# Diel gene expression improves software prediction of cyanobacterial operons

**Philip Heller** [Corresp. 1]

[1] 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 *Crocosphaera watsonii*. 22 operon enhancements are proposed.

1

# Diel gene expression improves software prediction of cyanobacterial operons

Philip Heller[1]

[1] Department of Computer Science, San Jose State University, San Jose, CA, USA

Corresponding Author:
Philip Heller[1]  1 Washington Sq., San Jose, CA 95192, USA
Email address: philip.heller@sjsu.edu

## Abstract

Add your abstract here. 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 *Crocosphaera watsonii*. 22 operon enhancements are proposed.

## Introduction

Photosynthesizing bacteria (Phylum *Cyanobacteria*) are significant participants in global biogeochemical cycles. They arose on Earth 3.5 billion years ago[1], and had oxygenated the atmosphere by 2.5 billion years ago[2]. Cyanobacteria participate in the ocean biological carbon pump[3], 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[4][5]. Cyanobacteria are used to produce medicines[6][7][8], biofuels[9][10], fertilizers[11][12], cosmetics[13], and food[14].

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[15]. In particular, identification of operons (consecutive genes controlled by a single promoter and expressed as a single transcript) appears to be incomplete. Operon identification provides clues for the

39   inference of regulatory pathways[16 17], supports interpretation of transcriptome
40   experiments[18], and can guide annotation of hypothetical genes. The expense of wetlab
41   operon discovery has prompted the development of algorithms for predicting operons
42   from assembled genomes[16 18 19]; predictions from one of these algorithms[19] for 1336
43   organisms are publicly available
44   (http://www.microbesonline.org/operons/OperonList.html). However, few of these
45   predictions have been experimentally verified and it is possible that operon sizes have
46   been underestimated.

47

48   Information for honing *in silico* operon predictions can be extracted from time-series
49   measurements of gene expression. Many cyanobacterial genes are not expressed at
50   constant rates, but rather exhibit fluctuating transcript abundance in repeating patterns
51   over a 24-hour cycle. For example, production of light-harvesting photosystem II
52   proteins, which are only useful during daylight and whose half-lives are generally less
53   than 12 hours[20 21], approximately coincides with available light[22]. Since oxygen disables
54   nitrogenase (the enzyme responsible for nitrogen fixation), diazotrophic cyanobacteria
55   segregate nitrogenase from the oxygen evolved by photosynthesis[23 24]; segregation is
56   sometimes temporal, with nitrogenase component proteins produced hours out of phase
57   from photosystem II proteins[25]. Diel cycling, defined as a transcript abundance change
58   of at least 2x over 24 hours, has been observed in 79% of genes of the diazotrophic
59   cyanobacterium *Crocosphaera watsonii*[26]. Since genes in an operon are expected to
60   have similar expression signatures[27 28], a high degree of diel expression similarity
61   among adjacent genes might indicate operon membership. Thus if two predicted
62   operons are adjacent, are on the same DNA strand, and exhibit similar diel expression,
63   then the predicted operons may in fact belong to a single common operon.
64
65   The approach presented here uses a Machine Learning classifier - specifically a
66   Logistic Model Tree[29 30] (LMT) - to determine when predicted operons in *Crocosphaera*
67   should be merged. A common metric for quantifying expression similarity is Pearson's
68   Correlation Coefficient (PCC); however, our earlier work[31] has determined that PCC has
69   deficiencies when applied to the current problem. The "Area Between Linear
70   Interpolations of Measurements" (ABLIM) metric, which we have presented elsewhere[31],
71   is more appropriate and is the basis of the research reported here. Based on the ABLIM
72   metric, positive and negative example operons were located in the *Crocosphaera*
73   *watsonii* genome. 48 kinds of classifier (Supplemental Table 1) were evaluated, and
74   LMT was selected due to its high accuracy. Adjacent predicted operons were identified
75   as candidates for merging, and the expression similarity of all genes was analyzed by
76   the classifier. 22 pairs of candidate operon predictions are recommended for merging
77   (Table 1).
78

79

80

81

## Materials & Methods

Computed operon predictions (hereafter the "prior predictions") for strain *Crocosphaera*
*watsonii* were downloaded from http://www.microbesonline.org/operons/. Log-expression
measurements for 4,407 *Crocosphaera* genes with 8 timepoints were retrieved from a study by
Shi et al.[26] 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 1 gene's expression exhibited diel variation. A negative training set for the classifier was
generated by identifying consecutive genes where at least 1 gene's expression exhibited diel
variation, and where each DNA strand is represented. (Since operons are transcribed as a single
unit, and transcription is restricted to a single strand, these sets of genes cannot be operons.)

