

pubs.acs.org/jpr Article

# Why Environmental Biomarkers Work: Transcriptome—Proteome Correlations and Modeling of Multistressor Experiments in the Marine Bacterium *Trichodesmium*

Nathan G. Walworth, Mak A. Saito, Michael D. Lee, Matthew R. McIlvin, Dawn M. Moran, Riss M. Kellogg, Fei-Xue Fu, David A. Hutchins, and Eric A. Webb\*



Cite This: J. Proteome Res. 2022, 21, 77-89



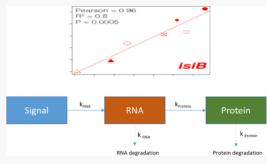
**ACCESS** 

III Metrics & More

Article Recommendations

Supporting Information

**ABSTRACT:** Ocean microbial communities are important contributors to the global biogeochemical reactions that sustain life on Earth. The factors controlling these communities are being increasingly explored using metatranscriptomic and metaproteomic environmental biomarkers. Using published proteomes and transcriptomes from the abundant colony-forming cyanobacterium *Trichodesmium* (strain IMS101) grown under varying Fe and/or P limitation in low and high CO<sub>2</sub>, we observed robust correlations of stress-induced proteins and RNAs (i.e., involved in transport and homeostasis) that yield useful information on the nutrient status under low and/or high CO<sub>2</sub>. Conversely, transcriptional and translational correlations of many other central metabolism pathways exhibit broad discordance. A cellular RNA and protein



production/degradation model demonstrates how biomolecules with small initial inventories, such as environmentally responsive proteins, achieve large increases in fold-change units as opposed to those with a higher basal expression and inventory such as metabolic systems. Microbial cells, due to their immersion in the environment, tend to show large adaptive responses in both RNA and protein that result in transcript—protein correlations. These observations and model results demonstrate multi-omic coherence for environmental biomarkers and provide the underlying mechanism for those observations, supporting the promise for global application in detecting responses to environmental stimuli in a changing ocean.

KEYWORDS: transcriptome-proteome, environmental biomarkers, marine microbes, Trichodesmium, metaproteomics

# INTRODUCTION

A recurring question in the interpretation of transcriptome and proteome datasets is the extent to which they co-vary. Messenger RNA (mRNA) and protein levels have been reported to generally be uncorrelated within a single cell, and only modestly correlated in populations of cells, due to differences in half-lives and degradation rates or phase variation within a population, respectively. It has also been shown that genes and their corresponding proteins associated with different cellular processes (e.g., core central vs stress metabolism) may retain varying degrees of correlation.<sup>2</sup> While observation of correlations in larger organisms may be challenging due to internal tissues being more remote environmental signals, microbes due to their small size and immersion within the environment often maintain multiple adaptive response capabilities to common environmental stimuli. Microbial datasets with transcriptomic and proteomic methods applied to the same experiment(s) are becoming more common and hence could aid in detection and interpretation of key environmental processes. In recent years, concurrent measurements of transcripts and proteins have been conducted in marine microbes such as Pelagibacter, Prochlorococcus, the marine diatom Thalassiosira pseudonana, the brown alga Aureococcus anophagefferens,  $^6$  the polar alga Phaeocystis antarctica,  $^7$  and the diazotrophic cyanobacterium Trichodesmium.  $^8$  These studies also observed varying extents of transcriptome—proteome coupling, with correlations observed particularly for genes involved in responding to environmental stresses, such as P, Fe, or vitamin  $B_{12}$  limitation. Despite these efforts, there is arguably a lack of consensus in the marine ecology community regarding the extent that RNA transcripts and proteins should correlate, with many having the opinion that correlations do not occur.

In this study, we synthesize the results from a number of *Trichodesmium* transcriptome—proteome datasets across a spectrum of conditions to understand their transcriptional and translational responses on a mechanistic level. *Trichodesmium* spp. are filamentous, buoyant micro-organisms that can

Received: June 23, 2021 Published: December 2, 2021





commonly grow in macroscopic colonies in close association with other microbes and have the capability to form massive blooms. Given their ability to fix both carbon and nitrogen, Trichodesmium spp. have relevance to both global productivity and biogeochemistry. 10,111 Together with other marine microbes, they can impact both ecosystem stability and climate feedbacks. 12 Trichodesmium is among one of the several oceanic cyanobacteria that are globally significant sources of N<sup>10,13</sup> as well as unicellular forms (Crocosphaera spp., Candidatus atelocyanobacterium thalassa, and Cyanothece) and other heterocystous, symbiotic forms<sup>11</sup> that can have relatively high cell numbers 14 and biogeochemical importance. 15,16 Trichodesmium is also considered one of the most important diazotrophs in many tropical and subtropical regimes. 10,17 As a result, research emphasis has been placed on developing field-ready, nutrient-limitation/stress markers to define the factors that control growth and N2 fixation in this important genus. These combined efforts have shown that iron (Fe) and phosphorus (P) primarily limit *Trichodesmium* across much of the global oceans. However, comparatively fewer studies have been conducted examining the interactions among future high CO<sub>2</sub> concentrations, Fe, and P in the context of holistic, molecular physiology.

A motivation in understanding the dynamics of RNA and protein responses is the interpretation of natural microbial populations and, in particular, the expression of genes that can provide clues to ecological or biogeochemical processes at play. Advances in sequencing (e.g., metagenomics/metatranscriptomics) and mass spectrometry (e.g., proteomics) technologies have enabled worldwide microbial characterizations from different biogeochemical regimes (e.g., Global Ocean Sampling Dataset<sup>20</sup>). While these data have revolutionized how we think about microbially mediated processes, extrapolating biogeochemistry from relative abundances in metagenome data may result in biased interpretation. For example, it is now recognized that globally abundant microbial species initially identified via metagenomics may not necessarily be dominant members of the transcriptionally active microbial community that is primarily responsible for biogeochemical turnover.<sup>21</sup> While metagenomic studies serve as powerful hypothesis-generating datasets and have unmasked previously underappreciated ecosystem processes (e.g., proteorhodopsin diversity<sup>20,22</sup>), they have also underscored the need to understand the dynamics of intracellular molecular processes in the context of organismal physiology both in the laboratory and field. Thus, it is important to investigate molecular dynamics of biogeochemically important microbes under differing nutrient and temporal regimes to resolve mRNA-protein dynamics that ultimately drive important biogeochemical processes.

Here, we focus primarily on  $N_2$  fixation. In past studies, we characterized the physiological and evolutionary responses of the globally important, photoautotrophic,  $N_2$ -fixing cyanobacterium, *Trichodesmium erythraeum* IMS101 (hereafter IMS101) to high  $CO_2$ . Using these cell lines as starting genetic backgrounds (i.e., replicates adapted to both 380 and 750  $\mu$ atm  $CO_2$ ), we then conditioned them to iron-limited, phosphorus-limited, and iron/phosphorus (Fe/P) co-limited conditions to define their predicted climate change-impacted Fe- and P-stress molecular responses. Consistent with another study, <sup>24</sup> this demonstration of the fitness advantage conferred by Fe/P "balanced limitation" compared to single limitation alone appears to contradict the long-standing Liebig limitation model and has implications for global biogeochemical cycles in both the current

and future ocean such as increased exogenous nitrogen scavenging coupled to a decrease in  $N_2$  fixation. <sup>23,25,26</sup>

In this study, we build on and synthesize our prior results to examine steady-state transcriptional/translational/physiological relationships in the context of environmentally relevant (e.g., Fe and/or P affected) nutrient regimes using our high and low CO<sub>2</sub>adapted, IMS101 cells lines. 23,25,27 As mentioned, numerous studies utilizing diverse techniques ranging from chemical quotas, enzyme activities, to gene- or protein-based molecular stress markers have pointed to the importance of Fe and P to Trichodesmium N<sub>2</sub> fixation in situ. 18,28,29 Despite the value of molecular markers for indicating species-specific bioavailability of nutrients, concerns persist that gene/protein expression can be decoupled with biogeochemically relevant rates like N<sub>2</sub> fixation (e.g., in ref 30) and thus give an inconsistent view of in situ limitation. We performed meta-analyses of our published, long-term Fe- and P-single limited, and co-limited proteomic and transcriptomic datasets<sup>23,25</sup> to determine which genes and proteins give congruent expression patterns during mid-day across eight experimental conditions vs those that do not. We chose mid-day for sampling as it is the period for peak photosynthesis, carbon fixation, and N2 fixation for Trichodesmium. This analysis shows that major components of core metabolic pathway genes and proteins show incongruent responses to environmental perturbations, whereas many wellclassified nutrient-responsive genes/proteins show consistent correlations under both low and high CO2 regimes. A model of the cellular inventories for RNA and protein biomolecules that compares the timing of biosynthesis and degradation was created that provides a mechanistic explanation for the observed coherence of responses to environmental perturbations. Finally, the environmental biomarkers that have been characterized thus far in marine microbes are briefly reviewed.

