

# Diel gene expression improves software prediction of cyanobacterial operons

**Philip Heller** <sup>Corresp. 1</sup>

<sup>1</sup> Department of Computer Science, San Jose State University, San Jose, CA, United States

Corresponding Author: Philip Heller  
Email address: philip.heller@sjsu.edu

Cyanobacteria are important participants in global biogeochemical process, but their metabolic processes and genomic functions are incompletely understood. In particular, operon structure, which can provide valuable metabolic and genomic insight, is difficult to determine experimentally, and algorithmic operon predictions probably underestimate actual operon extent. A software method is presented for enhancing current operon predictions by incorporating information from whole-genome time-series expression studies, using a Machine Learning classifier. Results are presented for the marine cyanobacterium *Crocospaera watsonii*. 15 operon enhancements are proposed. Source code is publicly available.

# Diel gene expression improves software prediction of cyanobacterial operons

Philip Heller<sup>1</sup>

<sup>1</sup> Department of Computer Science, San Jose State University, San Jose, CA, USA

Corresponding Author:

Philip Heller<sup>1</sup> 1 Washington Sq., San Jose, CA 95192, USA

Email address: philip.heller@sjsu.edu

## Abstract

Cyanobacteria are important participants in global biogeochemical process, but their metabolic processes and genomic functions are incompletely understood. In particular, operon structure, which can provide valuable metabolic and genomic insight, is difficult to determine experimentally, and algorithmic operon predictions probably underestimate actual operon extent. A software method is presented for enhancing current operon predictions by incorporating information from whole-genome time-series expression studies, using a machine learning classifier. Results are presented for the marine cyanobacterium *Crocospaera watsonii*. 15 operon enhancements are proposed. Source code is publicly available.

## Introduction

Photosynthesizing bacteria (Phylum *Cyanobacteria*) are significant participants in global biogeochemical cycles. They arose on Earth 3.5 billion years ago<sup>1</sup>, and had oxygenated the atmosphere by 2.5 billion years ago<sup>2</sup>. Cyanobacteria participate in the ocean biological carbon pump<sup>3</sup>, which transports atmospheric greenhouse carbon dioxide to sequestration in the deep ocean. Nitrogen reducing cyanobacteria (diazotrophs) annually convert approximately 200 Tg of atmospheric dinitrogen to bioavailable form<sup>4 5</sup>. Cyanobacteria are used to produce medicines<sup>6 7 8</sup>, biofuels<sup>9 10</sup>, fertilizers<sup>11 12</sup>, cosmetics<sup>13</sup>, and food<sup>14</sup>.

Despite their ecological and commercial importance, the metabolic processes of many cyanobacteria have not been fully characterized; this is especially true for marine cyanobacteria, which are difficult to cultivate<sup>15</sup>. In particular, identification of operons (consecutive genes on the same DNA strand, controlled by a single promoter and expressed as a single transcript) appears to be incomplete. Operon identification provides clues for the inference of regulatory

pathways<sup>16 17</sup>, supports interpretation of transcriptome experiments<sup>18</sup>, and can guide annotation of hypothetical genes. The expense of wet lab operon discovery has prompted the development of algorithms for predicting operons from assembled genomes<sup>16 18 19</sup>; predictions from one of these algorithms<sup>19</sup> for 1336 organisms are publicly available (<http://www.microbesonline.org/operons/OperonList.html>). However, few of these predictions have been experimentally verified and it is possible that operon sizes have been underestimated.

Information for honing *in silico* operon predictions can be extracted from time-series measurements of gene expression. Many cyanobacterial genes are not expressed at constant rates, but rather exhibit fluctuating transcript abundance in repeating patterns over a 24-hour cycle. For example, production of light-harvesting photosystem II proteins, which are only useful during daylight and whose half-lives are generally less than 12 hours<sup>20 21</sup>, approximately coincides with available light<sup>22</sup>. Since oxygen disables nitrogenase (the enzyme responsible for nitrogen fixation), diazotrophic cyanobacteria segregate nitrogenase from the oxygen evolved by photosynthesis<sup>23 24</sup>; segregation is sometimes temporal, with nitrogenase component proteins produced hours out of phase from photosystem II proteins<sup>25</sup>. Diel cycling, defined as a transcript abundance change of at least 2x over 24 hours, has been observed in 79% of genes of the diazotrophic cyanobacterium *Crocospaera watsonii*<sup>26</sup>. Since genes in an operon are expected to have similar expression signatures<sup>27 28</sup>, a high degree of diel expression similarity among adjacent genes might indicate operon membership. Thus if two predicted operons are adjacent, are on the same DNA strand, and exhibit similar diel expression, then the predicted operons may in fact belong to a single common operon.

