

# The Microbiologist's Guide to Metaproteomics

Tim Van Den Bossche<sup>1,2</sup>, Jean Armengaud<sup>3</sup>, Dirk Benndorf<sup>4,5,6</sup>, Jose Alfredo Blakeley-Ruiz<sup>7</sup>, Madita Brauer<sup>8,9</sup>, Kai Cheng<sup>10</sup>, Marybeth Creskey<sup>11</sup>, Daniel Figeys<sup>10,12,13</sup>, Lucia Grenga<sup>3</sup>, Timothy J Griffin<sup>14</sup>, Céline Henry<sup>15</sup>, Robert L Hettich<sup>16</sup>, Tanja Holstein<sup>1,2,17</sup>, Pratik D Jagtap<sup>14</sup>, Nico Jehmlich<sup>18</sup>, Manuel Kleiner<sup>7</sup>, Benoit J Kunath<sup>7,19</sup>, Xuxa Malliet<sup>1,2</sup>, Lennart Martens<sup>1,2</sup>, Subina Mehta<sup>14</sup>, Bart Mesuere<sup>20</sup>, Zhibin Ning<sup>10</sup>, Alessandro Tanca<sup>21,22</sup>, Sergio Uzzau<sup>21,22</sup>, Pieter Verschaffelt<sup>1,2,20</sup>, Jing Wang<sup>23</sup>, Paul Wilmes<sup>19,24</sup>, Xu Zhang<sup>11</sup>, Xin Zhang<sup>23</sup>, and Leyuan Li<sup>23\*</sup>, on behalf of the Metaproteomics Initiative

(1) Department of Biomolecular Medicine, Faculty of Medicine and Health Sciences, Ghent University, 9052 Ghent, Belgium

(2) VIB-UGent Center for Medical Biotechnology, VIB, 9052 Ghent, Belgium

(3) Département Médicaments Et Technologies Pour La Santé (DMTS), Université Paris-Saclay, CEA, INRAE, SPI, 30200, Bagnols-Sur-Cèze, France

(4) Applied Biosciences and Process Engineering, Anhalt University of Applied Sciences, Köthen, Germany

(5) Bioprocess Engineering, Max Planck Institute for Dynamics of Complex Technical Systems, Magdeburg, Germany

(6) Bioprocess Engineering, Otto von Guericke University, Magdeburg, Germany

(7) Department of Plant and Microbial Biology, North Carolina State University, Raleigh, USA

(8) Molecular Disease Mechanisms Group, Department of Life Sciences and Medicine, University of Luxembourg, Esch-sur-Alzette, Luxembourg

(9) Institute for Advanced Studies, University of Luxembourg, Esch-sur-Alzette, Luxembourg

(10) School of Pharmaceutical Sciences, Faculty of Medicine, University of Ottawa, Ottawa, Canada

(11) Biologic and Radiopharmaceutical Drugs Directorate, Health Products and Food Branch, Health Canada, Ottawa, Canada

(12) Quadram Institute Bioscience, Norwich Research Park, Norwich, Norfolk, UK

(13) University of East Anglia, Norwich, Norfolk, UK

(14) Department of Biochemistry, Molecular Biology, and Biophysics, University of Minnesota. Minneapolis, MN 55455

(15) Université Paris-Saclay, INRAE, AgroParisTech, Micalis Institute, PAPPSO, 78350, Jouy-en-Josas, France

(16) Biosciences Division, Oak Ridge National Laboratory, Oak Ridge, TN, USA

(17) Data Competence Center MF2, Robert-Koch-Institut, Berlin, Germany

(18) Helmholtz-Centre for Environmental Research - UFZ GmbH, Department of Molecular Toxicology, Leipzig, Germany

(19) Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Esch-sur-Alzette, Luxembourg

(20) Department of Applied Mathematics, Statistics and Computer Science, Faculty of Sciences, Ghent University, 9000 Ghent, Belgium

42 (21) Department of Biomedical Sciences, University of Sassari, 07100 Sassari, Italy  
43 (22) Unit of Microbiology and Virology, University Hospital of Sassari, 07100 Sassari, Italy  
44 (23) State Key Laboratory of Medical Proteomics, Beijing Proteome Research Center, National Center for  
45 Protein Sciences (Beijing), Beijing Institute of Lifeomics, 102206 Beijing, China  
46 (24) Department of Life Sciences and Medicine, Faculty of Science, Technology and Medicine, University  
47 of Luxembourg, Belvaux, Luxembourg  
48  
49 \* Correspondence should be addressed to Leyuan Li ([lilevuan@ncpsb.org.cn](mailto:lilevuan@ncpsb.org.cn))  
50

## 51 Abstract

52 Metaproteomics is an emerging approach for studying microbiomes, offering the ability to  
53 characterize proteins that underpin microbial functionality within diverse ecosystems. As  
54 the primary catalytic and structural components of microbiomes, proteins provide unique  
55 insights into the active processes and ecological roles of microbial communities. By  
56 integrating metaproteomics with other omics disciplines, researchers can gain a  
57 comprehensive understanding of microbial ecology, interactions, and functional dynamics.  
58 This review, developed by the Metaproteomics Initiative ([www.metaproteomics.org](http://www.metaproteomics.org)), serves  
59 as a practical guide for both microbiome and proteomics researchers, presenting key  
60 principles, state-of-the-art methodologies, and analytical workflows essential to  
61 metaproteomics. Topics covered include, among others, experimental design, sample  
62 preparation, mass spectrometry techniques, data analysis strategies, and statistical  
63 approaches.

64  
65  
66  
67

## 68 Table of Contents

69	<b>Table of Contents</b> .....	<b>3</b>
70	<b>1. Why metaproteomics?</b> .....	<b>5</b>
71	<b>2. Basics of proteomics</b> .....	<b>8</b>
72	<b>3. Experimental methods in metaproteomics</b> .....	<b>9</b>
73	3.1. Experiment design .....	11
74	3.1.1. Aligning experimental design with the scientific question .....	11
75	3.1.2. Reproducibility & statistics .....	14
76	3.2. Sample collection, preservation and storage prior to preprocessing .....	15
77	3.2.1 Sample collection and preservation .....	15
78	3.2.2. Storage conditions to maintain sample integrity .....	16
79	3.3. Sample preprocessing .....	17
80	3.4 Protein sample preparation: from extraction to digestion .....	19
81	3.4.1 Cell lysis and protein extraction .....	19
82	3.4.2 Protein clean-up: precipitation and alternative methods .....	21
83	3.4.3 Measuring protein concentration .....	22
84	3.4.4 Protein digestion .....	23
85	3.5 Separation and fractionation techniques .....	24
86	3.5.1 On-line and off-line peptide fractionation .....	25
87	3.5.2 Enrichment of peptides with post-translational modifications .....	26
88	3.5.3 Protein, cell-level and functional fractionation techniques .....	27
89	3.6 Automation .....	28
90	3.6.1 Microbial cell disruption and protein extraction .....	28
91	3.6.2 Protein digestion and peptide clean-up .....	29
92	3.6.3 Multiplexing .....	29
93	3.7 Mass spectrometry data acquisition methods .....	30
94	3.7.1 DDA .....	31
95	3.7.2 DIA .....	32
96	3.7.3 Critical parameters to optimize the HPLC and MS methods .....	33
97	3.7.4 Quality control of LC-MS/MS .....	36
98	3.7.5 Data management and data sharing .....	38

99	<b>4. Computational analysis of metaproteomics data.....</b>	<b>39</b>
100	4.1 Peptide identification, protein inference and quantification .....	39
101	4.1.1 Peptide identification with proteomics search engines .....	39
102	4.1.2 Database construction or selection .....	44
103	4.1.3 PSM FDR control.....	49
104	4.1.4 Protein inference .....	51
105	4.1.5 Protein quantification.....	53
106	4.1.6 DIA data analysis.....	56
107	4.2 Taxonomic and functional analysis .....	57
108	4.2.1 Taxonomic analysis.....	58
109	4.2.2 Functional analysis.....	58
110	4.2.3 Peptide-centric vs protein-centric approach .....	59
111	4.2.4 Metaproteomics tools for taxonomic and functional analysis .....	60
112	4.3 Downstream statistics .....	62
113	4.3.1 Identifying relevant scientific questions .....	63
114	4.3.2 Selecting appropriate levels of analytical insights .....	63
115	4.3.3 Data preprocessing strategies .....	65
116	4.3.4 Choosing data analysis methods .....	67
117	<b>5. A collaborative effort: writing a comprehensive review with members of the</b>	
118	<b>Metaproteomics Initiative .....</b>	<b>70</b>
119	<b>6. Conclusion .....</b>	<b>71</b>
120	<b>Author contributions.....</b>	<b>71</b>
121	<b>Abbreviations.....</b>	<b>73</b>
122	<b>Acknowledgements.....</b>	<b>74</b>
123	<b>Conflicts of Interest.....</b>	<b>75</b>
124	<b>References .....</b>	<b>75</b>
125		
126		

## 127 1. Why metaproteomics?

128 The importance of microbiomes in nearly all processes within the biosphere is increasingly  
129 clear. Composed of bacteria, bacteriophages, archaea, yeasts, fungi, protozoa, and viruses,  
130 microbiomes are highly diverse in taxonomic composition. A microbiome and its theater of  
131 activity—including microbial elements such as genes, transcripts, proteins, and  
132 metabolites—together form a microbiome (Berg et al. 2020). Microbiomes are, in most  
133 cases, highly structured in both membership and function. This underscores the need to  
134 understand microbiomes and their interactions with their environment or eukaryotic hosts,  
135 whether beneficial or harmful. However, the complexity of these systems challenges  
136 traditional research tools, particularly cultivation-dependent approaches, which, given the  
137 wealth of intra-organism interactions, are not scalable for large-scale microbiome studies.

138 The rapid advancement of omics-based approaches has opened new avenues for systems  
139 biology-based research into the complexity of microbiomes. Shotgun metagenomics, in  
140 particular, has proven to be a powerful tool, offering much deeper insights than older  
141 techniques such as 16S rRNA gene amplicon sequencing. Metagenomics enables the  
142 discovery of complete genomic inventories, even for uncultured microorganisms, revealing  
143 the metabolic and physiological capabilities of a microbiome. However, it is limited to  
144 predicting functions rather than identifying active processes. To overcome this limitation,  
145 omics approaches such as metatranscriptomics, metaproteomics, and metabolomics  
146 provide essential insights into actual gene expression and activity under specific conditions.  
147 Together, these techniques bridge the gap from taxonomic structure to genomic potential  
148 and dynamic, context-dependent functions.

149 Among these tools, metaproteomics enables the comprehensive analysis of the proteins  
150 expressed and functional in a microbiome, quantifies their abundances, and characterizes  
151 their modifications, interactions, and localizations (**Figure 1**). Proteins serve as the primary  
152 catalytic units and structural elements of microbiomes, making metaproteomics a direct  
153 reflection of the microbiome's phenotype. This approach provides a detailed functional  
154 description and examines specific protein changes associated with structure, homeostasis,  
155 and enzymatic activity. Differences in protein sequences allow researchers to determine  
156 the taxonomic origins of particular enzyme sets, linking functions to taxonomic units.

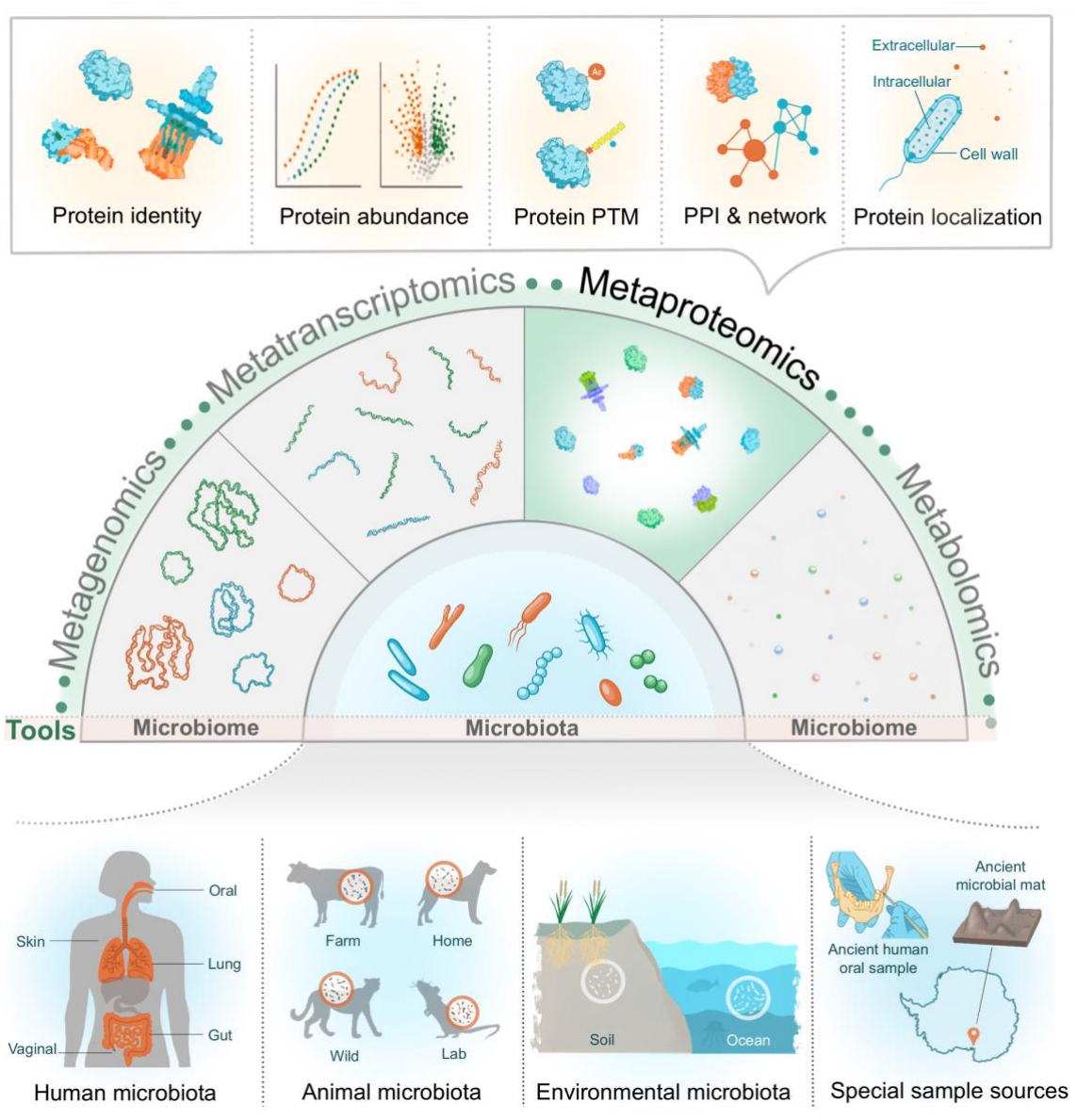
157 Metaproteomics can address several important questions such as:

- 158 • What are the metabolic and physiological processes of microorganisms in diverse  
159 habitats, including environmental, technical, and host-associated systems?

- 160 • How do microbiomes respond to changing conditions, as reflected by differential  
161 protein expression?
- 162 • How do microbes interact with their environment, including extracellular and  
163 intracellular protein dynamics?
- 164 • What post-translational modifications (PTMs) regulate protein activity and structure?
- 165 • How do microbiome phenotypes change over time or across spatial scales?
- 166 • How can stable isotope information from metaproteomes represent microbial activity  
167 and substrate utilization (Justice et al. 2014; Kleiner et al. 2023)?

168 While ongoing technological advancements are driving rapid progress, metaproteomics has  
169 already been successfully applied in the context of many impactful studies. It has  
170 contributed to fundamental understanding of microbial ecology, host-microorganism  
171 interactions, and disease mechanisms (Wolf et al. 2023). It has also improved  
172 biotechnological processes such as anaerobic digestion and wastewater treatment  
173 (Kleikamp et al. 2023; Justice et al. 2014; Heyer et al. 2024; Francesco Delogu et al. 2024),  
174 supported environmental monitoring (Pan, Wattiez, and Gillan 2024), and improved  
175 agricultural productivity (Andersen et al. 2021; Xue et al. 2024). Furthermore, it has  
176 applications in describing historical heritage and solving forensic questions (Jarman et al.  
177 2018). Readers interested in further details on the benefits of metaproteomics can explore  
178 several recommended reviews (Heintz-Buschart and Wilmes 2018; Sun, Ning, and Figeys  
179 2024; Herbst et al. 2016; Hettich et al. 2013; Kleiner 2019) and perspectives on its future  
180 (Van Den Bossche, Arntzen, et al. 2021; Wilmes, Heintz-Buschart, and Bond 2015; X.  
181 Zhang and Figeys 2019; Armengaud 2023).

182 This review, prepared by the Metaproteomics Initiative, aims to serve as a practical and  
183 accessible guide to metaproteomics. A detailed overview of the organization and  
184 presentation of this collaborative work is provided in **Section 5**, highlighting our dedication  
185 to delivering a comprehensive and valuable resource for the microbiome research  
186 community.



