1    MetaDIA: A Novel Database Reduction Strategy for DIA Human Gut Metaproteomics

2    **Abstract**

3    **Background**: Microbiomes, especially within the gut, are complex and may comprise hundreds

4    of species. The identification of peptides in metaproteomics presents a significant challenge, as

5    it involves matching peptides to mass spectra within an enormous search space for complex

6    and unknown samples. This poses difficulties for both the accuracy and the speed of

7    identification. Specifically, analysis of data-independent acquisition (DIA) datasets has relied on

8    libraries constructed from prior data-dependent acquisition (DDA) results. This approach

9    requires running the samples in DDA mode to construct a library from the identified results,

10   which can then be used for the DIA data. However, this method is resource-intensive, consumes

11   samples, and limits identification to peptides previously identified by DDA. These limitations

12   restrict the application of DIA in metaproteomics research.

13   **Results**: We introduced a novel strategy to reduce the search space by utilizing species

14   abundance and functional abundance information from the microbiome to score each peptide

15   and prioritize those most likely to be detected. Employing this strategy, we have developed and

16   optimized a workflow called MetaDIA for analysis of microbiome DIA data, which operates

17   independently of DDA assistance. Our method demonstrated strong consistency with the

18   traditional DDA-based library approach at both protein and functional levels.

19   **Conclusion**: Our approach successfully created a smaller, yet sufficient database for DIA data

20   search requirements in metaproteomics, showing high consistency with results from the

21   conventional DDA-based library. We believe this method can facilitate the application of DIA in

22   metaproteomics.

23   **Key**: Metaproteomics, Human gut microbiome, DIA, DDA-free, diaPASEF

## Introduction

The microbiome encompasses a diverse array of microorganisms residing in different organisms, ecosystems, and environmental settings such as the human body, animals, plants, soil, water bodies, and various ecological niches[1, 2]. Metaproteomics serves as a tool for understanding the roles of proteins within these microbial communities[3]. Mass spectrometry-based proteomics aims to study all proteins in a sample. However, applying these techniques to the microbiome is challenged by its complexity. Without prior knowledge of the microbes present in a sample, metaproteomics relies on searching mass spectra against a large database, making the task of matching peptides and spectra notably challenging. Employing an iterative search strategy significantly reduces the search complexity in which the final search is against a database generated from previous searching results [4, 5]. The iterative strategy has been successfully used but only for the data acquired by data-dependent acquisition (DDA) mode[6, 7], Unfortunately, in DDA mode, only the most abundant precursor ions are selected for further inquiry, and lower abundant ones are overlooked[8].

In contrast, data-independent acquisition (DIA) uses a set of precursor isolation windows to collect all the fragments ions indiscriminately[9]. It has shown remarkable robustness, sensitivity, and reproducibility with fewer missing values[10]. DIA can be coupled with microLC enabling high-throughput analysis[11]. This makes it particularly suitable for conducting large-scale analyses. The DIA-PASEF[12] method integrates ion mobility separation with the DIA workflow, adding a fourth dimension of analyzing ion mobility to the traditional three-dimensional data set. This not only enriches the structural information of analytes but also enhances ion utilization efficiency leveraging the linear relation between ion mobility and mass-to-charge ratio. Another improvement in mass spectrometer scanning speed enables the utilization of smaller isolation windows in DIA, termed as narrow-window DIA[13]. This approach achieves comprehensive peptide precursor coverage and high quantitative precision and accuracy. In bioinformatics, the

49　development of prediction software for peptide properties (theoretically predicted spectrum[14,

50　15], retention times[16-18]) enables the querying of DIA datasets without dependence on

51　libraries generated by DDA. Those predicted libraries even showed better performance than the

52　measured libraries[19]. Moreover, DIA-specific searching software such as DIA-NN[20, 21],

53　MaxDIA[22], and Spectronaut have shown reliable results for the identification and quantification

54　of peptides. The above advantages make DIA increasingly popular in proteomics. However, it is

55　noteworthy that the benefits conferred by these techniques have not yet been fully extended to

56　the field of metaproteomics. The main reason is that the inherent complexity of DIA data

57　requires a much more constrained searching space compared with DDA data. To date, only a

58　few metaproteomics studies have been done, and they were all compelled to use a spectral

59　library derived from DDA data[23-25]. The DDA-derived method involves creating a spectral

60　library from DDA runs for each sample, which is then used to interpret complex mass spectra

61　from subsequent analyses. This approach requires multiple sample aliquots, extensive mass

62　spectrometry resources and is limited to detecting peptides previously identified by DDA.

63　Gladiator[26] uses DIA-Umpire[27] to assemble pseudo-DDA spectra from DIA data for

64　microbiome samples. The method does not require a DDA-based spectral library for its

65　operation, however, it still relies on spectrum-centric algorithms and does not fully exploit the

66　potential advantages of DIA data.

67　Therefore, to leverage the benefits of DIA in metaproteomics, the searching space needs to be

68　further reduced. In the previous DDA iterative strategy[7, 28], the high-abundant proteins (HAP)

69　were used for the first search to infer the species that exist in the sample then all the proteins

70　belonging to those species were then used for the subsequent search.  However, this database

71　remains overly extensive when compared to the number of identified peptides. Since the

72　abundance of species within the microbiome shows significant disparity[29], the species

73　identified should not be considered equally. The same applies to proteins and peptides. Proteins

74  with high abundance and peptides with high detectability[30] or shared among various species

75  are more likely to be detected. Here we report on a DIA workflow for metaproteomics, called

76  MetaDIA, that relies on an annotated peptide database. This database comprises peptides that

77  are anticipated to be detected, leveraging information on species abundance and protein

78  abundance to score each peptide. We conducted a proof-of-concept experiment on human gut

79  microbiome data generated by diaPASEF mode[23]. The peptide identification number and

80  quantitative results obtained through our peptide library are comparable to those from the DDA-

81  based library. Moreover, the species and functional information obtained from both methods are

82  highly consistent.

83

## Materials and methods

**Reference peptide sequence with detectability score for human gut microbiome**

86  The Unified Human Gastrointestinal Protein (UHGP) catalog, encompassing 4744 assembled

87  genomes from the human gut microbiome, served as the reference database for this study[31].

88  Within this catalog, each protein sequence is uniquely associated with a distinct genome and is

89  accompanied by detailed taxonomic and functional annotations. The detectability of peptides

90  derived from these protein sequences was predicted using DeepDetect[30], a deep learning

91  algorithm specifically designed for this purpose. This process involved in silico digestion of the

92  protein sequences and subsequent assignment of a detectability score to each resultant peptide.

93  Consequently, the peptide sequence reference database was enhanced by annotating each

94  peptide with three key pieces of information: the genome identifier, the protein identifier, and the

95  peptide's detectability score. Please note that the database is structured on an identifier-centric

96  organization. This means that peptides with identical sequences may be present within the

97      database; however, as long as they are not from same genome and protein, they are

98      distinguished by unique identifiers.

99      **Generation of FuncTax score**

100     Firstly, identified peptides by MetaPep[32] are mapped to the UHGP database to establish

101     peptide-genome associations. Subsequently, a greedy algorithm is employed to identify the

102     minimal set of genomes that encompasses all peptide sequences, effectively reducing the

103     complexity of the dataset. Following this, the intensity of each peptide is aggregated to infer

104     genome abundance. The relative genome abundance will be used as the taxonomic score. To

105     address the assignment of shared peptides, a razor strategy is adopted, analogous to the

106     MaxQuant approach for protein inference[33]. Specifically, when a peptide is found in multiple

107     genomes, it is attributed to the genome with the greater number of associated peptides.

108     However, this typically results in approximately 1,000 genomes remaining, with many containing

109     only a single peptide. The number substantially larger than that is found in a typical human gut

110     microbiome which is around 200[29]. So, we only choose the most abundant species for

111     subsequent analysis. The selection of species for consideration is further explored in the

112     optimization section of the study.

113     For the functional score, we constructed a fixed table from the MetaPep project [32]. While

114     building the database Metapep, the peptide identification was performed by the software

115     MetaLab MAG[7],  which provides quantifications of protein abundance. Those proteins are well

116     annotated. Subsequently, the relative abundance of each Clusters of Orthologous Groups (COG)

117     accession was computed. Samples comprising fewer than 1000 COG accessions were

118     considered to be of low quality and consequently were omitted from the analysis. A total of

119     1,031 high-quality samples were retained for further evaluation. The mean of non-zero relative

120     abundance of the COG accessions was then determined across these 1,031 samples,

121     establishing a metric referred to the functional score.

122 The FuncTax score was obtained by multiplying two scores. In the case of peptides with the

123 same sequence, their FuncTax scores were combined to give higher priority to shared peptides;

124 the highest detectability score among them was utilized to ensure the inclusion of all possible

125 peptides.