#### MATERIALS AND METHODS

# Culturing, Molecular Extractions, Sequencing, and Proteomic Analyses

A depiction of the prior experimental design can be found in Figure S1a and in the work of Walworth et al.<sup>23</sup> in addition to physiological and proteomic results. Detailed culturing,<sup>23</sup> protein extraction,<sup>27</sup> RNA extraction,<sup>31</sup> and analysis methods can be found in our previously published papers. Briefly, proteomic data was generated using data-dependent acquisition on a Thermo Fusion mass spectrometer and mapped to the IMS101 genome using SEQUEST with Proteome Discoverer. Here, 1908 Trichodesmium proteins were detected across the sample set using a 0.3% false discovery rate and requiring two peptides per protein using Scaffold 3.0 (Proteome Software) and normalized spectral counts were reported (see Supplemental Dataset 1) and matched to transcript abundances. Spectral counts provided a robust relative abundance of the proteome, although future efforts could focus on peptide precursor or fragment ion peak areas, which have a broader dynamic range, although the former is subject to interference from sequence variants in metaproteomic samples. Transcript data were generated using Illumina Hi-Seq after rRNA removal, yielding single-end 50 bp read libraries (see Supplemental Dataset 1).

# Differential Expression Analysis

Differential expression data were from previously described datasets.<sup>25</sup>

Table 1. Parameters in Transcriptome and Proteome Expression Model

constant	description	value	reference <sup>a</sup>
$k_{ m RNA}$	production of RNA	100	
$k_{-\mathrm{RNA}}$	degradation of RNA	$0.102\mathrm{min^{-1}}$ (6.8 min half-life)	1 (E. coli)
$d_{ m RNA}$	delay between the signal and RNA	1 min	
$k_{ m protein}$	production of protein	3.5	3
$k_{ m -protein}$	degradation of protein	$0.0017 \text{ min}^{-1} (6.9 \text{ h half-life})$	2 (yeast)
$d_{ m protein}$	delay between RNA and protein	180 min	
$s_0$	basal signal level 1%	0.01	
$s_1$	increased signal due to environmental cue	0.2	
$s_{0alt}$	basal signal 20%	0.2	
$s_{1\text{alt}}$	increased signal due to environmental cue (+0.2)	0.4	
$t_{ m start}, t_{ m plateau1}, t_{ m plateau2}, \ t_{ m end}$	time of onset of signal $(t_{\text{start}})$ , begin plateau $(t_{\text{plateau1}})$ , end plateau $(t_{\text{plateau2}})$ , end of decrease in signal $t_{\text{end;}}$ ("long slow ramp")	600, 650, 700, 1200 min	
$t_0$	time point denominator for fold change	9 h (540 min)	
t <sub>0</sub> by day		540, 1980, 3420 min	

<sup>a1</sup>Selinger et al., 2003. <sup>2</sup>Pratt et al. 6.9–34 h protein half-lives reported. <sup>3</sup>Tuned using a 1 min signal to be 6–7 proteins per transcript based on Moran et al.

#### **Multivariate and Pairwise Analyses**

Redundancy analysis was conducted on TMM-normalized transcript levels and nutrient concentrations using the vegan package.<sup>32</sup> Redundancy analysis (RDA) ordinations were compared via the Procrustes test with the "protest" function in vegan using default settings.

Nonmetric multidimensional scaling (NMDS) was conducted using the "metaMDS" function from the vegan package with default settings (Bray—Curtis dissimilarities computed using Wisconsin double standardization) except for "autotransform = FALSE".

Hierarchical clustering with multiscale bootstrap resampling was conducted on Bray—Curtis dissimilarities calculated from TMM-normalized transcript levels using the pvclust package<sup>33</sup> with the following settings: method.dist = "euclidean", method.hclust = "average", nboot = 1000.

The "cor.test" function in "R" was used to calculate Pearson correlation coefficients on log2 fold changes (relative to the r380 condition; see below) of average TMM-normalized transcript levels vs log2 fold changes of average normalized spectral counts of the corresponding protein product detected in the work of Walworth et al. 23 Significance tests between mean correlation coefficients were carried out via the "perm" package 34 using a two-sample permutation test with Monte Carlo simulation and 2000 permutations.

# Gene Ontology (GO) Enrichment Analysis

As previously described in the work of Walworth et al.,<sup>31</sup> Gene Ontology (GO) mappings for *Trichodesmium* were downloaded from the Genome2D web server (http://genome2d.molgenrug.nl). The "phyper" function in "R" (R Core Team 2014) was used to test for significant enrichment of GO pathways, and *p*-values were corrected with the Benjamin and Hochberg method<sup>35</sup> using the "p.adjust" function ( $p \le 0.05$ ). Finally, genes in enriched GO categories were manually checked.

## **RNA-Protein Model**

A model simulating the biosynthesis of RNA transcripts and proteins was created in MATLAB (MathWorks) using parameters listed in Table 1. One-minute time steps over a duration of 3 days (4320 min) within a population of 100 cells were modeled. Expression was modeled with a signal parameter that varied between 0 and 1, representing the gene being turned

on or off with a simple linear ramping function. Basal signal level  $(s_0)$  was set to 0.01, representing 1% of the population or 1 cell within 100 cells with gene expression to avoid a zero denominator in fold-change estimates. In an alternate scenario representing metabolic functions in continuous use, a basal expression of 20% (0.2) was used. Rate constants for production of RNA and protein  $(k_{\rm RNA}$  and  $k_{\rm protein})$  as well as their decay  $(k_{\rm -RNA}$  and  $k_{\rm -protein})$  were obtained from experimental measurements in *Escherichia coli* and *Saccharomyces cerevisiae* where the degradation of each type of biomolecule was inhibited and measured. <sup>36,37</sup> Change in the inventories of RNA and protein biomolecules was estimated using simple production and loss terms (eqs 1 and 2):

$$\frac{\mathrm{dRNA}}{\mathrm{d}t} = (k_{\mathrm{RNA}} \times s) - (k_{-\mathrm{RNA}} \times [\mathrm{RNA}]) \tag{1}$$

$$\frac{\text{dprotein}}{\text{d}t} = (k_{\text{protein}} \times [\text{RNA}]) - (k_{\text{-protein}} \times [\text{protein}])$$
(2)

Turnover time estimates of these biomolecules are, to our knowledge, unavailable for marine organisms. While they may be different, presumably slower based on observed growth rates, it seems a reasonable assumption that the large difference between RNA and protein decay would be maintained. Hence, the trends observed here are likely to be consistent or further enhanced, particularly if the RNA degradation activity in marine microbes is slower than that of *E. coli* as might be expected.

This model has similarities to that of Moran et al.'s work in 2013,  $^{38}$  but it is applied here to focus on any coordination between RNA and protein and how results affect fold-change estimates. Moran et al. instead explored random triggers for transcript production using 10 random events per day of one transcript per event, and then simulated environmental stimuli by increasing to three transcripts per event. In contrast, the model in the present study estimates production of RNA and protein molecules based on their upstream signal (signal and transcripts) as described in eqs 1 and 2 rather than assigning a constant number of seven protein copies per transcript. Moran et al. observed that the total number of mRNA molecules in a typical marine bacterium was only  $\sim\!\!200$  per cell.  $^{38}$  These observations were incorporated into the model by allowing  $k_{\rm RNA}$ 

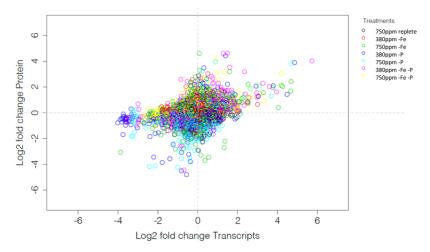


Figure 1. Overall comparison of proteome—transcriptome data in log2 fold space for eight experimental treatments, where treatments varied  $CO_{\mathcal{D}}$  iron, and phosphorus (see Figure S1 for full experimental design) in cultures of the marine cyanobacterium T. erythraeum IMS101. Legend shows symbols of seven treatments relative to the 380 ppm  $CO_2$  replete treatment. Transcript data is the average of biological duplicates, and protein data is the average of biological triplicates.

to be 100 to enable signal input to be transcribed into 100 transcripts across a population of 100 cells and then tuning the protein production  $k_{\text{protein}}$  parameter to produce 6–7 proteins per transcript also based on Moran et al.'s work. Similarly, Moran et al. used an RNA half-life of 1.5 min, whereas this model uses a more conservative RNA half-life of 6.8 min based on a tiling microarray study in *E. coli.*<sup>36</sup>

## RESULTS AND DISCUSSION

In this section, we describe statistical analyses of the marine cyanobacterium *Trichodesmium* experiments including (1) differential expression and cluster analysis of the transcriptome dataset across eight environmental treatments, (2) multivariate analysis of the whole RNA and protein datasets to detect influence of environmental parameters, and (3) pairwise transcript—protein correlation analyses. Subsequently, a simple RNA—protein model is used to reproduce the dynamics observed in transcriptome—proteome comparison experiments.

# Global Comparison of CO<sub>2</sub>-Impacted, Nutrient-Limited Molecular Dynamics

To investigate shifts in molecular metabolism underlying Felimited, P-limited, and Fe/P co-limited metabolisms as they interact with predicted future CO<sub>2</sub> atmospheric conditions, we examined the global transcriptional output in each scenario from Walworth et al.<sup>25</sup> and their corresponding translational correlations with proteins identified in Walworth et al. 23 Across all treatments, we detected a transcription of ≥96% and translation of 37% (see below) of the IMG-annotated genes. Instances where comparisons between transcripts and proteins were not possible due to one of those biomolecules not being detected were excluded. While not the goal of our study, future studies could apply deeper proteomic methods to explore correlative relationships in rarer biomolecules. Initial examination of the dataset in the log2 fold change space revealed a cluster of points around the origin with a subset of data extending into the northeast and southwest quadrants, implying a co-varying relationship between some transcript—protein pairs (Figure 1). It is important to note that nontranslated mRNAs can also potentially represent "irrelevant expression" if these transcripts are not regulatory. Thus, there is much to be learned from future studies reconciling the constantly developing techniques of RNA- and protein-based 'omics, with relatively few studies to date comparing multicondition proteomes with concomitantly taken transcriptomes.

