

1
2

3 Large Quantities of Bacterial DNA and Protein in Common Dietary
4 Protein Source Used in Microbiome Studies

5

6

7 Authors:

8 Alexandria Bartlett^{1,2}, J. Alfredo Blakeley-Ruiz¹, Tanner Richie¹ Casey M. Theriot³, Manuel
9 Kleiner¹

10

11 Affiliations:

12 1: Department of Plant and Microbial Biology, North Carolina State University, Raleigh NC

13 2: Department of Molecular Genetics and Microbiology, Duke University, Durham, NC

14 3: Department of Population Health and Pathobiology, North Carolina State University, Raleigh,

15 NC

16

17 Current Address:

18 Tanner Richie, Division of Biology, Kansas State University Manhattan, KS

19

20 Correspondence:

21 Manuel Kleiner – [manuel_kleiner\(at\)ncsu\(dot\)edu](mailto:manuel_kleiner@ncsu.edu)

22 **Abstract**

23 Diet has been shown to greatly impact the intestinal microbiota. To understand the role of
24 individual dietary components, defined diets with purified components are frequently used in
25 diet-microbiota studies. Many of the frequently used defined diets use purified casein as the
26 protein source. Previous work indicated that this casein contains microbial DNA potentially
27 impacting results of microbiome studies. Other diet-based microbially derived molecules that
28 may impact microbiome measurements, such as proteins detected by metaproteomics, have not
29 been determined for casein. Additionally, other protein sources used in microbiome studies have
30 not been characterized for their microbial content. We used metagenomics and metaproteomics
31 to identify and quantify microbial DNA and protein in a casein-based defined diet to better
32 understand potential impacts on metagenomic and metaproteomic microbiome studies. We
33 further tested six additional defined diets with purified protein sources with an integrated
34 metagenomic-metaproteomic approach and show that contaminating microbial protein is unique
35 to casein within the tested set as microbial protein was not identified in diets with other protein
36 sources. We also illustrate the contribution of diet-derived microbial protein in diet-microbiota
37 studies by metaproteomic analysis of stool samples from germ-free mice (GF) and mice with a
38 conventional microbiota (CV) following consumption of diets with casein and non-casein
39 protein. This study highlights a potentially confounding factor in diet-microbiota studies that
40 must be considered through evaluation of the diet itself within a given study.

41

42 **Importance**

43 Many diets used in diet-microbiota studies use casein as the source of dietary protein. We found
44 large quantities of microbial DNA and protein in casein-based diets. This microbial DNA and

45 protein are resilient to digestion as it is present in fecal samples of mice consuming casein-based
46 diets. This contribution of diet-derived microbial DNA and protein to microbiota measurements
47 may influence results and conclusions and must therefore be considered in diet-microbiota
48 studies. We tested additional dietary protein sources and did not detect microbial DNA or
49 protein. Our findings highlight the necessity of evaluating diet samples in diet-microbiota studies
50 to ensure that potential microbial content of the diet can be accounted for in microbiome
51 measurements.

52

53 **OBSERVATION**

54 The intestinal microbiota is highly influential to the health of the host (1, 2). Diet has been
55 shown to shape the intestinal microbiota, yet we are still unraveling how individual dietary
56 components impact the functioning of the microbiota (1, 3, 4). To understand the role of
57 individual dietary components, many studies use purified dietary components, for example as a
58 supplement in human feeding trials (5, 6) or as part of a defined diet in animal studies (1, 7).
59 Diet-microbiota studies frequently use casein as the purified protein component of a defined diet,
60 as casein is the protein source in the defined AIN-93 laboratory rodent diet (8).

61

62 Previous studies have shown that sterilized, purified casein protein used in defined diets contains
63 significant amounts of microbial DNA. 16S rRNA gene sequencing identified the Gram-positive
64 bacterium *Lactococcus lactis*, which is used in casein production, as the source of the microbial
65 DNA (9-12). *Lactococcus*, and in some cases specifically *L. lactis*, have been identified by
66 sequencing-based studies as a key microbiota member responding to a specific treatment, for
67 example a high fat diet, in studies using casein (13, 14). While in some of these studies the

68 dietary source of *L. lactis* was recognized (9-11), other studies do not show any indication of
69 considering a potential dietary source of *Lactococcus*. In the studies that did recognize the
70 potential issue with *L. lactis* contamination from the diet, two methods for addressing the issue
71 have been used; bioinformatic removal of *L. lactis* reads during analysis (10, 12) or use of
72 ethanol-washed casein in diets (12).

73 While the detection of *L. lactis* DNA in studies with casein-based diets has highlighted the
74 importance of knowing the microbial content of dietary protein sources, our understanding of the
75 breadth of the issue is limited. First, we currently only know about the presence of microbial
76 DNA in casein, however, the presence of other microbially derived biomolecules such as protein
77 could also critically impact microbiota measurements. Specifically, metaproteomics, which is
78 used to study functional interactions in the microbiota by identifying and quantifying host and
79 microbial proteins (15), would be impacted by microbial proteins introduced through the diet.

80 The contribution of diet-derived microbial protein to metaproteomic measurements of the
81 microbiota has, however, not been previously investigated. Second, there is an increasing interest
82 in studying the effect of different sources of dietary protein on the microbiota and resulting
83 effects on host health (3). The microbial content of other dietary protein sources that are used in
84 diet-microbiota studies has not yet been investigated.