The approach presented here uses a machine learning classifier - specifically a Logistic Model Tree<sup>29 30</sup> (LMT) - to determine when predicted operons in *Crocospaera* should be merged. A common metric for quantifying expression similarity is Pearson's Correlation Coefficient (PCC); however, our earlier work<sup>31</sup> has determined that PCC has deficiencies when applied to the current problem. The "Area Between Linear Interpolations of Measurements" (ABLIM) metric, which we have presented elsewhere<sup>31</sup>, is more appropriate and is the basis of the research reported here. Based on the ABLIM metric, positive and negative example operons were located in the *Crocospaera watsonii* genome. 45 kinds of classifier (Supplemental Table 1) were evaluated, and LMT was selected due to its high accuracy. Adjacent predicted operons were identified as candidates for merging, and the expression similarity of all genes was analyzed by the classifier. 15 pairs of candidate operon predictions are recommended for merging (Table 1).

Source code and instructions are available at <https://github.com/PhilipHeller/Operons> (DOI 10.5281/zenodo.5759925).

# Materials & Methods

Computed operon predictions (hereafter the “prior predictions”) for strain *Crocospaera watsonii* were downloaded from <http://www.microbesonline.org/operons/>. Log-expression measurements for 4,407 *Crocospaera* genes with 8 timepoints were retrieved from a study by Shi et al.<sup>26</sup> For each gene, log-expressions were normalized to a mean of zero. A positive training set of operons for the classifier was collected by identifying all prior predicted operons in which at least one gene’s expression exhibited diel variation. A negative training set for the classifier was generated by identifying consecutive genes where at least one gene’s expression exhibited diel variation, and where at least one gene is on each DNA strand. (All genes of an operon must be on the same strand, to allow correct translation of the transcribed operon.)

The classifier requires training and evaluation instances to be represented by vectors of numbers. For each prior in the training sets (and, later, for each merge candidate to be classified), the ABLIM distance between every pair of genes was computed; the instance was represented by a 4-vector consisting of the minimum, mean, standard deviation, and maximum of the ABLIM distances. Gene distances within priors were not considered, as they did not fit a Gaussian distributions and are therefore problematic. 45 classifiers (Supplemental Table 1) in the WEKA software suite<sup>32 33</sup> were evaluated on the positive and negative sets using five-fold cross-validation. The Logistic Model Tree (LMT) classifier gave the best accuracy on both the positive and negative data, and was therefore selected for the remainder of the study. The classifier was trained using all the positive and negative instances.

Pairs of prior predictions were identified as candidates for merging (Table 1) if there were no intervening genes, if all genes lay on the same DNA strand and in the same contig, and if each prior contained at least one gene whose expression exhibited diel variation. A 4-vector representation of each candidate was computed as described above, and the representations were evaluated on the trained LMT classifier to generate classification scores (Figure 1). A candidate was accepted (i.e. all its genes are predicted to be in a single operon) if classifier score was > 0.5. Note that this score is not to be interpreted as a probability that the classification is correct.

Software was developed on Eclipse 2020-12 (4.18.0) in Java SE-15 (and is compatible with Java 1.8), using version 3.5 of the WEKA library and version 3.6.1 of the Apache Commons Mathematics Library.

# Results

The positive training set consists of the 1195 operon predictions at <http://www.microbesonline.org/operons/>. The negative training set is listed in Supplemental

Table 1. 45 classifiers in the WEKA software were evaluated on the training data. The Logistical Model Tree (LMT) had the highest accuracy (Supplemental Table 1).

63 pairs of prior operon predictions were identified as candidates for merging. Each prior consisted of two genes, at least one of which exhibited diel expression variation; all genes were on the same DNA strand and in the same contig, and there were no intervening genes between the priors. 15 pairs of priors were classified as belonging to a common operon (Table 1).

## Discussion

Diel expression data was combined with prior operon predictions to compute 15 pairs of priors (Table 1) that appear to belong to common operons. It is recommended that each of these pairs be merged into a single prediction.

One reason for honing operon predictions is to gain insight into the function of unknown genes. When unknown genes share an operon with genes of known function, the known function can reasonably be hypothesized to relate to the unknown functions. In Table 1, unknown genes are marked in bold. Six prior predictions include operons where no gene has known function; in all these cases, the present analysis predicts that the prior prediction should be merged with another prior containing at least one gene of known function. Predicted operon membership *per se* may not be strong enough evidence to infer gene function, but it can provide the basis for hypothesizing function, and the hypothesis can be strengthened by other evidence.

