



HAL
open science

An Integrated Workflow for Enhanced Taxonomic and Functional Coverage of the Mouse Faecal Metaproteome

Nicolas Nalpas, Lesley Hoyles, Viktoria Anselm, Tariq Ganief, Laura Martinez-Gili, Cristina Grau, Irina Droste-Borel, Laetitia Davidovic, Xavier Altafaj, Marc-Emmanuel Dumas, et al.

► To cite this version:

Nicolas Nalpas, Lesley Hoyles, Viktoria Anselm, Tariq Ganief, Laura Martinez-Gili, et al.. An Integrated Workflow for Enhanced Taxonomic and Functional Coverage of the Mouse Faecal Metaproteome. 2020. hal-03090527

HAL Id: hal-03090527

<https://cnrs.hal.science/hal-03090527>

Preprint submitted on 29 Dec 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

An Integrated Workflow for Enhanced Taxonomic and Functional Coverage of the Mouse Faecal Metaproteome

Nicolas Nalpas¹, Lesley Hoyles^{2,3}, Viktoria Anselm¹, Tariq Ganief¹, Laura Martinez-Gili², Cristina Grau⁴,
Irina Droste-Borel¹, Laetitia Davidovic⁵, Xavier Altafaj^{4,6}, Marc-Emmanuel Dumas^{2,7,8}, Boris Macek¹

¹ Proteome Center Tuebingen, University of Tuebingen, Germany; ² Biomolecular Medicine Section, Division of systems Medicine, Department of Metabolism, Digestion and Reproduction, Imperial College London, Sir Alexander Fleming building, London SW7 2AZ, UK; ³ Department of Biosciences, Nottingham Trent University, UK; ⁴ Bellvitge Biomedical Research Institute, Spain; ⁵ Université Côte d'Azur, CNRS, Inserm, IPMC, France; ⁶ Neurophysiology Unit, University of Barcelona - IDIBAPS, Spain; ⁷ Genomic and Environmental Medicine, National Heart & Lung Institute, Faculty of Medicine, Imperial College London, London, SW3 6KY, United Kingdom. ⁸ European Genomic Institute for Diabetes, INSERM UMR 1283, CNRS UMR 8199, Institut Pasteur de Lille, Lille University Hospital, University of Lille, 59045 Lille, France.

Short title: Enhanced workflow for metaproteomics.

Keywords: Metaproteomics; Microbiome; Mus musculus; Mass spectrometry; Proteogenomics.

*To whom correspondence should be addressed:

Prof. Dr. Boris Macek
Proteome Center Tuebingen
Interfaculty Institute for Cell Biology
Auf der Morgenstelle 15
72076 Tuebingen
Germany
Phone: +49/(0)7071/29-70558
Fax: +49/(0)7071/29-5779
E-Mail: boris.macek@uni-tuebingen.de

39 Abstract

40 The intestinal microbiota plays a key role in shaping host homeostasis by regulating
41 metabolism, immune responses and behaviour. Its dysregulation has been associated with
42 metabolic, immune and neuropsychiatric disorders and is accompanied by changes in bacterial
43 metabolic regulation. Although proteomics is well suited for analysis of individual microbes,
44 metaproteomics of faecal samples is challenging due to the physical structure of the sample,
45 presence of contaminating host proteins and coexistence of hundreds of species. Furthermore,
46 there is a lack of consensus regarding preparation of faecal samples, as well as downstream
47 bioinformatic analyses following metaproteomic data acquisition. Here we assess sample
48 preparation and data analysis strategies applied to mouse faeces in a typical LC-MS/MS
49 metaproteomic experiment. We show that low speed centrifugation (LSC) of faecal samples
50 leads to high protein identification rates and a balanced taxonomic representation. During
51 database search, protein sequence databases derived from matched mouse faecal metagenomes
52 provided up to four times more MS/MS identifications compared to other database construction
53 strategies, while a two-step database search strategy led to accumulation of false positive
54 protein identifications. Comparison of matching metaproteome and metagenome data revealed
55 a positive correlation between protein and gene abundances, as well as significant overlap and
56 correlation in taxonomic representation. Notably, nearly all functional categories of detected
57 protein groups were differentially abundant in the metaproteome compared to what would be
58 expected from the metagenome, highlighting the need to perform metaproteomics when
59 studying complex microbiome samples.

60

61 Introduction

62 The prokaryotic component of the gut microbiota has multiple roles, contributing to
63 carbohydrate fermentation, maintenance of gut barrier integrity, as well as antimicrobial and
64 immunomodulation activities [1,2]. In metabolically healthy humans and mice, the gut
65 microbiota is predominated by two to three bacterial enterotypes [3-5]. These enterotypes
66 display significant heterogeneity in terms of species number, composition and relative
67 abundances depending on the location of the sample (upper vs lower GI tract) or the timing
68 (circadian variations) [6,7]. The gut microbiome has recently been associated in a number of
69 conditions ranging from inflammatory bowel syndrome to Parkinson's disease [8-11].
70 Interestingly, an increasing number of studies have reported a correlation between the gut
71 microbiome and neurodevelopmental disorders [12-14]. Notably, this includes changes in the
72 gut microbiome of Down syndrome individuals in comparison to non-trisomic individuals [15].
73 Given the established interaction between the host and the gut microbiome, a functional
74 analysis of the gut microbiome may help in understanding its contribution to pathophysiology.
75 To study the gut microbiome, approaches relying on nucleotide sequencing have so far been
76 preferred by the scientific community due to lower experimental costs, higher data throughput
77 and proven analytical workflows. However, metagenomics can only assess the genetic
78 potential, whereas metaproteomic investigates produced proteins (and therefore functions). In
79 particular, microbiome functional analysis can be performed using high-resolution mass
80 spectrometry (MS), to measure either protein abundance or metabolite production [16-18].
81 Although bacterial MS-based proteomic is well established, metaproteomic sample preparation
82 is hindered by many challenges, such as physical structure of the sample, the presence of
83 contaminating proteins and the coexistence of hundreds of microorganisms. Several studies in
84 human have shown that different sample preparation methodologies can result in significant
85 changes in the taxonomic composition and functional activities representation [19,20]. Beyond

86 sample preparation, the bioinformatic processing of metaproteomic data remains challenging,
87 due to the choice of representative protein sequence database, the redundancy in protein
88 functional annotation and elevated false discovery rate for peptide identification. Some of these
89 challenges have already been addressed by published software packages, such as
90 MetaProteomeAnalyzer [21] and MetaLab [22], which are all-in-one metaproteomic analytical
91 workflows, or UniPept [23], which allows peptide-based taxonomic representation.

92 However, a number of bottlenecks remain unaddressed, namely: (1) the lack of appropriate
93 sample preparation methods for optimal protein identification and taxonomic representation,
94 (2) high false positive rates in searches involving very large databases, (3) the impact of protein
95 sequence database constructions on protein identification, (4) the shortcomings of taxonomic
96 representation derived from MS-based peptide identification and (5) the lack of unbiased
97 assessment of the functional enrichment provided by the metaproteome compared to its
98 matching metagenome potential.

99 Here, we present a state-of-the-art LC-MS/MS-based workflow for the optimal metaproteome
100 characterisation of murine faecal samples. In terms of sample preparation, we achieve the
101 highest protein identification and biological correlation when combining LSC with in-solution
102 digestion. In terms of data processing, we demonstrate that the false discovery rate can only be
103 controlled using a single-step database search and that protein sequence database derived from
104 matching metagenome provides superior identification rate compared to publicly available
105 databases. We show that accurate taxonomic representation can be derived from peptide
106 spectral match (PSM) count. And while we overall observed a positive correlation between
107 protein and gene abundances, the metaproteome showed an enrichment in functionally-active
108 pathways compared to matching metagenomic potential.

109 Results

110 Low-speed centrifugation increases peptide identification rates

111 Our initial experiment involved the establishment of an optimal sample preparation workflow
112 applied to the mouse faecal metaproteome. In this context, we assessed several critical steps
113 within our sample preparation method: 1) the usage of LSC; 2) in-solution digestion versus
114 filter-aided sample preparation (FASP); and 3) cell lysis by bead-beating, nitrogen or
115 sonication (**Figure S1A**).

116 The number of peptides identified per MS raw file in the LSC group was significantly higher
117 with nearly 23 % more peptide identifications (**Figure 1A**). This was also observed at the
118 protein group level, but to a lower extent. Approximately 15 % of protein groups were
119 identified by a single peptide, while the median protein sequence coverage was 18.7 %. Such
120 metrics are usually indicative of highly complex samples that are not completely covered by a
121 single MS measurement under the stated parameters.

122 In-solution digestion consistently outperformed FASP based on peptide and protein groups
123 identification (**Figure 1B**). Compared to other methods, in-solution digestion combined with
124 LSC procedure provided nearly twice as many peptide identifications and 30 % more protein
125 groups. Furthermore, there was much less variability in the number of peptides and protein
126 groups identified across samples with this method.

127 When assessing the different lysis methods (in combination with LSC procedure), the methods
128 employed showed similar efficacies to identify peptides or protein groups (**Figure 1C**). This
129 suggests that any of these lysis methods could be used for preparation of mouse faecal samples.
130 However, the time necessary for the bead-beating approach was 30 min, while the nitrogen
131 grinding approach required 150 min and the sonication approach needed 90 min. The bead-

132 beating method would therefore provide a better time optimisation, particularly for large-scale
133 metaproteomic experiments.

134 LSC aids in recovery of Bacteroidetes proteins, whereas nLSC favours Firmicutes
135 proteins

136 Peptides that were identified after LSC and nLSC, were analysed to identify their phylogenetic
137 origin. The lowest common ancestor was determined using the Unipept interface [24], which
138 assigns peptide sequences to operational taxonomic units (OTUs). Based on the PSM count,
139 the most abundant superkingdom was bacteria, representing 77.7 % and 82.2 % of all PSMs
140 for LSC and nLSC, respectively (**Figure 1D-E, Table S1**). The second most represented
141 superkingdom was eukaryota in both the LSC and nLSC procedures. Under the assumption
142 that eukaryotic peptide sequences originated from the host, the proportion of *Mus musculus*
143 proteins was investigated further using intensity-based absolute quantification (iBAQ) values.
144 The LSC samples contained on average nearly two-fold more murine proteins (20.4 %) in
145 comparison to nLSC samples (14.6 %) (**Figure S1B**). Such findings were surprising since the
146 use of the LSC method was reported in a previous study to help with the removal of human
147 cells [19]. We also investigated the presence of peptides from host diet and found very low
148 levels of dietary peptides contamination (approximately 2 %), which was higher among LSC-
149 prepared samples (**Figure S1C**). This suggests that the majority of dietary proteins are absent
150 or depleted during the initial solubilisation step of the faecal pellet, a step common to both
151 procedures.