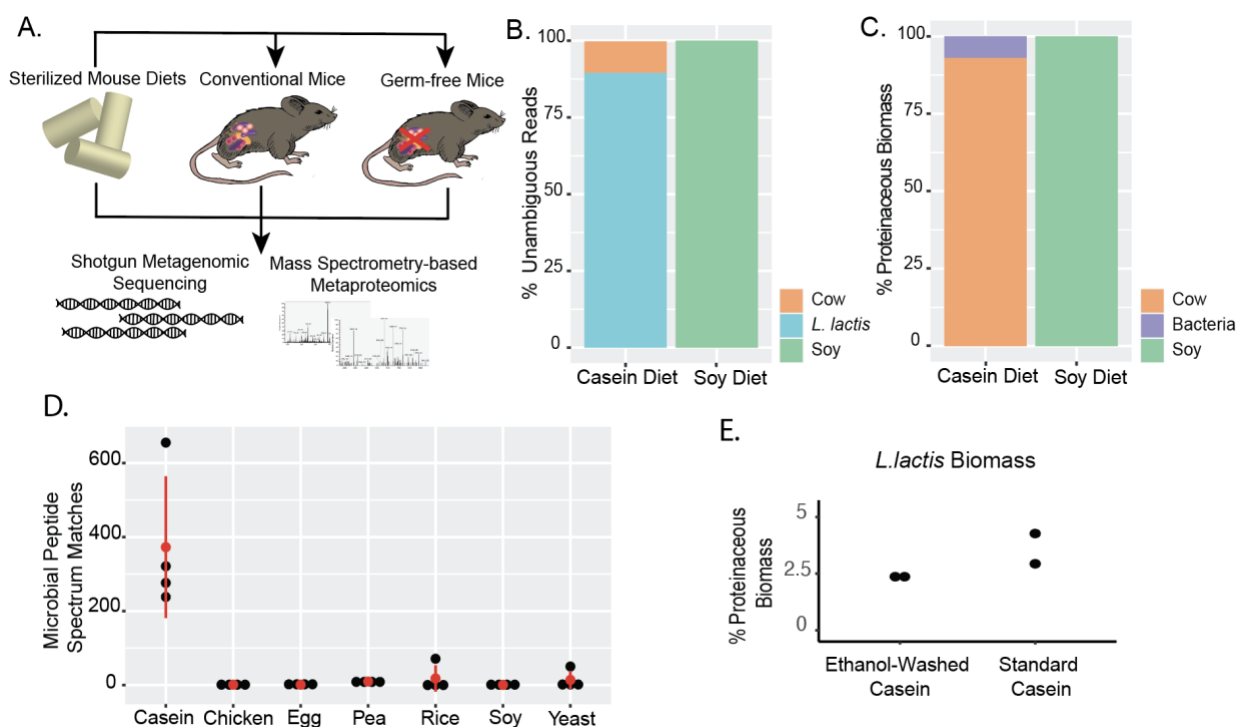
85

86 **Massive quantities of microbial DNA and protein in purified casein diet**

87

88 To assess the microbial protein content and its relevance in purified diets, we performed
89 metagenomic sequencing and metaproteomics of 1) defined mouse diets with a single source of
90 purified dietary protein (lactic casein or soy protein isolate), conventional (CV) and germ-free

91 (GF) mouse stool samples following consumption of the two diets (Fig. 1A). The metagenomic
92 sequencing allowed us to not only build a sample-matched protein sequence database for our
93 metaproteomic analyses (16), but also to expand on previous reports which primarily relied on
94 16S rRNA gene sequencing and only investigated casein-based diets and casein-fed mice. We
95 mapped the metagenomic reads to six reference genomes: *Mus musculus*, *Glycine max*, *Bos*
96 *taurus*, *Thioflaviccoccus mobilis* (control), *L. lactis* sequences retrieved from the assembled
97 metagenome and the assembled metagenome with *L. lactis* sequences removed. While the reads
98 from the soy diet mapped overwhelmingly to the soy genome, approximately 90% of reads from
99 the casein diet mapped to the *L. lactis* reference (Fig. 1B).



100
101
102 **Figure 1. Experimental design and quantification of bacterial DNA and protein in diets.** A)
103 Metagenomic sequencing and mass spectrometry-based metaproteomics was performed on
104 defined diets containing a single dietary protein source and on stool samples from germ-free
105 mice and mice with a conventional microbiota. The stool samples were collected from the mice
106 following consumption of diets containing purified soy protein or casein as the only protein
107 source. The diets were sterilized prior to measurement and feeding using gamma irradiation.
108 B) Average relative abundance of shotgun metagenomic reads from defined diet samples (n=3
109 per diet) that mapped unambiguously to diet reference genomes and assembled microbiota

110 metagenomes. C) Proteinaceous biomass composition of sterilized diets (n=4 per diet)
111 determined with metaproteomics according to the approach described by Kleiner et al (17).
112 Bacterial proteins were identified using protein sequences predicted from the assembled shotgun
113 metagenome to have a protein sequence database that matches the actual samples as described in
114 Blakeley-Ruiz and Kleiner (16). Approximately 4.5% of the bacterial proteinaceous biomass in
115 the casein diet samples was classified as *L. lactis* (Supplementary Table 3). The remaining 2.4%
116 of the bacterial proteinaceous biomass in the casein diet samples were unbinned or unambiguous
117 proteins in our microbiota database, which likely also represent *L.lactis* proteins. D) Number of
118 bacterial peptide spectrum matches (PSMs) in various defined diets with 20% of a purified
119 protein as the only protein source. The mean and standard deviation are indicated in red. E)
120 Bacterial PSMs in ethanol-washed casein vs unwashed casein.
121

122 To quantify the microbial protein in the purified diets, we identified and quantified proteins in
123 the soy (n=4) and casein (n=4) diets by LC-MS/MS and assessed the relative proteinaceous
124 biomass (17). Soy protein represented the entirety of the protein in the soy diet but microbial
125 proteins represented 4.7% of the proteinaceous biomass of the casein diet (Fig. 1C).

126

127 To investigate the microbial protein content of other dietary protein sources used in defined
128 diets, we performed metaproteomic analysis on seven different defined (20% protein by weight)
129 diets. In addition to the casein and soy diets described above, we analyzed diets formulated with
130 alternative purified protein sources including Egg White Solids, Torula Yeast, Chicken Bone
131 Broth, Yellow Pea, and Brown Rice (n=4 per diet). We observed more than 200 Peptide
132 Spectrum Matches (PSMs) to bacterial proteins for all replicates for only the casein diet (Fig.
133 1D). We did observe 71 bacterial PSMs for one replicate of the rice diet and 50 bacterial PSMs
134 for one replicate of the yeast diet, which we attributed to carryover from a sample with high
135 bacterial content run on the LC-MS/MS system immediately prior to these two samples.