Each operon (training priors and merge candidates) was represented by a 4-vector consisting of the minimum, mean, standard deviation, and maximum of the ABLIM distances among all gene pairs in the operon. None of these statistics alone was sufficient for training an accurate classifier. The LMT classifier had the best accuracy among the 45 classifiers that were evaluated (Supplemental Table 1). However, this does not imply that LMT should be used when analyzing other organisms. Future work on other organisms should repeat the classifier evaluation reported here, and should choose the best classifier for the organism at hand.

## Conclusions

The work presented here demonstrates that machine learning analysis of diel expression studies can improve *in silico* predictions of operons. When a prior prediction is extended to include genes of unknown function, the function of the known genes in the prior might elucidate the function of the new unknown genes.

The approach presented here can be applied to other cyanobacteria for which diel studies and prior predicted operons are available. Since the method is based on similarity of diel signatures,

best results should be expected from organisms whose genes exhibit strong and diverse diel variation. Organisms with weak diel variation can be expected to perform poorly, because the 4-vectors that describe operons to the classifier would all be similar. Experiments with a diel study<sup>34</sup> of the minimal bacterium *Prochlorococcus marinus* produced poor results with the approach presented here, possibly because the circadian clock mechanism is simplified in *Prochlorococcus*<sup>35</sup> and its diel genes fluctuate more weakly than those of *Crocospaera*.

## Acknowledgements

The author is grateful to Jonathan Zehr, Josh Stuart, Irina Shilova, Rex Malmstrom, and Laurence Nedelec for valuable discussions.

## References

1. Schopf, J. In: Whitton B, Potts M (eds) The fossil record: tracing the roots of the cyanobacterial lineage. in *The Ecology of Cyanobacteria* 13–35 (Kluwer, 2002).
2. Tomitani, A., Knoll, A. H., Cavanaugh, C. M. & Ohno, T. The evolutionary diversification of cyanobacteria: Molecular-phylogenetic and paleontological perspectives. *Proceedings of the National Academy of Sciences* **103**, 5442–5447 (2006).
3. Tréguer, P. *et al.* Influence of diatom diversity on the ocean biological carbon pump. *Nature Geosci* **11**, 27–37 (2018).
4. Wang, W.-L., Moore, J. K., Martiny, A. C. & Primeau, F. W. Convergent estimates of marine nitrogen fixation. *Nature* **566**, 205–211 (2019).
5. Tang, W., Li, Z. & Cassar, N. Machine Learning Estimates of Global Marine Nitrogen Fixation. *Journal of Geophysical Research: Biogeosciences* **124**, 717–730 (2019).
6. Tan, L. T. Bioactive natural products from marine cyanobacteria for drug discovery. *Phytochemistry* **68**, 954–979 (2007).
7. Soni, B., Trivedi, U. & Madamwar, D. A novel method of single step hydrophobic interaction chromatography for the purification of phycocyanin from *Phormidium fragile* and its characterization for antioxidant property. *Bioresource Technology* **99**, 188–194 (2008).

