

1 **Supplementary Material for: “Long-read metagenomics of soil communities reveals**  
2 **phylum-specific secondary metabolite dynamics”**

3 Marc W. Van Goethem<sup>1</sup>, Andrew R. Osborn<sup>1</sup>, Benjamin P. Bowen<sup>1</sup>, Peter F. Andeer<sup>1</sup>, Tami L.  
4 Swenson<sup>1</sup>, Alicia Clum<sup>2</sup>, Robert Riley<sup>2</sup>, Guifen He<sup>2</sup>, Maxim Koriabine<sup>2</sup>, Laura Sandor<sup>2</sup>, Mi Yan<sup>2</sup>,  
5 Chris G. Daum<sup>2</sup>, Yuko Yoshinaga<sup>2</sup>, Thulani P. Makhalanyane<sup>3</sup>, Ferran Garcia-Pichel<sup>4,5</sup>, Ronan C.  
6 O’Malley<sup>1,2</sup>, Axel Visel<sup>2</sup>, Len A. Pennacchio<sup>2</sup>, Trent R. Northen<sup>1,2\*</sup>

7 <sup>1</sup> Molecular EcoSystems Biology Division, Lawrence Berkeley National Laboratory, 1 Cyclotron  
8 Rd, Berkeley, CA, 94720, USA

9 <sup>2</sup> DOE Joint Genome Institute, Lawrence Berkeley National Laboratory, 1 Cyclotron Rd, Berkeley,  
10 CA, 94720, USA

11 <sup>3</sup> Centre for Microbial Ecology and Genomics, Department of Biochemistry, Genomics and  
12 Microbiology, University of Pretoria, Lynnwood Rd, Hatfield. Pretoria, 0028, South Africa

13 <sup>4</sup> Center for Fundamental and Applied Microbiomics, Biodesign Institute, Arizona State University,  
14 Tempe, Arizona, USA

15 <sup>5</sup> School of Life Sciences, Arizona State University, Tempe, Arizona, USA

16 \* Correspondence and requests for materials should be addressed to T.R.N. (email:  
17 TRNorthen@lbl.gov)

18 Prepared for: *Communications Biology*

19 **Supplementary Note 1**

20 *Expression of Cyanobacterial Siderophore BGCs at Night*

21 Differential gene expression analysis (DESeq2 ;  $P < 0.05$ ; FDR = 5%) using mapped transcripts  
22 revealed that 10 BGCs contained biosynthetic genes that underwent significantly more  
23 transcription at night, all of which were cyanobacterial in origin. The most dramatic of these  
24 enrichments involved two cyanobacterial NRPS-PKS hybrid BGCs, identified on Node\_81 and  
25 Node\_86. These BGCs likely encode for a novel siderophore, putatively assigned based on genes  
26 encoding predicted membrane proteins involved in siderophore and iron transport located in the  
27 clusters (Fig. 2c). Additionally, a subset of cation acquisition genes was upregulated at night,  
28 suggesting a multifaceted approach for cation import at night that is under significant control of  
29 native regulatory constraints.

30 Node\_81 (68 kb), appears to form a novel heptapeptide, while Node\_86 (66 kb) is 2 kb shorter  
31 and has one fewer NRPS module, thus forming a hexapeptide (Fig. 2c). We speculated these to  
32 be rearranged BGCs based on the presence of transposases within Node\_86, which were  
33 supported by differences in G+C content flanking the transposases, potentially indicating recent  
34 transposition. Overall, 80% of transposases located in BGCs were active across all time points.  
35 We were also able to recover BGCs from the assembled metatranscriptomic data (Supplementary  
36 Data 5).

37

38 **Supplementary Note 2**

39 *Cation acquisition at night*

40 To further investigate the transcriptional activity of cation acquisition genes at night, we mapped  
41 reads to all the co-assembled metagenomes. A subset of putative cation acquisition and  
42 sequestration gene were differentially expressed, specifically *hemH* (Ferrochelatase), *hxuB*  
43 (Heme/hemopexin transporter protein), *pacS* (putative copper-transporting ATPase) and *idiA*  
44 (iron ABC transporter) genes (Supplementary Data 6). The differential expression of siderophore

45 and cation acquisition genes at night suggests a common 'night-time' strategy of cation import in  
46 biocrust, which was consistent across most cyanobacterial BGCs (Figs. S6, S7).