To this end, we first conducted exploratory analysis on the global transcriptional output without constraining for environmental variables (see below) in order to observe how replicates clustered without imparting treatment information (i.e., no assumptions). We used both NMDS and hierarchical clustering with multiscale bootstrap resampling (replicates = 1000). NMDS revealed consistent nutrient-limited patterns across replicates (Figure S1b), which was further supported by hierarchical clustering (Figure S1c) that generated two highconfidence clusters with AU-bootstrap p-values >0.95 (i.e., rejection of the null hypothesis that the clusters do not exist at the 0.05 significance level). We labeled samples based on their environmental condition at the time of sampling with CO<sub>2</sub> concentration followed by the limiting nutrient. For example, 380-Fe signifies low CO<sub>2</sub> and low iron. For replete conditions, we chose to put an "r" in front of the CO2 concentration (i.e., r380 and r750). In these unconstrained analyses, the 380-Fe condition demonstrated greater inter-replicate transcriptional variation, although both replicates were indeed Fe-limited as evidenced by both reduced growth<sup>23</sup> and several upregulated Festress genes (Figure S1b; see below and Walworth et al.8). Interestingly, Fe-limited treatments (380- and 750-Fe) clustered with the nutrient replete ones (r380 and r750), supporting the notion that P limitation evokes a more varied molecular response than Fe-limitation relative to replete nutrient regimes (Figure S1c). This trend was further corroborated through hierarchical clustering of transcript levels of all differentially expressed (DE) genes (i.e., those DE in any treatment relative to r380; n = 1943; Figure 2a). To better resolve sources of variation producing these groupings, we further clustered upregulated (*n* = 951; Figure 2b) and downregulated (n = 1074; Figure 1c) transcripts, which yielded strikingly different trends among Felimited and replete conditions. Upregulated genes exhibited significantly different (bootstrap confidence level > 0.95) clusters among Fe-limited and CO2 regimes relative to those produced by clustering all DE genes, while those of downregulated genes retained similar groupings to those of DE genes. These patterns suggest that Fe-limited and replete clusters

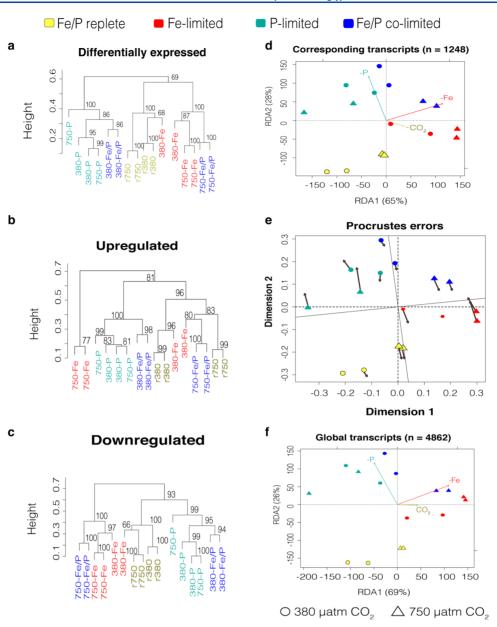


Figure 2. Global molecular analysis. Hierarchical clustering of Bray—Curtis dissimilarities with multiscale bootstrap resampling calculated from normalized transcripts of differentially expressed ((a) top left panel), upregulated ((b) middle left), and downregulated ((c) bottom left) genes. Numbers at dendrogram nodes are approximately unbiased p-values calculated from multiscale bootstrap resampling. The top right panel (d) shows an RDA of normalized transcript abundances of genes whose protein products were also detected in the work of Walworth et al. (2016a). The middle right panel (e) shows a Procrustes analysis between corresponding-RNA-RDA and prot-RDA ordinations, indicating close agreement evidenced by small vector residuals (r = 0.75; p < 0.001). The bottom right panel (f) displays the same redundancy analysis as in the top panel except now including all detected gene transcripts.

produced from all DE genes (Figure 2a) may have been driven more so by downregulated transcriptional variation of the DE gene pool. Hence, mechanistic adaptations to high CO<sub>2</sub> and responses to Fe limitation both seem to be more similar to each other than to molecular responses under P limitation. Interestingly, both up- and downregulated P-limited clusters were statistically analogous to those produced by clustering all DE genes (Figure 2a) and total genes (Figure S1c), which further corroborates the notion that P limitation evokes a highly conserved yet more varied transcriptional rearrangement than either high CO<sub>2</sub> or Fe limitation.

Next, we employed constrained multivariate methods that attempts to explain variation based on the given environmental

variables (e.g., RDA) for both transcripts and proteins. When relating nutrient-limited proteome variation identified via RDA in the work of Walworth et al.  $^{23}$  (n=1248; 37% of genome coding potential) to that of their corresponding gene transcripts (corresponding-RNA-RDA; Figure 2d), similar nutrient-limited patterns across replicates were observed with most of the proportional variance explained by the environmental variables for both transcripts (adjusted  $R^2=0.56$ ) and proteins (adjusted  $R^2=0.59$ ), respectively. Moreover, a statistically significant correlation was observed (r=0.75; P<0.001) when comparing the corresponding-RNA-RDA and prot-RDA ordinations via permutational, least-squares orthogonal analysis (i.e., Procrustes analysis; Materials and Methods) indicating that overall relative

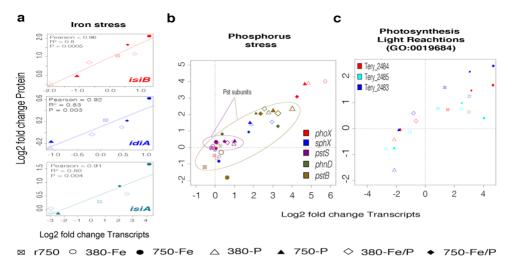


Figure 3. Scatterplots of log2 fold changes of transcript and protein abundances for select, highly correlated genes. (a) Shown are log2 fold changes of normalized protein (*y*-axis) and transcript (*x*-axis) abundances for Fe-stress genes. (b) Shown are log2 fold changes of normalized protein (*y*-axis) and transcript (*x*-axis) abundances for P-stress genes. Ellipses denote different subunits of the Pst transport system. (c) Shown are log2 fold changes of normalized protein (*y*-axis) and transcript (*x*-axis) abundances for photosynthesis light reaction genes. Different symbols denote different treatments.

transcriptional variation among treatments was captured in the corresponding proteome variation (Figure 2e).

Performing the same RDA on the global transcriptional output (global-RNA-RDA; Figure 2f; adjusted  $R^2 = 0.50$ ; n =4862) displayed a strikingly consistent ordination relative to that of both the prot-RDA (represented in Figure 1c in the work of Walworth et al.<sup>23</sup>) and corresponding-RNA-RDA (Figure 2d). In fact, when analyzing congruence of the global-RNA-RDA to both the prot-RDA<sup>23</sup> and the corresponding-RNA-RDA, respectively, a robust, statistically significant correlation was observed relative to both the prot-RDA (r = 0.75, p < 0.001) and the correlated-RNA-RDA (r = 0.98; p < 0.001). Taken together, these analyses indicate that both transcriptional and translational variations inherent in the overlapping 37% of the detected proteome and 96% of the detected transcriptome are representative not only of one another across nutrient-limited conditions but also of global transcriptional variation. In other words, this congruence implicates that the detected proteome variation is indeed representative of overall treatment-specific, transcriptional variation in IMS101 for Fe, P, Fe/P, and CO<sub>2</sub>.

Another striking observation is that 380-Fe/P replicates were more closely associated to P-limited treatments (380- and 750-P), while 750-Fe/P replicates were more correlated with Felimited ones (380- and 750-Fe) for both ordination-based (RDA, Figure 2d,f; NMDS, Figure S1b) and clustering analyses (Figure 2a–c and Figure S1c). Hence, these trends imply that regardless of the amount of variation explained by the given environmental variables, the global transcriptional profile of Fe/P co-limited metabolism is more similar to that of P-limitation when adapted to current  $\mathrm{CO}_2$  levels (380-Fe/P) but is more similar to Fe limitation when adapted to high  $\mathrm{CO}_2$  (750-Fe/P).