187

188 **Figure 1. Overview of metaproteomics within the multi-meta-omics toolbox applied to diverse**  
 189 **microbiome research domains.** This figure highlights the role of metaproteomics in identifying  
 190 proteins, quantifying their abundances, detecting post-translational modifications (PTMs), mapping  
 191 protein-protein interactions (PPIs), and determining protein localizations. Metaproteomics  
 192 complements other omics approaches, including metagenomics, metatranscriptomics, and  
 193 metabolomics, to provide a comprehensive understanding of microbial systems. Examples of  
 194 microbiome research domains include the human microbiome (oral, skin, gut, lung, vaginal), animal  
 195 microbiomes (farm, wild, and laboratory animals), environmental microbiomes (soil, ocean), and  
 196 special sample sources (e.g., ancient microbiome samples).

## 197 2. Basics of proteomics

198 Proteins are the essential structures and machinery that execute the instructions encoded  
199 in DNA, performing tasks ranging from catalyzing biochemical reactions to providing  
200 structural support. The term "proteome" refers to the complete set of proteins expressed in  
201 a cell, tissue, or organism (Wilkins et al. 1996). Proteomics, as a field, seeks to uncover the  
202 identities, quantities, structures, interactions, and modifications of proteins to better  
203 understand their roles in biological systems.

204 Although the term "proteome" was coined in the mid-1990s, its foundations lie in decades  
205 of protein biochemistry research that continues to shape modern proteomics. One of the  
206 earliest applications of proteomics combined gel electrophoresis (1D and 2D) with mass  
207 spectrometry techniques such as MALDI and ESI-MS/MS (James et al. 1993). Initially,  
208 protein samples were separated on a combination of 1D and 2D gels. One gel was electro-  
209 blotted onto a nitrocellulose membrane and stained using amido black, while the other gel  
210 was silver-stained for higher sensitivity. Protein bands or spots were excised from the  
211 nitrocellulose membrane, digested with trypsin, and identified using mass spectrometry.  
212 Aligning the nitrocellulose membrane with the silver-stained gel allowed researchers to  
213 locate bands that were difficult to visualize on the less-sensitive stain. Subsequent  
214 improvements, such as in-gel digestion, eliminated the need for electro-blotting. Early  
215 proteomics efforts also gave rise to software tools that automated protein identification, and  
216 therefore replacing manual annotation of peptide sequences. Many of these early  
217 innovations however, formed the basis for modern proteomics workflows.

218 The development of gel-free proteomics marked a significant advancement in the field. This  
219 approach bypasses gel-based separation, proceeding directly from protein extraction to  
220 digestion and mass spectrometry. Gel-free methods catalyzed a wave of new techniques,  
221 reagents (e.g., SILAC, ICAT, ITRAQ), and software, which collectively improved protein  
222 identification, PTM analysis, quantitation, and multiplexing. Tasks that were once labor-  
223 intensive with 2D gel MS became faster and more accessible through gel-free workflows.  
224 Moreover, mass spectrometers, which were initially optimized for small molecule research,  
225 were adapted for proteomics. Over the past 15 years, proteomics-dedicated mass  
226 spectrometers have been developed, offering greater speed, sensitivity, and accuracy in  
227 peptide identification and quantitation.

228 Proteomics today falls into two broad methodological categories: shotgun (or bottom-up)  
229 proteomics (Diz and Sánchez-Marín 2021) and top-down proteomics (Habeck et al. 2024).  
230 Shotgun proteomics, the more widely used approach, involves enzymatic digestion of



231 proteins into peptides, which are analyzed by mass spectrometry. This method is robust  
232 and effective for protein identification and quantification. In contrast, top-down proteomics  
233 directly analyzes intact proteins, providing insights into sequences, structures, and  
234 modifications. Although top-down proteomics offers unique advantages, it is technically  
235 demanding, less commonly used in single-species proteomics, and not currently applied in  
236 metaproteomics.

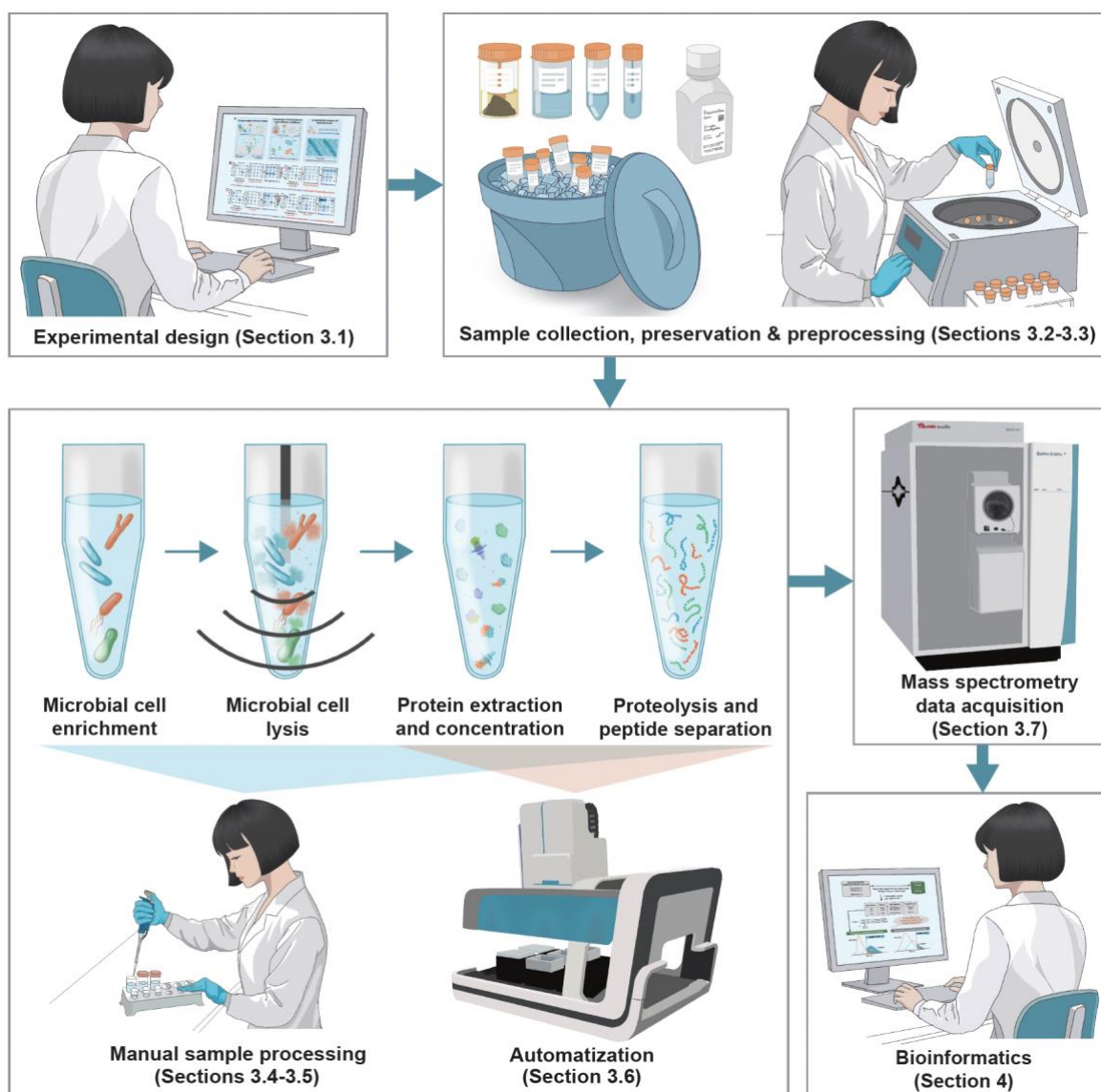
237 A typical bottom-up proteomics workflow begins with the enzymatic digestion of proteins,  
238 most commonly using trypsin, into smaller peptides. These peptides are separated through  
239 liquid chromatography and analyzed by tandem mass spectrometry (LC-MS/MS). In the  
240 mass spectrometer, the peptides are ionized, and their intact forms are detected to generate  
241 MS1 spectra. The peptides are further fragmented to produce MS2 spectra, which are  
242 analyzed by proteomics software. In most cases, database searches match these spectra  
243 to theoretical spectra derived from protein databases. This approach enables the  
244 identification and quantification of peptides and their corresponding proteins. For those  
245 seeking a deeper understanding of proteomics, numerous resources and reviews provide  
246 detailed insights into the field (Matthiesen and Bunkenborg 2013; Shuken 2023; Jiang et al.  
247 2024; Sinitcyn, Rudolph, and Cox 2018).

### 248 3. Experimental methods in metaproteomics

249 Metaproteomics expands upon proteomics techniques, leveraging high-resolution LC-  
250 MS/MS instruments (Gómez-Varela et al. 2023; Dumas et al. 2024) and accompanying  
251 software tools for mass spectra identification. However, metaproteomics goes beyond the  
252 straightforward application of proteomics to microbiome research. Its added complexity  
253 arises from the requirement to consider both species-specific and functional annotations for  
254 each protein. Additionally, the presence of protein homologs across phylogenetically related  
255 species within a single sample further complicates protein inference.

256 The key distinctions between proteomics and metaproteomics lie in the taxonomic and  
257 functional complexity of microbiomes, the vast size of microbiome databases, and the  
258 challenges associated with sample processing, as well as the identification and quantitation  
259 of peptides and proteins. Additionally, specialized bioinformatic and statistical tools are  
260 required to track both the taxonomic and functional annotations of peptides and proteins.  
261 These aspects, which are unique to metaproteomics, will be discussed in detail throughout  
262 the remainder of this article.

263 This section provides an essential foundational guide to start with metaproteomics studies  
 264 (**Figure 2**). We outline the basic principles for each step, starting with experimental design  
 265 (**Section 3.1**), followed by sample collection, preservation, and preprocessing (**Sections**  
 266 **3.2–3.3**). Protein sample preparation is then described, covering both manual workflows  
 267 (**Sections 3.4–3.5**) and automated workflows (**Section 3.6**). Next, we explain the basics of  
 268 MS data acquisition (**Section 3.7**), before delving into the detailed bioinformatics workflows  
 269 used in metaproteomics (**Sections 4.1–4.3**).



270

271 **Figure 2. Overview of key principles and workflows in metaproteomics, aligned with**  
 272 **corresponding subsections in this review.**

273

## 274 3.1. Experiment design

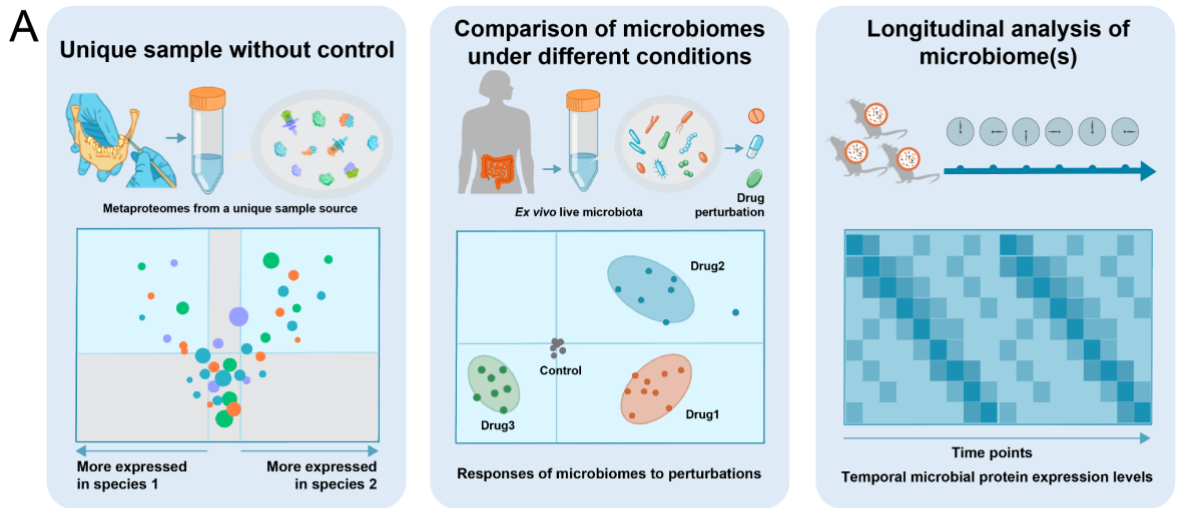
### 275 3.1.1. Aligning experimental design with the scientific question

276 A well-designed metaproteomics experiment forms the basis for generating meaningful  
277 insights that directly address the scientific question being studied. Most importantly, the  
278 experimental design must align with the specific scientific question being addressed and  
279 the resources available to answer that question. Broadly, three experimental scenarios can  
280 be outlined (**Figure 3A**):

281 **i) Unique sample without a control:** The goal here is to provide a comprehensive  
282 description of the taxonomic and functional units present in the sample, although  
283 comparisons with a control are not possible. Examples include desiccated material from a  
284 historical Antarctic ice core (Lezcano et al. 2022), a unique biofilm from an industrial storage  
285 pool (Pible et al. 2023), residues from an ancient tomb (Charlier et al. 2024), or medieval  
286 dental calculus (Jersie-Christensen et al. 2018) were analyzed using metaproteomics.  
287 Differential functional abundances among the identified microorganisms can reveal their  
288 metabolic specialization.

289 **ii) Comparison of microbiomes under different conditions:** This common approach  
290 highlights differences between conditions. Comparisons may involve two conditions (i.e.,  
291 condition A vs. condition B) or more complex setups with multiple conditions. Specific cases  
292 include dose-response analyses, where a single parameter such as stress intensity is  
293 modified, or spatial comparisons. Examples include characterizing microbial communities  
294 along a 5,000 km Pacific Ocean transect (Saunders et al. 2022) or analyzing microbiome  
295 responses to various xenobiotics *in vitro* (L. Li et al. 2020).

296 **iii) Longitudinal analysis of a single microbiome or multiple microbiomes:** This  
297 strategy captures temporal dynamics within a microbial community, and potentially the  
298 host's response, by analyzing the same microbiome at different time points. A more  
299 complex approach examines temporal changes across multiple conditions or sampling sites.  
300 Examples include monitoring gut microbiomes in Crohn's disease patients post-resection  
301 surgery over one year (Blakeley-Ruiz et al. 2019) or monthly analyses of specialized  
302 microbiomes in a two-stage anaerobic digester for lignocellulose breakdown, tracking the  
303 dynamics between hydrolytic and methanogenic subsystems (Heyer et al. 2024).



**B**

<table border="1" style="font-size: small;"> <thead> <tr><th>Taxa</th><th>Abundance</th></tr> </thead> <tbody> <tr><td>3</td><td>0 0 0 0</td></tr> <tr><td>0</td><td>2 0 0 0</td></tr> <tr><td>0</td><td>0 5 0 0</td></tr> <tr><td>0</td><td>0 0 0 1</td></tr> </tbody> </table> <p>Taxonomic Composition <math>t_0</math> (variable)</p>	Taxa	Abundance	3	0 0 0 0	0	2 0 0 0	0	0 5 0 0	0	0 0 0 1	<table border="1" style="font-size: small;"> <thead> <tr><th>Gene1</th><th>Gene2</th><th>Gene3</th><th>Gene4</th></tr> </thead> <tbody> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Genome contents (nearly constant)</p>	Gene1	Gene2	Gene3	Gene4	1	1	1	1	1	1	0	1	1	1	1	0	1	1	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Gene1</th><th>Gene2</th><th>Gene3</th><th>Gene4</th></tr> </thead> <tbody> <tr><td>3</td><td>3</td><td>3</td><td>3</td></tr> <tr><td>2</td><td>2</td><td>0</td><td>2</td></tr> <tr><td>5</td><td>5</td><td>5</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Metagenome <math>t_0</math></p>	Gene1	Gene2	Gene3	Gene4	3	3	3	3	2	2	0	2	5	5	5	0	1	1	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Taxa</th><th>Abundance</th></tr> </thead> <tbody> <tr><td>5</td><td>0 0 0 0</td></tr> <tr><td>0</td><td>2 0 0 0</td></tr> <tr><td>0</td><td>0 3 0 0</td></tr> <tr><td>0</td><td>0 0 0 2</td></tr> </tbody> </table> <p>Taxonomic Composition <math>t_1</math> (variable)</p>	Taxa	Abundance	5	0 0 0 0	0	2 0 0 0	0	0 3 0 0	0	0 0 0 2	<table border="1" style="font-size: small;"> <thead> <tr><th>Gene1</th><th>Gene2</th><th>Gene3</th><th>Gene4</th></tr> </thead> <tbody> <tr><td>1</td><td>1</td><td>1</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Genome contents (nearly constant)</p>	Gene1	Gene2	Gene3	Gene4	1	1	1	1	1	1	0	1	1	1	1	0	1	1	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Gene1</th><th>Gene2</th><th>Gene3</th><th>Gene4</th></tr> </thead> <tbody> <tr><td>5</td><td>5</td><td>5</td><td>5</td></tr> <tr><td>2</td><td>2</td><td>0</td><td>2</td></tr> <tr><td>3</td><td>3</td><td>3</td><td>0</td></tr> <tr><td>2</td><td>2</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Metagenome <math>t_1</math></p>	Gene1	Gene2	Gene3	Gene4	5	5	5	5	2	2	0	2	3	3	3	0	2	2	0	0
Taxa	Abundance																																																																																																								
3	0 0 0 0																																																																																																								
0	2 0 0 0																																																																																																								
0	0 5 0 0																																																																																																								
0	0 0 0 1																																																																																																								
Gene1	Gene2	Gene3	Gene4																																																																																																						
1	1	1	1																																																																																																						
1	1	0	1																																																																																																						
1	1	1	0																																																																																																						
1	1	0	0																																																																																																						
Gene1	Gene2	Gene3	Gene4																																																																																																						
3	3	3	3																																																																																																						
2	2	0	2																																																																																																						
5	5	5	0																																																																																																						
1	1	0	0																																																																																																						
Taxa	Abundance																																																																																																								
5	0 0 0 0																																																																																																								
0	2 0 0 0																																																																																																								
0	0 3 0 0																																																																																																								
0	0 0 0 2																																																																																																								
Gene1	Gene2	Gene3	Gene4																																																																																																						
1	1	1	1																																																																																																						
1	1	0	1																																																																																																						
1	1	1	0																																																																																																						
1	1	0	0																																																																																																						
Gene1	Gene2	Gene3	Gene4																																																																																																						
5	5	5	5																																																																																																						
2	2	0	2																																																																																																						
3	3	3	0																																																																																																						
2	2	0	0																																																																																																						

A microbiota's metagenomic response to a perturbation is driven by the changes in **taxonomic composition alone**.