47

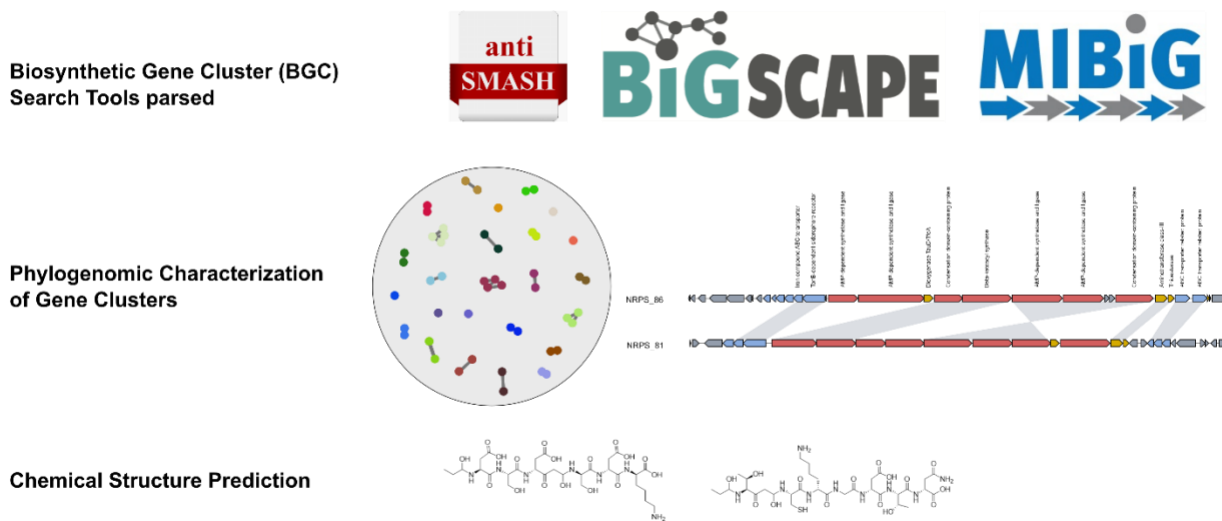
### 48 **Supplementary Note 3**

#### 49 *Transposases are prevalent among biocrust BGCs*

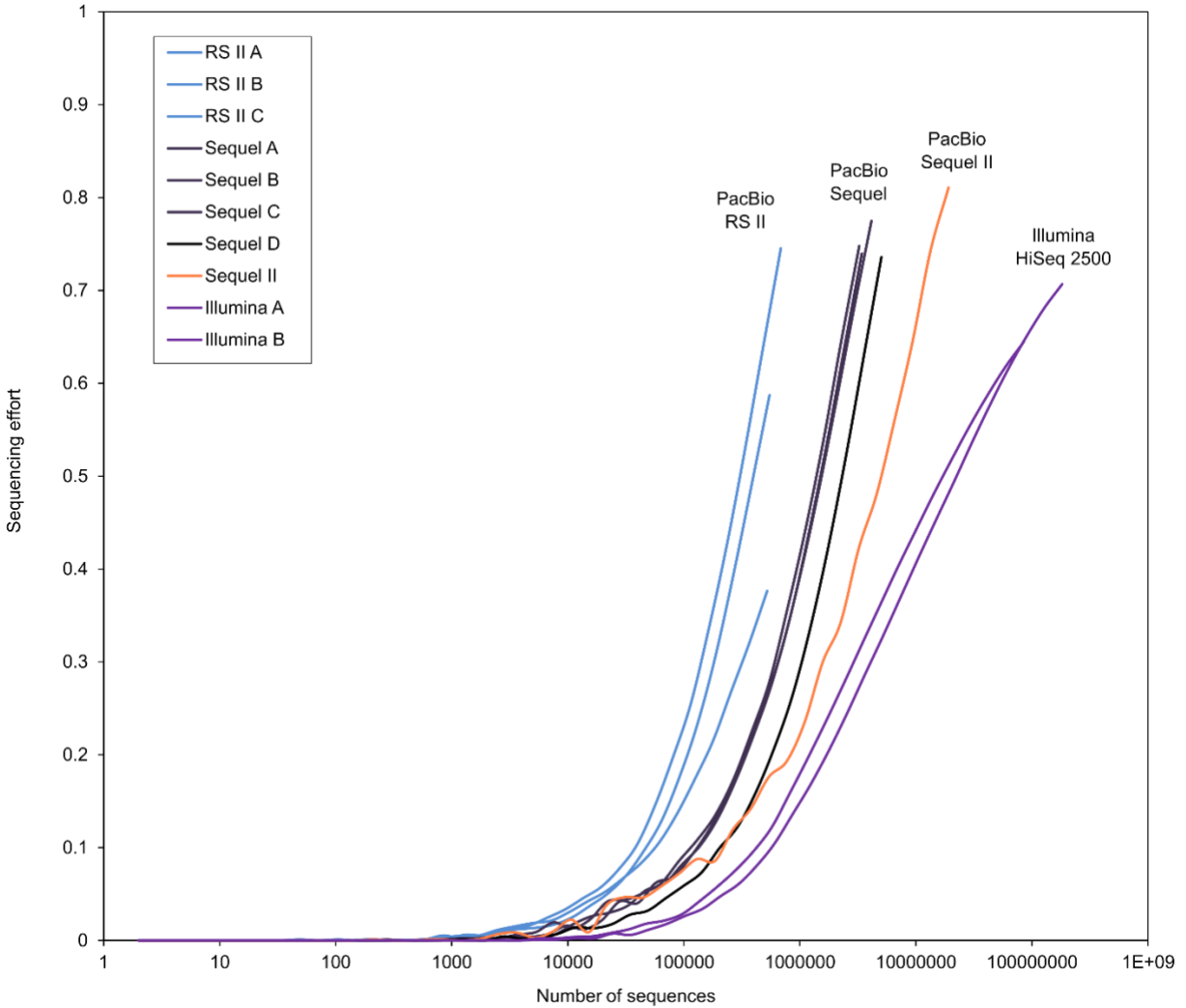
50 Two transposases and a phage integrase were found on Node\_86, with the transposases located  
51 upstream of the final NRPS gene in the cluster. G+C content surrounding the transposases within  
52 the BGC had high levels of G-C skew (>1.5% deviation from the mean) potentially indicative of a  
53 recent gene rearrangement. We hypothesized that these mobile elements were responsible for  
54 the rearrangement of Node\_81 into Node\_86. This is supported by the perpetual expression of  
55 the two transposases and phage integrase in Node\_86, implicating a mechanism for BGC  
56 rearrangement throughout the experiment and suggesting recombination events may still be  
57 occurring. We speculate that BGC re-arrangements may be on-going, long-term processes rather  
58 than brief, one-time events. Overall, 20% of transposases (Fig. S8) within BGCs were  
59 constitutively transcribed, hinting that similar long-term recombination events may be occurring in  
60 the biocrust community. Transposases located outside of BGCs, i.e., those more relevant to  
61 primary metabolism, showed less constitutive transcription (~7% of transposases) compared to  
62 those involved in secondary metabolism. Moreover, 26% of non-BGC transposases were never  
63 transcribed while only 19% of BGC transposases were never transcribed.

