are defined as promoters active *in vivo* in all circumstances, and, on
2 the other hand, inducible promoters are switched ON and OFF by transcription factors depending on the *in vivo*
3 conditions ⁶⁵. The Logo pattern ⁶⁶ generated from the alignment of the 33 identified metagenomic promoters (Figure
4 4C) indicated that positions -35 and -34 (-35 box) and positions -8, -7 and -3 (-10 box) were highly conserved.
5 Although this logo pattern was distant from the proposed for the RpoD-dependent constitutive promoters identified
6 *in vitro* (Figure 4A ⁶⁵), was very similar to previously described consensus ⁶⁷ from experimentally validated promoter
7 sets from RegulonDB ³ and EcoCyc ⁶⁸ databases (Figure 4B). To conclude, the results presented here has allowed
8 us to identify a consensus for exogenous promoter recognition in *E. coli*, which can be an important resource for
9 defining host-dependent restrictions in functional metagenomics.

10 **Estimating the accessibility of promoters in random metagenomic libraries**

11 In the present work we provide, for the first time in literature, a quantitative estimation regarding the
12 accessibility of natural promoters in random metagenomic libraries, supported by the integration of both *in silico*
13 and experimental results. We have estimated the existence of at least 553,300 promoters virtually recognized by
14 *E. coli* in a standard functional soil metagenomic library, from which approximately 4,961 promoters (~1%) were
15 readily accessible by our methodologies (see Experimental Procedures for calculations details). For the sake of
16 comparison, we have estimated an average rate of 1.1 promoters/Kb potentially recognizable by *E. coli* host in our
17 metagenomic fragments, which seems reasonable when compared to the rate of reported promoters in the well-
18 studied genome of *E. coli* K-12, ranging from 0.5 to 2.7 promoters/Kb, depending on both chosen datasets and *in*
19 *silico* prediction parameters ^{61,62,65,69}. We have also assessed the genomic features from 24 bacterial species
20 catalogued on the DOOR2 database (Database of prOkaryotic OpeRons) ⁷⁰ for a broader promoter rate estimation,
21 resulting in 2 promoters/Kb, which is also in concordance to our estimation and to others reported in literature
22 ^{61,62,65,69}. It should be mentioned that the average promoter rate from the present study (1.1 promoters/Kb) is
23 probably an underestimation of the whole set available as it is restricted to clones expressing GFP under specific
24 experimental conditions. Consequently, modifications of laboratory conditions during the screenings (such as
25 growth-phase, temperature, exposure to different chemicals and substrates, to cite some) would probably reveal
26 novel promoter elements ⁷¹.

1 A seminal study in functional metagenomics provided by Gabor *et al*⁶³(2004) estimated on a theoretical
2 basis, using 32 prokaryotic genomes, that 40% of the enzymatic activities present in a soil metagenomic library
3 could be readily accessed using *E. coli* as a host in an independent gene expression mode (in which both the
4 promoter and the ribosome binding sites (RBS) are provided by the metagenomics insert). Moreover, it was
5 predicted that Firmicutes, instead of Proteobacteria, would present the largest fraction of independently expressible
6 genes (73%). Contrastingly, recent empirical studies on *E. coli* and other hosts have shown that functional
7 expression faces a myriad of challenges that were not taken into account in previous mathematical models, such
8 as codon usage, improper promoter and RBS recognition⁷², missing initiation factors, protein misfolding, missing
9 co-factors, breakdown of product; improper secretion of product, toxicity of product or intermediates and formation
10 of inclusion bodies^{24,25}. Since it is impossible to predict the effect of the previously described difficulties in unknown
11 metagenomic fragments, the actual fraction of genes that can be successfully expressed in *E. coli* is probably
12 significantly lower than the proposed by Gabor and collaborators⁶³ (2004). In this context, our work supports the
13 previous arguments^{24,25} highlighting the large gap between theoretical predictions and experimental data as we
14 have shown only a small portion of the whole set of promoters is accessible for *E. coli* in metagenomics libraries
15 (~1%). Thus, we stress the importance of feeding mathematical models with empirical data in a continuous iterative
16 process for improving its predictive power.

17 **CONCLUSIONS**

18 In summary, we have developed a novel methodology for prospecting, characterising and estimating the
19 accessibility of promoter sequences in metagenomic samples by combining experimental and *in silico* approaches.
20 The expression profiling of fluorescent clones was used for the first time as a direct approach to analyse the
21 regulatory dynamics of an environmental sample, bearing great potential for revealing insightful trends regarding
22 the transcriptional diversity of microbial communities. It has already been computationally demonstrated by
23 Fernandez *et al.* (2014)⁷³ that the microbial metaregulome – the whole set of regulons of an environmental sample
24 – is shaped by the physicochemical conditions of the environment as an adaptive process. Thus, future studies
25 systematically applying our methodology to a range of environmental samples will greatly contribute to
26 understanding this relationship between regulatory diversity and environmental adaptation in bacteria. At the same

1 time, it can also be further applied to the design of efficient microbial communities for therapeutic or ecological
2 needs ^{73–76}.

3 Through the generation of a small-DNA insert library approach combined to *in silico* promoter prediction
4 we were able to provide both (i) a consensus of recognizable exogenous regulatory sequences in an *E. coli* host
5 and (ii) an estimation of the accessible promoter sequences in a plasmid-based functional metagenomic library,
6 which was close to 1% of the whole set of available promoters. These are resourceful data for building a concise
7 framework regarding the accessibility of genetic features from metagenomic libraries and how it can be influenced
8 by the choice of different microbial hosts ^{23,63,64} or by the tinkering of the host's transcription systems ^{72,77,78}.

9 Although this work provided seminal information regarding promoter accessibility in metagenomics
10 libraries, further high-throughput studies optimizing the proposed methods (e.g. application of automated screening
11 methods; exploration of the whole set of fluorescent clones in a metagenomics library by Next-Generation-
12 Sequencing) will be essential for expanding our current estimation into a more holistic landscape. Finally, we
13 highlight that besides providing novel approaches for studying the regulatory diversity underlying environmental
14 microbial communities, this work should be extremely useful for expanding the current Synthetic Biology toolbox
15 through the discovery and characterisation of novel regulatory features.

16

17 **EXPERIMENTAL PROCEDURES**

18 **Bacterial strains, primers, plasmids and general growth conditions**

19 *E. coli* DH10B (Invitrogen) cells were used for cloning and experimental procedures. *E. coli* strains were routinely
20 grown at 37°C in Luria-Broth medium or M9 minimal medium ⁷⁹ (6.4 g/L Na₂HPO₄·7H₂O, 1.5 g/L KH₂PO₄, 0.25
21 g/L NaCl, and 0.5 g/L NH₄Cl) supplemented with 2 mM MgSO₄, 0.1 mM casamino acid, and 1% glycerol as the
22 sole carbon source. When required, chloramphenicol (Cm) (34 µg/mL) was added to the medium to ensure plasmid
23 retention. When cells were grown in minimal medium, antibiotics were used at half concentrations. Transformed
24 bacteria were recovered on LB (Luria–Bertani) liquid medium for 1 hour at 37°C and 180 r.p.m, followed by plating
25 on LB-agar plates at 37°C for at least 18 hours. All constructions were cloned into the pMR1 bi-directional-reporter
26 vector ³⁶, which carries mCherry and GFP_{Iva}, a short-lived variant of GFP.