**C**

<table border="1" style="font-size: small;"> <thead> <tr><th>Taxa</th><th>Abundance</th></tr> </thead> <tbody> <tr><td>3</td><td>0 0 0 0</td></tr> <tr><td>0</td><td>2 0 0 0</td></tr> <tr><td>0</td><td>0 5 0 0</td></tr> <tr><td>0</td><td>0 0 0 1</td></tr> </tbody> </table> <p>Taxonomic Composition <math>t_0</math> (variable)</p>	Taxa	Abundance	3	0 0 0 0	0	2 0 0 0	0	0 5 0 0	0	0 0 0 1	<table border="1" style="font-size: small;"> <thead> <tr><th>Protein1</th><th>Protein2</th><th>Protein3</th><th>Protein4</th></tr> </thead> <tbody> <tr><td>1</td><td>1</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> <tr><td>1</td><td>0</td><td>1</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Proteome contents <math>t_0</math> (variable)</p>	Protein1	Protein2	Protein3	Protein4	1	1	0	1	1	1	0	0	1	0	1	0	1	1	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Protein1</th><th>Protein2</th><th>Protein3</th><th>Protein4</th></tr> </thead> <tbody> <tr><td>3</td><td>3</td><td>0</td><td>3</td></tr> <tr><td>2</td><td>2</td><td>0</td><td>0</td></tr> <tr><td>5</td><td>0</td><td>5</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Metaproteome <math>t_0</math></p>	Protein1	Protein2	Protein3	Protein4	3	3	0	3	2	2	0	0	5	0	5	0	1	1	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Taxa</th><th>Abundance</th></tr> </thead> <tbody> <tr><td>5</td><td>0 0 0 0</td></tr> <tr><td>0</td><td>2 0 0 0</td></tr> <tr><td>0</td><td>0 3 0 0</td></tr> <tr><td>0</td><td>0 0 0 2</td></tr> </tbody> </table> <p>Taxonomic Composition <math>t_1</math> (variable)</p>	Taxa	Abundance	5	0 0 0 0	0	2 0 0 0	0	0 3 0 0	0	0 0 0 2	<table border="1" style="font-size: small;"> <thead> <tr><th>Protein1</th><th>Protein2</th><th>Protein3</th><th>Protein4</th></tr> </thead> <tbody> <tr><td>3</td><td>2</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>0</td><td>0</td><td>1</td></tr> <tr><td>1</td><td>2</td><td>1</td><td>0</td></tr> <tr><td>1</td><td>2</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Proteome contents <math>t_1</math> (variable)</p>	Protein1	Protein2	Protein3	Protein4	3	2	0	1	1	0	0	1	1	2	1	0	1	2	0	0	<table border="1" style="font-size: small;"> <thead> <tr><th>Protein1</th><th>Protein2</th><th>Protein3</th><th>Protein4</th></tr> </thead> <tbody> <tr><td>15</td><td>10</td><td>0</td><td>5</td></tr> <tr><td>2</td><td>0</td><td>0</td><td>2</td></tr> <tr><td>3</td><td>6</td><td>3</td><td>0</td></tr> <tr><td>2</td><td>4</td><td>0</td><td>0</td></tr> </tbody> </table> <p>Metaproteome <math>t_1</math></p>	Protein1	Protein2	Protein3	Protein4	15	10	0	5	2	0	0	2	3	6	3	0	2	4	0	0
Taxa	Abundance																																																																																																								
3	0 0 0 0																																																																																																								
0	2 0 0 0																																																																																																								
0	0 5 0 0																																																																																																								
0	0 0 0 1																																																																																																								
Protein1	Protein2	Protein3	Protein4																																																																																																						
1	1	0	1																																																																																																						
1	1	0	0																																																																																																						
1	0	1	0																																																																																																						
1	1	0	0																																																																																																						
Protein1	Protein2	Protein3	Protein4																																																																																																						
3	3	0	3																																																																																																						
2	2	0	0																																																																																																						
5	0	5	0																																																																																																						
1	1	0	0																																																																																																						
Taxa	Abundance																																																																																																								
5	0 0 0 0																																																																																																								
0	2 0 0 0																																																																																																								
0	0 3 0 0																																																																																																								
0	0 0 0 2																																																																																																								
Protein1	Protein2	Protein3	Protein4																																																																																																						
3	2	0	1																																																																																																						
1	0	0	1																																																																																																						
1	2	1	0																																																																																																						
1	2	0	0																																																																																																						
Protein1	Protein2	Protein3	Protein4																																																																																																						
15	10	0	5																																																																																																						
2	0	0	2																																																																																																						
3	6	3	0																																																																																																						
2	4	0	0																																																																																																						

A microbiota's metaproteome response is driven by changes in both **taxonomic composition** and **proteome contents**.

304

305 **Figure 3. Metaproteomic experimental designs and their comparison with metagenomics in**

306 **studying microbiome dynamics.** (A) Overview of common metaproteomic experimental designs.

307 The left panel illustrates the comparison of microbial protein expression between species within a

308 unique sample source, lacking a control. The middle panel compares microbiomes under varying

309 conditions, such as drug treatments, using *ex vivo* microbiomes to assess microbial responses. The

310 right panel shows longitudinal studies that monitor temporal changes in microbial protein expression

311 over time. (B) Metagenomic responses to perturbations, showing shifts in taxonomic composition

312 while assuming genome content remains relatively constant. (C) Metaproteomic responses to

313 perturbations, showing changes in both taxonomic composition and proteome content. This

314 approach captures microbial abundances and their functional contributions, providing deeper

315 insights into microbiome dynamics.

316

317 Some readers may already have experience designing experiments for metagenomics and  
318 understand its principles. In contrast, metaproteomics offers a different perspective on  
319 microbiome changes (**Figure 3B**). Metagenomics captures shifts driven by changes in  
320 taxonomic composition, as genomic content within a sample is relatively constant. This  
321 approach reveals species abundance and diversity but does not provide functional insights.  
322 Metaproteomics, on the other hand, measures not only taxonomic changes through taxon-  
323 specific peptide intensities but also dynamic functional responses through proteome  
324 variations across taxa. This makes metaproteomics particularly well-suited for comparing  
325 microbiomes under different conditions or for longitudinal studies.

326 When selecting conditions or time points for a kinetic analysis, careful consideration is  
327 essential. Comparisons between vastly different samples, such as a soil microbiome versus  
328 a human gut microbiome, are in general uninformative, while overly similar samples may  
329 show no significant differences. Selection should be guided by a clear rationale and  
330 preliminary observations. The reference condition or time point depends on the scientific  
331 question but may involve using a mixture of all samples as a reference. While this approach  
332 increases peptide diversity in the reference sample, it can complicate analysis if the full  
333 diversity is not captured by the analytical workflow (Armengaud 2023) as further detailed in  
334 **Sections 3.5 and 3.7**.

335 Potential confounding factors must also be accounted for during experimental design.  
336 Comprehensive metadata collection is critical, including information on sampling location,  
337 timing, storage, processing conditions, and data acquisition. Additional metadata, such as  
338 weather conditions on sampling days, patient medication, or health status, may also be  
339 essential for interpreting results. Additionally, researchers should also consider using  
340 additional material to create appropriate databases for matching spectra to peptides and  
341 for testing methodologies before processing all samples. More details on proteomics  
342 software and database creation are provided in **Section 4.1.1 and 4.1.2**, respectively.

343 Finally, while a limited number of metaproteomics studies have used metabolic labeling  
344 (e.g., to study host-microbiome or plant rhizosphere interactions (Z. Li et al. 2019; Smyth et  
345 al. 2020; Sachsenberg et al. 2015), this approach is often impractical for environmental or  
346 human microbiome samples. Metabolic labeling, as briefly mentioned in **Section 2**, involves  
347 incorporating heavy isotopes like  $^{15}\text{N}$  or  $^{13}\text{C}$  into proteins through labeled substrates,  
348 enabling the study of metabolic crosstalk and protein production rates. However, its limited  
349 applicability means that it is not further discussed in this review.

### 350 3.1.2. Reproducibility & statistics

351 The high complexity and heterogeneity of metaproteomics samples necessitate careful  
352 consideration of statistical power and steps to ensure reproducibility during experimental  
353 design. Biological, technical, and analytical replicates are key to producing reliable data and  
354 accurate interpretations. Increasing the number of biological replicates improves the ability  
355 to detect smaller differences, even in the presence of high variability. When only slight  
356 differences between conditions are expected, the use of pooled samples may also be  
357 considered. Technical and analytical replicates are necessary to account for noise  
358 introduced during measurement. It is advisable to first evaluate the variability of sample  
359 preparation and the analytical workflow using a representative sample. Additionally,  
360 randomizing the order of samples before LC-MS/MS analysis reduces the risk of bias due  
361 to the sequence in which they are processed (Nakayasu et al. 2021). For cases where  
362 specific sources of variability, such as batch effects, are known, blocked randomization is  
363 preferable to further minimize bias (Oberg and Vitek, 2009). Rigorous quality control (QC)  
364 is essential during the LC-MS/MS phase of the metaproteomics workflow to ensure data  
365 reliability and consistency. **Section 3.7.4** provides further details on these QC procedures.

366 Determining the appropriate number of biological replicates is essential to detect  
367 meaningful biological differences, such as variations in taxon biomasses, protein  
368 abundances, or metabolic pathways. Power analysis is typically used to calculate the  
369 required sample size, but it can be challenging in metaproteomics due to the complexity of  
370 experimental designs and the inherent variability of samples. When precise endpoints are  
371 unavailable, rough estimates from similar studies can serve as a guide. Power analysis  
372 considers several key factors: the effect size, which reflects the expected magnitude of  
373 differences between groups and helps determine the necessary sample size; the  
374 significance level ( $\alpha$ ), usually set at 0.05 to allow a 5% risk of false positives; statistical  
375 power ( $1 - \beta$ ), often set at 0.8 or higher to reduce the likelihood of failing to detect a true  
376 effect; and the variability in the data, which can be estimated from pilot studies or previous  
377 literature on comparable experiments. In studies involving complex microbial communities,  
378 deriving precise sample size estimates may be impractical, but approximate estimates  
379 remain a valuable approach (Ferdous et al. 2022). Conducting power analysis is critical for  
380 avoiding underpowered studies and ensuring efficient use of resources (Levin 2011;  
381 Ferdous et al. 2022).

## 382 3.2. Sample collection, preservation and storage prior to 383 preprocessing

### 384 3.2.1 Sample collection and preservation

385 Metaproteomics has been applied to a variety of samples, including microbial communities  
386 from environmental niches such as water, soil, sewage, aerosols, and rocks (Starke,  
387 Jehmlich, and Bastida 2019; Nebauer, Pearson, and Neilan 2024). It has also been used  
388 to analyze microbiomes in fermented foods and beverages (L. Yang, Fan, and Xu 2020;  
389 Okeke et al. 2021) and in associations with various higher eukaryotes, including arachnids,  
390 insects, worms, mollusks, fish, plants, birds, and mammals (Ezzeldin et al. 2019; Andersen  
391 et al. 2021). In mammals and other vertebrates, metaproteomics has been applied to  
392 numerous body sites across the digestive, respiratory, and urogenital systems (Y. Wang et  
393 al. 2020; Wolf et al. 2023). However, many microbiomes remain unexplored by  
394 metaproteomics.

395 The choice of collection method significantly influences the resulting metaproteomic profile  
396 by altering the ratios of microbial to non-microbial components and the relative abundances  
397 of microbial taxa. Collection strategies also introduce operator-dependent variability,  
398 making user-friendly devices especially valuable for self-sampling of clinical specimens.  
399 Microbiome samples are often collected directly into sterile tubes or containers. This  
400 method is common for non-invasive clinical samples, such as feces, saliva, sputum, and  
401 urine, which can often be self-collected by study participants (Long et al. 2020; Arıkan et al.  
402 2023; Graf et al. 2021; XiaoLian Xiao et al. 2022). For clinical specimens requiring surface  
403 sampling, swabs, spatulas, or syringes are often used for oral, nasal, and cervicovaginal  
404 samples (Chen et al. 2024; Bihani et al. 2023; Berard et al. 2023), while periodontal curettes  
405 or paper strips are used for tooth- and gingiva-associated microbiomes (Rabe et al. 2022;  
406 Xiaolian Xiao et al. 2023). Invasive procedures, such as bronchoalveolar lavage,  
407 endotracheal aspiration (Pathak et al. 2020), intestinal biopsies (Jabbar et al. 2021), colonic  
408 luminal aspirates (X. Zhang et al. 2020), and surgical collection of colonic contents (Tanca  
409 et al. 2022), are necessary for some specimens. Similarly, gastrointestinal fistulation  
410 (Deusch et al. 2017) and post-mortem dissection (Haange et al. 2019) are used for  
411 collecting samples from laboratory or field animals. For environmental samples, specialized  
412 devices such as quartz filters for bioaerosols (Meyer et al. 2023) and large-volume water  
413 transfer/filter systems for aquatic environments (L.-F. Kong et al. 2021; S. Wang et al. 2024)  
414 are commonly employed. More complex ecosystems may require multi-step collection  
415 protocols (Aylward et al. 2012).

416 Microbiome sampling inherently involves the translocation of microbial communities from  
417 their native environment to laboratory conditions. During this transition, microbial  
418 communities are highly sensitive to environmental changes such as temperature, humidity,  
419 and exposure to chemical or biological agents. These factors can induce substantial  
420 alterations in the metaproteome profile. To minimize artifacts, protein extraction should  
421 ideally occur immediately after sampling. However, immediate processing is often  
422 impractical, particularly in large-scale studies or field collections. In such cases, proper  
423 transport and storage procedures are crucial to preserving the microbiome's original  
424 biological functions. This is especially important for low-biomass or low-diversity  
425 microbiomes, which are more vulnerable to rapid shifts in their composition and activity due  
426 to external stimuli.

### 427 3.2.2. Storage conditions to maintain sample integrity

428  
429 Proper storage is critical to preserving the integrity of microbial proteins and ensuring  
430 reliable downstream analyses. Exposure to environmental changes, such as air exposure,  
431 temperature fluctuations, or nutrient depletion, can significantly alter protein profiles,  
432 leading to misleading results. For instance, air exposure can introduce oxidative stress and  
433 enrich bacterial superoxide dismutase enzymes, which may bias colorectal cancer studies  
434 by mimicking disease-specific characteristics (Long et al. 2020). Therefore, appropriate  
435 storage immediately after sample collection is essential to maintain the microbiome's  
436 original state.