# From Global to Pairwise mRNA-Protein Correlations

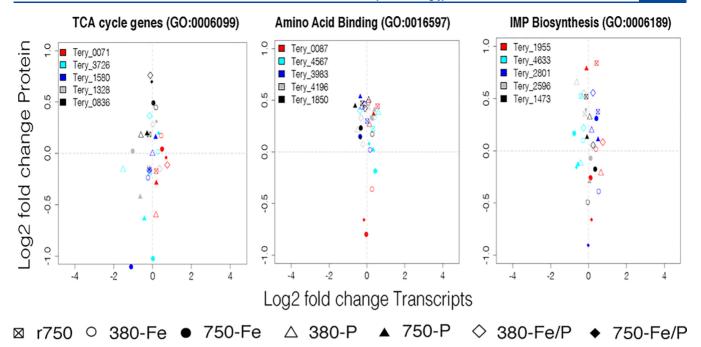
Bacterial messenger RNAs (mRNAs) typically have much shorter half-lives than proteins and can be rapidly transcribed in response to quickly changing environmental conditions. For example, in an *E. coli* cell, many mRNAs are typically degraded within minutes, whereas their corresponding proteins can have longer lifetimes than the cell cycle. These data indicate that intracellular mRNA copy number at any instant may typically reflect recent transcriptional activity, while protein levels in the

same instant can represent a relatively longer history of accumulated expression. Indeed, simultaneous measurements of mRNA and their corresponding protein concentrations have been shown to be uncorrelated within a single cell at a single time point, while mRNA levels integrated over many cells in a population generally correlate with protein concentrations. Different environmental variables invoking specific metabolic pathways can also impact overall mRNA—protein correlations. <sup>1,2</sup>

Unlike prior other microbial studies investigating mRNAprotein correlations across one or two treatments, we examined correlations across eight different growth conditions. As in most organisms to date, IMS101 transcript levels modestly correlated to their corresponding protein abundances across all treatments  $(n = 1246; squared Pearson correlation coefficient, <math>r^2 = 0.36)$ suggest that other forms of regulation might need to be invoked to explain the majority of variation. Upon calculating Pearson correlation coefficients across all treatments for genes that were DE in at least one experimental condition relative to r380 versus those that were not called as differentially expressed (NDE) under any condition, we observed a significantly greater mean correlation coefficient for DE  $(r^2 = 0.35)$  vs NDE (r = 0.03)genes (p < 0.0005; two-sample permutation test with Monte Carlo simulation). This same robust trend was also observed when comparing average correlation coefficients for DE vs NDE genes for each pairwise treatment comparison (Figure S2). Hence, this greater correlation suggests that transcription and translation of environmentally responsive genes and their protein products may be more closely coupled on average than those that do not respond to changing environmental conditions.

# **Pathway-Specific Molecular Relationships**

To search for metabolic pathways harboring highly correlated transcript and protein abundances, we calculated Pearson's correlation coefficients (PC) and kept genes retaining PC values  $\geq$ 0.7 and *p*-value <0.05. We then mapped genes onto their GO pathways and searched for significant GO enrichment using the hypergeometric test with Benjamini–Hochberg correction (FDR  $\leq$  0.05; Dataset S1). Enriched GO pathways included translation, nitrogen fixation, oxidation–reduction, calcium ion



**Figure 4.** Scatterplots of log2 fold changes of transcript and protein abundances for select core metabolism genes. Shown are log2 fold changes of normalized protein (y-axis) and transcript (x-axis) abundances for TCA cycle, amino acid binding, and IMP biosynthesis genes. Different symbols denote different treatments.

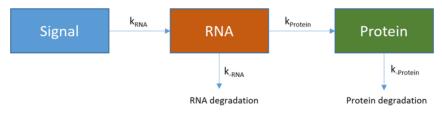


Figure 5. Representation of the transcriptome-proteome model. Parameters are described in Table 1.

binding, proteolysis, iron-sulfur cluster binding, and photosynthetic electron transport/light reactions. Hence, these data suggest that biosynthesis of these protein products increases with increasing transcript accumulation. Therefore, their abundances may be more representative of instantaneous cellular activity rather than those with low PC values, in which the latter may be a product of constitutive activity and/or reduced protein degradation. For example, well-characterized Fe- and P-stress genes exhibited high PC values in conjunction with genes involved in photosynthetic light reactions (Figure 3). In all cases, the interaction of high CO2 with limiting Fe irrespective of fluctuating P (750-Fe and 750-Fe/P) yields increases in both transcript and protein levels for all Fe-stress genes relative to their low CO<sub>2</sub>/low Fe conditions (i.e., 380-Fe and 380-Fe/P). Hence, the expression of these Fe-stress biomarkers increases under limiting Fe in future ocean conditions, suggesting their profiles to be robust for in situ profiling. The genes involved in photosynthesis light reactions (Figure 3c) exhibit parallel trends and could be related to midday sampling and increased transcription and translation of photosynthetic genes during the light period (Held et al., in press). Interestingly, P-stress transcript and protein levels robustly respond to limiting P irrespective of Fe and CO2, but increased CO2 has little to no effect on P-stress molecular machinery (Figure 3b). Additionally, these data suggest that different subunits of the same phosphate transport system (i.e., Pst) could respond differently to fluctuating P as transcript and

protein levels of, for example, the PstB subunit increased more considerably than those of PstS under limiting P (Figure 3b).

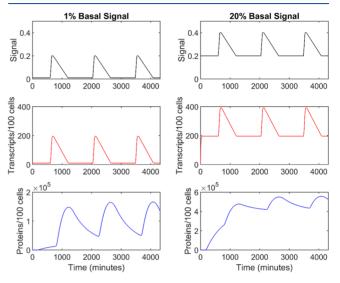
Conversely, GO pathways specific to the NDE pool with low PC values involve general biosynthetic processes, magnesium ion binding, proteolysis, respiration (e.g., TCA cycle), nucleotide metabolism, inosine monophosphate (IMP) biosynthesis, and amino acid metabolism (Figure 4). Of note, decoupling of TCA gene expression and enzyme abundance in yeast has also been previously observed. Interestingly, the GO pathway assigned to the photosystem I reaction center (GO: 0003989) experiences significant enrichment of genes with low PC values (Dataset S1) while conversely, the overall photosynthesis light reaction GO pathway (GO: 0019684) was enriched with genes with high PC values (Figure 3c). Hence, it is important to distinguish between potentially constitutive vs responsive components of large, multifaceted pathways with many components that can have widely differing roles (e.g., energy flow vs structure). Overall, it seems that the persistent activity of certain core pathways to maintain basic cellular homeostasis, irrespective of changing environmental conditions, may result in decoupling of transcription and translation on average (Figure 4), relative to pathways that synthesize new transcripts to generate protein products specific to particular environmental triggers (e.g., Figure 3).

# RNA-Protein Modeling: Influence of Biomolecule Inventories on Fold-Change Dynamics

In order to explore the potential causality of RNA-protein coupling and decoupling, we generated a model for production and decay of both biomolecules (Figure 5). Model parameters include previously published values for RNA and protein degradation based on studies of whole transcriptome and proteome degradation measurements from model organisms (E. coli<sup>36</sup> and yeast,<sup>37</sup> respectively) as well as delays between RNA and protein expression observed in *Prochlorococcus*<sup>4</sup> (Table 1). Notably, RNA decay rates are approximately two orders of magnitude higher (~60-275 fold) than protein degradation. While differences in degradation rates occur between different specific transcripts and protein molecules (for example, misfolded apo-metalloproteins being more susceptible to degradation), there is clearly a large difference between transcript and protein degradation rates, as expected based on the transient information transmission role of RNA versus the long-term functional roles of proteins. The choice of *E. coli* and yeast as organisms for obtaining decay parameters will likely vary from those of various marine microbes, with both degradation processes likely being slower within slower growing marine microbes such as Trichodesmium. Moreover, misfolded and mismetallated proteins are also known to be targeted for degradation by intracellular recycling processes. However, the large difference in cellular decay rates of RNA and protein rates is almost certainly a universal phenomenon, and hence we expect the general trends of temporal overlap between transcripts and protein inventories modeled here to be relevant even if future taxon-specific activities vary somewhat from those used here.

In particular, if the RNA degradation constant is lower in slower-growing organisms due to lower abundance and/or efficiency of RNases, then an increased temporal overlap between transcripts and proteins would result in greater opportunity for correlation between the two types of biomolecules. Additional parameters in the model include production of RNA and protein as well as constants for temporal offsets between signal transduction and transcription ( $d_{RNA}$ ), and transcription and translation processes  $(d_{protein})$ . For the latter, 3 h was used, similar to the range observed in Prochlorococcus. Finally, the model includes initial condition parameters whose influence can be explored including basal expression levels (1 and 20%) and time settings for onset, duration, and decay of signals to initiate transcription. The notion of the transcriptional signal is generalized here, whether they occur from riboswitches, two-component regulatory systems, sigma factors, or other mechanisms. Moreover, the onset and decay of the regulatory signaling are also idealized here as a linear ramp (using a "long slow ramp" that repeats daily), although it could be further qualified to have positive feedbacks to increase sensory capabilities such as those observed in increased expression of two-component regulatory systems, e.g., phosphate regulatory systems. 43 The influences of posttranslational modifications (PTMs) were purposely not included in this model in order to allow an examination of potential range of differences in RNA and protein fold-change expression levels prior to often invoking additional PTM regulatory controls. While PTMs likely contribute to protein regulation and enzyme activity in Trichodesmium and other marine microbes, those effects would likely not supersede the protein inventory dynamics measured and modeled here.

