1 **Normalization of metatranscriptomic and metaproteomic data for**

2 **differential gene expression analyses: The importance of accounting**

3 **for organism abundance**

4

5

6 **Author:** Manuel Kleiner

7

8 **Affiliations:**

9 Energy Bioengineering and Geomicrobiology Group, Department of Geoscience, University of Calgary,

10 Calgary, Canada

11 Department of Plant and Microbial Biology, North Carolina State University, Raleigh, NC, USA

12

13 **Correspondence:** manuel.kleiner@ucalgary.ca

14

15 **Article Format**: Perspective Piece

16

## Abstract

Metatranscriptomics and metaproteomics make it possible to measure gene expression in microbial communities. So far these approaches were mostly used to get a general overview of the dominant metabolism and physiologies of community members. Recently, environmental microbiologists have started using metatranscriptomics and metaproteomics to look at gene expression differences between different environments or conditions. This has been mostly done by using makeshift adaptations of pure culture focused differential transcriptomics and proteomics approaches. However, since meta-omics data has many more variables attached to it as compared to pure culture derived data, such makeshift adaptations are problematic at best. One particular challenge is posed by the data normalization strategies used to account for technical and biological variables in meta-omic data. Here I discuss the most common normalization strategy for transcriptomic and proteomic data and why it is not valid by itself for meta-omic data. I provide logical proof that variation in species abundances between samples is an additional variable that must be accounted for during normalization of meta-omic data. Finally, I show how the existing normalization methods for transcriptomic and proteomic data can be augmented to be applicable to meta-omic data.

2

## Main Text

34    **Main Text**

35    In the last decade technological advances in sequencing technology and mass spectrometry have made it

36    possible to measure gene expression in microbial communities on a large scale. The respective approaches

37    have been termed metatranscriptomics and metaproteomics [1, 2]. Metatranscriptomics is an umbrella

38    term for methods that measure the transcription levels in microbial communities and metaproteomics is

39    the corresponding term for methods that measure the protein abundances in microbial communities. The

40    outputs of both methods are tables which list gene expression values for individual genes (rows) across

41    multiple samples (columns). See the first worksheet in the supplemental table for a simulated example

42    (Supplementary Table S1). For metatranscriptomics, the expression values are usually based on the

43    counting of reads mapped to a set of reference genes/genomes. For metaproteomics, the expression values

44    are based on the number of spectra matching to reference protein sequences (spectral counting based

45    methods) or on the chromatographic peak intensities of peptides that match to reference protein sequences

46    [3]. These raw counts or intensities are usually converted into proportional (relative) data that gives

47    individual gene expression values as a fraction of 1. This conversion process is part of the data

48    normalization discussed below.

49    Initially metatranscriptomics and metaproteomics were mostly used for discovery based studies that

50    addressed the question which genes are expressed in the community and which proteins are the most

51    important players [4, 5]. In more recent years researchers have started to use these methods for a more in

52    depth investigation of how gene expression differs between different environmental sites, seasons or real

53    or artificially induced changes (e.g. [6-8]). So we are now entering an era in which we start applying

54    differential metatranscriptomics and metaproteomics. So far most differential meta-omics studies have

55    used makeshift adaptations of well-established differential transcriptomics and proteomics methods that

56    were developed for single-organism applications.

57    Metatranscriptomics and metaproteomics come with their own specific set of methodological challenges

58    including, for example, sample extraction biases, contaminants, the construction of suitable reference

59    databases and problems with database redundancies. These challenges are or will be discussed elsewhere

60    [9-11].