437 The standard practice for preserving metaproteomic samples involves flash-freezing in  
438 liquid nitrogen, followed by storage at  $-80^{\circ}\text{C}$ . This approach minimizes molecular  
439 degradation and prevents alterations in protein abundances. While this method is highly  
440 effective, some experimental setups do not allow for immediate freezing. In such cases,  
441 alternative preservation methods may be employed. Solutions like PBS (Delgado-Diaz et  
442 al. 2022), Amies liquid medium (Bankvall et al. 2023), NAP buffer (Mordant and Kleiner  
443 2021), and other commercially available liquid reagents (Birse et al. 2020) have been tested  
444 for their ability to enhance storage conditions or enable room-temperature preservation in  
445 metaproteomics. Protease inhibitors are often added to biological fluids such as saliva to  
446 prevent uncontrolled proteolysis (Ruan et al. 2021). RNAlater or RNAlater-like treatments  
447 have shown potential for preserving protein profiles in intestinal and marine samples,  
448 although with conflicting results (Mordant and Kleiner 2021; Jensen, Wippler, and Kleiner  
449 2021; Saito et al. 2011). Regardless of the method used, compatibility with downstream  
450 protein extraction, digestion, and analysis steps is crucial. Common pitfalls, including



451 polyethylene glycol (PEG) contamination from plasticware, keratin contamination from  
452 handling, and interference from detergents or salts, must be carefully managed.

453 Alternative long-term storage strategies, such as freeze-drying or storing samples at  $-20^{\circ}\text{C}$ ,  
454 in liquid nitrogen tanks, or as lyophilized powders, also require careful evaluation. These  
455 approaches may be suitable for some sample types but may not consistently maintain  
456 protein integrity. For example, frozen intact stool material has been shown to be more stable  
457 than extracted proteins when stored at  $-80^{\circ}\text{C}$ , underscoring the importance of selecting  
458 storage strategies tailored to the specific sample type (Morris and Marchesi 2016).

459 It is important to note that the stability of proteins during storage is highly dependent on the  
460 sample type and storage conditions. For example, the activity and stability of soil proteins  
461 are influenced by temperature, duration of storage, and soil organic matter content (Bandick  
462 and Dick 1999; Keiblinger et al. 2016). For studies involving prolonged transport or storage,  
463 incorporating a straightforward mock community can provide valuable controls to assess  
464 sample stability and detect potential storage-induced changes (Nebauer, Pearson, and  
465 Neilan 2024).

### 466 3.3. Sample preprocessing

467 Sample preprocessing ensures the removal of contaminants and debris, which can hinder  
468 protein extraction, degrade analytical quality (Heyer et al. 2019), and dilute biologically  
469 relevant signals. This step, as in other gene expression measurement workflows, ensures  
470 the enrichment of microbial fractions and improves the quality of downstream analysis.  
471 Ideally, preprocessing should involve minimal, rapid, and reproducible steps. Since no  
472 standardized protocols for metaproteomics (or metagenomics) currently exist, methods  
473 must be tailored to the specific sample type and evaluated based on the study's objectives  
474 (Tanca et al. 2015; Salvato, Hettich, and Kleiner 2021; Pettersen et al. 2022). While the  
475 breadth of samples processed for metaproteomics remains limited, this field is rapidly  
476 evolving, and many more methods are expected to emerge.

477 For soil samples, humic substances derived from decomposed organic material often co-  
478 extract with proteins, interfering with MS measurements (Benndorf et al. 2007; Waibel et al.  
479 2023). To address this, several methods have been developed to remove humic  
480 compounds while preserving protein integrity before digestion (Keiblinger et al. 2012;  
481 Giagnoni et al. 2011; Chourey et al. 2010; Bastida, Hernández, and García 2014).  
482 Alternatively, filter-aided sample preparation (FASP) can directly digest proteins within  
483 humic complexes. This method uses acidification to precipitate humic compounds and

484 undigested proteins while peptides are extracted via centrifugation through molecular  
485 weight cut-off filters (Qian and Hettich 2017).

486 For human gut microbiome samples, non-microbial proteins from host cells and food debris  
487 are often much more abundant than microbial proteins, reducing the efficiency of microbial  
488 metaproteome identification (X. Zhang and Figeys 2019). Techniques such as double  
489 filtering (Xiong et al. 2015) and differential centrifugation (Tanca et al. 2014) can enrich  
490 microbial cells to improve identification. However, these methods may introduce biases and  
491 depend on the study's goals (Tanca et al. 2015). For example, double filtration can remove  
492 host cells and exoproteins, while differential centrifugation may non-specifically remove  
493 microbial cells and proteins (Speda et al. 2017; Armengaud et al. 2012; A. Wang et al.  
494 2024). Moreover, these methods are time-consuming and may be influenced by fecal  
495 variability, such as texture, fiber, and water content. Automation technologies, including  
496 solid-phase extraction clean-up, have been proposed to streamline processing for large  
497 longitudinal studies, reducing variability and improving reliability (Gonzalez et al. 2020).

498 In studies analyzing heterogeneous samples with high host protein content, such as viscous  
499 sputum of cystic fibrosis patients, certain plant tissues or environmental samples, a  
500 homogenization step can improve sample consistency. This step should be performed  
501 under conditions (temperature and duration) that minimize alterations to the *in vivo*  
502 metaproteome. Various mechanical strategies can achieve homogenization, including  
503 laboratory mills (Graf et al. 2021) and glass homogenizers (Salvato et al. 2022). The  
504 addition of protease inhibitors and DNase I to prevent protein degradation and disrupt DNA-  
505 based aggregates may also be beneficial, yet should be carefully evaluated based on the  
506 sample type and study objective.

507 For clinical samples containing bacterial or viral pathogens, inactivation is required before  
508 further processing outside appropriate biosafety level (BSL) containment. Since no  
509 standardized pipeline exists for this step, protocols must be tailored to the specific pathogen  
510 and sample type. Methods such as heat inactivation in lithium dodecyl sulfate buffers  
511 (Grenga et al. 2022) and MPLEx extraction, which uses chloroform, methanol, and water  
512 (8:4:3) for simultaneous pathogen inactivation and fractionation into metabolite, protein, and  
513 lipid phases, are commonly used (Burnum-Johnson et al. 2017). These approaches ensure  
514 both safety and compatibility with downstream metaproteomics workflows.

## 515 3.4 Protein sample preparation: from extraction to digestion

516 Preparing protein samples from biological material involves a series of interconnected steps,  
517 each essential for obtaining high-quality metaproteomic data. The term "protein extraction"  
518 is often used broadly to describe the entire workflow of isolating proteins from a biological  
519 sample. This process typically begins with cell lysis using extraction buffers and may also  
520 include subsequent protein clean-up steps, such as precipitation, filtration, or other methods.  
521 In some workflows, however, protein clean-up is treated as a distinct step, especially in  
522 protocols where extraction, clean-up, and digestion are streamlined into a single process.  
523 This section provides an overview of the key stages in protein sample preparation: cell lysis  
524 and extraction (**Section 3.4.1**), protein clean-up (**Section 3.4.2**), protein concentration  
525 (**Section 3.4.3**), and protein digestion (**Section 3.4.4**).

### 526 3.4.1 Cell lysis and protein extraction

527 Cell lysis releases the proteome from microbial cells, with a variety of methods available,  
528 each with distinct advantages (Hansmeier, Sharma, and Chao 2022). Mechanical disruption  
529 methods, such as direct ultrasonication, non-contact ultrasonication, and bead beating, are  
530 commonly used. Ultrasonication usually involves direct ultrasonication, where the probe is  
531 directly inserted into the sample, or non-contact ultrasonication, where the sample in a tube  
532 receives sonication energy from a cup horn through a coupling fluid. An advanced non-  
533 contact method termed Adaptive Focused Acoustic (AFA) technique provides precise  
534 control over parameters like amplitude and duration, achieving efficient lysis while  
535 minimizing protein denaturation (Dhabaria et al. 2015). Bead beating, which uses zirconia  
536 or silica beads, is effective for cell disruption, with bead size modulating efficiency (X. Zhang  
537 et al. 2018).

538 Chemical lysis methods use detergents such as urea buffers containing Triton X-100 or  
539 sodium dodecyl sulfate (SDS) to disrupt microbial cell membranes, often in combination  
540 with mechanical disruption/ultrasonication (X. Zhang et al. 2018). Notably, when combining  
541 urea-containing buffers with mechanical disruption or ultrasonication, one should be aware  
542 of the risk of urea-induced carbamylation caused by sample overheating (Kollipara and  
543 Zahedi 2013). Physical methods, including freeze-thaw cycles or high-pressure  
544 homogenization, are also effective, with pressure settings tailored to specific sample types  
545 (Cai et al. 2022). Since microbial cell structures vary significantly, for example between  
546 Gram-positive bacteria, Gram-negative bacteria, and fungi, optimizing lysis conditions is

547 crucial to preserve protein integrity, maximize yield, and ensure unbiased protein extraction  
548 (Starke et al. 2019; J. Wang et al. 2020).

549 Recently, some of the above approaches have been compared and found that a urea- and  
550 SDS-containing lysis buffer coupled to ultrasonication yielded higher protein recovery than  
551 bead beating in microbiome samples, with minimal sample loss, though both methods  
552 achieved similar peptide and protein identifications (X. Zhang et al. 2018). Careful selection  
553 of lysis buffers is also critical to avoid interference with downstream MS analysis. For  
554 example, ion suppression-inducing detergents like Tween-20 should be avoided unless  
555 they are removed during cleanup, as in methods like suspension trapping (S-trap) or FASP.

556 **Table 1** compares commonly used protein sample preparation methods, summarizing their  
557 key advantages and disadvantages. The choice of lysis method depends on factors such  
558 as sample type, desired protein yield, and sensitivity of proteins to denaturation or  
559 degradation. The listed lysis methods can also be combined, for example, detergent-  
560 containing urea lysis buffers are often coupled with ultrasonication to achieve fast and  
561 unbiased bacterial cell lysis in complex microbiome samples.

562 **Table 1. Comparison of standard protein sample preparation methods.** This table summarizes  
563 commonly used protein sample preparation techniques, outlining their key advantages and potential  
564 disadvantages.

Method	Description	Advantages	Disadvantages
Chemical Lysis	Disrupts cell membranes with chemicals like urea or guanidine hydrochloride.	Can unfold complex proteins.	If not removed or sufficiently diluted, it can interfere with protease activity. Risk of urea-induced carbamylation.
Detergent Lysis	Uses detergents (e.g., SDS, Triton X-100) to solubilize cell membranes.	Mild, preserves protein function, ideal for membrane proteins.	If a detergent is not removed or sufficiently diluted, it can interfere with protease activity.
Freeze-Thaw	Repeatedly freezes	Simple, no special	Time-consuming,

Cycles	and thaws the sample to rupture cell membranes.	equipment needed.	may not fully lyse cells, risk of protein degradation.
Bead beating	Physical force such as using bead beating to break cell walls.	Effective for bacterial cell lysis.	Requires specific instrument, sample loss due to contact with beads, can generate heat, risk of protein degradation.
Ultrasonication	Uses ultrasound waves to break cell membranes/walls and release proteins.	Fast, effective and can be non-contact for small samples, no need for harsh chemicals.	Can denature proteins if overused, heat generation requires sample cooling.

565

### 566 3.4.2 Protein clean-up: precipitation and alternative methods

567 Protein precipitation addresses the challenges of complex environmental and fecal samples  
568 by removing contaminants such as lipids, nucleic acids, and polysaccharides that can  
569 interfere with downstream MS analysis. Following microbial cell lysis, effective separation  
570 of proteins from cellular debris and contaminants is essential to ensure high protein yield  
571 and purity. Removing contaminants not only improves protein recovery but also enhances  
572 MS sensitivity, enabling more accurate and reliable protein identification.

573 The trichloroacetic acid (TCA)/acetone precipitation method is widely employed for this  
574 purpose. This method involves adding cold (-20°C) TCA or acetone, or both, to the protein  
575 lysate to precipitate proteins, followed by centrifugation to pellet the proteins. The pellets  
576 are then washed with cold acetone (-20°C) to remove residual contaminants and insoluble  
577 particles (Nickerson and Doucette 2020). This approach has proven effective for high-yield  
578 protein precipitation in diverse sample types, including marine sediment and forest soil  
579 samples, which contain complex organic matrices (Niu et al. 2018). Similarly, acidified  
580 acetone/ethanol buffer has also been used in metaproteomics (X. Zhang et al. 2016).

581 An alternative method, phenol extraction, separates proteins into the organic phase while  
582 partitioning nucleic acids into the aqueous phase. This approach is particularly beneficial  
583 for "dirty" samples, such as soil and wastewater sludge, which are rich in organic and  
584 inorganic contaminants. Phenol extraction can reduce the interference caused by  
585 contaminants, thus improving the downstream analysis of target proteins (Benndorf et al.  
586 2009). Phenol extraction also enables the simultaneous extraction of nucleic acids from the  
587 same sample, making it highly suitable for integrated omics studies, especially in  
588 microbiome research (Baldrian 2017).

589 For samples with low microbial load, such as fecal samples, river sediment, or air filters,  
590 maximizing protein recovery is critical. Organic solvent systems, such as  
591 chloroform/methanol or chloroform/methanol/water mixtures, have proven effective for  
592 enhancing protein recovery and minimizing the loss of low-abundance proteins by  
593 optimizing solvent ratios and conditions (Vertommen et al. 2010). Biphasic systems, such  
594 as phenol/chloroform or Triton X-114, can also be used to selectively partition proteins and  
595 facilitate the removal of contaminants (Wessel and Flügge 1984).

596 Traditional protein precipitation methods, while effective, can be labor-intensive and may  
597 not always completely eliminate contaminants that interfere with downstream analyses. To  
598 address these limitations, alternative methods have been developed to improve protein  
599 clean-up and digestion efficiency. Techniques such as FASP, SP3, and suspension  
600 trapping (S-Trap) have shown promise for processing challenging samples like human fecal  
601 protein extracts (Tanca et al. 2024). Solid-phase alkylation, a novel strategy designed for  
602 low-loss and anti-interference sample preparation, utilizing covalent binding and purification  
603 of proteins, has also been proved effective for marine microbiome samples (S. Wang et al.  
604 2024). These approaches integrate clean-up and digestion steps into a single workflow,  
605 facilitating high-throughput applications.

### 606 3.4.3 Measuring protein concentration

607 Accurate protein concentration measurement ensures uniform loading in downstream LC-  
608 MS/MS analyses and for facilitating reliable data interpretation (Sapan and Lundblad 2015).  
609 Consistent peptide loading in LC-MS/MS is essential for accurate peptide quantification, as  
610 it maintains signal intensity and ensures reliable peptide detection across samples. Uniform  
611 loading also optimizes column performance, reducing variability in peak shapes and  
612 retention times. This consistency minimizes technical artifacts, enabling clearer biological  
613 insights when comparing samples.

614 Various methods are commonly used to determine protein concentration. The Bradford  
615 Assay, which utilizes Coomassie Brilliant Blue dye, measures protein concentration through  
616 a colorimetric change, requiring a standard curve prepared with known protein  
617 concentrations to ensure precision. The bicinchoninic acid (BCA) assay forms a purple-  
618 colored complex for protein quantification, with sensitivity optimized by adjusting reagent  
619 ratios and incubation conditions. Fluorescence-based assays, such as the Qubit Protein  
620 Assay, use dye-binding technology for highly sensitive quantification with minimal  
621 interference, making them suitable for samples with low protein concentrations.

622 The 2-D Quant Kit is another option, which quantitatively precipitates proteins while leaving  
623 interfering substances in solution. This method produces a color density inversely related  
624 to protein concentration, with a linear response in the range of 0–50 µg and a volume range  
625 of 1–50 µL. When selecting a protein concentration method, it is important to consider the  
626 required sensitivity, dynamic range, and compatibility with buffer components, as some  
627 assays show varying tolerance to substances like SDS or protease inhibitors, including  
628 PMSF.

629 If no suitable quantification assay is available, running SDS-PAGE gels can provide a rough  
630 estimate of protein abundance. While less precise, this approach can offer a practical  
631 alternative for assessing protein concentrations in certain scenarios.

632 This systematic approach to protein concentration measurement ensures consistency and  
633 reliability in downstream analyses, particularly when dealing with complex microbial  
634 samples containing proteins spanning a wide range of abundances.

#### 635 3.4.4 Protein digestion

636 Bottom-up (shotgun) metaproteomic studies involve the enzymatic digestion of proteins into  
637 peptides, a process known as proteolysis, for untargeted protein identification. This method  
638 requires several preparatory steps to ensure efficient proteolysis. Initially, proteins are  
639 denatured using agents such as urea or guanidine hydrochloride to expose cleavage sites.  
640 Disulfide bonds are then reduced using reducing agents like dithiothreitol (DTT) or tris(2-  
641 carboxyethyl)phosphine (TCEP). To prevent the re-formation of disulfide bonds, cysteine  
642 residues are alkylated with agents like iodoacetamide, which react with sulfhydryl groups to  
643 form stable thioether adducts (Sechi and Chait 1998). This alkylation introduces mass  
644 changes that must be accounted for during peptide identification, as discussed in **Section**  
645 **4.1.1**.

646 Following these preparatory steps, proteins are enzymatically cleaved into peptides suitable  
647 for downstream LC-MS/MS analysis (Hustoft et al. 2012). The most commonly used  
648 protease is trypsin due to its high specificity and efficiency. It cleaves proteins at the C-  
649 terminal side of lysine and arginine residues, producing peptides ideal for shotgun MS  
650 analysis. Lys-C, another commonly used protease, complements trypsin digestion by  
651 cleaving at the C-terminal side of lysine residues, particularly in high urea concentrations  
652 (8 M), enhancing peptide coverage. Alternative proteases such as chymotrypsin, Glu-C,  
653 and Asp-N may also be used to increase peptide diversity or for specific applications.  
654 However, the combination of trypsin and Lys-C is often the most practical and widely  
655 applied choice.