1 **Metagenomic libraries construction and screening for fluorescent clones**

2 The metagenomic libraries used in this work were generated in our laboratory from two soil samples of a Secondary
3 Atlantic Forest at the University of Sao Paulo, Ribeirão Preto, Brazil. Each sample was differentially enriched
4 regarding tree species abundance on plant-litter composition: (i) enriched in leaves from *Phytolacca dioica* and (ii)
5 from *Anadenanthera spp.* DNA was extracted from soil samples using the UltraClean™ Soil DNA isolation Kit (Mo
6 Bio Laboratories, Solana Beach, CA, USA). For the construction of the libraries, metagenomic DNA was partially
7 digested using Sau3AI, and fragments from 1.5 kb to 7 kb were extracted from an agarose gel for ligation into the
8 dephosphorylated and BamHI-digested pMR1 vector. Ligation mixtures were transformed by electroporation into
9 *E. coli* DH10B cells. To amplify the libraries, they were grown on LB agar plates containing Cm and incubated for
10 18 h at 37°C. Both green and red clones were manually isolated from LB-agar plates exposed to blue light
11 wavelength (at approximately 470 nm) by a transilluminator (Safe Imager™ 2.0 Blue Light Transilluminator). Ten
12 fluorescent and twenty non-fluorescent clones were randomly picked from each library and had their plasmids
13 extracted, following digestion with EcoRI and SmaI enzymes for checking presence/absence of inserts and their
14 sizes. Cells from the same library were collected and pooled together in LB supplemented with 10% (wt/vol)
15 glycerol for storing at -80°C. The plasmids from the 10 selected clones were isolated from individual clones and
16 transformed into new *E. coli* DH10B cells to reconfirm expression patterns.

17 **Nucleic acid techniques**

18 DNA preparation, digestion with restriction enzymes, analysis by agarose gel electrophoresis, isolation of DNA
19 fragments, ligations, and transformations were done by standard procedures (Ausubel et al., 1994). Plasmid DNA
20 was sequenced on both strands by primer walking using the ABI PRISM Dye Terminator Cycle Sequencing Ready
21 Reaction kit (PerkinElmer) and an ABI PRISM 377 sequencer (Perkin-Elmer) according to the manufacturer's
22 instructions.

23 **GFP fluorescence assay and data processing**

24 To measure promoter activity, freshly plated single colonies were grown overnight in M9 medium supplemented
25 with required antibiotics. Samples were diluted 1:20 (v/v) in M9 medium for a final volume of 200uL in 96-well
26 microplates. Cell growth and GFP fluorescence were quantified using a Victor X3 plate reader (PerkinElmer,

1 Waltham, MA, USA). Promoter activities were expressed as the emission of fluorescence at 535 nm upon excitation
2 with 485 nm light and then normalised with the optical density at each point (reported as fluorescence/OD600)
3 after background correction. Background signal was evaluated with non-inoculated M9 medium and used as a
4 blank for adjusting the baseline of measurements. *E. coli* DH10B harbouring the pMR1 empty plasmid was used
5 as a negative control. Three different positive controls were used, consisting in *E. coli* DH10B harbouring pMR1
6 plasmid with one of the following synthetic constitutive promoters from the iGEM BBa_J23104 Anderson's
7 catalogue (<http://parts.igem.org/Promoters/Catalog/Anderson>)⁸⁰ upstream a GFP_{Iva} reporter: J23100, J23106 and
8 J23114 (referred here as p100, p106 and p114, respectively). Unless otherwise indicated, measurements were
9 taken at 30 min intervals over 8 h. All experiments were performed with both technical and biological replicates,
10 being biological triplicates evaluated as independent measurements on different dates. Raw data were processed
11 and plots were constructed using Microsoft Excel. All data was normalised by background values and transformed
12 to a log₂ scale for better data visualisation. Heatmap dendrograms with expression profiles were generated by
13 using MeV2 (<http://mev.tm4.org/>) software.

14 **Small-DNA inserts libraries generation and screening**

15 In order to experimentally find and validate the promoter regions from each of the ten selected metagenomic
16 clones, an experimental technique was developed based on the previously described methodology of
17 metagenomic library construction. All selected clones had their plasmids extracted and pooled together in an
18 equimolar ratio. The pooled sample was amplified through a single PCR reaction using high-fidelity polymerase
19 enzyme (Phusion) and previously described primers flanking the MCS region (Multiple Cloning Site) of the pMR1
20 vector, into which the metagenomic inserts were cloned. The resulting amplicons were firstly submitted to an
21 analytical digestion followed by electrophoretic analysis for finding the optimal concentration of Sau3AI enzyme for
22 obtaining fragments size ranging from 0.1Kb to 0.5Kb. Then, the purified pooled samples were fragmented by
23 Sau3AI in preparative digestion and thereafter punctured from a 1% agarose gel in the region between 0.1 Kb and
24 0.5 Kb. These small DNA fragments, in turn, were ligated to pMR1 vector. Aliquots of electrocompetent *E. coli*
25 DH10B cells were transformed with ligated DNA. A total of 100 fluorescent clones (80 expressing GFP and 20
26 expressing mCherry) were isolated under blue light excitation screening and had their plasmids extracted for

1 sequencing reactions. Fluorescent clones were stored at -80°C in LB medium supplemented with required
2 antibiotics and 10% glycerol (v/v).

3 ***In silico* analysis of ORFs and promoter regions**

4 The inserts of selected clones were sequenced on both strands as previously described. Sequences were manually
5 assembled for the generation of 10 contigs. Putative ORFs were identified and analysed using the online ORF
6 Finder platform, available at the NCBI website (<http://www.ncbi.nlm.nih.gov/gorf/gorf.html>). Comparisons of
7 nucleotide and transcribed amino acid sequences were performed against public databases (NCBI) using BlastN,
8 BlastX and BlastP (BLAST, basic local alignment search tool) at the NCBI on-line server. For translation to protein
9 sequences, the bacterial code was selected, allowing ATG, GTG, and TTG as alternative start codons. All the
10 predicted ORFs longer than 270 bp were translated and used as queries in BlastP. Sequences with significant
11 matches were further analysed with psiBlast, and their putative function was annotated based on their similarities
12 to sequences in the COG (Clusters of Orthologous Groups) and Pfam (Protein Families) databases. Predicted
13 general cellular functions were annotated only for known ORFs based on the MultiFam classification (Serres et al,
14 2006). All sequences with an E-value higher than 0.001 in the BlastP searches and longer than 300 bp were
15 considered to be unknown. Transmembrane helices were predicted with TMprep ([http://www.ch.
16 embnet.org/software/TMPRED_form.html](http://www.ch.embnet.org/software/TMPRED_form.html)) and signal peptides with Signal P3.0 server ([http://www.cbs.
17 dtu.dk/services/SignalP/](http://www.cbs.dtu.dk/services/SignalP/)). A complete table can be found at Supplementary Table S1 online. Promoter prediction
18 was based on the analysis of the ten contigs by using both BPROM
19 (<http://www.softberry.com/berry.phtml?topic=bprom&group=programs&subgroup=gfindb>) and bTSSfinder
20 (<http://www.cbrc.kaust.edu.sa/btssfinder/>) web-based platforms. Both methods searched for rpoD-related
21 sequences and we have only considered as valid predictions the ones matched on both approaches. Those filtered
22 sequences were used to cross-validate 23 out of 33 experimentally defined regulatory regions by comparing the
23 positions between predicted and experimental sequences in metagenomic fragments. The positions of the 33 small
24 DNA fragments were obtained by a multiple alignment of the original contigs (queries) against those selected
25 sequences, which has also allowed the validation of the promoter's directionality – forward or reverse - by observing
26 the matched strands (Plus/Plus or Plus/Minus). The consensus Logo sequence was based on the alignment of the
27 33 experimentally validated promoters, using the WebLogo platform (<http://weblogo.berkeley.edu/logo.cgi>).