8. Zanchett, G. & Oliveira-Filho, E. Cyanobacteria and Cyanotoxins: From Impacts on Aquatic Ecosystems and Human Health to Anticarcinogenic Effects. *Toxins* **5**, 1896–1917 (2013).
9. Sakurai, H., Masukawa, H., Kitashima, M. & Inoue, K. How Close We Are to Achieving Commercially Viable Large-Scale Photobiological Hydrogen Production by Cyanobacteria: A Review of the Biological Aspects. *Life* **5**, 997–1018 (2015).
10. Farrokh, P., Sheikhpour, M., Kasaeian, A., Asadi, H. & Bavandi, R. Cyanobacteria as an eco-friendly resource for biofuel production: A critical review. *Biotechnol. Prog.* **35**, (2019).
11. Singh, H., Khattar, J. S. & Ahluwalia, A. S. Cyanobacteria and agricultural crops. *Vegetos- An Inter. Jour. of Plnt. Rese.* **27**, 37 (2014).
12. Chittapun, S., Limbipichai, S., Amnuaysin, N., Boonkerd, R. & Charoensook, M. Effects of using cyanobacteria and fertilizer on growth and yield of rice, Pathum Thani I: a pot experiment. *J Appl Phycol* **30**, 79–85 (2018).
13. Morone, J., Alfeus, A., Vasconcelos, V. & Martins, R. Revealing the potential of cyanobacteria in cosmetics and cosmeceuticals — A new bioactive approach. *Algal Research* **41**, 101541 (2019).
14. Khan, Z., Bhadouria, P. & Bisen, P. Nutritional and Therapeutic Potential of Spirulina. *CPB* **6**, 373–379 (2005).
15. Zehr, J. P. Nitrogen fixation by marine cyanobacteria. *Trends in Microbiology* **19**, 162–173 (2011).
16. Zheng, Y., Szustakowski, J. D., Fortnow, L., Roberts, R. J. & Kasif, S. Computational Identification of Operons in Microbial Genomes. *Genome Research* **12**, 1221–1230 (2002).
17. Westover, B. P., Buhler, J. D., Sonnenburg, J. L. & Gordon, J. I. Operon prediction without a training set. *Bioinformatics* **21**, 880–888 (2005).
18. Moreno-Hagelsieb, G. & Collado-Vides, J. A powerful non-homology method for the prediction of operons in prokaryotes. *Bioinformatics* **18**, S329–S336 (2002).
19. Price, M. N. A novel method for accurate operon predictions in all sequenced prokaryotes. *Nucleic Acids Research* **33**, 880–892 (2005).
20. Yao, D. C. I., Brune, D. C., Vavilin, D. & Vermaas, W. F. J. Photosystem II Component Lifetimes in the Cyanobacterium *Synechocystis* sp. Strain PCC 6803: SMALL Cab-LIKE PROTEINS STABILIZE BIOSYNTHESIS INTERMEDIATES AND AFFECT EARLY STEPS IN CHLOROPHYLL SYNTHESIS. *Journal of Biological Chemistry* **287**, 682–692 (2012).
21. Renger, G. *et al.* ON THE MECHANISM OF PHOTOSYSTEM II DETERIORATION BY UV-B IRRADIATION. *Photochemistry and Photobiology* **49**, 97–105 (1989).
22. Dodd, A. N. Plant Circadian Clocks Increase Photosynthesis, Growth, Survival, and Competitive Advantage. *Science* **309**, 630–633 (2005).
23. Bergersen, F. J. The Effects of Partial Pressure of Oxygen upon Respiration and Nitrogen Fixation by Soybean Root Nodules. *Journal of General Microbiology* **29**, 113–125 (1962).
24. Fay, P. Oxygen Relations of Nitrogen Fixation in Cyanobacteria. *MICROBIOL. REV.* **56**, 34 (1992).

25. Tuit, C., Waterbury, J. & Ravizza, G. Diel variation of molybdenum and iron in marine diazotrophic cyanobacteria. *Limnology and Oceanography* **49**, 978–990 (2004).
26. Shi, T., Ilikchyan, I., Rabouille, S. & Zehr, J. P. Genome-wide analysis of diel gene expression in the unicellular N<sub>2</sub>-fixing cyanobacterium *Crocospaera watsonii* WH 8501. *The ISME Journal* **4**, 621–632 (2010).
27. Sabatti, C. Co-expression pattern from DNA microarray experiments as a tool for operon prediction. *Nucleic Acids Research* **30**, 2886–2893 (2002).
28. Lercher, M. J. Coexpression of Neighboring Genes in *Caenorhabditis Elegans* Is Mostly Due to Operons and Duplicate Genes. *Genome Research* **13**, 238–243 (2003).
29. Landwehr, N., Hall, M. & Frank, E. Logistic Model Trees. *Mach Learn* **59**, 161–205 (2005).
30. Sumner, M., Frank, E. & Hall, M. Speeding Up Logistic Model Tree Induction. in *Knowledge Discovery in Databases: PKDD 2005* (eds. Jorge, A. M., Torgo, L., Brazdil, P., Camacho, R. & Gama, J.) vol. 3721 675–683 (Springer Berlin Heidelberg, 2005).
31. Heller, P. & Baiju, Bharath. An Improved Distance Metric For Clustering Gene Expression Time-Series Data. *Am J Adv Res* (2018) doi:10.5281/zenodo.1419822.
32. Hall, M. *et al.* The WEKA Data Mining Software: An Update. **11**, 9 (2009).
33. Frank, E., Hall, M. & Witten, I. The WEKA Workbench. in *Data Mining: Practical Machine Learning Tools and Techniques* (Morgan Kaufmann, 2016).
34. Zinser, E. R. *et al.* Choreography of the Transcriptome, Photophysiology, and Cell Cycle of a Minimal Photoautotroph, *Prochlorococcus*. *PLoS ONE* **4**, e5135 (2009).
35. Holtzendorff, J. *et al.* Genome Streamlining Results in Loss of Robustness of the Circadian Clock in the Marine Cyanobacterium *Prochlorococcus marinus* PCC 9511. *J Biol Rhythms* **23**, 187–199 (2008).