656 The enzyme-to-substrate ratio is another important factor, with typical ratios ranging from  
657 1:50 to 1:100 (w/w). Digestion time is also critical and usually involves incubating the  
658 proteome mixture at an appropriate temperature (e.g., 37°C) for several hours to overnight,  
659 depending on sample complexity and enzyme properties. Digestion is quenched by  
660 acidification, commonly using formic acid or trifluoroacetic acid to achieve a pH of 2–3. In  
661 methods such as S-trap or FASP, peptides may also be eluted without an acidification step.

662 Peptide lysates are subsequently desalted or purified to remove salts and contaminants.  
663 Solid-phase extraction (SPE), C18 ZipTips (Millipore), or ultrafiltration are commonly used  
664 for this purpose. In some cases, the desalting step can be omitted if peptides are desalted  
665 on a trap column in the LC system.

666 Direct in-solution protein digestion methods have been developed to streamline the  
667 workflow, offering efficient and high-throughput options. Notable examples include SP3  
668 (Hughes et al. 2014), FASP (Wiśniewski et al. 2009), S-trap (HaileMariam et al. 2018) and  
669 a commercial kit based on the in-StageTip (iST) (Kulak et al. 2014). These methods are  
670 designed to ensure high protein recovery and compatibility with downstream MS analysis,  
671 even when working with low protein amounts.

## 672 3.5 Separation and fractionation techniques

673 Separation and fractionation enable researchers to reduce sample complexity and enhance  
674 the depth and sensitivity of protein identification and quantification. These processes can  
675 be performed at multiple levels, including the peptide, protein, and cellular stages,  
676 depending on the specific goals of the analysis (Cheng et al. 2018). Techniques such as  
677 peptide fractionation are frequently used to enhance LC-MS/MS performance (**Section**  
678 **3.5.1**), while enrichment approaches allow for the targeted analysis of PTMs (**Section 3.5.2**).



679 At the protein or cellular level, fractionation strategies can further refine sample complexity  
680 or enrich specific components of interest (**Section 3.5.3**).

### 681 3.5.1 On-line and off-line peptide fractionation

682 Peptide separation workflows can generally be categorized into one-dimensional (1D) and  
683 two-dimensional (2D) or multi-dimensional approaches. In 1D liquid chromatography (LC),  
684 which is widely used in metaproteomics, reverse-phase (RP) nano-high-performance liquid  
685 chromatography (nanoHPLC, mostly just abbreviated as LC or HPLC) employs C18  
686 columns to separate peptides based on their hydrophobicity and is coupled directly with  
687 mass spectrometry for peptide analysis. 2D-LC, often based on multidimensional protein  
688 identification technology (MudPIT) (Washburn, Wolters, and Yates 2001), combines strong  
689 cation exchange (SCX) with RP-HPLC. Peptides are first fractionated on the SCX column  
690 based on their charge using salt or pH gradients for elution, and then further separated  
691 based on hydrophobicity on an RP-HPLC column (Verberkmoes et al. 2009). The 2D-LC  
692 strategy has been applied in metaproteomic analyses to improve identification depth, with  
693 online 2D LC-MS setups used for shotgun proteomics in studies of human gut and  
694 environmental microbiomes (Verberkmoes et al. 2009).

695 Off-line pre-fractionation, although less commonly used in metaproteomics due to its labor-  
696 intensive nature and the increased MS time required, offers potential for deeper peptide  
697 and protein identification (X. Zhang et al. 2017). High-pH RP chromatography is one such  
698 method and is orthogonal to low-pH RP-LC-MS gradients. This fractionation can be  
699 achieved using either stage-tip methods or HPLC systems. Stage-tip-based fractionation is  
700 straightforward to implement and is supported by commercially available kits (e.g., Pierce™  
701 High pH Reversed-Phase Peptide Fractionation Kit). On the other hand, micro-flow HPLC  
702 systems enable higher-resolution fractionation through continuous collection of numerous  
703 fractions and stepwise concatenation.

704 While extensive fractionation can significantly enhance the depth of metaproteomic analysis,  
705 it also increases costs, sample requirements, and instrument time, making it less feasible  
706 for large cohort studies. The adoption of multiplexing techniques, such as tandem mass  
707 tags (TMT) (Creskey et al. 2023), has mitigated these limitations by reducing MS time and  
708 the required sample quantity per condition. The combination of off-line peptide fractionation  
709 and multiplexing presents a promising and accessible option for researchers, particularly  
710 beginners, aiming to conduct in-depth metaproteomic analyses to investigate microbiome  
711 functionality.

### 712 3.5.2 Enrichment of peptides with post-translational modifications

713 PTMs are critical regulators of protein activity and function, and their study is uniquely  
714 possible through metaproteomics. Unlike other omics approaches, metaproteomics  
715 provides the direct capability to identify and quantify PTMs in microbial proteins, offering  
716 unparalleled insights into microbiome functionality. While analyzing PTMs at the  
717 metaproteome level is particularly challenging, several studies have successfully performed  
718 metaPTMomics on environmental and human gut microbiomes (Z. Li et al. 2014; W. Zhang  
719 et al. 2016; X. Zhang et al. 2021; 2020). These studies identified various PTMs, including  
720 methylation, hydroxylation, acylations, citrullination, deamination, phosphorylation, and  
721 nitrosylation, among others, with abundances varying across different microbiome types.  
722 Understanding the diversity and distribution of PTMs is essential for uncovering microbiome  
723 functionality. Recent advancements in the field have been detailed in two comprehensive  
724 reviews (Duchovni, Shmunis, and Lobel 2024; Duan, Zhang, and Figeys 2023).

725 Microbiome PTMs can be analyzed using non-enriched samples combined with tailored  
726 bioinformatics workflows (Z. Li et al. 2014; W. Zhang et al. 2016) or quantitatively profiled  
727 using enrichment techniques at the peptide or protein level (X. Zhang et al. 2021; 2020).  
728 Depending on the type of PTM, specific enrichment strategies may be employed to facilitate  
729 detection during MS analysis.

730 Immuno-affinity enrichment is widely used for protein acylations, such as lysine acetylation,  
731 propionylation, and succinylation, and has recently been applied to human gut microbiomes  
732 (X. Zhang et al. 2021). This technique uses antibodies bound to agarose or magnetic beads  
733 to selectively enrich acylated peptides, improving MS sensitivity and specificity. However,  
734 this approach can be limited by the availability of motif-specific antibodies and the inability  
735 to capture the full spectrum of modified peptides.

736 Immobilized metal affinity chromatography (IMAC) is a commonly used strategy in  
737 proteomics to enrich phosphorylated peptides for phosphoproteomic studies. Ti-IMAC and  
738 Fe-IMAC are typical examples, offering robust enrichment prior to LC-MS/MS analysis (Low  
739 et al. 2021).

740 Hydrophilic interaction liquid chromatography (HILIC) is another effective technique,  
741 particularly for enriching glycopeptides. This method capitalizes on its high selectivity and  
742 specificity for hydrophilic glycan moieties (Mysling et al. 2010). These enrichment  
743 approaches have been extensively applied to mammalian cells, tissues, and single bacterial  
744 strains, and they show potential for broader applications in microbiome studies.

### 745 3.5.3 Protein, cell-level and functional fractionation techniques

746 The high complexity of microbiomes often necessitates cellular and protein-level  
747 separations to complement peptide-level fractionation, enhancing the depth and resolution  
748 of metaproteomic analysis. Although high-speed, high-resolution mass spectrometers have  
749 made peptide fractionation sufficient for many proteomics workflows, the added complexity  
750 of microbiomes can still benefit from upstream fractionation approaches.

751 Capillary zone electrophoresis (CZE), a technique used to separate charged particles,  
752 shows promise for separating intact proteins and even bacterial cells (Cheng et al. 2018).  
753 Another method for separating proteomes from different bacteria is differential lysis, which,  
754 despite its relatively low granularity, can distinguish between bacterial types based on cell  
755 wall structure (J. Wang et al. 2020). In this approach, sequential lysis is achieved using  
756 buffers of increasing strength, such as those containing urea or varying concentrations of  
757 SDS. This method can separate the proteomes of Gram-negative bacteria, which have  
758 thinner cell walls, from those of Gram-positive bacteria with thicker, multilayered cell walls  
759 (J. Wang et al. 2020).

760 For host-associated microbiomes, removing abundant host cells is often critical to  
761 improving microbial signal detection. Techniques such as differential centrifugation and  
762 density gradient centrifugation (Hinzke, Kleiner, and Markert 2018) are commonly used to  
763 enrich microbial cells. Following lysis, additional separation of cellular components can be  
764 achieved through methods like ultracentrifugation (Henry et al. 2022), further increasing  
765 protein identification coverage.

766 Functional fractionation techniques, such as Activity-Based Protein Probing (ABPP), can  
767 be used to study enzymatic functions at the proteome level (Cravatt, Wright, and Kozarich  
768 2008). ABPP employs small-molecule probes that covalently bind to active sites of proteins  
769 with specific functions or residues. These labeled proteins can then be captured or enriched  
770 for LC-MS/MS analysis, enabling detailed profiling of protein functions and aiding in drug  
771 target discovery. ABPP is particularly useful for annotating proteins with unknown functions  
772 (Barglow and Cravatt 2007), making it a relevant approach in microbiome studies. Recent  
773 applications of ABPP in both host-associated and environmental microbiomes have  
774 uncovered diverse microbial enzymes, including thiol-containing proteases, bile salt  
775 hydrolases (BSHs), glycoside hydrolases (GHs), and  $\beta$ -glucuronidases (Han and Chang  
776 2023).

## 777 3.6 Automation

778 High-throughput techniques have transformed sample preparation, simplifying labor-  
779 intensive steps and revolutionizing workflows in proteomics, especially as datasets continue  
780 to grow in scale and complexity (Fu et al. 2023; Burns et al. 2021). These advancements  
781 have facilitated applications such as chemical proteomics (Lin et al. 2023), biomarker  
782 detection (Paramasivan et al. 2023) and drug target discovery (Qiong Wu et al. 2024).  
783 Although automation in metaproteomics has not advanced as rapidly as in proteomics, its  
784 potential for transforming the field is immense.

785 Automating metaproteomics workflows offers multiple benefits, including reduced sample  
786 handling time, minimized operator-induced variability, and enhanced reproducibility. These  
787 improvements provide broader coverage of microbiome responses to environmental factors  
788 within limited experimental timeframes. Furthermore, high-throughput automated workflows  
789 allow researchers to scale up the discovery of microbiome-associated biomarkers and  
790 explore dynamic functional landscapes across diverse microbiomes. Automation also  
791 generates large datasets, enabling the application of artificial intelligence (AI) to uncover  
792 hidden patterns within metaproteomic profiles.

793 Automated sample processing in metaproteomics can be broadly divided into four key steps:  
794 microbial cell disruption and protein extraction (**Section 3.6.1**), protein digestion and  
795 peptide clean-up (**Section 3.6.2**), and multiplexing (**Section 3.6.3**).

### 796 3.6.1 Microbial cell disruption and protein extraction

797 In certain scenarios, such as working with complex clinical samples like human stool or  
798 saliva, microbial cell enrichment is often required but poses significant challenges. Sample  
799 properties can vary greatly within a dataset, complicating efforts to standardize technical  
800 parameters for automated microbial cell purification. As a result, current automated  
801 metaproteomics workflows often exclude fully automated raw sample handling steps. For  
802 example, the RapidAIM 2.0 pipeline (L. Li et al. 2024) includes manual bacterial enrichment  
803 and cell washing, with a 96-channel liquid handler accelerating pipetting steps. In contrast,  
804 the SHT-Pro protocol (Gonzalez et al. 2020), the first high-throughput pipeline specifically  
805 designed for large-scale stool sample processing, begins with the lysis of raw stool samples  
806 without prior microbial enrichment. This approach is particularly beneficial when both host  
807 and microbial proteins are of interest.

808 Microbial cell disruption for protein extraction can be effectively automated in a 96-well  
809 format using ultra-sonication devices designed for high-throughput workflows. These  
810 instruments facilitate efficient protein extraction, enabling downstream high-throughput  
811 protein clean-up. Several methods, including FASP, SP3, and S-Trap, have been  
812 successfully adapted to microplate-based formats, with studies showing that the  
813 combination of FASP and SP3 with iST yields the most robust results for high-throughput  
814 protein processing (Tanca et al. 2024).

### 815 3.6.2 Protein digestion and peptide clean-up

816 Similar to manual metaproteomics workflows, automated protein preparation typically  
817 involves protein denaturation, reduction, alkylation, and protease digestion. These steps  
818 are relatively straightforward to automate and can be performed using liquid handling  
819 platforms equipped with low-volume pipetting accuracy and heater-shaker capabilities.  
820 Therefore, protein digestion is often considered one of the least complex steps to automate  
821 in metaproteomic workflows.

822 Peptide clean-up, however, presents greater challenges. Typically, this step is carried out  
823 manually by skilled personnel using solid-phase extraction (SPE), C18 ZipTips, or  
824 ultrafiltration, as described in **Section 3.4.4**. During automation, sample heterogeneity at  
825 this stage can introduce variability, complicating experimental parameter control. A  
826 promising solution involves replacing centrifugation through reverse-phase columns with  
827 pipette-based mixing of reverse-phase resins. This approach has been incorporated into  
828 workflows like RapidAIM 2.0 (L. Li et al. 2024) and is supported by established proteomics  
829 automation protocols. For example, the autoSISPROT system offers all-in-tip sample  
830 preparation capabilities, demonstrating compatibility with automated platforms (Qiong Wu  
831 et al. 2024).

### 832 3.6.3 Multiplexing

833 The integration of automated sample handling with techniques like TMT labeling  
834 significantly enhances throughput and accelerates the discovery process in  
835 metaproteomics. However, the high cost of TMT reagents might be a challenge for broader  
836 application. One solution involves pre-aliquoting and drying TMT reagents in a 96-well plate  
837 format, a strategy that reduces reagent waste and preparation time. This approach is  
838 compatible with automated workflows, such as those used in the RapidAIM 2.0 platform,  
839 and facilitates more efficient reagent utilization (L. Li et al. 2024).

840 While advancements in automation have enabled notable progress in metaproteomics,  
841 most current systems are semi-automated rather than fully automated. Continued  
842 development of automation technologies is essential to further streamline workflows,  
843 enhance sample processing speed, and achieve higher throughput.

### 844 3.7 Mass spectrometry data acquisition methods

845 Mass spectrometry analysis of (meta)proteomes is predominantly carried out using (HP)LC-  
846 MS/MS. A fundamental limitation of mass spectrometers, even when combined with  
847 multidimensional separations, is their inability to generate fragmentation spectra (or MS/MS  
848 spectra) for all peptides in a sample within a single run. This constraint has led to the  
849 widespread adoption of data-dependent acquisition (DDA) as the dominant approach in  
850 proteomics over the past 25 years.

851 DDA, as discussed in **Section 3.7.1**, involves selecting the most abundant precursor ions  
852 from the MS1 spectra for fragmentation in the MS2 (or MS/MS) stage, dynamically  
853 excluding previously fragmented ions to prioritize unfragmented targets. This strategy  
854 increases the diversity of identified peptides and proteins. In metaproteomics, however, the  
855 complexity of the samples presents significant challenges for DDA, particularly in achieving  
856 comprehensive sequencing depth and coverage. Even with the latest high-resolution and  
857 highly sensitive mass spectrometers, DDA is inherently biased toward the most abundant  
858 ions, leaving many lower-abundance peptides uncharacterized. Nevertheless, DDA  
859 remains the most widely used method due to its extensive validation, established workflows,  
860 and compatibility with a broad range of analytical tools.

861 Data-independent acquisition (DIA), as discussed in **Section 3.7.2**, is a more recent  
862 advancement that offers an alternative approach by fragmenting all peptide ions within  
863 predefined mass-to-charge ( $m/z$ ) windows, rather than selectively targeting the most  
864 abundant ones. DIA addresses some of the limitations of DDA, particularly in terms of  
865 peptide coverage and reproducibility, making it increasingly attractive for metaproteomics.  
866 However, the broader data capture in DIA results in significantly more complex datasets  
867 that require advanced computational tools for processing and analysis. While progress has  
868 been made in developing such tools, further validation and optimization are needed before  
869 DIA can become a routine method for metaproteomics.

870 Both DDA and DIA have distinct advantages and limitations, and their choice depends on  
871 the specific goals of the experiment, the complexity of the sample, and the available  
872 computational resources.