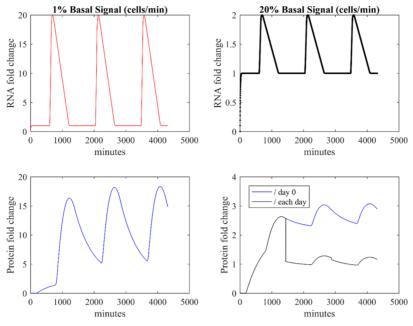
The RNA–protein expression model results imply that changes in basal expression level can have a major influence on the extent of RNA–proteome coupling and decoupling observed experimentally. For example, when there is no significant basal expression (1% or  $s_0 = 0.01$ ; where a nonzero value was chosen to allow fold-change ratio estimate) as would be expected prior to an environmental stimulus such as phosphate or iron stress, both RNA and protein signals respond with a large increase in inventory and fold-change (Figures 6 and



**Figure 6.** RNA—protein expression model results. Low basal expression of 1% (left panels) and moderate basal expression (right panels), with a simulated daily expression pattern. Short and long half-lives of RNA and protein, respectively, result in different inventories over time. When basal expression is increased to 20%, protein inventories begin to accumulate between successive daily cycles.

7, left panels). In contrast, when there was a 20% basal level of expression ( $s_0 = 0.2$ ), as expected for routinely used metabolic machinery, a significant inventory of protein product accumulates (Figure 6, right panels), causing the fold change signal to become muted (Figure 6). These results are also displayed as a movie format (Supplemental Movie S1) to display the temporal unfolding of these transcript—protein inventories and their influence on fold change results. The fold change decrease is accentuated if calculated each day as the protein inventory accumulates (e.g., the choice of what to normalize the expression to in the fold change estimate). Notably, proteins that respond to environmental stress are consistent with the lower basal expression model that yields correlations between RNA and protein.

In *Trichodesmium*, photosynthetic transcripts and proteins were also observed to be correlated (Figure 4c). While generally considered metabolic proteins with high inventories, this observed correlation can be particularly driven by iron limitation scenarios, which can remodel overall machinery by reducing the iron-rich photosystem I (PSI) in favor of PSII.<sup>23</sup> Iron stress transcripts and proteins like *isiA* can also concurrently increase in abundance solely in times of enhanced iron stress to form PSIIsiA-chlorophyll—protein—antenna super-complexes,<sup>44</sup> while under other conditions like phosphorus stress, they show little correlation due to the lack of an environmental trigger. Additionally, diazotrophs have been shown to have strong diel cycles that affect both transcripts and proteins of photosynthetic and nitrogen fixation enzymes in order to promote intracellular



**Figure 7.** RNA and protein expression in fold change space. Model output from Figure 6 converted to fold-change space results in high values with a low basal expression (1%; left panels) and muted values with some basal expression (20%; right panels). Fold change was calculated relative to day 0 or each successive day in the right panels.

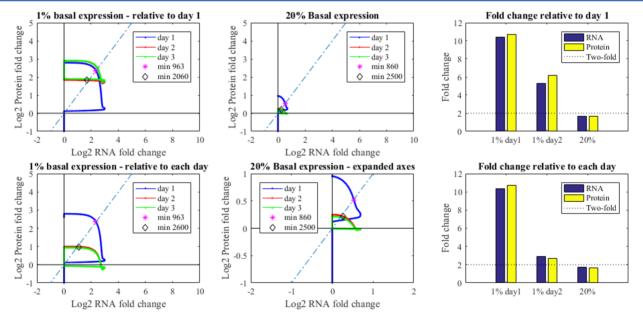


Figure 8. Expression patterns in transcriptome—proteome space. The succession of RNA and protein biomolecules resulted in a counterclockwise path, with more pronounced fold change signals with low basal expression (left panels) compared to higher basal expression (center panels) where the signal circular progression falls close to the origin (top middle). Coherence between biomolecules was observed as in the northeast or southwest quadrants (dotted line). The resulting maximum fold change (right panels) observed was high with low basal expression (top right) and when relative to day 1, rather than successive days as the protein inventory accumulates (Figure 6). Environmental biomarkers tend to have lower basal expression and are expressed based on regulation by environmental cues.

oxygen protection and iron conservation efforts. <sup>45,46</sup> Hence, the choice of when to sample transcriptomic and proteomic experiments can also be consequential. The short-reported half-lives of RNA molecules cause a rapid decrease in transcripts if they are not continually being produced. <sup>38</sup> The expected large two orders of magnitude difference in half-lives of RNA and protein molecules results in a counterclockwise progression on RNA and protein space plots commonly used for assessing coherence between these biomolecules (Figure 8). Choice of

sampling time is typically estimated to allow sufficient time for RNA and protein biosynthesis, where if too early then there is insufficient time for protein synthesis, and if too late then RNA transcripts may have begun degradation. Notably, the magnitude of these counterclockwise traces is affected by the basal expression level, where low basal expression creates an easily observable pattern due to the high fold change values for transcripts and proteins, while high basal expression results in the circular trace to be compressed near the origin. This is

indeed what many experimental datasets look like with metabolic proteins found near the origin in the log2 transcript—protein space (Figure 1), while genes that respond strongly to environmental cues have stronger signals that project into the northeast quadrant of these plots. The model results here may be overly conservative, as it seems likely that marine microbes with much slower reproduction rates than *E. coli* or yeast would also have lower concentrations of RNases and concurrent longer residence times of RNA, as mentioned above. If lower RNase activity were common among marine microbes like *Trichodesmium*, then the overlapping time window for coexistence of RNA and proteins would be enhanced.

Based on this simple modeling approach, the (often negative) conclusions about the coherence of RNA and protein expression in microbes appear to have been clouded by use of relative units of fold change that discriminate against metabolic proteins with larger inventories. In these cases, there may be highly abundant proteins that change in inventory by 20% for example, which would be a large change in copy number but a small change in fold change ratios. Use of absolute units for both transcripts and proteins and efforts to calibrate with standards to achieve copy number per cell estimates would resolve this problem. Importantly, however, for environmental applications, the successful use of environmentally responsive genes in both RNA and protein spaces can be explained by their large signals relative to basal expression that allows easy detection of large fold changes. In short, the lack of observed correlations of metabolic systems represented in transcripts and proteins can be explained, at least partially, by an over-reliance on ratios of data, and those observations should not detract from the statistically significant correlations of environmentally responsive genes in the RNA-protein space and their application as research and monitoring tools in environmental microbiology and biogeochemistry.

While each experimentalist can determine which biomole-cule(s) (proteins and/or transcripts) would complement its research program, these matters often default to practical decisions based on prior training or the availability of sequencing and mass spectrometry facilities. The results and model presented here could help in guiding experimental design and perhaps inspiring future multi-omic studies. In circumstances of meta-omics where environmental samples containing highly diverse biological communities are measured, the measurement of transcripts is particularly useful when environmental eukaryotic organisms have large genomes, as the transcripts allow for the generation of a sequence database without introns for peptide-to-spectrum mapping in data-dependent acquisition approaches.

Together, these empirical observations and model results validate the utility and provide a theoretical basis behind proteins that respond to environmental stimuli in marine microbes. Because organisms that comprise natural ecosystems are frequently limited in their productivity by nutritional elements, they often have evolved specific adaptive responses to nutritional scarcity. As such, these proteins and their mRNA transcripts represent ideal environmental biomarkers to detect controls on ecosystem productivity. Biomarkers for nutritional stress in marine phytoplankton and bacteria that have been identified are shown in Table 2, including those for nitrogen, <sup>43,48–50</sup> phosphorus, <sup>5,6,51–53</sup> iron, <sup>18,19,54–60</sup> vitamin B<sub>12</sub>, <sup>61,62</sup> and zinc. <sup>63</sup> Notably, proteins involved in adaptive responses in marine microbes are often distinct from those model organisms because of their very different environmental conditions. As a

Table 2. Example Biomarkers for Nutrient Stress in Marine Microbes

protein name	nutrient	relevant taxa	example refs
urea transporter	nitrogen	cyanobacteria	43, 48-50
NtcA and P-II N regulatory proteins	nitrogen	cyanobacteria	43 48, and 65
NtrX regulatory protein	nitrogen	Pelagibacter	49
ammonia transporter	nitrogen	cyanobacteria and diatoms	43 and 50
alkaline phosphatase, PhoA	phosphorus	cyanobacteria and Diatoms	5 6, and 51
phosphate transporter, PstS	phosphorus	cyanobacteria	51-53
iron transporter IdiA/ FutA	iron	cyanobacteria	18 19 58, and 59
flavodoxin IsiB and ferredoxin	iron	cyanobacteria and eukaryotic algae	57 43, and 7
iron transporter SfuC	iron	Pelagibacter	60
ISIP2	iron	eukaryotic algae	7 54, and 55
ISIP3	iron	eukaryotic algae	7 and 54
B <sub>12</sub> transporter CBA1	vitamin $B_{12}$	eukaryotic algae	61 and 62
methioine synthase MetE and MetH	vitamin $B_{12}$	eukaryotic algae	61 and 62
Zn chaperone ZCRP-A	zinc and cobalt	eukaryotic algae	63
Zn membrane protein ZCRP-B	zinc and cobalt	eukaryotic algae	63
unknown	Co, Ni, Mn, vitamins, etc.		64

result, several of these marine biomarkers have only recently been discovered and it is almost certain that more remain to be discovered. In some cases, biomarkers for required elements such as cobalt do not appear to have adaptive responses (or the chemical species they respond to has not yet been identified)<sup>64</sup> and corresponding biomarkers perhaps due to the nutritional burden of such systems in highly streamlined genomes. With the natural environment undergoing rapid and unprecedented changes, the ability to observe and understand how the oceans are changing on a global scale is possible through the application of environmental biomarkers. Demonstration that RNA and protein do correlate for environmentally responsive genes and the theory behind when and why these correlations occur is valuable in increasing confidence for study of metaproteomic and other meta-omic studies to the global oceans. <sup>39,65–69</sup>

