

This work was written as part of one of the author's official duties as an Employee of the United States Government and is therefore a work of the United States Government. In accordance with 17 U.S.C. 105, no copyright protection is available for such works under U.S. Law. Access to this work was provided by the University of Maryland, Baltimore County (UMBC) ScholarWorks@UMBC digital repository on the Maryland Shared Open Access (MD-SOAR) platform.

Please provide feedback

Please support the ScholarWorks@UMBC repository by emailing [scholarworks-group@umbc.edu](mailto:scholarworks-group@umbc.edu) and telling us what having access to this work means to you and why it's important to you. Thank you.



# Reply to: Examining microbe–metabolite correlations by linear methods

James T. Morton<sup>1,2</sup>, Daniel McDonald<sup>1,3</sup>, Alexander A. Aksenov<sup>4,5</sup>, Louis Felix Nothias<sup>4,5</sup>, James R. Foulds<sup>6</sup>, Robert A. Quinn<sup>7</sup>, Michelle H. Badri<sup>8</sup>, Tami L. Swenson<sup>9</sup>, Marc W. Van Goethem<sup>9</sup>, Trent R. Northen<sup>9,10</sup>, Yoshiki Vazquez-Baeza<sup>3,11</sup>, Mingxun Wang<sup>4,5</sup>, Nicholas A. Bokulich<sup>12,13</sup>, Aaron Watters<sup>14</sup>, Se Jin Song<sup>1,3</sup>, Richard Bonneau<sup>8,14,15,16</sup>, Pieter C. Dorrestein<sup>4,5</sup> and Rob Knight<sup>1,2,17,18</sup> ✉

REPLYING TO T. P. Quinn & I. Erb *Nature Methods* <https://doi.org/10.1038/s41592-020-01006-1> (2020)

Quinn and Erb<sup>1</sup> propose to apply a centered log-ratio (CLR) transform before performing correlation analysis and make the case that, when used correctly, correlation and proportionality can outperform MMvec in identifying microbe–metabolite interactions. While this may be an appealing strategy, it is important to note that the correlations estimated from CLR-transformed data will have a fundamentally different interpretation than the true correlations in the environment, namely:

$$\text{Cov}(x_i, y_j) \neq \text{Cov}(\text{clr}(x)_i, \text{clr}(y)_j)$$

where  $x_i$  and  $y_j$  are the absolute abundances for microbe abundances  $x$  and metabolite abundances  $y$  in taxon  $i$  and metabolite  $j$ . Because the absolute abundances are often not available, inferring the true correlations between microbes and metabolites is not tractable (Supplementary Note 1). This phenomenon has been extensively studied in refs. 2–4, and one of our recent studies provides the intuition behind this in the case of differential abundance<sup>5</sup>. Because of this discrepancy, we proposed to use co-occurrence probabilities instead of correlation.

We relied on simulated data in the original paper<sup>6</sup> as an artificial ground truth, as is common in the evaluation of omics tools. However, simulated data will always have limitations because of the inability to model unknown features of the real system or because of deliberate simplifications that clarify key points in the model system. Furthermore, it is possible to identify simulations where a proposed model is optimal. In Fig. 1, we used Bayesian Optimization<sup>7</sup> to identify simulations where MMvec was able to accurately estimate the correct parameters and Pearson underperformed. If the appropriate assumptions are satisfied, MMvec can correctly estimate the co-occurrence probabilities with machine precision.

Therefore, a crucial aspect of the MMvec manuscript was to test performance both on simulations and on real data. Performance on real data is the ultimate test of methods, and we recommend that simulated datasets be complemented with experimentally validated datasets where possible. Accordingly, we applied the same proportionality-based scripts described by Quinn and Erb<sup>1</sup> and evaluated them on one of the real datasets we used in the MMvec paper.