1 **Calculations for promoter/Kb rates from databases and for promoter accessibility estimation**

2 Data from predicted sequences of promoter sites, TSS (Transcriptional Start Site) and TUs (Transcription Unit)
3 reported in different studies and databases regarding *E. coli* and other bacteria ^{3,61,65,70} were used as proxies for
4 the total number of predicted promoters. Those values were divided by their respective genome sizes (or average
5 genome sizes when calculating an average rate of multiple species at once) in order to provide promoter/Kb rates
6 (i.e. 8,000 predicted promoters, TSS or TUs on a genome of 4.6 Mb would result in a rate of 1.7 promoters/Kb). The
7 promoter accessibility estimation followed the same rationale and was based on the combination of the data from
8 both metagenomics libraries presented in Table 1 and the rate of experimentally discovered promoters per Kb (33
9 promoters found in 30 Kb of metagenomic DNA, resulting in a rate of 1.1 promoters/Kb). Firstly, we have merged
10 data from both metagenomic libraries and calculated the predicted number of promoters in a metagenomics library
11 with an effective size of 503 Mb (combined effective sizes of USP1 and USP3) – the “effective size” takes into
12 account only the percentage of clones with an insert -. Thus, we have multiplied the effective library size by the
13 23previously obtained promoter rate (1.1 promoters/Kb), resulting in a total estimated set of 553,300 promoter
14 sequences. Secondly, we have calculated the predicted set of accessible promoters by multiplying the number of
15 fluorescent clones (1,100 clones, considering both libraries) by the average insert size (4.1 Kb) and by the rate of
16 observed promoters per Kb (1.1 promoters/Kb), resulting in 4,961 potentially accessible promoters. Lastly, we have
17 calculated the proportion of accessible promoters among the total number of predicted promoters, which
18 represents approximately ~1% of the whole available set.

19

20 **Data Availability**

21 The nucleotide sequences obtained for the plasmid inserts have been deposited in the GenBank database under
22 the Accession numbers (KY939589-KY939597), which are also shown in Table 2.

23

24 **REFERENCES:**

25 1. Ishihama, A. Prokaryotic genome regulation: Multifactor promoters, multitarget regulators and hierarchic
26 networks. *FEMS Microbiol. Rev.* **34**, 628–645 (2010).

- 1 2. Browning, D. F. & Busby, S. J. W. Local and global regulation of transcription initiation in bacteria. *Nat. Rev. Microbiol.* **14**, 638–650 (2016).
- 3 3. Gama-Castro, S. *et al.* RegulonDB version 9.0: high-level integration of gene regulation, coexpression, motif clustering and beyond. *Nucleic Acids Res.* **44**, D133–D143 (2016).
- 5 4. Covert, M. W., Knight, E. M., Reed, J. L., Herrgard, M. J. & Palsson, B. O. Integrating high-throughput and computational data elucidates bacterial networks. *Nature* **429**, 92–96 (2004).
- 7 5. Martínez-Antonio, A., Collado-Vides, J., Martínez-Antonio, A. & Collado-Vides, J. Identifying global regulators in transcriptional regulatory networks in bacteria. *Curr. Opin. Microbiol.* **6**, 482–489 (2003).
- 9 6. Shen-Orr, S. S., Milo, R., Mangan, S. & Alon, U. Network motifs in the transcriptional regulation network of *Escherichia coli*. *Nat. Genet.* **31**, 64–68 (2002).
- 11 7. Shimada, T., Fujita, N., Maeda, M. & Ishihama, A. Systematic search for the Cra-binding promoters using genomic SELEX system. *Genes to Cells* **10**, 907–918 (2005).
- 13 8. Land, M. *et al.* Insights from 20 years of bacterial genome sequencing. *Funct. Integr. Genomics* **15**, 141–61 (2015).
- 15 9. Locey, K. J. & Lennon, J. T. Scaling laws predict global microbial diversity. *Proc. Natl. Acad. Sci.* **113**, 5970–5975 (2016).
- 17 10. Handelsman, J., Rondon, M. R., Brady, S. F., Clardy, J. & Goodman, R. M. Molecular biological access to the chemistry of unknown soil microbes: a new frontier for natural products. *Chem. Biol.* **5**, R245–R249 (1998).
- 20 11. Amann, R. I., Ludwig, W. & Schleifer, K. H. Phylogenetic identification and in situ detection of individual microbial cells without cultivation. *Microbiol. Rev.* **59**, 143–69 (1995).
- 22 12. Venter, J. C. Environmental Genome Shotgun Sequencing of the Sargasso Sea. *Science (80-.)*. **304**, 66–74 (2004).
- 24 13. Torsvik, V. & Øvreås, L. Microbial diversity and function in soil: from genes to ecosystems. *Curr. Opin. Microbiol.* **5**, 240–5 (2002).
- 26 14. Tringe, S. G. Comparative Metagenomics of Microbial Communities. *Science (80-.)*. **308**, 554–557 (2005).
- 27 15. Singh, J. *et al.* Metagenomics: Concept, methodology, ecological inference and recent advances. *Biotechnol. J.* **4**, 480–494 (2009).
- 29 16. Cowan, D. *et al.* Metagenomic gene discovery: past, present and future. *Trends Biotechnol.* **23**, 321–329 (2005).
- 31 17. Li, X. & Qin, L. Metagenomics-based drug discovery and marine microbial diversity. *Trends Biotechnol.* **23**, 539–543 (2005).
- 33 18. Guazzaroni, M. E., Morgante, V., Mirete, S. & González-Pastor, J. E. Novel acid resistance genes from the metagenome of the Tinto River, an extremely acidic environment. *Environ. Microbiol.* **15**, 1088–1102 (2013).
- 36 19. Uchiyama, T., Abe, T., Ikemura, T. & Watanabe, K. Substrate-induced gene-expression screening of environmental metagenome libraries for isolation of catabolic genes. *Nat. Biotechnol.* **23**, 88–93 (2005).
- 38 20. Uchiyama, T. & Miyazaki, K. Product-Induced Gene Expression, a Product-Responsive Reporter Assay Used To Screen Metagenomic Libraries for Enzyme-Encoding Genes. *Appl. Environ. Microbiol.* **76**, 7029–7035 (2010).
- 41 21. Williamson, L. L. *et al.* Intracellular Screen To Identify Metagenomic Clones That Induce or Inhibit a Quorum-Sensing Biosensor. *Appl. Environ. Microbiol.* **71**, 6335–6344 (2005).
- 42