61    Here I will discuss data normalization for differential gene expression analyses of metatranscriptomes and

62    metaproteomes, which differs in part from the normalization steps required for differential transcriptomics

63    and proteomics. To make samples comparable on a gene expression level for transcriptomics and

64    proteomics the necessity for two normalization steps is widely accepted [12-15]: (i) In the first

65    normalization step, the expression values are adjusted for the gene/protein sequence length, which can for

66    example be done by simple division of the expression values by gene length. This normalization step is

3

67      justified by the fact that both metatranscriptomics and the metaproteomics will yield higher raw

68      expression values (read counts, spectral counts or summed peptide intensities) for larger

69      transcripts/proteins. (ii) In the second normalization step, the expression values are adjusted for variations

70      in the sum of expression values for each sample (column). After this normalization step the sum of

71      expression values for each sample should be identical across all samples (e.g. if you normalize to %, the

72      sum of each column should be 100). This normalization step is justified and needed because of technical

73      variations between sample runs. In nextSeq based metatranscriptomics each sample will for example yield

74      a different number of total reads, while in metaproteomics variation between runs can lead to difference in

75      total spectral counts or peptide intensities. These normalization steps have been implemented in many

76      different forms for both transcriptomics and proteomics and are reviewed elsewhere [12-14]. Suitable

77      implementations of this normalization scheme for transcriptomics are the transcripts per million (TPMs)

78      metric [12] and for proteomics either normalized spectral abundance factors (NSAFs) [14, 16] or for

79      peptide intensities MaxLFQ [17].

80      For metatranscriptomes and metaproteomes an additional level of variation needs to be considered when

81      comparing expression differences between genes of individual organisms. This additional level is

82      variation of organism abundances between samples. Here an important differentiation has to be made, as

83      the kind of normalization required in meta-omes very much depends on the exact question asked:

84          (a) If your question is of the type: "Does the expression of *gene*A contribute a higher number of

85                transcripts/protein mass to COMMUNITY1 as compared to COMMUNITY2?" OR "Which genes

86                differ in contribution to total community transcript number or protein mass between COMMUNITY1

87                and COMMUNITY2?", then the above described two-step normalization scheme for transcriptomics

88                and proteomics is perfectly adequate.  To give two concrete example for such questions (1) "Is the

89                human structural protein collagen enriched in the intestinal microbiome samples of PERSON1

90                versus PERSON2?" and (2) "Are Carbohydrate-Active Enzymes (*CAZymes*) overall more

91                abundantly expressed in COMMUNITY1 as compared to COMMUNITY2

92         (b) If your question is of the type: "Is the expression of *gene*A from SPECIESX higher in SPECIESX in

93                COMMUNITY1 as compared to SPECIESX in COMMUNITY2?" OR "Which genes differ in expression

94                between COMMUNITY1 and COMMUNITY2 on the species level?", then the above described two-

95                step normalization scheme for transcriptomics and proteomics by itself is not valid. As I will

96                prove here, an additional normalization step is needed after the two-step normalization to account

97                for variation in species/strain abundances between samples.

98      Generally, there are at least two ways to provide evidence or proof for this. First, one could generate

99      empirical data using two or more mock communities made with the same species, but different species

4

100 abundances. This data could then be used to validate normalizations methods. Such mock community

101 studies have helped to validate other omics methods for environmental microbiology in the past e.g.

102 methods for quantitative metagenomic sequencing [18, 19]. The second approach that one can use in this

103 case is to do a thought experiment to show that the comparison of expression values is invalid if the data is

104 not corrected for variation in species abundance in each sample (and valid if the correction is done). I will

105 use simulated datasets that represent two extreme cases for this thought experiment.

106 **To re-iterate the assumptions:**

107     (1) Gene expression is measured for a microbial community with >1 species.

108     (2) Gene expression values have been normalized to gene length and the sum of expression values in

109         each sample (column).

110     (3) We ask a question of the type (b) above.

111 **Proof:**

112 In the first worksheet of the supplemental table the simplest case of a microbial community is shown: one

113 with only two community member species. To keep it simple, I assume that for each of the two species

114 gene expression was detected for 50 genes and that the expression of all genes is identical. To emphasize

115 the importance of replication for differential omics [20], I show 6 replicate columns; although for the

116 purpose of this proof replication is not really relevant.