64 Nine of the 10 genes comprising Node\_86 showed differential expression, with significantly higher  
65 transcription at night. As a putative siderophore, the function of this metabolite could be to acquire  
66 iron at night in preparation for photosynthesis the following day. In contrast, Node\_81 only had  
67 one differentially expressed gene, but still tended towards nighttime activation.

		PACIFIC BIOSCIENCES®			illumina®	
		Metagenomics				
Sequencer						
		RS II	Sequel	Sequel II	HiSeq 2500	
No. of Metagenomes		3	4	1	2	
Assembler		Canu	metaFlye	metaFlye	metaSPAdes	
Co-assembler	Co-assemble Sequel (n=4) with metaFlye	Co-assemble Sequel (n=4) and HiSeq 2500 (n=2) with metaSPAdes		Co-assemble Sequel (n=4) and Sequel II with metaFlye		



68 **Supplementary Figure 1 | Workflow overview.** 8 long-read metagenomes were generated  
 69 using Pacific Biosciences instruments that span 3 generations (RS II -> Sequel -> Sequel II). 2  
 70 short-read metagenomes were generated from the same biocrust samples using an Illumina  
 71 HiSeq 2500 instrument.



72

73 **Supplementary Figure 2 | Estimates of sequencing effort.** Nonpareil estimates sequence  
 74 coverage by calculating read redundancy. The generations of long-read technology (RS II ->  
 75 Sequel -> Sequel II) lead to improvements in sequencing depth. Short-read sequencing required  
 76 orders of magnitude more sequencing to achieve a similar sequencing effort as RS II.

contig 3185  
(111 kb)