### 873 3.7.1 DDA

874 DDA is the most widely used method in proteomics, particularly in shotgun proteomics, for  
875 identifying peptides in biological samples. In DDA mode, the mass spectrometer  
876 dynamically selects a specified number of the most abundant precursor ions (commonly  
877 referred to as the "topN") for fragmentation. This prioritization ensures that the most intense  
878 ions within each acquisition cycle are fragmented into smaller ions, generating MS/MS  
879 spectra that serve as unique fingerprints for peptide identification. To enhance the detection  
880 of lower-abundance peptides, DDA incorporates a process known as dynamic exclusion.  
881 Previously selected precursor ions are temporarily excluded from subsequent  
882 fragmentation, increasing the diversity of peptides analyzed within a single run. These  
883 MS/MS spectra are then analyzed using proteomics software packages (**Section 4.1.1**).

884 DDA has several advantages, making it a popular choice for metaproteomics workflows. It  
885 is relatively simple to configure and analyze compared to more complex approaches like  
886 DIA, making it accessible for both beginners and experienced researchers. The one-to-one  
887 relationship between spectra and peptides reduces computational demands during data  
888 analysis, particularly when a well-curated protein database is available. More information  
889 on creating a protein database is provided in **Section 4.1.2**. Furthermore, DDA supports  
890 relative quantification of proteins using both label-free quantification (LFQ) and labeling  
891 approaches, offering flexibility for various experimental designs (**Section 4.1.5**). Its  
892 longstanding use in proteomics has also led to the development of numerous software tools  
893 and well-established workflows, enhancing its reliability and versatility.

894 Despite its strengths, DDA has notable limitations. Its reliance on selecting the most intense  
895 precursor ions means that low-abundance proteins may go undetected, especially in  
896 complex samples. Additionally, DDA often fails to identify the same peptides consistently  
897 across multiple runs, resulting in missing values for low-abundance proteins and  
898 complicating large-scale quantitative studies.

899 Overall, while DDA is not without its limitations, it remains the most widely used and  
900 versatile technique in metaproteomics (Van Den Bossche, Kunath, et al. 2021). For studies  
901 requiring deeper proteome coverage or greater reproducibility, alternative methods like DIA  
902 may offer complementary advantages.

### 903 3.7.2 DIA

904 DIA mass spectrometry has emerged as a powerful approach in proteomics, providing  
905 broad protein coverage, high reproducibility, and quantitative accuracy. Unlike DDA, which  
906 focuses on fragmenting a limited number of the most intense precursor ions, DIA fragments  
907 all ions within predefined  $m/z$  windows. These windows are repeatedly scanned across the  
908 entire  $m/z$  range, generating complex MS/MS spectra that provide a more comprehensive  
909 view of the proteome. This inclusivity is particularly advantageous in metaproteomics,  
910 where samples contain an overwhelming diversity of peptides and low-abundance proteins  
911 that might be missed by DDA (E. Wu et al. 2024).

912 DIA has demonstrated significant potential in metaproteomics applications. Its application  
913 in metaproteomics was first evaluated in gut microbiome studies (Aakko et al. 2020) and  
914 has since expanded to various contexts, including Chinese liquor fermenter starters (Zhao,  
915 Yang, Chen, et al. 2023), and multicenter diagnostic research on tongue coating samples  
916 for gastric cancer (Chen et al. 2024). Recent advances in MS instrumentation, such as DIA-  
917 PASEF (Gómez-Varela et al. 2023) and the Orbitrap Astral (Dumas et al. 2024), have  
918 significantly improved DIA's sensitivity and resolution, enabling deeper proteome coverage  
919 in highly complex microbial communities.

920 One of DIA's key advantages lies in its ability to capture a broader range of peptides  
921 compared to DDA, enabling deeper proteome coverage and improved detection of low-  
922 abundance proteins (Chen et al. 2024; Gómez-Varela et al. 2023; Pietilä, Suomi, and Elo  
923 2022; Zhao, Yang, Xu, et al. 2023; Zhao, Yang, Chen, et al. 2023; Aakko et al. 2020; Zhao,  
924 Yang, Teng, et al. 2023). Another significant advantage is its reproducibility across samples,  
925 as it is less susceptible to variations in ionization efficiency (Fernández-Costa et al. 2020).  
926 This consistency makes DIA particularly well-suited for large-scale quantitative studies.

927 Despite its advantages, DIA also comes with challenges, particularly in data analysis.  
928 Indeed, analyzing the complex MS/MS spectra generated by DIA requires advanced  
929 computational tools and specialized expertise which is further discussed in **Section 4.1.1**.  
930 Additionally, because DIA fragments all ions within a given  $m/z$  window simultaneously, the  
931 resulting spectra are more complex and less specific to individual peptides compared to  
932 DDA. This reduced specificity can make it challenging to confidently resolve detailed  
933 structural or sequence-level information for single peptides, limiting DIA's utility for  
934 applications that require precise characterization, such as studying PTMs or differentiating  
935 highly similar peptide sequences. These inherent trade-offs highlight the importance of  
936 carefully tailoring DIA workflows to specific research objectives.



937 Nevertheless, DIA's rapid advancements make it a promising tool for metaproteomics,  
938 providing the depth and reproducibility required to explore the functional landscape of  
939 microbial communities comprehensively.

### 940 3.7.3 Critical parameters to optimize the HPLC and MS methods

941 Optimization of HPLC and MS methods is crucial for obtaining high-quality data in  
942 metaproteomics workflows. Each parameter below plays a significant role in ensuring  
943 accurate peptide separation, identification, and quantification. Metaproteomics, with its  
944 added complexity compared to standard proteomics workflows, requires specific  
945 adjustments to many of these parameters.

#### 946 **i) Analytical column quality, gradient and flow rates**

947 Peptides are commonly separated using HPLC, which is directly coupled to the MS, using  
948 either commercial or in-house analytical HPLC columns. These separations are achieved  
949 with a mobile phase composed of increasing concentrations of acetonitrile (ACN). For  
950 laboratories using in-house columns, stringent QC checks are crucial to ensure consistent  
951 column performance, as explained in **Section 3.7.4**.

952 Metaproteomics samples present significantly greater chromatographic challenges than  
953 single-species proteomics due to their inherent complexity (Duan et al. 2022). To address  
954 this, typical mobile phase gradients of 5–35% of 80% ACN or 5–30% of 100% ACN over  
955 1–2 hours are generally sufficient for tryptic peptide elution. However, adjustments may be  
956 required for specific experimental setups. For example, chemically labeled digests with  
957 increased hydrophobicity often require a steeper gradient with a higher final concentration  
958 of ACN for complete peptide elution.

959 Efficient gradient design is essential to optimize runtime and achieve an even distribution  
960 of peptide elution across the gradient. Since fewer peptides elute at the beginning and end  
961 of the gradient, tailoring the gradient can improve separation and detection (Xu, Duong, and  
962 Peng 2009). Accurate peptide quantification requires sufficient sampling points per LC peak,  
963 making short gradients (e.g., 10-minute gradients) generally unsuitable for metaproteomics  
964 in data-dependent acquisition (DDA) mode. Comprehensive tutorials on gradient  
965 optimization are available for general proteomics (Lenčo et al. 2022), and metaproteomics  
966 specifically (Hinze et al. 2019).

967 LC flow rates typically range from 200–300 nL/min. Recently, higher flow rates have gained  
968 popularity to accelerate sample duty cycles. However, these higher flow rates compromise  
969 sensitivity. Strategies to offset this limitation include increasing the sample loading amount

970 or using dimethyl sulfoxide (DMSO) to boost signal intensity, making higher flow rates more  
971 viable for metaproteomics workflows.

## 972 **ii) MS Settings in DDA workflows**

973 Optimizing MS parameters plays a key role in obtaining high-quality data in metaproteomics.  
974 While those new to the field are generally not expected to configure MS settings,  
975 understanding key optimization steps can provide valuable context for interpreting data and  
976 troubleshooting issues.

977 Accurate mass measurements require regular calibration of the mass spectrometer, which  
978 is crucial for reliable peptide identification and quantification. Additionally, source  
979 parameters such as source temperature, flow rates, and nebulizer gas pressure must be  
980 optimized to enhance ionization efficiency and maximize signal intensity. The specific  
981 optimization steps vary depending on the type of mass analyzer used, such as time-of-flight  
982 (TOF) or Orbitrap instruments. Key parameters for these analyzers include scan range,  
983 resolution, and scan speed, which must be fine-tuned to ensure precise mass  
984 measurements and resolve closely spaced peptide ions. Similarly, collision energy settings  
985 for peptide fragmentation need careful adjustment to generate high-quality fragment  
986 spectra for peptide identification.

987 Dynamic exclusion is a critical parameter in DDA workflows, requiring careful calibration to  
988 align with the chromatographic gradient and peak width. This setting prevents repeated  
989 fragmentation of the same peptide by excluding it temporarily after its initial fragmentation,  
990 thereby increasing peptide diversity. However, this approach poses challenges, particularly  
991 in metaproteomics. Many researchers rely on spectral counting for relative quantification,  
992 as it has been shown robust for metaproteomic datasets with significant differences in cell  
993 numbers and total protein amounts between community members (Kleiner et al. 2017).  
994 Nonetheless, dynamic exclusion can limit the number of spectra acquired for abundant  
995 peptides, leading to fewer spectral counts than expected and potentially skewing  
996 quantification accuracy. This issue is exacerbated with modern high-resolution instruments,  
997 where the correlation between peptide abundance and peptide-spectrum matches (PSMs)  
998 becomes less relevant due to faster scan rates and increased resolving power. Dynamic  
999 exclusion times must therefore strike a balance, ensuring high-quality fragmentation  
1000 spectra while maximizing the diversity of peptides analyzed. The choice between spectral  
1001 counting and MS1-based quantification methods like area under the curve (AUC) remains  
1002 a topic of debate in metaproteomics.

1003 In DDA, selecting the isolation window width for precursor ions is a critical optimization step.  
1004 A wider isolation window, up to 2 Da, allows the collection of more ions, resulting in higher-  
1005 quality MS spectra. However, this increases the risk of generating chimeric spectra, where  
1006 fragments from multiple precursor ions are combined, complicating peptide identification.  
1007 Conversely, narrower isolation windows, down to 0.7 Da, reduce the likelihood of chimeric  
1008 spectra but limit the number of ions isolated, potentially impacting signal intensity. In  
1009 metaproteomics, the high density and diversity of precursor ions in certain mass ranges  
1010 complicates this balance, as even narrow windows can capture multiple ions. Advances in  
1011 mass spectrometers, such as faster scan speeds, now enable higher topN settings in DDA  
1012 workflows, helping to address this challenge by acquiring more fragmentation spectra within  
1013 a given run.

### 1014 **iii) MS Settings in DIA workflows**

1015 Optimizing data-independent acquisition (DIA) workflows requires careful calibration of  
1016 several key parameters to achieve accurate and comprehensive peptide identification. The  
1017 width of mass isolation windows is particularly critical, as narrower windows, such as 2 m/z,  
1018 provide higher resolution and more precise fragmentation spectra, which are essential for  
1019 resolving complex peptide mixtures. However, narrower windows can reduce proteome  
1020 coverage, as fewer ions are isolated in each cycle. Balancing resolution with proteome  
1021 coverage is thus a central challenge in DIA optimization. Recent advancements, such as  
1022 the Orbitrap Astral mass spectrometer, support exceptionally narrow isolation windows  
1023 while maintaining high scanning speeds, effectively bridging the gap between DDA and DIA  
1024 methodologies.

1025 In addition to tuning isolation windows, optimizing collision energy is required for generating  
1026 high-quality fragment ions, while chromatographic conditions, including gradient length and  
1027 flow rate, must be carefully calibrated to align with the DIA cycle time. Ensuring sufficient  
1028 acquisition points across peptide elution peaks is essential for accurate quantification and  
1029 peptide identification. DIA workflows in metaproteomics are advancing rapidly, providing  
1030 enhanced resolution and deeper proteome coverage in complex microbial samples (E. Wu  
1031 et al. 2024; Dumas et al. 2024). Detailed guidelines for these optimization strategies can  
1032 be found in recent studies exploring advancements in DIA methodologies (Ishikawa et al.  
1033 2022; Demichev et al. 2022; Gu et al. 2024).

### 1034 **3.7.4 Quality control of LC-MS/MS**

1035 A comprehensive QC workflow begins with a blank injection of solvent without any sample  
1036 to check for background contamination. Ideally, a blank run should produce minimal

1037 identifications, which can be verified visually or through database searches. Contamination  
1038 sources can include transport solvents used in HPLC systems, so these should be carefully  
1039 monitored. Next, a standard injection of a known peptide mixture, such as cytochrome C or  
1040 BSA digest, is performed to confirm instrument calibration and performance. Simple  
1041 mixtures like these are useful for testing HPLC performance, while more complex peptide  
1042 mixtures, such as HeLa digest, assess the mass spectrometer's ability to analyze complex  
1043 samples. A representative microbiome sample digest can also be injected to refine the LC  
1044 gradient profile, and such standards should be injected regularly throughout the run.  
1045 Additionally, using reference microbiome material as a positive control can help verify the  
1046 efficiency of protein extraction protocols. This ensures that the extraction method reliably  
1047 captures a representative set of proteins from the sample, which is particularly important  
1048 for metaproteomic studies. Database searches on complex standards should be used to  
1049 monitor metrics like number of PSMs, peptide and protein identifications. Consistently  
1050 tracking these values over time helps detect performance declines, signaling when the  
1051 instrument requires cleaning or recalibration.

1052 During the LC-MS/MS run, retention times for known peaks should be monitored closely,  
1053 as significant shifts compared to previous runs may indicate issues such as column  
1054 blockage, connector leakage, or valve wear. Similarly, column back pressure should be  
1055 monitored as a potential indicator of problems. Peak shape should also be evaluated for  
1056 symmetry and sharpness; tailing or broadening peaks may suggest problems with  
1057 chromatography or ionization efficiency. Signal intensity is another important parameter,  
1058 and any significant drop compared to expected values may point to reduced instrument  
1059 sensitivity or ionization issues.

1060 After the run, each raw file must be carefully reviewed to identify potential issues. Failed  
1061 runs should be rerun immediately to avoid batch effects caused by delayed reanalysis. The  
1062 total ion current (TIC) chromatogram provides valuable information on instrument  
1063 performance, and it should be examined for unexpected peaks or a noisy baseline, both of  
1064 which may point to contamination or hardware issues. The base peak chromatogram  
1065 provides additional insights into LC resolution. Comparing the TIC-to-base peak intensity  
1066 ratio is also informative, as higher values often reflect increased sample complexity or poor  
1067 chromatographic performance. Retention times and peak intensities across samples should  
1068 be consistent, indicating good repeatability. Additional QC checks, such as PCA clustering  
1069 or heatmaps, can help pinpoint variations between runs and ensure data quality.

1070 Metrics collected after protein identification and quantification are also essential for  
1071 evaluating QC (Bielow, Mastrobuoni, and Kempa 2016). For example, the number of

1072 identified PSMs to the total number of MS2 spectra, the PSM identification rate, serves as  
1073 a key indicator of data quality. Using a 1-hour gradient on a Q-Exactive mass spectrometer  
1074 with optimized conditions and high-quality sample preparation, metaproteomic samples can  
1075 achieve an ID rate of approximately 50%, meaning that 50% of spectra yield identified  
1076 peptide sequences after 1% FDR filtering. Note that for samples in less trivial environments,  
1077 such as soil, the PSM identification rate will be much lower. It is crucial to analyze high-  
1078 quality QC samples using the same LC-MS/MS methods, as the identification rate depends  
1079 heavily on both the instrument's performance and sample preparation.

1080 In large-scale projects lasting several weeks, retention time drift and signal drops are  
1081 common. Blocking and randomizing samples during analysis is recommended to reduce  
1082 systematic biases caused by these performance variations (Oberg and Vitek 2009).  
1083 Implementing rigorous QC procedures at each step of LC-MS/MS is essential to maintain  
1084 data reliability and consistency, with standardized QC samples serving as valuable  
1085 benchmarks for long-term performance evaluation.

1086 Several dedicated QC tools, such as MaCProQC (Rozanova et al. 2023), QCloud2 (Olivella  
1087 et al. 2021), Rawtools (Cortay et al. 1988) are available to evaluate the quality of LC-MS/MS  
1088 data. These tools provide a range of functionalities, from tracking performance metrics to  
1089 generating clustering analyses for data quality evaluation. However, more recently, the  
1090 HUPO-PSI Quality Control working group has introduced the mzQC file format, a JSON-  
1091 based standard designed to streamline the reporting and exchange of mass spectrometry  
1092 (MS) quality control metrics. To facilitate adoption, they have also developed open-source  
1093 software libraries in Python (pymzqc), R (rmzqc), and Java (jmzqc), which provide  
1094 functionalities for creating, validating, and analyzing mzQC files. These libraries enable  
1095 researchers to integrate mzQC into diverse workflows for proteomics, metabolomics, and  
1096 other MS applications, ensuring consistent data quality assessment and fostering  
1097 interoperability across different analytical platforms (Bielow et al. 2024).