152 At the phylum level, three main taxa were represented in both LSC and nLSC: *Bacteroidetes*,
153 *Firmicutes* and *Metazoa*. There were large differences in the number of PSMs assigned to the
154 two main bacterial phyla when comparing LSC and nLSC methods. *Bacteroidetes* accounted
155 for 33 % and 13 % of PSMs, whereas Firmicutes amounted to 24 % and 48 % of PSMs in LSC
156 and nLSC procedures, respectively. The PSM counts were also calculated for each individual

157 sample across the two procedures in order to investigate the proportion of *Firmicutes* and
158 *Bacteroidetes*. We observed a significant inversion in the *Firmicutes*:*Bacteroidetes* ratio
159 between LSC and nLSC (**Figure S1D**). Interestingly, large sample-to-sample variation in the
160 *Firmicutes* to *Bacteroidetes* ratio was observed in the nLSC group. Such variation seemed
161 biological rather than technical, since technical replicates (samples measured multiple times by
162 MS) displayed relatively similar ratios. Overall, our results suggest different bacterial
163 representation, whereby LSC favours *Bacteroidetes* and corresponds more to genome-derived
164 taxonomy.

165 LSC and nLSC methods are characterised by different protein abundance profiles
166 We further investigated the overlap between the peptides or protein groups identified following
167 either LSC and nLSC procedures (**Figure 2A**). In terms of peptides, only 27.7 % were
168 identified with both procedures, the rest of the peptides being split equally into unique to LSC
169 and nLSC methods. Similar results were observed at the protein groups level with 38.7 % of
170 protein groups being identified in both procedures. Label-free quantitative (LFQ) comparison
171 between LSC and nLSC procedures revealed an intermediate correlation ($\rho = 0.44$) (**Figure**
172 **S2A**). Notably, the correlation between technical replicates was very high. This was illustrated
173 further through a principal component analysis (PCA) (**Figure 2B**). Clustering of technical
174 replicates confirmed the high technical reproducibility. Our findings indicate that while the two
175 procedures have a poor identification overlap and quantification correlation, the main
176 differences may still result from biological variations.

177 Using LFQ intensities, we then performed a *t*-test to identify which protein groups have
178 different abundances between the two procedures. Out of 2,589 quantified protein groups, 365
179 and 267 were significantly up-regulated and down-regulated in LSC versus nLSC samples,
180 respectively ($FDR \leq 0.01$ and absolute fold-change ≥ 2.5) (**Figure 2C, Table S2**). We gained
181 functional insights into these differences by performing an over-representation analysis of

182 KEGG pathways. The over-represented pathways based on the up- or down-regulated protein
183 groups were mostly similar ($FDR \leq 0.1$) and were associated with core microbial functions,
184 such as ribosome, carbon metabolism and carbon fixation pathways (**Figure 2D, Table S2**).
185 The protein groups unique to LSC or nLSC showed over-representation of protein export in
186 the LSC samples, whereas biosynthesis of amino acid, fatty acid degradation and bacterial
187 chemotaxis were over-represented in the nLSC samples (**Figure S2B**). Mapping of the
188 individual protein groups on the glycolysis-gluconeogenesis KEGG pathway highlighted the
189 redundancy in protein function (i.e. protein with identical function found across multiple
190 OTUs) and discrepancy in abundance (i.e. OTU depletion leads to reduced protein abundance,
191 while OTU enrichment increases protein abundance) (**Figure S2C**). Protein differential
192 abundance testing confirmed the divergence between LSC and nLSC procedures and was
193 suggestive of broad taxonomic changes, rather than variation in functional activities.

194 MS instrument selection critically determines the identification rate

195 Following assessment of sample preparation, we investigated the identification rates obtained
196 through LC-MS/MS measurements with two different mass spectrometers, namely Orbitrap
197 Elite and Q Exactive HF. In this context, we prepared samples using the LSC method from
198 faeces collected in a cohort of 38 mice. The newer generation Orbitrap instrument, namely Q
199 Exactive HF, provided a median of 229 MS/MS spectra identification per minute, while the
200 Orbitrap Elite resulted in less than half that number (**Figure S3A-B**). A similar significant trend
201 was also observed at the peptide and protein group levels, despite the fact that the measurement
202 time was halved on the Q Exactive HF versus the Orbitrap Elite instrument. Of note,
203 measurement on the Q Exactive HF required half of the sample material as compared to the
204 Orbitrap Elite (0.5 $\mu\text{g/h}$ vs 1 $\mu\text{g/h}$), a critical point when processing limited amount of material.
205 While expected, our findings highlight the crucial impact of mass spectrometer speed and

206 sensitivity in a typical metaproteomic measurements. Indeed, the choice of MS instrument was
207 among the parameters with the greatest impact on identification rates.

208 Two-step database search strategy shows a dramatic increase in false positive rate
209 After acquisition of LC-MS/MS raw data, the MS/MS spectra are searched against a protein
210 sequence database. One aspect of database search is the controversial use of a two-step search
211 strategy, whereby LC-MS/MS measurements are processed initially against a large protein
212 sequence database with no FDR control ($FDR \leq 1$). Subsequently, the original database is
213 filtered to retain only protein sequences that were identified during the first search. During the
214 second database search, the measurements are processed against the reduced database with
215 FDR control (e.g., $FDR \leq 0.01$) [25]. To assess the false discovery rates in such approach, we
216 searched a single HeLa cell LC-MS/MS file using MaxQuant software against a *Homo sapiens*
217 protein sequence database supplemented with different number of bacterial protein sequences
218 (**Figure S3C**). The HeLa measurement is used here as a proxy for a complex microbiome
219 measurement, with the exception that the sample composition is known (i.e. of human origin).
220 We initially established a gold standard by processing the HeLa measurement only against an
221 *H. sapiens* database, which resulted in approximately 5,000 human (eukaryota) protein groups
222 identified for the single-step search at $FDR \leq 0.01$ (**Figure 3A, Table S3**). Notably, the same
223 database used in a two-step search identified less than 1 % additional protein groups in
224 comparison to a single-step search, despite nearly twice as much processing time. We then
225 processed our HeLa measurement against *H. sapiens* database supplemented with 1:1, 1:2, 1:5,
226 1:10 and 1:20 *H. sapiens*:bacteria protein sequences, resulting in increasingly large databases
227 (**Figure S3C, Table S3**). For the single-step database search against the 1:20 database, we
228 observed a 10 % decline in the number of human protein groups identified, while 132 bacterial
229 protein groups were identified (false positives). On the contrary, the 1:20 two-step database
230 search resulted only in a 1 % decrease compared to the gold standard. This processing also

231 revealed a large number of bacterial protein groups identification (980 protein groups).
232 Furthermore, the two-step search led to large number of MS/MS spectra to be assigned to
233 different sequences (or newly assigned) in comparison to the gold standard (**Figure S3D, Table**
234 **S3**); this phenomenon was much less pronounced when performing the single-step search.
235 We then calculated the actual FDR for each processing approach using either the reverse hits
236 or the reverse hits plus the bacterial hits (which in our case are false positives). For both the
237 single-step and the two-step search, we obtained an FDR of 2.6 % when using only the reverse
238 hits for FDR calculation (**Figure 3B**). However, when using the reverse hits plus the bacterial
239 hits, we calculated an actual FDR of 8 % and 34 % for the single- and two-step search with
240 1:20 database, respectively. This represents a dramatic increase in the rate of false positive
241 identification when using two-step search, despite controlling for 1 % FDR. Notably, these
242 false positive hits would remain unnoticed in a microbiome sample of unknown composition,
243 thus highlighting the inherent problem associated with the two-step database search.

244 Metagenome-derived protein sequence database outperforms databases from public
245 resources

246 We then assessed the importance of the protein sequence databases and their impact on the
247 identification of microbiome sample proteins. In this context, we used faeces collected in a
248 cohort of 38 mice that were measured on Q Exactive HF instrument. We compared four
249 different databases, namely UniProt bacterial reference proteomes (5,408,622 protein entries),
250 UniProt bacterial pan proteomes (18,541,701 entries), protein sequences from the Mouse Gut
251 Metagenome catalogue (2,626,630 entries) [26] and protein sequences obtained from our
252 matched-metagenome data (1,595,268 entries) (**Figure 3C**). Notably, the number of protein
253 groups identified was inversely proportional to the number of protein sequences in the
254 databases, with the matched-metagenome database resulting in 25,230 identified protein
255 groups—that is seven-times more than the UniProt reference proteome database. A similar

256 trend could be observed at the peptide and MS/MS levels, with on average 23.23 % MS/MS
257 spectra identified using the matched-metagenome database. Not surprisingly, the matched-
258 metagenome database yielded the most identifications; however, the murine microbiome
259 catalogue also performed surprisingly well, making it an excellent substitute in cases where the
260 matched metagenome is not available.

261 Taxonomic representation correlates significantly between metaproteome and
262 metagenome down to species level

263 We also investigated the correlation in taxonomic representation between the metagenomic and
264 metaproteomic datasets. Because metaproteomics is not generally used as the method of choice
265 to infer taxonomy in a sample, there are fewer software tools available for this purpose in
266 comparison to metagenomic approaches. Initially, we retrieved taxonomic assignment using
267 Kraken2 software [27] and compared six methods to calculate taxon abundance using our
268 metaproteome data for our cohort of 38 mice (**Figure 4A, Table S4**). A range of filtering
269 thresholds were also implemented to remove low abundant taxa. The Kraken2 software in
270 combination with PSM count resulted in high OTU abundance correlation between
271 metagenome and metaproteome datasets even without any taxonomic representation filtering.
272 This correlation increased when a minimum threshold of 0.4 % total taxonomic representation
273 was used; however, it also reduced the number of identified OTUs by more than half. Aside
274 from the comparison to the metagenome, we investigated the taxonomic assignment obtained
275 using Diamond, Unipept and Kraken2 software tools [23,27,28]. This revealed far superior
276 results when using Kraken2, with between 20 and 30 % PSM assigned to a species, as opposed
277 to approximately 5 % with Diamond or Unipept (**Figure S4A**).

278 Using Kraken, we inspected the lowest taxonomic level that could be reached while still
279 displaying significant correlation against metagenomes (**Figure 4B**). Highly significant
280 correlations in taxonomic abundances were observed at all taxonomic levels. At the phylum

281 level, the ratio of *Firmicutes* to *Bacteroidetes* (or *Bacteroidota*) was significantly different
282 between metagenomes and metaproteomes, nevertheless ratios were consistent between the
283 two omics approaches and ranged between 0 and 1 (**Figure S4B**). When assessing the
284 identification of different phyla between omics datasets, we did not observe technical bias
285 relating to sample preparation or MS sensitivity (**Figure S4C**). Indeed, the low number of
286 identifications in metagenomic data did not necessarily translate into reduced or missing
287 identification in metaproteomic data; in addition, Gram-positive phyla were not necessarily
288 over-represented among the missing phyla in the metaproteomes. Under the current conditions,
289 metaproteomic can be used to assess taxonomic biomass [29] from phylum down to species
290 level for the mouse faecal microbiota.

291 Metaproteome to metagenome correlation highlights an over-representation in the
292 core microbiome functions

293 Due to the availability of matching metagenomic and metaproteomic data for our cohort of 38
294 mice, we assessed the correlation between gene and protein abundances. To deal with the
295 intrinsic difference between the two datasets, the gene entries were grouped in a similar fashion
296 as the protein groups (i.e. based on peptide identification) and the maximum expression was
297 calculated per gene group. Here, we show that a majority of gene-protein pairs (91 %) have a
298 positive correlation, with a median of 0.39, the rest having a median negative correlation of
299 -0.09 (**Figure 4C, Table S4**). Notably, 3,519 gene-protein pairs displayed a significant positive
300 correlation.

