# Interpreting UniFrac with Absolute Abundance: A Conceptual and Practical Guide

Augustus Pendleton[1]* & Marian L. Schmidt[1]*

[1]Department of Microbiology, Cornell University, 123 Wing Dr, Ithaca, NY 14850, USA

**Corresponding Authors:** Augustus Pendleton: arp277@cornell.edu; Marian L. Schmidt: marschmi@cornell.edu

**Author Contribution Statement:** Both authors contributed equally to the manuscript.

**Data Availability:** All data and code used to produce the manuscript are available at https://github.com/MarschmiLab/Pendleton_2025_Absolute_Unifrac_Paper, in addition to a reproducible `renv` environment. All packages used for analysis are listed in Table S1.

Microbial ecologists routinely compare communities using -diversity metrics derived from relative abundances. Yet this approach overlooks a critical ecological dimension: microbial load. Across ecosystems, ranging from gut to soil to water, total microbial abundance can vary by orders of magnitude, even among healthy systems. Relative data can misrepresent true shifts in biomass, making taxa appear to change when they have not. High-throughput sequencing produces compositional data, in which each taxon's abundance is constrained by all others [1]. In low-biomass samples, this distortion is especially pronounced, allowing contaminants to appear biologically meaningful when absolute counts are unknown.

Recognition of these limitations has reshaped how we interpret microbial variation. Quantitative profiling studies show that cell abundance, not only composition, can drive major community differences and alter co-occurrence networks [2]. To overcome compositional constraints, researchers increasingly use flow cytometry, qPCR, and genomic spike-ins to quantify microbial load [3, 4]. These tools improve detection of functionally relevant taxa and mitigate the compositional constraints imposed by sequencing [1, 2]. Most studies using absolute data rely on Bray-Curtis dissimilarity, which does not require normalization to proportions [e.g. 3, 5]. Weighted UniFrac remains the dominant phylogenetic -diversity metric, yet its behavior under absolute abundance frameworks remains unclear. Here, we present *Absolute UniFrac*, a direct extension of Weighted UniFrac that incorporates total abundance, and evaluate its impact across simulated and real-world datasets.

The UniFrac distance was first introduced by Lozupone & Knight (2005), and has since become enormously popular as a measure of $\beta$-diversity within the field of microbial ecology [6]. A benefit of the UniFrac distance is that is considers phylogenetic information when estimating the distance between two communities. After first generating a phylogenetic tree representing species (or amplicon sequence variants, "ASVs") from all samples, the UniFrac distance computes the fraction of branch-lengths which is *shared* between communities, relative to the total branch length represented in the tree. UniFrac can be both unweighted, in which only the incidence of species is considered, or weighted, wherein a branch's contribution is weighted by the proportional abundance of taxa on that branch [7]. The weighted UniFrac is derived:

$$U^R = \frac{\sum_{i=1}^{n} b_i |p_i^a - p_i^b|}{\sum_{i=1}^{n} b_i (p_i^a + p_i^b)}$$

where we weight the length of each branch, $b_i$, by the difference in the relative abundance of all species ($p_i$) descended from that branch in sample *a* or sample *b*. Here, we denote this distance as $U^R$, for "Relative Unifrac". Popular packages which calculate weighted Unifrac-including `diversity-lib` QIIME plug-in and the R packages `phyloseq` and `GUniFrac`-run this normalization by default.

Because $U^R$ is most sensitive to changes in abundant lineages, it can sometimes obscure compositional differences driven by rare to moderately-abundant taxa [8]. To address this weakness, Chen et al. (2012) introduced the generalized UniFrac distance ($GU^R$), in which the impact of abundant lineages can be mitigated by decreasing the

parameter $\alpha$:

$$GU^R = \frac{\sum\limits_{i=1}^{n} b_i(p_i^a + p_i^b)^\alpha \left| \frac{p_i^a - p_i^b}{p_i^a + p_i^b} \right|}{\sum\limits_{i=1}^{n} b_i(p_i^a + p_i^b)^\alpha}$$

where $\alpha$ ranges from 0 (close to unweighted UniFrac) up to 1 (identical to $U^R$, above).
However, if one wishes to use absolute abundances, both $U^R$ and $GU^R$ can be derived
without normalizing to proportions:

$$U^A = \frac{\sum\limits_{i=1}^{n} b_i|c_i^a - c_i^b|}{\sum\limits_{i=1}^{n} b_i(c_i^a + c_i^b)} \qquad GU^A = \frac{\sum\limits_{i=1}^{n} b_i(c_i^a + c_i^b)^\alpha \left| \frac{c_i^a - c_i^b}{c_i^a + c_i^b} \right|}{\sum\limits_{i=1}^{n} b_i(c_i^a + c_i^b)^\alpha}$$