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. 48 classifiers (Supplemental Table 1) in the WEKA software suite[32][33] were evaluated
on the positive and negative sets using 5-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 (Supplemental Table 2) 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 1 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.

To estimate the accuracy of the classifier's predictions, each negative training example in turn
was censored from the training set; the model was then re-trained on the remaining data, and the
censored example's classification score was computed. A Gaussian distribution was computed
for the classification scores thus generated. Given a candidate with score s, the cumulative
probability of scores >= s is an estimate of the probability of erroneously accepting the
candidate. Table 1 lists the accepted predictions, with their classifier scores and estimated error
probabilities.

118

## Results

119
120 The positive training set consists of the 1195 operon predictions at
121 http://www.microbesonline.org/operons/. The negative training set is listed in Supplemental
122 Table 1. 48 classifiers in the WEKA software were evaluated on the training data. The Logistical
123 Model Tree (LMT) had the highest accuracy (Supplemental Table 1).
124
125 79 pairs of prior operon predictions were identified as candidates for merging. Each prior
126 consisted of 2 genes, at least 1 of which exhibited diel expression variation; all genes were on the
127 same DNA strand and in the same contig, and there were no intervening genes between the 2
128 priors. 22 pairs of priors were classified as belonging to a common operon (Table 1).
129

## Discussion

131 Diel expression data was combined with prior operon predictions to compute 22 pairs of priors
132 (Table 1) that appear to belong to common operons. It is recommended that each of these pairs
133 be merged into a single prediction.
134
135 One reason for honing operon predictions is to gain insight into the function of unknown genes.
136 When unknown genes share an operon with genes of known function, the known function can
137 reasonably be hypothesized to relate to the unknown functions. In Table 1, unknown genes are
138 marked in underlined bold. 8 prior predictions include operons where no gene has known
139 function; in all these cases, the present analysis predicts that the prior prediction should be
140 merged with another prior containing at least 1 gene of known function. Predicted operon
141 membership *per se* may not be strong enough evidence to infer gene function, but it can provide
142 the basis for hypothesizing function, and the hypothesis can be strengthened by other evidence.
143
144 Each operon (training priors and merge candidates) was represented by a 4-vector consisting of
145 the minimum, mean, standard deviation, and maximum of the ABLIM distances among all gene
146 pairs in the operon. None of these statistics alone was sufficient for training an accurate
147 classifier. The LMT classifier had the best accuracy among the 44 classifiers that were evaluated
148 (Supplemental Table 1). However, this does not imply that LMT should be used when analyzing
149 other organisms. Future work on other organisms should repeat the classifier evaluation reported
150 here, and should choose the best classifier for the organism at hand.
151
152 The false-positive probability (column "P(false +)" in Supplemental Table 2) is a rough estimate.
153 It has much in common with a p-value: the null hypothesis is that the prior operons should not be
154 merged; the alternative hypothesis is that they should be merged; the statistic is the cumulative
155 probability of the null hypothesis when a score is at least as strong as the score at hand.
156 However, the cumulative probability is based on a negative training set of non-operons which is
157 specific but not sensitive. No members of the negative set can possibly be operons, because both
158 DNA strands are present. However no same-strand non-operons are present in the negative

159  training set, because these are difficult to ascertain. Thus there is a bias in the negative set, and
160  the resulting P(false +) values should not be viewed as rigorous.
161
162
163

164  ## Conclusions

165  The work presented here demonstrates that Machine Learning analysis of diel expression studies
166  can improve *in silico* predictions of operons. When a prior prediction is extended to include
167  genes of unknown function, the function of the known genes in the prior might elucidate the
168  function of the new unknown genes.
169