#### CONCLUSIONS

In summary, these data elucidate broad, nutrient-limited mRNA-protein dynamics following long-term high CO<sub>2</sub> adaptation under multiple nutrient limitation regimes in a biogeochemically important organism. Overall, nutrient-limited proteome variation representing 37% of the IMS101 coding potential was congruent with that of its corresponding transcripts as well as with global transcriptional variation. Pairwise mRNA-protein transcript levels were modestly correlated to protein abundances, yet environmentally responsive transcripts were significantly more correlated than those that did not respond to any of the eight conditions. Wellcharacterized stress mRNA and protein biomarkers responded robustly to their respective limiting nutrients following longterm CO<sub>2</sub> adaptation, thereby supporting their future use under ocean acidification conditions. Other responsive pathways harboring genes with high PC values suggest that they could

be good indicators of cellular physiology, although certain major pathways (e.g., photosynthesis) contain both environmentally responsive (e.g., light reactions) and unresponsive (e.g., photosystem I) genes. Hence, it is important to consider the functions of different components that make up larger biochemical pathways when investigating energy and/or nutrient flux through the cell. It is as equally important to consider the nature of the environmental conditions being tested. As discussed above, environmentally responsive iron stress genes that form complexes with high-inventory photosystems (e.g., IsiA) in iron limitation can yield high RNAprotein correlations across treatments, which may not be the case if another set of conditions was tested that lacked iron stress. These data fill knowledge gaps relating not only to longterm, nutrient-limited mRNA-protein dynamics but also to differences and commonalities among molecular metabolisms of interrelated single- and co-limited physiologies under increasing CO<sub>2</sub>. A model of transcript and protein production and decay was able to reproduce correlations (in fold change units) observed for environmentally responsive genes with small initial inventories, consistent with the observations, demonstrating the mechanics of how both transcripts and proteins can be useful in assessing ecosystem and biogeochemical functions. Although these data highlight molecular dynamics fueling important biogeochemical processes, future studies examining mRNAprotein correlations over longer temporal periods (e.g., diel to several cell cycles) will help in elucidating transcriptional and corresponding protein synthesis times in the face of environmental change.

# ASSOCIATED CONTENT

# Supporting Information

The Supporting Information is available free of charge at https://pubs.acs.org/doi/10.1021/acs.jproteome.1c00517.

(Supplementary Figures S1 and S2) Experimental design and Pearson correlation coefficients, respectively (PDF) (Supplemental Dataset 1) Combined transcriptome—proteome dataset (tabs within the sheet include normalized biological duplicate transcript counts, biological triplicate protein spectral counts, the resulting log2 fold change dataset for transcripts and proteins, and a key for gene annotations) (XLSX)

(Supplemental Movie S1) Time course of model results of transcript and protein inventories and the influence on log2 fold change results (AVI)

## AUTHOR INFORMATION

## **Corresponding Authors**

Mak A. Saito – Marine Chemistry and Geochemistry Department, Woods Hole Oceanographic Institution, Woods Hole, Massachusetts 02543, United States; oocid.org/ 0000-0001-6040-9295; Email: msaito@whoi.edu

Eric A. Webb — Marine and Environmental Biology, Department of Biological Sciences, University of Southern California, Los Angeles, California 90089, United States; Email: eawebb@usc.edu

# **Authors**

Nathan G. Walworth – Marine and Environmental Biology, Department of Biological Sciences, University of Southern California, Los Angeles, California 90089, United States Michael D. Lee — Blue Marble Space Institute of Science, Seattle, Washington 98104, United States; Exobiology Branch, NASA Ames Research Center, Moffett Field, California 94035, United States

Matthew R. McIlvin – Marine Chemistry and Geochemistry Department, Woods Hole Oceanographic Institution, Woods Hole, Massachusetts 02543, United States; orcid.org/ 0000-0002-5301-8365

Dawn M. Moran – Marine Chemistry and Geochemistry Department, Woods Hole Oceanographic Institution, Woods Hole, Massachusetts 02543, United States

Riss M. Kellogg – Marine Chemistry and Geochemistry Department, Woods Hole Oceanographic Institution, Woods Hole, Massachusetts 02543, United States

Fei-Xue Fu — Marine and Environmental Biology, Department of Biological Sciences, University of Southern California, Los Angeles, California 90089, United States

David A. Hutchins – Marine and Environmental Biology, Department of Biological Sciences, University of Southern California, Los Angeles, California 90089, United States

Complete contact information is available at: https://pubs.acs.org/10.1021/acs.jproteome.1c00517

#### **Author Contributions**

\*N.G.W. and M.A.S. are co-first authors.

# **Author Contributions**

F.-X.F., D.A.H., M.A.S., and E.A.W. designed the research; N.G.W., M.D.L., M.A.S., D.M.M., M.R.M., and F.-X.F. performed the research; N.G.W., M.D.L., F.-X.F., D.A.H., and E.A.W. analyzed data; M.A.S. developed the RNA—protein model; and N.G.W., M.A.S., M.D.L., R.M. K., F.-X.F., D.A.H., and E.A.W. wrote the paper.

# Notes

The authors declare no competing financial interest. All RNA-Seq data used in this study have been deposited as raw fastq files in NCBI's Gene Expression Omnibus<sup>47</sup> and are accessible through GEO Series accession number GSE94951. All protein spectral count data used in the above analyses can be found in Supplementary Data 4 of Walworth et al.<sup>23</sup> (DOI: 10. 1038/ncomms12081). Combined transcriptome—proteome dataset and gene annotations are provided as Supplemental Dataset 1. Raw mass spectral datasets used in the above analyses are available at ProteomeXchange under the identifier PXD010515 (DOI: 10.6019/PXD010515). The code for the RNA—protein model is available on GitHub under M.A.S.'s repository: https://github.com/maksaito/RNA-Protein\_MATLAB model.

#### ACKNOWLEDGMENTS

This work was supported by US National Science Foundation Grants OCE 1851222 and OCE 1657757 and (to D.A.H., E.A.W., and F.-X.F.), OCE 1924554, OCE 1850719, and OCE 2019589 (Center for Chemical Currencies on a Microbial Planet), G.B. Moore Foundation, and NIH R01 GM135709 grants (to M.A.S.).

#### REFERENCES

(1) Vogel, C.; Marcotte, E. M. Insights into the regulation of protein abundance from proteomic and transcriptomic analyses. *Nat. Rev. Genet.* **2012**, *13*, 227–232.