77

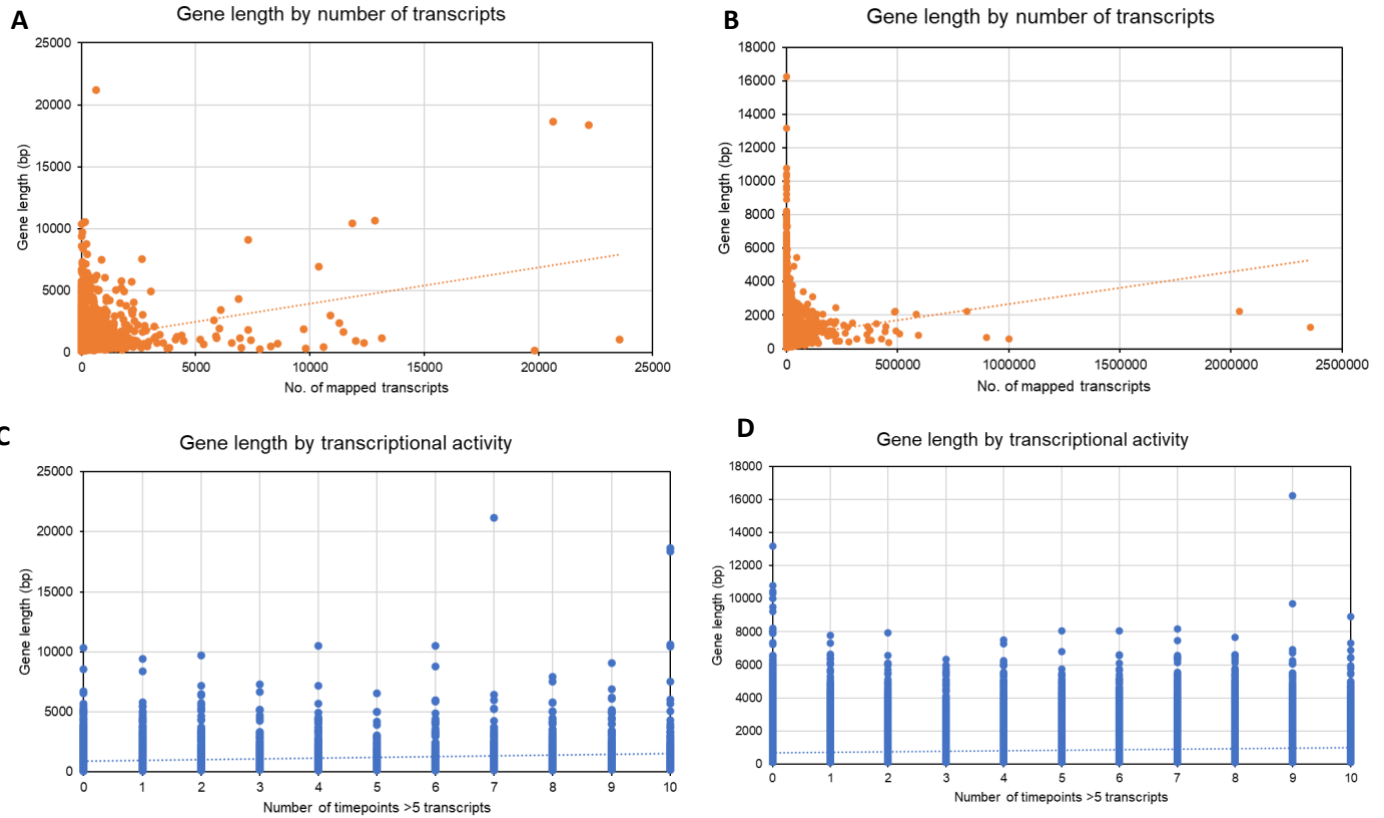
**Cluster 1**

polyketide  
nonribosomal peptide