- 1 22. Han, S. S., Lee, J. Y., Kim, W. H., Shin, H. J. & Kim, G. J. Screening of promoters from metagenomic DNA
2 and their use for the construction of expression vectors. *J. Microbiol. Biotechnol.* **18**, 1634–1640 (2008).
- 3 23. Guazzaroni, M.-E., Silva-Rocha, R. & Ward, R. J. Synthetic biology approaches to improve biocatalyst
4 identification in metagenomic library screening. *Microb. Biotechnol.* **8**, 52–64 (2015).
- 5 24. Ekkers, D. M., Cretoiu, M. S., Kielak, A. M. & van Elsas, J. D. The great screen anomaly—a new frontier
6 in product discovery through functional metagenomics. *Appl. Microbiol. Biotechnol.* **93**, 1005–1020 (2012).
- 7 25. Vester, J. K., Glaring, M. A. & Stougaard, P. Improved cultivation and metagenomics as new tools for
8 bioprospecting in cold environments. *Extremophiles* **19**, 17–29 (2015).
- 9 26. Silander, O. K. *et al.* A genome-wide analysis of promoter-mediated phenotypic noise in *Escherichia coli*.
10 *PLoS Genet.* **8**, (2012).
- 11 27. Silva-Rocha, R. & de Lorenzo, V. Mining logic gates in prokaryotic transcriptional regulation networks.
12 *FEBS Lett.* **582**, 1237–1244 (2008).
- 13 28. Boyle, P. M. & Silver, P. A. Harnessing nature's toolbox: regulatory elements for synthetic biology. *J. R.*
14 *Soc. Interface* **6**, S535–S546 (2009).
- 15 29. Blount, B. A., Weenink, T., Vasylechko, S. & Ellis, T. Rational Diversification of a Promoter Providing Fine-
16 Tuned Expression and Orthogonal Regulation for Synthetic Biology. *PLoS One* **7**, e33279 (2012).
- 17 30. Kubota, M., Yamazaki, Y. & Ishihama, A. Random screening of promoters from *Escherichia coli* and
18 classification based on the promoter strength. *Japanese J. Genet.* **66**, 399–409 (1991).
- 19 31. Dunn, A. K. & Handelsman, J. A vector for promoter trapping in *Bacillus cereus*. *Gene* **226**, 297–305
20 (1999).
- 21 32. Lu, C., Bentley, W. E. & Rao, G. A High-Throughput Approach to Promoter Study Using Green Fluorescent
22 Protein. *Biotechnol. Prog.* **20**, 1634–1640 (2004).
- 23 33. Chen, S., Bagdasarian, M., Kaufman, M. G. & Walker, E. D. Characterization of Strong Promoters from an
24 Environmental *Flavobacterium hibernum* Strain by Using a Green Fluorescent Protein-Based Reporter
25 System. *Appl. Environ. Microbiol.* **73**, 1089–1100 (2007).
- 26 34. Mao, X., Zhang, H., Yin, Y. & Xu, Y. The percentage of bacterial genes on leading versus lagging strands
27 is influenced by multiple balancing forces. *Nucleic Acids Res.* **40**, 8210–8218 (2012).
- 28 35. Mao, X. *et al.* Revisiting operons: an analysis of the landscape of transcriptional units in *E. coli*. *BMC*
29 *Bioinformatics* **16**, 356 (2015).
- 30 36. Guazzaroni, M.-E. & Silva-Rocha, R. Expanding the Logic of Bacterial Promoters Using Engineered
31 Overlapping Operators for Global Regulators. *ACS Synth. Biol.* **3**, 666–675 (2014).
- 32 37. Jiménez, D. J., Montaña, J. S., Álvarez, D. & Baena, S. A novel cold active esterase derived from
33 Colombian high Andean forest soil metagenome. *World J. Microbiol. Biotechnol.* **28**, 361–370 (2012).
- 34 38. Pushpam, P., Rajesh, T. & Gunasekaran, P. Identification and characterization of alkaline serine protease
35 from goat skin surface metagenome. *AMB Express* **1**, 3 (2011).
- 36 39. Gabor, E. M., de Vries, E. J. & Janssen, D. B. Construction, characterization, and use of small-insert gene
37 banks of DNA isolated from soil and enrichment cultures for the recovery of novel amidases. *Environ.*
38 *Microbiol.* **6**, 948–958 (2004).
- 39 40. Hebisch, E., Knebel, J., Landsberg, J., Frey, E. & Leisner, M. High Variation of Fluorescence Protein
40 Maturation Times in Closely Related *Escherichia coli* Strains. *PLoS One* **8**, e75991 (2013).
- 41 41. Andersen, J. B. *et al.* New unstable variants of green fluorescent protein for studies of transient gene
42 expression in bacteria. *Appl. Environ. Microbiol.* **64**, 2240–2246 (1998).