- (2) Newman, J. R. S.; Ghaemmaghami, S.; Ihmels, J.; Breslow, D. K.; Noble, M.; DeRisi, J. L.; Weissman, J. S. Single-cell proteomic analysis of *S.cerevisiae* reveals the architecture of biological noise. *Nature* **2006**, 441, 840–846.
- (3) Sowell, S. M.; Norbeck, A. D.; Lipton, M. S.; Nicora, C. D.; Callister, S. J.; Smith, R. D.; Barofsky, D. F.; Giovannoni, S. J. Proteomic Analysis of Stationary Phase in the Marine Bacterium "Candidatus Pelagibacter ubique". Appl. Environ. Microbiol. 2008, 74, 4091–4100.
- (4) Waldbauer, J. R.; Rodrigue, S.; Coleman, M. L.; Chisholm, S. W. Transcriptome and Proteome Dynamics of a Light-Dark Synchronized Bacterial Cell Cycle. *PLoS One* **2012**, *7*, No. e43432.
- (5) Dyhrman, S. T.; Jenkins, B. D.; Rynearson, T. A.; Saito, M. A.; Mercier, M. L.; Alexander, H.; Whitney, L. P.; Drzewianowski, A.; Bulygin, V. V.; Bertrand, E. M.; Wu, Z.; Benitez-Nelson, C.; Heithoff, A. The Transcriptome and Proteome of the Diatom *Thalassiosira pseudonana* Reveal a Diverse Phosphorus Stress Response. *PLoS One* **2012**, 7, No. e33768.
- (6) Wurch, L. L.; Bertrand, E. M.; Saito, M. A.; Van Mooy, B. A. S.; Dyhrman, S. T. Proteome Changes Driven by Phosphorus Deficiency and Recovery in the Brown Tide-Forming Alga *Aureococcus anophagefferens*. *PLoS One* **2011**, *6*, No. e28949.
- (7) Bender, S. J.; Moran, D. M.; McIlvin, M. R.; Zheng, H.; McCrow, J. P.; Badger, J.; DiTullio, G. R.; Allen, A. E.; Saito, M. A. Colony formation in *Phaeocystis antarctica*: connecting molecular mechanisms with iron biogeochemistry. *Biogeosciences* **2018**, *15*, 4923–4942.
- (8) Frischkorn, K. R.; Haley, S. T.; Dyhrman, S. T. Transcriptional and proteomic choreography under phosphorus deficiency and re-supply in the  $N_2$  fixing cyanobacterium *Trichodesmium erythraeum*. Front. Microbiol. **2019**, 10, 330.
- (9) Capone, D. G.; Zehr, J. P.; Paerl, H. W.; Bergman, B.; Carpenter, E. J. *Trichodesmium*, a Globally Significant Marine Cyanobacterium. *Science* 1997, 276, 1221–1229.
- (10) Hutchins, D. A.; Fu, F. Microorganisms and ocean global change. *Nat. Microbiol.* **2017**, *2*, 17058.
- (11) Sohm, J. A.; Webb, E. A.; Capone, D. G. Emerging patterns of marine nitrogen fixation. *Nat. Rev. Microbiol.* **2011**, *9*, 499–508.
- (12) Hutchins, D. A.; Mulholland, M. R.; Fu, F. Nutrient cycles and Marine Microbes in a CO<sub>2</sub>-Enriched Ocean. *Oceanography* **2009**, 22, 128–145.
- (13) Zehr, J. P. Nitrogen fixation by marine cyanobacteria. *Trends Microbiol.* **2011**, *19*, 162–173.
- (14) Luo, Y. W.; Doney, S. C.; Anderson, L. A.; Benavides, M.; Berman-frank, I.; Bode, A.; Bonnet, S.; Boström, K. H.; Böttjer, D.; Capone, D. G.; et al. Database of diazotrophs in global ocean: abundance, biomass and nitrogen fixation rates. *Earth Syst. Sci. Data* **2012**, *4*, 47–73.
- (15) Montoya, J. P.; Holl, C. M.; Zehr, J. P.; Hansen, A.; Villareal, T. A.; Capone, D. G. High rates of  $N_2$  fixation by unicellular diazotrophs in the oligotrophic Pacific Ocean. *Nature* **2004**, *430*, 1027–1031.
- (16) Karl, D. M.; Church, M. J.; Dore, J. E.; Letelier, R. M.; Mahaffey, C. Predictable and efficient carbon sequestration in the North Pacific Ocean supported by symbiotic nitrogen fixation. *Proc. Natl. Acad. Sci. U. S. A.* 2012, 109, 1842–1849.
- (17) Bergman, B.; Sandh, G.; Lin, S.; Larsson, J.; Carpenter, E. J. *Trichodesmium* a widespread marine cyanobacterium with unusual nitrogen fixation properties. *FEMS Microbiol. Rev.* **2013**, *37*, 286–302.
- (18) Webb, E. A.; Jakuba, R. W.; Moffett, J. W.; Dyhrman, S. T. Molecular assessment of phosphorus and iron physiology in *Trichodesmium* populations from the western Central and western South Atlantic. *Limnol. Oceanogr.* **2007**, *52*, 2221–2232.
- (19) Chappell, P. D.; Moffett, J. W.; Hynes, A. M.; Webb, E. A. Molecular evidence of iron limitation and availability in the global diazotroph *Trichodesmium*. *ISME J.* **2012**, *6*, 1728–1739.
- (20) Venter, J. C.; Remington, K.; Heidelberg, J. F.; Halpern, A. L.; Rusch, D.; Eisen, J. A.; Wu, D.; Paulsen, I.; Nelson, K. E.; Nelson, W.; et al. Environmental Genome Shotgun Sequencing of the Sargasso Sea. *Science* **2004**, *304*, *66*–74.
- (21) Dupont, C. L.; McCrow, J. P.; Valas, R.; Moustafa, A.; Walworth, N.; Goodenough, U.; Roth, R.; Hogle, S. L.; Bai, J.; Johnson, Z. I.;

- Mann, E.; Palenik, B.; Barbeau, K. A.; Craig Venter, J.; Allen, A. E. Genomes and gene expression across light and productivity gradients in eastern subtropical Pacific microbial communities. *ISME J.* **2015**, *9*, 1076–1092.
- (22) Moran, M. A. The global ocean microbiome. Science 2015, 350, aac8455.
- (23) Walworth, N. G.; Fu, F.-X.; Webb, E. A.; Saito, M. A.; Moran, D.; Mcllvin, M. R.; Lee, M. D.; Hutchins, D. A. Mechanisms of increased *Trichodesmium* fitness under iron and phosphorus co-limitation in the present and future ocean. *Nat. Commun.* **2016**, *7*, 12081.
- (24) Garcia, N. S.; Fu, F.; Sedwick, P. N.; Hutchins, D. A. Iron deficiency increases growth and nitrogen-fixation rates of phosphorus-deficient marine cyanobacteria. *ISME J.* **2015**, *9*, 238–245.
- (25) Walworth, N. G.; Fu, F.-X.; Lee, M. D.; Cai, X.; Saito, M. A.; Webb, E. A.; Hutchins, D. A. Nutrient-Colimited *Trichodesmium* as a Nitrogen Source or Sink in a Future Ocean. *Appl. Environ. Microbiol.* **2018**, *84*, e02137–e02117.
- (26) Hutchins, D. A.; Boyd, P. W. Marine phytoplankton and the changing ocean iron cycle. *Nat. Clim. Change* **2016**, *6*, 1072–1079.
- (27) Hutchins, D. A.; Walworth, N. G.; Webb, E. A.; Saito, M. A.; Moran, D.; McIlvin, M. R.; Gale, J.; Fu, F.-X. Irreversibly increased nitrogen fixation in *Trichodesmium* experimentally adapted to elevated carbon dioxide. *Nat. Commun.* **2015**, *6*, 8155.
- (28) Sañudo-Wilhelmy, S. A.; Kustka, A. B.; Gobler, C. J.; Hutchins, D. A.; Yang, M.; Lwiza, K.; Burns, J.; Capone, D. G.; Raven, J. A.; Carpenter, E. J. Phosphorus limitation of nitrogen fixation by *Trichodesmium* in the central Atlantic Ocean. *Nature* **2001**, *411*, 66–69.
- (29) Mills, M. M.; Ridame, C.; Davey, M.; La Roche, J.; Gelder, R. J. Iron and phosphorus co-limit nitrogen fixation in the eastern tropical North Atlantic. *Nature* **2004**, 429, 292–294.
- (30) Levitan, O.; Sudhaus, S.; LaRoche, J.; Berman-Frank, I. The Influence of pCO2 and Temperature on Gene Expression of Carbon and Nitrogen Pathways in *Trichodesmium* IMS101. *PLoS One* **2010**, *S*, e15104.
- (31) Walworth, N. G.; Lee, M. D.; Fu, F.-X.; Hutchins, D. A.; Webb, E. A. Molecular and physiological evidence of genetic assimilation to high CO<sub>2</sub> in the marine nitrogen fixer *Trichodesmium*. *Proc. Natl. Acad. Sci. U. S. A.* **2016**, *113*, E7367–E7374.
- (32) Oksanen, J.; Blanchet, F. G.; Kindt, R.; Legendre, P.; Minchin, P. R.; O'hara, R. B.; Simpson, G.L.; Solymos, P.; Stevens, M. H. H.; Wagner, H. *Package 'vegan'*. *Community ecology package, version* 2.; R Core Team: 2013. 2(9), 1–295.
- (33) Suzuki, R.; Shimodaira, H. Pvclust: an R package for assessing the uncertainty in hierarchical clustering. *Bioinformatics* **2006**, 22, 1540–1542.
- (34) Fay, M. P.; Shaw, P. A. Exact and Asymptotic Weighted Logrank Tests for Interval Censored Data: The interval R package. *J. Stat. Softw.* **2010**, *36*, 1–34.
- (35) Benjamini, Y.; Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Series B Stat. Methodol.* **1995**, *57*, 289–300.
- (36) Selinger, D. W.; Saxena, R. M.; Cheung, K. J.; Church, G. M.; Rosenow, C. Global RNA half-life analysis in *Escherichia coli* reveals positional patterns of transcript degradation. *Genome Res.* **2003**, *13*, 216–223
- (37) Pratt, J. M.; Petty, J.; Riba-Garcia, I.; Robertson, D. H. L.; Gaskell, S. J.; Oliver, S. G.; Beynon, R. J. Dynamics of Protein Turnover, a Missing Dimension in Proteomics. *Mol. Cell. Proteomics* **2002**, *1*, 579–591.
- (38) Moran, M. A.; Satinsky, B.; Gifford, S. M.; Luo, H.; Rivers, A.; Chan, L.-K.; Meng, J.; Durham, B. P.; Shen, C.; Varaljay, V. A.; et al. Sizing up metatranscriptomics. *ISME J.* **2013**, *7*, 237–243.
- (39) Santoro, A. E.; Dupont, C. L.; Richter, R. A.; Craig, M. T.; Carini, P.; McIlvin, M. R.; Yang, Y.; Orsi, W. D.; Moran, D. M.; Saito, M. A. Genomic and proteomic characterization of "Candidatus Nitrosopelagicus brevis": An ammonia-oxidizing archaeon from the open ocean. Proc. Natl. Acad. Sci. U. S. A. 2015, 112, 1173–1178.
- (40) Ramette, A. Multivariate analyses in microbial ecology. *FEMS Microbiol. Ecol.* **2007**, *62*, 142–160.