126 **Taxonomic and functional analysis**

127 The taxonomic analysis is similar to the generation of genomic abundance score.  The identified

128 peptides are mapped to a database to establish peptide-genome associations. The database

129 contains only the top 50 genomes. In our workflow, the peptide database was filtered out from

130 the top 50 genomes. So, all the identified peptides were from the top genomes and thus can be

131 used for the taxonomic analysis (the peptides added from MetaPep may not be used). In the

132 DDA-based method, the peptide identified by the DDA library can be annotated to over 1,000

133 genomes even after using the greedy algorithm described above (method: generation of

134 FuncTax score). However, we found that the top 50 genomes accounted for 79%-90% of

135 peptides and 87%-92% of peptide intensity (Supplementary Figure 1). To simplify the

136 comparison between the two methods, we discarded the small number of peptides that cannot

137 be annotated to the top 50 genomes. Similarly, the razor strategy is used to process peptides

138 shared by multiple genomes. Finally, the intensity of each peptide is aggregated to infer genome

139 abundance.

140 For functional analysis, a protein abundance was firstly generated using the same strategy as

141 taxonomic analysis. The proteins in the UHGP database have been extensively annotated thus

142 the protein abundance can be further interpreted into functional abundance.

143 **Deepdetect software configuration**

144 Protein digestion was simulated using Trypsin with the following parameters: a maximum of two

145 missed cleavages, and peptide lengths ranging from 7 to 50 amino acids. Default settings were

146 applied for all other parameters.

**DIA software configuration**

148 DIA-NN (version 1.8.1) was used to process all the DIA data in this study. Maximum mass

149 accuracy tolerances were set to 10 ppm for both MS1 and MS2 spectra. The --relaxed-prot-inf

150 option was used for library-free searching. The --no-maxlfq option was used to disable the

151 normalization for the quantification benchmark experiment. All other settings were left default.

152 The precursor matrix containing the peptide information was used for taxonomic and functional

153 analysis.

**Metaproteomic datasets**

155 The dataset used for optimizing workflow is sourced from a published study and shared by the

156 authors[23]. The dataset for evaluating accuracy is from in-house samples. *Blautia*

157 *hydrogenotrophica* (DSM 101114; Leibniz Institute DSMZ- German collection of microorganisms

158 and cell cultures) was cultured in LB broth. The human stool was collected from a healthy adult

159 volunteer at the University of Ottawa, Ottawa, ON, CAN. The protocol (# 20160585-01H) was

160 approved by Ottawa Health Science Network Research Ethics. The protein extraction and

161 digestion were performed as described previously[34]. Peptide concentrations were measured

162 using Thermo Scientific Pierce Quantitative Colorimetric Peptide Assays according to the

163 manufacturer's directions.

164 The in-house samples were then analysed using an UltiMate 3000 RSLCnano system (Thermo

165 Fisher Scientific, USA) coupled to an Orbitrap Exploris 480 mass spectrometer (Thermo Fisher

166 Scientific, USA). Peptides were loaded onto a tip column (75 µm inner diameter ×15 cm) packed

167 with reverse phase beads (3 µm/120 Å ReproSil-Pur C18 resin, Dr. Maisch HPLC GmbH). A 60

168 min gradient of 5 to 35% (v/v) from buffer A (0.1% (v/v) formic acid) to B (0.1% (v/v) formic acid

169 with 80% (v/v) acetonitrile) at a flow rate of 300 μL/min was used. The mass spectrometer was

170 in data-independent mode covering the mass range of 380–980 m/z with 10 m/z isolation

171 windows. 

**Availability of the pipeline**

173 The whole pipeline is available for use at https://github.com/northomics/MetaDIA

174

# Result

**MetaDIA Workflow overview: Taxonomy- and function-guided construction of peptide database for metaproteomics**

178 Here we propose a new workflow for DIA based metaproteomics called MetaDIA. MetaDIA is a

179 multistep workflow that systematically reduces the search space for DIA searching. At its basis,

180 it relies on a combination of taxonomic abundance, functional abundance as a proxy of protein

181 levels, and peptide detectability ultimately enabling DIA searching without the need for DDA

182 results. Briefly, in the first step, we created a new database of peptides, called MetaPepDetec,

183 obtained by in silico digestion and detectability prediction of the Unified Human Gastrointestinal

184 Protein (UHGP, 4744 genomes) database into peptides[30, 31]. Then each peptide is annotated

185 with a FuncTax score (Figure 1). Both the FuncTax and the detectability scores are used to

186 reduce the peptide database.

187 The FuncTax scores for each peptide in the MetaPepDetec are calculated using information

188 from the MetaPep database [32]. MetaPep is a core peptide database compiling peptides

189 previously identified in the published human gut metaproteomics studies. The information from

190 MetaPep was used to create a static table of COG relative functional abundances and a

191   sample-specific table of taxonomic relative abundances (method). We noted that despite

192   significant differences in the species composition of gut bacteria among different individuals,

193   their functions are remarkably similar[35]. Therefore, functional abundance hierarchy

194   information could act to estimate the likelihood of a protein being observed. We analyzed the

195   search results used to construct the MetaPep database which contains 2,134 raw files and 415

196   individuals[32] (Method). The functional ranking among various samples exhibits a strong

197   correlation (Supplementary Figure 2a and 2b). We observed a stable pattern in the functional

198   hierarchy of human gut bacteria: abundant functions consistently remain high, while scarce

199   functions persistently stay low across all samples (Supplementary Figure 2c and Supplementary

200   File 1). The sample-specific table of taxonomic relative abundances was generated by

201   searching the DIA data against MetaPep[32]. The identified peptides and their quantitation were

202   used to create the table (Method). The FuncTax score for each peptide is calculated by

203   multiplying the taxonomic score for its taxonomic annotation and the functional score of its

204   functional annotation. For peptides with identical sequences, their FuncTax scores were

205   aggregated thereby leading to shared peptides having a higher ranking.

206   In the last step, sample-specific reduced peptide database is generated by filtering

207   MetaPepDetec using the FuncTax score and the detectability score (Figure 1). The final search

208   of the DIA data is done against the reduced peptide database. To validate the efficiency of our

209   peptide ranking method, peptides were sorted by FuncTax score and partitioned into equal-

210   sized subsets based on their percentile rank (e.g., top 0-5%, 5-10%, ..., 35-40%). Each subset

211   was subjected to database searching with uniform parameters. We observed a decline in the

212   number of peptides identified as the percentile ranking of the subsets decreased (Figure 2). The

213   decreasing trend suggested our ranking method effectively prioritizes peptides with a higher

214   probability of detection.

215   **Optimized MetaDIA parameters reduces the database size**

216  We explored whether the number of microbes in the reduced peptide database, the threshold

217  for FuncTax score and the threshold for detectability score influenced the identification of

218  peptides. We explored the impacts of the parameters using the DIA data from 10 different

219  human gut microbiome samples previously reported[23] (Supplementary File 2, sample

220  information).

221  In particular, we first tested effect of the number of microbes (genomes) ranging from 50 to 150

222  and FuncTax score ranging from top 1% to 40% (Figure 3). We keep the detectability threshold

223  at the top 40% in this experiment which is suggested by the author of Deepdetect[30].

224  Interestingly, no mater how many genomes we choose, the size of the reduced peptide

225  database had the strongest effect on the number of identified peptides. The identification

226  number plateaued once the reduced peptide database size reached around 1.6 million entries

227  (Figure 3a and 3b, Supplementary Figure 3 and 4), corresponding to a FuncTax score threshold

228  of 40% for 50 genomes, 20% for 100 genomes and 15% for 150 genomes respectively. We

229  compared the three different reduced peptide databases, which led to consistent peptide

230  identification results (Figure 2c and 2d, Supplementary Figure 5 and 6). In our previous studies,

231  we observed that low-abundance species were underrepresented[8]. In this context, we chose

232  to focus on the top 50 genomes to prioritize high-abundance genomes. It is important to note

233  that this cut-off is a variable parameter that can be adjusted according to the specific objectives

234  of different studies.

235  Subsequently, we explored whether the detectability threshold impacted the number of peptides

236  identified. While the recommended threshold by the author of Deepdetect is 40%, we explored

237  thresholds ranging from 40% to 10%. We observed that a threshold of 25% was the point at

238  which the number of identifications began to decrease significantly (Figure 4a). However, both

239  the database size and the search time decreased substantially (Figure 4b). Comparing the

240  identification results at thresholds of 25% and 40%, we found a substantial overlap (Figure 4c).

241    Therefore, we selected a 25% threshold for detectability. Based on this analysis, we proceeded

242    with peptides ranking in the top 40% by FuncTax score and the top 25% by detectability. Given

243    that these two scores are entirely uncorrelated, applying both filters effectively reduced the

244    database to one-tenth of its original size (25% times 40%, Supplementary Figure 7). After

245    applying these optimized parameters, approximately 1 million peptide sequences remain in the

246    reduced database.