136

137 A previous study showed that ethanol-washed casein contains 1,000-fold less *L. lactis* DNA and
138 suggested ethanol-washed casein as a potential alternative to the standard preparation of casein

139 currently used in purified diets (12). To assess if ethanol-washed casein has a similar reduction in
140 *L. lactis* protein, we assessed the proteinaceous biomass contribution of *L. lactis* in ethanol-
141 washed casein and standard casein. Although the *L. lactis* protein content in ethanol-washed
142 casein (2.4%) was on average lower than the standard casein (3.6%), the reduction was only
143 minor and not comparable to the previously reported reduction in *L. lactis* DNA (Fig. 1E). This
144 suggests that ethanol washing of casein is not a viable strategy for reducing *L. lactis* protein
145 content of casein-based diets.

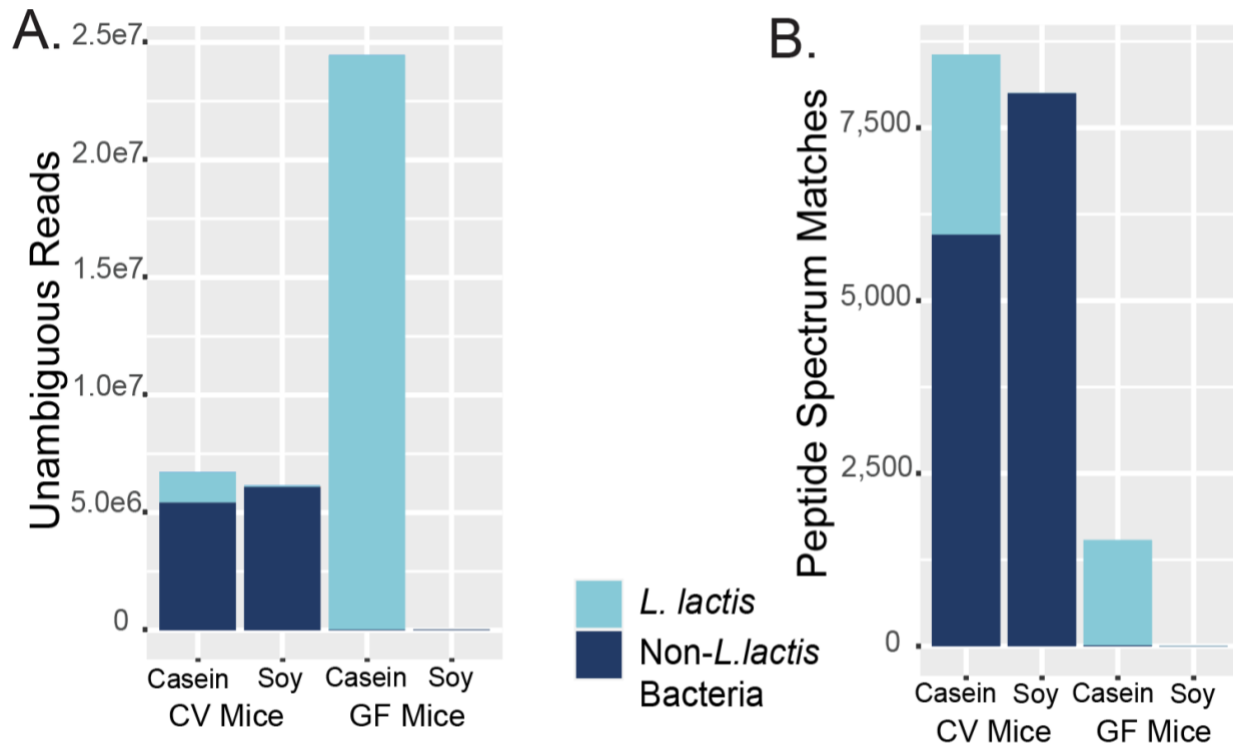
146

147 **Massive quantities of diet-derived microbial DNA and protein in stool of CV and GF mice** 148 **fed casein diet**

149

150 To examine the contribution of diet-derived microbial DNA and protein to metagenomic and
151 metaproteomic measurements of the microbiota, we performed shotgun metagenomics and
152 metaproteomics on stool collected from four cages of CV and GF mice after mice were fed the
153 casein or soy diet for 7 days (Fig. 1A). We mapped the raw metagenomic reads to the 6
154 references described above. Over 24 million reads from the casein-fed GF mice mapped to the *L.*
155 *lactis* reference (Fig. 2A). The majority of reads from the soy and casein CV mouse samples
156 mapped to the Non-*L.lactis* microbiota reference, however, the casein CV mouse samples had an
157 average of 1,291,589 reads that mapped to the *L. lactis* reference. In the metaproteomic analyses,
158 we measured an average of 1,527 *L. lactis* PSMs in the casein-fed GF mice (Fig. 2B). As
159 expected, the microbial PSMs we identified in the CV mice were primarily from Non-*L. Lactis*
160 microbiota proteins (5,713 PSMs in casein-fed and 7,999 PSMs in soy-fed mice). However, in
161 the casein-fed CV mice we measured on average, an additional 2,444 *L. Lactis* PSMs. Our data

162 suggests that a large quantity of diet-derived *L. lactis* DNA and protein withstands the passage
163 through the intestinal tract.