301 To identify the core pathways within our mice cohort, we performed an over-representation
302 analysis of the significantly correlated gene-protein pairs (**Figure 4D, Table S4**). Among these
303 pairs, there was an over-representation in carbon fixation, glycolysis-gluconeogenesis, citrate
304 cycle and carbon metabolism pathways (KEGG) [30]. We further characterised the correlating
305 genes and proteins based on gene ontology (GO) and identified 178 biological processes

306 (GOBP), 51 cellular components (GOCC) and 20 molecular functions (GOMF) that were over-
307 represented (**Figure S4D-F, Table S4**). Our results confirm the central role of carbon fixation
308 and general metabolism, which are associated with bacterial energy production, in the murine
309 faecal microbiome under the analysed conditions.

310 The metaproteome is enriched in functionally active pathways compared to the
311 matching potential encoded in the metagenome

312 The metagenome corresponds to the microbiome genetic potential, whereas the metaproteome
313 represents its truly expressed functional activities. Thereby, we compared the functional
314 abundance derived from the metagenomic versus metaproteomic datasets within our cohort of
315 38 mice. To allow comparison, the KEGG level 2 categories were quantified and normalised
316 separately for each omics datasets (**Figure S5A, Table S5**). Out of 55 KEGG categories, we
317 found 15 and 37 to be significantly increased and decreased in abundance at the metaproteome
318 level in comparison to the metagenome ($FDR \leq 0.05$). In general, the metagenome-based
319 quantification of KEGG categories was stable across categories, whereas large differences
320 were observed for the metaproteome.

321 To prioritise the KEGG categories, we selected eight categories differing significantly in terms
322 of gene-protein correlation in comparison to the overall correlation (**Figure 5A and S5B**).

323 Among the KEGG categories displaying higher abundance in the metaproteome compared to
324 the metagenome were the membrane transport, translation, signalling and cellular processes,
325 and genetic information processing. Conversely, transcription, carbohydrate metabolism and
326 antimicrobial drug resistance exhibited lower abundance. The KEGG Orthology (KO) entries
327 differing significantly in abundance between the metagenomes and metaproteomes were
328 identified via *t*-test and used for gene set enrichment analysis (GSEA). GSEA revealed an
329 enrichment of a number of overlapping KEGG pathways, with 19 and 6 pathways positively
330 and negatively enriched, respectively (**Figure 5B, Table S5**). Interestingly, we found the

331 ribosome pathway enriched in protein with increased abundance (between metaproteome and
332 metagenome datasets), therefore highlighting the functional activation of this pathway (**Figure**
333 **5C** and **S5C**). Conversely, homologous recombination, DNA replication and mismatch repair
334 were enriched in protein with decreased abundance, suggesting no or low activation of these
335 pathways. Overall, our findings highlight the critical importance of metaproteomic to
336 characterise microbiome samples particularly when it comes to their functional activity.

337 Discussion

338 Here, we provide solutions to some key bottlenecks hindering metaproteomic of murine faecal
339 samples in order to enhance protein identification, taxonomic and functional coverage. These
340 solutions include (1) an adequate sample preparation method, (2) the best strategy to control
341 for false positive rates, (3) the ideal protein sequence database construction, (4) an accurate
342 MS-derived taxonomic representation and (5) the leverage provided by metaproteomic to
343 determine functionally enriched pathways.

344 An integrated workflow that provide the highest identification rate for metaproteomic
345 of murine faecal samples and is amenable to other hosts

346 To the best of our knowledge this is one of the largest and most extensive comparisons
347 undertaken to date, comprising over 40 different biological samples and over 200 LC-MS/MS
348 runs. Overall, we reached identification rates that are similar to bacterial shotgun proteomics
349 (ca. 20-40 %). In comparison to previous murine faecal metaproteomic studies, we identified
350 more non-redundant peptides per samples (approximately 20,000 non-redundant peptides on a
351 60 min gradient) [31,32]. Several parameters may have influenced our greater performance,
352 among which are the use of a faster and more sensitive Orbitrap instrument (i.e. Q Exactive

353 HF) [33,34], an optimised LC gradient [35] and a more representative protein sequence
354 database (i.e. mouse metagenome catalogue or mouse matching metagenome) [26].

355 It should be noted that a number of aspects detailed herein would be directly applicable to
356 metaproteomic study of human samples. Indeed, many of our conclusions are not connected to
357 the taxonomic composition of species-specific samples and should therefore be transferrable
358 [36]. Regarding the sample preparation, while the choice of LSC was in part based on
359 taxonomic representation that is host-specific, the superior identification rate still encourages
360 its usage in other host organisms. Noteworthy, the selected approaches used to control false
361 discovery rate, construct protein sequence database and derive peptide-based taxonomic
362 representation will likely hold true in many diverse metaproteomic samples.

363 Increased identifications are obtained when using LSC with in-solution digestion

364 Our study confirms previous observation with regard to the depletion or enrichment of several
365 major bacterial phyla, which is dependent on laboratory preparation method and specifically
366 the usage of differential centrifugation [19]. In this context our results do not match with the
367 study from Tanca and colleagues, who reached opposite conclusions. However, there are
368 several possible explanations for such discrepancy, such as the host organism under study (*i.e.*
369 *Mus musculus* versus *Homo sapiens*) and different protein sequence database construction (*i.e.*
370 mouse microbiome catalogue versus UniProtKB custom microbiome). Nonetheless, to
371 preserve the trend in phylum distribution observed at the metagenome level in *Mus musculus*,
372 the LSC approach seems more appropriate based on previous studies, as well as our direct
373 comparison of the *Firmicutes* to *Bacteroidetes* proportion between metaproteomes and
374 matching metagenomes [37,38]. The LSC approach also leads to more consistent
375 identifications and as a result fewer missing values, which is a general and extensive problem
376 in metaproteomic datasets. Therefore, we recommend the use of LSC when preparing murine
377 faecal samples for measurement by mass spectrometry (Table 1).

378 While we observed a significant increase in identification when using in-solution digestion in
379 comparison to FASP, the lysis methods using either bead-beating, nitrogen grinding or
380 sonication did not impact performance. Our results partially confirm the study from Zhang and
381 colleagues regarding the superior performance of in-solution digestion [20] as opposed to other
382 protein digestion strategies [39], possibly due to limited protein loss.

383 Here, we identified numerous protein groups showing significant changes in abundance
384 between the LSC and nLSC approaches, as well as the over-representation of key KEGG
385 pathways. Notably, degradation of carbohydrates and proteins is over-represented in LSC,
386 while biosynthesis of amino acids is characteristic of the nLSC approach. While, these
387 observations show an opposite trend compared to the study by Tanca and colleagues [19],
388 possibly due to difference in host organism (*i.e.* *M. musculus* versus *H. sapiens*). Our results
389 are similarly indicative of broad taxonomic changes more so than variation in functional
390 activities.

391 Single-step search against matching-metagenome protein sequence database allows
392 control of false discovery rate and highest protein identification

393 Currently, many metaproteomic studies use two-step database searches as ways to boost
394 identification rates [25]. However, we demonstrate that this type of search dramatically
395 underrepresents the number of false positives, due to the use of a decoy search strategy that is
396 unsuitable in this context. Our results elaborate on a previous study by Muth and co-workers,
397 who also emphasised the drawbacks of using a two-step search together with decoy strategy
398 [40]. Here, our findings were so extreme that the number of false positives was equal or greater
399 to the number of false negatives, with FDR outside of the accepted range (*i.e.* $FDR > 0.1$). We
400 argue that the use of a two-step search should be avoided whenever possible and replaced by
401 alternative strategies, such as taxonomic foreknowledge or using matching metagenomes
402 (Table 1) [41].

403 We also show that the choice of protein sequence database had a serious impact on the
404 identification rates, leading to nearly five-fold differences. While the best results were obtained
405 using protein sequences derived from matched metagenomes, it was surprising that
406 concatenation of bacterial protein sequences from online resources performed so poorly. This
407 highlights the large diversity in proteins and peptides, which are not accounted among well
408 characterised bacterial species, such as sequences from UniProtKB Reference or Pan
409 proteomes. However, our results also reveal the importance of microbiome characterisation
410 studies performed in different organisms (e.g. *H. sapiens*, *M. musculus*) or in specific tissues
411 (e.g. oral, nasal) [26,42,43]. Indeed, the *M. musculus* microbiome catalogue led to
412 identification rates similar to our matched metagenomes without the associated costs. In the
413 future, more comprehensive microbiome catalogues may completely alleviate the need for
414 matched metagenome (Table 1).

415 Accurate taxonomic annotation and quantification are obtained via Kraken2 software
416 and PSM count

417 Based on the approaches tested here and in other studies, it is possible to derive taxonomic
418 representation and abundance from MS-based peptide identification [29,44]. In this study,
419 taxonomic representation correlated significantly between metaproteomic and metagenomic
420 data at all taxonomic levels tested (*i.e.* phylum down to species), thus confirming observations
421 from Erickson and colleagues [45]. Interestingly, the metaproteome PSM count is the
422 calculation method that is closest to metagenome read count, it is therefore unsurprising that
423 this calculation method showed the best correlation for taxonomic representation between the
424 two omics datasets. While the Kraken2 software [27] provided the best taxonomic annotation
425 (Table 1), it should be mentioned that the comparison to Diamond or Unipept softwares [23,28]
426 is not entirely fair, since the metagenome taxonomy was solely derived from Kraken2 and may
427 thus result in a favouring bias. Importantly, the proportion of *Firmicutes* and *Bacteroidetes*

428 displayed a similar trend (despite a significant difference) between omics, indicating that
429 overall murine faecal microbiota contains more *Bacteroidetes* than *Firmicutes* [37,38].

430 While, the majority of OTUs were commonly identified across omics datasets (especially at
431 higher taxonomic levels), there were several OTUs quantified exclusively at the metaproteome
432 or metagenome level, e.g. *Firmicutes_H*, *Euryarchaeota*, *Thermoplasmata*. Some of these
433 discrepancies might be explained by different analytical artefacts (i.e. different sensitivity
434 between Thermo Orbitrap versus Illumina HiSeq instruments). Yet it is important to state that
435 bacterial activity and bacterial presence are different, therefore it is unsurprising to report only
436 a medium overlap between metaproteome and metagenome [29].

437 The metaproteome shows an enrichment in functionally-active pathways compared to
438 the matching metagenomic potential

439 Here, we observed an overall positive correlation between gene and protein abundances derived
440 from metaproteome and matching-metagenome analysis. This was previously reported in a
441 longitudinal study of metaproteome/metagenome fluctuations from one individual with
442 Crohn's Disease [46]. In our case the significantly correlated entries were associated with core
443 bacterial metabolic functions, such as carbon and energy metabolism or electron transfer
444 activity [47]. Despite such correlations, we also reported extensive differences in quantified
445 functions between metagenomic and metaproteomic. Notably, with regard to genetic
446 information processing (KEGG level 2), the ribosome pathway was over-represented in entries
447 with higher abundance in metaproteomes, whereas pathways associated with DNA repair,
448 replication or recombination were over-represented in entries with increased abundance in
449 metagenomes. This greatly highlights the main advantage of metaproteomic, which capture
450 functionally active pathways, as opposed to the genetic potential represented by metagenomic
451 [48]. Thus, these approaches are complementary to each other and can provide a more
452 comprehensive understanding of a biological system.

