

# 1 **MetaDIA: A DDA-free Database Reduction Strategy for DIA Human Gut**

## 2 **Metaproteomics**

3 Haonan Duan (段浩楠)<sup>1,2</sup>, Zhibin Ning (宁志斌)<sup>2</sup>, Zhongzhi Sun (孙中智)<sup>2</sup>, Tiannan Guo (郭天  
4 南)<sup>3,4,5</sup>, Yingying Sun (孙莹莹)<sup>3,4,5</sup>, Daniel Figeys<sup>2,6,7,\*</sup>

5  
6 <sup>1</sup> *Nanjing Women and Children's Healthcare Hospital, Women's Hospital of Nanjing Medical*  
7 *University, Nanjing 210004, China*

8 <sup>2</sup> *School of Pharmaceutical Sciences, Faculty of Medicine, University of Ottawa, Ottawa ON K1H*  
9 *8M5, Canada*

10 <sup>3</sup> *Westlake Center for Intelligent Proteomics, Westlake Laboratory of Life Sciences and*  
11 *Biomedicine, Hangzhou 310030, China*

12 <sup>4</sup> *School of Medicine, School of Life Sciences, Westlake University, Hangzhou 310030, China*

13 <sup>5</sup> *Research Center for Industries of the Future, Westlake University, Hangzhou 310030, China*

14 <sup>6</sup> *Quadram Institute Bioscience, Norwich Research Park, Norwich NR4 7UQ, UK*

15 <sup>7</sup> *University of East Anglia, Norwich NR4 7TJ, UK*

16

17 \* Corresponding author.

18 E-mail: Daniel.Figeys@quadram.ac.uk (Figeys D).

19

20 **Running title:** *Duan H et al / MetaDIA: Database Reduction*

21

22 The number of words: 9027

23 The number of references: 44

24 The number of figures: 8

25 The number of tables: 0

26 The number of supplementary files: 3

27 The number of supplementary figures: 20

28 The number of supplementary tables: 2

29 The number of letters in the article title: 68

30 The number of letters in the running title: 33

31 The number of keywords: 5  
32 The number of words in abstract: 199

33

## 34 **Abstract**

35 Microbiomes, especially within the gut, are complex and may comprise hundreds of species. The  
36 identification of peptides in metaproteomics presents a substantial challenge, as it involves  
37 matching peptides to mass spectra within an enormous search space for complex and unknown  
38 samples. This poses difficulties for both the accuracy and the speed of identification. Specifically,  
39 analysis of data-independent acquisition (DIA) datasets has relied on libraries constructed from  
40 prior data-dependent acquisition (DDA) results. However, this method is resource-intensive,  
41 consumes samples, and limits identification to peptides previously identified. These limitations  
42 restrict the application of DIA in metaproteomics research. We introduced a novel strategy to  
43 reduce the search space by utilizing species abundance and functional abundance information from  
44 the microbiome to score each peptide and prioritize those most likely to be detected. Using this  
45 strategy, we have developed and optimized a workflow called MetaDIA for the analysis of  
46 microbiome data generated by DIA, which operates independently of DDA assistance. Our  
47 approach successfully created a smaller, yet sufficient database for DIA data search in  
48 metaproteomics. The results demonstrated strong consistency with the traditional DDA-based  
49 library approach at both protein and functional levels. MetaDIA is readily accessible as an open-  
50 source project hosted on GitHub (<https://github.com/northomics/MetaDIA>).

51

52 **Key words:** Metaproteomics; Human gut microbiome; Data independent acquisition; Data-  
53 dependent acquisition-free; diaPASEF

54

55

## 56 **Introduction**

57 The microbiome encompasses a diverse array of microorganisms residing in different organisms,  
58 ecosystems, and environmental settings such as the human body, animals, plants, soil, water bodies,  
59 and various ecological niches [1,2]. Metaproteomics serves as a tool for understanding the roles of  
60 proteins within these microbial communities [3]. Mass spectrometry-based proteomics aims to  
61 study all proteins in a sample. However, applying these techniques to the microbiome is challenged

62 by its complexity. Without prior knowledge of the microbes present in a sample, metaproteomics  
63 relies on searching mass spectra against a large database, making the task of matching peptides  
64 and spectra notably challenging. Employing an iterative search strategy significantly reduces the  
65 search complexity in which the final search is against a database generated from previous searching  
66 results [4,5]. The iterative strategy has been successfully used for the data acquired by data-  
67 dependent acquisition (DDA) mode [6,7]. Unfortunately, in DDA mode, only the most abundant  
68 precursor ions are selected for further inquiry, and lower abundant ones are overlooked [8].

69 In contrast, data-independent acquisition (DIA) uses a set of precursor isolation windows to  
70 collect all the fragment ions indiscriminately [9]. It has shown remarkable robustness, sensitivity,  
71 and reproducibility with fewer missing values [10]. DIA can be coupled with microLC enabling  
72 high-throughput analysis [11]. This makes it particularly suitable for conducting large-scale  
73 analyses. The diaPASEF [12] method integrates ion mobility separation with the DIA workflow,  
74 adding a fourth dimension of analyzing ion mobility to the traditional three-dimensional data set.  
75 This not only enriches the structural information of analytes but also enhances ion utilization  
76 efficiency, leveraging the linear relation between ion mobility and mass-to-charge ratio. Another  
77 improvement in mass spectrometer scanning speed enables the utilization of smaller isolation  
78 windows in DIA, termed as narrow-window DIA [13]. This approach achieves comprehensive  
79 peptide precursor coverage and high quantitative precision and accuracy. In bioinformatics, the  
80 development of prediction software for peptide properties (theoretically predicted spectrum  
81 [14,15], retention times [16–18]) enables the querying of DIA datasets without dependence on  
82 libraries generated by DDA. Those predicted libraries even showed better performance than the  
83 measured libraries [19]. Moreover, DIA-specific searching software such as DIA-NN [20,21],  
84 MaxDIA [22], and Spectronaut have shown reliable results for the identification and quantification  
85 of peptides. The above advantages make DIA increasingly popular in proteomics. However, it is  
86 noteworthy that the benefits conferred by these techniques have not yet been fully extended to the  
87 field of metaproteomics. The main reason is that the inherent complexity of DIA data requires a  
88 much more constrained searching space compared with DDA data. To date, only a few studies  
89 employed DIA for metaproteomics analysis, and the majority of them were compelled to use a  
90 spectral library derived from DDA data [23–25]. The DDA-derived method involves creating a  
91 spectral library from DDA runs for each sample, which is then used to interpret complex mass  
92 spectra from subsequent analyses. This approach requires multiple sample aliquots, extensive mass

93 spectrometry resources and is limited to detecting peptides previously identified by DDA. To  
94 surmount these limitations, several endeavors have been undertaken. Gladiator [26] uses DIA-  
95 Umpire [27] to assemble pseudo-DDA spectra from DIA data for microbiome samples. The  
96 method does not require a DDA-based spectral library for its operation; however, it still relies on  
97 spectrum-centric algorithms and does not fully exploit the potential advantages of DIA data.  
98 Dumas et al. [28] have refined the database by selecting the most abundant species from DDA data  
99 analysis. This approach, while reducing the size of database, preserves the entirety of species'  
100 proteins, thereby retaining peptides that may be overlooked by a DDA-based library. Similarly,  
101 Wu et al. [29] have derived species composition information of samples by analyzing DIA data,  
102 laying the foundation for subsequent database construction. This strategy, entirely independent of  
103 DDA data, reduces the demand on sample and instrumental resources. Nevertheless, the databases  
104 constructed via these two methodologies, encompassing the complete protein sequences of the  
105 selected species, remain vast, necessitating substantial computational resources.

106 Therefore, to leverage the benefits of DIA in metaproteomics, the searching space needs to be  
107 further reduced. In the previous iterative strategy [7,29,30], the high-abundant proteins were used  
108 for the first search to infer the species that existed in the sample then all the proteins belonging to  
109 those species were then used for the subsequent search. However, this database remains overly  
110 extensive when compared to the number of identified peptides. Since the abundance of species  
111 within the microbiome shows significant disparity [31], the species identified should not be  
112 considered equally. The same applies to proteins and peptides. Proteins with high abundance and  
113 peptides with high detectability [32] or shared among various species are more likely to be detected.  
114 Here we report on a DIA workflow for metaproteomics, called MetaDIA, that relies on an  
115 annotated peptide database. This database comprises peptides that are anticipated to be detected,  
116 leveraging information on species abundance and protein abundance to score each peptide. We  
117 conducted a proof-of-concept experiment on human gut microbiome data generated by diaPASEF  
118 mode [25]. The peptide identification number and quantitative results obtained through our peptide  
119 library are comparable to those from the DDA-based library. Moreover, the species and functional  
120 information obtained from both methods are highly consistent.

121

## 122 **Method**

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

124 The Unified Human Gastrointestinal Protein (UHGP v2.0.2) catalog, encompassing 4744  
125 assembled genomes from the human gut microbiome, served as the reference database for this  
126 study [33]. Within this catalog, each protein sequence is uniquely associated with a distinct genome  
127 and is accompanied by detailed taxonomic and functional annotations. The detectability of  
128 peptides derived from these protein sequences was predicted using DeepDetect [32], a deep  
129 learning algorithm specifically designed for this purpose. This process involved *in silico* digestion  
130 of the protein sequences and subsequent assignment of a detectability score to each resultant  
131 peptide. Consequently, the peptide sequence reference database was enhanced by annotating each  
132 peptide with three key pieces of information: the genome identifier, the protein identifier, and the  
133 detectability score of the peptide. Please note that the database is structured on an identifier-centric  
134 organization. This means that peptides with identical sequences may be present within the database;  
135 however, as long as they are not from same genome and protein, they are distinguished by unique  
136 identifiers.