**MetaDIA maintains accuracy in DIA peptide identification**

248    We next explored whether the enrichment of high abundant and highly detectable peptides in

249    our reduced database impacted the accuracy of peptide identification when applying the false

250    discovery rate strategy. To evaluate this, we conducted a benchmark experiment using three

251    samples: a human gut microbiome sample (A), a *Blautia hydrogenotrophica* sample (C), and a

252    50:50 mixed sample of the two (B) (Figure 5a). *Blautia hydrogenotrophica* was selected due to

253    its absence in the microbiome sample used here and its minimal peptide overlap with the

254    microbiome sample. Each sample was subjected to triplicate DIA measurements. Sample A and

255    B were analyzed using the reduced peptide database generated by our workflow with optimized

256    parameters of top 40% FuncTax score and top 25% detectability score, whereas sample B and

257    C were searched against species-specific protein databases derived from NCBI (Genome

258    assembly ASM15797v1). In the first search against the peptide database, 32,624 unique

259    peptides were identified. Of these, 1,952 peptides also present in the *Blautia hydrogenotrophica*

260    database were excluded. Further, peptides unique to each sample were removed, leaving

261    27,830 peptides identified in both sample A and B. Ideally, the peptide abundance ratio between

262    samples A and B should approximate 2. In the second search, 14,821 unique peptides were

263    identified. Among these, 10,995 peptides were unique to the *Blautia hydrogenotrophica*

264    database and were found in both samples B and C. The expected ratio between samples B and

265    C should be around 0.5. We found whether using a protein database or a peptide database, the

266 ratios of peptides identified in both searches closely aligned with the expected values (Figure

267 5b). This suggests that the employment of our reduced peptide database does not significantly

268 affect the accuracy of peptide identification, thereby supporting its use in peptide identification

269 workflows with a controlled FDR.

270 **MetaDIA yields consistent peptide and protein identification results with DDA based**

271 **strategies.**

272 We next evaluated whether MetaDIA performed similarly to a conventional DDA-based library

273 for DIA data analysis. The DIA data and corresponding DDA-based library were obtained from a

274 published study[23]. We found that the MetaDIA provided identification numbers comparable to

275 those obtained through the DDA-based library (Figure 6a). Notably, in certain instances, such as

276 with samples 8 and 9, the MetaDIA surpassed DDA library in the number of identifications. The

277 initial step in our workflow involves searching the raw data against MetaPep, which leverages

278 the results from an open search strategy, thereby encompassing modified peptides not included

279 in the original database. Subsequent integration of peptides identified by MetaPep into a refined

280 peptide database resulted in a marked increase in identification rates (Figure 6a).

281 Over 50% of peptides identified from the DDA-based library were also identified by MetaDIA

282 (Figure 6b and Supplementary Figure 8). The divergence in unique identifications between the

283 two methods may be attributed to inherent differences between DDA acquisition and DIA

284 acquisition. Upon examining the quantification results of those peptides found by both methods,

285 we observed a significant consistency in the outcomes, with a Pearson coefficient above 0.9

286 (Figure 6d and Supplementary Figure 9). It is worth noting that the fragment ions used for

287 quantification in the DDA-based library correspond to actual DDA acquisitions. In contrast,

288 MetaDIA uses theoretical spectra that are predicted from peptide sequences. The high degree

289 of agreement between the quantification results underscores the reliability of the MS-Simulator

290 algorithm which is employed by DIA-NN for spectra prediction [14].

291    At the protein level, our findings revealed greater consistency in identification compared to the

292    peptide level (Figure 6c). Around 70% of proteins found by the DDA-based library can be found

293    by MetaDIA. The overlap on protein level reinforces the reliability of the identifications and

294    indicates that a significant subset of proteins is consistently identified by both methods despite

295    differences at the peptide level (Supplementary Figure 10). Proteins like

296    GYG000002545_00035 had greater sequence coverage and higher detection intensity with the

297    DDA library, while others like MGYG000002272_00452 showed higher coverage and intensity

298    with MetaDIA. Given that the quantification of a protein is derived from different subsets of

299    peptides in these two methods, we observed reduced consistency of quantification in the protein

300    level between the methods, as reflected by Pearson correlation coefficients of approximately 0.7

301    (Figure 6e and Supplementary Figure 11). However, it is important to note that in most

302    proteomic studies, the primary interest lies in the differential abundance of the same protein

303    across various samples. Therefore, it is crucial that we use the same fragment ions to quantify a

304    protein. In this regard, the inconsistencies in protein quantification between the two methods do

305    not undermine the utility of either approach. The substantial overlap in peptide and protein

306    identification by both methods suggests a robust cross-validation of both methods. Then we

307    annotated the proteins using COG accessions and calculated their relative abundances. Our

308    analysis revealed that approximately 90% of the COG accessions identified by the DDA-based

309    method were also covered by our MetaDIA (Figure 6c). Furthermore, the Pearson correlation

310    coefficient for the relative abundance of COG accessions exceeded 0.9, with a stronger

311    correlation for those COG accessions that were highly abundant (Figure 6f and Supplementary

312    Figure 12).

313    **MetaDIA provides taxonomic profiles highly similar to those obtained from searching**

314    **DDA-libraries.**

315    We verified whether both methods had a high degree of similarity in the taxonomic composition.

316    We did comparative analysis of microbiome composition across different taxonomic levels using

317    the result from both methods. Our findings indicate that there is a significant linear correlation

318    between the compositions identified by both methods, with the degree of correlation

319    strengthening at higher taxonomic levels (Figure 7a and b, Supplementary Figure 13). The two

320    methods showed remarkably consistent taxonomic composition at the genus level with a

321    Pearson coefficient above 0.98 across all the samples tested. Even Sample 9, which displayed

322    the lowest correlation, demonstrated a substantial degree of consistency between the two

323    methods. To underscore the consistency, we have provided a detailed visualization of the

324    taxonomic composition for Sample 9 (Figure 7c and d, Supplementary Figure 14)

325    The species compositions observed by MetaDIA in these ten samples differed significantly as

326    expected, indicating that our database and taxonomic analysis have the capability to identify a

327    diverse range of microbiota (Supplementary Figure 14 and Supplementary File 3: searching

328    result). The most abundant species identified in the ten samples have been previously reported

329    as high-abundance species in the human gut microbiome[36-40]. Except for *Phocaeicola dorei*

330    which were identified as the top species in sample 2, 5 and 10, the other top species were all

331    unique to each sample.

332    **MetaDIA is universally applicable to different types of DIA, including DIA-PASEF**

333    To further validate the versatility and applicability of our proposed metaproteomic workflow, we

334    extended our analysis to a diverse set of 79 DIA datasets obtained from a published study[23].

335    This dataset encompasses samples from 62 individuals, featuring replicate injections, quality

336    control (QC) samples, and pooled samples (Supplementary File 2: Sample information). We

337    applied MetaDIA to this extensive dataset and compared the results with the conventional DDA-

338    based approach. Remarkably, the number of peptides identified by both methods demonstrated

339    a close equivalence, reinforcing the robustness and universal applicability of our metaproteomic

340  workflow (Figure 8). Validating our method across diverse samples enhances confidence in its

341  effectiveness and consistency, demonstrating its potential for widespread adoption in

342  metaproteomics research.

343

## Discussion

345  We propose a novel workflow for DIA data analysis from human gut microbiome called MetaDIA.

346  The approach aims at prioritizing peptides with a higher likelihood of detection based on their

347  detectability, taxonomic and functional scores.

348  MetaDIA is entirely devoid of DDA, thereby circumventing the drawbacks of DDA-based

349  methods. Not only does this approach save time and resources, but it also enables the creation

350  of a tailored database for each sample. In contrast, DDA-based methods typically rely on a

351  single pooled sample to generate a library. For instance, Gomez *et al.*[25] used a pooled sample

352  to represent 12 individual mice, while Sun *et al.*[23] did so for a cohort of 62 individuals.

353  However, such a pooled sample may not effectively represent every sample. In our study, the

354  ten samples showed highly diverse taxonomic composition (Supplementary Figure 13). To

355  increase the sampling depth for the pooled sample, Sun *et al.*[23] had to fractionate the pooled

356  sample into 30 portions and Gomez *et al.*[25] repeatedly injected the pooled sample 10 times.

357  Moreover, utilizing a static library to search various samples may potentially compromise the

358  accuracy of peptide identification, as it includes peptides from the pooled samples that are

359  absent in the specific sample under investigation.

360  In MetaDIA, we pre-defined the range of genomes for each microbiome sample (50 genomes in

361  this study). This approach not only enabled us to narrow the search space but also to mitigate

362  the issues associated with protein inference that arise from common peptides. When assigning

363  peptides to proteins, we confined our consideration to the genomes within the predefined range

364 rather than the entire dataset. This strategy significantly reduced the incidence of common

365 peptides.

366 Although MetaDIA is currently focused on the human gut microbiome, we foresee that it can be

367 extended to other types of microbiomes, such as those in animal intestines, environmental

368 microbiomes when using an appropriate bait database. A database similar to MetaPep could be

369 constructed for other microbiomes.

370

## Conclusion
371

372 In conclusion, we introduced a new strategy to prioritize peptides with a high probability of

373 detection. This strategy simulates protein digestion procedures in silico and uses taxonomic and

374 functional information to infer the peptide abundance. MetaDIA is a fully DDA-free workflow and

375 provides a user interface to change the different parameters. We compared the performance of

376 MetaDIA with the DDA-based library and observed a high degree of consistency. We further

377 validated our method across a DIA-PASEF dataset with 79 samples, thereby confirming its wide

378 applicability. We believe that our approach will help the application of DIA in metaproteomics.