Where $c_i^a$ and $c_i^b$ stands for the absolute counts of species descended from branch $b_i$ in
community $a$ and $b$, respectively. We refer to these distances as "Absolute Unifrac" and
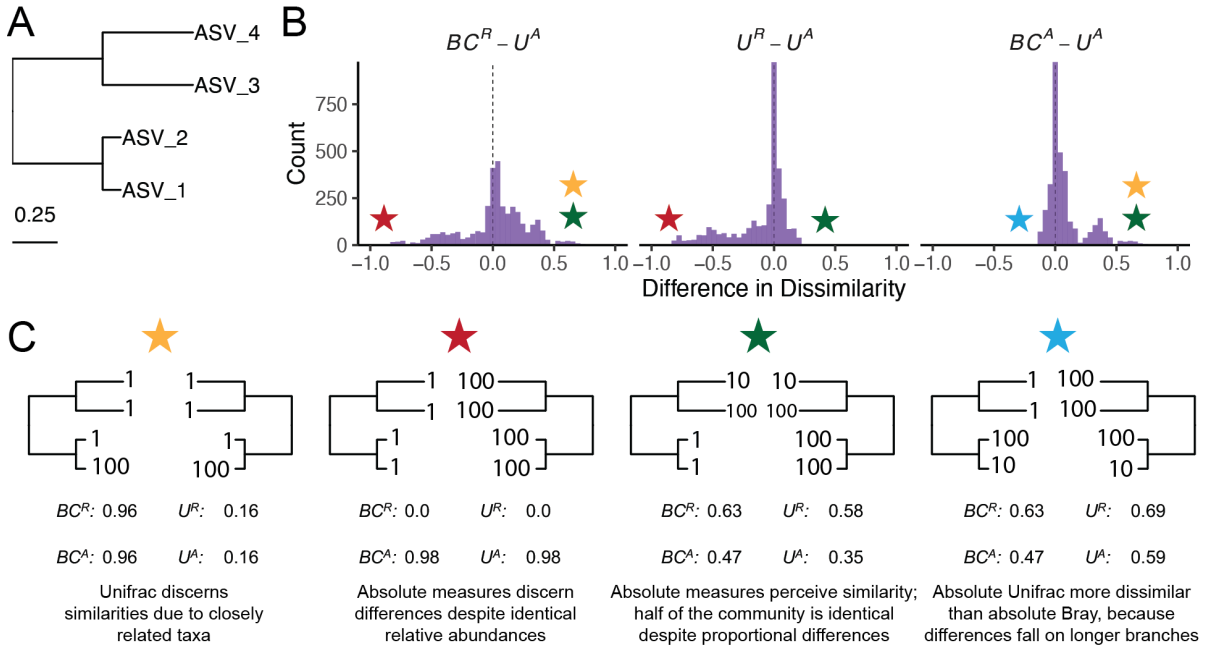"Generalized Absolute Unifrac" ($U^A$ and $GU^A$) .

To illustrate how $U^A$ behaves, we constructed a simulated community of four ASVs
arranged in a simple, balanced phylogeny (Fig. 1A). By varying the absolute abundance
of each ASV (1, 10, or 100), we generated 81 samples and 3,240 pairwise comparisons. For
each pair, we computed four dissimilarity metrics: Bray-Curtis with relative abundance
($BC^R$), Bray-Curtis with absolute abundance ($BC^A$), Weighted UniFrac with relative
abundance ($U^R$), and Weighted UniFrac with absolute abundance ($U^A$).

$U^A$ does not consistently yield higher or lower distances but instead varies depending
on how abundance and phylogeny intersect (Fig. 1B). In the improbable scenario that
all branch lengths are equal, $U^A$ is always less than or equal to $BC^A$ (Fig. S1), but
real-world trees rarely behave this way. These comparisons emphasize that incorporating
phylogeny and absolute abundance reshapes distance estimates in nontrivial ways.