164
165 **Figure 2: Quantification of bacterial proteins and DNA in fecal samples from germ-free**
166 **(GF) mice and mice with a conventional microbiota (CV) following consumption of a defined**
167 **diet with either purified soy protein or casein as the sole protein source. A) Average number of**
168 **shotgun metagenomic reads from GF (n=3 per diet) and CV (n=3 per diet) mice fecal samples**
169 **that mapped unambiguously to diet reference genomes and assembled microbiota metagenomes.**
170 **B) Average microbial PSMs in GF (n=12) and CV mice (n=11-12).**
171

172 **Conclusions**

173
174 Here we show that one purified protein source (casein) used in diet-microbiota studies contains
175 high amounts of microbial DNA and protein. This diet-derived microbial DNA and protein
176 withstand passage through the intestinal tract and are thus present in metagenomic and
177 metaproteomic measurements of the microbiota. Additionally, we show that a diversity of other
178 purified dietary protein sources do not contain measurable amounts of microbial DNA and

179 protein, based on the metagenomic and metaproteomic approaches used in this study. It is,
180 however, to be expected that other purified dietary components (lipids, fiber, other protein
181 sources etc.) and unpurified foods and drinks contain microbial DNA and protein to various
182 extents, which may be present in microbiome measurements. In particular, foods and drinks that
183 involve microbial fermentation in their preparation such as, bread, cheese, yogurt, kimchi and
184 beer may contain significant amounts of microbial DNA (e.g. Fig. 4 in (4)) and protein, even if
185 the microbes from which the DNA and protein originate are non-viable in the intestinal tract.
186 Therefore, in any study that seeks to evaluate the impact of diet on the microbiota, the content of
187 microbial compounds in the diet that might confound the measurements should be evaluated
188 using the same measurement approach as the one applied to the intestinal or fecal samples. The
189 obtained information on the microbial content of dietary components can then be used to either
190 choose diets with a lower amount of microbial compounds or alternatively to bioinformatically
191 remove/account for the known microbial compounds from the diet in the data obtained from the
192 intestinal/fecal sample measurements.

193 **NCBI bioproject, ProteomeXchange Consortium and Dryad Data Repository.** Sequencing
194 data available under bioprojects PRJNA1026909 (microbial database) and PRJNA1026974
195 (metagenomic dataset). Mass spectrometry data and protein sequence databases available via
196 PRIDE repository under PXD041586 and PXD040649. Proteins identified at 5% FDR in diet and
197 mouse samples available via Dryad data repository with DOI: 10.5061/dryad.nvx0k6dzq. See
198 data accessibility in materials and methods for details.

199

200 **Materials and Methods**

201

202 ***Animals and housing***

203 12 conventional C57BL/J6 mice (6 males, 6 females, Jackson Labs Bar Harbor) and 12 germ-
204 free C57BL/J6 mice (6 males, 6 females, NCSU gnotobiotic core) were used in this study. All
205 mice were 3-6 months in age and housed in groups of two or three by sex. The food, bedding and
206 water were autoclaved. All conventional mouse cage changes were performed in a laminar flow
207 hood. The mice were subjected to a 12 h light and 12 h dark cycle and were housed at an average
208 temperature of 70F and 35% humidity. Animal experiments were conducted in the Laboratory
209 Animal Facilities located on the NCSU CVM campus. The animal facilities are managed by full-
210 time animal care staff coordinated by the Laboratory Animal Resources (LAR) division at
211 NCSU. The NCSU CVM is accredited by the Association for the Assessment and Accreditation
212 of Laboratory Animal Care International (AAALAC). Trained animal handlers in the facility fed
213 and assessed the status of animals several times per day. Those assessed as moribund were
214 humanely euthanized by CO₂ asphyxiation. This protocol is approved by NC State's Institutional
215 Animal Care and Use Committee (IACUC).

216 ***Animal diets and sample collection***

217 All diets used in this study were irradiated, not supplemented with amino acids and contained a
218 single source of dietary protein (Supplementary Table 1). The 20% soy diet (Envigo Teklad
219 Diets) was fed *ad libitum* to both conventional and germ-free mice for 7 days. After 7 days of the
220 soy diet, fecal samples were collected and the soy diet was replaced with a 20% casein diet.
221 After 7 days of the casein diet, fecal samples were collected again. Fecal samples were collected
222 into NAP buffer preservation solution at a ratio of approximately 1:10 Sample Weight:
223 Preservation Solution Volume, and roughly homogenized with a sterilized disposable pestle (18-

224 20). All animal protocols were approved by the Institutional Animal Care and Use Committee of
225 North Carolina State University.

226 ***Metagenomic sequencing***

227 Two different shotgun metagenomic datasets were generated in this study. The first round of
228 sequencing was done to prepare a matched metagenomic-based metaproteomic database for the
229 identification of microbiota proteins in the CV mice (PRJNA1026909). In the methods described
230 below, this dataset is referred to as the microbial database. The microbial database was used 1) to
231 identify the microbial proteins via metaproteomic analysis and 2) as a reference for the read
232 mapping of the second metagenomic dataset. The second round of metagenomic sequencing was
233 conducted to identify the DNA present in the diet samples, GF mice and CV mice
234 (PRJNA1026974). This dataset is referred to as the metagenomic dataset.

235 ***DNA extraction for the Microbial Database***

236 To generate a matched metagenomic-based proteomic database of the conventional mice
237 microbiota, DNA was extracted from 16 different CV mouse samples. These 16 samples
238 represented multiple sampling points from four different cages of conventional mice, collected 7
239 days after mice consumed different purified protein diets. The protocol described by Knudsen *et*
240 *al.*, which is based on the QIAamp DNA stool mini kit (Qiagen), was used to extract DNA, with
241 minor modifications (21). To remove the preservation solution from the fecal samples, 5 mL of
242 1X Phosphate Buffered Saline solution (VWR) was added to samples to dilute the preservation
243 solution, followed by centrifugation (17,000 x g, 5 min) to pellet solids and bacterial cells.
244 Pellets were lysed by beadbeating (3.1 m/s for 3 cycles of 30 sec. with 1 min of cooling on ice in
245 between each cycle) in 2 ml bead beating tubes (Lysing Matrix E, MP Biomedicals) using a
246 Bead Ruptor Elite 24 (Omni International). DNA concentration of eluates was assessed using a

247 DS-11 FX+ Spectrophotometer (Denovix) using a Qubit™ dsDNA High Sensitivity Assay Kit
248 (Invitrogen). 200 ng of each individually extracted sample was used to pool by cage (4 samples
249 sent for sequencing).

250 ***DNA extraction for the Metagenomic Dataset***

251 DNA was extracted from fecal samples of 3 GF mice and 3 CV mice (2 male, 1 female)
252 both after feeding on the soy diet and the casein diet (12 samples in total), 3 replicates of the
253 irradiated soy diet, 3 replicates of the casein diet and a *Thioflavicoccus mobilis* (*T. mobilis*)
254 culture as a control. We used the same protocol for DNA extraction as above.