379

## Author information
380

### Authors and Affiliations
381

382 School of Pharmaceutical Sciences, Faculty of Medicine, University of Ottawa, Ottawa, ON K1H

383 8M5, Canada

384 Haonan Duan, Zhibin Ning, Zhongzhi Sun & Daniel Figeys

385    Department of Biochemistry, Microbiology and Immunology, Faculty of Medicine, University of

386    Ottawa, Ottawa, ON K1H 8M5, Canada

387    Haonan Duan & Zhongzhi Sun

388    Westlake Center for Intelligent Proteomics, Westlake Laboratory of Life Sciences and

389    Biomedicine, Hangzhou, Zhejiang Province, 310030, China

390    Tiannan Guo & Yingying Sun

391    School of Medicine, School of Life Sciences, Westlake University, Hangzhou, Zhejiang Province,

392    310030, China

393    Tiannan Guo & Yingying Sun

394    Research Center for Industries of the Future, Westlake University, 600 Dunyu Road, Hangzhou,

395    Zhejiang, 310030, China

396    Tiannan Guo & Yingying Sun

## Contributions

398    HD, ZN, ZS, TG, YS and DF designed the study. HD performed the experiments. HD, ZN, DF

399    analyzed the data. TG and YS shared the raw data. HD and DF wrote the manuscript. All

400    authors read and approved the final manuscript.

## Corresponding authors

402    Correspondence to Daniel Figeys

403

## **Acknowledgements**

405   The authors acknowledge the assistance of OpenAI's ChatGPT in grammar correction and

406   clarification of this article.

407

## Funding

409   This work was funded by the Natural Sciences and Engineering Research Council of Canada

410   (NSERC) discovery grant to D.F. H.D. was funded by a stipend from the NSERC CREATE in

411   Technologies for Microbiome Science and Engineering (TECHNOMISE) Program.

412

## Ethics declarations

### Competing interests

415   DF is the founder of MedBiome Inc. a microbiome nutrition and therapeutic company.

### Consent for publication

417   Not applicable

### Competing interests

419   The Human stool was collected from a healthy adult volunteer at the University of Ottawa,

420   Ottawa, ON, CAN. The protocol (# 20160585-01H) was approved by Ottawa Health Science

421   Network Research Ethics.

422

## Figure Captions

424

**Figure 1**. The flowchart for the MetaDIA. All proteins in Unified Human Gastrointestinal Protein (UHGP) database were firstly in silico digested into peptides. The detectability of each peptide was predicated by DeepDetect algorithm. Following this prediction, each peptide was assigned a functional score and a taxonomic score, derived from a predetermined functional relative abundance table and a sample-specific taxonomic relative abundance table, respectively (method). The FuncTax score was calculated by multiplying the two scores. For peptides with identical sequence, their FuncTax scores were aggregated to prioritize shared peptides; the maximum of their detectability scores was used to ensure the inclusion of all potential peptides. The detectability and FuncTac scores are both used for filtering peptides. The reduced peptide database was used for a second search.

435

**Figure 2**. The number of peptides identified from each subset. Ten samples were tested in the experiment. For constructing the peptide database, the top 100 genomes were considered; the detectability threshold was set at 40%. Each subset contains around 400,000 peptides. Peptide identification was performed by DIA-NN under same conditions. The maximum identification from the last subset was heighted in the figure.

**Figure 3**. Optimization for genome number and FuncTax score (Sample 1 was shown. For the other samples, please see Supplementary figures). Peptides from n (50, 100, 150) genomes were ranked by the FuncTax score and top x% (1-40 for 50 and 100 genomes; 1–35 for 150 genomes) peptides was used as database. (a) Number of identified peptides against database percent. (b) Number of identified peptides against database size. The inflection point has been highlighted with a red box. (c) The overlap of the reduced peptide database and (d) identified peptide when taking top 40% peptides for 50 genomes, top 20% for 100 genomes and top 15%

449  for 150 genomes as database. Peptide identification was performed by DIA-NN under same

450  conditions.



451

452  **Figure 4**. Optimization for detectability threshold. (a) The number of peptides identified and (b)

453  the searching time under detectability threshold from 10% to 40%. (c) The overlap of peptides

454  identified by top 25% and top 40% of the database. Peptide identification was performed by

455  DIA-NN under same conditions.

**Figure 5**. Benchmark experiment for peptide identification. (a) The experimental design. Each sample was subjected to triple-run measurements (b) Log-transformed ratios are plotted as a function of peptide intensity for n = 27,830 (green) microbiome peptides and n = 10,995 (purple) *Blautia hydrogenotrophica* peptides. The point density for ratio was plotted at right. Dashed lines indicate the expected ratio. Peptide identification was performed by DIA-NN under same conditions. The intensities derived from various charge states of the same peptide were aggregated.

**Figure 6**. Comparison between the DDA-based method and DDA-free method. (a) The peptide identified by each method. (b) The overlap of peptide identified in sample 1 by each method. (c) Coverage of peptides, proteins and cog accessions identified by DDA-based method with those found using DDA-free method. The intensity correlation of the overlapped peptides (d), proteins (e) and COG accessions in sample 1. The dashed line indicates y = x. For DDA-based method, the peptides identified as derived from human proteins are removed.

**Figure 7**. Comparison of the taxonomic composition between the DDA-based method and DDA-free method. Pearson correlation (a) and Bray-Curtis distance (b) analysis between DDA-based method and DDA-free method on different taxonomic levels from Phylum to Species. The relative taxonomic abundance was used for the analysis. In the correlation analysis, taxonomic categories that were unique to one method were imputed with a value of zero. The taxonomic composition (Phylum to Family) of sample 9 derived from DDA-based method (c) and DDA-free method (d).
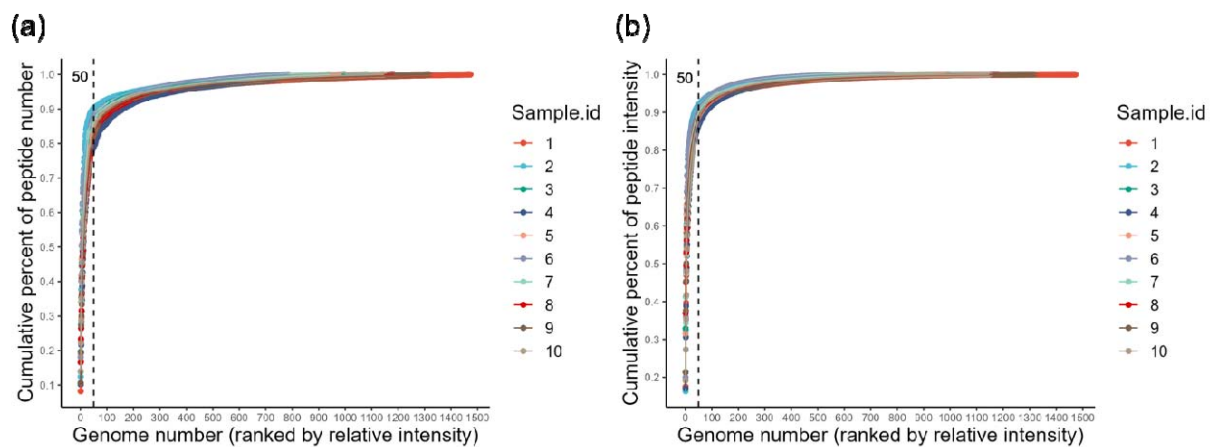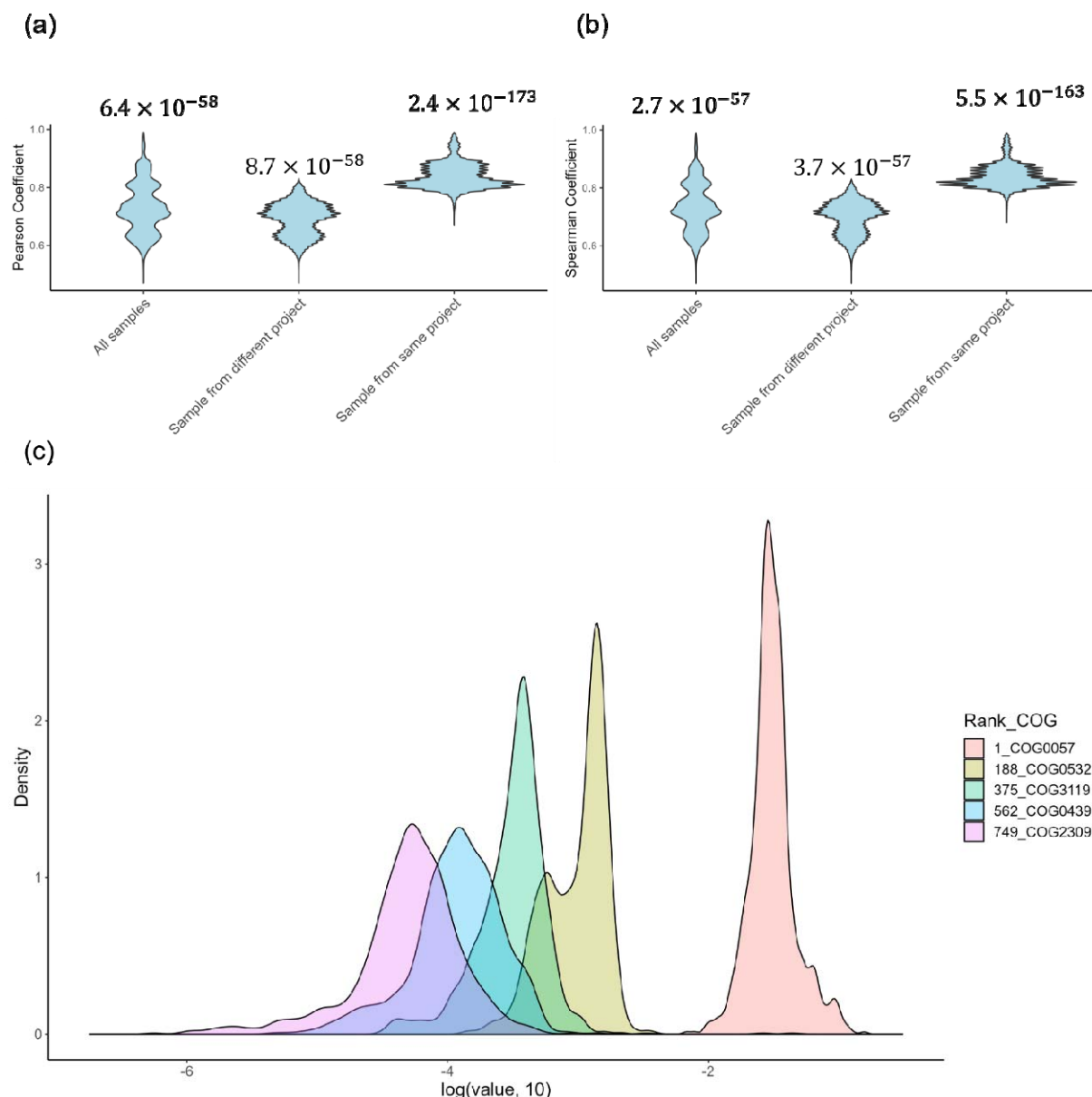
479