To better understand how these metrics diverge, we examined individual sample pairs
(Fig. 1C). Scenario 1 (gold star) illustrates the classic advantage of UniFrac: ASV_1
and ASV_2 are phylogenetically close, so $U^R$ and $U^A$ discern greater similarity between
samples than $BC^R$ and $BC^A$, which ignore phylogenetic structure. Scenario 2 (red star)
highlights a limitation of relative metrics: two samples with identical relative composition
but a 100-fold difference in biomass appear identical to $BC^R$ and $U^R$, but not to their
absolute counterparts. In Scenario 3 (green star), incorporating absolute abundance
decreases dissimilarity. $BC^A$ and $U^A$ are lower than their relative counterparts because
half the community is identical in absolute terms, despite proportional differences. In
contrast, Scenario 4 (blue star) shows that $U^A$ can increase dissimilarity relative to $BC^A$
when abundance differences occur on long branches, amplifying phylogenetic dissimilarity.

Across all 3,240 pairwise comparisons, $U^A$ is usually smaller than and strongly corre-
lated with $BC^A$ (Pearsons $r = 0.82$, $p < 0.0001$) than with $BC^R$ ($r = 0.41$) and $U^R$ ($r = 0.55$), reflecting the effect illustrated in Scenario 1. However, exceptions like Scenario 4
show that $U^A$ can also yield larger distances than $BC^A$ when abundance differences occur
on long branches. These scenarios demonstrate that $U^A$ integrates ecological realism by

capturing differences in both lineage identity and total biomass, offering a nuanced view of community structure grounded in biological context, while remaining sensitive to long branches.



*Figure 1. Simple simulations to compare the impact of absolute abundance on phylogenetic and non-phylogenetic distance measures.* (A) We created a small community represented by 4 ASVs. We explored all permutations of each ASV holding an absolute count of 1, 10, or 100 to create 81 communities, producing 3240 comparisons for which we calculated distance metrics. (B) Distributions of the difference between weighted Unifrac using absolute abundance ($U^A$), Bray-Curtis using relative abundance ($BC^R$), weighted Unifrac using relative abundance ($U^R$), and Bray-Curtis using absolute abundance ($BC^A$). (C) Scenarios (comparisons) highlighting occasions when $U^A$ is larger or lesser compared to other metrics. The stars indicate what region of the distributions in (B) this scenario would fall. The actual values for each metric are displayed beneath the scenario.

We next explored the practical effect $U^A$ has on separating groups within a real dataset. We used a previously published 16S dataset from Lake Ontario, represented by 66 samples and >7,000 ASVs, where samples clustered into three main groups defined by depth and month due to changes in both composition and abundance (Fig. S2, [9]). We sought to demonstrate how weighting by absolute abundance can affect the interpretation and statistical power to separate samples between these depth/month groups.
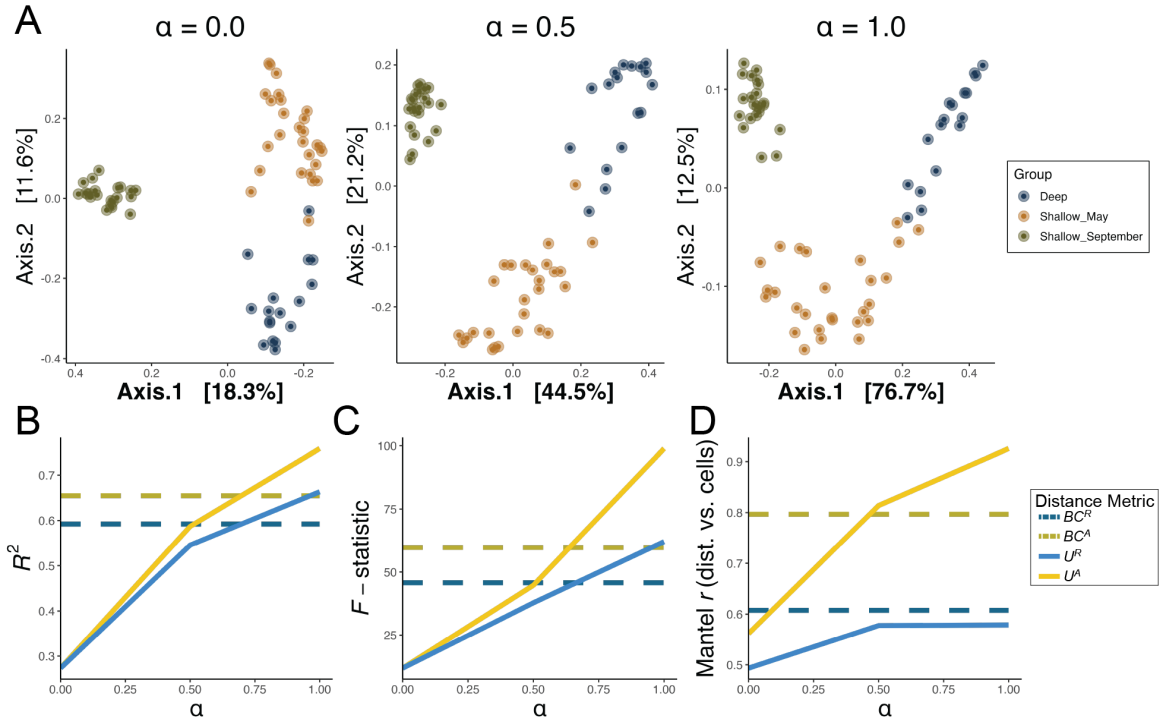
We calculated the generalized absolute UniFrac ($GU^A$) across three levels of $\alpha$: 0.0 (close to unweighted UniFrac), 0.5, and 1.0 (identical to $U^A$). Fig. 2A shows PCoA ordinations of the same samples across this $\alpha$-range. As $\alpha$ increases, we see a greater recognition of similarity between Shallow September and Shallow May, reflecting their relatively higher cell counts compared to the Deep samples. Note that an increasing amount of variation is also assigned to the first axis, going from 18.3% up to 76.7%. This trend was also true for $U^R$ across multiple $\alpha$, but to a much weaker degree (Fig. S3).