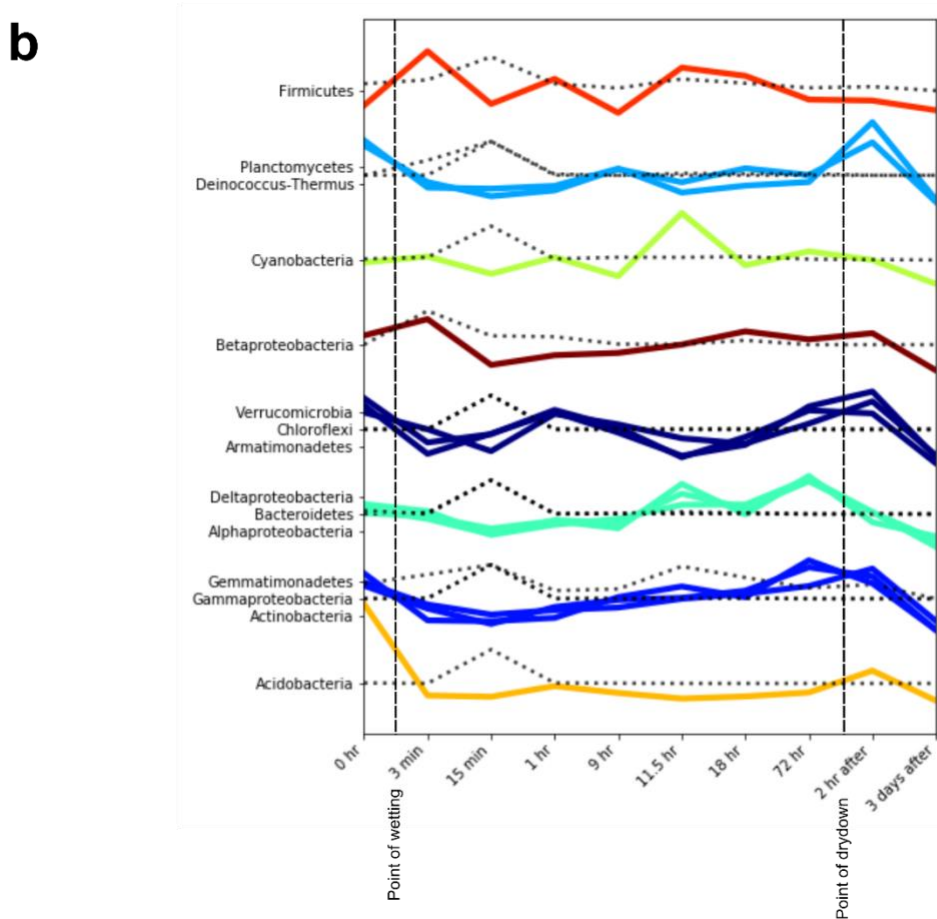
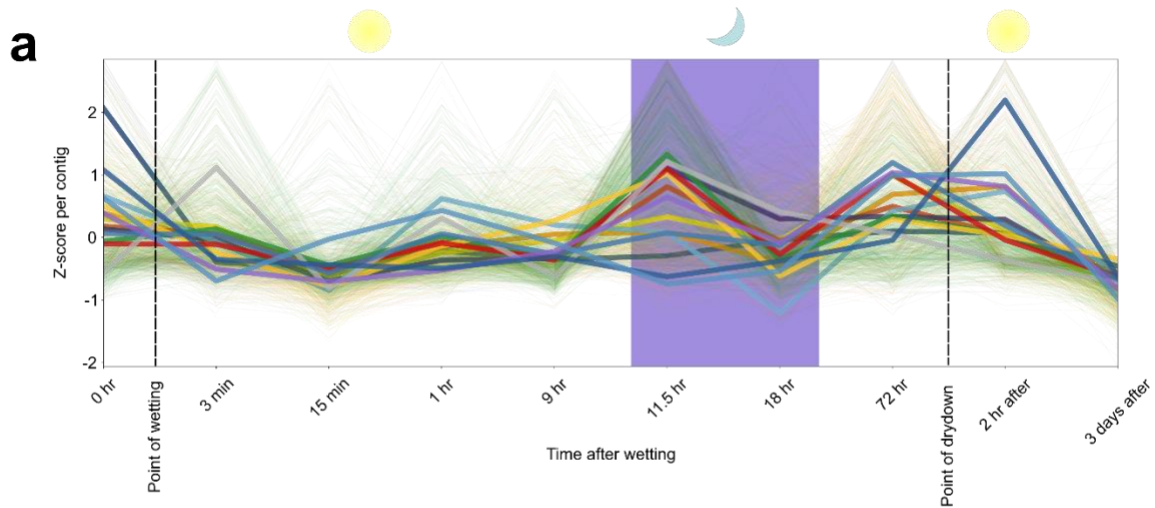