170  The approach presented here can be applied to other cyanobacteria for which diel studies and
171  prior predicted operons are available. Since the method is based on similarity of diel signatures,
172  best results should be expected from organisms whose genes exhibit strong and diverse diel
173  variation. Organisms with weak diel variation can be expected to perform poorly, because the 4-
174  vectors that describe operons to the classifier would all be similar. Experiments with a diel
175  study[34] of the minimal bacterium *Prochlorococcus marinus* produced poor results with the
176  approach presented here, possibly because the circadian clock mechanism is simplified in
177  *Prochlorococcus*[35] and its diel genes fluctuate more weakly than those of *Crocosphaera.*
178
179
180
181
182
183
184

185  ## Acknowledgements

186  The author is grateful to Jonathan Zehr, Josh Stuart, Irina Shilova, Rex Malmstrom, and
187  Laurence Nedelec for valuable discussions.
188
189
190
191

192  ## References

193  1.      Schopf, J. In: Whitton B, Potts M (eds) The fossil record: tracing the roots of the
194  cyanobacterial lineage. in *The Ecology of Cyanobacteria* 13–35 (Kluwer, 2002).

195   2.      Tomitani, A., Knoll, A. H., Cavanaugh, C. M. & Ohno, T. The evolutionary
196   diversification of cyanobacteria: Molecular-phylogenetic and paleontological perspectives.
197   *Proceedings of the National Academy of Sciences* **103**, 5442–5447 (2006).
198   3.      Tréguer, P. *et al.* Influence of diatom diversity on the ocean biological carbon pump.
199   *Nature Geosci* **11**, 27–37 (2018).
200   4.      Wang, W.-L., Moore, J. K., Martiny, A. C. & Primeau, F. W. Convergent estimates of
201   marine nitrogen fixation. *Nature* **566**, 205–211 (2019).
202   5.      Tang, W., Li, Z. & Cassar, N. Machine Learning Estimates of Global Marine Nitrogen
203   Fixation. *Journal of Geophysical Research: Biogeosciences* **124**, 717–730 (2019).
204   6.      Tan, L. T. Bioactive natural products from marine cyanobacteria for drug discovery.
205   *Phytochemistry* **68**, 954–979 (2007).
206   7.      Soni, B., Trivedi, U. & Madamwar, D. A novel method of single step hydrophobic
207   interaction chromatography for the purification of phycocyanin from Phormidium fragile and its
208   characterization for antioxidant property. *Bioresource Technology* **99**, 188–194 (2008).
209   8.      Zanchett, G. & Oliveira-Filho, E. Cyanobacteria and Cyanotoxins: From Impacts on
210   Aquatic Ecosystems and Human Health to Anticarcinogenic Effects. *Toxins* **5**, 1896–1917
211   (2013).
212   9.      Sakurai, H., Masukawa, H., Kitashima, M. & Inoue, K. How Close We Are to Achieving
213   Commercially Viable Large-Scale Photobiological Hydrogen Production by Cyanobacteria: A
214   Review of the Biological Aspects. *Life* **5**, 997–1018 (2015).
215   10.     Farrokh, P., Sheikhpour, M., Kasaeian, A., Asadi, H. & Bavandi, R. Cyanobacteria as an
216   eco-friendly resource for biofuel production: A critical review. *Biotechnol. Prog.* **35**, (2019).
217   11.     Singh, H., Khattar, J. S. & Ahluwalia, A. S. Cyanobacteria and agricultural crops.
218   *Vegetos- An Inter. Jour. of Plnt. Rese.* **27**, 37 (2014).
219   12.     Chittapun, S., Limbipichai, S., Amnuaysin, N., Boonkerd, R. & Charoensook, M. Effects
220   of using cyanobacteria and fertilizer on growth and yield of rice, Pathum Thani I: a pot
221   experiment. *J Appl Phycol* **30**, 79–85 (2018).
222   13.     Morone, J., Alfeus, A., Vasconcelos, V. & Martins, R. Revealing the potential of
223   cyanobacteria in cosmetics and cosmeceuticals — A new bioactive approach. *Algal Research* **41**,
224   101541 (2019).
225   14.     Khan, Z., Bhadouria, P. & Bisen, P. Nutritional and Therapeutic Potential of Spirulina.
226   *CPB* **6**, 373–379 (2005).
227   15.     Zehr, J. P. Nitrogen fixation by marine cyanobacteria. *Trends in Microbiology* **19**, 162–
228   173 (2011).
229   16.     Zheng, Y., Szustakowski, J. D., Fortnow, L., Roberts, R. J. & Kasif, S. Computational
230   Identification of Operons in Microbial Genomes. *Genome Research* **12**, 1221–1230 (2002).
231   17.     Westover, B. P., Buhler, J. D., Sonnenburg, J. L. & Gordon, J. I. Operon prediction
232   without a training set. *Bioinformatics* **21**, 880–888 (2005).
233   18.     Moreno-Hagelsieb, G. & Collado-Vides, J. A powerful non-homology method for the
234   prediction of operons in prokaryotes. *Bioinformatics* **18**, S329–S336 (2002).