We used PERMANOVAs to quantify the variance ($R^2$) and statistical power ($F$-statistic) our three depth-month groups can explain across different distance metrics ($GU^R$, $GU^A$, $BC^R$ and $BC^A$) and different alpha values (Fig. 2B and 2C). Absolute abundance metrics consistently allow for greater $R^2$ values ($GU^A$ reaching a maximum of 75.8%) and F-statistics (1.56X higher for $GU^A$ than $GU^R$) than their relative counter-

parts. For UniFrac distances, this separation increases as alpha increases. As such, the incorporation of absolute abundances can be a powerful tool to emphasize differences in microbial load between groups.

That said, we identify a major caveat; $U^A$ can become strongly correlated to the cell-counts alone compared to other distance metrics (Fig. 2D). We used a Mantel test to assess the correlation between each distance metric to the absolute difference in cell counts between each sample. Absolute abundance measures are more sensitive to differences in cell counts compared to their relative counterparts. This is intuitive, and to a degree, intentional; the utility of these metrics is to identify changes in absolute abundance even in the case of compositional similarity. However, at $\alpha = 1$, the value of $U^A$ becomes almost entirely based on the overall abundance of a sample. We see this reflected in the ordinations as well; as $\alpha$ increases, the first Axis increasingly represents and correlates to the absolute abundance (Spearman's $\rho$ = -1.0 for $\alpha = 1.0$ versus Spearman's $\rho = 0.58$ for $\alpha = 0.0$ versus ). It is possible that any 2D structure we observed in the third panel of Fig. 2A is simply the result of the horseshore effect [10].

With this in mind, we encourage others to think critically about how heavily they want their distance metric to correspond to absolute abundance, and modulate this effect using $GU^A$ rather than $U^A$. As a simple start from these data, we recommended using an $\alpha$ value of 0.5, consistent with the original recommendation of Chen et al. (2012) but of heightened importance when using absolute abundances.



*Figure 2. Impact of absolute abundance on separating freshwater microbial communities* (A) Samples represent microbial communities from Lake Ontario taken at different times (May and September) and at different depths (Shallow and Deep), as originally described in [9]. We used Principal Coordinates Analysis (PCoA) to ordinate samples based on the $GU^A$ at $\alpha$ values of 0.0, 0.5, and 1.0. The proportion of variance assigned to each axis is provided in brackets. Note that the x-axis is reversed in the first panel, to provide visual symmetry across ordinations. We then used PERMANOVAs to quantify the variance ($R^2$) (B) and statistical power ($F$-statistic) (C) our three depth-month groups can explain across different distance metrics ($GU^R$, $GU^A$, $BC^R$ and $BC^A$) and different alpha values. (D) Correlations between

151 different distance metrics and differences in over-all cell abundances as assessed by a Mantel test.