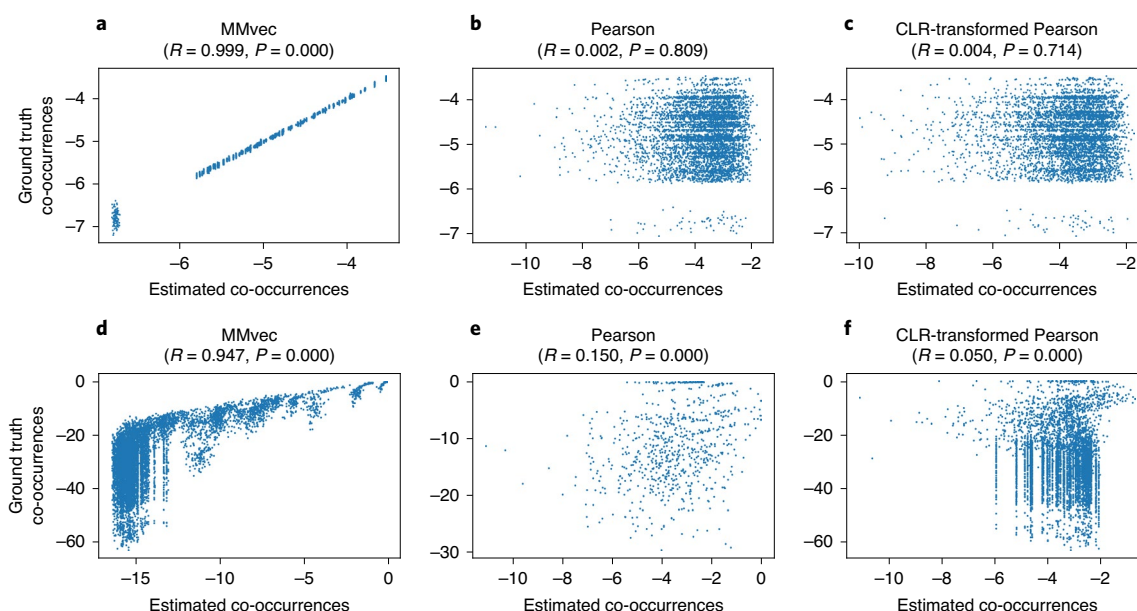
A major obstacle to analyzing real-world microbiome and metabolomics data is sparsity. Traditional compositional methods such as the proposed CLR transform cannot automatically deal with zeros and require imputation as a preprocessing step. This imputation adds bias and is impractical for the sparse datasets typically encountered<sup>8,9</sup>. Microbiome and untargeted metabolomics datasets are generally sparse: in large studies, such as the American Gut Project<sup>10</sup>, the sparsity for stool samples alone is 99.946%. MMvec was designed to handle sparse data. In the desert biocrust soils dataset (sparsity of 51%; ref. 11) that was used in the MMvec publication, we observe that MMvec dramatically outperformed the newly proposed linear methods (Fig. 2).

Contrary to the argument by Quinn and Erb<sup>1</sup> regarding the complexity of neural networks, the MMvec model<sup>6</sup> is not much more complex than the proposed regression techniques. It is a simple one-layer neural network, which is in effect a two-stage log–bilinear regression.

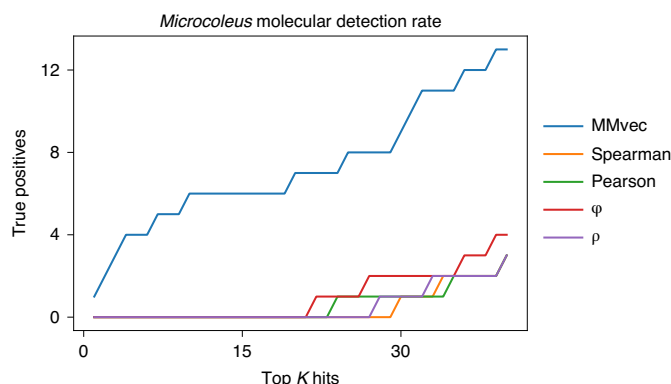
Methods similar to MMvec have been successful at the task of learning word co-occurrences. Since Mikolov et al.<sup>12</sup>, these models have been designed with an emphasis on practical methods for learning useful word representations at scale, rather than on perfectly modeling the data distribution.