**Figure 8**. Number of peptides identified by both methods with mean value from 79 diaPASEF samples. For DDA-based method, the peptides identified as derived from human proteins are removed.
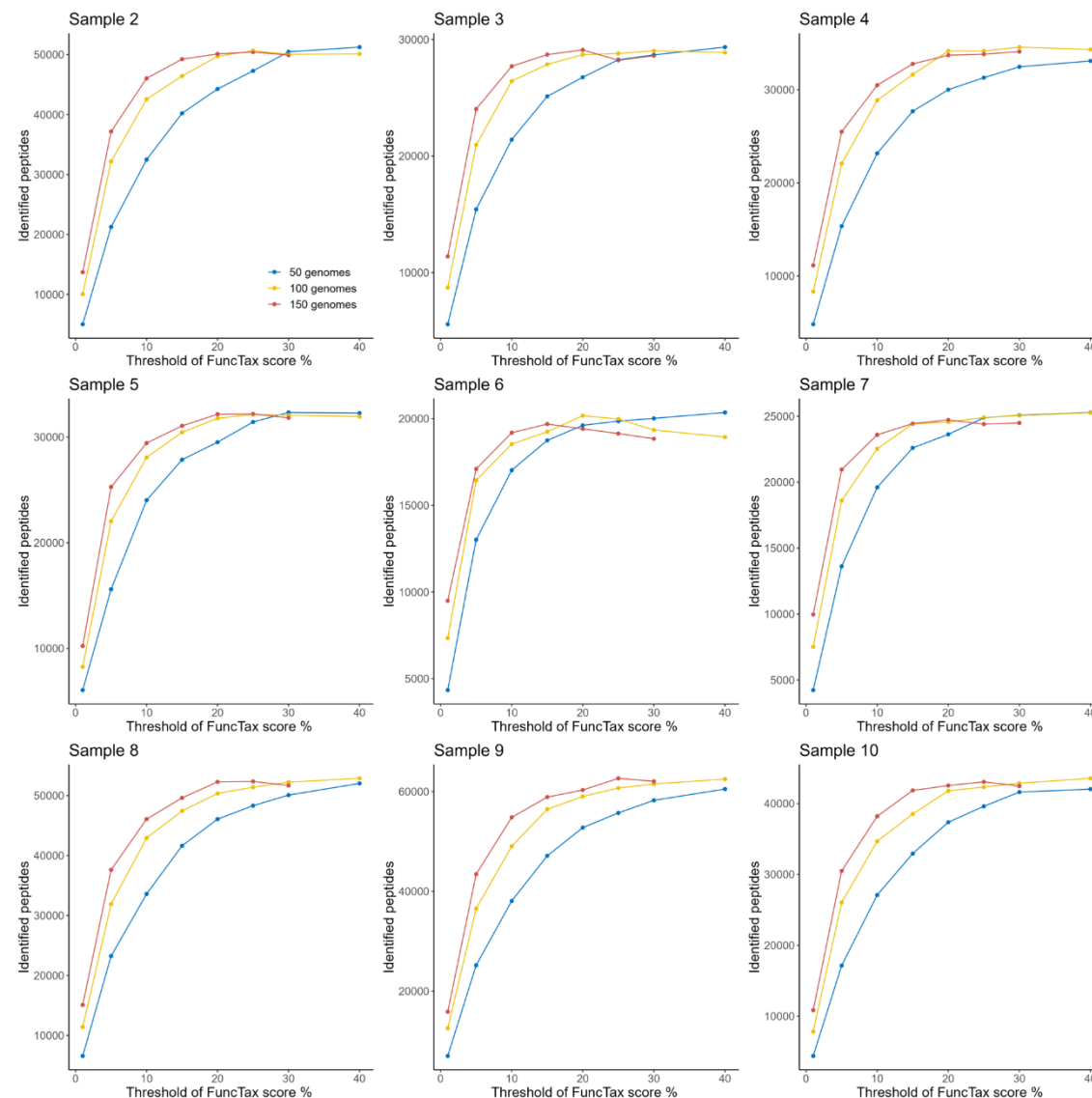