453 Conclusion

454 To conclude, in this study we present an integrated analytical and bioinformatic workflow to
455 improve protein identification, taxonomic and functional coverage of the murine faecal
456 metaproteome. Notably, this workflow should be easily amenable to studying the human faecal
457 metaproteome. LSC combined with in-solution digestion provided the highest identification
458 rates. We also show that fast and accurate MS data processing can be achieved using a single-
459 step database search against publicly available metagenome catalogues or matching
460 metagenomes. Taxonomic representation can be generated directly from MS-based peptide
461 identification. While protein and gene abundances show an overall positive correlation, the
462 metaproteome showed a significant functional enrichment compared to its metagenomic
463 potential; thus, emphasizing the need for more metaproteomic studies for adequate functional
464 characterisation of the microbiome.

465 Methods

466 Animals and faecal samples collection

467 Mouse faecal pellets obtained from a small cohort of six wild-type B6EiC3SnF1/J mice were
468 used to compare sample purification and protein extraction methodologies (**Figure S1A**). A
469 larger cohort of 38 mice (euploid and trisomic Ts65Dn) was used to obtain mouse faeces, for
470 further sample preparation described below (**Figure S1A**). Mice were housed and faeces were
471 collected following the experimental procedures evaluated by the local Ethical Committee
472 (Barcelona Biomedical Research Park, Spain). After collection, faecal pellets were frozen and
473 stored at -80 °C until analysis.

474 DNA extraction and whole-genome sequencing

475 Whole genome analysis was performed on the mouse cohort used for data analysis assessment.

476 In brief, DNA was extracted from faecal samples using the FastDNA SPIN Kit (MP

477 Biochemicals) and following manufacturer's instructions. DNA concentration was measured

478 using a Qubit fluorometer (Invitrogen) and samples were shipped frozen to the Quantitative

479 Biology Centre (QBiC) at the University of Tuebingen for whole genome sequencing.

480 Sequence data were generated on an Illumina HiSeq 2500 instrument (chemistry SBS v3 plus

481 ClusterKit cBot HS) and processed as described previously [49] but with minor modifications

482 that follow. Supplied sequence data were checked using fastQC v0.11.5 [50]. Data were

483 trimmed with Trim Galore! (--clip_R1 10 --clip_R2 10 --three_prime_clip_R1 10 --

484 three_prime_clip_R2 10 --length 50; Babraham Bioinformatics). Mouse DNA within samples

485 was detected by mapping reads against the mouse genome (GRCm38). Mouse-filtered read

486 files (with an average of 3.58 ± 0.08 Gb sequence data per sample) were used for all subsequent

487 analyses. Kraken2 2.0.8-beta [51] with the pre-compiled Genome Taxonomy Database [52]

488 Kraken2 GTDB_r89_54k index (downloaded on 3 May 2020) available from

489 <https://bridges.monash.edu/ndownloader/files/16378439> [53] was used to determine the

490 bacterial and archaeal taxonomic composition/abundance for each sample. Functional

491 annotation was achieved by mapping centroid protein sequences generated as described before

492 [49,51] using the eggNOG-mapper software (v.1.0.3) [54] and associated database (v.4.5).

493 Microbial gene richness was determined as previously described [49]. Data were downsized to

494 adjust for sequencing depth and technical variability by randomly selecting 20 million reads

495 mapped to the merged gene catalogue (of 1,540,712 genes) for each sample and then computing

496 the mean number of genes over 30 random drawings.

497 Sample treatment before protein extraction

498 Mouse faecal pellets obtained from wild-type B6EiC3SnF1/J mice were used to compare
499 sample initial preparation methodologies. For the LSC procedure, faeces (~50 mg) were
500 resuspended in phosphate buffer (50 mM Na₂HPO₄/NaH₂PO₄, pH 8.0, 0.1 % Tween 20, 35x
501 volume per mg) by vortexing vigorously for 5 min using 4 mm glass beads (ColiRollers™
502 Plating beads, Novagen), followed by incubation in a sonication bath for 10 min and shaking
503 at 1,200 rpm for 10 min in a Thermomixer with a thermo block for reaction tubes. Insoluble
504 material was removed by centrifugation at 200 × g at 4 °C for 15 min. The supernatant was
505 removed and the remaining pellet was subjected to two additional rounds of microbial cell
506 extraction. After merging supernatants, microbial cells were collected by centrifugation at
507 13,000 × g at 4 °C for 30 min. The pellet was resuspended in 80 µL sodium dodecyl sulfate
508 (SDS) buffer (2 % SDS, 20 mM Tris, pH 7.5; namely pellet extraction buffer) and heated at
509 95 °C for 30 min in a Thermomixer.

510 For the nLSC procedure, mouse faeces (~25 mg) were homogenised in 150 µL pellet extraction
511 buffer as described above with the following changes. A bead mixture of 0.1 mm glass beads
512 (100 mg), 5 × 1.4 mm ceramic beads (Biolab products), and 1 × 4 mm glass bead was used for
513 five cycles of homogenisation.

514 Cell lysis and protein extraction

515 To compare protein extraction by bead beating, two bead beaters were used to disrupt bacterial
516 cell pellets derived from LSC or nLSC procedures. Samples were split in two and were
517 homogenised using 0.1 mm glass beads (100 mg, Sartorius™ Glass Beads) together with the
518 FastPrep-24 5G instrument (MP) at 4 m/s or the BeadBug microtube homogeniser (BeadBug)
519 at 4,000 rpm. Three homogenisation cycles were performed and consisted of 1 min bead
520 beating, 30 sec incubation at 95 °C and 30 sec centrifugation at 13,000 × g. The homogenate
521 was diluted with 800 µL MgCl₂ buffer (0.1 mg/mL MgCl₂, 50 mM Tris, pH 7.5) and

522 centrifuged at 13,000 rpm for 15 min. Proteins from the supernatant were precipitated
523 overnight in acetone/methanol at -20 °C (acetone:methanol:sample with 8:1:1 ratio). Protein
524 pellets were resuspended in 120 µL denaturation buffer (6 M urea, 2 M thiourea, 10 mM Tris,
525 pH 8.0) for downstream use.

526 Additional protein extraction methods were compared only using three biological samples that
527 were subjected to the LSC procedure (as described above). Microbial pellets were prepared in
528 order to compare (1) bead beating, (2) ultrasonication and (3) grinding on liquid
529 nitrogen/ultrasonication. Bead beating was performed as described above by resuspension in
530 100 µL pellet extraction buffer and by using the MP bead beater. The microbial pellets for
531 ultrasonication were resuspended in 120 µL pellet extraction buffer and incubated at 60 °C and
532 1,400 rpm for 10 min. After addition of 1 mL pellet extraction buffer, benzonase was added
533 for DNA removal (final concentration of 1 µL/mL). Ultrasonication was performed using an
534 ultrasonicator with an amplitude of 50 % and a cycle time of 0.5 for 2 min on ice. The samples
535 were incubated at 37 °C and 1,400 rpm for 10 min, followed by centrifugation at 4 °C and
536 10,000 × g for 15 min. The third procedure included grinding of the microbial pellet in liquid
537 nitrogen for 1-2 min with a pestle in reaction tube, followed by ultrasonication as described
538 above.

539 Protein digestion

540 Following extraction, protein amount was quantified using Bradford assay (Bio-Rad, Munich,
541 Germany) [55] and two methods were compared to digest proteins extracted from LSC or nLSC
542 procedures.

543 The in-solution digestion method was performed as follows. Proteins (20 µg starting material)
544 were reduced in 1 mM dithiothreitol (DTT) and alkylated in 5.5 mM iodoacetamide at room
545 temperature (RT) for 1 h each. Proteins were pre-digested with LysC at RT for 3 h using a
546 protein to protease ratio of 75:1. Samples were diluted nine-fold with 50 mM ammonium

547 bicarbonate and digested overnight with trypsin (Sequencing Grade Modified Trypsin,
548 Promega) at pH 8.0 using a protein to protease ratio of 75:1.
549 Filter-aided sample preparation (FASP) was performed as previously published [56]. Briefly,
550 proteins (10 µg starting material) were reduced in 0.1 M DTT for 40 min at RT. The reduced
551 samples were added to the filter units (30 kDa membrane cut off) and centrifuged at 14,000 × g
552 for 15 min. All further centrifugation steps were performed similarly unless otherwise noted.
553 Samples were then washed with 2X 200 µL urea buffer (100mM Tris/HCl, pH 8.5, 8M urea)
554 and centrifuged. Proteins were incubated in 50 mM IAA for 20 min at RT in the dark. After
555 alkylation, samples were centrifuged and washed three times with 100 µL urea buffer. This was
556 followed by three wash steps with 50 mM ammonium bicarbonate (ABC) for 10 min. Proteins
557 were digested overnight at 37 °C using trypsin digestion (Sequencing Grade Modified Trypsin,
558 Promega) at pH 8.0 using a protein to protease ratio of 100:1. On the following day, the peptides
559 were centrifuged into fresh tubes at 14,000 × g for 10 min. An additional 40 µL ABC buffer
560 was added to the filter units and this solution was also centrifuged to increase the peptide yield.
561 Following digestion either in-solution or FASP, samples were acidified to pH 2.5 with formic
562 acid and cleaned for LC-MS/MS measurement using Empore C18 disks in StageTips [57].

563 LC-MS/MS measurements

564 Samples were measured on an EASY-nLC 1200 (Thermo Fisher Scientific) coupled to a Q
565 Exactive HF mass spectrometer (Thermo Fisher Scientific). Peptides were
566 chromatographically separated using 75 µm (ID), 20 cm packed in-house with reversed-phase
567 ReproSil-Pur 120 C18-AQ 1.9 µm resin (Dr. Maisch GmbH).
568 Peptide samples generated as part of the laboratory method optimisation were eluted over
569 43 min using a 10 to 33 % gradient of solvent B (80 % ACN in 0.1 % formic acid) followed
570 by a washout procedure. Peptide samples generated as part of the data analysis assessment were

571 eluted over 113 min using a 10 to 33 % gradient of solvent B (80 % ACN in 0.1 % formic acid)
572 followed by a washout procedure.

573 MS1 spectra were acquired between 300-1,650 Thompson at a resolution of 60,000 with an
574 AGC target of 3×10^6 within 25 ms. Using a dynamic exclusion window of 30 sec, the top 12
575 most intense ions were selected for HCD fragmentation with an NCE of 27. MS2 spectra were
576 acquired at a resolution of 30,000 and a minimum AGC of 4.5×10^3 within 45 ms.