### 1098 3.7.5 Data management and data sharing

1099 Effective data management and sharing are essential to advancing metaproteomics  
1100 research, ensuring data integrity, reproducibility, and collaboration. A robust data  
1101 management plan should include secure, redundant storage solutions to protect against  
1102 data loss, particularly for large-scale studies conducted over extended periods.  
1103 Implementing version control for raw and processed data facilitates systematic tracking of  
1104 updates and reanalyses, improving reproducibility and transparency.

1105 Adhering to community standards, such as those established by the Human Proteome  
1106 Organization Proteomics Standards Initiative (HUPO-PSI) (Deutsch, Vizcaíno, et al. 2023),  
1107 is crucial for consistency and interoperability. The HUPO-PSI defines data representation  
1108 standards in proteomics to facilitate data comparison, exchange, and verification. Using  
1109 standardized formats like mzML for mass spectrometry data (Martens et al. 2011),  
1110 mzIdentML for identification results (Combe et al. 2024), and the Universal Spectrum  
1111 Identifier (USI) for referring to any mass spectrum in publicly deposited proteomics datasets  
1112 (Deutsch et al. 2021), ensures compatibility across platforms and tools, thereby  
1113 streamlining collaborative efforts and enabling more efficient data use.

1114 Metadata plays a critical role in making datasets interpretable, reusable, and comparable  
1115 across studies. Comprehensive metadata should capture sample origins, preparation  
1116 protocols, instrument settings, and data processing workflows, ideally using standardized  
1117 ontologies like PSI-MS Ontology. In proteomics, this information is collected in the Sample  
1118 and Data Relationship Format for Proteomics (SDRF-Proteomics) format, which provides a  
1119 structured, tab-delimited format for describing the relationships between samples and data  
1120 files, mirroring the experimental workflow in proteomics (Dai et al. 2021). Tools like  
1121 lesSDRF offer user-friendly interfaces to annotate metadata in SDRF format, facilitating  
1122 standardization (Claeys et al. 2023). Recognizing the added complexity of microbial  
1123 environments, the Metaproteomics Initiative is developing SDRF-Proteomics templates  
1124 tailored for metaproteomics, as current formats for single-species proteomics do not fully  
1125 address the nuances of microbial data. Standardized metadata not only supports  
1126 computational analyses but also ensures structured inputs for machine learning models,  
1127 advancing reproducibility and consistency across the field.

1128 Depositing both data and metadata in recognized international ProteomeXchange  
1129 repositories (Deutsch, Bandeira, et al. 2023), such as PRIDE (Perez-Riverol et al. 2024),  
1130 aligns with the FAIR (Findable, Accessible, Interoperable, and Reusable) principles,  
1131 promoting open science and innovation. These repositories make data accessible to the  
1132 broader research community, enabling others to validate findings, conduct systematic  
1133 reviews, and perform large-scale analyses. Sharing practices in metaproteomics help with  
1134 benchmarking studies, development of new interpretation tools, and the ability to draw  
1135 broader conclusions, significantly improving the field's collaborative potential and impact.

## 1136 4. Computational analysis of metaproteomics 1137 data

### 1138 4.1 Peptide identification, protein inference and quantification

1139 After acquiring MS/MS spectra from mass spectrometry, the next step is to identify the  
1140 peptides present in the sample. This involves analyzing the fragmentation patterns in the  
1141 MS/MS spectra to determine the specific amino acid sequences of the peptides. This  
1142 process is performed using search engines, often integrated into comprehensive  
1143 proteomics software packages (**Section 4.1.1**). Typically, these algorithms match the  
1144 experimental MS/MS spectra to a theoretical protein sequence database, and the success  
1145 of this step depends heavily on the selection or construction of an appropriate database, as  
1146 outlined in **Section 4.1.2**. The search engine then applies a false discovery rate (FDR)  
1147 threshold to filter out potential false positives (**Section 4.1.3**). Peptides passing this filter  
1148 are subsequently used for protein inference (**Section 4.1.4**) and quantification (**Section**  
1149 **4.1.5**). All these sections focus on DDA MS, while **Section 4.1.6** is dedicated to tools  
1150 specifically designed for analyzing DIA MS data.

#### 1151 4.1.1 Peptide identification with proteomics search engines

1152 Shotgun metaproteomics experiments generate large datasets of MS1 and MS2 spectra,  
1153 which form the basis for downstream analysis. With advancements in high-throughput MS,  
1154 these datasets now range from thousands to millions of spectra, making manual  
1155 interpretation impractical. To address this challenge, search engines are essential for  
1156 interpreting the data and identifying peptides. Peptide identification relies on three main  
1157 strategies: (i) sequence database searching, where experimental spectra are matched to  
1158 theoretical spectra derived from protein or peptide sequences in a database; (ii) *de novo*  
1159 sequencing, which directly infers peptide sequences from spectra without a reference  
1160 database; and (iii) spectral library searching, where experimental spectra are compared to  
1161 curated libraries of previously validated spectra. These methods are often complemented  
1162 by post-processing steps to enhance accuracy and confidence in peptide identification, as  
1163 outlined in **Section 4.1.3**. Additionally, most proteomics software packages integrate  
1164 peptide identification with protein inference and quantification, a topic discussed in **Section**  
1165 **4.1.4** and **Section 4.1.5**. Some specific metaproteomics software also integrates taxonomic  
1166 and functional analyses, as outlined in **Section 4.2**.

## 1167 i) **Protein sequence database searching**

1168 Database search algorithms are fundamental for interpreting mass spectrometry data,  
1169 particularly in metaproteomics, where the complexity of microbial communities poses  
1170 significant analytical challenges. These algorithms match experimental MS/MS spectra to  
1171 theoretical spectra generated from protein sequence databases. The success of this  
1172 process depends on the choice of search engine, the search parameters used, and the  
1173 composition of the database, all of which influence the number and type of peptides and  
1174 proteins detected.

1175 Database search engines start by using a selected reference protein sequence database,  
1176 which is *in silico* digested to emulate the cleavage rules of the enzyme used during protein  
1177 digestion, most commonly trypsin. From these digested sequences, theoretical MS/MS  
1178 spectra are generated and compared to the experimental MS/MS spectra obtained during  
1179 mass spectrometry. Each combination of theoretical peptide and spectrum (peptide-  
1180 spectrum match, PSM) is assigned a similarity score, with the search engine ranking and  
1181 filtering potential PSMs based on the score and peptide properties. The exact method of  
1182 score calculation varies between search engines, and these differences can affect both  
1183 sensitivity and specificity. An in-depth explanation of the various scoring algorithms used in  
1184 database search engines can be found in this comprehensive review (Verheggen et al.  
1185 2020).

1186 Each database search engine offers unique advantages and limitations, including variations  
1187 in processing speed, compatibility with input and output formats, support for post-  
1188 processing tools, and overall user-friendliness. These factors significantly influence their  
1189 performance in metaproteomics workflows, where the complexity and scale of datasets  
1190 demand highly efficient and reliable analysis tools. A detailed discussion of these tools and  
1191 their applications is available in a comprehensive review (Schiebenhoefer et al. 2019). A  
1192 selection of database search engines and proteomics software commonly used in  
1193 metaproteomics research is highlighted below:

- 1194 • SearchGUI (Vaudel et al. 2011) provides simultaneous access to multiple  
1195 complementary search algorithms, including X!Tandem (Craig and Beavis 2004),  
1196 Comet (Eng, Jahan, and Hoopmann 2013), Andromeda (Cox et al. 2011), OMSSA  
1197 (Geer et al. 2004), Sage (Lazear 2023), and others. Its companion tool,  
1198 PeptideShaker (Vaudel et al. 2015), seamlessly imports SearchGUI output and  
1199 offers a comprehensive, user-friendly interface for interpreting and visualizing  
1200 results. Additionally, PeptideShaker includes a direct export feature to Unipept,



1201 enabling streamlined downstream taxonomic and functional analysis (Vande  
1202 Moortele, Devlaminck, et al. 2024; Van Den Bossche et al. 2020). A detailed  
1203 tutorial is available on the CompOmics web page to guide users through these  
1204 workflows (Vaudel et al. 2014).

- 1205 • Andromeda (Cox et al. 2011), used in MaxQuant (Cox and Mann 2008), is widely  
1206 used for its ease of use and MS1 quantitative capabilities. Users benefit from a  
1207 well-established community, including annual user meetings and a dedicated  
1208 forum for support.
- 1209 • Mascot (Matrix Science) and Proteome Discoverer (Thermo Fisher Scientific) are  
1210 popular commercial tools with extensive user bases.
- 1211 • FragPipe, using MSFragger (A. T. Kong et al. 2017), and pFind (Le-heng Wang et  
1212 al. 2007) incorporate open search strategies, which improve sensitivity by enabling  
1213 the identification of PTMs.
- 1214 • Sipros (Guo et al. 2018), ProteoStorm (Beyter et al. 2018) and COMPIL 2.0 (Park  
1215 et al. 2019) are tailored specifically for metaproteomics but are perceived less  
1216 user-friendly than mainstream software.
- 1217 • Tools such as Sage (Lazear 2023) and MSFragger (A. T. Kong et al. 2017)  
1218 leverage advanced spectral and sequence indexing strategies to significantly  
1219 accelerate database searches, making them highly promising for improving the  
1220 speed of metaproteomics analysis.

1221 For researchers that want more integrated solutions, several software suites can simplify  
1222 metaproteomics workflows by consolidating multiple steps and managing the high density  
1223 of information inherent to the field.

- 1224 • Galaxy for Proteomics (Galaxy-P) is another versatile platform offering numerous  
1225 tools and workflows tailored to metaproteomics, including database generation,  
1226 discovery analysis, verification, quantitation, and statistical analysis (Blank et al.  
1227 2018; P. D. Jagtap et al. 2015; Do et al. 2024). With public gateway availability (The  
1228 Galaxy Community 2024) and access to training resources via the Galaxy Training  
1229 Network (Hiltemann et al. 2023), Galaxy-P is a valuable resource for researchers  
1230 seeking an open and user-friendly platform for users to access metaproteomic  
1231 workflows.
- 1232 • The MetaProteomeAnalyzer (MPA) software suite (Muth, Behne, et al. 2015) offers  
1233 modules for protein database creation, database searching, protein grouping,  
1234 annotation, and results visualization. Its user-oriented design makes it a suitable  
1235 option for both beginners and experienced researchers.

1236 • MetaLab (Cheng et al. 2017) is an integrated data processing pipeline that includes  
1237 tools for sample-specific database generation, peptide determination, taxonomic  
1238 and functional profiling, and abundance analysis. Its open search strategy enables  
1239 comprehensive profiling of PTMs and improved sensitivity. Additionally, MetaLab  
1240 offers workflows for taxonomic analysis based on metagenome-assembled genome  
1241 (MAG) databases, allowing peptide-to-genome linkages for improved specificity  
1242 compared to traditional lowest common ancestor (LCA) methods.

1243 In these tools, selecting appropriate search parameters is essential for reliable and  
1244 meaningful results. The choices regarding modifications, enzyme specificity, and mass  
1245 tolerance significantly impact the identification of PSMs. Below are key considerations:

1246 • Selection of modifications: It is important to distinguish between modifications  
1247 introduced by the experimental workflow and biological modifications. Fixed  
1248 modifications, like carbamidomethylation of cysteine, are commonly applied across  
1249 all peptides to account for standard sample preparation artifacts, as discussed in  
1250 **Section 3.4.4**. Variable modifications, such as methionine oxidation, are applied  
1251 selectively to explore biologically relevant modifications. However, including too  
1252 many variable modifications can expand the search space excessively, reducing  
1253 identification rates. It is often best to limit variable modifications to the most  
1254 biologically relevant ones.

1255 • Enzyme specificity and number of missed cleavages: Choosing the correct enzyme  
1256 and setting an appropriate number of allowed missed cleavages affects the range  
1257 of detectable peptides. For instance, trypsin, the most commonly used enzyme in  
1258 proteomics, may occasionally miss cleavages after lysine (K) or arginine (R).  
1259 Allowing one or two missed cleavages is generally a good compromise in  
1260 metaproteomics, as it accounts for incomplete digestion without excessively  
1261 broadening the search. Semi-specific or non-specific cleavage settings might be  
1262 useful in some cases but can lead to longer processing times and a lower  
1263 identification rate due to the expanded search space.

1264 • Mass tolerance: Mass tolerance settings should match the resolution capabilities of  
1265 the mass spectrometer. For example, on a high-resolution Q Exactive instrument  
1266 with HCD data, setting a precursor mass tolerance of 10 ppm (for MS1) and a  
1267 fragment mass tolerance of 0.02 Da (for MS2) can balance accuracy and  
1268 computational efficiency, restricting the search to relevant matches while taking  
1269 advantage of the instrument's resolution.

1270 Thoughtful parameter selection helps balance sensitivity and specificity, leading to high-  
1271 quality data that accurately reflects the sample's biological characteristics. Parameter  
1272 adjustments should consider the mass spectrometer type, sample complexity, and specific  
1273 research objectives.

#### 1274 **ii) *de novo* searching**

1275 *De novo* peptide sequencing assigns amino acid sequences to MS/MS spectra without  
1276 requiring a protein sequence database for spectral matching. This approach provides an  
1277 unbiased method for detecting peptides, independent of the quality and completeness of  
1278 the protein sequence database. Several *de novo* sequencing algorithms have been  
1279 introduced in recent years, including PEAKS, Casanovo (Yilmaz et al. 2024), PepNovo  
1280 (Frank and Pevzner 2005), and the newly developed  $\pi$ -HelixNovo (T. Yang et al. 2024),  
1281 metaSpectraST (Hao et al. 2023), and NovoBridge (Kleikamp et al. 2021).

1282 When applied effectively, *de novo* sequencing can sensitively and accurately estimate the  
1283 taxonomic composition and functional content of the microbiome without prior knowledge  
1284 of the system under study. It also has the potential to identify unsequenced members of the  
1285 microbial community. Furthermore, *de novo* sequencing can be used to evaluate the  
1286 completeness and suitability of a protein sequence database for metaproteomics research  
1287 (R. S. Johnson et al. 2020). Recently, the progress and opportunities in *de novo* sequencing  
1288 for metaproteomics were reviewed, emphasizing its potential for unsequenced species  
1289 detection and deeper functional insights into microbial communities (Van Den Bossche,  
1290 Beslic, et al. 2024).

1291 Despite its promise, there remains a need for systematic benchmarking of *de novo*  
1292 sequencing tools to assess their applicability to metaproteomics. In particular, most tools  
1293 and approaches for *de novo* metaproteomic analysis still require some input from databases  
1294 either to help selecting peptides or to gain information from the identified peptides.  
1295 Evaluating their performance in terms of sensitivity, accuracy, and throughput is essential  
1296 to ensure their effectiveness in the complex and diverse datasets characteristic of  
1297 microbiome studies.

#### 1298 **iii) Spectral library searching**

1299 Spectral library search engines operate on principles similar to database searching but  
1300 differ by directly comparing experimental MS/MS spectra to pre-existing libraries of  
1301 validated spectra. These libraries consist of MS/MS spectra previously acquired through  
1302 the analysis of complex peptide mixtures and conventional sequence database searches

1303 or generated using predictive deep-learning algorithms. Unlike sequence database  
1304 searching, spectral library searching can incorporate additional parameters, such as  
1305 retention time on the LC column and the relative intensities of fragment peaks within the  
1306 spectra, enhancing both accuracy and confidence in peptide identification.

1307 The development of AI-based tools like MS<sup>2</sup>PIP (Degroeve, Maddelein, and Martens 2015)  
1308 and Prosit (Gessulat et al. 2019) has made it possible to generate high-quality spectral  
1309 libraries from protein sequence databases (Lautenbacher et al. 2024). These  
1310 advancements have expanded the applicability of spectral library searches by enabling the  
1311 generation of predictive libraries tailored to specific experiments. Newer spectral library  
1312 search tools designed for data-dependent acquisition (DDA) data, such as Mistle (Nowatzky  
1313 et al. 2023) and Scribe (Searle, Shannon, and Wilburn 2023), have also emerged for  
1314 metaproteomics research.

1315 Spectral library searching offers a fast and efficient approach to match peptide sequences  
1316 to MS/MS data, often outperforming traditional database searching in terms of speed and  
1317 precision for well-curated libraries. However, despite its potential, spectral library tools for  
1318 metaproteomics require further evaluation, particularly regarding their usability and  
1319 effectiveness for highly complex microbial datasets.

#### 1320 4.1.2 Database construction or selection

1321 For single-organism proteomics, constructing a protein sequence database is relatively  
1322 straightforward, as it can be derived directly from the organism's genome. In  
1323 metaproteomics, however, the complexity of microbial communities, the diversity of  
1324 organisms, and the prevalence of unknown proteins present significant challenges.  
1325 Selecting or generating an appropriate database is crucial, as the database must balance  
1326 comprehensiveness and specificity. An incomplete database risks missing or falsely  
1327 identifying proteins, while an excessively large database decreases the sensitivity of the  
1328 analysis and inflates the FDR, as detailed in **Section 4.1.3** (Nesvizhskii and Aebersold 2005;  
1329 Blakeley-Ruiz and Kleiner 2022).