483

**Supplementary Figure 1.** Cumulative contribution of identified peptides (a) and intensity (b). Figure were plotted against the number of genomes. Genomes are ordered by decreasing peptide count.
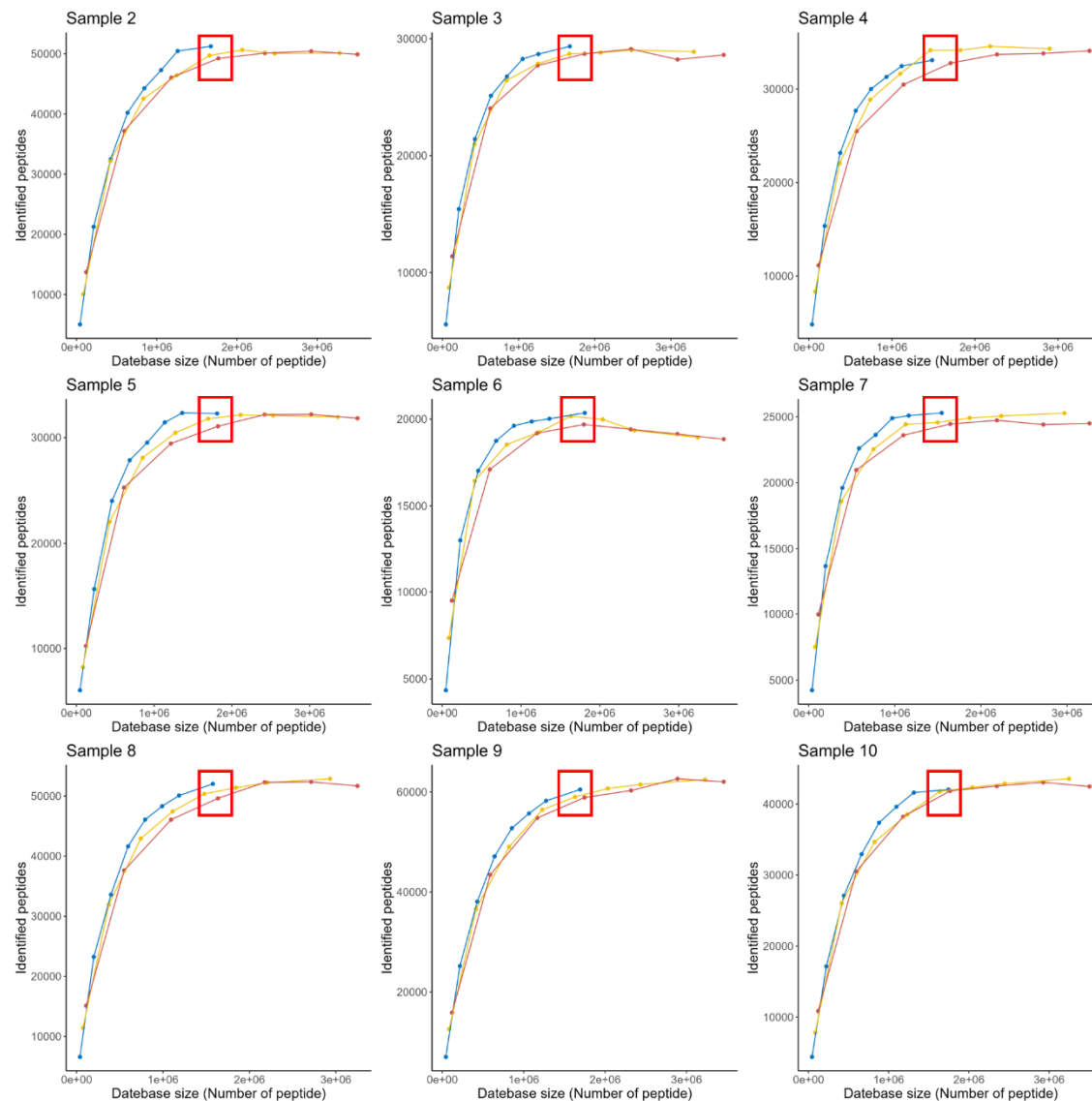
**Supplementary Figure 2**. The functional correlation across samples in MetaPep. The Pearson coefficient (a) and Spearman coefficient (b) on COG accessions. Sample pairs from same project are likely from the same individuals and plotted in different groups. The average p value are plotted above each groups. (c) The distribution of relative abundance of COG accession across samples in MetaPep. The legends shows the rank and name of COG accessions. COG accessions present in more than 95% of the samples were retained. A total of 750 COG accessions remained and were subsequently sorted based on their functional abundance
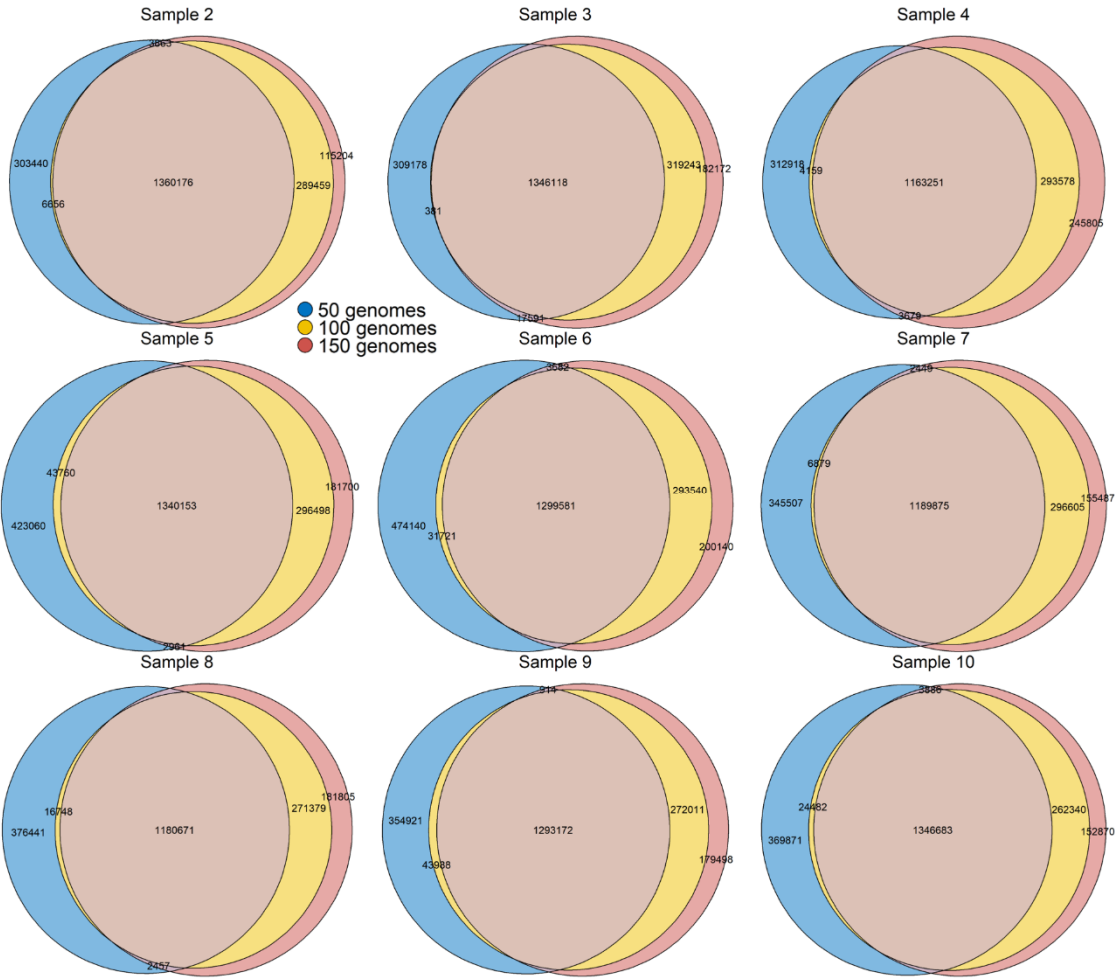
495    scores. Five COG accessions were selected for density plot at evenly spaced intervals from this

496    ordered list.



**Supplementary Figure 3**. Optimization for genome number and FuncTax score (Sample 2-10).
Number of identified peptides against database percent. Peptides from n (50, 100, 150)
genomes were ranked by the FuncTax score and top x% (1-40 for 50 and 100 genomes; 1–35
for 150 genomes) peptides was used as database. Peptide identification was performed by DIA-
NN under same conditions.

**Supplementary Figure 4**. Optimization for genome number and FuncTax score (Sample 2-10). Number of identified peptides against database size. The inflection point has been highlighted with a red box. Peptides from n (50, 100, 150) genomes were ranked by the FuncTax score and top x% (1-40 for 50 and 100 genomes; 1–35 for 150 genomes) peptides was used as database. Peptide identification was performed by DIA-NN under same conditions.
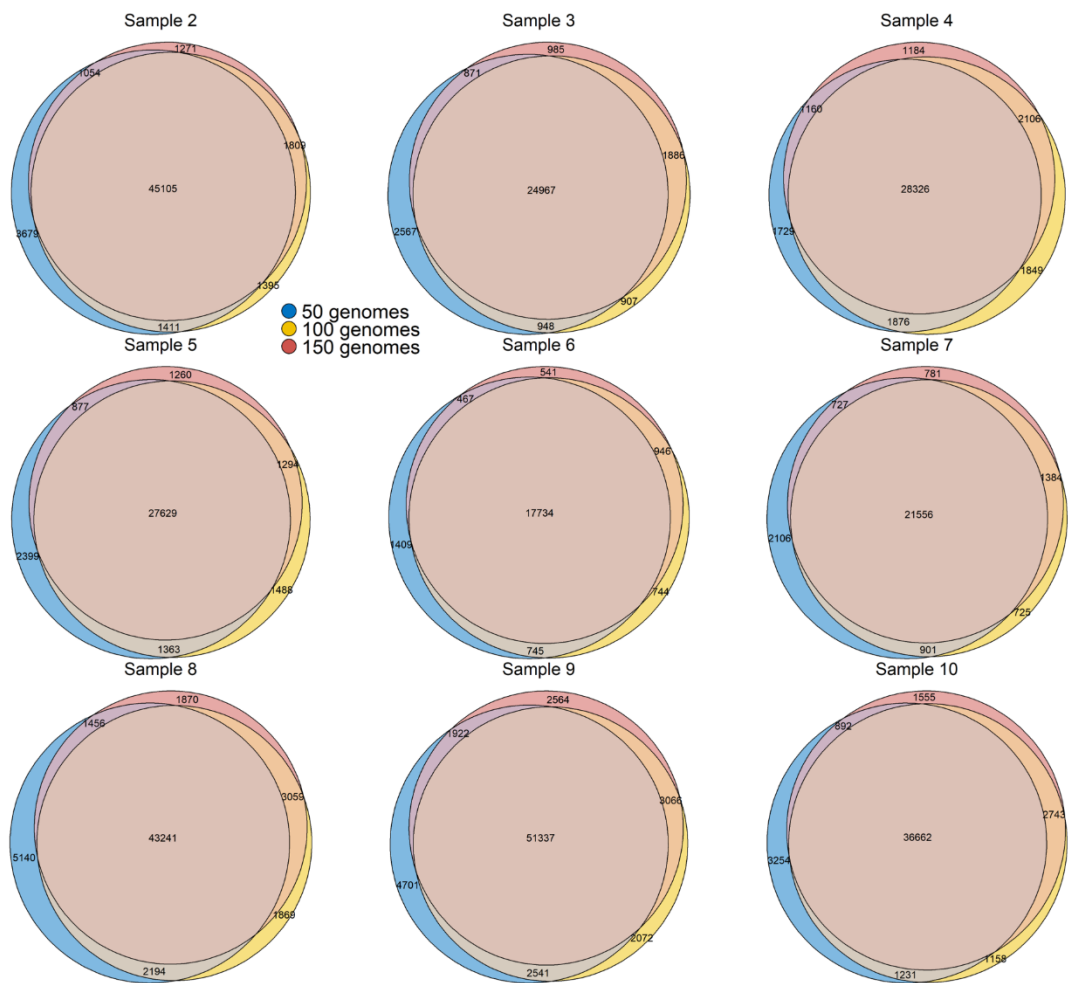
509