577 LC-MS/MS data processing

578 Raw data obtained from the instrument were processed using MaxQuant (version 1.5.2.8) [58].
579 The protein sequence databases used for database search consisted of the complete *Mus*
580 *musculus* Uniprot database (54,506 sequences) and frequently observed contaminants (248
581 entries), as well as the mouse microbiome catalogue (~2.6 million proteins) [26] for the raw
582 data from laboratory method optimisation samples or the matching metagenome gene
583 translation (~1.5 million proteins) for the raw data from data analysis assessment samples. A
584 FDR of 1 % was required at the peptide and protein levels. A maximum of two missed
585 cleavages was allowed and full tryptic enzyme specificity was required. Carbamidomethylation
586 of cysteines was defined as fixed modification, while methionine oxidation and N-terminal
587 acetylation were set as variable modifications. Match between runs was enabled where
588 applicable. Quantification was performed using label-free quantification (LFQ) [59] and a
589 minimum ratio of 1. All other parameters were left to MaxQuant default settings.

590 Comparison of sample preparation methods

591 Unless stated otherwise, the analyses described below were performed in the R environment
592 [60]. To compare the different centrifugation, digestion and lysis methods, we counted for each
593 sample the number of peptide and protein groups with intensities and LFQ intensities superior
594 to zero, respectively. We tested for significant differences between methods using paired t-tests

595 via the ggplot2 package [61]. Quantified peptides and protein groups were checked for overlap
596 between the centrifugation methods using the VennDiagram package. The proportion of host
597 (*Mus musculus*) proteins was computed by summing up all host proteins iBAQ values and then
598 dividing by the total iBAQ per sample. The centrifugation methods were evaluated using a
599 paired t-test.

600 The taxonomy representation for the centrifugation methods was done via the Unipept software
601 [24]. The quantified peptides (intensity superior to zero) were imported into Unipept with I-L
602 not equal and advanced missed cleavages handling selected only for peptides with 1 or 2 miss-
603 cleavages. The Unipept result were used to count the number of non-redundant peptides
604 assigned to each taxonomic node. The Firmicutes to Bacteroidetes ratio was calculated by
605 summing the spectral count for each phylum and sample. The centrifugation methods were
606 compared on this basis using a paired t-test.

607 For the differential protein abundance analysis (between LSC and nLSC), the MSnBase
608 package was used as organisational framework for the protein groups LFQ data [62]. Host
609 proteins, reverse hit and potential contaminant proteins were filtered out. Protein groups were
610 retained for further analysis only if more than 90 % of samples within either LSC or nLSC
611 group had an LFQ superior to the first quartile overall LFQ. Significantly changing proteins
612 were identified using paired t-test. Significance was set at an adjusted p-value of 0.01 following
613 Benjamini-Hochberg multiple correction testing, as well as a minimum LSC/nLSC fold-change
614 of ± 1.5 . The over-representation and GSEA testing of KEGG pathways were done for the
615 significantly up- and down-regulated proteins as well as for the proteins uniquely identified per
616 group via the clusterProfiler package based on hypergeometric distribution ($p\text{-adj.} \leq 0.05$) [63].
617 Selected over-represented KEGG pathways were displayed in context of protein quantification
618 using the pathview package.

619 MS instruments comparison

620 The samples generated as part of the data analysis assessment were measured on a Q Exactive
621 HF mass spectrometer as described above, as well as on an Orbitrap Elite mass spectrometer,
622 as described below.

623 Samples were measured on an EASY-nLC 1200 (Thermo Fisher Scientific) coupled to an
624 Orbitrap Elite mass spectrometer (Thermo Fisher Scientific). Peptides were
625 chromatographically separated using 75 μm (ID), 20 cm packed in-house with reversed-phase
626 ReproSil-Pur 120 C18-AQ 1.9 μm resin (Dr. Maisch GmbH).

627 Peptide samples were eluted over 213 min using a 10 to 33 % gradient of solvent B (80 % ACN
628 in 0.1 % formic acid) followed by a washout procedure. An optimised gradient was also tested
629 whereby peptide samples were eluted over 225 min using a 5 % (0 min), 10 % (5 min), 13 %
630 (80 min), 15 % (170 min) and 30 % (225 min) gradient of solvent B (80 % ACN in 0.1 %
631 formic acid) followed by a washout procedure.

632 MS1 spectra were acquired between 300-2,000 Thomson at a resolution of 120,000 with an
633 AGC target of 1×10^6 within 100 ms. Using a dynamic exclusion window of 30 sec, the top 15
634 most intense ions were selected for HCD fragmentation with an NCE of 35. MS2 spectra were
635 acquired at a resolution of 120,000 and a minimum AGC of 5×10^3 within 150 ms.

636 Raw data were processed, together with raw data measured on Q Exactive mass spectrometer,
637 as described in the LC-MS/MS data processing section (the matching metagenome gene
638 translation (1,595,268 entries) was used as microbial database). To compare the different
639 measurement methods, we counted for each sample the number of identified MS/MS, non-
640 redundant peptides and protein groups. These were then tested for significant differences
641 between measurement methods using paired t-tests via the ggplot2 package [61].

642 Single- versus two-step assessment

643 HeLa cells were prepared for LC-MS/MS measurements using published method [64]. Briefly,
644 cells were grown in DMEM medium and harvested at 80 % confluence. Proteins were
645 precipitated using acetone and methanol. Proteins were reduced with DTT and digested with
646 Lys-C and trypsin. Peptides were purified on Sep-Pak C18 Cartridge.

647 Sample was measured as described in the LC-MS/MS measurements section but for a few
648 changes. Peptide sample was eluted over 213 min using a 7 % (0 min), 15 % (140 min) and
649 33 % (213 min) gradient of solvent B (80 % ACN in 0.1 % formic acid) followed by a washout
650 procedure. The top 10 most intense ions were selected for HCD fragmentation.

651 Raw data were processed using MaxQuant (version 1.5.2.8) [58] as described in the LC-
652 MS/MS data processing section with a few alterations. The protein sequence databases used
653 for database search consisted of the complete *Homo sapiens* Uniprot database (93,799
654 sequences) and frequently observed contaminants (248 entries), as well as the mouse
655 microbiome catalogue (~2.6 million proteins) [26]. Several processings were performed
656 differing in the number of microbiome catalogue entries included, such as 1:0, 1:1, 1:2, 1:5,
657 1:10 and 1:20 *H. sapiens*:bacteria protein sequences and using a single- or two-step database
658 search [25].

659 Identified MS/MS, peptides and protein groups were assigned to kingdom of origin (conflicts
660 were resolved to Eukaryota by default). To compare the different database search strategies,
661 we counted the number of identified MS/MS, non-redundant peptides and protein groups
662 associated to each kingdom (as well as reverse hits and potential contaminants). We also
663 calculated the FDR based solely on reverse hits or together with bacterial hits in order to
664 investigate the true number of false positives.

665 Microbiome protein sequences databases comparison

666 The samples generated as part of the data analysis assessment were measured on a Q Exactive
667 HF mass spectrometer as described above, but were processed using additional databases and
668 search strategies in MaxQuant. Only the microbiome sequence databases differed and consisted
669 of one of (1) UniProt bacterial reference proteomes (5,408,622 protein entries), (2) UniProt
670 bacterial pan proteomes (18,541,701 entries), (3) protein sequences from the Mouse Gut
671 Metagenome catalogue (2,626,630 entries) [26], or (4) the matching metagenome gene
672 translation (1,595,268 entries). Processing also involved comparison between single- and two-
673 step database search [25]. To compare the different databases and search strategies, we counted
674 the number of sequences and OTUs per database, as well as reported the identification rates,
675 non-redundant peptides count and protein groups count.

676 Metagenome to metaproteome correlation

677 All subsequent sections use the samples generated as part of the data analysis assessment that
678 were measured on a Q Exactive HF mass spectrometer and processed against the matching
679 metagenome gene translation as described above. For direct comparison between metagenome
680 and metaproteome, the identified genes were collapsed into groups identical to protein groups
681 composition from mass spectrometry. Each gene groups abundance was calculated as the
682 highest gene abundance within that group. Each gene groups and corresponding protein groups
683 abundances were correlated across samples using Spearman's rank correlation from the stats
684 package. Significance was set at an adjusted p-value of 0.05 following Benjamini-Hochberg
685 multiple correction testing. The GSEA testing of KEGG pathways and Gene ontologies were
686 performed via the clusterProfiler package based on hypergeometric distribution ($p\text{-adj.} \leq 0.05$)
687 [63] following z-scoring of Spearman rho estimate per KEGG orthologies.

688 Metaproteome-based taxonomic representation

689 Metaproteome-derived taxonomic assignments were obtained using (1) Diamond (v. 0.9.23)
690 [28], (2) Unipept online (v. Dec. 2018) [23], or (3) Kraken2 (2.0.8-beta) [51] softwares for
691 either the identified peptides or all matched-metagenome protein sequences. The Diamond
692 alignment was performed against NCBI non-redundant protein sequences database using
693 sensitive and taxonomic classification mode. The Unipept online analysis was done via the
694 metaproteome analysis function with I-L not equal and advanced missed cleavages handling
695 selected against all UniProt entries. The Kraken2 k-mer analysis was carried out at the
696 nucleotide level (corresponding to identified peptides or proteins) against the Genome
697 Taxonomy Database (v. 89) [65] obtained from Struo software [66]. For each software
698 approach, the complete taxonomic lineage (NCBI or GTDB) was retrieved per peptide or
699 protein groups and the lowest common ancestor was determined. OTUs were quantified per
700 sample based on the different software approaches by summing either (1) PSM count, (2)
701 peptide intensity, (3) protein groups iBAQ, (4) protein groups intensity, (5) protein groups
702 LFQ, (6) protein groups MS/MS count. OTUs quantification were normalised on a per sample
703 basis as percentage of total to get OTUs representation. OTUs were filtered on a per sample
704 basis based on minimum representation threshold (representation $\geq X$, with X equal from 0 %
705 to 10 %), while the correlation between metagenome and metaproteome was calculated using
706 Spearman's rank correlation. The optimal representation filtering threshold was identified as
707 the threshold that maximises number of identified OTUs and correlation between omics (i.e.
708 count \times spearman rho). Using the optimal representation threshold and quantification approach,
709 the overlap in identified OTUs between omics was calculated, as well as the Spearman's rank
710 correlation at each taxonomic level.

711 Functional KEGG categories representation

712 For each sample, the protein groups iBAQ values were summed per KEGG category (level 2)
713 on the basis of KEGG orthology annotation. The same approach was also undertaken for gene
714 count. The KEGG category abundance were normalised for differing number of KO entries per
715 category and for variation between samples; this was done separately for metagenome and
716 metaproteome. Differences in KEGG category abundance between metagenome and
717 metaproteome were tested using paired t-tests from the stats package. Significance was set at
718 an adjusted p-value of 0.01 following Benjamini-Hochberg multiple correction testing.
719 Significantly changing KEGG categories were prioritised based on gene groups to protein
720 groups correlation (see section Metagenome to metaproteome correlation), whereby the
721 Wilcoxon rank-sum test was used to identify KEGG category containing KO entries whose
722 correlation differ from overall distribution (adjusted p-value ≤ 0.05).

723 To investigate further these selected KEGG categories, the protein groups iBAQ and gene
724 count were used as described in the previous paragraph to derive KO normalised abundance
725 and t-test results. Using the KO entries from each selected KEGG categories, separate GSEA
726 testing of KEGG pathways were performed via the clusterProfiler package based on
727 hypergeometric distribution (p-adj. ≤ 0.05).