78

79 **Supplementary Figure 3 | Longest BGC recovered from a metagenome.** The longest BGC  
80 found across the metagenomes encodes a 111 kb transAT-PKS-NRPS. The domain architecture  
81 is provided by PRISM.



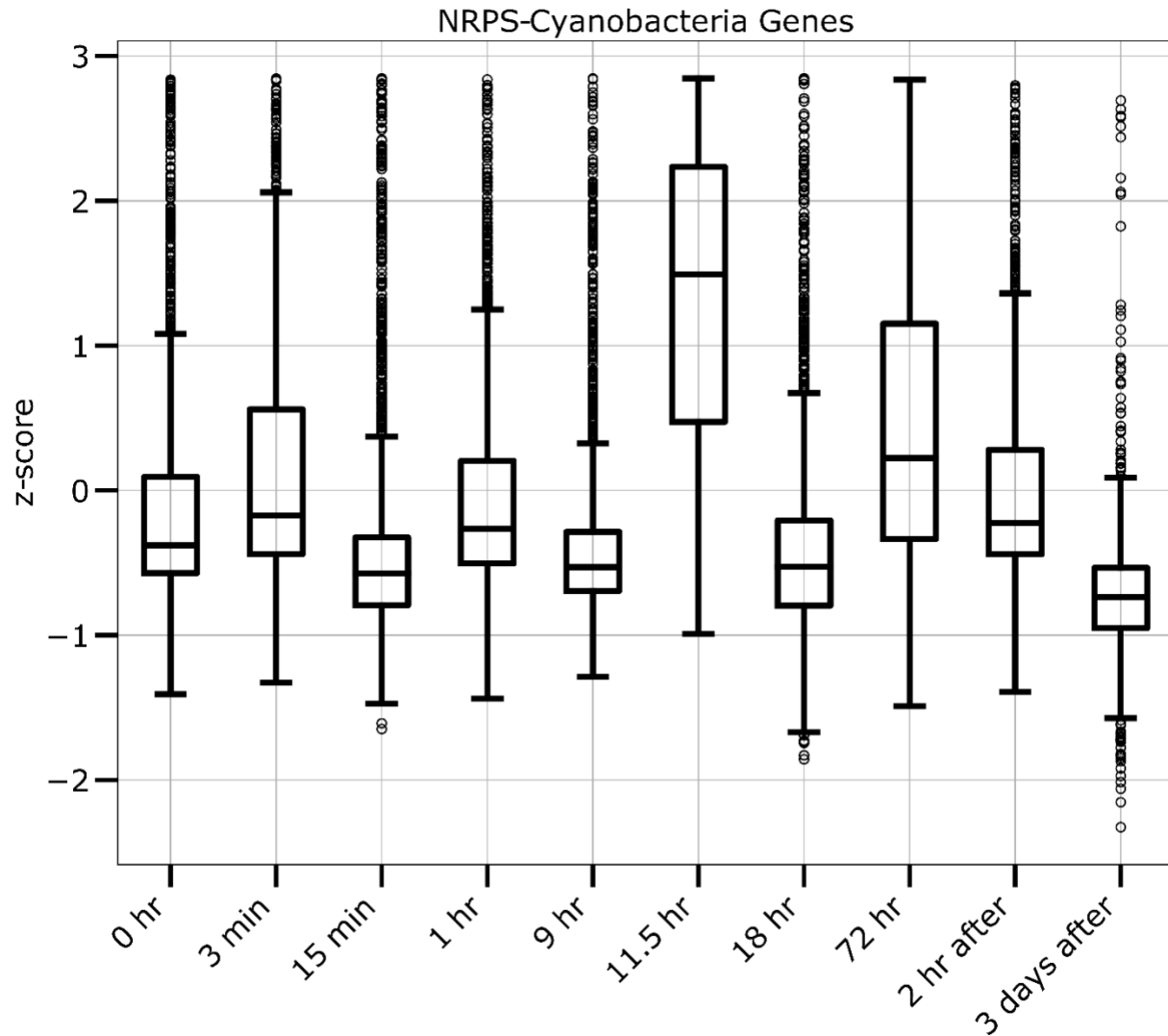
82 **Supplementary Figure 4 | Mapping transcripts as a function of gene length.** To test whether  
 83 gene length influenced mapping rates we compared how read recruitment differed for the longer  
 84 secondary metabolite genes (average gene length=1153 bp) compared to the primary metabolic  
 85 genes (average gene length=688 bp). Visualizing the number of mapped transcripts by gene  
 86 length showed no correlation for either **a)** secondary or **b)** primary metabolic genes. Similarly,  
 87 comparing the number of timepoints with 5 or more transcripts (where 0 means never expressed  
 88 and 10 means constitutive expression) to gene length indicated no trend that followed an increase  
 89 in transcript recruitment onto longer genes for either **c)** secondary or **d)** primary metabolic genes.



90 **Supplementary Figure 5 | Diurnal trends in BGCs.** a, Trends across the 3-day phase were  
 91 detected in 12,470 expressed biosynthetic genes (cutoff > 20 mapped transcripts across time  
 92 points). Read counts of transcript relative abundances per gene were Z-score normalized for the



93 purpose of visualization. Each gene trend is color-coded by its taxonomic affiliation. The purple  
94 background indicates night-time transcription. **b**, Clusters of bacterial phyla based on their  
95 average Z-score from all contigs with BGCs. The single-colored lines indicate secondary  
96 metabolism over time, while the dotted lines indicate the number of 16S rRNA transcripts at each  
97 timepoint.



98

99 **Supplementary Figure 6 | Cyanobacterial transcription occurs primarily at night.** All NRPS

100 gene clusters belonging to *Cyanobacteria* showed significantly more transcriptional activity 11.5

101 hours after wetting ( $P < 0.05$ ). This first dark time point indicates a shunt of secondary metabolism

102 by *Cyanobacteria* that also includes the increased transcription of RiPPs, T1PKS and T3PKS

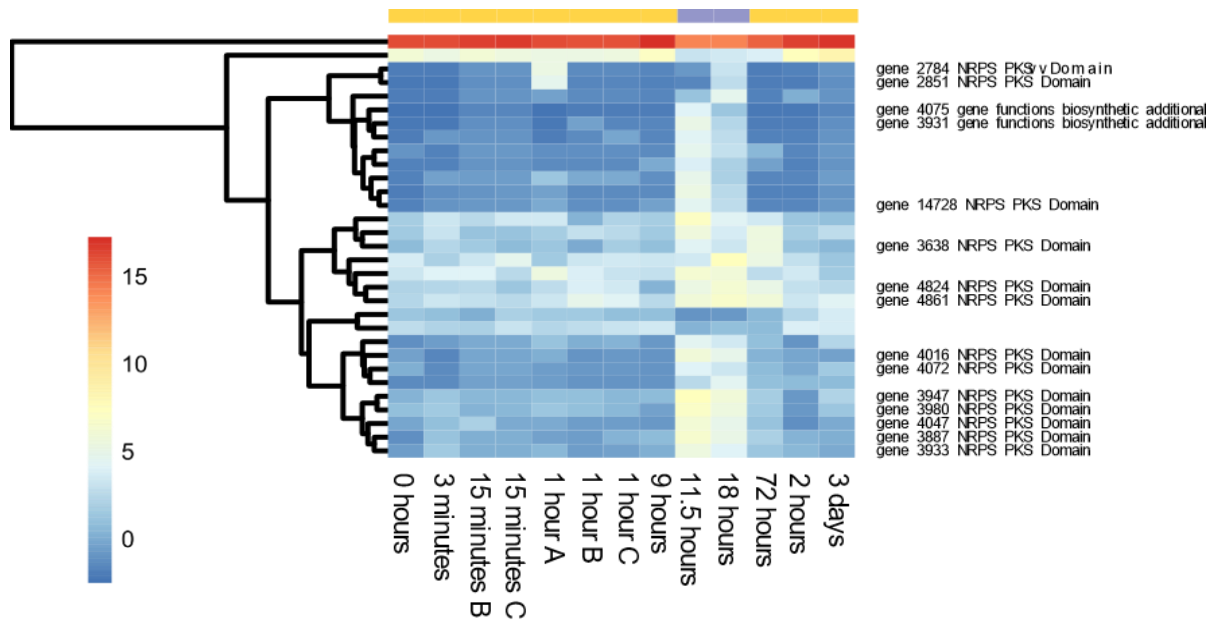
103 gene clusters. The boxes extend from Q1 to the Q3 quartile values of the data. The line is at the