255 ***DNA sequencing for the Microbial Database***

256 Metagenomic DNA was submitted to the North Carolina State Genomic Sciences
257 Laboratory (Raleigh, NC, USA) for Illumina library construction and sequencing to produce
258 between 51,152,549 and 74,618,259 150 bp paired-end reads for each of the 4 samples. Library
259 construction was performed using an Illumina TruSeq Nano Library kit according to
260 manufacturer's instructions. Libraries were sequenced on an Illumina NovaSeq 6000 sequencer.

261 ***DNA sequencing for the Metagenomic Dataset***

262 Metagenomic DNA was submitted to Diversigen for Illumina library preparation and
263 sequencing using a single lane of a NovaSeq to generate 100 bp single-end reads. To mitigate
264 index hopping, dual indexing was performed and a control sample was included (DNA extracted
265 from *T. mobilis*) for assessment between samples in a single lane. The number of reads that
266 unambiguously mapped to the *T. mobilis* reference was less than 0.25% of all unambiguously
267 mapped reads for diet and mouse samples (Supplementary Table 2).

268 ***Read processing, assembly, binning and annotation (Microbial Database)***

269 The BBSplit algorithm was used to remove phix174 (NCBI GenBank accession
270 CP004084.1) and mouse genome (mm10) reads, followed by quality trimming with BBDuk

271 (BBMap, Version 38.06) using the following settings: $mink = 6$, $minlength = 20$. MetaSPAdes
272 (version 3.12.0) with error correction and k-mer lengths 33, 55, 99 was used to assemble the 4
273 sequenced samples individually (22). In addition to assembling the samples individually, a co-
274 assembly of the four samples was performed using MEGAHIT (Version 1.2.4) to increase the
275 number of high-quality Metagenome Assembled Genomes (MAGs). MetaBAT (version 2.12.1)
276 was used to bin the assembled contigs and the resulting MAGs were evaluated using CheckM
277 (version 1.1.2) (23, 24). MAGs with a CheckM quality score of >50 completeness <10 were
278 considered medium quality and accepted for further consideration. MAGs with >30
279 completeness and <5 contamination were also included to avoid missing small genomes or
280 taxonomic groups that did not assemble well. dRep (Version 2.6.2) was used to cluster the
281 MAGs into species groups at 95% average nucleotide identity (25). PROKKA (Version 1.14.6)
282 was used for gene prediction of the MAGS and unbinned contigs. Taxonomy of the MAGs was
283 predicted using GTDB-Tk (Version 1.3.0) using reference database r95 and BAT(Version 5.0.3)
284 (26, 27).

285 *Protein sequence database construction for metaproteomics*

286 A non-redundant microbiota protein sequence database was constructed with the
287 annotated protein sequences from the microbial database to identify microbial proteins via
288 metaproteomic analysis. To remove redundant protein sequences, protein sequences from MAGs
289 were clustered with an identity threshold of 95% using cd-hit (Version 4.7) (28). Protein
290 sequences from unbinned contigs were separately clustered at 95% similarity. Cd-hit-2d was
291 used to identify sequences from unbinned contigs and low-quality MAGs that were not
292 represented in the set of binned sequences with at least 90% similarity. Sequences with less than
293 90% similarity to binned sequences were added to the microbiota protein sequence database.

294 The microbiota protein sequence database was combined with the *Mus musculus* reference
295 proteome (UP000000589, Downloaded 19Feb20) and one of the respective dietary proteomes to
296 generate six different databases. The dietary reference proteomes included *Glycine max*
297 (UP000008827, Downloaded 19Feb20), *Bos taurus* (UP000009136, Downloaded 19Feb20),
298 *Cyberlindnera jadinii* (UP000094389, Downloaded 25May20), *Oryza sativa* (UP000059680,
299 Downloaded 25May20) and *Gallus gallus* (UP000000539, Downloaded 25May20). Due to the
300 lack of a reference proteome for the yellow pea diet, we created a custom pea reference with all
301 available UniProtKB protein sequences for *Pisum sativum* (Taxon ID: 388 Downloaded
302 25Apr20) and the reference proteome of *Cajanus cajan* (UP000075243, Downloaded 25May20).
303 Each reference proteome and the *Pisum sativum* protein sequences were clustered individually
304 with an identity threshold of 95% using cd-hit (28).

305 ***Read mapping of metagenomic dataset***

306 The BBsplit algorithm was used to map raw reads from the metagenomic dataset to six
307 references. The references consisted of the *Mus musculus* genome (GCA_000001635.9), *Glycine*
308 *max* genome (GCA_000004515.4), *Bos taurus* genome (GCF_002263795.1), *T. mobilis*
309 (GCF_000327045.1), *L. lactis* and Non-*L. lactis* microbiota. The *L. lactis* reference comprised
310 the seven MAGs from our microbial dataset that were taxonomically classified as *L. lactis*. The
311 Non-*L. lactis* microbiota reference consisted of all other high-quality MAGs not classified as *L.*
312 *lactis* from our metagenomic database. The following BBsplit parameters were used:
313 ambiguous2=toss, qtrim=lr, minid=0.97. Plots were made using ggplot2 (Version 3.4.0) in
314 Rstudio (Version 4.1.1) (29, 30).