# Table 1 (on next page)

Pairs of predicted operons recommended for merging.

Each row presents a consecutive pair of previously predicted operons. The score is generated by a Logistic Model Tree classifier. Each pair presented here has classifier score  $> 0.5$ , and therefore should likely be merged into a single longer predicted operon. The numbers in the “Prior 1” and “Prior 2” columns are gene identifiers, truncated for formatting; prepend “CwatDRAFT\_” to the numbers to generate the full identifier. In the “Gene Functions” columns, a question mark denotes a gene of unknown function; known function of other genes in a predicted operon can provide clues to the unknown function.

Score	Prior 1	Prior 1 Gene Functions	Prior 2	Prior 2 Gene Functions
0.868	2207 <b>2208</b>	HAD-superfamily hydrolase, subfamily IA, var <b>No predicted function</b>	<b>2209</b> <b>2210</b>	<b>No predicted function</b> <b>No predicted function</b>
0.834	<b>3078</b> <b>3079</b>	<b>No predicted function</b> <b>No predicted function</b>	3080 3081	Glucose-6-phosphate dehydrogenase OpcA
0.673	4526 4527	Hydrogenase expression/synthesis, HypA Hydrogenase accessory protein HypB	<b>4528</b> <b>4529</b>	<b>No predicted function</b> <b>No predicted function</b>
0.672	3989 3990	extracellular solute-binding protein, family 3 Amino acid ABC transporter, permease protein	<b>3991</b> <b>3992</b>	<b>No predicted function</b> <b>No predicted function</b>
0.631	<b>5385</b> 5386	<b>No predicted function</b> Cytochrome c oxidase, subunit II: Cytochrome	5388 5389	Cytochrome-c oxidase Cytochrome c oxidase, subunit 3
0.618	1168 1169	3-isopropylmalate dehydratase small subunit DegT/DnrJ/EryC1/StrS aminotransferase	<b>1165</b> 1166	<b>No predicted function</b> Glutathione S-transferase, N-term
0.599	4216 <b>4217</b>	Pentapeptide repeat <b>No predicted function</b>	<b>4214</b> 4215	<b>No predicted function</b> HEAT:PBS lyase HEAT-like rpt
0.576	2640 2641	Phosphopantethiene-protein transferase ATP-binding region, ATPase-like:Histidin	2638 2639	K+ channel, pore region K+ channel, pore region
0.565	<b>6744</b> <b>6745</b>	<b>No predicted function</b> <b>No predicted function</b>	6742 6743	Competence-damaged protein:CinA, C-trmnl Uracil phosphoribosyltransferase
0.56	4082 <b>4083</b>	Carbamoyltransferase <b>No predicted function</b>	<b>4084</b> <b>4085</b>	<b>No predicted function</b> <b>No predicted function</b>
0.552	2855 <b>2856</b>	Glycosyl transferase, group 1 <b>No predicted function</b>	2853 2854	Phycobilisome linker polypeptide Ferredoxin-dependent bilin reductase
0.537	5116 5117	GTP cyclohydrolase I Cobalamin synthesis protein/P47K:Cobalami	5114 5115	Cobalamin synthesis protein/P47K:Cobalami Dihydrouridine synthase TIM-barrel protein yjbN
0.537	2158 2159	Ribosomal protein L33 Ribosomal protein S18	2160 <b>2161</b>	Exoribonuclease II <b>No predicted function</b>
0.518	3464 <b>3465</b>	Porphobilinogen synthase <b>No predicted function</b>	<b>3466</b> 3467	<b>No predicted function</b> Metallophosphoesterase
0.514	3440 <b>3441</b>	Hemolysin-type calcium-binding region <b>No predicted function</b>	<b>3438</b> 3439	<b>No predicted function</b> Quinate/Shikimate 5-dehydrogenase

1

# Figure 1

Classification algorithm.

Prior predicted operons (red, blue) are candidates for merging if their genes are consecutive, if all genes are on the same DNA strand, and if at least one gene in each prior exhibits diel expression.