117 To show the effect of relative species abundance in the community on gene expression data I have

118 simulated the gene expression data for two distinct species abundance profiles. Samples 1 through 6 come

119 from a community in which both species have the same abundance (1:1 abundance ratio). In samples 7 to

120 12 the same exact gene expression patterns are shown, but expression values have been adjusted to be

121 coming from a very different species abundance profile (species ratio is 20:1). Without the need for

122 statistical tests, it becomes immediately clear that the expression of individual genes would be considered

123 to be different between the two community types. This proves that for type (b) questions two-step

124 normalized data is not sufficient.

125 **How to normalize expression data for species abundances?**

126 Now the question is of course how to actually normalize the data to species abundance. The simplest way

127 is to normalize the expression values for each sample and species to a constant value (i.e. the sum of

128 expression values for each species in each sample should be the same after normalization), which make

129 expression values comparable across samples as the effect of different species abundance profiles is

130 removed. A simple implementation of this is shown in the second worksheet of the supplemental table. An

5

131 implementation of this procedure for spectral counting based metaproteomics was published by Mueller et

132 al. [21] and has been used in many other metaproteomics studies [7, 22, 23]. One important thing to check

133 before normalizing to species/strains is that there are enough measurements (e.g. read counts, spectral

134 counts) for the species/strain to be normalized to. This is crucial to avoid skewing the data simply because

135 there are only very few transcripts/proteins to be considered for the respective species/strain.

136 Normalization to species could actually be abolished if only the reference genome/protein sequences of

137 the organism of interest were used for generating the expression profile data by read mapping or spectral

138 counting. However, using only a subset of reference sequences for the generation of expression data

139 carries the danger of reads or spectra falsely mapping to this reference due to the absence of the

140 potentially better matching reference sequences of the other community members. For

141 metatranscriptomics this can be alleviated by using very strict read mapping criteria, i.e. only use counts

142 from reads mapped with very high identity. For metaproteomics, the strategy of only using the target

143 organism reference genome cannot be recommended, because spectra that would match non-uniquely to

144 multiple sequences if the complete database were used, may match uniquely to a single protein sequence if

145 a limited set of sequences is used (for more details on the so called protein inference problem see [24]).

146 There are several alternative approaches to data intrinsic normalization that could be used. First,

147 abundance profile data obtained with other methods, e.g. 16S rRNA amplicon sequencing or metagenomic

148 sequencing, could be used to correct expression values for each sample. However, this kind of data might

149 bring its own skews and biases into the normalization procedure. Second, spiking in of known amounts of

150 mRNA or protein into samples prior to extraction allows estimating transcript or protein abundances in

151 relation to the standard. This spike in strategy can provide absolute per cell quantification if cell numbers

152 are determined prior to extraction [10, 25].

153 A normalization of expression values to housekeeping genes, which is sometimes used for transcriptomic

154 and proteomic data [26] can currently not be used for metatranscriptomics and metaproteomic data. A

155 housekeeping gene based normalization requires that the housekeeping gene in question is quantified as a

156 function of cell number or cell mass for all conditions that will be considered in a differential expression

157 experiment. In theory, such a correlation of cell number with housekeeping gene expression could be

158 measured for members of a microbial community e.g. by using a combination of mRNA FISH with 16S

159 rRNA FISH, however, the effort required for this seems prohibitive, particularly since much simpler

160 methods are already available.

161 **What comes after normalization?**

6

162     Of course, the normalization steps are only a small part of the workflow for looking at gene expression

163     differences. After normalization of the data, simple checks should be done to test the overall validity of

164     the data and to discover potential sample mixups and alike. This can, for example, be done by hierarchical

165     clustering or principal component analysis of samples based on expression values. Here you should see a

166     separation of samples based on the sampling sites or conditions used. If all seems in order, one can

167     proceed with statistical testing for differential gene expression. Statistical methods for differential gene

168     expression analyses, which have to account for the multiple testing problem, missing values and the fact

169     that the normalized gene expression data represents compositional data, are discussed elsewhere [15, 27-

170     29].

## Acknowledgments

# References

1.  Warnecke F, Hess M: **A perspective: Metatranscriptomics as a tool for the discovery of novel biocatalysts**. *J Biotechnol* 2009, **142**:91-95.