728 Acknowledgments

729 MED, BM, XA and LD are grateful to the European Community 7th Framework Program under
730 Coordinated Action NEURON-ERANET (grant agreement 291840). BM was supported by
731 grants from the Deutsche Forschungsgemeinschaft (German Research Foundation Cluster of
732 Excellence EXC 2124). BM and NN acknowledge support by the High Performance and Cloud
733 Computing Group at the Center for Data Processing of the University of Tübingen, the state of
734 Baden-Wuerttemberg through bwHPC. The metagenomic work detailed herein used the

735 computing resources of the UK MEDical BIOinformatics partnership – aggregation,
736 integration, visualization and analysis of large, complex data (UK Med-Bio) – which was
737 supported by the Medical Research Council (grant number MR/L01632X/1). XA acknowledge
738 support from the MINECO, Spain (grant number PCIN-2014-105).

739 Author contributions

740 LD, XA, MD and BM designed the study. XA and CG generated the mouse cohorts and
741 collected the murine faecal material. VA, TG and ID prepared the murine faecal samples for
742 proteomic measurement by mass spectrometry. LH processed the metagenomic data,
743 generating the taxonomic and gene abundance outputs. NN processed the metaproteomic
744 datasets and performed the proteogenomic integration. NN wrote the manuscript with the input
745 from all authors.

746 Data Access

747 The complete metaproteomic bioinformatic workflow is available online [67]. The mass
748 spectrometry proteomics data have been deposited to the ProteomeXchange Consortium via
749 the PRIDE [68] partner repository with the dataset identifiers PXD020695, PXD020738,
750 PXD021928 and PXD021932. Trimmed whole genome sequence data with mouse reads
751 removed have been deposited with GenBank, EMBL and DDBJ databases under the BioProject
752 accession PRJNA473429.

753 References

- 754 1. Jandhyala SM, Talukdar R, Subramanyam C, et al. Role of the normal gut microbiota.
755 World J Gastroenterol. 2015;21(29):8787-8803.
- 756 2. Valdes AM, Walter J, Segal E, et al. Role of the gut microbiota in nutrition and health.
757 BMJ. 2018;361:k2179.

- 758 3. Arumugam M, Raes J, Pelletier E, et al. Enterotypes of the human gut microbiome.
759 Nature. 2011 2011/05//;473(7346):174-180.
- 760 4. Vieira-Silva S, Falony G, Belda E, et al. Statin therapy is associated with lower
761 prevalence of gut microbiota dysbiosis. Nature. 2020 May;581(7808):310-315.
- 762 5. Wang J, Linnenbrink M, Künzel S, et al. Dietary history contributes to enterotype-like
763 clustering and functional metagenomic content in the intestinal microbiome of wild
764 mice. Proceedings of the National Academy of Sciences. 2014;111(26):E2703.
- 765 6. Ladau J, Eloje-Fadrosch EA. Spatial, Temporal, and Phylogenetic Scales of Microbial
766 Ecology. Trends in Microbiology. 2019 2019/08/01//;27(8):662-669.
- 767 7. Parfrey LW, Knight R. Spatial and temporal variability of the human microbiota.
768 Clinical Microbiology and Infection. 2012 2012/07/01//;18:5-7.
- 769 8. Sampson TR, Debelius JW, Thron T, et al. Gut Microbiota Regulate Motor Deficits and
770 Neuroinflammation in a Model of Parkinson's Disease. Cell. 2016 Dec 1;167(6):1469-
771 1480 e12.
- 772 9. Ley RE, Turnbaugh PJ, Klein S, et al. Microbial ecology: human gut microbes
773 associated with obesity. Nature. 2006 Dec 21;444(7122):1022-3.
- 774 10. Xiao L, Sonne SB, Feng Q, et al. High-fat feeding rather than obesity drives
775 taxonomical and functional changes in the gut microbiota in mice. Microbiome. 2017
776 Apr 8;5(1):43.
- 777 11. Zhang X, Deeke SA, Ning Z, et al. Metaproteomics reveals associations between
778 microbiome and intestinal extracellular vesicle proteins in pediatric inflammatory
779 bowel disease. Nat Commun. 2018 Jul 20;9(1):2873.
- 780 12. Hsiao EY, McBride SW, Hsien S, et al. Microbiota modulate behavioral and
781 physiological abnormalities associated with neurodevelopmental disorders. Cell. 2013
782 Dec 19;155(7):1451-63.
- 783 13. Tengeler AC, Dam SA, Wiesmann M, et al. Gut microbiota from persons with
784 attention-deficit/hyperactivity disorder affects the brain in mice. Microbiome. 2020
785 2020/04/01//;8(1):44.
- 786 14. Wang L, Christophersen CT, Sorich MJ, et al. Increased abundance of *Sutterella* spp.
787 and *Ruminococcus torques* in feces of children with autism spectrum disorder. Mol
788 Autism. 2013 Nov 4;4(1):42.
- 789 15. Biagi E, Candela M, Centanni M, et al. Gut Microbiome in Down Syndrome. PLOS
790 ONE. 2014;9(11):e112023.
- 791 16. Hettich RL, Pan C, Chourey K, et al. Metaproteomics: harnessing the power of high
792 performance mass spectrometry to identify the suite of proteins that control metabolic
793 activities in microbial communities. Anal Chem. 2013 May 7;85(9):4203-14.
- 794 17. Li X, LeBlanc J, Truong A, et al. A metaproteomic approach to study human-microbial
795 ecosystems at the mucosal luminal interface. PloS one. 2011;6(11):e26542-e26542.
- 796 18. Ram RJ, Verberkmoes NC, Thelen MP, et al. Community proteomics of a natural
797 microbial biofilm. Science. 2005 Jun 24;308(5730):1915-20.
- 798 19. Tanca A, Palomba A, Pisanu S, et al. Enrichment or depletion? The impact of stool
799 pretreatment on metaproteomic characterization of the human gut microbiota.
800 Proteomics. 2015 Oct;15(20):3474-85.
- 801 20. Zhang X, Li L, Mayne J, et al. Assessing the impact of protein extraction methods for
802 human gut metaproteomics. J Proteomics. 2018 May 30;180:120-127.
- 803 21. Muth T, Behne A, Heyer R, et al. The MetaProteomeAnalyzer: a powerful open-source
804 software suite for metaproteomics data analysis and interpretation. J Proteome Res.
805 2015 Mar 6;14(3):1557-65.
- 806 22. Cheng K, Ning Z, Zhang X, et al. MetaLab: an automated pipeline for metaproteomic
807 data analysis. Microbiome. 2017 Dec 2;5(1):157.

- 808 23. Mesuere B, Devreese B, Debyser G, et al. Unipept: tryptic peptide-based biodiversity
809 analysis of metaproteome samples. *J Proteome Res.* 2012 Dec 7;11(12):5773-80.
- 810 24. Mesuere B, Van der Jeugt F, Willems T, et al. High-throughput metaproteomics data
811 analysis with Unipept: A tutorial. *J Proteomics.* 2018 Jan 16;171:11-22.
- 812 25. Jagtap P, Goslinga J, Kooren JA, et al. A two-step database search method improves
813 sensitivity in peptide sequence matches for metaproteomics and proteogenomics
814 studies. *Proteomics.* 2013 Apr;13(8):1352-7.
- 815 26. Xiao L, Feng Q, Liang S, et al. A catalog of the mouse gut metagenome. *Nature*
816 *biotechnology.* 2015 Oct;33(10):1103-8.
- 817 27. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using
818 exact alignments. *Genome Biol.* 2014 Mar 3;15(3):R46.
- 819 28. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND.
820 *Nat Methods.* 2015 Jan;12(1):59-60.
- 821 29. Kleiner M, Thorson E, Sharp CE, et al. Assessing species biomass contributions in
822 microbial communities via metaproteomics. *Nat Commun.* 2017 Nov 16;8(1):1558.
- 823 30. Kanehisa M, Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids*
824 *Res.* 2000 Jan 1;28(1):27-30.
- 825 31. Zhang X, Ning Z, Mayne J, et al. MetaPro-IQ: a universal metaproteomic approach to
826 studying human and mouse gut microbiota. *Microbiome.* 2016 Jun 24;4(1):31.
- 827 32. Tanca A, Manghina V, Fraumene C, et al. Metaproteogenomics Reveals Taxonomic
828 and Functional Changes between Cecal and Fecal Microbiota in Mouse. *Front*
829 *Microbiol.* 2017;8:391.
- 830 33. Michalski A, Damoc E, Hauschild JP, et al. Mass spectrometry-based proteomics using
831 Q Exactive, a high-performance benchtop quadrupole Orbitrap mass spectrometer. *Mol*
832 *Cell Proteomics.* 2011 Sep;10(9):M111 011015.
- 833 34. Williamson JC, Edwards AV, Verano-Braga T, et al. High-performance hybrid
834 Orbitrap mass spectrometers for quantitative proteome analysis: Observations and
835 implications. *Proteomics.* 2016 Mar;16(6):907-14.
- 836 35. Shishkova E, Hebert Alexander S, Coon Joshua J. Now, More Than Ever, Proteomics
837 Needs Better Chromatography. *Cell Systems.* 2016 2016/10/26;3(4):321-324.
- 838 36. Nguyen TLA, Vieira-Silva S, Liston A, et al. How informative is the mouse for human
839 gut microbiota research? *Disease Models & Mechanisms.* 2015;8(1):1.
- 840 37. Hart ML, Meyer A, Johnson PJ, et al. Comparative Evaluation of DNA Extraction
841 Methods from Feces of Multiple Host Species for Downstream Next-Generation
842 Sequencing. *PLOS ONE.* 2015;10(11):e0143334.
- 843 38. Nagpal R, Wang S, Solberg Woods LC, et al. Comparative Microbiome Signatures and
844 Short-Chain Fatty Acids in Mouse, Rat, Non-human Primate, and Human Feces. *Front*
845 *Microbiol.* 2018;9:2897-2897.
- 846 39. Speicher KD, Kolbas O, Harper S, et al. Systematic analysis of peptide recoveries from
847 in-gel digestions for protein identifications in proteome studies. *J Biomol Tech.*
848 2000;11(2):74-86.
- 849 40. Muth T, Kolmeder CA, Salojarvi J, et al. Navigating through metaproteomics data: A
850 logbook of database searching. *Proteomics.* 2015 Oct;15(20):3439-53.
- 851 41. Heyer R, Schallert K, Zoun R, et al. Challenges and perspectives of metaproteomic data
852 analysis. *J Biotechnol.* 2017 Nov 10;261:24-36.
- 853 42. Li J, Jia H, Cai X, et al. An integrated catalog of reference genes in the human gut
854 microbiome. *Nature biotechnology.* 2014 2014/08/01;32(8):834-841.
- 855 43. Chen T, Yu WH, Izard J, et al. The Human Oral Microbiome Database: a web
856 accessible resource for investigating oral microbe taxonomic and genomic information.
857 *Database (Oxford).* 2010 Jul 6;2010:baq013.