- 1 42. Frias-Lopez, J. *et al.* Microbial community gene expression in ocean surface waters. *Proc. Natl. Acad. Sci.*
2 **105**, 3805–3810 (2008).
- 3 43. Dupont, C. L. *et al.* Genomes and gene expression across light and productivity gradients in eastern
4 subtropical Pacific microbial communities. *ISME J.* **9**, 1076–92 (2015).
- 5 44. Fortunato, C. S. & Crump, B. C. Microbial gene abundance and expression patterns across a river to ocean
6 salinity gradient. *PLoS One* **10**, 1–22 (2015).
- 7 45. Stewart, F. J., Ulloa, O. & Delong, E. F. Microbial metatranscriptomics in a permanent marine oxygen
8 minimum zone. *Environ. Microbiol.* **14**, 23–40 (2012).
- 9 46. Koonin, E. V. Evolution of genome architecture. *Int. J. Biochem. Cell Biol.* **41**, 298–306 (2009).
- 10 47. Serres, M. H. & Riley, M. MultiFun, a Multifunctional Classification Scheme for Escherichia coli K-12 Gene
11 Products. *Microb. Comp. Genomics* **5**, 205–222 (2000).
- 12 48. Bohlin, J. *et al.* Analysis of intra-genomic GC content homogeneity within prokaryotes. *BMC Genomics* **11**,
13 464 (2010).
- 14 49. Mann, S. & Chen, Y. P. P. Bacterial genomic G + C composition-eliciting environmental adaptation.
15 *Genomics* **95**, 7–15 (2010).
- 16 50. Foerster, K. U., von Mering, C., Hooper, S. D. & Bork, P. Environments shape the nucleotide composition
17 of genomes. *EMBO Rep.* **6**, 1208–13 (2005).
- 18 51. Fierer, N. *et al.* Cross-biome metagenomic analyses of soil microbial communities and their functional
19 attributes. *Proc. Natl. Acad. Sci.* **109**, 21390–21395 (2012).
- 20 52. Fierer, N., Bradford, M. A. & Jackson, R. B. Toward an ecological classification of soil bacteria. *Ecology*
21 **88**, 1354–1364 (2007).
- 22 53. Janssen, P. H. Identifying the Dominant Soil Bacterial Taxa in Libraries of 16S rRNA and 16S rRNA Genes
23 MINIREVIEWS Identifying the Dominant Soil Bacterial Taxa in Libraries of 16S rRNA and 16S rRNA
24 Genes. *Appl. Environ. Microbiol.* **72**, 1719–1728 (2006).
- 25 54. Solovyev, V. & Salamov, A. in *Metagenomics and its Applications in Agriculture, Biomedicine and*
26 *Environmental Studies* (ed. R.W., L.) 61–78 (Nova Science Publishers, 2011).
- 27 55. Gruber, T. M. & Gross, C. A. Multiple Sigma Subunits and the Partitioning of Bacterial Transcription Space.
28 *Annu. Rev. Microbiol.* **57**, 441–466 (2003).
- 29 56. Paget, M. S. B. & Helmann, J. D. The sigma70 family of sigma factors. *Genome Biol.* **4**, 203 (2003).
- 30 57. Lonetto, M., Gribskov, M. & Gross, C. A. The sigma 70 family: sequence conservation and evolutionary
31 relationships. *J. Bacteriol.* **174**, 3843–9 (1992).
- 32 58. Vanet, A., Marsan, L. & Sagot, M. F. Promoter sequences and algorithmical methods for identifying them.
33 *Res. Microbiol.* **150**, 779–799 (1999).
- 34 59. de Jong, A., Pietersma, H., Cordes, M., Kuipers, O. P. & Kok, J. PePPER: a webserver for prediction of
35 prokaryote promoter elements and regulons. *BMC Genomics* **13**, 299 (2012).
- 36 60. Shahmuradov, I. A., Mohamad Razali, R., Bougouffa, S., Radovanovic, A. & Bajic, V. B. bTSSfinder: a
37 novel tool for the prediction of promoters in Cyanobacteria and *Escherichia coli*. *Bioinformatics* **33**, btw629
38 (2016).
- 39 61. Conway, T. *et al.* Unprecedented high-resolution view of bacterial operon architecture revealed by RNA
40 sequencing. *MBio* **5**, 1–12 (2014).
- 41 62. Cho, B.-K. *et al.* Elucidation of the transcription unit architecture of the Escherichia coli K-12 MG1655
42 genome. *Nat. Biotechnol.* **27**, 1043–1049 (2009).

- 1 63. Gabor, E. M., Alkema, W. B. L. & Janssen, D. B. Quantifying the accessibility of the metagenome by
2 random expression cloning techniques. *Environ. Microbiol.* **6**, 879–886 (2004).
- 3 64. Liebl, W. *et al.* Alternative hosts for functional (meta)genome analysis. *Appl. Microbiol. Biotechnol.* **98**,
4 8099–8109 (2014).
- 5 65. Shimada, T., Yamazaki, Y., Tanaka, K. & Ishihama, A. The whole set of constitutive promoters recognized
6 by RNA polymerase RpoD holoenzyme of Escherichia coli. *PLoS One* **9**, (2014).
- 7 66. Crooks, G. E. WebLogo: A Sequence Logo Generator. *Genome Res.* **14**, 1188–1190 (2004).
- 8 67. Mitchell, J. E. Identification and analysis of ‘extended -10’ promoters in Escherichia coli. *Nucleic Acids*
9 *Res.* **31**, 4689–4695 (2003).
- 10 68. Keseler, I. M. *et al.* The EcoCyc database: reflecting new knowledge about Escherichia coli K-12. *Nucleic*
11 *Acids Res.* **45**, D543–D550 (2017).
- 12 69. Rangel, C. P., Galán, E. & Martínez, A. Consensus architecture of promoters and transcription units in
13 Escherichia coli: design principles for synthetic biology. *Mol. BioSyst.* (2017). doi:10.1039/C6MB00789A
- 14 70. Mao, X. *et al.* DOOR 2.0: presenting operons and their functions through dynamic and integrated views.
15 *Nucleic Acids Res.* **42**, D654–D659 (2014).
- 16 71. Chang, D. E., Smalley, D. J. & Conway, T. Gene expression profiling of Escherichia coli growth transitions:
17 An expanded stringent response model. *Mol. Microbiol.* **45**, 289–306 (2002).
- 18 72. Bernstein, J. R., Bulter, T., Shen, C. R. & Liao, J. C. Directed Evolution of Ribosomal Protein S1 for
19 Enhanced Translational Efficiency of High GC Rhodospseudomonas palustris DNA in Escherichia coli. *J.*
20 *Biol. Chem.* **282**, 18929–18936 (2007).
- 21 73. Fernandez, L., Mercader, J. M., Planas-Fèlix, M. & Torrents, D. Adaptation to environmental factors shapes
22 the organization of regulatory regions in microbial communities. *BMC Genomics* **15**, 877 (2014).
- 23 74. Solé, R. Bioengineering the biosphere? *Ecol. Complex.* **22**, 40–49 (2015).
- 24 75. Johns, N. I., Blazejewski, T., Gomes, A. L. & Wang, H. H. Principles for designing synthetic microbial
25 communities. *Curr. Opin. Microbiol.* **31**, 146–153 (2016).
- 26 76. Fredrickson, J. K. Ecological communities by design. *Science (80-.).* **348**, 1425–1427 (2015).
- 27 77. Lämmle, K. *et al.* Identification of novel enzymes with different hydrolytic activities by metagenome
28 expression cloning. *J. Biotechnol.* **127**, 575–592 (2007).
- 29 78. Gaida, S. M. *et al.* Expression of heterologous sigma factors enables functional screening of metagenomic
30 and heterologous genomic libraries. *Nat. Commun.* **6**, 7045 (2015).
- 31 79. Sambrook, J.; Fritsch, E. F.; Maniatis, T. *Molecular cloning: a laboratory manual.* (Cold Spring Harbor
32 Laboratory Press, 1989).
- 33 80. Kelly, J. R. *et al.* Measuring the activity of BioBrick promoters using an in vivo reference standard. *J. Biol.*
34 *Eng.* **3**, 4 (2009).
- 35 81. Raes, J., Korbøl, J. O., Lercher, M. J., von Mering, C. & Bork, P. Prediction of effective genome size in
36 metagenomic samples. *Genome Biol.* **8**, R10 (2007).
- 37

38 Acknowledgements

1 This work was supported by the National Council for Technological and Scientific Development (CNPq
2 472893/2013-0 and 441833/2014-4) and by Young Research Awards by the Sao Paulo State Foundation
3 (FAPESP, award numbers 2015/04309-1 and 2012/21922-8). CAW and LFA are beneficiaries of FAPESP
4 fellowships (award numbers 2016/05472-6 and 2016/06323-4, respectively). Authors have no conflict of interest
5 to declare.