2.  Hettich RL, Sharma R, Chourey K, Giannone RJ: **Microbial metaproteomics: identifying the repertoire of proteins that microorganisms use to compete and cooperate in complex environmental communities**. *Curr Opin Microbiol* 2012, **15**:373-380.

3.  Nahnsen S, Bielow C, Reinert K, Kohlbacher O: **Tools for label-free peptide quantification**. *Mol Cell Proteomics* 2013, **12**:549-556.

4.  Williams TJ, Cavicchioli R: **Marine metaproteomics: deciphering the microbial metabolic food web**. *Trends Microbiol* 2014, **22**:248-260.

5.  VerBerkmoes NC, Denef VJ, Hettich RL, Banfield JF: **Systems Biology: Functional analysis of natural microbial consortia using community proteomics**. *Nat Rev Microbiol* 2009, **7**:196-205.

6.  Hawley AK, Brewer HM, Norbeck AD, Paša-Tolić L, Hallam SJ: **Metaproteomics reveals differential modes of metabolic coupling among ubiquitous oxygen minimum zone microbes**. *Proc Natl Acad Sci USA* 2014, **111**:11395-11400.

7.  Ponnudurai R, Kleiner M, Sayavedra L, Petersen JM, Moche M, Otto A, Becher D, Takeuchi T, Satoh N, Dubilier N, Schweder T, Markert S: **Metabolic and physiological interdependencies in the *Bathymodiolus azoricus* symbiosis**. *ISME J* 2017, **11**:463-477.

8.  Shi W, Moon CD, Leahy SC, Kang D, Froula J, Kittelmann S, Fan C, Deutsch S, Gagic D, Seedorf H, Kelly WJ, Atua R, Sang C, Soni P, Li D, Pinares-Patiño CS, McEwan JC, Janssen PH, Chen F, Visel A, Wang Z, Attwood GT, Rubin EM: **Methane yield phenotypes linked to differential gene expression in the sheep rumen microbiome**. *Genome Res* 2014, **24**:1517-1525.

9.  Wilmes P, Heintz-Buschart A, Bond PL: **A decade of metaproteomics: Where we stand and what the future holds**. *Proteomics* 2015, **15**:3409-3417.

10. Moran MA, Satinsky B, Gifford SM, Luo H, Rivers A, Chan L-K, Meng J, Durham BP, Shen C, Varaljay VA, Smith CB, Yager PL, Hopkinson BM: **Sizing up metatranscriptomics**. *ISME J* 2013, **7**:237-243.

11. Kuske CR, Hesse CN, Challacombe JF, Cullen D, Herr JR, Mueller RC, Tsang A, Vilgalys R: **Prospects and challenges for fungal metatranscriptomics of complex communities**. *Fungal Ecology* 2015, **14**:133-137.

12. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, Szcześniak MW, Gaffney DJ, Elo LL, Zhang X, Mortazavi A: **A survey of best practices for RNA-seq data analysis**. *Genome Biol* 2016, **17**:13.

13. Dillies M-A, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, Keime C, Marot G, Castel D, Estelle J, Guernec G, Jagla B, Jouneau L, Laloë D, Le Gall C, Schaëffer B, Le Crom S, Guedj M, Jaffrézic F: **A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis**. *Brief Bioinform* 2012, **14**:671-683.

14. McIlwain S, Mathews M, Bereman MS, Rubel EW, MacCoss MJ, Noble WS: **Estimating relative abundances of proteins from shotgun proteomics data**. *BMC Bioinf* 2012, **13**:308.

15. Rapaport F, Khanin R, Liang Y, Pirun M, Krek A, Zumbo P, Mason CE, Socci ND, Betel D: **Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data**. *Genome Biol* 2013, **14**:3158.

16. Florens L, Carozza MJ, Swanson SK, Fournier M, Coleman MK, Workman JL, Washburn MP: **Analyzing chromatin remodeling complexes using shotgun proteomics and normalized spectral abundance factors**. *Methods* 2006, **40**:303-311.