315 ***Protein extraction, peptide preparation and determination of diet and fecal samples***

316 We removed the NAP buffer preservation solution from fecal samples by centrifugation
317 (21,000 x g, 5 min). As diet samples were not collected in preservation solution, 100 mg was
318 placed directly into Lysing Matrix E tubes (MP Biomedicals) for each replicate. Cells were lysed
319 and proteins solubilized with SDT lysis buffer [4% (w/v) SDS, 100 mM Tris-HCl pH 7.6, 0.1 M
320 DTT] and bead beating in Lysing Matrix E tubes (MP Biomedicals) (5 cycles of 45s at 6.45 m/s,
321 1 min between cycles). Following bead beating, samples were heated to 95°C for 10 min., then
322 centrifuged (21,000 x g, 5 min). Tryptic digests were prepared (16 hour digestion) using the
323 filter-aided sample preparation protocol (31). In brief, 60 µl of lysate was combined with 400 µl
324 of UA solution (8 M urea in 0.1 M Tris/HCl pH 8.5) in 10 kDa MWCO 500 µl centrifugal filters
325 (VWR International). Samples were centrifuged at 14,000 g for 30 min. Depending on sample
326 concentration, this step was repeated up to three times to load the filter to capacity. Filters were
327 washed with 200 µl of UA solution at 14,000 g for 40 min. 100 µl IAA (0.05 M iodoacetamide in
328 UA solution) was added to filters, incubated for 20 min, and centrifuged at 14,000 g for 20 min.
329 Filters were washed with 100 µl of UA three times, followed by three washes of 100 µl ABC
330 (50 mM Ammonium Bicarbonate). For digestion, 0.95 µg of MS grade trypsin (Thermo
331 Scientific Pierce, Rockford, IL, USA) in 40 µl of ABC was added to filters and incubated for 16
332 hours in a wet chamber at 37 °C. Following digestion, samples were centrifuged at 14,000 g for
333 20 min. To elute peptides, 50 µl of 0.5 M NaCl was added and samples were centrifuged for 20
334 min. Peptide concentrations were determined with the Pierce Micro BCA assay (Thermo
335 Scientific Pierce) according to the manufacturer's instructions.

336 ***LC-MS/MS of diet and fecal samples***

337 We analyzed the peptides from fecal and diet samples by LC-MS/MS as previously
338 described with small modifications (18). The samples were blocked and randomized to control

339 for batch effects as previously described (32). An UltiMate™ 3000 RSLCnano Liquid
340 Chromatograph (Thermo Fisher Scientific) was used to load peptides (600 ng for mouse fecal
341 samples, 300 ng for diet samples) onto a 5 mm, 300 µm ID C18 Acclaim® PepMap100 pre-
342 column (Thermo Fisher Scientific) with loading solvent A (2% acetonitrile, 0.05% TFA).
343 Peptides were then separated on an EASY-Spray analytical column heated to 60°C (PepMap
344 RSLC C18, 2 µm material, 75 cm× 75 µm, Thermo Fisher Scientific) using a 140 min gradient at
345 a flow rate of 300 nl/min. The first 102 minutes of the gradient went from 95% eluent A (0.1%
346 formic acid) to 31% eluent B (0.1% formic acid, 80% acetonitrile), followed by 18 min from 31
347 to 50% B, and 20 min at 99% B. To reduce carryover, a wash run with 100% acetonitrile was
348 inserted between samples. Eluting peptides were ionized by electrospray ionization and analyzed
349 in a Q Exactive HF hybrid quadrupole-Orbitrap mass spectrometer (Thermo Fisher Scientific)
350 with the following parameters: m/z 445.12003 lock mass, normalized collision energy equal to
351 24, 25 s dynamic exclusion, and exclusion of ions of +1 charge state. A full MS scan from 380 to
352 1600 m/z was performed at a resolution of 60,000 and a max IT of 200 ms. Data-dependent MS²
353 was performed for the 15 most abundant ions at a resolution of 15,000 and max IT of 100 ms.

354 ***Protein identification and analysis of diet and fecal samples***

355 A protein sequence database containing the matched metagenomic database and multiple
356 reference proteomes was used to search the MS² spectra. The Proteome Discoverer software
357 version 2.3 (Thermo Fisher Scientific) was used for protein identification using run calibration
358 and the Sequest HT node, with the following settings: trypsin (Full), maximum 2 missed
359 cleavages, 10 ppm precursor mass tolerance, 0.1 Da fragment mass tolerance and maximum 3
360 equal dynamic modifications per peptide. The following dynamic modifications were
361 considered: oxidation on M (+15.995 Da), deamidation on N,Q,R (0.984 Da) and acetyl on the
362 protein N terminus (+42.011 Da). The static modification carbamidomethyl on C (+57.021 Da)

363 was also included. The percolator node in Proteome Discoverer was used to calculate peptide
364 false discovery rate (FDR) with the following parameters: maximum Delta Cn 0.05, a strict
365 target FDR of 0.01, a relaxed target FDR of 0.05 and validation based on q-value. The Protein-
366 FDR Validator node in Proteome Discoverer was used for protein inference with a strict target
367 FDR of 0.01 and a relaxed target FDR of 0.05 to restrict protein FDR to below 5%.

368 To assess the contribution of microbial proteins in the diet and mouse samples, identified
369 proteins were filtered for 5% FDR and 2 protein unique peptides and summed by organism, as
370 described in Kleiner et al (17). Only organisms with at least 10 PSMs for all 4 replicates of each
371 diet sample were included in the biomass assessment. Plots were made using ggplot2 (Version
372 3.4.0) in Rstudio (Version 4.1.1) (29, 30).

373 *Comparison of ethanol washed casein with standard lactic casein*

374 Envigo gifted us two lots of ethanol washed vitamin-free casein and standard casein used in their
375 purified protein diets. Proteins were extracted from 250 µg of sample and tryptic digests were
376 prepared using the filter-aided sample preparation protocol described above (31). LC-MS/MS
377 analysis was similar to the method described above for diet and fecal samples the only
378 modification being that 400 ng of peptides were loaded onto the analytical column. To identify
379 peptides and proteins, MS/MS spectra were searched against the same protein sequence database
380 as described above. Proteins were considered *Lactococcus lactis* if they 1) were present in the
381 *Lactococcus lactis* annotated MAGs from the matched metagenomic database or 2) matched to
382 the Uniprot *Lactococcus lactis* reference proteomes UP000002196 and UP000015854 at a 95%
383 identity threshold using diamond blastp (33). We then calculated percent proteinaceous biomass
384 as described in Kleiner et al (17).