104 median (Q2) while the whiskers extend from the edges of the box (Q1 and Q3) to show the range

105 of the data. The whiskers extend no more than 1.5 times the interquartile range (IQR = Q3 - Q1)

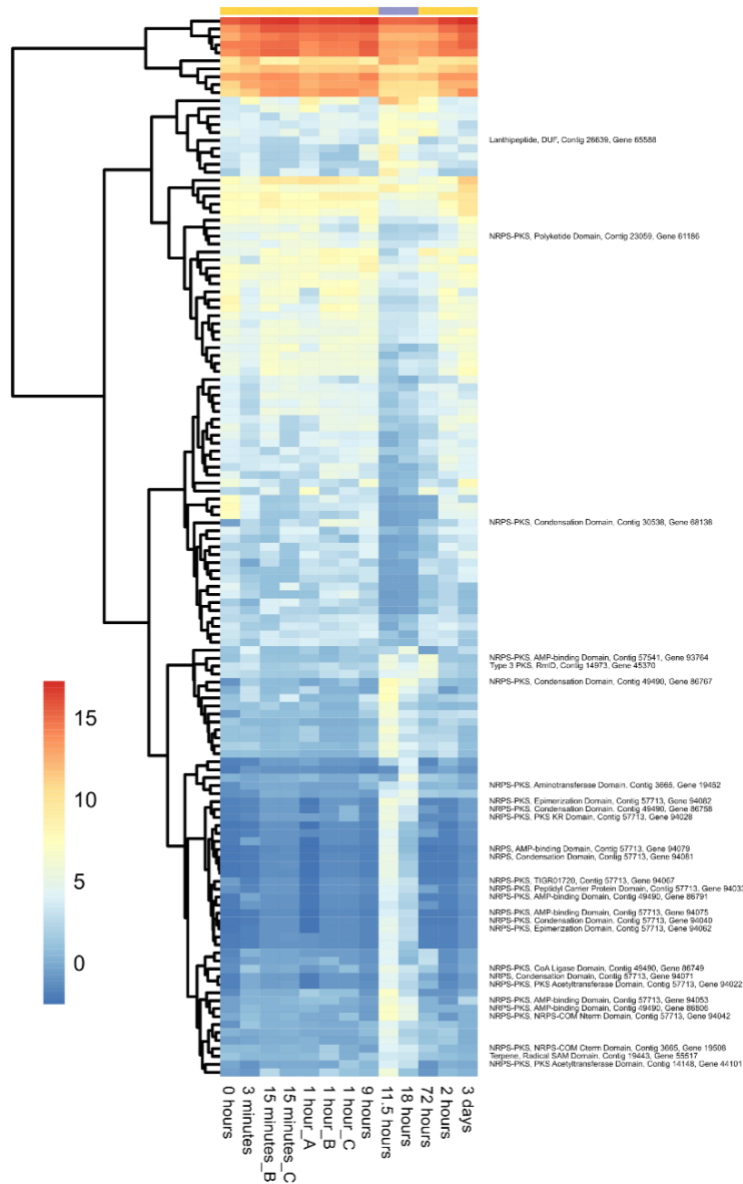
106 from the edges of the box and end at the farthest data point within that interval. All outliers are

107 plotted as separate dots beyond the whiskers.