17. Cox J, Hein MY, Luber CA, Paron I, Nagaraj N, Mann M: **Accurate proteome-wide label-free quantification by delayed normalization and maximal peptide ratio extraction, termed MaxLFQ**. *Mol Cell Proteomics* 2014, **13**.

8

225   18.   Bowers RM, Clum A, Tice H, Lim J, Singh K, Ciobanu D, Ngan CY, Cheng J-F, Tringe SG,
226         Woyke T: **Impact of library preparation protocols and template quantity on the**
227         **metagenomic reconstruction of a mock microbial community**. *BMC Genomics* 2015, **16**:856.
228   19.   Kleiner M, Hooper LV, Duerkop BA: **Evaluation of methods to purify virus-like particles for**
229         **metagenomic sequencing of intestinal viromes**. *BMC Genomics* 2015, **16**:7.
230   20.   Prosser JI: **Replicate or lie**. *Environ Microbiol* 2010, **12**:1806-1810.
231   21.   Mueller RS, Denef VJ, Kalnejais LH, Suttle KB, Thomas BC, Wilmes P, Smith RL, Nordstrom
232         DK, McCleskey RB, Shah MB, VerBerkmoes NC, Hettich RL, Banfield JF: **Ecological**
233         **distribution and population physiology defined by proteomics in a natural microbial**
234         **community**. *Mol Syst Biol* 2010, **6**:374.
235   22.   Hamann E, Gruber-Vodicka H, Kleiner M, Tegetmeyer HE, Riedel D, Littmann S, Chen J,
236         Milucka J, Viehweger B, Becker KW, Dong X, Stairs CW, Hinrichs K-U, Brown MW, Roger AJ,
237         Strous M: **Environmental Breviatea harbour mutualistic *Arcobacter* epibionts**. *Nature* 2016,
238         **534**:254-258.
239   23.   Justice NB, Pan C, Mueller R, Spaulding SE, Shah V, Sun CL, Yelton AP, Miller CS, Thomas
240         BC, Shah M, VerBerkmoes N, Hettich R, Banfield JF: **Heterotrophic archaea contribute to**
241         **carbon cycling in low-pH, suboxic biofilm communities**. *Appl Environ Microbiol* 2012,
242         **78**:8321-8330.
243   24.   Nesvizhskii AI, Aebersold R: **Interpretation of shotgun proteomic data: The protein inference**
244         **problem**. *Mol Cell Proteomics* 2005, **4**:1419-1440.
245   25.   Zauber H, Schüler V, Schulze W: **Systematic evaluation of reference protein normalization in**
246         **proteomic experiments**. *Frontiers in Plant Science* 2013, **4**:25.
247   26.   Ferguson RE, Carroll HP, Harris A, Maher ER, Selby PJ, Banks RE: **Housekeeping proteins: A**
248         **preliminary study illustrating some limitations as useful references in protein expression**
249         **studies**. *Proteomics* 2005, **5**:566-571.
250   27.   Tyanova S, Temu T, Sinitcyn P, Carlson A, Hein MY, Geiger T, Mann M, Cox J: **The Perseus**
251         **computational platform for comprehensive analysis of (prote)omics data**. *Nat Meth* 2016,
252         **13**:731-740.
253   28.   Fernandes AD, Reid JNS, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB: **Unifying the**
254         **analysis of high-throughput sequencing datasets: characterizing RNA-seq, 16S rRNA gene**
255         **sequencing and selective growth experiments by compositional data analysis**. *Microbiome*
256         2014, **2**:15-15.
257   29.   Weiss SJ, Xu Z, Amir A, Peddada S, Bittinger K, Gonzalez A, Lozupone C, Jesse R Zaneveld,
258         Vazquez-Baeza Y, Birmingham A, Knight R: **Effects of library size variance, sparsity, and**
259         **compositionality on the analysis of microbiome data**. *PeerJ PrePrints* 2015, **3**:e1157v1151

260

261