# HTSSIP: an R package for analysis of high throughput sequencing data from nucleic acid stable isotope probing (SIP) experiments

Nicholas D. Youngblut*[1], Samuel E. Barnett[2], and Daniel H. Buckley[2]

[1] Max Planck Institute for Developmental Biology, Spemannstraße 35, 72076 Tübingen, Germany
[2] School of Integrative Plant Science, Cornell University, Ithaca, NY 14843, USA

*Corresponding author:
    Nicholas D. Youngblut
    Max Planck Institute for Developmental Biology
    Spemannstraße 35, 72076
    Tübingen, Germany
    nyoungblut@tuebingen.mpg.de

Running title: *HTSSIP for analysis of HTS-SIP data*

1 **Abstract**

2       Combining high throughput sequencing with stable isotope probing (HTS-SIP) is a

3 powerful method for mapping *in situ* metabolic processes to thousands of microbial taxa.

4 However, accurately mapping metabolic processes to taxa is complex and challenging. Multiple

5 HTS-SIP data analysis methods have been developed, including high-resolution stable isotope

6 probing (HR-SIP), multi-window high-resolution stable isotope probing (MW-HR-SIP),

7 quantitative stable isotope probing (q-SIP), and ΔBD. Currently, the computational tools to

8 perform these analyses are either not publicly available or lack documentation, testing, and

9 developer support. To address this shortfall, we have developed the *HTSSIP* R package, a

10 toolset for conducting HTS-SIP analyses in a straightforward and easily reproducible manner.

11 The *HTSSIP* package, along with full documentation and examples, is available from CRAN at

12 https://cran.r-project.org/web/packages/HTSSIP/index.html and Github at

13 https://github.com/nick-youngblut/HTSSIP.

14

15 **Introduction**

16       Stable isotope probing of nucleic acids (DNA- and RNA-SIP) is a powerful method for

17 mapping *in situ* metabolic processes, such as nitrogen and carbon cycling, to microbial taxa.

18 Historically the sensitivity of nucleic acid SIP has been limited by the low throughput of DNA

19 sequencing and the low taxonomic resolution of DNA fingerprinting techniques [1,2]. Recently,

20 DNA- and RNA-SIP have been combined with high throughput sequencing of PCR amplicons

21 (HTS-SIP), which allows researchers to map *in situ* metabolic processes to thousands of taxa

22 resolved at a fine taxonomic resolution [3–5].

23       While HTS-SIP is proving to be a very useful method for exploring *in situ* metabolic

24 processes in complex microbial communities, the accurate analysis of HTS-SIP datasets is

25 complex [6,7]. Multiple strategies have been developed for analyzing HTS-SIP data, including

26    high-resolution stable isotope probing (HR-SIP) [5], multi-window high-resolution stable isotope

27    probing (MW-HR-SIP) [7], quantitative stable isotope probing (q-SIP) [3], and ΔBD [5]. The

28    goals of these methods differ, with HR-SIP and MW-HR-SIP designed to accurately identify taxa

29    that have incorporated isotopically labeled substrate (*i.e.* 'incorporators'), while the main goal of

30    q-SIP and ΔBD is to quantify the amount of isotopic enrichment for each taxon (*i.e.* atom %

31    excess). While all methods use amplicon sequences (*e.g.* 16S rRNA or fungal ITS sequences)

32    from multiple fractions of each isopycnic gradient, HR-SIP, MW-HR-SIP, and ΔBD solely use

33    sequence data while q-SIP additionally requires qPCR derived estimations of gene copy

34    number from each gradient fraction. Recently, Youngblut and Buckley developed a HTS-SIP

35    simulation model and showed that MW-HR-SIP is more accurate for identifying incorporators

36    than HR-SIP and q-SIP, while q-SIP is generally more precise than ΔBD for quantifying isotopic

37    enrichment [7].

38         The code for performing each of these HTS-SIP analyses is limited in availability,

39    documentation, and developer support; all of which severely limit the ease of use and

40    reproducibility of HTS-SIP analyses. To address this deficiency, we developed the *HTSSIP* R

41    package, which includes the following features:

42    ● Functions for conducting HR-SIP, MW-HR-SIP, q-SIP, and ΔBD to analyze data from

43       DNA-SIP and RNA-SIP experiments

44    ● Functions for performing HTS-SIP dataset simulation, as described [7]

45    ● Functions for exploratory analysis of simulated HTS-SIP data, useful for predicting how

46       different experimental designs can alter experimental outcomes

47    ● Functions for exploratory analysis of real HTS-SIP data, useful for conducting post-hoc

48       analyses

49    ● Ability to run analyses with parallel processing

50    ● Extensive documentation and tutorials (see the *HTSSIP* vignettes)

51

52    **Package description**

53    *Input data*

54    Dataset input is handled by the *Phyloseq* R package, a feature-rich package for general

55    microbiome data analysis that can be used to import many common microbiome data formats

56    [8]. *HTSSIP* includes convenience functions to easily and flexibly designate the experimental

57    design of the SIP experiment for downstream HTS-SIP analyses (Figure 1).

58

59    **Figure 1.** *A diagram depicting the possible analyses available in the* HTSSIP *R package.* The R functions to conduct

60    each workflow step are italicized, and the figure references refer to example data produced by these workflow steps.

61    *HTS-SIP dataset exploratory analyses*

62    A common first step in analyzing nucleic acid SIP data is to quantify the total nucleic acid

63    concentration or gene copy number (estimate by qPCR) across density gradients in order to

64    determine the buoyant density (BD) "shift" of nucleic acids in isotopically labeled treatments

65    versus unlabeled controls [9,10]. The general expectation is that a "shift" of nucleic acid BD from

66    "light" towards "heavy" densities is indicative of isotope incorporation. However, in a well

67    designed SIP experiment, the ratio of exogenous to indigenous substrate should be small, and

68    this can produce an imperceptible BD shift [4]. In addition, an extensive shift may indicate

69    excessive cross-feeding [11]. HTS-SIP methods can detect taxa that have incorporated low

70    levels of isotope, or occur at frequencies that are so low that they do not cause a shift in the

71    overall BD of community nucleic acids [5]. As a result, analysis of the BD distribution of total

72    nucleic acids within density gradients is of little utility in assessing the results of nucleic acid SIP

73    experiments performed on complex communities.

74    As a simpler alternative, which leverages the power of high-throughput sequencing

75    techniques, BD "shifts" can be inferred solely from sequence data [4,5]. Given that incorporators

76    will be more abundant in "heavy" gradient fractions of the labeled treatment versus the

77    unlabeled control, a BD shift can be inferred by assessing the beta-diversity between treatment

78    and control gradient fractions. This approach is more sensitive for detecting community-level

79    isotope incorporation than the approach of quantifying total nucleic acid concentration across

80    the density gradient [7]. *HTSSIP* implements two methods for using beta-diversity to assess

81    isotope incorporation at the community-level: an ordination approach and an approach that

82    expresses beta-diversity between corresponding treatment and control fractions as a function of

83    their BD (Figure 1).

84        The ordination approach simply involves pairwise calculations of a beta-diversity metric

85    between all gradient fractions from isotopically labeled treatments and corresponding unlabeled

86    controls, followed by visualizing the distance matrix with either principal coordinates analysis

87    (PCoA) or non-metric multidimensional scaling (NMDS). An increase in beta-diversity between

88    corresponding gradient fractions of labeled samples and controls is expected if isotope

89    incorporation causes a change in the BD of OTUs (Figure 2A & 2B).

90

91    **Figure 2.** *Examples of the ordination and BD-shift analyses for assessing community-level incorporation.* Plots A and

92    B are non-metric multidimensional scaling (NMDS) ordinations of beta-diversity (16S rRNA OTUs; 97% sequence

93    identity; weighted Unifrac) calculated between gradient fractions from a HTS-DNA-SIP experiment conducted with

94    agricultural soil. Plot A compares fractions from replicate unlabeled control gradients, with different symbols (circles

95    and triangles) used to distinguish different replicates, and with symbol diameter scaled in relation to fraction buoyant

96    density as indicated in the accompanying scale. Plot B compares fractions from labeled treatments ("13C-Cel" for

97    [13]C-cellulose  or "13C-Xyl" for [13]C-xylose) versus their corresponding unlabeled controls ("12C-Con") at 3 or 14 days

98    after substrate addition ("D03" and "D14", respectively). The NMDS stress values ranged from 0.06 to 0.07. An

99    increase in beta-diversity is expected between labeled and unlabeled "heavy" fractions in response to isotope

100    incorporation. Plots C and D depict the same data as in Plots A and B, but the beta-diversity comparisons between

101    labeled treatment and unlabeled control are indicated only for fractions that correspond in BD. To account for partial

102    overlap between labeled and unlabeled fractions, the weighted mean beta-diversity value is calculated based on

103    percent overlap in BD ranges. "BD shift windows" indicate regions defined by ≥3 consecutive fractions with

104    significantly high beta-diversity resulting from isotope incorporation, with significance defined by permuting OTU

105    abundances and recalculating beta-diversity values (100 bootstrap replicates; $P < 0.05$). The dataset used is a subset

106    from the dataset from Youngblut and Buckley [7].

107

108        While the ordination approach provides a useful overview of community-wide isotope

109    incorporation, the extent of incorporation is difficult to compare among multiple treatments (*e.g.*

110    $^{13}$C-cellulose versus $^{13}$C-xylose). The second approach implemented in *HTSSIP* visualizes DNA

111    BD shifts by calculating pairwise beta-diversity of corresponding gradient fractions between

112    treatment and control gradients. To deal with partially overlapping gradient fractions between

113    gradients, the weighted mean beta-diversity is calculated from all treatment gradient fractions

114    that overlap each control gradient fraction, with weights defined as the percent overlap in the BD

115    range of each fraction (Figure 2C & 2D). A permutation test is used to identify BD ranges of high

116    beta-diversity resulting from BD shifts ("BD shift windows"). The permutation test involves

117    constructing bootstrap confidence intervals (CI) of beta-diversity by permuting OTU abundances

118    among labeled treatments (*i.e.* a null model where OTUs in treatment are randomly dispersed

119    relative to the control). A note in interpreting these data is that isotope incorporation will cause

120    DNA to shift out of "light" gradient fractions and into "heavy" gradient fractions. Hence, in the

121    presence of isotope incorporation, high beta-diversity can be observed in both "heavy" and

122    "light" gradient fractions. Alternatively, in the absence of isotope incorporation, beta-diversity will

123    remain low across all gradient fractions.

124    *Identifying incorporators*

125        HR-SIP, MW-HR-SIP, and q-SIP can all be used to identify incorporators. To illustrate

126    the application of HR-SIP, MW-HR-SIP, and q-SIP in the *HTSSIP* R package, we simulated a

127    simplified HTS-SIP dataset consisting of 10 OTUs (Figure 3A). Our purpose here is merely to

128    illustrate functions of the *HTSSIP* R package; comprehensive assessment of the accuracy of

129    these techniques is available elsewhere [7]. Briefly, HR-SIP identifies incorporators by utilizing

130    DESeq2 to identify OTUs that have high differential relative abundance in "heavy" fractions of

131    labeled treatment versus unlabeled control [12]. MW-HR-SIP takes the same relative

132    abundance based approach as HR-SIP but uses multiple overlapping "heavy" BD windows

133    (while correcting for multiple hypotheses). In contrast, q-SIP uses qPCR data to transform OTU

134    relative abundance distributions into pseudo-absolute abundance distributions (Figure 3A), and

135    then BD shifts are determined from these transformed distributions by calculating the difference

136    in center of mass for each OTU in treatment versus control gradients. Atom fraction excess can

137    thus be calculated for specific isotopes (*e.g.* $^{13}$C or $^{15}$N) based on the calculations described in

138    the work of Hungate and colleagues [3]. In order to identify incorporators, a permutation test is

139    used to construct bootstrap confidence intervals of atom fraction excess. Sensitivity in

140    identifying incorporators can depend on the methods used (Figure 3B; and see [7]). SIP

141    experiments can be simulated using the SIPSim toolset [7], and these data analyzed using the

142    *HTS-SIP* R package. Such *in silico* evaluation is valuable for predicting possible experimental

143    outcomes and the expected analytical accuracy of SIP experiments based on details of

144    experimental design prior to conducting experiments.

145        *HTSSIP* implements HR-SIP based on the code provided in the work of Pepe-Ranney

146    and colleagues [5]. MW-HR-SIP is implemented in *HTSSIP* based on the R code provided in the

147    SIPSim HTS-SIP dataset simulation toolset [7]. The *HTSSIP* implementation of q-SIP is based

148    on the method's description in the work of Hungate and colleagues [3]. Implementations of each

149    method include the option for parallel processing of each algorithm. Parallelization is

150    implemented through the *plyr* R package [13], which allows for various parallel backends to be

151    used such as *doSNOW* and *doParallel*.

152 *Quantifying isotopic enrichment*

153       Unlike HR-SIP and MW-HR-SIP, the main goal of q-SIP and ΔBD is to quantify isotopic

154    enrichment. To illustrate the use of *HTSSIP* for conducting q-SIP and ΔBD, we applied both

155    analyses to the simplified HTS-SIP dataset described above (Figure 3C). ΔBD is implemented

156    in *HTSSIP* as described in the work of Pepe-Ranney and colleagues [5]. As shown in Youngblut

157    and Buckley [7], q-SIP and ΔBD can produce substantially different estimates of isotope

158    incorporation.

159

160    **Figure 3.** *Examples of using the HTSSIP R package for data processing, data exploration, incorporator identification,*

161    *and quantification of BD shifts (Z).* The SIPSim toolset was used to simulate the relative abundances of 10 OTUs

162    across 24 gradient fractions in an experiment that includes a single $^{13}$C-treatment ("13C-Treat") and a single $^{12}$C-

163    control ("12C-Con") condition, each with 3 experimental replicates. Half of the OTUs had an atom fraction excess of

164    30 to 100%, while the others were 0%. qPCR estimates of total community 16S rRNA copy numbers were also

165    simulated with SIPSim, and qPCR analytical error was modeled based on error estimated from Hungate and

166    colleagues [3]. Plot A depicts the raw abundances ("Counts"), fractional relative abundance ("Rel. Abund."), and

167    relative abundances transformed by simulated qPCR data ("Rel. Abund. qPCR-trans."). For clarity, only 1 of the 3

168    experimental replicates is shown. Plot B shows which OTUs were identified as incorporators by the statistical

169    methods described for HR-SIP, MW-HR-SIP, or q-SIP. A Benjamini-Hochberg corrected p-value cutoff of 0.1 was

170    used for HR-SIP and MW-HR-SIP, and 100 bootstrap replicates were used to calculate confidence intervals for q-

171    SIP, as described [3]. Plot C shows the mean BD shift of each OTU as quantified by ΔBD or q-SIP. The dashed line

172    signifies a BD shift (Z) of 0.0 g ml$^{-1}$, and the red bars show the true theoretical BD shift resulting from $^{13}$C isotope

173    incorporation.

174 *Simulating datasets*

175       *HTSSIP* provides functions to simulate simple HTS-SIP datasets for use in software

176    testing, analysis pipeline development, and gaining familiarity with software and data formats.

177    However, the SIPSim toolset is recommended for evaluating possible SIP experimental designs

178    and for testing the accuracy of HTS-SIP analyses, because the simulation framework for

179 SIPSim is based on the physics of isopycnic centrifugation, unlike the simulations possible with

180 *HTSSIP* [7]. *HTSSIP* utilizes *coenocliner*, an R package designed for simulating taxon

181 abundance across environmental gradients, to simulate taxon abundance distributions across

182 buoyant density gradient fractions [13].

183

184 **Availability**

185 The *HTSSIP* package and the data used in this work are available from CRAN at

186 https://cran.r-project.org/web/packages/HTSSIP/index.html and Github at

187 https://github.com/nick-youngblut/HTSSIP.

188

189 **Future work**

190 Future development of the *HTSSIP* package will include i) functions for mapping

191 incorporator status to phylogenies and visualizing the results ii) direct integration with the

192 SIPSim toolset for rapid HTS-SIP experimental design and assessment of accuracy iii) functions

193 analyzing shotgun metagenome data derived from SIP experiments.

194

195 **Conclusions**

196 Given the power of HTS-SIP for mapping *in situ* metabolism to taxonomic identity,

197 adoption of the technique by researchers will greatly help to resolve connections between

198 microbial ecology and taxonomy. Currently, HTS-SIP data analysis is complex, with few existing

199 computational tools to aid researchers. The R package *HTSSIP* provides a single, standardized

200 analysis pipeline that facilitates reproducible analyses on HTS-SIP datasets and direct cross-

201    study comparisons. Moreover, *HTSSIP* can be combined with the SIPSim toolset to simulate

202    and evaluate possible DNA-SIP experimental designs.
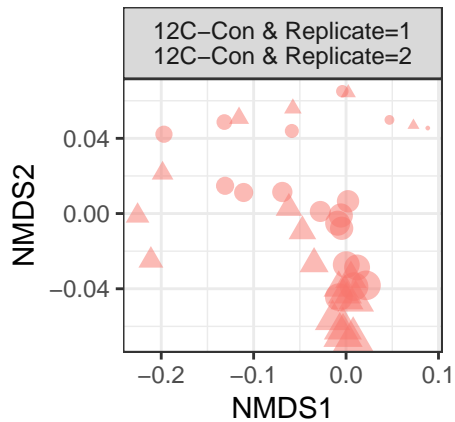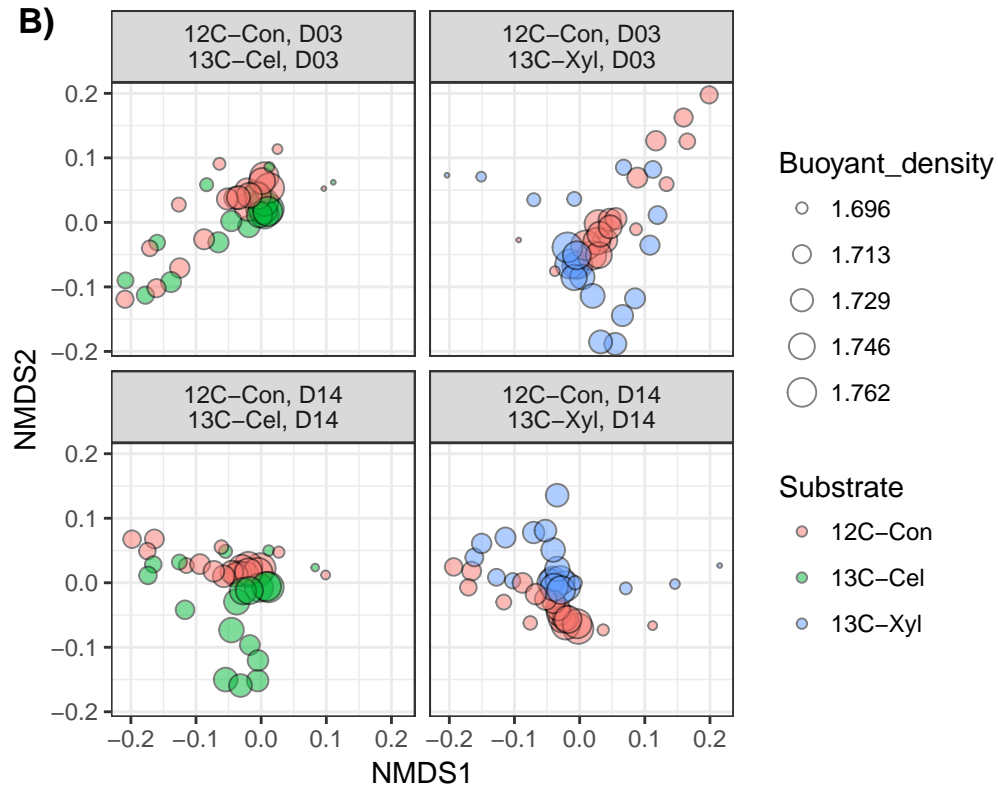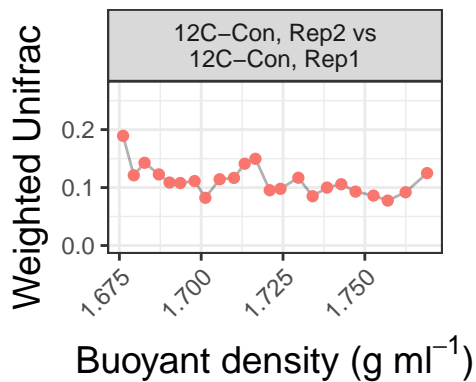
203

**Acknowledgements**

208

**References**

210    1.    Whiteley AS, Thomson B, Lueders T, Manefield M. RNA stable-isotope probing. Nat Protoc.

211        2007;2: 838–844.

212    2.    Uhlík O, Jecná K, Leigh MB, Macková M, Macek T. DNA-based stable isotope probing: a

213        link between community structure and function. Sci Total Environ. 2009;407: 3611–3619.

214    3.    Hungate BA, Mau RL, Schwartz E, Caporaso JG, Dijkstra P, van Gestel N, et al.

215        Quantitative microbial ecology through stable isotope probing. Appl Environ Microbiol.

216        2015;81: 7570–7581.

217    4.    Pepe-Ranney C, Koechli C, Potrafka R, Andam C, Eggleston E, Garcia-Pichel F, et al. Non-

218        cyanobacterial diazotrophs mediate dinitrogen fixation in biological soil crusts during early

219        crust formation. ISME J. 2016;10: 287–298.

220    5.    Pepe-Ranney C, Campbell AN, Koechli CN, Berthrong S, Buckley DH. Unearthing the

221        ecology of soil microorganisms using a high resolution DNA-SIP approach to explore

222        cellulose and xylose metabolism in soil. Front Microbiol. 2016;7: 703.

223    6.  Youngblut ND, Buckley DH. Intra-genomic variation in G + C content and its implications for DNA stable isotope probing. Environ Microbiol Rep. 2014;6: 767–775.

224

225    7.  Youngblut ND, Buckley DH. Evaluating the accuracy of DNA stable isotope probing.

226        bioRxiv. 2017. p. 138719. doi:10.1101/138719

227    8.  McMurdie PJ, Holmes S. phyloseq: An R package for reproducible interactive analysis and

228        graphics of microbiome census data. PLoS ONE. 2013. p. e61217. Available:

229        http://dx.plos.org/10.1371/journal.pone.0061217

230    9.  Lueders T, Kindler R, Miltner A, Friedrich MW, Kaestner M. Identification of bacterial

231        micropredators distinctively active in a soil microbial food web. Appl Environ Microbiol.

232        2006;72: 5342–5348.

233    10. El Zahar Haichar F, Achouak W, Christen R, Heulin T, Marol C, Marais M-F, et al.

234        Identification of cellulolytic bacteria in soil by stable isotope probing. Environ Microbiol.

235        2007;9: 625–634.

236    11. DeRito CM, Pumphrey GM, Madsen EL. Use of field-based stable isotope probing to

237        identify adapted populations and track carbon flow through a phenol-degrading soil

238        microbial community. Appl Environ Microbiol. 2005;71: 7858–7865.

239    12. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-

240        seq data with DESeq2. Genome Biol. 2014;15: 550.

241    13. Wickham H. The split-apply-combine strategy for data analysis. Journal of Statistical

242        Software, Articles. 2011;40: 1–29.

|  | q-SIP | HR-SIP | MW-HR-SIP |
|---|---|---|---|

**Input Data**

Phyloseq object
+ qPCR dataset

Phyloseq object

**Data Modification (Figure 3A)**

Multiply OTU counts by qPCR values.
*OTU_qPCR_trans()*

Parse data into a list based on treatments. Each list entry should contain a phyloseq object containing data for a treatment sample and it's corresponding control sample (*i.e.* $^{13}$C-Cellulose Day 14 and $^{12}$C-Control Day 14). *phyloseq_subset()*

**Community-wide incorporation (Figure 2)**

**Ordination:** ordination of beta-diversity among fractions from labeled treatment versus unlabeled control. *SIP_betaDiv_ord()*

**BD-shift:** ordination of beta-diversity among fractions from labeled treatment versus unlabeled control. BD shift windows defined by permuting OTU abundances. *BD_shift()*

**Identifying incorporators (Figure 3B)**

1. Calculate BD shift and atom fraction excess for each OTU, comparing between treatment and control samples. *qSIP_atom_excess()*
2. Produce atom fraction excess confidence intervals for for each OTU. *qSIP_bootstrap()*

Calculate change in abundance between "heavy" window of treatment and control samples for each OTU. *HRSIP()*

Calculate change in abundance between treatment and control sample using multiple overlapping "heavy" windows for each OTU. *HRSIP()*

**Quantifying incorporation (Figure 3C)**

Calculate the difference in weighted mean BD between labeled treatment and unlabeled control for each OTU. *qSIP_atom_excess()*

ΔBD analysis

Calculate the difference in weighted mean BD between labeled treatment and unlabeled control for each OTU (Note: using untransformed relative abundances). *delta_BD()*