6 **Author Contributions**

7 CAW, LFA, MEG and RSR designed the experiments. CAW and LFA performed the experiments. CAW analyzed
8 the data. CAW and RSR prepared the figures. CAW and MEG wrote the manuscript. All authors reviewed the
9 manuscript.

10

11 **Additional Information**

12 **Competing financial interests**

13 The author(s) declare no competing financial interests.

14

15 **Tables and figures**

16

17 **Table 1.** Features of the generated metagenomic libraries.

Metagenomic Library	USP 1	USP 3
Total number of clones	100,000	90,000
Percentage of clones with insert	60%	70%
Number of clones with insert	60,000	63,000
Total number and rate* of fluorescent clones	400 (1:150)	700 (1:90)
Total number and rate* of green clones	270 (1:220)	400 (1:157)
Total number and rate* of red clones	130 (1:460)	300 (1:210)
Average insert size	4,5 kb	3,7 kb
Library Size	270 Mb	233 Mb
Estimated number of genomes**	60	52

18

19 * Rate represented by the number of fluorescent clones divided by the total number of clones with inserts.

20 ** Assuming 4.5 Mb per genome ⁸¹.

21

22

1 **Table 2.** Description of the ORFs contained in plasmids from the selected clones (pCAW1 to pCAW10) and their
 2 sequence similarities.

3

Clone_Sample [insert bp]	% G + C	GenBank accession No.	ORF ^a	Strand	Length (aa ^b)	Closest similar protein ^c (Length in aa)	Organism	Identity (%)	Putative function
pCAW1 (2367bp)	55%	KY939589	1	Minus	131	hypothetical protein (416)	Bacterioidetes bacterium	68%	Alginate lyase
			2	Plus	271	hypothetical protein (261)	Acidobacteria bacterium	73%	17-B-hydroxysteroid dehydrogenase
			3 ^a	Plus	295	beta-glucosidase (777)	Caulobacter sp. OV484	66%	beta-glucosidase
pCAW2 (2069bp)	52%	KY939590	1	Plus	304	Unkonwn ^c	Hyphomicrobium sp. NDB2Meth4	33%	Unknown
			2	Plus	249	Unkonwn	Hungatella hathewayi	33%	Unknown
pCAW3 (4404bp)	53%	KY939591	1	Minus	318	IS4 family transposase (320)	Escherichia coli	96%	IS4 family transposase
			2	Minus	1011	DNA-directed RNA polymerase subunit beta' (1430)	Sphingobacteriales bacterium 44-61	83%	RNA polymerase - Beta Subunit
			3	Plus	120	Uncharacterised protein (135)	Bordetella pertussis	47%	Unknown
			4	Plus	151	Uncharacterised protein (130)	Bordetella pertussis	37%	Unknown
			5	Plus	94	Uncharacterised protein (64)	Bordetella pertussis	82%	Unknown
			6	Plus	96	Uncharacterised protein (86)	Vibrio cholerae	48%	Unknown
			7	Plus	173	predicted protein (585)	Ruminococcus sp. CAG:403	26%	Unknown
pCAW4 (4002bp)	61%	KY939592	1	Minus	245	nosine monophosphate cyclohydrolase (246)	Ktedonobacter racemifer	63%	IMP cyclohydrolase
			2	Minus	214	phosphodiesterase (498)	candidate division NC10 bacterium	40%	phosphodiesterase
			3	Minus	402	hypothetical protein A2Y08_02680 (625)	Planctomycetes bacterium GWA2_40_7	43%	Unknown
			4 ^a	Plus	142	gentisate 1,2-dioxygenase (349)	Pseudomonas sp. 21C1	60%	gentisate 1,2-dioxygenase
pCAW5 (2724bp)	54%	KY939593	1 ^a	Plus	642	pyruvate:ferredoxin oxidoreductase (1565)	uncultured bacterium HF770_11D24] / Acidobacterium	80%	pyruvate:ferredoxin oxidoreductase
pCAW6 (2125bp)	57%	KY939594	1	Plus	159	hypothetical protein BGO39_33875 (215)	Chloroflexi bacterium 54-19	65%	MerR family
			2	Plus	336	hypothetical protein BGO39_33870 (347)	Chloroflexi bacterium 54-19	78%	PrsW intramembrane metalloprotease
			3 ^a	Plus	163	hypothetical protein BGO39_33865 (173)	Chloroflexi bacterium 54-19	75%	chromate transporter
pCAW7 (2558bp)	46%	KY939595	1 ^a	Minus	391	hypothetical protein A2X07_06330 (480)	Flavobacteria bacterium GWF1_32_7	45%	Por secretion system sorting domain
			2	Minus	250	hypothetical protein (586)	Chitinophagaceae bacterium PMP191F	65%	Polysaccharide Lyase
pCAW8 (4480bp)	57%	KY939596	1	Plus	508	hypothetical protein AUH20_02325 (597)	Rokubacteria bacterium	76%	5-oxoprolinase / Hydantoinase_B
			2	Minus	348	Oxidoreductase (336)	Rokubacteria bacterium	61%	Flavin-utilizing monooxygenases
			3	Plus	314	hypothetical protein ETSY1_46935 (279)	Candidatus Entotheonella sp. TSY1	76%	Cellulose biosynthesis BcsQ
pCAW9 (2573bp)	43%	KY939597	1 ^a	Minus	81	hypothetical protein (129)	Janthinobacterium	50%	Unknown
			2	Minus	303	Formylglycine-generating enzyme (379)	Mucilagibacter sp.	65%	Formylglycine-generating enzyme
			3	Minus	457	acetylglucosamine-6-sulfatase (504)	Flaviumibacter solisilvae	67%	acetylglucosamine-6-sulfatase

pCAW10 (2076bp)	56%	Submitted - Waiting for accession no.	1	Plus	204	hypothetical (195)	protein	Luminiphilus sylvensis	50%	Unknown
-----------------	-----	---	---	------	-----	-----------------------	---------	------------------------	-----	---------

1 ^aTruncated proteins.

2 ^baa, amino acids.

3 ^cSequences with an *E*-value higher than 0.001 in Blastp searches were considered to be unknown proteins.