- (41) Rauhut, R.; Klug, G. mRNA degradation in bacteria. FEMS Microbiol. Rev. 1999, 23, 353–370.
- (42) Taniguchi, Y.; Choi, P. J.; Li, G.-W.; Chen, H.; Babu, M.; Hearn, J.; Emili, A.; Xie, X. S. Quantifying *E. coli* Proteome and Transcriptome with Single-Molecule Sensitivity in Single Cells. *Science* **2010**, 329, 533–538.
- (43) Saito, M. A.; McIlvin, M. R.; Moran, D. M.; Goepfert, T. J.; DiTullio, G. R.; Post, A. F.; Lamborg, C. H. Multiple nutrient stresses at intersecting Pacific Ocean biomes detected by protein biomarkers. *Science* **2014**, 345, 1173–1177.
- (44) Morrissey, J.; Bowler, C. Iron utilization in marine cyanobacteria and eukaryotic algae. *Front. Microbiol.* **2012**, *3*, 43.
- (45) Mohr, W.; Intermaggio, M. P.; LaRoche, J. Diel rhythm of nitrogen and carbon metabolism in the unicellular, diazotrophic cyanobacterium *Crocosphaera watsonii* WH8501. *Environ. Microbiol.* **2010**, *12*, 412–421.
- (46) Saito, M. A.; Bertrand, E. M.; Dutkiewicz, S.; Bulygin, V. V.; Moran, D. M.; Monteiro, F. M.; Follows, M. J.; Valois, F. W.; Waterbury, J. B. Iron conservation by reduction of metalloenzyme inventories in the marine diazotroph *Crocosphaera watsonii. Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108*, 2184–2189.
- (47) Edgar, R.; Domrachev, M.; Lash, A. E. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Res.* **2002**, *30*, 207–210.
- (48) Lindell, D.; Padan, E.; Post, A. F. Regulation of ntcA expression and nitrite uptake in the marine *Synechococcus sp.* strain WH 7803. *J. Bacteriol* 1998, 180, 1878–1886.
- (49) Smith, D. P.; Thrash, J. C.; Nicora, C. D.; Lipton, M. S.; Burnum-Johnson, K. E.; Carini, P.; Smith, R. D.; Giovannoni, S. J. Proteomic and transcriptomic analyses of "Candidatus Pelagibacter ubique" describe the first  $P_{II}$ -independent response to nitrogen limitation in a free-living Alphaproteobacterium. *MBio* **2013**, *4*, e001336—e001312.
- (50) Bender, S. J.; Durkin, C. A.; Berthiaume, C. T.; Morales, R. L.; Armbrust, E. V. Transcriptional responses of three model diatoms to nitrate limitation of growth. *Front. Mar. Sci.* **2014**, *1*, 3.
- (51) Orchard, E. D.; Webb, E. A.; Dyhrman, S. T. Molecular analysis of the phosphorus starvation response in *Trichodesmium spp. Environ. Microbiol.* **2009**, *11*, 2400–2411.
- (52) Scanlan, D. J.; Silman, N. J.; Donald, K. M.; Wilson, W. H.; Carr, N. G.; Joint, I.; Mann, N. H. An immunological approach to detect phosphate stress in populations and single cells of photosynthetic picoplankton. *Appl. Environ. Microbiol.* **1997**, *63*, 2411–2420.
- (53) Pereira, N.; Shilova, I. N.; Zehr, J. P. Use of the high-affinity phosphate transporter gene, pstS, as an indicator for phosphorus stress in the marine diazotroph *Crocosphaera watsonii* (Chroococcales, Cyanobacteria). J. Phycol. 2019, 55, 752–761.
- (54) Morrissey, J.; Sutak, R.; Paz-Yepes, J.; Tanaka, A.; Moustafa, A.; Veluchamy, A.; Thomas, Y.; Botebol, H.; Bouget, F. Y.; McQuaid, J. B.; Tirichine, L.; Allen, A. E.; Lesuisse, E.; Bowler, C. A novel protein, ubiquitous in marine phytoplankton, concentrates iron at the cell surface and facilitates uptake. *Curr. Biol.* **2015**, *25*, 364–371.
- (55) McQuaid, J. B.; Kustka, A. B.; Oborník, M.; Horák, A.; McCrow, J. P.; Karas, B. J.; Zheng, H.; Kindeberg, T.; Andersson, A. J.; Barbeau, K. A.; Allen, A. E. Carbonate-sensitive phytotransferrin controls high-affinity iron uptake in diatoms. *Nature* **2018**, *555*, 534–537.
- (56) Behnke, J.; LaRoche, J. Iron uptake proteins in algae and the role of Iron Starvation-Induced Proteins (ISIPs). *Eur. J. Phycol.* **2020**, *55*, 339–360.
- (57) Roche, J. L.; Murray, H.; Orellana, M.; Newton, J. Flavodoxin Expression as an Indicator of Iron Limitation in Marine Diatoms<sup>1</sup>. *J. Phycol.* **1995**, *31*, 520–530.
- (58) Chappell, P. D.; Webb, E. A. A molecular assessment of the iron stress response in the two phylogenetic clades of *Trichodesmium*. *Environ. Microbiol.* **2010**, 12, 13–27.
- (59) Held, N. A.; Webb, E. A.; McIlvin, M. M.; Hutchins, D. A.; Cohen, N. R.; Moran, D. M.; Kunde, K.; Lohan, M. C.; Mahaffey, C.; Woodward, E. M. S.; Saito, M. A. Co-occurrence of Fe and P stress in natural populations of the marine diazotroph *Trichodesmium*. *Biogeosciences* **2020**, *17*, 2537–2551.

- (60) Smith, D. P.; Kitner, J. B.; Norbeck, A. D.; Clauss, T. R.; Lipton, M. S.; Schwalbach, M. S.; Steindler, L.; Nicora, C. D.; Smith, R. D.; Giovannoni, S. J. Transcriptional and translational regulatory responses to iron limitation in the globally distributed marine bacterium *Candidatus* Pelagibacter ubique. *PLoS One* **2010**, *5*, No. e10487.
- (61) Bertrand, E. M.; Allen, A. E.; Dupont, C. L.; Norden-Krichmar, T. M.; Bai, J.; Valas, R. E.; Saito, M. A. Influence of cobalamin scarcity on diatom molecular physiology and identification of a cobalamin acquisition protein. *Proc. Natl. Acad. Sci. U. S. A.* **2012**, *109*, E1762—E1771.
- (62) Bertrand, E. M.; Moran, D. M.; McIlvin, M. R.; Hoffman, J. M.; Allen, A. E.; Saito, M. A. Methionine synthase interreplacement in diatom cultures and communities: Implications for the persistence of  $B_{12}$  use by eukaryotic phytoplankton. *Limnol. Oceanogr.* **2013**, *58*, 1431-1450.
- (63) Kellogg, R. M., M. A., Moosburner, N. R., Cohen, N. J., Hawco, M. R., McIlvin, D. M., Moran, A. E., Allen, M. A., Saito Two novel zinc and cobalt responsive proteins in marine eukaryotic phytoplankton and implications for Zn scarcity in the surface oceans. *Nat. Commun.*, **2021**. Submitted.
- (64) Hawco, N. J.; McIlvin, M. M.; Bundy, R. M.; Tagliabue, A.; Goepfert, T. J.; Moran, D. M.; Valentin-Alvarado, L.; DiTullio, G. R.; Saito, M. A. Minimal cobalt metabolism in the marine cyanobacterium *Prochlorococcus. Proc. Natl. Acad. Sci. U. S. A.* **2020**, *117*, 15740–15747.
- (65) Saito, M. A.; Dorsk, A.; Post, A. F.; McIlvin, M. R.; Rappé, M. S.; DiTullio, G. R.; Moran, D. M. Needles in the blue sea: Sub-species specificity in targeted protein biomarker analyses within the vast oceanic microbial metaproteome. *Proteomics* **2015**, *15*, 3521–3531.
- (66) Saito, M. A.; Bertrand, E. M.; Duffy, M. E.; Gaylord, D. A.; Held, N. A.; Hervey, W. J., IV; Hettich, R. L.; Jagtap, P. D.; Janech, M. G.; Kinkade, D. B.; Leary, D. H.; McIlvin, M. R.; Moore, E. K.; Morris, R. M.; Neely, B. A.; Nunn, B. L.; Saunders, J. K.; Shepherd, A. I.; Symmonds, N. I.; Walsh, D. A. Progress and challenges in ocean metaproteomics and proposed best practices for data sharing. *J. Proteome Res.* **2019**, *18*, 1461–1476.
- (67) Saito, M. A.; Saunders, J. K.; Chagnon, M.; Gaylord, D. A.; Shepherd, A.; Held, N. A.; Dupont, C.; Symmonds, N.; York, A.; Charron, M.; Kinkade, D. B. Development of an Ocean Protein Portal for Interactive Discovery and Education. *J. Proteome Res.* **2021**, *20*, 326–336.
- (68) Villar, E.; Vannier, T.; Vernette, C.; Lescot, M.; Cuenca, M.; Alexandre, A.; Bachelerie, P.; Rosnet, T.; Pelletier, E.; Sunagawa, S.; Hingamp, P. The Ocean Gene Atlas: exploring the biogeography of plankton genes online. *Nucleic Acids Res.* **2018**, 46, W289–W295.
- (69) Ustick, L. J.; Larkin, A. A.; Garcia, C. A.; Garcia, N. S.; Brock, M. L.; Lee, J. A.; Wiseman, N. A.; Moore, J. K.; Martiny, A. C. Metagenomic analysis reveals global-scale patterns of ocean nutrient limitation. *Science* **2021**, *372*, 287–291.