**Supplementary Figure 5**. Optimization for genome number and FuncTax score (Sample 2-10).

The overlap of reduced peptide database when taking top 40% peptides for 50 genomes, top 20%

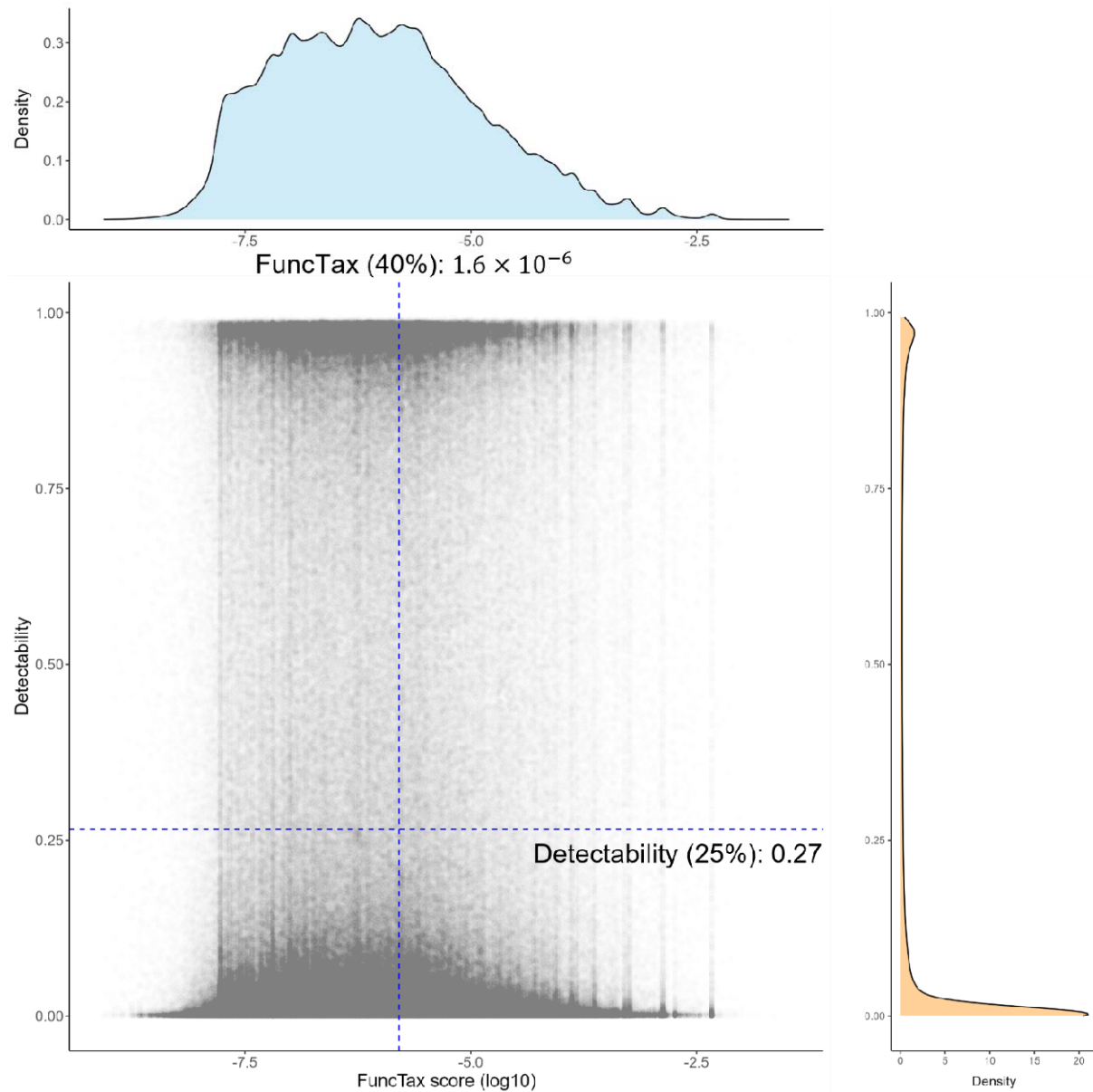for 100 genomes and top 15% for 150 genomes as database.

513

**Supplementary Figure 6**. Optimization for genome number and FuncTax score (Sample 2-10). The overlap of identified peptide when taking top 40% peptides for 50 genomes, top 20% for 100 genomes and top 15% for 150 genomes as database. Peptide identification was performed by DIA-NN under same conditions.

518

**Supplementary Figure 7.** The distribution of detectability score and FuncTax score (Sample 1, top 50 genomes). The dotted blue line shows the cutoff for detectability (top 25%) and FuncTax (top 40%).

**Supplementary Figure 8**. The overlap of peptide identified by each DDA-based method and DDA-free method (Sample 2-10)

**Supplementary Figure 9**. The intensity correlation of the overlapped peptides found by both DDA-based and DDA-free method (Sample 2-10). The dashed line indicates y = x. For DDA-based method, the peptides identified as derived from human proteins are removed.

**Supplementary Figure 10**. The sequence coverage of representative proteins that exhibit the largest relative difference in intensity DDA-based and DDA-free method.

**Supplementary Figure 11**. The intensity correlation of the overlapped proteins found by both DDA-based and DDA-free method (Sample 2-10). The dashed line indicates y = x.



**Supplementary Figure 12**. The intensity correlation of the overlapped COG accessions found by both DDA-based and DDA-free method (Sample 2-10). The dashed line indicates y = x.

**Supplementary Figure 13**. Correlation analysis between DDA-based method and DDA-free method on different taxonomic levels from Species to Phylum. The relative abundance was used for the analysis. Taxonomic categories that were unique to one method were imputed with a value of zero.



**Supplementary Figure 14.** Comparative Taxonomic Composition of the Microbiome at Levels from Phylum to Species. (a) DDA-Based. (b) DDA-Free.

546

**Supplementary Figure 15.** UpSet plot illustrating the overlap in genomes identified by the DDA-free method across ten microbiome samples. The top 40 intersections were plotted.

549

# Reference

551  1.  Human Microbiome Project C: **Structure, function and diversity of the healthy human**

552      **microbiome**. *Nature* 2012, **486**(7402):207-214.

553  2.  Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, Prill RJ, Tripathi A, Gibbons SM,

554      Ackermann G *et al*: **A communal catalogue reveals Earth's multiscale microbial diversity**.

555      *Nature* 2017, **551**(7681):457-463.

556   3.      Wilmes P, Bond PL: **Metaproteomics: studying functional gene expression in microbial**

557          **ecosystems**. *Trends Microbiol* 2006, **14**(2):92-97.

558   4.      Jagtap P, Goslinga J, Kooren JA, McGowan T, Wroblewski MS, Seymour SL, Griffin TJ: **A two-step**

559          **database search method improves sensitivity in peptide sequence matches for**

560          **metaproteomics and proteogenomics studies**. *Proteomics* 2013, **13**(8):1352-1357.

561   5.      Zhang X, Ning Z, Mayne J, Moore JI, Li J, Butcher J, Deeke SA, Chen R, Chiang CK, Wen M *et al*:

562          **MetaPro-IQ: a universal metaproteomic approach to studying human and mouse gut**

563          **microbiota**. *Microbiome* 2016, **4**(1):31.

564   6.      Cheng K, Ning Z, Zhang X, Li L, Liao B, Mayne J, Figeys D: **MetaLab 2.0 Enables Accurate Post-**

565          **Translational Modifications Profiling in Metaproteomics**. *J Am Soc Mass Spectrom* 2020,

566          **31**(7):1473-1482.

567   7.      Cheng K, Ning Z, Li L, Zhang X, Serrana JM, Mayne J, Figeys D: **MetaLab-MAG: A Metaproteomic**

568          **Data Analysis Platform for Genome-Level Characterization of Microbiomes from the**

569          **Metagenome-Assembled Genomes Database**. *J Proteome Res* 2023, **22**(2):387-398.

570   8.      Duan H, Cheng K, Ning Z, Li L, Mayne J, Sun Z, Figeys D: **Assessing the Dark Field of**

571          **Metaproteome**. *Anal Chem* 2022, **94**(45):15648-15654.

572   9.      Gillet LC, Navarro P, Tate S, Rost H, Selevsek N, Reiter L, Bonner R, Aebersold R: **Targeted data**

573          **extraction of the MS/MS spectra generated by data-independent acquisition: a new concept**

574          **for consistent and accurate proteome analysis**. *Mol Cell Proteomics* 2012, **11**(6):O111 016717.

575   10.    Collins BC, Hunter CL, Liu Y, Schilling B, Rosenberger G, Bader SL, Chan DW, Gibson BW, Gingras

576          AC, Held JM *et al*: **Multi-laboratory assessment of reproducibility, qualitative and quantitative**

577          **performance of SWATH-mass spectrometry**. *Nat Commun* 2017, **8**(1):291.