235    19.     Price, M. N. A novel method for accurate operon predictions in all sequenced
236    prokaryotes. *Nucleic Acids Research* **33**, 880–892 (2005).
237    20.     Yao, D. C. I., Brune, D. C., Vavilin, D. & Vermaas, W. F. J. Photosystem II Component
238    Lifetimes in the Cyanobacterium *Synechocystis* sp. Strain PCC 6803: SMALL Cab-LIKE
239    PROTEINS STABILIZE BIOSYNTHESIS INTERMEDIATES AND AFFECT EARLY STEPS
240    IN CHLOROPHYLL SYNTHESIS. *Journal of Biological Chemistry* **287**, 682–692 (2012).
241    21.     Renger, G. *et al.* ON THE MECHANISM OF PHOTOSYSTEM II DETERIORATION
242    BY UV-B IRRADIATION. *Photochemistry and Photobiology* **49**, 97–105 (1989).
243    22.     Dodd, A. N. Plant Circadian Clocks Increase Photosynthesis, Growth, Survival, and
244    Competitive Advantage. *Science* **309**, 630–633 (2005).
245    23.     Bergersen, F. J. The Effects of Partial Pressure of Oxygen upon Respiration and Nitrogen
246    Fixation by Soybean Root Nodules. *Journal of General Microbiology* **29**, 113–125 (1962).
247    24.     Fay, P. Oxygen Relations of Nitrogen Fixation in Cyanobacteria. *MICROBIOL. REV.* **56**,
248    34 (1992).
249    25.     Tuit, C., Waterbury, J. & Ravizza, G. Diel variation of molybdenum and iron in marine
250    diazotrophic cyanobacteria. *Limnology and Oceanography* **49**, 978–990 (2004).
251    26.     Shi, T., Ilikchyan, I., Rabouille, S. & Zehr, J. P. Genome-wide analysis of diel gene
252    expression in the unicellular N2-fixing cyanobacterium Crocosphaera watsonii WH 8501. *The
253    ISME Journal* **4**, 621–632 (2010).
254    27.     Sabatti, C. Co-expression pattern from DNA microarray experiments as a tool for operon
255    prediction. *Nucleic Acids Research* **30**, 2886–2893 (2002).
256    28.     Lercher, M. J. Coexpression of Neighboring Genes in Caenorhabditis Elegans Is Mostly
257    Due to Operons and Duplicate Genes. *Genome Research* **13**, 238–243 (2003).
258    29.     Landwehr, N., Hall, M. & Frank, E. Logistic Model Trees. *Mach Learn* **59**, 161–205
259    (2005).
260    30.     Sumner, M., Frank, E. & Hall, M. Speeding Up Logistic Model Tree Induction. in
261    *Knowledge Discovery in Databases: PKDD 2005* (eds. Jorge, A. M., Torgo, L., Brazdil, P.,
262    Camacho, R. & Gama, J.) vol. 3721 675–683 (Springer Berlin Heidelberg, 2005).
263    31.     Heller, P. & Baiju, Bharath. An Improved Distance Metric For Clustering Gene
264    Expression Time-Series Data. *Zenodo* (2018) doi:10.5281/zenodo.1419822.
265    32.     Hall, M. *et al.* The WEKA Data Mining Software: An Update. **11**, 9 (2009).
266    33.     Frank, E., Hall, M. & Witten, I. The WEKA Workbench. in *Data Mining: Practical
267    Machine Learning Tools and Techniques* (Morgan Kaufmann, 2016).
268    34.     Zinser, E. R. *et al.* Choreography of the Transcriptome, Photophysiology, and Cell Cycle
269    of a Minimal Photoautotroph, Prochlorococcus. *PLoS ONE* **4**, e5135 (2009).
270    35.     Holtzendorff, J. *et al.* Genome Streamlining Results in Loss of Robustness of the
271    Circadian Clock in the Marine Cyanobacterium *Prochlorococcus marinus* PCC 9511. *J Biol
272    Rhythms* **23**, 187–199 (2008).
273