385 *Data accessibility*

386 All sequencing data has been submitted to NCBI as part of bioprojects PRJNA1026909
387 (reviewer link:
388 <https://dataview.ncbi.nlm.nih.gov/object/PRJNA1026909?reviewer=a6ipv6iees1fq136cqt5rmq2q>
389 4) and PRJNA1026974 (reviewer link:
390 <https://dataview.ncbi.nlm.nih.gov/object/PRJNA1026974?reviewer=t94m1sut42r15j9himabksub>
391 [hp](#)). The mass spectrometry data and protein sequence databases were deposited to the
392 ProteomeXchange Consortium via the PRIDE (34) partner repository with the data set identifier
393 PXD041586 [reviewer access at this link <https://www.ebi.ac.uk/training/user/login>, with
394 credentials user:reviewer_pxd041586@ebi.ac.uk, password:V9Jz2n4h] and PXD040649
395 [reviewer access at this link <https://www.ebi.ac.uk/training/user/login>, with credentials
396 user:reviewer_pxd040649@ebi.ac.uk, password:ucnMkbYg]. All proteins identified at 5% FDR
397 in diet and mouse samples have been submitted to the Dryad data repository with DOI:
398 10.5061/dryad.nvx0k6dzq.

399 **Supplemental Material**

400 Supplementary Table 1: Composition of diets used in this study.
401 Supplementary Table 2: Raw read numbers and unambiguously mapped read numbers from
402 metagenomic dataset samples, including 3 replicates of soy diet, casein diet, GF and CV mice.
403

404 **3.7 Acknowledgements**

405 We are grateful to Fernanda Salvato for providing extensive training for the
406 metaproteomic analysis, Alissa Rivera for assistance with conventional mouse sample collection,
407 Karen Flores for technical assistance with manipulations of germ-free mice, Angie Mordant for
408 assistance with mouse experiments and Susan Tonkonogy for guidance on germ-free mouse
409 experiments. We thank Envigo for the kind gift of ethanol-washed casein samples. We thank the
410 Genomic Sciences Laboratory at North Carolina State University for performing Metagenomic

411 sequencing and the Bioinformatic and Statistical Consulting Cores at North Carolina State
412 University for performing statistical analysis and consultation. All LC-MS/MS measurements
413 were made in the Molecular Education, Technology, and Research Innovation Center (METRIC)
414 at North Carolina State University. The Gnotobiotic Core at the College of Veterinary Medicine,
415 North Carolina State University is supported by the National Institutes of Health funded Center
416 for Gastrointestinal Biology and Disease, NIH-NIDDK P30 DK034987. This work was
417 supported by the Foundation for Food and Agriculture Research Grant ID: 593607 and by the
418 National Institute Of General Medical Sciences of the National Institutes of Health under Award
419 Number R35GM138362.

420 **Declaration of Interests**

421 The authors declare no competing interests.

422 **References**

- 423
- 424 1. Desai MS, Seekatz AM, Koropatkin NM, Kamada N, Hickey CA, Wolter M, Pudlo NA,
425 Kitamoto S, Terrapon N, Muller A, Young VB, Henrissat B, Wilmes P, Stappenbeck TS,
426 Nunez G, Martens EC. 2016. A Dietary Fiber-Deprived Gut Microbiota Degrades the
427 Colonic Mucus Barrier and Enhances Pathogen Susceptibility. *Cell* 167:1339-1353 e21.
 - 428 2. Ma C, Han M, Heinrich B, Fu Q, Zhang Q, Sandhu M, Agdashian D, Terabe M,
429 Berzofsky JA, Fako V, Ritz T, Longerich T, Theriot CM, McCulloch JA, Roy S, Yuan
430 W, Thovarai V, Sen SK, Ruchirawat M, Korangy F, Wang XW, Trinchieri G, Greten TF.
431 2018. Gut microbiome-mediated bile acid metabolism regulates liver cancer via NKT
432 cells. *Science* 360.
 - 433 3. Bartlett A, Kleiner M. 2022. Dietary protein and the intestinal microbiota: An
434 understudied relationship. *iScience* 25:105313.
 - 435 4. David LA, Maurice CF, Carmody RN, Gootenberg DB, Button JE, Wolfe BE, Ling AV,
436 Devlin AS, Varma Y, Fischbach MA, Biddinger SB, Dutton RJ, Turnbaugh PJ. 2014.
437 Diet rapidly and reproducibly alters the human gut microbiome. *Nature* 505:559-63.
 - 438 5. Beaumont M, Portune KJ, Steuer N, Lan A, Cerrudo V, Audebert M, Dumont F,
439 Mancano G, Khodorova N, Andriamihaja M, Airinei G, Tome D, Benamouzig R, Davila
440 AM, Claus SP, Sanz Y, Blachier F. 2017. Quantity and source of dietary protein
441 influence metabolite production by gut microbiota and rectal mucosa gene expression: a
442 randomized, parallel, double-blind trial in overweight humans. *Am J Clin Nutr* 106:1005-
443 1019.
 - 444 6. Han ND, Cheng J, Delannoy-Bruno O, Webber D, Terrapon N, Henrissat B, Rodionov
445 DA, Arzamasov AA, Osterman AL, Hayashi DK, Meynier A, Vinoy S, Desai C, Marion