152 The incorporation of absolute abundance allows microbial ecologists to assess more
153 realistic, ecologically-relevant differences in microbial communities in situations for where
154 microbial load matters. For example, the temporal development of the infant microbiome
155 is captured in both a rise in absolute abundance in addition to compositional changes
156 [5]; bacteriophage predation in wastewater bioreactors can be understood only when
157 incorporating microbial load [11]; and the negative impact of antibiotics on the abundance
158 of specific taxa in the swine gut was missed using relative abundance approaches [12].
159 As $\beta$-diversity metrics are an essential tool for every microbial ecologist-and UniFrac
160 distances are highly popular within microbiome research specifically-we encourage other
161 researchers to adopt the usage of $GU^A$ in comparing their samples.

162 Our data show, however, that the interpretation of these metrics can be complex, and
163 researchers should consider (and justify) the relative importance of absolute abundance
164 in their estimates of $\beta$-diversity. We also do not address other concerns relevant to
165 estimating $\beta$-diversity, including the impacts of sampling effort on richness and whether
166 rarefaction should be applied when calculating $GU^A$ [13]. As ever, interpreting the results
167 of these metrics critically and exploring the sensitivity of your conclusions across multiple
168 metrics is encouraged [14]. Through this work, we hope to facilitate the adoption of $GU^A$
169 and its interpretation in an approachable manner.

## References

171 1. Gloor GB et al. Microbiome Datasets Are Compositional: And This Is Not Optional. *Front Microbiol* 2017;**8**:2224. https://doi.org/10.3389/fmicb.2017.02224

172 2. Vandeputte D et al. Quantitative microbiome profiling links gut community variation to microbial load. *Nature* 2017;**551**:507–511. https://doi.org/10.1038/nature24460

173 3. Props R et al. Absolute quantification of microbial taxon abundances. *ISME J* 2017;**11**:584–587. https://doi.org/10.1038/ismej.2016.117

174 4. Wang X et al. Current Applications of Absolute Bacterial Quantification in Microbiome Studies and Decision-Making Regarding Different Biological Questions. *Microorganisms* 2021;**9**:1797. https://doi.org/10.3390/microorganisms9091797

175 5. Rao C et al. Multi-kingdom ecological drivers of microbiota assembly in preterm infants. *Nature* 2021;**591**:633–638. https://doi.org/10.1038/s41586-021-03241-8

176 6. Lozupone C, Knight R. UniFrac: a New Phylogenetic Method for Comparing Microbial Communities. *Applied and Environmental Microbiology* 2005;**71**:8228–8235. https://doi.org/10.1128/AEM.71.12.8228-8235.2005

177 7. Lozupone CA et al. Quantitative and Qualitative Diversity Measures Lead to Different Insights into Factors That Structure Microbial Communities. *Applied and Environmental Microbiology* 2007;**73**:1576–1585. https://doi.org/10.1128/AEM.01996-06

178 8. Chen J et al. Associating microbiome composition with environmental covariates using generalized UniFrac distances. *Bioinformatics* 2012;**28**:2106–2113. https://doi.org/10.1093/bioinformatics/bts342

179 9. Pendleton A, Wells M, Schmidt ML. Upwelling periodically disturbs the ecological assembly of microbial communities in the Laurentian Great Lakes.

180 10. Morton JT et al. Uncovering the Horseshoe Effect in Microbial Analyses. *mSystems* 2017;**2**:10.1128/msystems.00166–16. https://doi.org/10.1128/msystems.00166-16

181 11. Shapiro OH, Kushmaro A, Brenner A. Bacteriophage predation regulates microbial abundance and diversity in a full-scale bioreactor treating industrial wastewater. *The ISME Journal* 2010;**4**:327–336. https://doi.org/10.1038/ismej.2009.118

182 12. Wagner S et al. Absolute abundance calculation enhances the significance of microbiome data in antibiotic treatment studies. *Front Microbiol* 2025;**16**. https://doi.org/10.3389/fmicb.2025.1481197

183 13. Schloss PD. Waste not, want not: revisiting the analysis that called into question the practice of rarefaction. *mSphere* 2023;**9**:e00355–23. https://doi.org/10.1128/msphere.00355-23

184 14. Kers JG, Saccenti E. The Power of Microbiome Studies: Some Considerations on Which Alpha and Beta Metrics to Use and How to Report Results. *Front Microbiol* 2022;**12**. https://doi.org/10.3389/fmicb.2021.796025