4

5

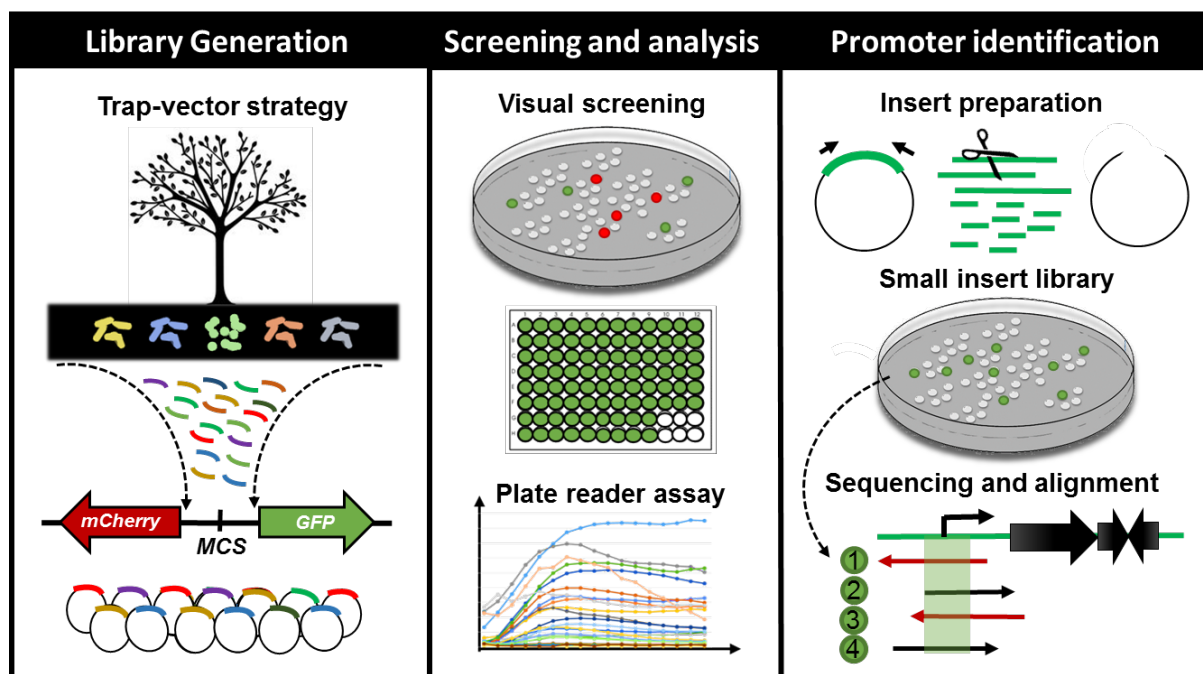
6

7

8

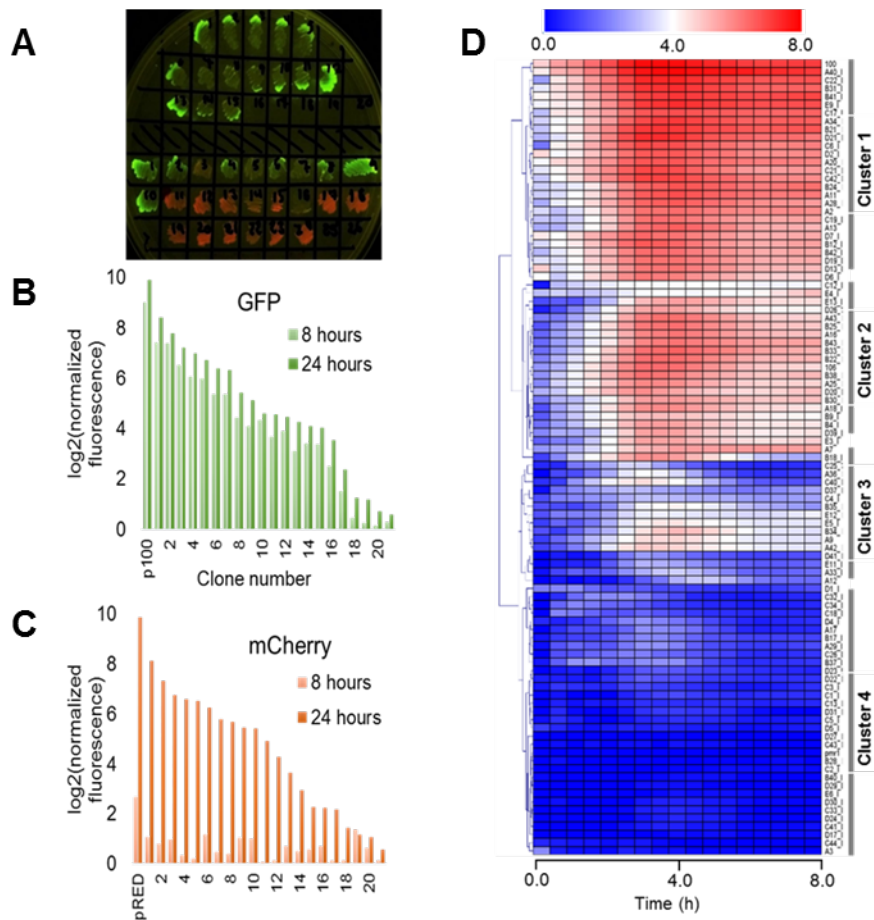
9

10 **Figures**



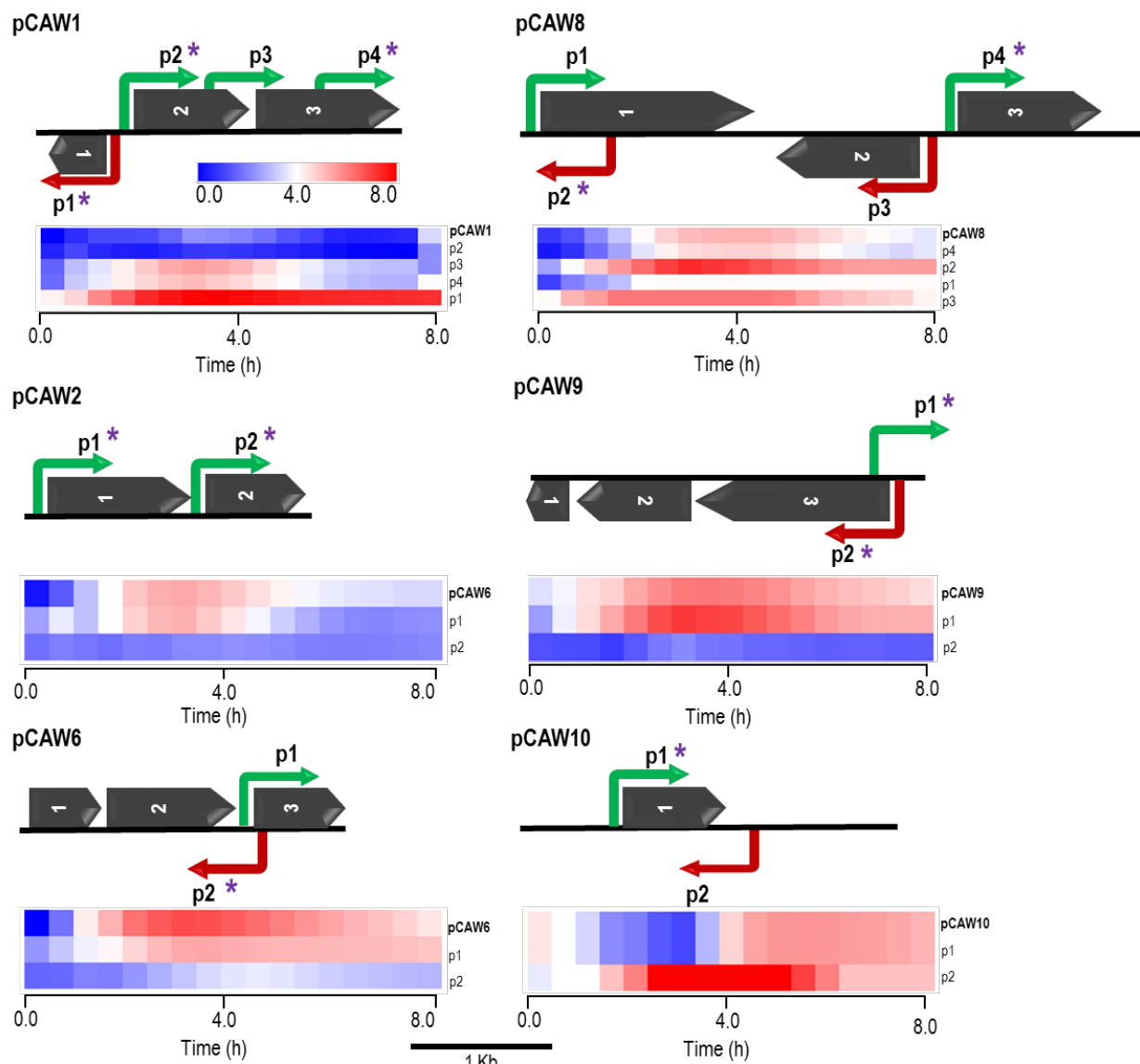
11

12 **Figure 1. Schematic representation of the workflow for finding, characterising and cross-validating novel**
 13 **bacterial cis-regulatory elements in environmental samples.** From left to right: firstly, we have generated
 14 metagenomic libraries from soil samples in *E. coli* DH10B. The DNA fragments were cloned into a bi-directional
 15 reporter trap-vector (bearing *mCherry* and *GFP* fluorescent reporters), pMR1, which allowed for the screening
 16 of promoters in both DNA strands. Secondly, we have manually screened all visible fluorescent clones from our
 17 metagenomic libraries and analysed the expression patterns of all green fluorescent clones on a microplate reader
 18 during 8 hours. Lastly, we have selected ten clones based on their *GFP* expression patterns for an in-depth
 19 analysis combining experimental (small DNA insert library generation) and *in silico* promoter prediction. This
 20 integrated strategy has allowed us to identify, validate and estimate the accessibility of novel promoter regions
 21 from metagenomic libraries.



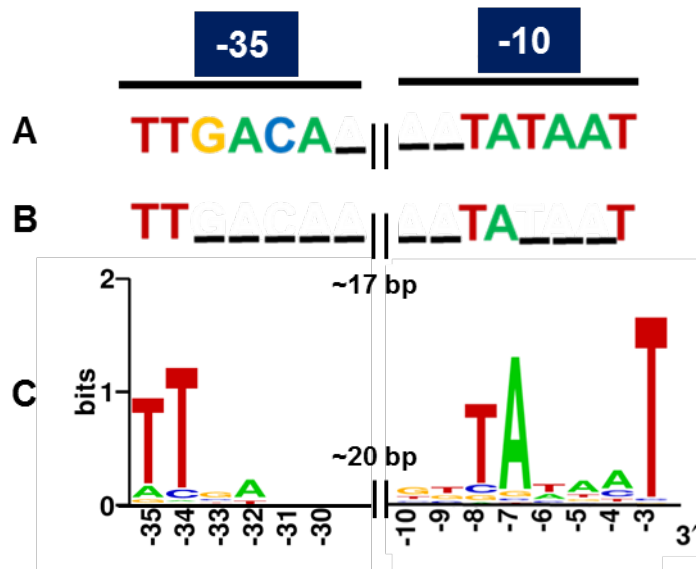
1

2 **Figure 2. Evaluating the expression dynamics of fluorescent clones.** A) LB-agar plate under blue light
3 light excitation comprising a subset of metagenomic isolated clones expressing GFP_{Iva} (top) and mCherry (bottom)
4 fluorescent reporters. A few clones were observed to express both reporters. All isolated clones were initially
5 considered to hold at least one endogenous promoter. B-C) Indirect assessment of maturation times from both
6 fluorescent reporters GFP_{Iva} (B) and mCherry (C) after 8 hours (light bars) and 24 hours (dark bars) of the
7 beginning of the experiment. Maturation times are substantially lower for mCherry than for GFP_{Iva}, which excluded
8 the former from further analyses. Positive controls for GFP and mCherry are represented by p100 and pRED,
9 respectively. Fluorescence data has been normalised by OD₆₀₀ values for each sample following normalisation by
10 values from the negative control (empty-pMR1). Data was transformed to log₂ scale to allow better visualisation of
11 fluorescence variation. D) Hierarchical representation of a metaconstitutome (i.e. all expression profiles from a
12 single metagenomic library). Fluorescence time-lapse dynamics were measured during 8 hours for each clone and