578    11.    Vowinckel J, Zelezniak A, Bruderer R, Mulleder M, Reiter L, Ralser M: **Cost-effective generation**

579            **of precise label-free quantitative proteomes in high-throughput by microLC and data-**

580            **independent acquisition**. *Sci Rep* 2018, **8**(1):4346.

581    12.    Meier F, Brunner AD, Frank M, Ha A, Bludau I, Voytik E, Kaspar-Schoenefeld S, Lubeck M, Raether

582            O, Bache N *et al*: **diaPASEF: parallel accumulation-serial fragmentation combined with data-**

583            **independent acquisition**. *Nat Methods* 2020, **17**(12):1229-1236.

584    13.    Guzman UH, Martinez-Val A, Ye Z, Damoc E, Arrey TN, Pashkova A, Renuse S, Denisov E, Petzoldt

585            J, Peterson AC *et al*: **Ultra-fast label-free quantification and comprehensive proteome coverage**

586            **with narrow-window data-independent acquisition**. *Nat Biotechnol* 2024.

587    14.    Sun S, Yang F, Yang Q, Zhang H, Wang Y, Bu D, Ma B: **MS-Simulator: predicting y-ion intensities**

588            **for peptides with two charges based on the intensity ratio of neighboring ions**. *J Proteome Res*

589            2012, **11**(9):4509-4516.

590    15.    Zeng WF, Zhou XX, Zhou WJ, Chi H, Zhan J, He SM: **MS/MS Spectrum Prediction for Modified**

591            **Peptides Using pDeep2 Trained by Transfer Learning**. *Anal Chem* 2019, **91**(15):9724-9731.

592    16.    Bouwmeester R, Gabriels R, Hulstaert N, Martens L, Degroeve S: **DeepLC can predict retention**

593            **times for peptides that carry as-yet unseen modifications**. *Nat Methods* 2021, **18**(11):1363-

594            1369.

595    17.    Ma C, Ren Y, Yang J, Ren Z, Yang H, Liu S: **Improved Peptide Retention Time Prediction in Liquid**

596            **Chromatography through Deep Learning**. *Anal Chem* 2018, **90**(18):10881-10888.

597    18.    Al Musaimi O, Valenzo OMM, Williams DR: **Prediction of peptides retention behavior in**

598            **reversed-phase liquid chromatography based on their hydrophobicity**. *J Sep Sci* 2023,

599            **46**(2):e2200743.

600    19.    Cox J: **Prediction of peptide mass spectral libraries with machine learning**. *Nat Biotechnol* 2023,

601            **41**(1):33-43.

602  20.  Demichev V, Messner CB, Vernardis SI, Lilley KS, Ralser M: **DIA-NN: neural networks and**

603  **interference correction enable deep proteome coverage in high throughput**. *Nat Methods* 2020,

604  **17**(1):41-44.

605  21.  Demichev V, Szyrwiel L, Yu F, Teo GC, Rosenberger G, Niewienda A, Ludwig D, Decker J, Kaspar-

606  Schoenefeld S, Lilley KS *et al*: **dia-PASEF data analysis using FragPipe and DIA-NN for deep**

607  **proteomics of low sample amounts**. *Nat Commun* 2022, **13**(1):3944.

608  22.  Sinitcyn P, Hamzeiy H, Salinas Soto F, Itzhak D, McCarthy F, Wichmann C, Steger M, Ohmayer U,

609  Distler U, Kaspar-Schoenefeld S *et al*: **MaxDIA enables library-based and library-free data-**

610  **independent acquisition proteomics**. *Nat Biotechnol* 2021.

611  23.  Sun Y, Xing Z, Liang S, Miao Z, Zhuo L-b, Jiang W, Zhao H, Gao H, Xie Y, Zhou Y: **metaExpertPro: a**

612  **computational workflow for metaproteomics spectral library construction and data-**

613  **independent acquisition mass spectrometry data analysis**. *bioRxiv* 2023:2023.2011.

614  2029.569331.

615  24.  Aakko J, Pietila S, Suomi T, Mahmoudian M, Toivonen R, Kouvonen P, Rokka A, Hanninen A, Elo LL:

616  **Data-Independent Acquisition Mass Spectrometry in Metaproteomics of Gut Microbiota-**

617  **Implementation and Computational Analysis**. *J Proteome Res* 2020, **19**(1):432-436.

618  25.  Gomez-Varela D, Xian F, Grundtner S, Sondermann JR, Carta G, Schmidt M: **Increasing taxonomic**

619  **and functional characterization of host-microbiome interactions by DIA-PASEF metaproteomics**.

620  *Front Microbiol* 2023, **14**:1258703.

621  26.  Pietilä S, Suomi T, Elo LL: **Introducing untargeted data-independent acquisition for**

622  **metaproteomics of complex microbial samples**. *ISME Communications* 2022, **2**(1):1-8.

623  27.  Tsou CC, Avtonomov D, Larsen B, Tucholska M, Choi H, Gingras AC, Nesvizhskii AI: **DIA-Umpire:**

624  **comprehensive computational framework for data-independent acquisition proteomics**. *Nat*

625  *Methods* 2015, **12**(3):258-264, 257 p following 264.

626   28.   Stamboulian M, Li S, Ye Y: **Using high-abundance proteins as guides for fast and effective**

627         **peptide/protein identification from human gut metaproteomic data**. *Microbiome* 2021, **9**(1).

628   29.   Yang J, Pu J, Lu S, Bai X, Wu Y, Jin D, Cheng Y, Zhang G, Zhu W, Luo X *et al*: **Species-Level Analysis**

629         **of Human Gut Microbiota With Metataxonomics**. *Front Microbiol* 2020, **11**:2029.

630   30.   Yang J, Cheng Z, Gong F, Fu Y: **DeepDetect: Deep Learning of Peptide Detectability Enhanced by**

631         **Peptide Digestibility and Its Application to DIA Library Reduction**. *Anal Chem* 2023,

632         **95**(15):6235-6243.

633   31.   Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, Pollard KS, Sakharova E, Parks

634         DH, Hugenholtz P *et al*: **A unified catalog of 204,938 reference genomes from the human gut**

635         **microbiome**. *Nat Biotechnol* 2021, **39**(1):105-114.

636   32.   Sun Z, Ning Z, Cheng K, Duan H, Wu Q, Mayne J, Figeys D: **MetaPep: A core peptide database for**

637         **faster human gut metaproteomics database searches**. *Comput Struct Biotechnol J* 2023,

638         **21**:4228-4237.

639   33.   Tyanova S, Temu T, Cox J: **The MaxQuant computational platform for mass spectrometry-based**

640         **shotgun proteomics**. *Nat Protoc* 2016, **11**(12):2301-2319.

641   34.   Zhang X, Li L, Mayne J, Ning Z, Stintzi A, Figeys D: **Assessing the impact of protein extraction**

642         **methods for human gut metaproteomics**. *J Proteomics* 2018, **180**:120-127.

643   35.   Li L, Wang T, Ning Z, Zhang X, Butcher J, Serrana JM, Simopoulos CMA, Mayne J, Stintzi A, Mack

644         DR *et al*: **Revealing proteome-level functional redundancy in the human gut microbiome using**

645         **ultra-deep metaproteomics**. *Nat Commun* 2023, **14**(1):3428.

646   36.   Davis-Richardson AG, Ardissone AN, Dias R, Simell V, Leonard MT, Kemppainen KM, Drew JC,

647         Schatz D, Atkinson MA, Kolaczkowski B *et al*: **Bacteroides dorei dominates gut microbiome prior**

648         **to autoimmunity in Finnish children at high risk for type 1 diabetes**. *Front Microbiol* 2014,

649         **5**:678.

650    37.    Hosomi K, Saito M, Park J, Murakami H, Shibata N, Ando M, Nagatake T, Konishi K, Ohno H,

651            Tanisawa K *et al*: **Oral administration of Blautia wexlerae ameliorates obesity and type 2**

652            **diabetes via metabolic remodeling of the gut microbiota**. *Nat Commun* 2022, **13**(1):4477.

653    38.    Scher JU, Sczesnak A, Longman RS, Segata N, Ubeda C, Bielski C, Rostron T, Cerundolo V, Pamer

654            EG, Abramson SB *et al*: **Expansion of intestinal Prevotella copri correlates with enhanced**

655            **susceptibility to arthritis**. *Elife* 2013, **2**:e01202.

656    39.    Lan PTN, Sakamoto M, Sakata S, Benno Y: **Bacteroides barnesiae sp. nov., Bacteroides**

657            **salanitronis sp. nov. and Bacteroides gallinarum sp. nov., isolated from chicken caecum**. *Int J*

658            *Syst Evol Microbiol* 2006, **56**(Pt 12):2853-2859.

659    40.    Ferreira-Halder CV, Faria AVS, Andrade SS: **Action and function of Faecalibacterium prausnitzii**

660            **in health and disease**. *Best Pract Res Clin Gastroenterol* 2017, **31**(6):643-648.

661