**Table 1**(on next page)

The 48 classifiers evaluated in this study.

48 classifiers, evaluated by 5-fold cross validation on the *Crocosphaera* time-series data, in descending order of accuracy. Best accuracy was achieved by the Logistic Model Tree (LMT) classifier, which was selected for this study.

1

| Classifier | Accuracy |
|---|---|
| trees.LMT | 95.28 |
| functions.MultilayerPerceptron | 95.03 |
| functions.Logistic | 92.8 |
| meta.MultiClassClassifier | 92.8 |
| meta.RandomizableFilteredClassifier | 91.89 |
| functions.SimpleLogistic | 91.56 |
| lazy.IBk | 91.31 |
| trees.RandomForest | 91.14 |
| lazy.KStar | 90.89 |
| meta.Bagging | 89.82 |
| meta.RandomCommittee | 89.32 |
| trees.RandomTree | 88.91 |
| functions.SGD | 88.49 |
| meta.MultiClassClassifierUpdateable | 88.49 |
| trees.REPTree | 88.49 |
| trees.J48 | 88.25 |
| rules.JRip | 87.83 |
| meta.RandomSubSpace | 87.67 |
| functions.VotedPerceptron | 87.33 |
| rules.DecisionTable | 87.17 |
| meta.LogitBoost | 86.51 |
| meta.FilteredClassifier | 86.26 |
| meta.IterativeClassifierOptimizer | 86.09 |
| rules.PART | 85.6 |
| meta.AdaBoostM1 | 85.51 |
| functions.SMO | 85.1 |
| bayes.BayesNet | 84.44 |
| trees.HoeffdingTree | 83.53 |
| bayes.NaiveBayesMultinomial | 83.44 |
| bayes.NaiveBayesMultinomialUpdateable | 83.44 |
| bayes.NaiveBayes | 82.7 |
| bayes.NaiveBayesUpdateable | 82.7 |
| meta.AttributeSelectedClassifier | 82.12 |
| lazy.LWL | 81.71 |
| trees.DecisionStump | 81.71 |
| rules.OneR | 81.29 |
| bayes.NaiveBayesMultinomialText | 66.14 |
| functions.SGDText | 66.14 |
| meta.CVParameterSelection | 66.14 |
| meta.MultiScheme | 66.14 |
| meta.Stacking | 66.14 |
| meta.Vote | 66.14 |
| meta.WeightedInstancesHandlerWrapper | 66.14 |
| misc.InputMappedClassifier | 66.14 |
| rules.ZeroR | 66.14 |

2  Supplemental Table 1 – 48 Classifiers evaluated by 5-fold cross validation, in descending order of accuracy. The LMT (Logistic
3  Model Tree) classifier was chosen for this study.
4

**Figure 1**(on next page)

Method for accepting/rejecting proposed merger of 2 prior predicted operons (red and green).

Both priors must lie on the same DNA strand, and each must contain at least 1 gene with diel variation. All pairwise ABLIM distances among the 4 genes are computed (blue). A 4-vector comprising the minimum, mean, standard deviation, and maximum of the ABLIM distances is computed (purple) and submitted to the LMT classifier. The classifier produces a score s. Gaussian distributions over scores of positive (upper curve) and negative (lower curve) are used to compute, respectively, the false negative and false positive probabilities.