13 represented as heat maps. Promoter activities (calculated as GFP/OD₆₀₀) were normalised by the negative control
14 (*E. coli* DH10B harbouring empty pMR1) and transformed to log₂ scale in order to facilitate the visualisation of
15 subtle activities. Data are representative of three independent experiments.



1

2 **Figure 3. Schematic representation of six metagenomic inserts (contigs) showing predicted ORFs and**
 3 **experimentally validated/characterised promoters.** Each contig is identified on the far left of each subfigure.
 4 Promoters are indicated by elbow-shaped arrows and name according to their relative position in the contig.
 5 Promoter directionality, regarding the leading and lagging strands, is represented by green and red colours,
 6 respectively. Asterisks over specific promoters indicate regulatory regions which were cross-validated by matching
 7 *in silico* predictions. Dark arrows represent predicted ORFs, according to their relative positions in each contig (see
 8 Table 2 for more information). All genetic features respect their original relative sizes, following the 1 Kb scale
 9 depicted at the bottom of this figure. Beneath each metagenomic insert, there is a heat map cluster representing
 10 the whole set of promoter activities measured during 8-hours fluorescence assays. The first line of each cluster
 11 shows the original expression profile initially measured for each metagenomic insert. All other lines represent
 12 expression activities from *de novo* experimentally validated promoters within each contig (small DNA fragments).
 13 The second line of each cluster represents the endogenous promoter showing the most similar activity with respect
 14 to the original expression profile for each contig. All expression profiles are properly identified at the most rightmost
 15 side of each line, following their respective contig/promoter name. For the supplementary set of analysed contigs,
 16 see Supplementary Figure S4 online.



1

2 **Figure 4. Consensus of RpoD-related metagenomic promoters.** A) Known consensus sequences of the RpoD-
3 dependent promoter determined in vitro, TTGAAC (-35) and TATAAT (-10) separated by 17 plus/minus 2 bp in E.
4 coli ⁶⁵. B) Known consensus sequences of 582 promoters experimentally validated in E. coli ^{3,65,68}. C) The
5 sequences of the 33 promoters experimentally validated in this study were aligned and subjected to Logo analysis
6 ⁶⁶. The consensus from the metagenomic set (C) is very similar to the one from the experimentally validated set
7 from E. coli (B).