MMvec is only one tool in the arsenal of correlative methods. It is not perfect for every correlation type or dataset and is not a one-size-fits-all solution. However, we have found that MMvec is a powerful discovery tool, as demonstrated by the other real datasets

<sup>1</sup>Department of Pediatrics, University of California, San Diego, La Jolla, CA, USA. <sup>2</sup>Department of Computer Science and Engineering, University of California, San Diego, La Jolla, CA, USA. <sup>3</sup>Center for Microbiome Innovation, University of California, San Diego, La Jolla, CA, USA. <sup>4</sup>Collaborative Mass Spectrometry Innovation Center, University of California, San Diego, La Jolla, CA, USA. <sup>5</sup>Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California, San Diego, La Jolla, CA, USA. <sup>6</sup>Department of Information Systems, University of Maryland–Baltimore County, Baltimore, MD, USA. <sup>7</sup>Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI, USA. <sup>8</sup>Department of Biology, New York University, New York, NY, USA. <sup>9</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. <sup>10</sup>DOE Joint Genome Institute, Walnut Creek, CA, USA. <sup>11</sup>Jacobs School of Engineering, University of California, San Diego, La Jolla, CA, USA. <sup>12</sup>The Pathogen and Microbiome Institute, Northern Arizona University, Flagstaff, AZ, USA. <sup>13</sup>Department of Biological Sciences, Northern Arizona University, Flagstaff, AZ, USA. <sup>14</sup>Flatiron Institute, Simons Foundation, New York, NY, USA. <sup>15</sup>Department of Computer Science, Courant Institute, New York, NY, USA. <sup>16</sup>Center for Data Science, New York University, New York, NY, USA. <sup>17</sup>Department of Bioengineering, University of California, San Diego, La Jolla, CA, USA. <sup>18</sup>Center for Microbiome Innovation, University of California, San Diego, La Jolla, CA, USA. ✉e-mail: [rknight@ucsd.edu](mailto:rknight@ucsd.edu)



**Fig. 1 | A simulation benchmark comparing MMvec to Pearson.** Simulations were obtained through Bayesian Optimization<sup>7</sup> to showcase scenarios where MMvec outperforms Pearson. **a–c**, Simulation of a scenario where the microbiome dataset is 99% dense. **d–f**, Simulation of a scenario where the microbiome dataset is 60% dense. All axes are represented on a log scale. Pearson's  $R$  is used to measure the agreement between the simulated ground truth co-occurrences and the estimated co-occurrences.



**Fig. 2 | Biocrust soils benchmark.** A comparison of MMvec to metrics proposed by Quinn and Erb<sup>1</sup>. These proposed metrics include Spearman, Pearson,  $\phi$  and  $\rho$  applied after a CLR transformation<sup>13</sup>.

we evaluated in the original article. It is critical that we provide accurate guidance to the community so that scenarios where one method works better than others are well understood. While there may be scenarios where linear methods outperform neural networks, we show that there are scenarios where neural networks outperform linear methods. We appreciate the communication on the topic to the extent that it helps the community better understand the advantages and limitations of the different approaches and prompts the community to continue to innovate in this area.

### Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41592-020-01007-0>.

Received: 17 March 2020; Accepted: 27 October 2020;  
Published online: 4 January 2021