- 446 S, Barratt MJ, Heath AC, Gordon JI. 2022. Microbial liberation of N-methylserotonin
447 from orange fiber in gnotobiotic mice and humans. *Cell* 185:3056-3057.
- 448 7. Faith JJ, McNulty NP, Rey FE, Gordon JI. 2011. Predicting a human gut microbiota's
449 response to diet in gnotobiotic mice. *Science* 333:101-4.
- 450 8. Reeves PG, Nielsen FH, Fahey GC, Jr. 1993. AIN-93 purified diets for laboratory
451 rodents: final report of the American Institute of Nutrition ad hoc writing committee on
452 the reformulation of the AIN-76A rodent diet. *J Nutr* 123:1939-51.
- 453 9. Carmody RN, Gerber GK, Luevano JM, Jr., Gatti DM, Somes L, Svenson KL,
454 Turnbaugh PJ. 2015. Diet dominates host genotype in shaping the murine gut microbiota.
455 *Cell Host Microbe* 17:72-84.
- 456 10. Dalby MJ, Ross AW, Walker AW, Morgan PJ. 2017. Dietary Uncoupling of Gut
457 Microbiota and Energy Harvesting from Obesity and Glucose Tolerance in Mice. *Cell*
458 *Rep* 21:1521-1533.
- 459 11. Dollive S, Chen YY, Grunberg S, Bittinger K, Hoffmann C, Vandivier L, Cuff C, Lewis
460 JD, Wu GD, Bushman FD. 2013. Fungi of the murine gut: episodic variation and
461 proliferation during antibiotic treatment. *PLoS One* 8:e71806.
- 462 12. Bisanz JE, Upadhyay V, Turnbaugh JA, Ly K, Turnbaugh PJ. 2019. Meta-Analysis
463 Reveals Reproducible Gut Microbiome Alterations in Response to a High-Fat Diet. *Cell*
464 *Host Microbe* 26:265-272 e4.
- 465 13. Chen G, Chen D, Zhou W, Peng Y, Chen C, Shen W, Zeng X, Yuan Q. 2021.
466 Improvement of Metabolic Syndrome in High-Fat Diet-Induced Mice by Yeast beta-
467 Glucan Is Linked to Inhibited Proliferation of *Lactobacillus* and *Lactococcus* in Gut
468 Microbiota. *J Agric Food Chem* 69:7581-7592.
- 469 14. Zhao F, Song S, Ma Y, Xu X, Zhou G, Li C. 2019. A Short-Term Feeding of Dietary
470 Casein Increases Abundance of *Lactococcus lactis* and Upregulates Gene Expression
471 Involving Obesity Prevention in Cecum of Young Rats Compared With Dietary Chicken
472 Protein. *Front Microbiol* 10:2411.
- 473 15. Salvato F, Hettich RL, Kleiner M. 2021. Five key aspects of metaproteomics as a tool to
474 understand functional interactions in host-associated microbiomes. *PLoS Pathog*
475 17:e1009245.
- 476 16. Blakeley-Ruiz JA, Kleiner M. 2022. Considerations for constructing a protein sequence
477 database for metaproteomics. *Comput Struct Biotechnol J* 20:937-952.
- 478 17. Kleiner M, Thorson E, Sharp CE, Dong X, Liu D, Li C, Strous M. 2017. Assessing
479 species biomass contributions in microbial communities via metaproteomics. *Nat*
480 *Commun* 8:1558.
- 481 18. Mordant A, Kleiner M. 2021. Evaluation of Sample Preservation and Storage Methods
482 for Metaproteomics Analysis of Intestinal Microbiomes. *Microbiol Spectr* 9:e0187721.
- 483 19. Camacho-Sanchez M, Burraco P, Gomez-Mestre I, Leonard JA. 2013. Preservation of
484 RNA and DNA from mammal samples under field conditions. *Mol Ecol Resour* 13:663-
485 73.
- 486 20. Menke S, Gillingham MA, Wilhelm K, Sommer S. 2017. Home-Made Cost Effective
487 Preservation Buffer Is a Better Alternative to Commercial Preservation Methods for
488 Microbiome Research. *Front Microbiol* 8:102.
- 489 21. Knudsen BE, Bergmark L, Munk P, Lukjancenko O, Prieme A, Aarestrup FM, Pamp SJ.
490 2016. Impact of Sample Type and DNA Isolation Procedure on Genomic Inference of
491 Microbiome Composition. *mSystems* 1.

- 492 22. Nurk S, Meleshko D, Korobeynikov A, Pevzner PA. 2017. metaSPAdes: a new versatile
493 metagenomic assembler. *Genome Res* 27:824-834.
- 494 23. Kang DD, Li F, Kirton E, Thomas A, Egan R, An H, Wang Z. 2019. MetaBAT 2: an
495 adaptive binning algorithm for robust and efficient genome reconstruction from
496 metagenome assemblies. *PeerJ* 7:e7359.
- 497 24. Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. 2015. CheckM:
498 assessing the quality of microbial genomes recovered from isolates, single cells, and
499 metagenomes. *Genome Res* 25:1043-55.
- 500 25. Olm MR, Brown CT, Brooks B, Banfield JF. 2017. dRep: a tool for fast and accurate
501 genomic comparisons that enables improved genome recovery from metagenomes
502 through de-replication. *ISME J* 11:2864-2868.
- 503 26. von Meijenfeldt FAB, Arkhipova K, Cambuy DD, Coutinho FH, Dutilh BE. 2019.
504 Robust taxonomic classification of uncharted microbial sequences and bins with CAT
505 and BAT. *Genome Biol* 20:217.
- 506 27. Chaumeil PA, Mussig AJ, Hugenholtz P, Parks DH. 2019. GTDB-Tk: a toolkit to classify
507 genomes with the Genome Taxonomy Database. *Bioinformatics* 36:1925-7.
- 508 28. Li W, Godzik A. 2006. Cd-hit: a fast program for clustering and comparing large sets of
509 protein or nucleotide sequences. *Bioinformatics* 22:1658-9.
- 510 29. Wickham H, Wickham H. 2016. *Data analysis*. Springer.
- 511 30. Team R. 2020. RStudio: Integrated Development for R., p <http://www.rstudio.com/>.
512 Boston, MA.
- 513 31. Wisniewski JR, Zougman A, Nagaraj N, Mann M. 2009. Universal sample preparation
514 method for proteome analysis. *Nat Methods* 6:359-62.
- 515 32. Oberg AL, Vitek O. 2009. Statistical design of quantitative mass spectrometry-based
516 proteomic experiments. *J Proteome Res* 8:2144-56.
- 517 33. Buchfink B, Reuter K, Drost HG. 2021. Sensitive protein alignments at tree-of-life scale
518 using DIAMOND. *Nat Methods* 18:366-368.
- 519 34. Perez-Riverol Y, Bai J, Bandla C, Garcia-Seisdedos D, Hewapathirana S,
520 Kamatchinathan S, Kundu DJ, Prakash A, Frericks-Zipper A, Eisenacher M, Walzer M,
521 Wang S, Brazma A, Vizcaino JA. 2022. The PRIDE database resources in 2022: a hub
522 for mass spectrometry-based proteomics evidences. *Nucleic Acids Res* 50:D543-D552.
523