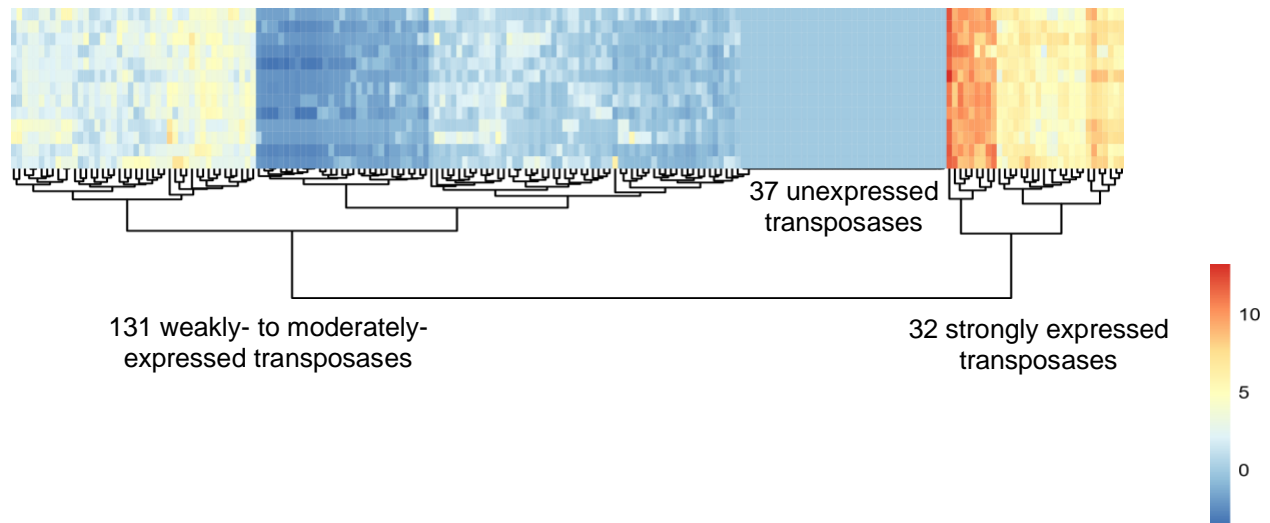


108 **Supplementary Figure 7 | DESeq2 showed significantly enriched gene transcription at**  
 109 **night.** Genes labelled on the heatmaps were those located within BGCs (Flye co-assembly), while  
 110 unlabeled rows are 'non-BGC' genes. Heatmap colors are based on DESeq2 comparisons  
 111 between night and day based on Log2Fold changes. Higher Log2Fold changes are shown in  
 112 warmer colors while cooler colors show less change between treatments. All genes are  
 113 significantly differentially transcribed at night. Left color axis indicates the condition, i.e. day  
 114 (yellow) or night (purple).

115



117 **Supplementary Figure 8 | DESeq2 showed significantly enriched gene transcription at**  
 118 **night.** Genes labelled on the heatmaps were those located within BGCs (Ultimate co-assembly),  
 119 while unlabeled rows are 'non-BGC' genes. Heatmap colors are based on DESeq2 comparisons  
 120 between night and day based on Log<sub>2</sub>Fold changes. Higher Log<sub>2</sub>Fold changes are shown in  
 121 warmer colors while cooler colors show less change between treatments. All genes are  
 122 significantly differentially transcribed at night. Left color axis indicates the condition, i.e. day  
 123 (yellow) or night (purple).



124

125 **Supplementary Figure 9 | Transposase transcription.** Heatmap showing the transcriptional  
 126 Log2Fold change of all transposases located in BGCs over time, with 0 hours the bottom row and  
 127 3 days after wetting the top row as in Supplementary Figure 7. We identified 3 broad categories  
 128 of expression: (i) unexpressed (no mapped transcripts), (ii) weak- to moderate-expression, and  
 129 (iii) strongly expressed transposases. Reds indicate higher levels of transcription while blues are  
 130 lowly-transcribed.

131 **Supplementary References**

- 132 1. Hiraoka, S. et al. Metaepigenomic analysis reveals the unexplored diversity of DNA methylation in  
133 an environmental prokaryotic community. *Nature communications* **10**, 159 (2019).
- 134 2. Frank, J.A. et al. Improved metagenome assemblies and taxonomic binning using long-read  
135 circular consensus sequence data. *Scientific reports* **6**, 25373 (2016).
- 136 3. Koren, S. et al. Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and  
137 repeat separation. *Genome research* **27**, 722-736 (2017).
- 138 4. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P.A. metaSPAdes: a new versatile  
139 metagenomic assembler. *Genome research* **27**, 824-834 (2017).
- 140 5. Kolmogorov, M., Rayko, M., Yuan, J., Pevnikov, E. & Pevzner, P. metaFlye: scalable long-read  
141 metagenome assembly using repeat graphs. *bioRxiv*, 637637 (2019).
- 142 6. Mikheenko, A., Saveliev, V. & Gurevich, A. MetaQUAST: evaluation of metagenome assemblies.  
143 *Bioinformatics* **32**, 1088-1090 (2015).
- 144 7. Blin, K. et al. antiSMASH 5.0: updates to the secondary metabolite genome mining pipeline.  
145 *Nucleic acids research* (2019).
- 146 8. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-  
147 seq data with DESeq2. *Genome biology* **15**, 550 (2014).