### References

- Quinn, T. P. & Erb, I. Examining microbe–metabolite correlations by linear methods. *Nat. Methods* <https://doi.org/10.1038/s41592-020-01006-1> (2020).
- Aitchison, J. A concise guide to compositional data analysis. [http://www.leg.ufpr.br/lib/exe/fetch.php/pessoais:abtmarins:a\\_concise\\_guide\\_to\\_compositional\\_data\\_analysis.pdf](http://www.leg.ufpr.br/lib/exe/fetch.php/pessoais:abtmarins:a_concise_guide_to_compositional_data_analysis.pdf) (2003).
- Filzmoser, P. & Hron, K. Correlation analysis for compositional data. *Math. Geosci.* **41**, 905 (2009).
- Friedman, J. & Alm, E. J. Inferring correlation networks from genomic survey data. *PLoS Comput. Biol.* **8**, e1002687 (2012).
- Morton, J. T. et al. Establishing microbial composition measurement standards with reference frames. *Nat. Commun.* **10**, 2719 (2019).
- Morton, J. T. et al. Learning representations of microbe–metabolite interactions. *Nat. Methods* **16**, 1306–1314 (2019).
- Nogueira, F. Bayesian Optimization: open source constrained global optimization tool for Python. <https://github.com/fmfn/BayesianOptimization> (2014).
- Martín-Fernández, J. A., Barceló-Vidal, C. & Pawłowsky-Glahn, V. Dealing with zeros and missing values in compositional data sets using nonparametric imputation. *Math. Geol.* **35**, 253–278 (2003).
- Silverman, J. D., Roche, K., Mukherjee, S. & David, L. A. Naught all zeros in sequence count data are the same. Preprint at *bioRxiv* <https://doi.org/10.1101/477794> (2018).
- McDonald, D. et al. American Gut: an open platform for citizen science microbiome research. *mSystems* **3**, e00031-18 (2018).
- Swenson, T. L., Karaoz, U., Swenson, J. M., Bowen, B. P. & Northen, T. R. Linking soil biology and chemistry in biological soil crust using isolate exometabolomics. *Nat. Commun.* **9**, 19 (2018).
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. & Dean, J. Distributed representations of words and phrases and their compositionality. in *Advances in Neural Information Processing Systems* 3111–3119 (2013).
- Aitchison, J. The statistical analysis of compositional data. *J. R. Stat. Soc. B* **44**, 139–160 (1982).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2021

## Methods

The simulations were created by using the generative form of MMvec; the microbe and metabolite factor loadings were randomly generated from a normal distribution to parameterize the MMvec parameters. Microbial counts were then drawn from a multinomial logistic normal distribution and fed into MMvec to generate the metabolite counts. To identify scenarios where CLR correlations underperformed in comparison to MMvec, we used Bayesian Optimization to tune the distributions used to generate the simulations.

The CLR-transformed correlations suggested by Quinn and Erb were benchmarked on the desert biocrust soils dataset using the R scripts provided in ref. <sup>1</sup>.

**Reporting Summary.** Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

The datasets to reproduce the results presented here can be found at <https://github.com/knightlab-analyses/multiomic-cooccurrences>.

## Code availability

The analysis software to reproduce the results presented here can be found at <https://github.com/knightlab-analyses/multiomic-cooccurrences>.

## Author contributions

J.T.M. performed all analyses and wrote the manuscript. All authors have contributed edits to the manuscript.

## Competing interests

M.W. is the founder of Omata Labs. The remaining authors declare no competing interests.

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41592-020-01007-0>.

**Correspondence and requests for materials** should be addressed to R.K.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- |                                     |   |
|-------------------------------------|---|
| n/a                                 | Confirmed   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly  |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided<br><i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i>  |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A description of all covariates tested   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. $F$ , $t$ , $r$ ) with confidence intervals, effect sizes, degrees of freedom and $P$ value noted<br><i>Give <math>P</math> values as exact values whenever suitable.</i>                 |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings   |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes   |
| <input type="checkbox"/>            | <input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's $d$ , Pearson's $r$ ), indicating how they were calculated  |

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

### Software and code

Policy information about [availability of computer code](#)

Data collection

Only simulation data was used.

Data analysis

All data analysis scripts can be found here: <https://github.com/knightlab-analyses/multiomic-cooccurrences>

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

### Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The biocrust soils data was retrieved from the supplemental section in Swenson et al

### Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

- ☒ Life sciences      ☐ Behavioural & social sciences      ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	We used 500 samples in our simulations, since this is larger than most of the studies that we have analyzed. For the biocrust soils study, there were 19 samples and after filtering there were 466 unique microbial taxa and 85 metabolite features.
Data exclusions	Taxa that appeared in less than 10 samples for each study were removed, since there are fewer samples than degrees of freedom in the model to infer these microbes co-occurrence patterns.
Replication	Replication was not necessary, since the data was simulated, not collected.
Randomization	Randomization was not necessary, since the data was simulated, not collected.
Blinding	Blinding was not necessary, since the data was simulated, not collected.

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

### Materials & experimental systems

### Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging