1  # Risk of Bias Assessment Tool for Systematic Review and Metanalysis of the Gut

2  # Microbiome

3  Thomas Lampeter[1*], Charles Love[2], Trien Trey Tang[2], Aditi Marella[2], Hayden Young Lee[2], Armani
4  Oganyan[2], Devin Moffat[2], Anisha Karim[2], Matthew Rusling DO[2*], Aubrey Massmann DO/MPH[2], Melanie
5  Orr[1], Christian Bongiorno[2], Li-Lian Yuan PhD[2*]

6  [1]New York Institute of Technology College of Osteopathic Medicine, 101 Northern Blvd, Glen Head, NY
7  11545, USA

8  [2]Des Moines University College of Osteopathic Medicine, 3200 Grand Ave, Des Moines, IA 50312, USA

9  *corresponding authors: Lampeter.thomasm@gmail.com; matthewrrusling@gmail.com;
10  lilian.yuan@dmu.edu

11

12  **Author Contributions**

23

## Abstract:

Risk of bias assessment is a critical step of any metanalysis or systematic review. Given the low sample count of many microbiome studies, especially observational or cohort studies involving human subjects, many microbiome studies have low power. This increases the importance of performing metanalysis and systematic review for microbiome research in order to enhance the relevance and applicability of microbiome results. This work proposes a method based on the ROBINS-I tool to systematically consider sources of bias in microbiome research seeking to perform metanalysis or systematic review for microbiome studies.

## Introduction:

The most common experimental design used to evaluate the effects of gut microbiome (GMB) genomic or taxonomic post-exposure remodeling has been cohort studies using either animal or human models. Randomized controlled trials (RCTs) for microbiome interventions are less common because we are still characterizing microbiome post-exposure remodeling to identify promising markers or targets for microbiome intervention that would warrant subsequent evaluation by RCTs. Therefore, results from a systematic review with quantitative or pooled metanalysis are essential in identifying candidates for RCTs.

A diligent risk of bias (ROB) assessment is a key step in systematic review or metanalysis to determine the likelihood that features of the study design or conduct of the study will give misleading results. GMB research is highly heterogeneous in its methods, reporting, and attempts to address bias. This manuscript and its associated rubric (**table 1**) are based on the Risk of Bias in Non-randomized Studies - of Interventions (ROBINS-I) tool, and are meant to be used as a GMB-specific adjunct to ROBINS-I. This manuscript and its associated rubric together form a tool that was developed to help standardize ROB assessment in metanalyses and systematic reviews of GMB studies. A small-scale validation test by first-time ROB assessors produced consistently similar ROB determinations, suggesting that this tool can successfully guide consistent ROB determinations. This tool may allow for improved ROB assessment when evaluating studies for metanalyses and systematic reviews of the GMB.

## Using This Tool:

This manuscript and its associated rubric provide a framework for assessing ROB specific to GMB research. This tool strives to provide insight and reduce variability between individual researchers and groups conducting systematic reviews of the GMB. We do not seek to suggest best practices. Instead, we aim to indicate potential sources of bias that may significantly impact GMB studies and are thus vital when considering the strength of evidence for systematic review and metanalysis. The essential criteria in this manuscript are summarized in **table 1,** which was compiled to act as a rubric in guiding ROB determination.

Table 1, "the rubric," guides the determination of low, moderate, or high ROB across seven domains. In each cell of the rubric, there are signaling statements to help guide low, moderate, or high ROB determination in that domain. Two additional ROB determinations are not included on the rubric as they are to be used at the judgement of the person assessing ROB in a study. They are "critical ROB" and "no information". Critical ROB can be determined when a reviewer believes a study to be too problematic to

65  provide useful evidence on the effect of an intervention. As such, a study determined to be of critical
66  ROB in any one domain should not be included in any synthesis. A determination of no information
67  applies to domains where there is no clear evidence of a critical ROB *and* a lack of information to judge
68  ROB otherwise.

69

# 70  1 – Confounding

## 71  1.1 Demographic Differences

72  Important demographic considerations in GMB studies are sex and age. Substantial differences in the
73  gut microbiota are attributable to sex differences in mammals (Org *et al.* 2016, Kim *et al.* 2020). Because
74  of this, any study which includes one sex in one arm and a different sex in another should be classified
75  as having a high risk of bias. In addition to the risk of bias from sex, other demographic factors may also
76  introduce confounding bias into the studies being examined. The GMB changes with age across
77  numerous conditions, disease models, and species impacting microbial diversity and biome composition
78  (Ticinesi *et al.* 2019, Liu *et al.* 2020). Therefore, age differences between cohorts and study arms should
79  be assessed. If the study being examined uses organisms of one age in one arm and a different age in a
80  second arm, it should be classified as having a high risk of bias. The age gap which introduces significant
81  confounding bias, varies by organism. An example of an age gap that would introduce a high risk of bias
82  is 8-week-old mice versus 1-year-old mice (Yoon *et al.* 2021).

83

## 84  1.2 Habitat Stability

85  The habitat in which organisms are kept substantially impacts their GMB (Singh *et al.* 2021). Mice,
86  common subjects of microbiome research, are known to have highly variable microbiomes on arrival at
87  a facility, likely because of transportation stress on the microbiome itself and the immune system and
88  hormonal functions of the host organism (Lipinski *et al.* 2021, Montonye *et al.* 2018, Capdevila *et al.*
89  2007). Studies that do not allow for microbiome stabilization before research begins risk confounding
90  bias due to a lack of habitat stability. Organisms should be acclimated to the study condition before
91  baseline measurements or interventions are performed. However, an extensive acclimation period risks
92  microbiome drift occurring due to the increasing age of the organism or other unknown factors, so
93  habitat stabilization must be time limited (Hoy *et al.* 2015). Additional bias would also be introduced if
94  the acclimation period is included in the interventional period of the research.

95

## 96  1.3 Genotype, Familial, & Source Differences:

97  Subject genotype, degree of familial relation, and in the case of animal models, the source can
98  significantly impact GMB composition. Differences in the genotype of animal models have been found to
99  impact the diversity and abundance of organisms (Campbell *et al.* 2012, McKnite *et al.* 2012, Leamy *et*
100 *al.* 2014). For this reason, if the study being evaluated uses organisms of significantly different
101 genotypes, such as the use of different strains of mice from the Collaborative Cross, where the effect of
102 genotype difference is not the target of the study, it should be classified as having a high risk of bias.
103 Suppose the study uses a similar genotype between treatment groups, such as the same strain of inbred
104 animal model or monozygotic twin subjects. In that case, it should be considered a low risk of bias for
105 confounding due to the genotype effect.

106  Regarding familial relation, genetically related subjects have been demonstrated to share a core of
107  similar GMB for up to three generations in the female line (Turnbaugh *et al.* 2008, Valles-Colomer *et al.*
108  2021). With animal models, breeding within familial relations is often used to maintain genotypically and
109  GMB homogeneity (Hufeldt *et al.* 2010 ). A caution regarding inbreeding is that while selective breeding
110  between siblings can create a more stable and uniform GMB composition, the effects of genetic drift can
111  also introduce confounders across multiple generations that may affect experimental reproducibility
112  with subsequent generations (Laukens *et al.* 2016).

113  Additionally, with animal models, an organism's litter of origin impacts the gut microbiota (Vilson *et al.*
114  2018, Fujiwara *et al.* 2008). This may relate not only to parent genetics but also to the host of maternal
115  factors that can affect the development of progeny GMB, including mode of delivery, maternal diet,
116  maternal stress, and maternal antibiotic use (Friwell *et al.* 2010, Walker *et al.* 2017, Stokholm *et al.*
117  2014, Golubeva *et al.* 2015, Bailey *et al.* 2004, Zhang *et al.* 2021). For these reasons, if the study being
118  examined utilizes organisms from differing litters (from separate mothers or separate deliveries from
119  the same mother) that have not yet reached their mature adult development and are not randomly
120  assorted between research arms, it should be classified as having a high risk of bias. Suppose a study
121  uses organisms from the same mother and litter or randomly assorts progeny from different mothers
122  and litters. In that case, it should be classified as having a low risk of bias.

123  Regarding sourcing of animal models, subjects sourced from different vendors have substantial
124  differences in GMB at baseline (Rasmussen *et al.* 2019, Long *et al.* 2021, Wolff *et al.* 2020). The
125  microbiological or physiological basis of these effects is unknown but may be due to differential
126  exposures to environmental or infectious factors between vendors (Mandal *et al.* 2020).

127

## 1.4 Extreme Diet

129  Dietary differences have been shown to alter the abundance of most gut microbes (Daniel *et al.* 2014,
130  Ang *et al.* 2020, Li *et al.* 2021, Do *et al.* 2018). Because of this, maintaining the diet of interest is
131  essential to avoid introducing confounding bias to the study. However, it may not always be possible to
132  strictly control diet. This is especially relevant to clinical studies involving humans. In this situation, an
133  evaluation of bias must note how a study documented these diet variations.

134

## 1.5 GMB Normalization

136  It is important to assure organisms being studied in research have similar baseline GMB. This allows for
137  more definitive inference as to the effect of the intervention. Several strategies have been used to make
138  the GMB as similar as possible over time. Removal of the entire GMB through the use of germ-free mice
139  can allow for artificial seeding of a select group of organisms (Kennedy *et al.* 2018, Yi and Li 2012).
140  However, the use of these mice necessarily limits the generalizability of a study. For this reason,
141  research often uses organisms with populated GMBs and rely instead on antibiotics to homogenize the
142  microbiome. The use of antibiotics introduces additional risks of bias which must be considered when
143  evaluating a study (Theriot *et al.* 2016).The most significant risk of bias arises from beginning the
144  intervention of interest before the gut microbiota has stabilized after normalization with antibiotics. The
145  GMB continues to fluctuate unpredictably for long periods following antibiotic administration

146  (Merenstein *et al.* 2021). This variance has been found for at least a year after antibiotic usage in
147  humans and for times ranging between one week and 16 weeks in mice depending on the length of the
148  course of antibiotics used (Elvers *et al.* 2020, Rashid *et al.* 2015, Zhu *et al.* 2021). However, short, or
149  single doses of antibiotics such as those often used to normalize the microbiome allow for substantial
150  stabilization of the GMB within 7 days (Gu *et al*. 2020).

151  A third method used to standardize the GMB is to intermix the bedding of multiple cages and then
152  redistribute it (Miyoshi *et al.* 2018). This method is less invasive than antibiotic usage and has a lower
153  risk of long-term impact on the GMB than the use of antibiotics. The use of homogenization of the
154  bedding allows for similar microbiomes to develop in more mice than can be practically housed in a
155  single cage, where the organisms also share all of their bedding (McCafferty *et al.* 2013).

156  Because of the impact of different methods of GMB normalization, it is critical to note the method that
157  was used to normalize the GMB and how long before the intervention this normalization was
158  completed.

159

## 160  2 - Selection Bias

### 161  2.1 Extreme genotype

162  Host genotype shows a stable and heritable impact on GMB composition (Goodrich *et al.* 2016). In the
163  context of GMB research, extreme genotype selection refers to the selection of GMB subjects with
164  genotypes that vary significantly between subjects within a study. Selection of subjects with identical or
165  similar genetic make-up limits genotype confounding effects. A subject with an established history of
166  use along with maximized genetic correlation can be considered a low risk of selection bias. For
167  example, while inbred Balb/C mice do have an extreme genotype, they also have a long-established
168  history of use in immune modulation studies with their known Th2 immune response wherein they
169  exhibit low IFNy and high IL-4 production (Khan *et al.* 2022, Mills *et al.* 2000, Watanabe *et al.* 2004).
170  Furthermore, prior literature has established the correlation between subject genetics and variation in
171  the GMB population and subsequent disease states (Xu *et al.* 2020).

172

### 173  2.2 Randomization or Demographic Balancing Sufficiently Applied

174  Randomization is essential in ensuring subject-level differences between participants in the intervention
175  and control groups can be attributed to chance alone. It is a standard method that attempts to create
176  the necessary pre-intervention equivalence between groups, allowing for conclusions based on the
177  effect of the intervention. In trials where randomization was not appropriately utilized, the outcome
178  was overestimated by up to 40% compared to trials where randomization was utilized (Suresh, 2011). If
179  randomization was not applied, implementing demographic balancing is an appropriate measure to
180  ensure adequate control and intervention arms distribution. Any demographic balancing performed
181  should be sufficiently described in the study. This method focuses on ensuring each group is
182  demographically balanced at baseline to lessen the difference between groups and utilize randomization
183  if no subject background information is available (Saint, 2015). Both randomization and demographic
184  balancing can be applied to human and animal model studies. For example, in studies utilizing syngeneic

185  mice, randomization must be performed outside the scope of human intervention in that random
186  number generators should assign mice numbers which can then correlate to intervention and control
187  groups, hence this places randomization outside the scope of human influence, limiting bias to a
188  maximum degree. In syngeneic animals, demographic balancing would have a limited impact on the
189  bias, however, wherein studies utilize genetically unrelated animals, the need for implementation of
190  both randomization and demographic balancing is necessary for limiting substantial bias (Hirst *et al.*
191  2014). Similar principles apply in human studies. Given a majority of human studies utilize genetically
192  unrelated subjects, randomization is required to avoid high risk of bias. In human studies, a step beyond
193  randomization should be taken, i.e., implementing blinded randomization with description of the
194  randomization protocol to give the reader the ability to discern breaks in randomization or similar bias
195  control methods within the study (Chalmers *et al.* 1983).

196

## 3 - Classification of Intervention

### 3.1 Intervention Bias

199  Bias in intervention can occur when interventions or outcomes are inappropriately selected for or
200  measured. In non-differential misclassification, test subjects' exposures are misidentified, and they are
201  categorized into the wrong group (McCoy, 2017). This misclassification can dilute the effect of the
202  intervention causing effect estimates to favor the null (LaMorfe, 2016). The probability of non-
203  differential misclassification is equal across all groups. Bias may be reduced by ensuring a proper
204  background check on test subjects and equalizing any differences. On the other hand, differential
205  misclassification occurs when misclassification of exposure or outcome is not equal between subjects
206  and is less easily predictable in whether it will bias results towards or away from the null. Therefore, the
207  probability of assigning subjects to the wrong group differs based on the individual. This may also
208  introduce recall bias towards recalling specific exposures because the subject has the disease state
209  versus a subject that does not. In GMB studies, this may present in the form of researchers explaining
210  results that show a significant effect as attributed to specific causes but leaving out explanations for
211  non-significant results.  Because this type of misclassification is more applicable in case studies, it is less
212  relevant for animal studies but can be prominent in human studies (Spencer *et al.* 2018).

213

### 3.2 Validation of Method

215  The establishment of an effective intervention is imperative for a successful study. Before the
216  experiment, researchers must verify that their chosen intervention method will produce the intended
217  effect. In studies where this is not done, the produced results may or may not be relied on because the
218  protocol was never validated. Verification can be internal (tested and proved by the researchers) or
219  external (via other established studies). If the study calls for a particular disease state to be expressed, it
220  must be validated that the test subjects have the disease state. In studies that call for a specific
221  procedure, there can be potential bias in how the readers know the procedure was correctly obtained if
222  it is not reported. For example, in microbiome hypertension studies, animal subjects were tested based
223  on blood pressure measurements by a well-established method, tail-cuff plethysmography (Marques *et*
224  *al.* 2019). If a lesser-known and validated method was used, it could introduce a high risk of bias if
225  researchers did not verify that their method was accurate.  When testing for the effect of a disease state

226 as influenced by the microbiome, it is helpful to transplant the experimental group microbiome into a
227 germ-free animal model to confirm the effect. This reduces an intermediate risk of bias by
228 demonstrating that the effect of the intervention is associated with the levels of change in the
229 microbiome (Gottfredson *et al.* 2015).

230

# 4 – Deviation from Intervention

232 It is well understood that experiments that deviate from their initial protocol have an increased
233 potential for bias in their study should they decide to include data prior to the deviation. Therefore, all
234 deviations from the protocol should be well documented with time stamps, and the data included in the
235 study should also include the time at which it was collected—either post-protocol or pre-protocol
236 addendum. Rationale and limitations should also be included should researchers decide to include data
237 from any time the protocol was different.

238

# 5 - Missing Data

240 Missing data is prevalent in many academic disciplines, from the social to biomedical sciences, and may
241 contribute to bias in any given study. GMB research likewise suffers from inadequate consideration of
242 missing data and the statistical methods to address it. To begin, two types of missing data should be
243 distinguished: missing data due to patient drop-out in clinical, longitudinal studies and missing data as a
244 result of inadequate sequencing depth leading to "false zeroes" in the microbiome genetic data. Both
245 have potential to increase ROB.

## 5.1 Cause/Category of Missing Data

247 Missing data falls into multiple categories based on the mechanism of missingness: Missing Completely
248 at Random (MCAR), Missing at Random (MAR), and Missing Not at Random (MNAR) (Groenwold and
249 Dekkers 2020). These categories apply assumptions to missing data based on the cause. MCAR assumes
250 that data is missing due to a factor entirely unrelated for the study. MAR assumes data is missing due to
251 observed variables relevant to the study. MNAR assumes data is missing based on unknown or not
252 quantifiable variables to the authors. MAR and MNAR are most relevant to clinical research, specifically
253 in regard to patient drop-out, including clinical GMB trials (Pugh *et al.* 2021). Sampling zeroes in
254 microbiome data are a more generalized form of missing data but are primarily reminiscent of MAR
255 (Kaul *et al.* 2017, Kaul *et al.* 2017). Each of these areas will be further discussed in the following sections.
256 Under MAR, studies may utilize various statistical imputation techniques to replace missing data, though
257 the most well-known and effective method is multiple imputations (Spineli *et al.* 2015). With MNAR,
258 various statistical modeling techniques may address missing data. Such techniques are further discussed
259 in relation to GMB studies in the section "Sequencing Depth and Sampling Zeroes." The distinction
260 between MAR and MNAR also indicates whether bias related to missing data is entirely removable in
261 analysis - the former can, while the latter cannot (Mack *et al.* 2018). This should not be confused with
262 the notion that MNAR assumptions immediately denote a study as biased. If the missingness in MNAR or
263 MAR is independent of the outcome, then the study may be unbiased in regard to missing data. Thus, a
264 study with MNAR data is not necessarily high ROB.

265    Notably, a significant number of studies do not clearly state the mechanism of missingness or adjust for
266    missing data (Carpenter and Smuk, 2021). It is important that studies distinguish mechanism of
267    missingness or explain relevant missing data. If a study does not acknowledge missingness in data or
268    ensures the absence of missing data, the study may be considered high ROB. If a study acknowledges
269    missing data but does not adequately address it through MAR/MNAR distinction and proper statistical
270    techniques related to its missing data category, then the study may be considered intermediate ROB. If a
271    study demonstrates all of this, it may be considered low ROB.

272

## 5.2 Subject Drop-out

274    Missing data in the form of patient drop-out has a marked effect on statistical power, type 1 error, and
275    various outcome measures (Fiero *et al.* 2016, Cai *et al.* 2020, Thompson *et al.* 2011). In traditional
276    clinical research, missing data has a clear effect on useful measures, such as relative risk and risk ratio
277    calculations. Further, although researchers attempt to minimize drop-out and its statistical effects, drop-
278    out ratios were reported to be greater than 40% depending on the study and the degree of
279    unpleasantness in medical interventions to the patient (Schnicker *et al.* 2013, Li *et al.* 2021).
280    Consequently, it has been proposed that a 20% drop-out ratio is reasonable (Furlan *et al.* 2009, Cramer
281    *et al.* 2016). Interestingly, it has been shown that fecal sampling of patients in GMB studies has not been
282    a significant reason for drop-out, suggesting typical sources of patient non-retention (Vandeputte *et al.*
283    2017). The effect of drop-out on statistical measures is expected to be the same in clinical GMB trials.
284    Despite drop-out being common in clinical studies, its effect on outcome measures involving microbial
285    compositional data (e.g., beta diversity) is not currently well described in clinical GMB studies. However,
286    it is expected that such measurements relying on consistent analysis from a wide array of samples will
287    be biased if there is inadequate sampling size.

288    The effect of bias comes into effect when there is interpretation between samples, in that missing data
289    prevents consistent interpretation of genetic data through a larger body of samples. For example,
290    microbiome samples stratified by disease state versus control should be held to higher statistical power,
291    similar to traditional clinical studies. Yet, the complexity of GMB genetic analysis often prevents large
292    sample sizes from being a practical implementation due to costs unless utilizing less-expensive protocols
293    such as those involving qPCR to monitor microbial composition at high taxonomic levels (i.e., phyla)
294    (Koliada *et al.* 2020). Some studies demonstrate shallow shotgun metagenomic sequencing as an
295    alternative methodology for large, longitudinal GMB studies (Xu *et al.* 2021). Nonetheless, making
296    interpretations in GMB data between samples stratified by host conditions may need to be more
297    consistent and accurate when samples are unavailable from a patient drop-out. Based on the literature
298    of other areas in clinical research as discussed, it is again reasonable to assert that drop-out will
299    influence outcome measures if authors make interpretations across hosts of varying condition states.

300    Due to few clinical studies analyzing the effect of drop-out on GMB outcomes, it is reasonable to use a
301    20% patient drop-out ratio, as many clinical trials traditionally utilize. GMB studies that have a high
302    patient dropout are considered high ROB. GMB studies that have low patient drop-out are considered
303    low ROB.

304

## 5.3 Sequencing Depth and Sampling Zeroes

GMB researchers should consider sequencing depth as a contributor to missing data and subsequent bias. It is established that low-sequencing depth (2000 single-end reads per sample) can adequately predict the same diversity patterns as high-depth sequencing (on the scale of millions of reads per sample) (Caporaso *et al.* 2011, Lundin *et al.* 2012, Xiao *et al.* 2018). Experiments that quantify GMB outcome measures (like alpha and beta diversity) should utilize the same depth for all samples. Bias would be introduced if different sequencing depths are used for a set of samples. It should be noted, however, that false zeroes influence microbiome genetic data at both high and low depth. While true zeroes (or biological zeroes) represent true taxonomic absences, false zeroes (or sampling zeroes) represent a lack of sequencing depth to adequately detect certain microbial taxa. Notably, low sequencing depth, as is often the case of 16S rRNA sequencing, may not detect low abundance taxa or low taxa (subspecies) due to lower resolution. Though whole genome sequencing (WGS), such as shotgun metagenomic sequencing, utilizes high sequencing depth to sequence entire genomes, sampling zeroes still persist (Pereira-Marques *et al.* 2019).

At the time of writing, this issue of zero-inflation – or the excess of sampling zeroes at high and low depth – and the resulting bias in GMB genetic data is an active area of research. Interestingly, relatively few studies utilize any statistical modeling to correct for such missing data. Yet, various modeling techniques were recently developed to address zero-inflation (Deek and Li, 2021, Zhang *et al.* 2020, Ha *et al.* 2020). Similar to modeling techniques, imputation is a method traditionally used to address missing data in the form of patient drop out, but a promising imputation method is recently available to also deal with GMB sampling zeroes. Previous studies showed an increase of Pearson correlation from 0.59 (between 16S and WGS in non-corrected data) to 0.64 (between 16S and WGS in corrected data) (Jiang *et al. 2021)*. There were also marked differences in mean and standard deviation of abundances per taxon between corrected and non-corrected data. This suggests greater homogeneity of samples across sequencing methods if imputation is utilized to correct data. However, as our article focuses on the role of bias in GMB research, we do not yet place best-practice recommendations for a particular method of missing data correction.

As of date, few GMB studies utilize statistical techniques to correct for sampling zeroes. Furthermore, common bioinformatics pipelines (such as QIIME2) do not incorporate such techniques into data-correction programs.

As such, the available literature suggests future GMB studies that do not consider sampling zeroes and lack a statistical technique for missing data correction may be considered high ROB. Studies that utilize missing data correction may be considered low ROB. These data correction methods, once more, include various modeling techniques or imputation.

# 6 – Measurement of Outcomes

## 6.1 Sample collection

Currently, there is no standard method for sample collection for GMB studies. While biopsy of the lower intestine provides a controlled sampling site and an accurate microbiota account, it is expensive, time-consuming, and unsuitable for healthy control groups. In contrast, fecal collection is non-invasive and

345  cost-effective (Tang *et al.* 2020). Thus, it is a standard sampling method in both clinical and research
346  applications. However, fecal collection introduces temporal inconsistency that is a risk of bias when
347  unaccounted for.

348  Fecal samples collected at different times of the day are at risk for inaccurate representation of the
349  absolute abundance of gut microbiota (Caporaso *et al.* 2011). Specifically for mouse studies, the
350  snapshots of the microbiota provided by fecal samples is more accurate and consistent within treatment
351  groups when collected in the morning due to the nocturnal feeding nature of mice (Jones *et al.* 2021).
352  For studies involving subjects with unpredictable and inconsistent bowel movements, samples should be
353  preserved immediately after defecation as oxidation of the outer layer can alter the microbiota (Pepper
354  and Rosenfeld, 2012). Specifically, Firmicutes and Bifidobacteria Spp. are two known phylum that are
355  unstable in the outer microenvironment when exposed to oxygen (Gorselak *et al.* 2015). Therefore, to
356  minimize the differential errors, the methods of measurement must be consistent between control and
357  intervention groups.

358

## 6.2 Blinding

360  In a GMB study, the primary outcome is based on definitive and objective genetic sequencing.
361  Therefore, assessor bias is typically negligible, and a low risk of bias is expected (Higgins *et al.* 2022).

362

# 7 – Reporting of Results

## 7.1 Selection of Reported Results
365  Selective reporting of results can lead to biased interpretations of significance and or non-significance
366  via particular selection of results from multiple outcome measures in estimating outcome effect. Bias in
367  selection of reported results can be difficult to detect without access to a protocol from which one can
368  compare pre-specified intended outcomes of interest to the outcomes analyzed in the published paper
369  (Heneghan *et al.* 2019). Often, results are selected for significance, omitted for non-significance, or
370  omitted for adverse effect of intervention (Dwan *et al.* 2013, Hedin *et al.* 2016, Van der Steen *et al.*
371  2019).

372

# Validation Test
374  Four medical students with no prior experience in ROB assessment were recruited to test this tool by
375  using it to independently assess ROB on three selected studies of similar length in a predetermined
376  sequence (Wu *et al.* 2017, Mohammed *et al.* 2020, Saunders *et al.* 2020). Subjects were provided with
377  the manuscript and ROB rubric. They were asked to track time to completion per study and complete
378  the ROB rubric for each study. Subjects assessed ROB in an average of 44.75 minutes per study with time
379  to completion generally decreasing from the first study assessed to the last study assessed.

380  Inter-rater variability was assessed by assigning values of 1, 2, and 3 to low, medium, and high ROB in
381  order to construct visual representations of rater scores in each sub-domain of bias and to compare
382  summed ROB scores between raters for each study. **Figures 1.1-1.3** demonstrate variability within a

383    study in each subdomain of bias assessed by this tool between raters. The figures demonstrate similar
384    ROB judgements between at least three of four raters in the majority of subdomains across the three
385    studies assessed.

386    **Figure 2** demonstrates variation in summed ROB score by rater for each of the three studies.  It shows
387    the decreasing magnitude of difference between raters' summed ROB scores with each subsequent use
388    of the tool from a max-score min-score difference of six points in study1 and study3, and of four points
389    in study2 out of 45 possible points. One way ANOVA test of rater subdomain scores across all
390    subdomains for each study returned p-values of 0.554, 0.568, and 0.399 for study1, study2, and study3
391    respectively indicating no significant difference between overall ROB assessment scores between raters
392    of the same study. First time ROB assessors using this tool showed a relatively high degree of
393    concordance in ROB determination at the subdomain level and in magnitude of summed ROB score.

394

## Conclusion

396    Risk of bias assessment is a crucial step in systematic review and metanalysis to assess quality of
397    information being collected. By outlining common sources of bias that can impact GMB research
398    following the structure of the ROBINS-I tool, this tool can serve as an adjunct to improve and
399    standardize ROB assessment of GMB studies. A standardized ROB assessment for GMB studies will
400    improve accuracy of risk assessment, improve reproducibility between researchers, and promote the
401    inclusion of high-quality information in systematic reviews and metanalyses of the GMB.

402

403 **Financial Support**

404 This work was supported by the Iowa Osteopathic Education and Research (IOER) Foundation.

405

406 **Conflicts of Interest declarations in manuscripts**

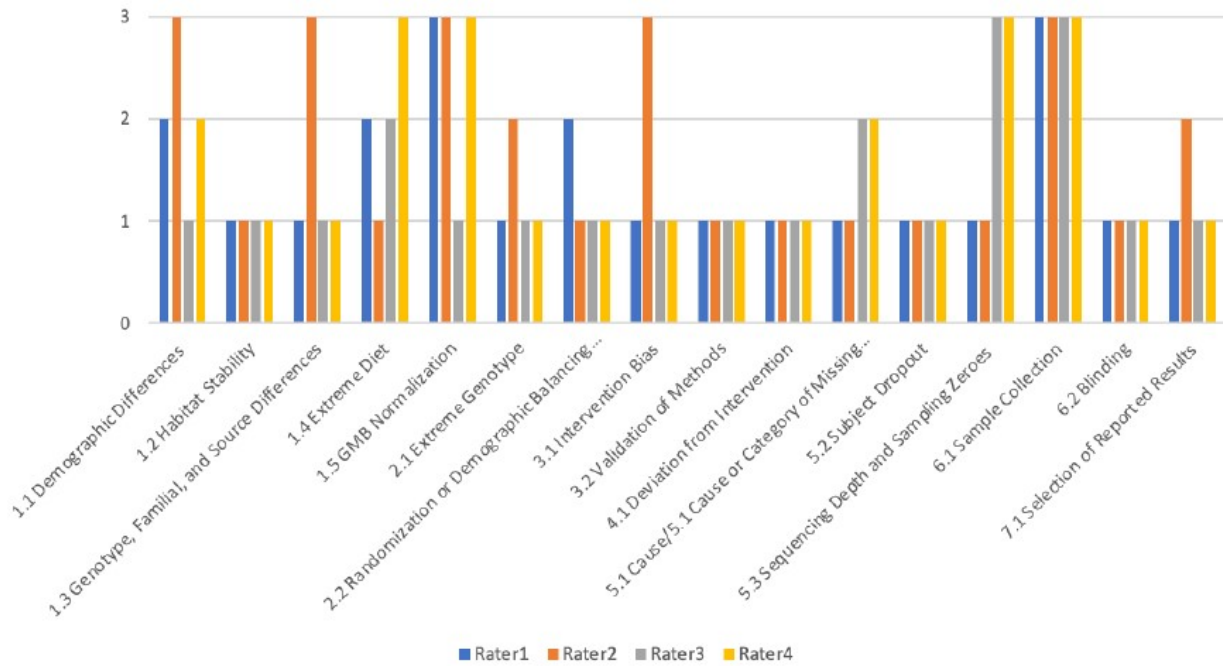407 Authors declare no conflicts of interest.

408

409 **Research Transparency and Reproducibility**

410 Following the journal's policy for supporting research transparency and reproducibility, we will make all
411 data and protocols available to readers.
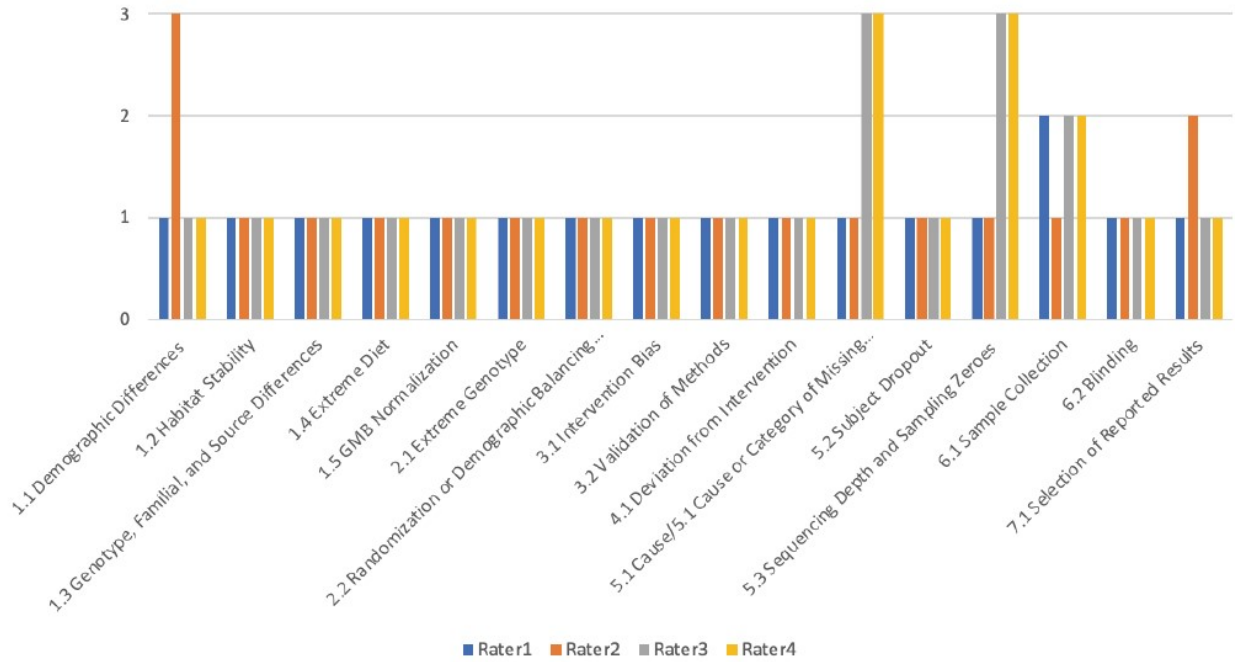
412 **Figure Captions:**

413 **Figure 1.1** - Inter-rater variability in ROB determinations by subdomain for validation test study 1,
414 *"Metformin alters the gut microbiome of individuals with treatment-naive type 2 diabetes, contributing*
415 *to the therapeutic effects of the drug"* by Wu *et al.* 2017, where "1" on the y-axis indicates that the rater
416 determined the study to be at low ROB for the subdomain indicated on the x-axis; "2" indicates medium
417 ROB and "3" indicates a high ROB determination by the individual rater.



418

419

420 **Figure 1.2** - Inter-rater variability in ROB determinations by subdomain for validation test on study 2,
421 *"Protective effects of Δ9-tetrahydrocannabinol against enterotoxin-induced acute respiratory distress*
422 *syndrome are mediated by modulation of microbiota"* by Mohammed *et al.* 2020, where "1" on the y-
423 axis indicates that the rater determined the study to be at low ROB for the subdomain indicated on the
424 x-axis; "2" indicates medium ROB and "3" indicates a high ROB determination by the individual rater.
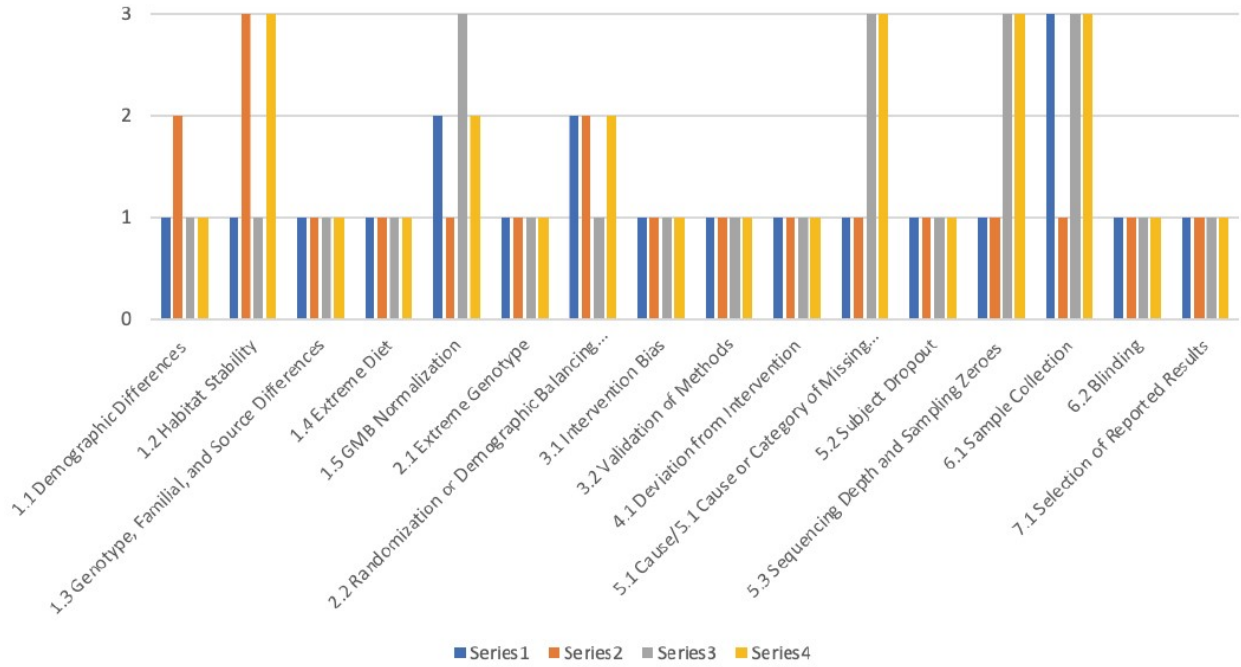


425

426

427

428    **Figure 1.3** - Inter-rater variability in ROB determinations by subdomain for validation test on study 3,

429    *"Gut microbiota manipulation during the prepubertal period shapes behavioral abnormalities in a mouse*

430    *neurodevelopmental disorder model"* by Saunders *et al.* 2020, where "1" on the y-axis indicates that the

431    rater determined the study to be at low ROB for the subdomain indicated on the x-axis; "2" indicates

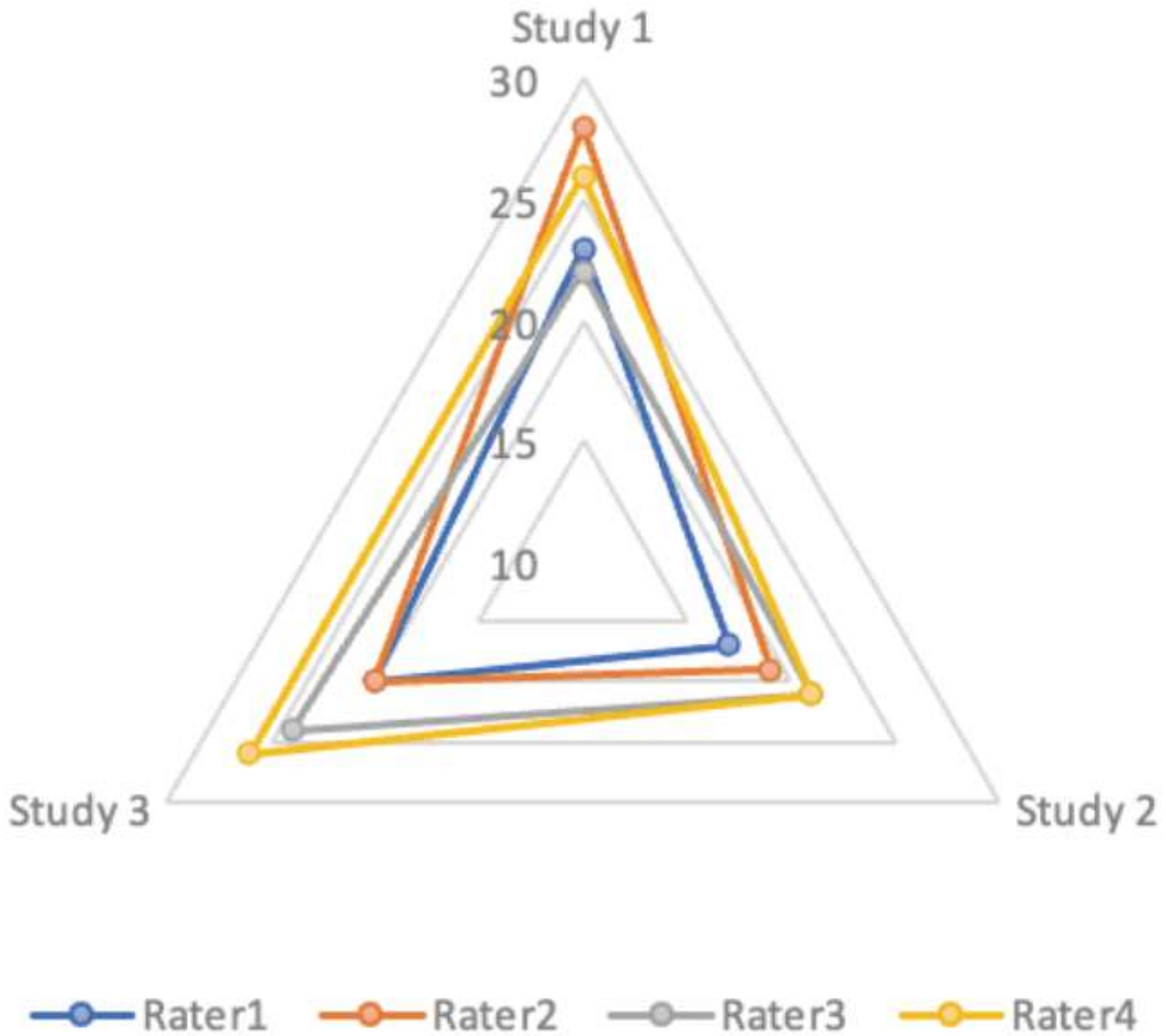432    medium ROB and "3" indicates a high ROB determination by the individual rater.



433

434

435 **Figure 2** – Visual representation comparing summed ROB score (as determined by assigning point values
436 of 1, 2, and 3 to low, medium, and high ROB respectively) by rater for each of the three studies assessed
437 in the validation test where each increasingly large concentric triangle indicates an increase of 5 points.



438

439

440 **References:**

441 Org *et al.* (2016). Sex differences and hormonal effects on gut microbiota composition in mice. *Gut*
442 *microbes*, *7*(4), 313–322.

443 Kim *et al.* (2020). Sex Differences in Gut Microbiota. *The world journal of men's health*, *38*(1), 48–60.
444 https://doi.org/10.5534/wjmh.190009

445 Ticinesi *et al.* (2019). Gut Microbiota, Muscle Mass and Function in Aging: A Focus on Physical Frailty and
446 Sarcopenia. *Nutrients*, *11*(7), 1633.

447 Liu *et al.* (2020). Aging Increases the Severity of Colitis and the Related Changes to the Gut Barrier and
448 Gut Microbiota in Humans and Mice. *The journals of gerontology. Series A, Biological sciences and*
449 *medical sciences*, *75*(7), 1284–1292.

450 Yoon *et al.* (2021). Association between aging-dependent gut microbiome dysbiosis and dry eye severity
451 in C57BL/6 male mouse model: a pilot study. *BMC microbiology*, *21*(1), 106.

452 Singh *et al.* (2021). Cage and maternal effects on the bacterial communities of the murine gut. *Scientific*
453 *reports*, *11*(1), 9841. https://doi.org/10.1038/s41598-021-89185-5

454 Lipinski *et al.* (2021). Cage Environment Regulates Gut Microbiota Independent of Toll-Like
455 Receptors. *Infection and immunity*, *89*(9), e0018721. https://doi.org/10.1128/IAI.00187-21

456 Montonye *et al.* (2018). Acclimation and Institutionalization of the Mouse Microbiota Following
457 Transportation. *Frontiers in microbiology*, *9*, 1085. https://doi.org/10.3389/fmicb.2018.01085

458 Capdevila *et al.* (2007). Acclimatization of rats after ground transportation to a new animal
459 facility. *Laboratory animals*, *41*(2), 255–261. https://doi.org/10.1258/002367707780378096

460 Hoy *et al.* (2015). Variation in Taxonomic Composition of the Fecal Microbiota in an Inbred Mouse Strain
461 across Individuals and Time. *PloS one*, *10*(11), e0142825. https://doi.org/10.1371/journal.pone.0142825

462 Campbell *et al.* (2012). Host genetic and environmental effects on mouse intestinal microbiota. *The*
463 *ISME journal*, *6*(11), 2033–2044. https://doi.org/10.1038/ismej.2012.54

464 McKnite *et al.* (2012). Murine gut microbiota is defined by host genetics and modulates variation of
465 metabolic traits. *PloS one*, *7*(6), e39191. https://doi.org/10.1371/journal.pone.0039191

466 Leamy *et al.* (2014). Host genetics and diet, but not immunoglobulin A expression, converge to shape
467 compositional features of the gut microbiome in an advanced intercross population of mice. *Genome*
468 *biology*, *15*(12), 552. https://doi.org/10.1186/s13059-014-0552-6

469 Turnbaugh *et al.* (2008). A core gut microbiome in obese and lean twins. *Nature*, *457*(7228), 480–484.
470 https://doi.org/10.1038/nature07540

471 Valles-Colomer *et al.* (2021). Variation and transmission of the human gut microbiota across multiple
472 familial generations. *Nature Microbiology*, *7*(1), 87–96. https://doi.org/10.1038/s41564-021-01021-8

473     Hufeldt *et al.* (2010). Family relationship of female breeders reduce the systematic inter-individual
474     variation in the gut microbiota of inbred laboratory mice. *Laboratory Animals, 44*(4), 283–9.
475     https://doi.org/10.1258/la.2010.010058

476     Laukens *et al.* (2016). Heterogeneity of the gut microbiome in mice: guidelines for optimizing
477     experimental design. *Fems Microbiology Reviews, 40*(1), 117–32.
478     https://doi.org/10.1093/femsre/fuv036

479     Vilson *et al.* (2018). Disentangling factors that shape the gut microbiota in German Shepherd dogs. *PloS*
480     *one*, *13*(3), e0193507. https://doi.org/10.1371/journal.pone.0193507

481     Fujiwara *et al.* (2008). Assessing changes in composition of intestinal microbiota in neonatal BALB/c mice
482     through cluster analysis of molecular markers. *The British journal of nutrition*, *99*(6), 1174–1177.
483     https://doi.org/10.1017/S0007114507862349

484     Friswell *et al.* (2010). Site and strain-specific variation in gut microbiota profiles and metabolism in
485     experimental mice. *Plos One, 5*(1), 8584. https://doi.org/10.1371/journal.pone.0008584

486     Walker *et al.* (2017). The prenatal gut microbiome: are we colonized with bacteria in utero ?:
487     colonization of the gut microbiome in utero. *Pediatric Obesity, 12*, 3–17.
488     https://doi.org/10.1111/ijpo.12217

489     Stokholm *et al.* (2014). Antibiotic use during pregnancy alters the commensal vaginal microbiota. *Clinical*
490     *Microbiology and Infection : The Official Publication of the European Society of Clinical Microbiology and*
491     *Infectious Diseases, 20*(7), 629–35. https://doi.org/10.1111/1469-0691.12411

492     Golubeva *et al.* (2015). Prenatal stress-induced alterations in major physiological systems correlate with
493     gut microbiota composition in adulthood. *Psychoneuroendocrinology, 60*, 58–74.
494     https://doi.org/10.1016/j.psyneuen.2015.06.002

495     Bailey *et al.* (2004). Prenatal Stress Alters Bacterial Colonization of the Gut in Infant Monkeys. *Journal of*
496     *Pediatric Gastroenterology and Nutrition*, *38*(4), 414–421. https://doi.org/10.1097/00005176-
497     200404000-00009

498     Zhang *et al.* (2021). Prenatal stress leads to deficits in brain development, mood related behaviors and
499     gut microbiota in offspring. *Neurobiology of Stress*, *15*, 100333.
500     https://doi.org/10.1016/j.ynstr.2021.100333

501     Rasmussen *et al.* (2019). Mouse Vendor Influence on the Bacterial and Viral Gut Composition Exceeds
502     the Effect of Diet. *Viruses*, *11*(5), 435. https://doi.org/10.3390/v11050435

503     Long *et al.* (2021). Shared and distinctive features of the gut microbiome of C57BL/6 mice from different
504     vendors and production sites, and in response to a new vivarium. *Lab animal*, *50*(7), 185–195.
505     https://doi.org/10.1038/s41684-021-00777-0

506     Wolff *et al.* (2020). Vendor effects on murine gut microbiota and its influence on lipopolysaccharide-
507     induced lung inflammation and Gram-negative pneumonia. *Intensive care medicine experimental*, *8*(1),
508     47. https://doi.org/10.1186/s40635-020-00336-w

509    Mandal *et al*. (2020). Temporospatial shifts within commercial laboratory mouse gut microbiota impact
510    experimental reproducibility. *BMC biology*, *18*(1), 83. https://doi.org/10.1186/s12915-020-00810-7

511    Daniel *et al.* (2014). High-fat diet alters gut microbiota physiology in mice. *The ISME journal*, *8*(2), 295–
512    308. https://doi.org/10.1038/ismej.2013.155

513    Ang *et al.* (2020). Ketogenic Diets Alter the Gut Microbiome Resulting in Decreased Intestinal Th17
514    Cells. *Cell*, *181*(6), 1263–1275.e16. https://doi.org/10.1016/j.cell.2020.04.027

515    Li *et al.* (2021). Ketogenic Diets Induced Glucose Intolerance and Lipid Accumulation in Mice with
516    Alterations in Gut Microbiota and Metabolites. *mBio*, *12*(2), e03601-20.
517    https://doi.org/10.1128/mBio.03601-20

518    Do *et al.* (2018). High-Glucose or -Fructose Diet Cause Changes of the Gut Microbiota and Metabolic
519    Disorders in Mice without Body Weight Change. *Nutrients*, *10*(6), 761.
520    https://doi.org/10.3390/nu10060761

521    Kennedy *et al.* (2018). Mouse Microbiota Models: Comparing Germ-Free Mice and Antibiotics Treatment
522    as Tools for Modifying Gut Bacteria. *Frontiers in physiology*, *9*, 1534.
523    https://doi.org/10.3389/fphys.2018.01534

524    Yi and Li (2012). The germfree murine animal: an important animal model for research on the
525    relationship between gut microbiota and the host. *Veterinary microbiology*, *157*(1-2), 1–7.
526    https://doi.org/10.1016/j.vetmic.2011.10.024

527    Theriot *et al.* (2016). Antibiotic-Induced Alterations of the Gut Microbiota Alter Secondary Bile Acid
528    Production and Allow for Clostridium difficile Spore Germination and Outgrowth in the Large
529    Intestine. *mSphere*, *1*(1), e00045-15. https://doi.org/10.1128/mSphere.00045-15

530    Merenstein *et al.* (2021). *Bifidobacterium animalis* subsp. *lactis* BB-12 Protects against Antibiotic-
531    Induced Functional and Compositional Changes in Human Fecal Microbiome. *Nutrients*, *13*(8), 2814.
532    https://doi.org/10.3390/nu13082814

533    Elvers *et al.* (2020). Antibiotic-induced changes in the human gut microbiota for the most commonly
534    prescribed antibiotics in primary care in the UK: a systematic review. *BMJ open*, *10*(9), e035677.
535    https://doi.org/10.1136/bmjopen-2019-035677

536    Rashid *et al.* (2015). Determining the Long-term Effect of Antibiotic Administration on the Human
537    Normal Intestinal Microbiota Using Culture and Pyrosequencing Methods. *Clinical infectious diseases :*
538    *an official publication of the Infectious Diseases Society of America*, *60 Suppl 2*, S77–S84.
539    https://doi.org/10.1093/cid/civ137

540    Zhu *et al.* (2021). Effects of long-term antibiotic treatment on mice urinary aromatic amino acid profiles.
541    Bioscience reports, 41(1), BSR20203498. https://doi.org/10.1042/BSR20203498

542    Gu *et al.* (2020). Effect of the Short-Term Use of Fluoroquinolone and β-Lactam Antibiotics on Mouse
543    Gut Microbiota. *Infection and drug resistance*, *13*, 4547–4558. https://doi.org/10.2147/IDR.S281274

544    Miyoshi *et al.* (2018). Minimizing confounders and increasing data quality in murine models for studies
545    of the gut microbiome. *PeerJ*, *6*, e5166. https://doi.org/10.7717/peerj.5166

546    McCafferty *et al.* (2013). Stochastic changes over time and not founder effects drive cage effects in
547    microbial community assembly in a mouse model. *The ISME journal*, *7*(11), 2116–2125.
548    https://doi.org/10.1038/ismej.2013.106

549    Goodrich *et al.* (2016). Genetic determinants of the gut microbiome in UK twins. Cell Host & Microbe,
550    19(5), 731–743. https://doi.org/10.1016/j.chom.2016.04.017

551    Khan *et al.* (2022). Effects of Shrimp Peptide Hydrolysate on Intestinal Microbiota Restoration and
552    Immune Modulation in Cyclophosphamide-Treated Mice. Molecules, 27(5), 1720.
553    https://doi.org/10.3390/molecules27051720

554    Mills *et al.* (2000). M-1/M-2 macrophages and the Th1/Th2 paradigm. Journal of immunology, 164(12),
555    6166–6173. https://doi.org/10.4049/jimmunol.164.12.6166

556    Watanabe *et al.* (2004). Innate immune response in Th1- and Th2-dominant mouse strains. Shock, 22(5),
557    460–466. https://doi.org/10.1097/01.shk.0000142249.08135.e9

558    Xu *et al.* (2020). The interplay between host genetics and the gut microbiome reveals common and
559    distinct microbiome features for complex human diseases. Microbiome, 8(1), 145.
560    https://doi.org/10.1186/s40168-020-00923-9

561    Suresh (2011). An overview of randomization techniques: An unbiased assessment of outcome in clinical
562    research. Journal of human reproductive sciences , 4 (1), 8 – 11. https://doi.org/10.4103/0974 -
563    1208.82352

564    Saint (2015). Randomization Does Not Help Much, Comparability Does. PloS one, 10 (7), e0132102.
565    https://doi.org/10.1371/journal.pone.0132102

566    Hirst *et al.* (2014) The Need for Randomization in Animal Trials: An Overview of Systematic Reviews.
567    PLOS ONE 9(6): e98856. https://doi.org/10.1371/journal.pone.0098856

568    Chalmers *et al.* (1983). Bias in treatment assignment in controlled clinical trials. *The New England journal*
569    *of medicine*, *309*(22), 1358–1361. https://doi.org/10.1056/NEJM198312013092204

570    McCoy  (2017). Understanding the Intention-to-treat Principle in Randomized Controlled Trials. *The*
571    *western journal of emergency medicine*, *18*(6), 1075–1078.
572    https://doi.org/10.5811/westjem.2017.8.35985

573    LaMorfe (2016). Information Bias (Observation Bias). Information bias(observation bias). Retrieved
574    March 26, 2022, fromhttps://sphweb.bumc.bu.edu/otlt/mph-
575    modules/ep/ep713_bias/EP713_Bias4.html

576    Spencer *et al.* **Misclassification bias.** In Catalogue Of Bias 2018.
577    http://www.catalogueofbiases.org/biases/misclassificationbias

578    Marques *et al.* (2019). Guidelines for transparency on GutMicrobiome Studies in essential and
579    experimental hypertension. Hypertension,74(6),1279–
580    1293.https://doi.org/10.1161/hypertensionaha.119.13079

581    Gottfredson *et al.* Standards of Evidence for Efficacy, Effectiveness, and Scale-up Research in Prevention
582    Science: Next Generation. Prev Sci. 2015 Oct;16(7):893-926. doi: 10.1007/s11121-015-0555-x.
583    PMID:25846268; PMCID: PMC4579256.

584    Groenwold and Dekkers (2020). Missing data: the impact of what is not there. *European journal of*
585    *endocrinology*, *183*(4), E7–E9. https://doi.org/10.1530/EJE-20-0732

586    Pugh *et al.* (2021). Missing repeated measures data in clinical trials. *Neuro-oncology practice*, *9*(1), 35–
587    42. https://doi.org/10.1093/nop/npab043

588    Kaul *et al.* (2017). Structural zeros in high-dimensional data with applications to microbiome studies.
589    *Biostatistics (Oxford, England)*, *18*(3), 422–433. https://doi.org/10.1093/biostatistics/kxw053

590    Kaul *et al.* (2017). Analysis of Microbiome Data in the Presence of Excess Zeros. *Frontiers in*
591    *microbiology*, *8*, 2114. https://doi.org/10.3389/fmicb.2017.02114

592    Spineli *et al.* (2015). Addressing missing participant outcome data in dental clinical trials. *Journal of*
593    *dentistry*, *43*(6), 605–618. https://doi.org/10.1016/j.jdent.2015.03.007

594    Mack *et al.* Managing Missing Data in Patient Registries: Addendum to Registries for Evaluating Patient
595    Outcomes: A User's Guide, Third Edition [Internet]. Rockville (MD): Agency for Healthcare Research and
596    Quality (US); 2018 Feb. Types of Missing Data. Available from:
597    https://www.ncbi.nlm.nih.gov/books/NBK493614/

598    Carpenter and Smuk (2021). Missing data: A statistical framework for practice. *Biometrical journal.*
599    *Biometrische Zeitschrift*, *63*(5), 915–947. https://doi.org/10.1002/bimj.202000196

600    Fiero *et al.* (2016). Statistical analysis and handling of missing data in cluster randomized trials: a
601    systematic review. *Trials*, *17*, 72. https://doi.org/10.1186/s13063-016-1201-z

602    Cai *et al.* (2020). Estimands and missing data in clinical trials of chronic pain treatments: advances in
603    design and analysis. *Pain*, *161*(10), 2308–2320. https://doi.org/10.1097/j.pain.0000000000001937

604    Thompson *et al.* (2011). Who's left? Symptoms of schizophrenia that predict clinical trial dropout.
605    *Human psychopharmacology*, *26*(8), 609–613. https://doi.org/10.1002/hup.1253,

606    Schnicker *et al.* (2013). Drop-out and treatment outcome of outpatient cognitive-behavioral therapy for
607    anorexia nervosa and bulimia nervosa. *Comprehensive psychiatry*, *54*(7), 812–823.
608    https://doi.org/10.1016/j.comppsych.2013.02.007

609    Li *et al.* (2021). Analysis of risk factors and construction of prediction model of drop out from peritoneal
610    dialysis. *Medicine*, *100*(3), e24195. https://doi.org/10.1097/MD.0000000000024195

611    Furlan *et al.* (2009). 2009 updated method guidelines for systematic reviews in the Cochrane Back
612    Review Group. *Spine*, *34*(18), 1929–1941. https://doi.org/10.1097/BRS.0b013e3181b1c99f

613    Cramer *et al.* (2016). A Systematic Review and Meta-Analysis Estimating the Expected Dropout Rates in
614    Randomized Controlled Trials on Yoga Interventions. *Evidence-based complementary and alternative*
615    *medicine : eCAM*, *2016*, 5859729. https://doi.org/10.1155/2016/5859729

616   Vandeputte *et al.* Practical considerations for large-scale gut microbiome studies, *FEMS Microbiology*
617   *Reviews*, Volume 41, Issue Supplement_1, August 2017, Pages S154–S167,
618   https://doi.org/10.1093/femsre/fux027

619   Koliada *et al.* (2020). Seasonal variation in gut microbiota composition: cross-sectional evidence from
620   Ukrainian population. *BMC microbiology*, *20*(1), 100. https://doi.org/10.1186/s12866-020-01786-8

621   Xu *et al.* (2021). Characterization of Shallow Whole-Metagenome Shotgun Sequencing as a High-
622   Accuracy and Low-Cost Method by Complicated Mock Microbiomes. *Frontiers in microbiology*, *12*,
623   678319. https://doi.org/10.3389/fmicb.2021.678319

624   Caporaso *et al.* (2011). Global patterns of 16S rRNA diversity at a depth of millions of sequences per
625   sample. *Proceedings of the National Academy of Sciences of the United States of America*, *108 Suppl*
626   *1*(Suppl 1), 4516–4522. https://doi.org/10.1073/pnas.1000080107

627   Lundin *et al.* (2012). Which sequencing depth is sufficient to describe patterns in bacterial α- and β-
628   diversity?. *Environmental microbiology reports*, *4*(3), 367–372. https://doi.org/10.1111/j.1758-
629   2229.2012.00345.x

630   Xiao *et al.* (2018). Necessary Sequencing Depth and Clustering Method to Obtain Relatively Stable
631   Diversity Patterns in Studying Fish Gut Microbiota. *Current microbiology*, *75*(9), 1240–1246.
632   https://doi.org/10.1007/s00284-018-1516-y

633   Pereira-Marques *et al.* (2019). Impact of Host DNA and Sequencing Depth on the Taxonomic Resolution
634   of Whole Metagenome Sequencing for Microbiome Analysis. *Frontiers in microbiology*, *10*, 1277.
635   https://doi.org/10.3389/fmicb.2019.01277

636   Deek and Li  (2021). A Zero-Inflated Latent Dirichlet Allocation Model for Microbiome Studies. *Frontiers*
637   *in genetics*, *11*, 602594. https://doi.org/10.3389/fgene.2020.602594

638   Zhang *et al.* (2020). Zero-Inflated gaussian mixed models for analyzing longitudinal microbiome
639   data. *PloS one*, *15*(11), e0242073. https://doi.org/10.1371/journal.pone.0242073

640   Ha *et al.* (2020). Compositional zero-inflated network estimation for microbiome data. *BMC*
641   *bioinformatics*, *21*(Suppl 21), 581. https://doi.org/10.1186/s12859-020-03911-w

642   Jiang *et al.* (2021). mbImpute: an accurate and robust imputation method for microbiome data. *Genome*
643   *biology*, *22*(1), 192. https://doi.org/10.1186/s13059-021-02400-4

644   Tang *et al.* (2020). Current Sampling Methods for Gut Microbiota: A Call for More Precise
645   Devices. *Frontiers in cellular and infection microbiology*, *10*, 151.
646   https://doi.org/10.3389/fcimb.2020.00151

647   Caporaso *et al.* (2011). Moving pictures of the human microbiome. Genome biology, 12(5), R50.
648   https://doi.org/10.1186/gb-2011-12-5-r50

649   Jones *et al.* (2021). Fecal sample collection methods and time of day impact microbiome composition
650   and short chain fatty acid concentrations. *Scientific reports*, *11*(1), 13964.
651   https://doi.org/10.1038/s41598-021-93031-z

652 Pepper and Rosenfeld (2012). The emerging medical ecology of the human gut microbiome. *Trends in*
653 *ecology & evolution*, *27*(7), 381–384. https://doi.org/10.1016/j.tree.2012.03.002

654 Gorzelak *et al.* (2015). Methods for Improving Human Gut Microbiome Data by Reducing Variability
655 through Sample Processing and Storage of Stool. *PloS one*, *10*(8), e0134802.
656 https://doi.org/10.1371/journal.pone.0134802

657 Higgins *et al. Cochrane Handbook for Systematic Reviews of Interventions* version 6.3 (updated
658 February 2022). Cochrane, 2022. Available from www.training.cochrane.org/handbook.

659 Heneghan (2019). *Outcome reporting bias*. Catalog of Bias. Retrieved September 6, 2022, from
660 https://catalogofbias.org/biases/outcome-reporting-bias/

661 Dwan *et al.* (2013). Systematic Review of the Empirical Evidence of Study Publication Bias and Outcome
662 Reporting Bias — An Updated Review. *PLoS ONE*, *8*(7), e66844.
663 https://doi.org/10.1371/journal.pone.0066844

664 Hedin *et al.* (2016). Publication Bias and Nonreporting Found in Majority of Systematic Reviews and
665 Meta-analyses in Anesthesiology Journals. *Anesthesia & Analgesia, 123* (4), 1018-1025. doi:
666 10.1213/ANE.0000000000001452.

667 Van der Steen *et al.* (2019). Causes of reporting bias: a theoretical framework. *F1000Research*, *8*, 280.
668 https://doi.org/10.12688/f1000research.18310.2

669 Wu *et al.* (2017). Metformin alters the gut microbiome of individuals with treatment-naive type 2
670 diabetes, contributing to the therapeutic effects of the drug. *Nature Medicine*, *23*(7), 850–858.
671 https://doi.org/10.1038/nm.4345

672 Mohammed *et al.* (2020). Protective effects of Δ9-tetrahydrocannabinol against enterotoxin-induced
673 acute respiratory distress syndrome are mediated by modulation of microbiota. *British Journal of*
674 *Pharmacology*, *177*(22), 5078–5095. https://doi.org/10.1111/bph.15226

675 Saunders *et al.* (2020). Gut microbiota manipulation during the prepubertal period shapes behavioral
676 abnormalities in a mouse neurodevelopmental disorder model. *Scientific Reports*, *10*(1).
677 https://doi.org/10.1038/s41598-020-61635-6

678

| Domain | High ROB | Moderate ROB | Low ROB |
|---|---|---|---|
| **1 - Confounding** | | | |
| 1.1 Demographic Differences | - Age: consistently different between study arms<br><br>- Sex: consistently different between study arms | - Age: mixed ages within a study arm, but equal in distribution between study arms<br><br>- Sex: mixed sexes within a study arm, but equal in distribution between study arms | - Age: consistently similar between study arms<br><br>- Sex: consistently similar between study arms |
| 1.2 Habitat Stability | - No acclimation period, or acclimation period <2 days<br><br>- Acclimation period included in the interventional period | - Acclimation period ≥2 days but <5 days | - Acclimation period ≥5 days and <9 weeks |
| 1.3 Genotype, Familial, and Source Differences | - Significantly different subject genotypes between study arms (where genotype effect is not the target of investigation)<br><br>- Non-matured animal models from different litters and/or mothers without random assortment into study arms<br><br>- Comparison of animal subjects from different source or vendor between study arms | - Animal subjects from same vendor, but from separate and temporally spaced orders without random assortment into study arms | - Adequately similar genotypes used between study arms (where host genotype effect is not the target of study)<br><br>- Animal subjects from same litter<br><br>- Animal subjects from same vendor and same order<br><br>- Adult animal subjects from different litters/mothers/vendors randomly assorted into study arms |

| | | | |
|---|---|---|---|
| 1.4 Extreme Diet | - No statement of dietary standards or documentation of dietary variation<br><br>- Major deviations from stated diet | - Study uses human subjects outside of a highly controlled environment (for example an inpatient healthcare setting) | - Use of identical diet between study arms where diet is not the target of study |
| 1.5 GMB Normalization | - No documented means of verified GMB normalization methods employed prior to intervention<br><br>- Use of different normalization methods between study arms or use of non-validated technique | - Antibiotic normalization employed < 7 days prior to intervention | - Antibiotic normalization employed ≥7 days prior to intervention<br><br>- Validated technique of GMB normalization employed |
| **2 – Selection Bias** | | | |
| 2.1 Extreme Genotype | - Subjects of known extremely different genotypes<br><br>- Subjects with no established history of use | - Syngeneic subjects with limited established history of use | - Syngeneic subjects with established history of use |
| 2.2 Randomization or Demographic Balancing Sufficiently Applied | - Absence of both RCT and implementation of consistent host demographic across study | - Utilization of RCT or implementation of consistent host demographics across study | - Utilization of RCT and implementation of consistent host demographics across study |
| **3 - Classification of Intervention** | | | |
| 3.1 Intervention Bias | - Differential misclassification of intervention or test subject based on exposures present or suspected | - n/a | - Differential misclassification of intervention or test subject based on exposures absent or not suspected |
| 3.2 Validation of Method | - No validation that treatment method produces intended effect | - n/a | - Documented use of validated methods |

| | | | |
|---|---|---|---|
| | - Use of new method without internal validation | | - Use of a new method with adequate internal validation |
| **4 – Deviation from Intervention** | | | |
| 4.1 Deviation from Intervention | - Large deviations to protocol without adequate time stamps, rationale, and limitations noted | - Slight deviations to protocol with adequate time stamps, rationale, and limitations noted | - Intervention successfully carried out without protocol deviation |
| **5 - Missing Data** | | | |
| 5.1 Cause or Category of Missing Data | - Does not address missing data qualitatively or quantitatively<br><br>- Or, does not ensure to readers the absence of missing data | - Acknowledges missing data qualitatively or quantitatively<br><br>- Inadequate MAR/MNAR distinction or proper statistical correction | - Addresses missing data qualitatively or quantitatively, or ensures absence of missing data.<br><br>- Adequate MAR/MNAR distinction or proper statistical correction |
| 5.2 Subject Dropout | - Subject drop-out exceeds 20% | n/a | - Subject drop-out is equal to or less than 20% |
| 5.3 Sequencing Depth and Sampling Zeroes | - Does not address sampling zeroes with statistical correction | n/a | - Addresses sampling zeroes with statistical correction |
| **6 – Measurement of Outcomes** | | | |
| 6.1 Sample Collection | - Inconsistent collection time | - Animal models: Collected at same time, not in the morning<br><br>- Human models: Inconsistent collection time, but preserved immediately after defecation | - Animal models: Collected at same time, in the morning<br><br>- Human models: Consistent collection time & preserved immediately after defecation |

| 6.2 Blinding | - No double blinding when the primary measurement is subjective | n/a | - Primary outcome is objective measure such as genetic sequencing not subject to bias by the subject or investigator<br><br>- Primary outcome is subjective and double or greater blinding employed |
| **7 – Reporting of Results** | | | |
| 7.1 Selection of Reported Results | Any of:<br>- Omission of stated outcomes that are unfavorable or statistically insignificant<br><br>- Addition of outcomes not in initial protocol<br><br>- Results reported are only on a subset of data<br><br>- Changing outcome(s) of interest | - Any of the above, but with valid and satisfactory explanation provided | - Inclusion of relevant null and significant findings as stated in protocol |

679 **Table 1** – Rubric of domains and subdomains of bias with signaling statements to guide risk of bias assessment of gut microbiome studies.

680

681