- 858 44. Grassl N, Kulak NA, Pichler G, et al. Ultra-deep and quantitative saliva proteome
859 reveals dynamics of the oral microbiome. *Genome medicine*. 2016;8(1):44.
- 860 45. Erickson AR, Cantarel BL, Lamendella R, et al. Integrated
861 metagenomics/metaproteomics reveals human host-microbiota signatures of Crohn's
862 disease. *PLoS One*. 2012;7(11):e49138.
- 863 46. Mills RH, Vazquez-Baeza Y, Zhu Q, et al. Evaluating Metagenomic Prediction of the
864 Metaproteome in a 4.5-Year Study of a Patient with Crohn's Disease. *mSystems*. 2019
865 Jan-Feb;4(1).
- 866 47. Edirisinghe JN, Weisenhorn P, Conrad N, et al. Modeling central metabolism and
867 energy biosynthesis across microbial life. *BMC Genomics*. 2016 Aug 8;17:568.
- 868 48. Sidoli S, Kulej K, Garcia BA. Why proteomics is not the new genomics and the future
869 of mass spectrometry in cell biology. *J Cell Biol*. 2017 Jan 2;216(1):21-24.
- 870 49. Hoyles L, Fernandez-Real JM, Federici M, et al. Molecular phenomics and
871 metagenomics of hepatic steatosis in non-diabetic obese women. *Nat Med*. 2018
872 Jul;24(7):1070-1080.
- 873 50. Andrews S. FastQC A Quality Control tool for High Throughput Sequence Data.
874 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>; 2010.
- 875 51. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. *Genome*
876 *Biol*. 2019 Nov 28;20(1):257.
- 877 52. Parks DH, Chuvochina M, Chaumeil PA, et al. A complete domain-to-species
878 taxonomy for Bacteria and Archaea. *Nature biotechnology*. 2020 Sep;38(9):1079-1086.
- 879 53. Méric G, Wick RR, Watts SC, et al. Correcting index databases improves metagenomic
880 studies. *bioRxiv*. 2019:712166.
- 881 54. Huerta-Cepas J, Forslund K, Coelho LP, et al. Fast Genome-Wide Functional
882 Annotation through Orthology Assignment by eggNOG-Mapper. *Molecular Biology*
883 *and Evolution*. 2017;34(8):2115-2122.
- 884 55. Bradford MM. A rapid and sensitive method for the quantitation of microgram
885 quantities of protein utilizing the principle of protein-dye binding. *Anal Biochem*. 1976
886 May 7;72:248-54.
- 887 56. Wisniewski JR, Zougman A, Nagaraj N, et al. Universal sample preparation method for
888 proteome analysis. *Nat Methods*. 2009 May;6(5):359-62.
- 889 57. Rappsilber J, Mann M, Ishihama Y. Protocol for micro-purification, enrichment, pre-
890 fractionation and storage of peptides for proteomics using StageTips. *Nat Protoc*.
891 2007;2(8):1896-906.
- 892 58. Cox J, Mann M. MaxQuant enables high peptide identification rates, individualized
893 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature*
894 *biotechnology*. 2008 Dec;26(12):1367-72.
- 895 59. Cox J, Hein MY, Lubner CA, et al. Accurate proteome-wide label-free quantification by
896 delayed normalization and maximal peptide ratio extraction, termed MaxLFQ. *Mol Cell*
897 *Proteomics*. 2014 Sep;13(9):2513-26.
- 898 60. R Core Team. R: A Language and Environment for Statistical Computing. Vienna,
899 Austria: R Foundation for Statistical Computing; 2018.
- 900 61. Wickham H. *ggplot2: Elegant Graphics for Data Analysis*. 2016.
- 901 62. Gatto L, Lilley KS. MSnbase-an R/Bioconductor package for isobaric tagged mass
902 spectrometry data visualization, processing and quantitation. *Bioinformatics*. 2012 Jan
903 15;28(2):288-9.
- 904 63. Yu G, Wang LG, Han Y, et al. clusterProfiler: an R package for comparing biological
905 themes among gene clusters. *OMICS*. 2012 May;16(5):284-7.

- 906 64. Schmitt M, Sinnberg T, Nalpas NC, et al. Quantitative Proteomics Links the
907 Intermediate Filament Nestin to Resistance to Targeted BRAF Inhibition in Melanoma
908 Cells. *Mol Cell Proteomics*. 2019 Jun;18(6):1096-1109.
- 909 65. Parks DH, Chuvochina M, Waite DW, et al. A standardized bacterial taxonomy based
910 on genome phylogeny substantially revises the tree of life. *Nature biotechnology*. 2018
911 Nov;36(10):996-1004.
- 912 66. de la Cuesta-Zuluaga J, Ley RE, Youngblut ND. Struo: a pipeline for building custom
913 databases for common metagenome profilers. *Bioinformatics*. 2020 Apr 1;36(7):2314-
914 2315.
- 915 67. Nalpas N, Macek B. Integrated metaproteomics workflow. 1.0. Zenodo; 2020.
- 916 68. Perez-Riverol Y, Csordas A, Bai J, et al. The PRIDE database and related tools and
917 resources in 2019: improving support for quantification data. *Nucleic Acids Res*. 2019
918 Jan 8;47(D1):D442-D450.
- 919

920 Tables

921 **Table 1: Performance comparison of different sample preparation and data analysis**
 922 **steps.** In bold are the best methods according to assessed criteria: peptide/protein count,
 923 host/dietary contamination, *Firmicutes* to *Bacteroidetes* ratio, time efficiency, sample amount,
 924 FDR, identification rate, OTU assigned peptides and number of OTUs. The performance status
 925 is displayed using minus sign for poor, equal sign for similar/no difference or plus sign for
 926 good performance.

		Peptide/protein count	Host/dietary contamination	Firmicutes:Bacteroidetes	Time efficiency	Sample amount	FDR	Identification rate	OTU assigned peptides	Number of OTUs
Centrifugation	LSC	+	-	+	-	=				
	nLSC	-	+	-	+	=				
Digestion	In-solution	+			-	=				
	FASP	-			+	=				
Lysis	Bead-beating	=			+	=				
	Nitrogen	=			-	=				
	Sonication	=			-	=				
MS instrument	Elite	-			-	-				
	Elite long LC	-			-	-				
	Q Exactive	++			+	+				
Search strategy	Single-step	-			+		+	-		
	Two-step	+			-		--	+		
Protein sequence database	UP Ref proteomes	--			-		-	--		
	UP Pan proteomes	--			-		-	--		
	Catalogue	+			+		+	+		
	Metagenome	++			+		+	++		
OTUs quantification	Diamond				--				--	-
	Unipept				+				-	-
	Kraken2				-				+	++

927

928 Figures

929 **Figure 1: Low speed centrifugation impacts protein identification and taxonomic**
930 **representation.** A) Number of identified peptides and protein groups per samples for the
931 comparison between LSC (red) and nLSC (blue) methods. B) Number of identified peptides
932 and protein groups per samples for the comparison between LSC-in solution digestion (red),
933 LSC-FASP (grey), nLSC-in solution digestion (blue) and nLSC-FASP (orange) methods. C)
934 Number of identified peptides and protein groups per samples for the comparison between
935 bead-beating (red), nitrogen (grey) and sonication (blue) lysis methods. A-C) Represented
936 significance results correspond to paired t-test on $N = 12$: * p -value ≤ 0.05 , ** ≤ 0.01 , *** \leq
937 0.01. D-E) Unipept-derived taxonomic representation (down to phylum level) for the peptide
938 identified in the LSC (D) and nLSC (E) samples. Number represent the Peptide Spectrum
939 Match count for each taxon.

940 **Figure 2: Functional representation of proteins detected in LSC- and nLSC-samples is**
941 **similar despite different peptide and protein abundances.** A) Overlap in the overall
942 identified peptides or protein groups between the LSC and nLSC methods. B) Principal
943 component analysis showing separation of the samples based on biological replicates and
944 sample preparation (LSC or nLSC) and the clustering of the technical replicates. C) Volcano
945 plot of the protein abundance comparison between LSC and nLSC approaches. Significant
946 protein groups based on paired t-test from $N = 12$ with $FDR \leq 0.01$ and absolute fold-change
947 ≥ 2.5 . D) KEGG pathways over-representation testing for the significantly up-regulated (red)
948 and down-regulated (blue) protein groups between LSC and nLSC sample preparation
949 approaches. Fisher exact-test threshold set to adjusted p -value ≤ 0.1 .

950 **Figure 3: Two-step database search in combination with target-decoy strategy leads to a**
951 **dramatic increase in false positive rate.** A) The protein groups count is shown for single- or
952 two-step search strategies across increasingly large protein sequence databases. Counts are

953 colour-coded per category, with eukaryote (grey), bacteria (red), contaminant (blue) and
954 reverse (orange) hits. B) The FDR is calculated for single- or two-step search strategies across
955 increasingly large protein sequence databases. The FDR is calculated based on reverse hits only
956 (circle shape) or reverse plus bacterial hits (triangle shape). C) The faeces metaproteome data
957 from a cohort of 38 mice were searched against different protein sequence databases and either
958 via single- or two-step search strategy. The comparison between databases focuses on the
959 number of sequences, number of OTUs, identification rate, number of peptides identified and
960 number of protein groups identified. The following databases were compared: UniProt
961 bacterial reference proteome (blue colours), UniProt bacterial pan-proteome (black colours),
962 mouse microbiome catalog (orange) and matching metagenome sequences (red).

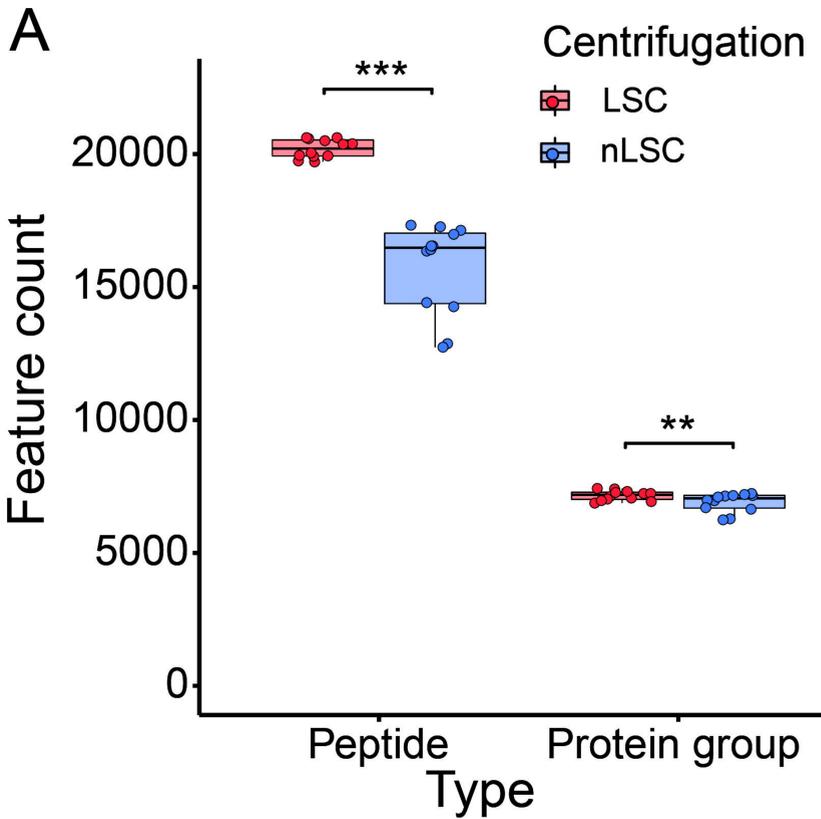
963 **Figure 4: Metaproteome and metagenome show positive correlation at the level of**
964 **protein, gene and taxonomic representation.** A) Assessment of the optimal filtering
965 threshold to maximise the total number of identified OTUs and the OTUs abundance
966 correlation between metaproteome- and metagenome-derived taxonomy representation. For
967 the filtering, OTUs are retained on a per sample basis when having a taxonomic representation
968 superior or equal to defined filtering value. Metaproteome-derived taxonomic representation
969 was calculated from Kraken2 software results using different calculation methods (i.e., peptide
970 and protein quantifications). B) The overlap in OTUs between metaproteome and metagenome
971 at different taxonomic levels. Numbers above bars correspond to the spearman rank correlation
972 between metaproteome- and metagenome-derived OTUs abundance for the corresponding
973 taxonomic level with $N = 38$: * p -value ≤ 0.05 , ** ≤ 0.01 , *** ≤ 0.01 . C) Correlation in
974 abundances is shown between each protein groups (metaproteome) and corresponding gene
975 “groups” (metagenome). Correlation was tested using Spearman’s rank correlation and p -value
976 was adjusted for multiple testing using Benjamini-hochberg correction. Significantly
977 correlating protein/gene groups are in red colours, while significantly anti-correlating