137

### 138 **Generation of FuncTax score**

139 Firstly, the peptides identified by MetaPep [34] are mapped to the UHGP database to establish  
140 peptide-genome associations. Subsequently, a greedy algorithm is employed to identify the  
141 minimal set of genomes that encompasses all peptide sequences, effectively reducing the  
142 complexity of the dataset. Following this, the intensity of each peptide is summed to infer genome  
143 abundance. The relative genome abundance will be used as the taxonomic score. To address the  
144 assignment of shared peptides, a razor strategy is adopted, analogous to the MaxQuant approach  
145 for protein inference [35]. Specifically, when a peptide is found in multiple genomes, it is attributed  
146 to the genome with the greater number of associated peptides

147 For the functional score, we constructed a fixed table from the MetaPep project [34]. While  
148 building the database Metapep, the peptide identification was performed by the software MetaLab  
149 MAG [7], which provided quantifications of protein abundance. Those proteins are well annotated.  
150 Subsequently, the relative abundance of each Clusters of Orthologous Groups (COG) accession  
151 was computed. The mean of non-zero relative abundance of the COG accessions was determined  
152 across samples, establishing a metric referred to the functional score.

153 The FuncTax score was obtained by multiplying two scores. In the case of peptides with the same  
154 sequence, their FuncTax scores were combined to give higher priority to shared peptides; the  
155 highest detectability score among them was utilized to ensure the inclusion of all possible peptides.  
156

### 157 **Taxonomic and functional analysis**

158 The taxonomic analysis is similar to the generation of genomic abundance score. The identified  
159 peptides are mapped to a database to establish peptide-genome associations. The database contains  
160 only the top 50 genomes. In our workflow, the peptide database was filtered out from the top 50  
161 genomes. So, all the identified peptides were from the top genomes and thus could be used for the  
162 taxonomic analysis (the peptides added from MetaPep might not be used). In DDA-based method,  
163 only the peptides mapping to the top 50 genomes were used for taxonomic analysis. Similarly, the  
164 razor strategy is used to process peptides shared by multiple genomes. Finally, the intensity of each  
165 peptide is summed to infer genome abundance.

166 For functional analysis, a protein abundance was firstly generated using the same strategy as  
167 taxonomic analysis. The proteins in the UHGP database have been extensively annotated thus the  
168 protein abundance can be further interpreted into functional abundance.  
169

### 170 **Deepdetect software configuration**

171 Protein digestion was simulated using Trypsin with the following parameters: a maximum of two  
172 missed cleavages, and peptide lengths ranging from 7 to 50 amino acids. Default settings were  
173 applied for all other parameters. This step was executed on a system equipped with a single Intel®  
174 Xeon® Gold 5118 processor (2.30 GHz, 12 cores), resulting in a total elapsed runtime of 126.5  
175 hours. The peptide database can be generated using Deepdetect by user or downloaded via the link  
176 on GitHub.  
177

### 178 **DIA software configuration**

179 DIA-NN (version 1.8.1) was used to process all the DIA data in this study. Maximum mass  
180 accuracy tolerances were set to 10 ppm for both MS1 and MS2 spectra for data generated from  
181 tims-TOF. For data generated from Exploris 480, mass accuracy tolerances were set to 0 to let  
182 DIA-NN optimise them automatically. False discovery rate (FDR) threshold of 1% was applied to  
183 ensure the reliability of identification. The --relaxed-prot-inf option was used for library-free

184 searching (MetaPep and MetaDIA). The --no-norm option was used to disable the normalization  
185 for the quantification benchmark experiment. All other settings were left default. The precursor  
186 matrix containing the peptide information was used for taxonomic and functional analysis.

187 FragPipe (version 23.1) was used for quantification benchmark experiment and samples 1–10  
188 for validation. For samples 1–10, DIA\_SpecLib\_Quant\_diaPASEF workflow was used in which  
189 diaTracer was used for spectrum deconvolution. For the mixed samples, DIA\_SpecLib\_Quant  
190 workflow was used, the --no-norm option was used to disable the normalization. FDR threshold  
191 of 1% was applied to ensure the reliability of identification. All other settings were left default.

192

### 193 **Metaproteomic datasets**

194 The dataset used for optimizing workflow is sourced from a published study and downloaded from  
195 ProteomeXchange Consortium (<http://www.proteomexchange.org>) with Dataset identifier  
196 PXD051104 [25]. The dataset for evaluating accuracy is from in-house samples. *Blautia*  
197 *hydrogenotrophica* (DSM 101114; Leibniz Institute DSMZ-German collection of microorganisms  
198 and cell cultures) was cultured in LB broth. The human stool was collected from a healthy adult  
199 volunteer at the University of Ottawa, Ottawa, ON, CAN. The protocol was approved by Ottawa  
200 Health Science Network Research Ethics (Approval No. 20160585-01H). The protein extraction  
201 and digestion were performed as described previously [36]. Peptide concentrations were measured  
202 using Thermo Scientific Pierce Quantitative Colorimetric Peptide Assays according to the  
203 manufacturer's directions.

204 The in-house samples were then analysed using an UltiMate 3000 RSLCnano system (Thermo  
205 Fisher Scientific, USA) coupled to an Orbitrap Exploris 480 mass spectrometer (Thermo Fisher  
206 Scientific, USA). Peptides were loaded onto an analytical column (75  $\mu\text{m}$  inner diameter  $\times$  15 cm)  
207 packed with reverse phase beads (3  $\mu\text{m}$ /120  $\text{\AA}$  ReproSil-Pur C18 resin, Dr. Maisch HPLC GmbH).  
208 The liquid chromatography system was directly connected to the analytical column, and no trap  
209 column was used in this analysis. A 60 min gradient of 5 to 35% (v/v) from buffer A (0.1% (v/v)  
210 formic acid) to B (0.1% (v/v) formic acid with 80% (v/v) acetonitrile) at a flow rate of 300  $\mu\text{L}/\text{min}$   
211 was used. The MS1 scan was conducted with a mass resolution of 60,000 (at 200 m/z) and covered  
212 a mass range of 380 to 985 m/z. The normalized automatic gain control (AGC) was set to 100%,  
213 with a maximum injection time of 100 ms. The MS2 scans were performed at a mass resolution of  
214 15,000 (at 200 m/z), within a mass range of 380 to 980 m/z. Isolation windows of 10 m/z with 1

215 m/z overlap were used. The normalized AGC for MS2 was set to 200%, with a maximum injection  
216 time of 40 ms. The normalized collision energy was maintained at 38%.

217

## 218 **Result**

### 219 **MetaDIA Workflow overview: taxonomy- and function-guided construction of peptide** 220 **database for metaproteomics**

221 Here we propose a new workflow for DIA based metaproteomics called MetaDIA. MetaDIA is a  
222 multistep workflow that systematically reduces the search space for DIA searching. At its basis, it  
223 relies on a combination of taxonomic abundance, functional abundance as a proxy of protein levels,  
224 and peptide detectability ultimately enabling DIA searching without the need for DDA results.  
225 Briefly, in the first step, we created a new database of peptides (UHGP peptide database), obtained  
226 by *in silico* digestion and detectability prediction of the UHGP database into 773,472,327 peptides  
227 [32,33]. Then each peptide is annotated with a FuncTax score (**Figure 1**). Both the FuncTax and  
228 the detectability scores are used to reduce the peptide database.

229 The FuncTax scores for each peptide in the UHGP peptide database are calculated using  
230 information from the MetaPep database [34]. MetaPep is a core peptide database compiling  
231 peptides previously identified in the published human gut metaproteomics studies. The  
232 information from MetaPep was used to create a static table of COG relative functional abundances  
233 and a sample-specific table of taxonomic relative abundances (Method). We noted that despite  
234 significant differences in the species composition of gut bacteria among different individuals, their  
235 functions are remarkably similar [37]. Therefore, functional abundance hierarchy information  
236 could act to estimate the likelihood of a protein being observed. We analyzed the search results  
237 used to construct the MetaPep database which contained 2134 raw files and 415 individuals [34]  
238 (Method). The functional ranking among various samples exhibits a strong correlation (Figure S1A  
239 and B). We observed a stable pattern in the functional hierarchy of human gut bacteria: abundant  
240 functions consistently remain high, while scarce functions persistently stay low across all samples  
241 (Figure S1C; Table S1). The sample-specific table of taxonomic relative abundances was generated  
242 by searching the DIA data against MetaPep [34]. The identified peptides and their quantitation  
243 were used to create the table (Method). However, this typically results in approximately 1000  
244 genomes remaining, with many containing only a single peptide. The number substantially larger  
245 than that is found in a typical human gut microbiome which is around 200 [31]. So, we only choose

246 the most abundant species for subsequent analysis. The selection of species for consideration is  
247 further explored in the optimization section of the study. The FuncTax score for each peptide is  
248 calculated by multiplying the taxonomic score for its taxonomic annotation and the functional  
249 score of its functional annotation. For peptides with identical sequences, their FuncTax scores were  
250 summed thereby leading to higher rankings for shared peptides.

251 In the last step, sample-specific reduced peptide database is generated by filtering UHGP peptide  
252 database using the FuncTax score and the detectability score (Figure 1). The final search of the  
253 DIA data is done against the reduced peptide database. To validate the efficiency of our peptide  
254 ranking method, peptides were sorted by FuncTax score and partitioned into equal-sized subsets  
255 based on their percentile rank (*e.g.*, top 0–5%, 5%–10%, ..., 35%–40%). Each subset was subjected  
256 to database searching with uniform parameters. We observed a decline in the number of peptides  
257 identified as the percentile ranking of the subsets decreased (**Figure 2**). The decreasing trend  
258 suggested our ranking method effectively prioritizes peptides with a higher probability of detection.

259

### 260 **Optimized MetaDIA parameters reduce the database size**

261 We explored whether the number of microbes in the reduced peptide database, the threshold for  
262 FuncTax score and the threshold for detectability score influenced the identification of peptides.  
263 To examine the impact of these parameters, we utilized DIA data from 10 previously reported  
264 human gut microbiome samples [25] (File S1, sample information).

265 In particular, we first tested the effect of the number of microbes (genomes) ranging from 50 to  
266 150 and FuncTax score ranging from top 1% to 40% (**Figure 3**). We keep the detectability  
267 threshold at the top 40% in this experiment which is suggested by the author of Deepdetect [32].  
268 Interestingly, no matter how many genomes we chose, the size of the reduced peptide database had  
269 the strongest effect on the number of identified peptides. The identification number plateaued once  
270 the reduced peptide database size reached around 1.6 million entries (Figure 3A and B, Figures S2  
271 and S3), corresponding to a FuncTax score threshold of 40% for 50 genomes, 20% for 100  
272 genomes, and 15% for 150 genomes, respectively. We compared the three different reduced peptide  
273 databases, which led to consistent peptide identification results (Figure 3C and D, Figures S4 and  
274 S5). In our previous studies, we observed that low-abundance species were underrepresented [8].  
275 In this context, we chose to focus on the top 50 genomes to prioritize high-abundance genomes. It

276 is important to note that this cut-off is a variable parameter that can be adjusted according to the  
277 specific objectives of different studies.

278 Subsequently, we explored whether the detectability threshold impacted the number of peptides  
279 identified. While the recommended threshold by the author of Deepdetect was 40%, we explored  
280 thresholds ranging from 40% to 10%. We observed that a threshold of 25% was the point at which  
281 the number of identifications began to decrease significantly (**Figure 4A**). However, both the  
282 database size and the search time decreased substantially (Figure 4B). Comparing the identification  
283 results at thresholds of 25% and 40%, we found a substantial overlap (Figure 4C). Therefore, we  
284 selected a 25% threshold for detectability. Based on this analysis, we proceeded with peptides  
285 ranking in the top 40% by FuncTax score and the top 25% by detectability. Given that these two  
286 scores are entirely uncorrelated, applying both filters effectively reduces the database to one-tenth  
287 of its original size (25% times 40%, Figure S6). After applying these optimized parameters,  
288 approximately 1 million peptide sequences remain in the reduced database. Compared to  
289 previously reported database reduction strategies, our workflow substantially reduces the size of  
290 the database. Wu et al. [29] employed a high-abundance protein-guided method to construct their  
291 database, incorporating 256 or 366 complete proteomes. Rough estimates indicate that the database  
292 generated by our approach is 50 or 70 times smaller. A smaller database reduces computational  
293 resource consumption (Figure S7). Additionally, smaller searching space potentially facilitates  
294 peptide identification. Although different searching engines were used, almost double peptides  
295 were identified when we apply MetaDIA to the same data (File S2; Table S2).

296

### 297 **MetaDIA maintains accuracy in DIA peptide identification**

298 We next explored whether the enrichment of high abundant and highly detectable peptides in our  
299 reduced database impacted the accuracy of peptide identification when applying the false  
300 discovery rate strategy. To evaluate this, we conducted a benchmark experiment using three  
301 samples: a human gut microbiome sample (sample A), a *Blautia hydrogenotrophica* sample  
302 (sample C), and a 50:50 mixed sample of the two (sample B) (**Figure 5A**). *Blautia*  
303 *hydrogenotrophica* was selected due to its absence in the microbiome sample used here and its  
304 minimal peptide overlap with the microbiome sample. Each sample was subjected to triplicate DIA  
305 measurements. Samples A and B were analyzed using the reduced peptide database generated by  
306 our workflow with optimized parameters of top 40% FuncTax score and top 25% detectability

307 score, whereas samples B and C were searched against species-specific protein databases derived  
308 from NCBI (Genome assembly ASM15797v1). In the first search against the peptide database,  
309 36,971 unique peptides were identified. Of these, 2072 peptides also present in the *Blautia*  
310 *hydrogenotrophica* database were excluded. Further, peptides unique to each sample were removed,  
311 leaving 32,210 peptides identified in both samples A and B. Ideally, the peptide abundance ratio  
312 between samples A and B should be approximately 2. In the second search, 16,173 unique peptides  
313 were identified. Among these, 12,109 peptides were unique to the *Blautia hydrogenotrophica*  
314 database and were found in both samples B and C. The expected ratio between samples B and C  
315 should be around 0.5. We found that, whether we used a protein database or a peptide database,  
316 the ratios of peptides identified in both searches closely aligned with the expected values (Figure  
317 5B). To ensure that the performance is not specific to DIA-NN, we independently re-analyzed the  
318 three samples using FragPipe [38,39], a widely adopted alternative to DIA-NN with distinct  
319 algorithmic foundations. The results were consistent with those achieved using DIA-NN (Figure  
320 S8). This suggests that the employment of our reduced peptide database does not significantly  
321 affect the accuracy of peptide identification, thereby supporting its use in peptide identification  
322 workflows with a controlled FDR.

323

### 324 **MetaDIA yields consistent peptide and protein identification results with DDA-based** 325 **strategies.**

326 We next evaluated whether MetaDIA performed similarly to a conventional DDA-based library  
327 for DIA data analysis. The DIA data and corresponding DDA-based library were obtained from a  
328 published study [25]. We found that the MetaDIA provided identification numbers comparable to  
329 those obtained through the DDA-based library (**Figure 6A**). Notably, in certain instances, such as  
330 with samples 8 and 9, the MetaDIA surpassed DDA library in the number of identifications. The  
331 initial step in our workflow involves searching the raw data against MetaPep. MetaPep was  
332 constructed by collecting searching result from an open search algorithm, Metalab-MAG [7]. Thus,  
333 it encompasses modified peptides not included in the original database. Subsequent integration of  
334 peptides identified by MetaPep into a refined peptide database resulted in a marked increase in  
335 identification rates (Figure 6A, Figure S9; File S3). Because of the augmentation in the number of  
336 peptide identifications attributable to MetaPep, we incorporate the identification results from  
337 MetaPep into the database in the following study.

338 Over 50% of peptides identified from the DDA-based library were also identified by  
339 MetaDIA+MetaPep (Figure 6B, Figure S10). The divergence in unique identifications between the  
340 two methods may be attributed to a variety of factors, including the size and composition of the  
341 data, as well as the inherent differences between DDA acquisition and DIA acquisition. For  
342 instance, we have found that approximately half of the unique peptides from the DDA-based  
343 library are indeed present within the database produced by MetaDIA+MetaPep; however, they  
344 remain unidentified in our workflow (Figure S11). Moreover, peptides provided by MetaDIA  
345 originate exclusively from a pre-selected group of 50 species, whereas the results from the DDA  
346 library encompass peptides from thousands of species. This, to some extent, elucidates the small  
347 overlap observed between the two methods. Upon examining the quantification results of those  
348 peptides found by both methods, we observed a significant consistency in the outcomes, with a  
349 Pearson coefficient above 0.9 (Figure 6D, Figure S12). It is worth noting that the fragment ions  
350 used for quantification in the DDA-based library correspond to actual DDA acquisitions. In  
351 contrast, our workflow leverages the library-free mode within DIA-NN which uses theoretical  
352 spectra that are predicted from peptide sequences. The high degree of agreement between the  
353 quantification results underscores the reliability of the MS-Simulator algorithm which is employed  
354 by DIA-NN for spectra prediction [14]. We further validated the MetaDIA using the FragPipe. In  
355 general, FragPipe identified less peptides (Figure S13A). This may potentially be attributed to its  
356 dependence on MS1 signals for assembling associated fragment ion intensities into pseudo MS2  
357 spectra, whereas DIA-NN does not require such an approach. By using FragPipe, MetaDIA  
358 demonstrated a distinct advantage over DDA-based method in terms of the number of  
359 identifications. Similar number of common peptides (those detected by both the DDA-based  
360 library and MetaDIA) were identified by FragPipe and DIA-NN (Figure S13B). The DDA-based  
361 library and MetaDIA demonstrated greater consistency when using the FragPipe on peptide  
362 quantification (Figure S13C and D).

363 At the protein level, our findings revealed greater consistency in identification compared to the  
364 peptide level (Figure 6C). Around 70% of proteins found by the DDA-based library can be found  
365 by MetaDIA+MetaPep. The overlap on protein level reinforces the reliability of the identifications  
366 and indicates that a significant subset of proteins is consistently identified by both methods despite  
367 differences at the peptide level (Figure S14). Proteins like GYG000002545\_00035 had greater  
368 sequence coverage and higher detection intensity with the DDA library, while others like

369 MGYG000002272\_00452 showed higher coverage and intensity with MetaDIA+MetaPep. Given  
370 that the quantification of a protein was derived from different subsets of peptides in these two  
371 methods, we observed reduced consistency of quantification in the protein level between the  
372 methods, as reflected by Pearson correlation coefficients of approximately 0.7 (Figure 6E, Figure  
373 S15). However, it is important to note that in most proteomic studies, the primary interest lies in  
374 the differential abundance of the same protein across various samples. Therefore, it is crucial that  
375 we use the same fragment ions to quantify a protein. In this regard, the inconsistencies in protein  
376 quantification between the two methods do not undermine the utility of either approach. The  
377 substantial overlap in peptide and protein identification by both methods suggests a robust cross-  
378 validation of both methods. Then we annotated the proteins using COG accessions and calculated  
379 their relative abundances. Our analysis revealed that approximately 90% of the COG accessions  
380 identified by the DDA-based method were also covered by our method (Figure 6C). Furthermore,  
381 the Pearson correlation coefficient for the relative abundance of COG accessions exceeded 0.9,  
382 with a stronger correlation for those COG accessions that were highly abundant (Figure 6F, Figure  
383 S16).

384

### 385 **MetaDIA provides taxonomic profiles highly similar to those obtained from searching DDA-** 386 **libraries.**

387 We verified whether both methods had a high degree of similarity in the taxonomic composition.  
388 We did comparative analysis of microbiome composition across different taxonomic levels using  
389 the results from both methods. In our workflow, taxonomic analysis was confined to the predefined  
390 database containing 50 genomes. In the DDA-based method, the peptide identified by the DDA  
391 library can be annotated to over 1000 genomes even after using the greedy algorithm (Method:  
392 generation of FuncTax score). However, we found that the top 50 genomes accounted for 79%–  
393 90% of peptides and 87%–92% of peptide intensity (Figure S17). To simplify the comparison  
394 between the two methods, we discarded the small number of peptides that could not be annotated  
395 to the top 50 genomes. Our findings indicate that there is a significant linear correlation between  
396 the compositions identified by both methods, with the degree of correlation strengthening at higher  
397 taxonomic levels (Figure 7A and B, Figure S18). The two methods showed remarkably consistent  
398 taxonomic composition at the genus level with a Pearson coefficient above 0.98 across all the  
399 samples tested. Even sample 9, which displayed the lowest correlation, demonstrated a substantial

400 degree of consistency between the two methods. To underscore the consistency, we have provided  
401 a detailed visualization of the taxonomic composition for sample 9 (Figure 7C and D, Figure S19)  
402 The species compositions observed by our method in these ten samples differ significantly as  
403 expected, indicating that our database and taxonomic analysis have the capability to identify a  
404 diverse range of microbiota (Figure S20; File S2). The most abundant species identified in the ten  
405 samples have been previously reported as high-abundance species in the human gut microbiome  
406 [40–44]. Except for *Phocaeicola dorei* which was identified as the top species in samples 2, 5, and  
407 10, the other top species were all unique to each sample.

408

### 409 **Performance validation of MetaDIA on large-scale diaPASEF datasets.**

410 To further validate the versatility and applicability of our proposed metaproteomic workflow, we  
411 extended our analysis to a diverse set of 79 DIA datasets obtained from a published study [25].  
412 This dataset encompasses samples from 62 individuals, featuring replicate injections, quality  
413 control samples, and pooled samples (File S1). We applied MetaDIA to this extensive dataset and  
414 compared the results with the conventional DDA-based approach. Remarkably, the number of  
415 peptides identified by both methods demonstrated a close equivalence, reinforcing the robustness  
416 and universal applicability of our metaproteomic workflow (**Figure 8**). Validating our method  
417 across diverse samples enhances confidence in its effectiveness and consistency, demonstrating its  
418 potential for widespread adoption in metaproteomics research.

419

## 420 **Discussion**

421 We propose a novel workflow for DIA data analysis from human gut microbiome called MetaDIA.  
422 The approach aims at prioritizing peptides with a higher likelihood of detection based on their  
423 detectability, as well as taxonomic and functional scores.

424 MetaDIA is entirely devoid of DDA, thereby circumventing the drawbacks of DDA-based  
425 methods. Not only does this approach save time and resources, but it also enables the creation of  
426 a tailored database for each sample. In contrast, DDA-based methods typically rely on a single  
427 pooled sample to generate a library. For instance, Gomez et al. [24] used a pooled sample to  
428 represent 12 individual mice, while Sun et al. [25] did so for a cohort of 62 individuals. However,  
429 such a pooled sample may not effectively represent every sample. In our study, the ten samples  
430 showed highly diverse taxonomic composition (Figure S20). To increase the sampling depth for

431 the pooled sample, Sun et al. [25] had to fractionate the pooled sample into 30 portions and Gomez  
432 et al. [24] repeatedly injected the pooled sample 10 times. Moreover, utilizing a static library to  
433 search various samples may potentially compromise the accuracy of peptide identification, as it  
434 includes peptides from the pooled samples that are absent in the specific sample under  
435 investigation.

436 In MetaDIA, we pre-defined the range of genomes for each microbiome sample (50 genomes in  
437 this study). This approach not only enabled us to narrow the search space but also to mitigate the  
438 issues associated with protein inference that arise from common peptides. When assigning peptides  
439 to proteins, we confined our consideration to the genomes within the predefined range rather than  
440 the entire dataset. This strategy significantly reduced the incidence of common peptides.

441 In the process of calculating the taxonomic composition of a sample, we employed a greedy  
442 algorithm to minimize the number of genomes required to account for all identified peptides.  
443 Practically, this was achieved through an iterative method, sequentially selecting the genome that  
444 encompasses the greatest number of peptides, until every peptide could be explained by the  
445 selected genomes. By using this, the matches between peptides and genomes are significantly  
446 simplified. However, this approach is prone to yielding locally optimal solutions that may not align  
447 with the global optimum, primarily due to the shared nature of peptides among different genomes.  
448 For example, a genome containing greatest number of identified peptides may not present in the  
449 sample, potentially compromising the precision of taxonomic analysis at the genomic level.  
450 Fortunately, genomes possessing a multitude of shared peptides typically align at higher taxonomic  
451 levels (*e.g.*, Genus, Family). Hence, the precision of taxonomic analysis is anticipated to be  
452 improved on higher taxonomic levels.

453 MetaDIA relies on MetaPep for the first search to acquire the taxonomic composition of the  
454 sample. Although MetaPep showed good performance in this study, it may exhibit bias due to the  
455 underrepresentation of peptides from rare species or functions, instrumentation-specific detection  
456 limits, and variations in sample processing methods, which compromise its comprehensiveness  
457 and reliability. Because MetaPep is for human gut microbiome, the current workflow can only be  
458 applied to the human gut microbiome. However, we foresee that its strategy could be extended to  
459 other types of microbiomes. In the bottom-up approach to metaproteomics, the abundance of a  
460 specific peptide is ascertained by its corresponding protein and species abundance, both of which  
461 are generally diminished during the processes of protein digestion and cellular lysis. Within this

462 strategy, efforts are made to digitally reconstruct these two crucial pieces of information using  
463 functional score and taxonomic score and employ them to prioritize the peptides with the highest  
464 probability of detection. In the future, a selected database containing proteins common to all  
465 taxonomy could be used to help us infer the taxonomic composition for other types of microbiome  
466 [28,29], and the functional score could be acquired by analyzing previously published data.  
467 Importantly, owing to the comprehensive nature of DIA data, researchers may also request the  
468 inclusion of specific functions or pathways of interest into the database.

469

## 470 **Conclusion**

471 In conclusion, we introduced a new strategy to prioritize peptides with a high probability of  
472 detection. This strategy simulates protein digestion procedures *in silico* and uses taxonomic and  
473 functional information to infer the peptide abundance. MetaDIA is a fully DDA-free workflow and  
474 provides a user interface to change the different parameters. We compared the performance of  
475 MetaDIA with the DDA-based library and observed a high degree of consistency. We further  
476 validated our method across a DIA-PASEF dataset with 79 samples, thereby confirming its wide  
477 applicability. We believe that our approach will help the application of DIA in metaproteomics.

478

## 479 **Code availability**

480 MetaDIA is publicly available at GitHub (<https://github.com/northomics/MetaDIA>). The code has  
481 also been submitted to BioCode at the National Genomics Data Center (NGDC), China National  
482 Center for Bioinformation (CNCB) (BioCode: BT007936), which is publicly accessible at  
483 <https://ngdc.cncb.ac.cn/biocode/tool/BT007936>.

484

## 485 **Data availability**

486 The datasets generated in this study were sourced from ProteomeXchange Consortium  
487 (<http://www.proteomexchange.org>) with dataset identifier PXD063632.

488

## 489 **CRedit author statement**

490 **Haonan Duan:** Conceptualization, Methodology, Software, Formal analysis, Writing – original  
491 draft p reparation. **Zhibin Ning:** Conceptualization, Methodology, Writing – review & editing.  
492 **Zhongzhi Sun:** Writing – review & editing. **Tiannan Guo:** Data Curation, Writing – review &

493 editing. **Yingying Sun**: Data Curation, Writing – review & editing, **Daniel Figeys**: Supervision,  
494 Writing – review & editing, Funding acquisition. All authors read and approved the final  
495 manuscript.

496

### 497 **Competing interests**

498 Dainel Figeys is the founder of MedBiome Inc. a microbiome nutrition and therapeutic company.  
499 The other authors have declared that they have no competing interests.

500

### 501 **Acknowledgments**

502 This work was funded by the Natural Sciences and Engineering Research Council of Canada  
503 (NSERC, Grant No. 2018-03905) discovery grant to Daniel Figeys. Haonan Duan was funded by  
504 a stipend from the NSERC CREATE in Technologies for Microbiome Science and Engineering  
505 (TECHNOMISE) Program (Grant No. 497995) and Jiangsu Funding Program for Excellent  
506 Postdoctoral Talent.

507

### 508 **Supplementary material**

509 Supplementary material is available at *Genomics, Proteomics & Bioinformatics* online  
510 (<https://doi.org/xxxxxxx>).

511

### 512 **ORCID**

513 0000-0002-2753-8594 (Haonan Duan)

514 0000-0003-2045-7596 (Zhibin Ning)

515 0000-0002-4862-8599 (Zhongzhi Sun)

516 0009-0008-6015-2109 (Yingying Sun)

517 0000-0003-3869-7651 (Tiannan Guo)

518 0000-0002-5373-7546 (Daniel Figeys)

519

### 520 **Reference**

521 [1] Human Microbiome Project C. Structure, function and diversity of the healthy human  
522 microbiome. *Nature* 2012;486:207–14.

- 523 [2] Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, et al. A communal  
524 catalogue reveals Earth's multiscale microbial diversity. *Nature* 2017;551:457–63.
- 525 [3] Wilmes P, Bond PL. Metaproteomics: studying functional gene expression in microbial  
526 ecosystems. *Trends Microbiol* 2006;14:92–7.
- 527 [4] Jagtap P, Goslinga J, Kooren JA, McGowan T, Wroblewski MS, Seymour SL, et al. A two-step  
528 database search method improves sensitivity in peptide sequence matches for metaproteomics and  
529 proteogenomics studies. *Proteomics* 2013;13:1352–7.
- 530 [5] Zhang X, Ning Z, Mayne J, Moore JI, Li J, Butcher J, et al. MetaPro-IQ: a universal  
531 metaproteomic approach to studying human and mouse gut microbiota. *Microbiome* 2016;4:31.
- 532 [6] Cheng K, Ning Z, Zhang X, Li L, Liao B, Mayne J, et al. MetaLab 2.0 enables accurate post-  
533 translational modifications profiling in metaproteomics. *J Am Soc Mass Spectrom* 2020;31:1473–  
534 82.
- 535 [7] Cheng K, Ning Z, Li L, Zhang X, Serrana JM, Mayne J, et al. MetaLab-MAG: a metaproteomic  
536 data analysis platform for genome-level characterization of microbiomes from the metagenome-  
537 assembled genomes database. *J Proteome Res* 2023;22:387–98.
- 538 [8] Duan H, Cheng K, Ning Z, Li L, Mayne J, Sun Z, et al. Assessing the dark field of  
539 metaproteome. *Anal Chem* 2022;94:15648–54.
- 540 [9] Gillet LC, Navarro P, Tate S, Rost H, Selevsek N, Reiter L, et al. Targeted data extraction of  
541 the MS/MS spectra generated by data-independent acquisition: a new concept for consistent and  
542 accurate proteome analysis. *Mol Cell Proteomics* 2012;11:O111.016717.
- 543 [10] Collins BC, Hunter CL, Liu Y, Schilling B, Rosenberger G, Bader SL, et al. Multi-laboratory  
544 assessment of reproducibility, qualitative and quantitative performance of SWATH-mass  
545 spectrometry. *Nat Commun* 2017;8:291.
- 546 [11] Vowinckel J, Zelezniak A, Bruderer R, Mülleler M, Reiter L, Ralser M. Cost-effective  
547 generation of precise label-free quantitative proteomes in high-throughput by microLC and data-  
548 independent acquisition. *Sci Rep* 2018;8:4346.
- 549 [12] Meier F, Brunner AD, Frank M, Ha A, Bludau I, Voytik E, et al. diaPASEF: parallel  
550 accumulation-serial fragmentation combined with data-independent acquisition. *Nat Methods*  
551 2020;17:1229–36.

- 552 [13] Guzman UH, Martinez-Val A, Ye Z, Damoc E, Arrey TN, Pashkova A, et al. Ultra-fast label-  
553 free quantification and comprehensive proteome coverage with narrow-window data-independent  
554 acquisition. *Nat Biotechnol* 2024;42:1855–66.
- 555 [14] Sun S, Yang F, Yang Q, Zhang H, Wang Y, Bu D, et al. MS-Simulator: predicting y-ion  
556 intensities for peptides with two charges based on the intensity ratio of neighboring ions. *J*  
557 *Proteome Res* 2012;11:4509–16.
- 558 [15] Zeng WF, Zhou XX, Zhou WJ, Chi H, Zhan J, He SM. MS/MS spectrum prediction for  
559 modified peptides using pDeep2 trained by transfer learning. *Anal Chem* 2019;91:9724–31.
- 560 [16] Bouwmeester R, Gabriels R, Hulstaert N, Martens L, Degroeve S. DeepLC can predict  
561 retention times for peptides that carry as-yet unseen modifications. *Nat Methods* 2021;18:1363–9.
- 562 [17] Ma C, Ren Y, Yang J, Ren Z, Yang H, Liu S. Improved peptide retention time prediction in  
563 liquid chromatography through deep learning. *Anal Chem* 2018;90:10881–8.
- 564 [18] Al Musaimi O, Valenzo OMM, Williams DR. Prediction of peptides retention behavior in  
565 reversed-phase liquid chromatography based on their hydrophobicity. *J Sep Sci* 2023;46:e2200743.
- 566 [19] Cox J. Prediction of peptide mass spectral libraries with machine learning. *Nat Biotechnol*  
567 2023;41:33–43.
- 568 [20] Demichev V, Messner CB, Vernardis SI, Lilley KS, Ralser M. DIA-NN: neural networks and  
569 interference correction enable deep proteome coverage in high throughput. *Nat Methods*  
570 2020;17:41–4.
- 571 [21] Demichev V, Szyrwił L, Yu F, Teo GC, Rosenberger G, Niewianda A, et al. dia-PASEF data  
572 analysis using FragPipe and DIA-NN for deep proteomics of low sample amounts. *Nat Commun*  
573 2022;13:3944.
- 574 [22] Sinitcyn P, Hamzeiy H, Salinas Soto F, Itzhak D, McCarthy F, Wichmann C, et al. MaxDIA  
575 enables library-based and library-free data-independent acquisition proteomics. *Nat Biotechnol*  
576 2021;39:1563–73.
- 577 [23] Aakko J, Pietila S, Suomi T, Mahmoudian M, Toivonen R, Kouvonen P, et al. Data-  
578 independent acquisition mass spectrometry in metaproteomics of gut microbiota-implementation  
579 and computational analysis. *J Proteome Res* 2020;19:432–6.

- 580 [24] Gomez-Varela D, Xian F, Grundtner S, Sondermann JR, Carta G, Schmidt M. Increasing  
581 taxonomic and functional characterization of host-microbiome interactions by DIA-PASEF  
582 metaproteomics. *Front Microbiol* 2023;14:1258703.
- 583 [25] Sun Y, Xing Z, Liang S, Miao Z, Zhuo LB, Jiang W, et al. metaExpertPro: a computational  
584 workflow for metaproteomics spectral library construction and data-independent acquisition mass  
585 spectrometry data analysis. *Mol Cell Proteomics* 2024;23:100840.
- 586 [26] Pietilä S, Suomi T, Elo LL. Introducing untargeted data-independent acquisition for  
587 metaproteomics of complex microbial samples. *ISME Commun* 2022;2:51.
- 588 [27] Tsou CC, Avtonomov D, Larsen B, Tucholska M, Choi H, Gingras AC, et al. DIA-Umpire:  
589 comprehensive computational framework for data-independent acquisition proteomics. *Nat*  
590 *Methods* 2015;12:258–64.
- 591 [28] Dumas T, Martinez Pinna R, Lozano C, Radau S, Pible O, Grenga L, et al. The astounding  
592 exhaustiveness and speed of the Astral mass analyzer for highly complex samples is a quantum  
593 leap in the functional analysis of microbiomes. *Microbiome* 2024;12:46.
- 594 [29] Wu E, Yang Y, Zhao J, Zheng J, Wang X, Shen C, et al. High-abundance protein-guided hybrid  
595 spectral library for data-independent acquisition metaproteomics. *Anal Chem* 2024;96:1029–37.
- 596 [30] Stambouliau M, Li S, Ye Y. Using high-abundance proteins as guides for fast and effective  
597 peptide/protein identification from human gut metaproteomic data. *Microbiome* 2021;9:80.
- 598 [31] Yang J, Pu J, Lu S, Bai X, Wu Y, Jin D, et al. Species-level analysis of human gut microbiota  
599 with metataxonomics. *Front Microbiol* 2020;11:2029.
- 600 [32] Yang J, Cheng Z, Gong F, Fu Y. DeepDetect: deep learning of peptide detectability enhanced  
601 by peptide digestibility and its application to DIA library reduction. *Anal Chem* 2023;95:6235–43.
- 602 [33] Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, et al. A unified catalog  
603 of 204,938 reference genomes from the human gut microbiome. *Nat Biotechnol* 2021;39:105–14.
- 604 [34] Sun Z, Ning Z, Cheng K, Duan H, Wu Q, Mayne J, et al. MetaPep: A core peptide database  
605 for faster human gut metaproteomics database searches. *Comput Struct Biotechnol J*  
606 2023;21:4228–37.
- 607 [35] Tyanova S, Temu T, Cox J. The MaxQuant computational platform for mass spectrometry-  
608 based shotgun proteomics. *Nat Protoc* 2016;11:2301–19.

- 609 [36] Zhang X, Li L, Mayne J, Ning Z, Stintzi A, Figeys D. Assessing the impact of protein  
610 extraction methods for human gut metaproteomics. *J Proteomics* 2018;180:120–7.
- 611 [37] Li L, Wang T, Ning Z, Zhang X, Butcher J, Serrana JM, et al. Revealing proteome-level  
612 functional redundancy in the human gut microbiome using ultra-deep metaproteomics. *Nat*  
613 *Commun* 2023;14:3428.
- 614 [38] Yu F, Teo GC, Kong AT, Frohlich K, Li GX, Demichev V, et al. Analysis of DIA proteomics  
615 data using MSFragger-DIA and FragPipe computational platform. *Nat Commun* 2023;14:4154.
- 616 [39] Li K, Teo GC, Yang KL, Yu F, Nesvizhskii AI. diaTracer enables spectrum-centric analysis of  
617 diaPASEF proteomics data. *Nat Commun* 2025;16:95.
- 618 [40] Davis-Richardson AG, Ardisson AN, Dias R, Simell V, Leonard MT, Kemppainen KM, et  
619 al. *Bacteroides dorei* dominates gut microbiome prior to autoimmunity in Finnish children at high  
620 risk for type 1 diabetes. *Front Microbiol* 2014;5:678.
- 621 [41] Hosomi K, Saito M, Park J, Murakami H, Shibata N, Ando M, et al. Oral administration of  
622 *Blautia wexlerae* ameliorates obesity and type 2 diabetes via metabolic remodeling of the gut  
623 microbiota. *Nat Commun* 2022;13:4477.
- 624 [42] Scher JU, Sczesnak A, Longman RS, Segata N, Ubeda C, Bielski C, et al. Expansion of  
625 intestinal *Prevotella copri* correlates with enhanced susceptibility to arthritis. *Elife* 2013;2:e01202.
- 626 [43] Lan PTN, Sakamoto M, Sakata S, Benno Y. *Bacteroides barnesiae* sp. nov., *Bacteroides*  
627 *salanitronis* sp. nov. and *Bacteroides gallinarum* sp. nov., isolated from chicken caecum. *Int J Syst*  
628 *Evol Microbiol* 2006;56:2853–9.
- 629 [44] Ferreira-Halder CV, Faria AVS, Andrade SS. Action and function of *Faecalibacterium*  
630 *prausnitzii* in health and disease. *Best Pract Res Clin Gastroenterol* 2017;31:643–8.

631

## 632 **Figure legend**

### 633 **Figure 1 The flowchart for the MetaDIA**

634 All proteins in Unified Human Gastrointestinal Protein (UHGP) database were firstly *in silico*  
635 digested into peptides. The detectability of each peptide was predicated by DeepDetect algorithm.  
636 Following this prediction, each peptide was assigned a functional score and a taxonomic score,  
637 derived from a predetermined functional relative abundance table and a sample-specific taxonomic  
638 relative abundance table, respectively (Method). The FuncTax score was calculated by multiplying

639 the two scores. For peptides with identical sequence, their FuncTax scores were summed to  
640 prioritize shared peptides; the maximum of their detectability scores was used to ensure the  
641 inclusion of all potential peptides. The detectability and FuncTax scores are both used for filtering  
642 peptides. The reduced peptide database was used for a second search.

643

## 644 **Figure 2 The number of peptides identified from each subset**

645 Ten samples were tested in the experiment. For constructing the peptide database, the top 100  
646 genomes were considered; the detectability threshold was set at 40%. Each subset contains around  
647 400,000 peptides. Peptide identification was performed by DIA-NN under same conditions. The  
648 maximum identification in the ten samples from the last subset was highlighted in the figure.

649

## 650 **Figure 3 Optimization for genome number and FuncTax score (Sample 1)**

651 Peptides from  $n$  (50, 100, and 150) genomes were ranked by the FuncTax score and top  $x\%$  (1–40  
652 for 50 and 100 genomes; 1–35 for 150 genomes) peptides was used as database. **A.** Number of  
653 identified peptides against database percentage. **B.** Number of identified peptides against database  
654 size. The inflection point has been highlighted with a red box. The overlap of the reduced peptide  
655 database (**C**) and identified peptide (**D**) when taking top 40% peptides for 50 genomes, top 20%  
656 for 100 genomes, and top 15% for 150 genomes as database. Peptide identification was performed  
657 by DIA-NN under same conditions.

658

## 659 **Figure 4 Optimization for detectability threshold**

660 The number of peptides identified (**A**) and the searching time (**B**) under detectability threshold  
661 from 10% to 40%. Only the searching time is included, time consumed by DeepDetect is not  
662 accounted for. **C.** The overlap of peptides identified by top 25% and top 40% of the database.  
663 Peptide identification was performed by DIA-NN under same conditions.

664

## 665 **Figure 5 Benchmark experiment for peptide identification**

666 **A.** The experimental design. Each sample was subjected to triple-run measurements. **B.** Log-  
667 transformed ratios are plotted as a function of peptide intensity for microbiome peptides ( $n =$   
668 32,210, in green) and *Blautia hydrogenotrophica* peptides ( $n = 12,109$ , in purple). The point  
669 density for ratio was plotted at right. Dashed lines indicate the expected ratio. Peptide identification

670 was performed by DIA-NN under same conditions. The intensities for x axis were derived from  
671 pure samples (Samples A and C).

672

### 673 **Figure 6 Comparison between the DDA-based method and DDA-free method**

674 **A.** The peptide identified by each method. **B.** The overlap of peptide identified in sample 1 by each  
675 method. **C.** Coverage of peptides, proteins and cog accessions identified by DDA-based method  
676 with those found using DDA-free method. The intensity correlation of the overlapped peptides (**D**),  
677 proteins (**E**), and COG accessions (**F**) in sample 1. The dashed line indicates  $y = x$ . For DDA-  
678 based method, the peptides identified as derived from human proteins are removed.

679

### 680 **Figure 7 Comparison of the taxonomic composition between the DDA-based method and** 681 **DDA-free method**

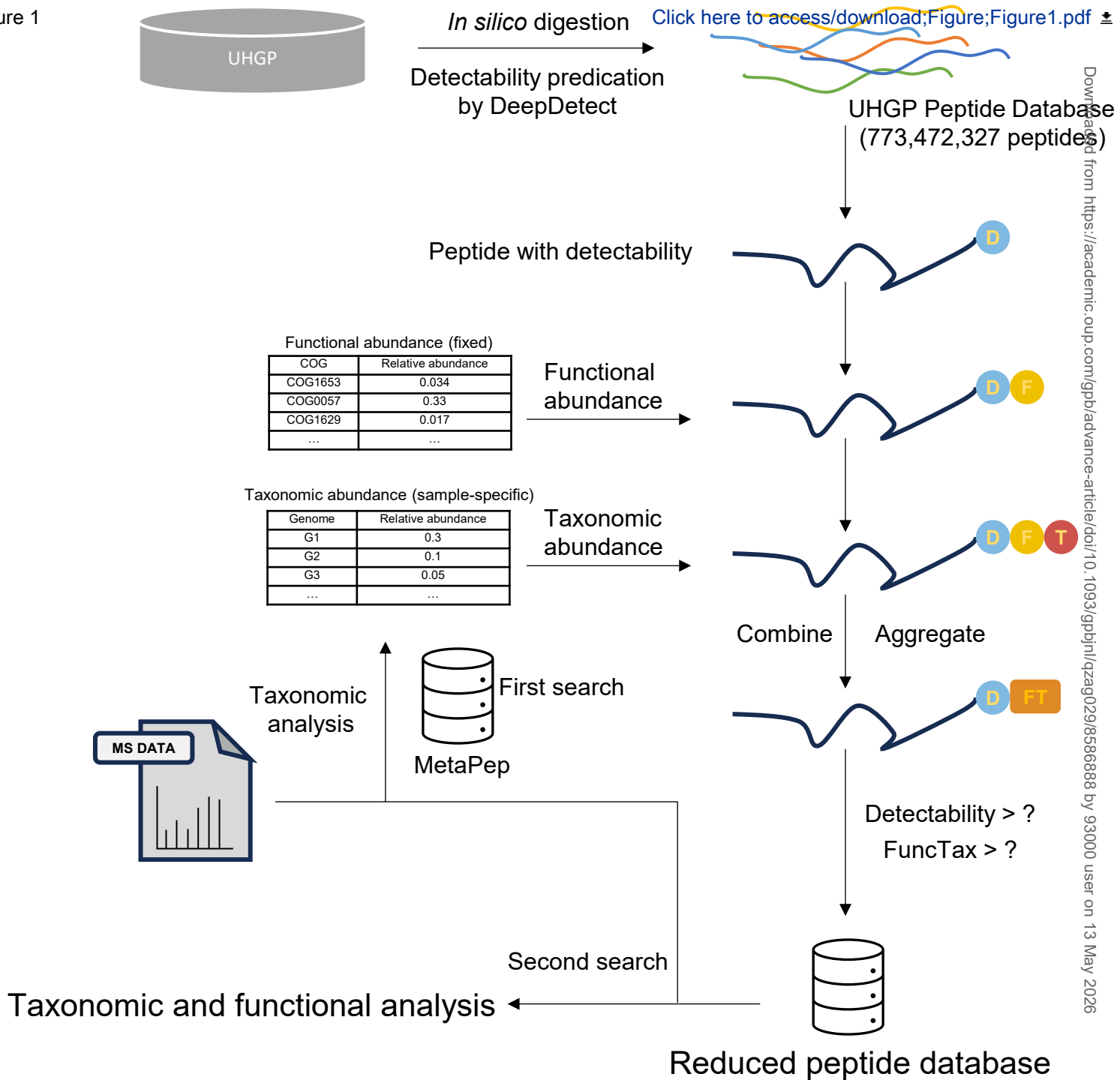
682 Pearson correlation (**A**) and Bray-Curtis distance (**B**) analysis between DDA-based method and  
683 DDA-free method on different taxonomic levels from Phylum to Species. The sunburst chart  
684 showing taxonomic composition (Phylum to Family) of sample 9 derived from DDA-based  
685 method (**C**) and DDA-free method (**D**). The relative taxonomic abundance was used for the  
686 analysis. In the correlation analysis, taxonomic categories that were unique to one method were  
687 imputed with a value of zero. In sunburst chart, some species names are omitted due to spatial  
688 constraints. Each color corresponds to a distinct species, and consistent coloring is applied for the  
689 same species across panels.

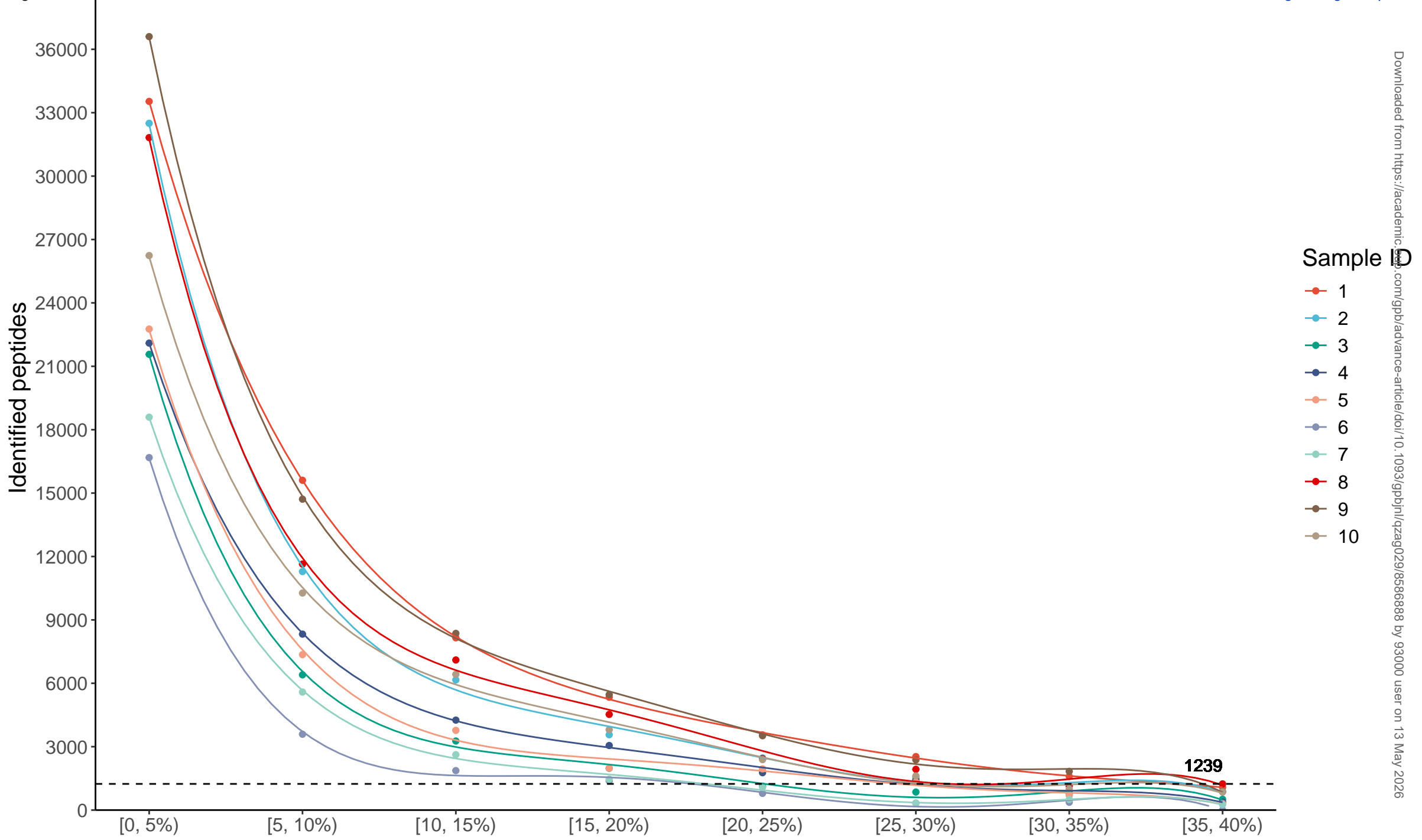
690

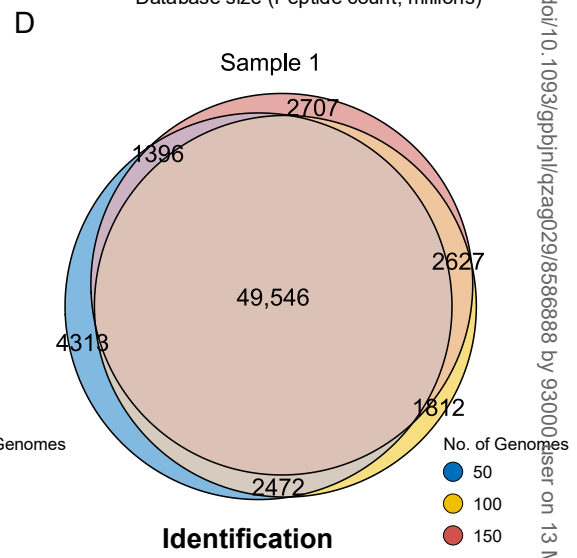
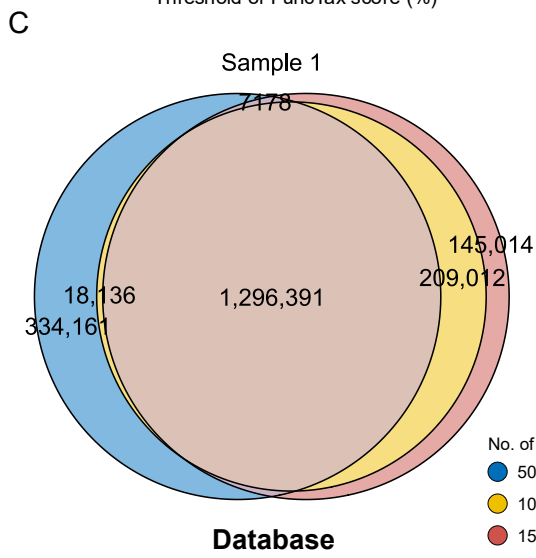
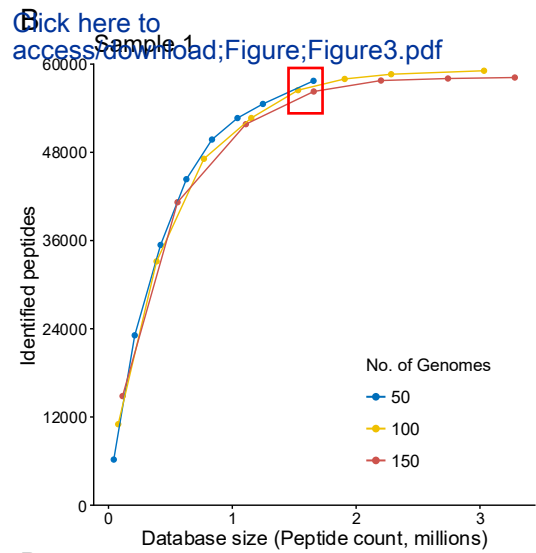
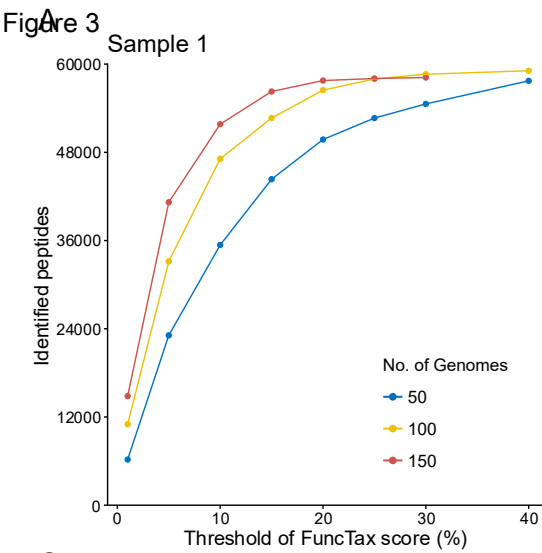
### 691 **Figure 8 Number of peptides identified by both methods**

692 Mean values from 79 diaPASEF samples were shown. For DDA-based method, the peptides  
693 identified as derived from human proteins are removed. The specific values can be found in File  
694 S3.

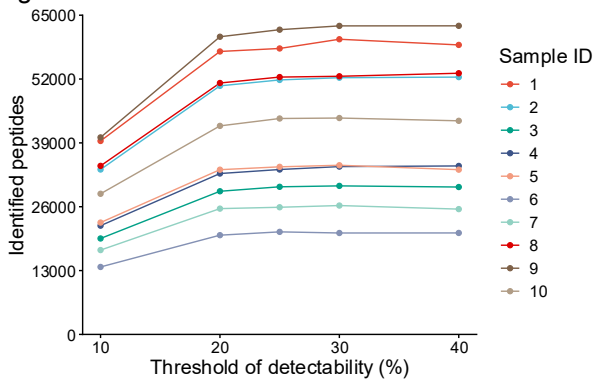
Figure 1



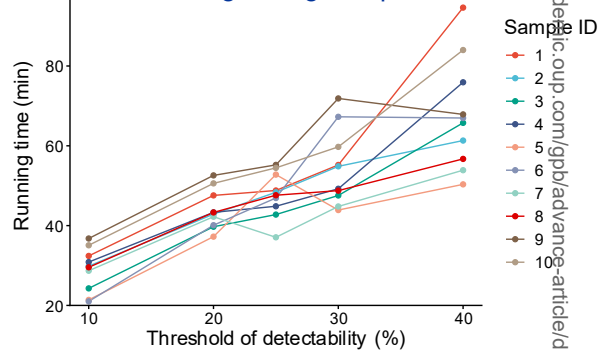




**Figure 4**



Click here to access/download;Figure;Figure4.pdf



**C**

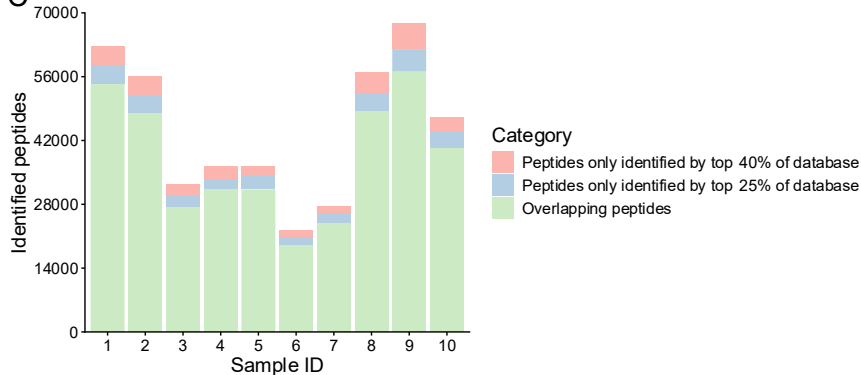
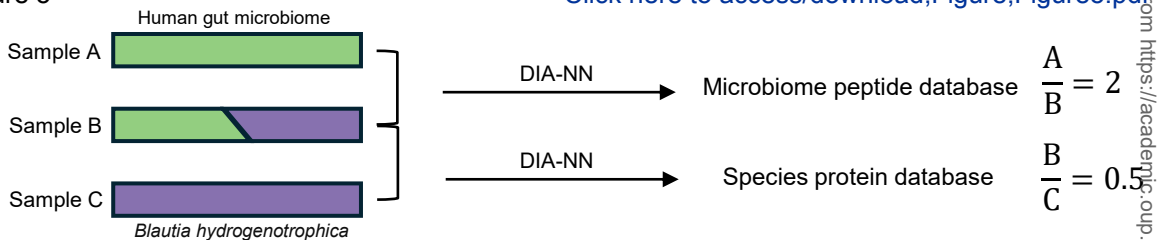
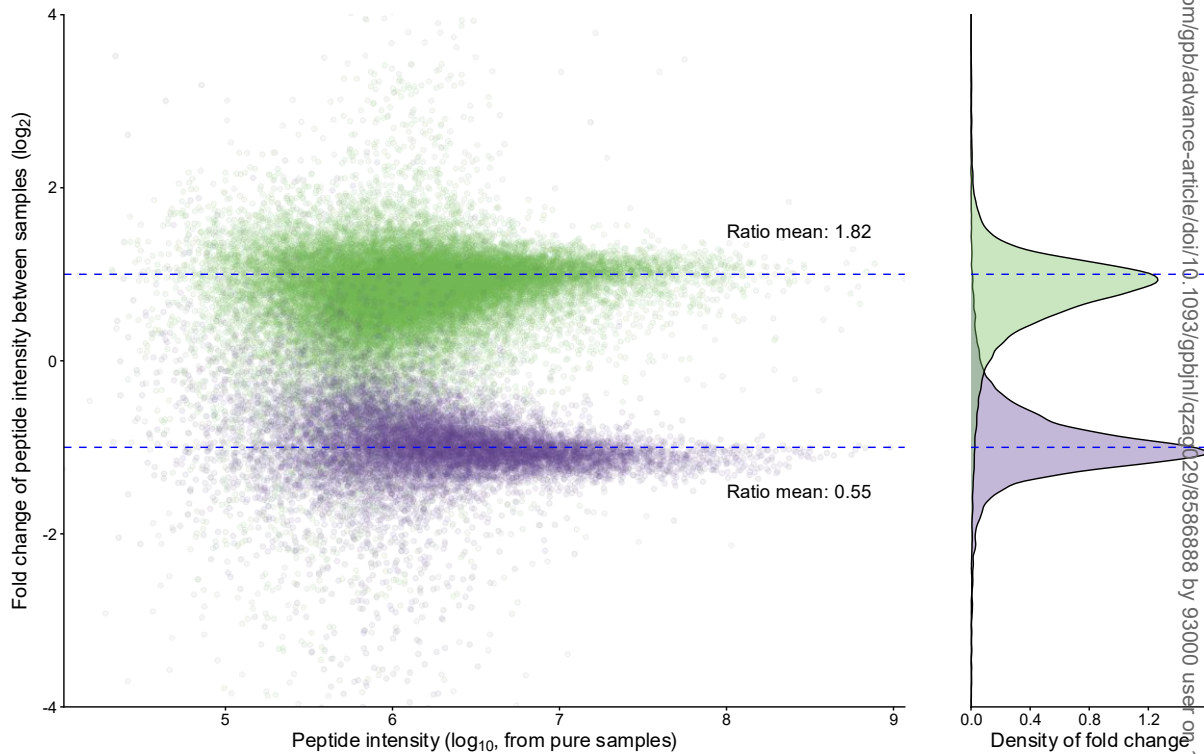


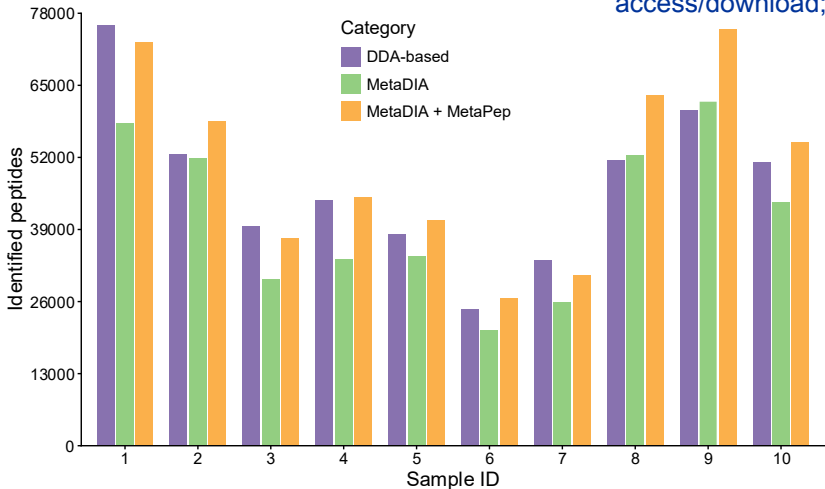
Figure 5



B



Downloaded from https://academic.oup.com/gpb/advance-article/doi/10.1093/gpb/ij/qzaa029/8586888 by 93000 user on 11 July 2020

**Figure 6**

Click here to  
access/download;Figure;Figure6.pdf

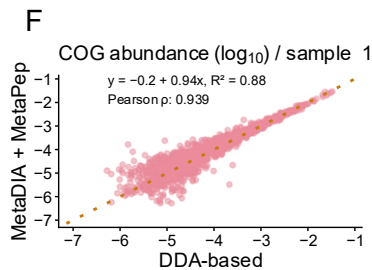
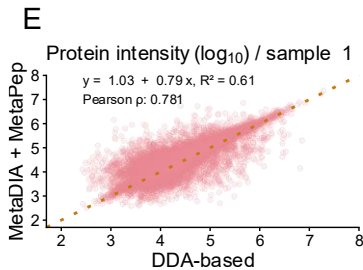
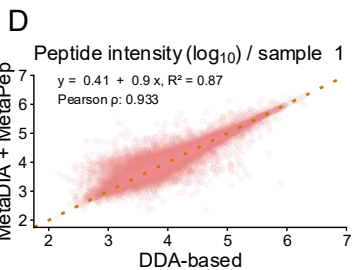
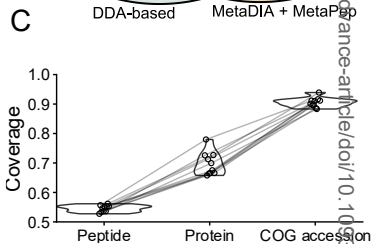
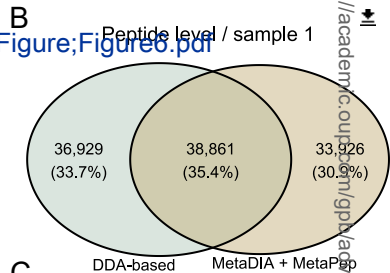
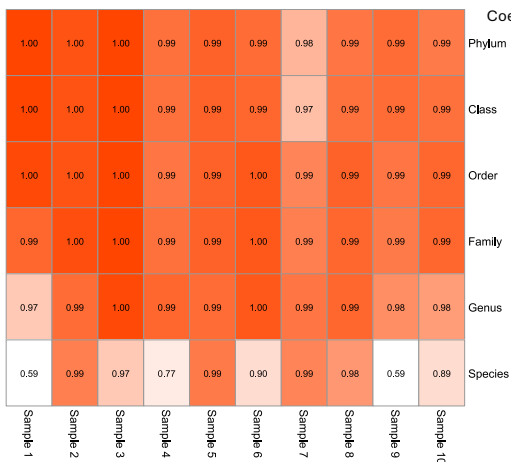
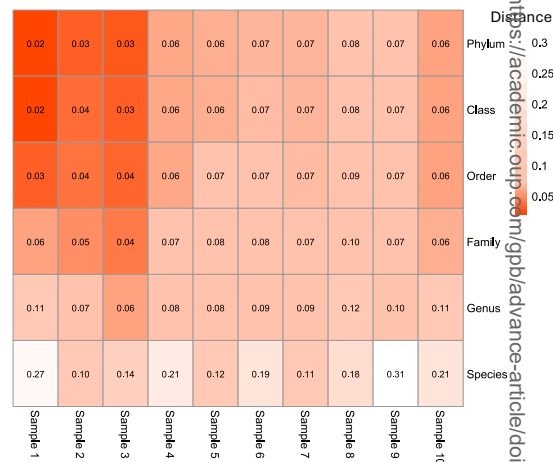


Figure 7

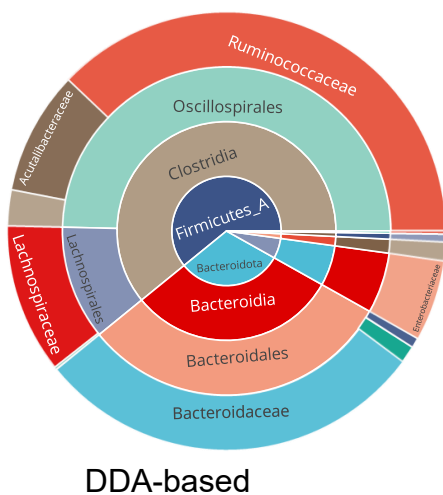
## Pearson coefficient



## Bray-Curtis distance



C



D

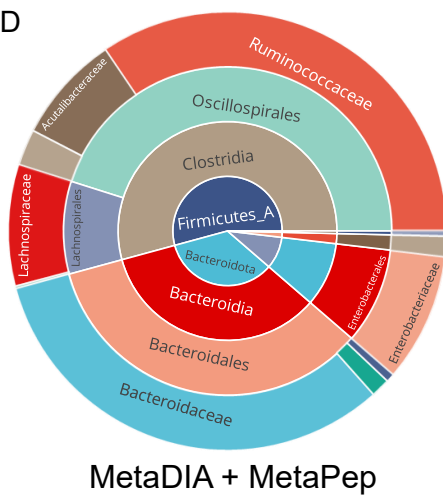


Figure 8