978 protein/gene groups are in green colours. D) GSEA of KEGG pathways based on ranking of
979 the protein/gene groups correlation. Pathway node colour corresponds to GSEA results
980 adjusted *p*-value and node size matches the number of protein/gene group assigned to the
981 pathway.

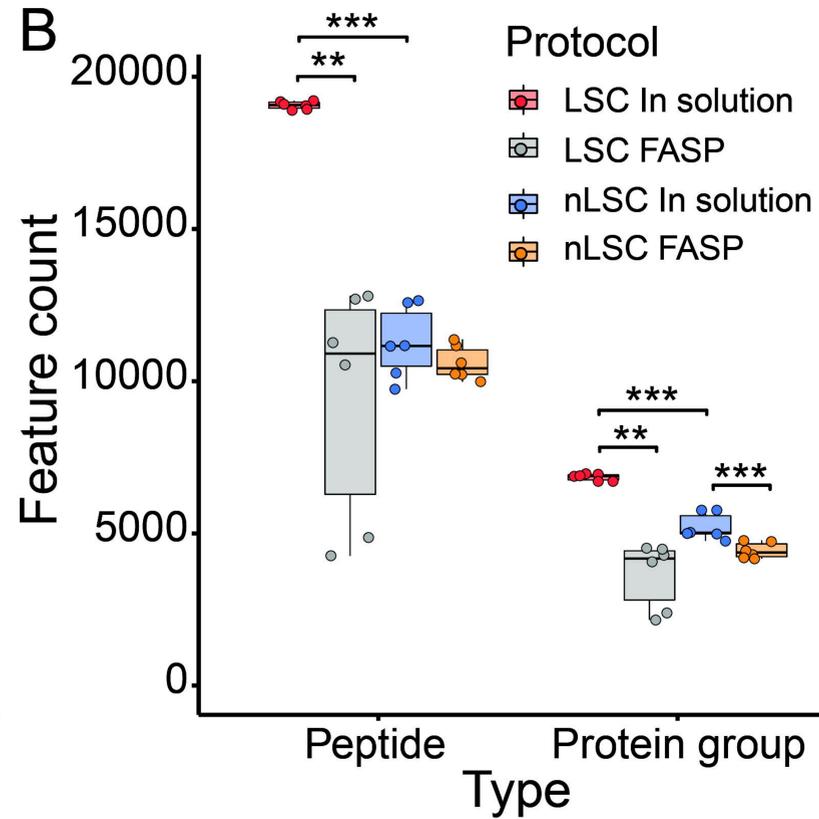
982 **Figure 5: Functionally active pathways derived from the metaproteome differs from the**
983 **metagenome potential.** A) Comparison in the proportion of selected KEGG functional
984 categories (level 2) between metaproteome (red) and metagenome (grey). Paired t-test *p*-values
985 are indicated (N = 38). B) GSEA of KEGG pathways based on ranking of t-test results from
986 KEGG orthology proportion between metaproteome and metagenome. KEGG pathways are
987 colour-coded based on KEGG functional categories (level 2). Only significantly over-
988 represented KEGG pathways are shown with adjusted *p*-value ≤ 0.05 . C) Interaction network
989 between KEGG orthologies and KEGG pathways for the KEGG functional category “Protein
990 families: genetic information processing”. Pathway node size corresponds to number of KEGG
991 orthologies associated to it. KEGG orthologies are colour-coded based on directional adjusted
992 *p*-value from the t-test comparison between metaproteome and metagenome.

993

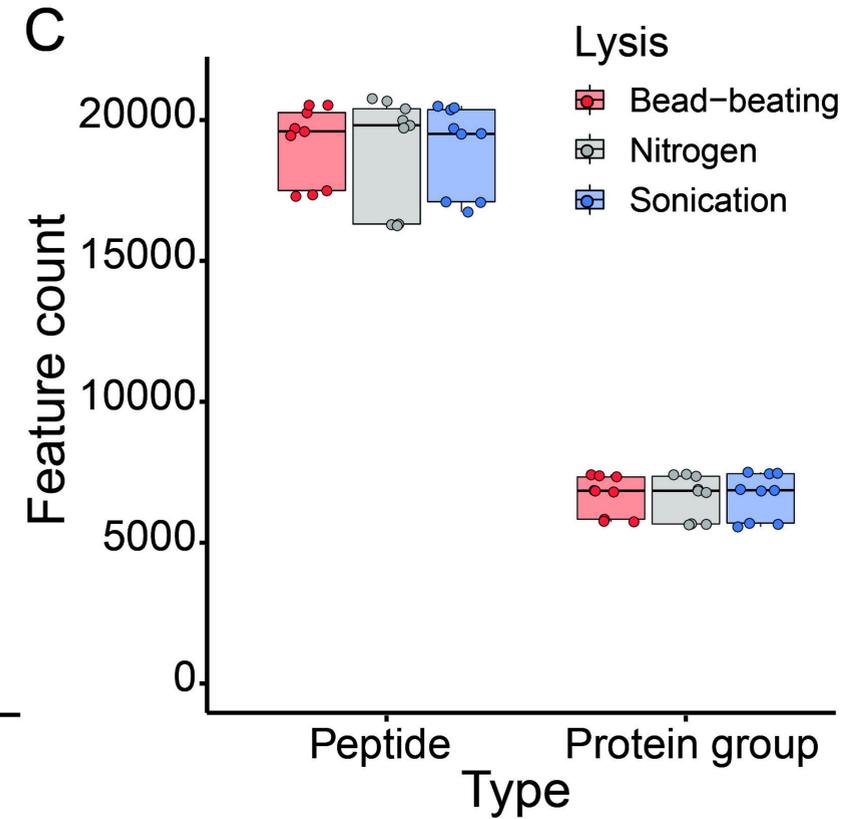
A



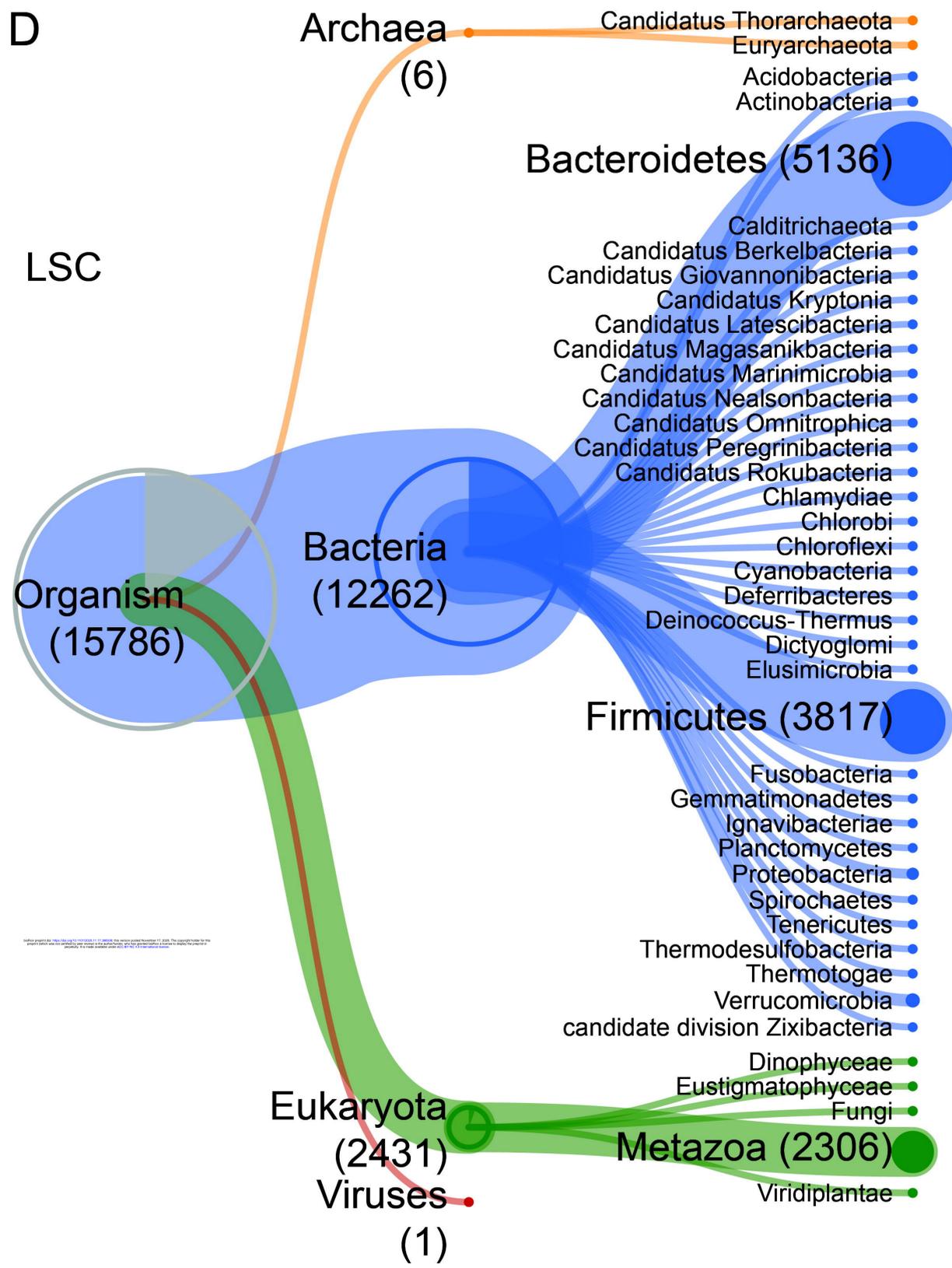
B



C



D



E