1330 An optimal database for metaproteomics should be both comprehensive and specific.  
1331 Comprehensive, as it should include all proteins potentially present in the sample. Missing  
1332 sequences lead to false negatives, reducing peptide and protein identification rates.  
1333 Specific, because it should exclude sequences unexpected to be present in the sample.  
1334 Including irrelevant sequences increases random matches, inflates the FDR, and therefore  
1335 negatively affects peptide (and protein) identification (see also **Section 4.1.3**). Additionally,

1336 metaproteomic analyses often include contaminants from sample processing, such as  
1337 leftover trypsin, BSA carry-over, or keratin from handling. Incorporating these contaminants  
1338 into the database, using resources like the common Repository of Adventitious Proteins  
1339 (cRAP, <https://www.thegpm.org/crap/>), allows for their accurate identification and prevents  
1340 misidentification with other proteins in the sample.

1341 To create a suitable database, prior knowledge of the community composition is essential.  
1342 This information can be derived from various sources, including prior literature, 16S rRNA  
1343 amplicon sequencing, or metagenomic and/or metatranscriptomic sequencing, each  
1344 offering different levels of resolution and success. Literature reviews provide only limited  
1345 insights, whereas meta-omics approaches offer the most comprehensive and detailed  
1346 characterization of the community (Kleiner et al. 2012; Blakeley-Ruiz et al. 2022; Minniti et  
1347 al. 2019). Additionally, depending on the sample's environment, host or dietary proteins  
1348 may need to be included in the database. While adding these proteins can improve  
1349 identification rates, it also increases database size and complexity, potentially complicating  
1350 the analysis. The inclusion of nearly identical sequences, often inevitable in large databases,  
1351 can further exacerbate protein inference issues (see **Section 4.1.5**). Sequence clustering  
1352 algorithms (W. Li, Jaroszewski, and Godzik 2001) or protein grouping tools (Audain et al.  
1353 2017; The et al. 2016) can address these challenges by consolidating redundant entries  
1354 while retaining essential taxonomic and functional annotations.

1355 The choice of database type depends on the sample type, the level of understanding of the  
1356 microbial community, and the available resources. Based on these factors, different types  
1357 of databases can be used, each with its own set of advantages and limitations (see **Table**  
1358 **2**). These include public repositories, reference catalogs, and meta-omic databases, as  
1359 detailed below.

#### 1360 **i) Public repositories**

1361 Public repositories like UniProtKB (The UniProt Consortium 2023) and NCBI RefSeq  
1362 (O'Leary et al. 2016) provide extensive reference collections of protein sequences.  
1363 However, these untailored (or unrestricted) databases often lack specificity and contain  
1364 many unrelated sequences, leading to reduced identification rates and increased FDR  
1365 (**Section 4.1.3**). Furthermore, public repositories are biased toward well-characterized  
1366 microbes, such as model organisms or pathogens, and heavily studied environments or  
1367 systems, such as clinical and human samples. This bias results in significant gaps for less-  
1368 studied environmental microbial communities, making these repositories incomplete for  
1369 many metaproteomics applications. Filtering (or restricting) these repositories based on

1370 16S rRNA analysis results can improve specificity, but the resolution of 16S rRNA  
1371 sequencing is limited. Entire genera or sets of species often need to be included, preventing  
1372 strain-level specificity (Odom et al. 2023; J. S. Johnson et al. 2019).

## 1373 **ii) Reference catalogs**

1374 Reference catalogs are curated collections of protein sequences tailored to specific  
1375 environments or systems. They are available for well-studied ecosystems such as the  
1376 human gut (J. Li et al. 2014; Almeida et al. 2021), the cow rumen (Stewart et al. 2019; Xie  
1377 et al. 2021), and the mouse gut (Kieser, Zdobnov, and Trajkovski 2022; Beresford-Jones et  
1378 al. 2022; Lesker et al. 2020). These catalogues are typically constructed by combining data  
1379 from isolated microbes and metagenomic studies (Gurbich et al. 2023). Although smaller  
1380 and more targeted than public repositories, reference catalogs can still be relatively large  
1381 for metaproteomic analyses and often aggregate data from many samples, including  
1382 different individuals and studies - yet, not from the study itself, therefore also called  
1383 unmatched meta-omics databases. This composite nature introduces challenges, as even  
1384 samples from similar environments can exhibit substantial variation in species composition  
1385 and strain diversity. Consequently, reference catalogs can suffer from inaccuracies,  
1386 incompleteness, and overrepresentation of certain subsamples (Van Den Bossche, Kunath,  
1387 et al. 2021; Abdill, Adamowicz, and Blekhman 2022). Like repositories, the specificity of  
1388 reference catalogs can be improved by incorporating prior knowledge of the microbial  
1389 community, such as results from 16S rRNA analysis, to narrow down the included  
1390 sequences to those most relevant to the sample.

1391

1392 Alternatively, to address the challenges posed by large and composite catalogs, database-  
1393 reduction methods have been developed. These methods include the two-step search  
1394 approach (P. Jagtap et al. 2013), iterative workflows such as MetaPro-IQ (X. Zhang et al.  
1395 2016) and MetaLab (Cheng et al. 2017), next to others. While these methods are often used  
1396 in the field and increase the number of identified PSMs and peptides, some have been  
1397 shown to significantly raise the number of false positives at both levels, exceeding the FDR  
1398 estimate (Muth, Kolmeder, et al. 2015). These methods should therefore be treated with  
1399 caution, and additional validation might be appropriate prior to drawing biological  
1400 conclusions.

1401

## 1402 **iii) (matched) meta-omics databases**

1403

1404 Meta-omic databases are constructed using metagenomic and/or metatranscriptomic data  
1405 collected from the same sample as the metaproteomic analysis, making them the most  
1406 specific databases available. These databases accurately reflect the species composition  
1407 and strain diversity of the sample (Blakeley-Ruiz and Kleiner 2022; Heintz-Buschart and  
1408 Wilmes 2018; B. J. Kunath et al. 2022). However, generating a high-quality meta-omic  
1409 database requires significant sequencing effort, cost, computational resources, and  
1410 technical expertise. Although the specific details of this process are beyond the scope of  
1411 this manuscript, they have been extensively covered elsewhere (Blakeley-Ruiz and Kleiner  
1412 2022; Benoit J. Kunath et al. 2019). Briefly, constructing a meta-omic database involves  
1413 four key steps: sequencing, assembly, binning, and annotation.

1414 To create a comprehensive database suitable for metaproteomic analysis, the sequencing  
1415 effort must be sufficiently deep to capture the complexity of the community. One major  
1416 advantage of meta-omic databases is their ability to provide precise insights into the species  
1417 and strain diversity of the sample, enabling direct linkage between genomes and identified  
1418 proteins. This requires genome reconstruction through binning, where contigs are grouped  
1419 into MAGs based on shared features. However, due to the complexity of microbial  
1420 communities and limitations in sequencing depth, some MAGs may remain incomplete.  
1421 Therefore, a robust meta-omic database should include both binned and unbinned  
1422 sequences to retain as much information as possible (Benoit J. Kunath et al. 2017;  
1423 Narayanasamy et al. 2016).

1424 Once reconstructed, MAGs and contigs are taxonomically annotated, and protein  
1425 sequences or open-reading frames (ORFs) are predicted and functionally annotated. The  
1426 choice of tools and resources for these steps depends on the study's objectives (Queirós  
1427 et al. 2021). Despite their specificity, meta-omic databases can still be incomplete due to  
1428 insufficient sequencing depth or the inability to recover all relevant MAGs from the sample.  
1429 This issue can be partially addressed by performing exploratory 16S rRNA gene  
1430 sequencing to assess the required sequencing depth for optimal metagenomic analysis  
1431 (Blakeley-Ruiz et al. 2022).

1432 Combining metagenomic data with metatranscriptomic data further improves the quality  
1433 and specificity of the database (Narayanasamy et al. 2016; F. Delogu et al. 2020). Since  
1434 metatranscriptomics focuses on mRNA, it captures the active portion of the community,  
1435 providing a gene-centric view that aligns closely with the functional content of interest for  
1436 metaproteomics.

1437 **Table 2:** Comparison of database types for metaproteomics: public repositories, reference  
 1438 catalogs, and meta-omic databases. The color indicates our preference: green represents  
 1439 favorable choices, yellow indicates intermediate choices, and red highlights unfavorable choices.  
 1440

	<b>Public repositories (*)</b>	<b>Reference catalogs</b>	<b>Meta-omic databases</b>
<b>Monetary cost</b>	Free	Free	Sample type dependent \$100-\$2,000/sample or pooled samples
<b>Time cost (labor &amp; computation)</b>	Days	Days	Genome-resolved month-year, otherwise weeks
<b>Comprehensiveness</b>	Low to Medium depending on the sample representation in the repository	Medium to High depending on sequencing effort and multi-omics integration	Medium to High depending on sequencing effort and multi-omics integration
<b>Identification probability</b>	Low	Medium	High
<b>Specificity</b>	Low due to high diversity of the repository	Medium due to lack of strains resolution	High due to sample specificity
<b>Misidentification probability</b>	High	Medium	Low
<b>Sequence Redundancy and Impact</b>	High and difficult to resolve due to high diversity of the repository	Medium but can be resolved depending the curation level	Low and can be resolved as part of the metagenomic processing
<b>Taxonomic Annotation and Resolution</b>	Taxonomy not curated and potentially outdated	Depends on curation level (potential for misidentification due to closely related taxa)	Possibility of <i>de novo</i> annotation and species resolution based on metagenomic processing
<b>Certainty/Applicability</b>	Easily available but lacks the guarantee of appropriate sequences	Available for few sample types only and lacks of accuracy	High accuracy but requires particular expertise and extra time/cost



(\*) Restricted repositories have similar characteristics to reference catalogs in terms of specificity and sequence redundancy.

1441

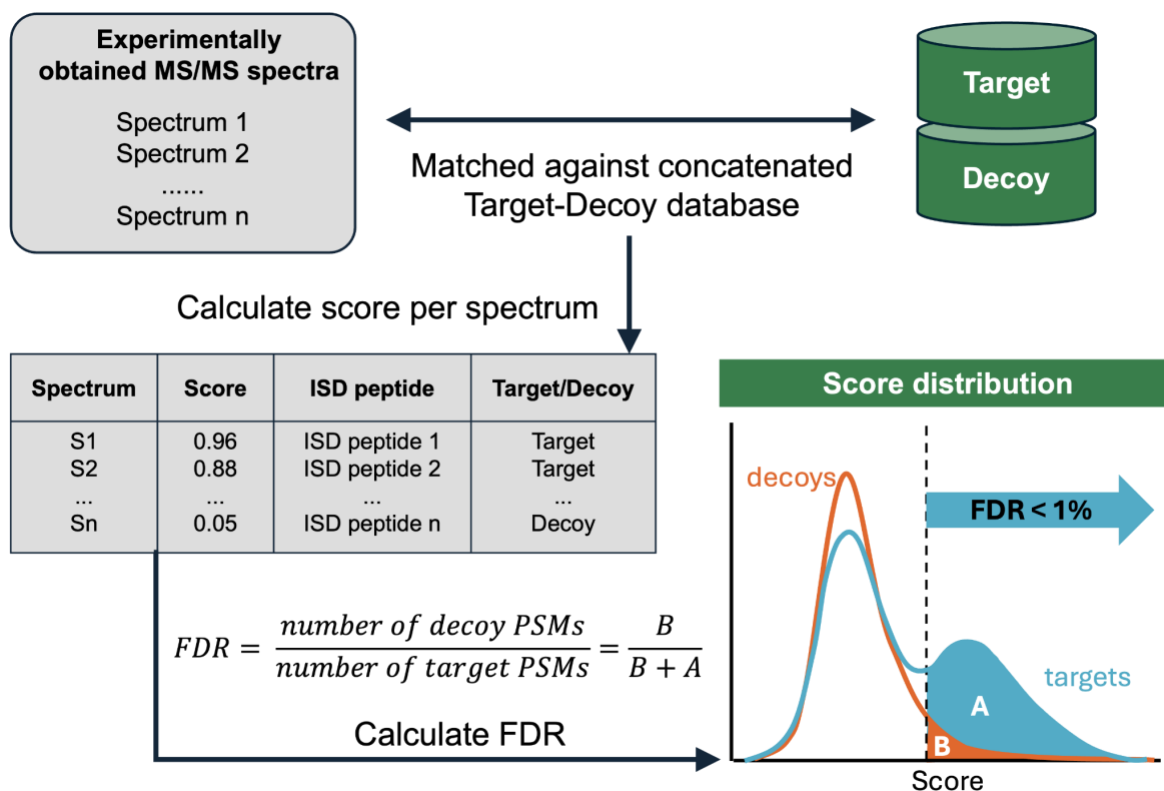
1442

#### 1443 4.1.3 PSM FDR control

1444 A critical step in the process of peptide identification is acquiring a set of reliable PSMs.  
1445 After PSMs are acquired, they are evaluated based on the scoring function of the search  
1446 engine, retaining the highest-ranked PSM for each spectrum — that is, the peptide  
1447 sequence whose theoretical spectrum most closely matches the experimental MS/MS  
1448 spectrum. However, regardless of the scoring algorithm used, some PSMs will inevitably  
1449 represent false matches, making robust control of false positives essential.

1450 The most commonly used strategy to manage false positives in (meta)proteomics is the  
1451 target-decoy approach (Elias and Gygi 2007). In this approach, the protein sequences in  
1452 the target database are processed *in silico* to emulate enzymatic digestion, generating  
1453 theoretical peptides. The same procedure is applied to the reversed or shuffled sequences  
1454 of a decoy database, ensuring that these decoy peptides are biologically implausible and  
1455 not present in the sample. During the search, the experimental spectra are matched to both  
1456 the target and decoy sequences in a concatenated target-decoy database. This process  
1457 results in PSMs being labeled as either target or decoy. The proportion of decoy PSMs in  
1458 the final result serves as an estimate of the FDR, calculated as the number of decoy PSMs  
1459 divided by the total number of accepted PSMs (**Figure 4**). The FDR is typically controlled  
1460 at 1% in proteomics and metaproteomics experiments, but for highly complex samples such  
1461 as soil microbiomes, the FDR threshold can be increased to 5% to retain a sufficient number  
1462 of identifications for biological interpretation.

1463



1464

1465 **Figure 4. Principle of target-decoy analysis and False Discovery Rate (FDR) calculation.** (Top)  
 1466 The experimentally obtained MS/MS spectra are matched to *in silico* generated spectra of the  
 1467 concatenated target/decoy protein sequence database. (Middle) For each obtained spectrum, the  
 1468 match with the highest score is retained, together with the assigned *in silico* digested peptide  
 1469 sequence and its target or decoy label. (Bottom) The score distribution is used to select which PSMs  
 1470 will be considered as true matches. The metric to control the false positives is the FDR, and is  
 1471 calculated as the number of decoy PSMs divided by the number of target PSMs (in the Figure  
 1472 depicted as area B divided by the sum of areas B and A). Figure of (schematic) target/decoy  
 1473 distribution adjusted from (Käll et al. 2008).

1474 The specific challenges of metaproteomics add complexity to FDR control. The larger, more  
 1475 diverse protein sequence databases required for metaproteomics often increase the search  
 1476 space significantly, leading to a greater overlap between the score distributions of target  
 1477 and decoy PSMs. This overlap reduces the resolution of FDR estimation and necessitates  
 1478 careful database construction to limit irrelevant sequences, as discussed in **Section 4.1.2**.  
 1479 Overly large but unspecific databases inflate the FDR by increasing random matches to  
 1480 both target and decoy sequences, resulting in fewer confident peptide identifications  
 1481 (Schiebenhoefer et al. 2019; Tanca et al. 2016). Conversely, overly restrictive databases  
 1482 risk excluding true target sequences, resulting in missed matches, false negatives, and  
 1483 reduced proteome coverage. Therefore, achieving an optimal balance between database

1484 specificity and comprehensiveness is crucial to minimize false positives from decoy  
1485 matches while maximizing target identifications, ensuring effective FDR control.

1486 Metaproteomics workflows often rely on advanced post-processing tools to improve the  
1487 accuracy and confidence of peptide identifications. Tools such as Percolator (Käll et al.  
1488 2007) and MS<sup>2</sup>Rescore (C. Silva et al. 2019) refine PSM scores using machine learning  
1489 algorithms that consider additional features beyond the initial search engine score, such as  
1490 precursor intensity, fragmentation patterns, and retention time. These tools can  
1491 substantially improve the separation between target and decoy PSMs, enabling more  
1492 accurate FDR estimation even for complex datasets.

1493 In metaproteomics, where samples often contain thousands of species, the challenge of  
1494 FDR control is even larger by the inherent complexity and diversity of the microbial  
1495 communities under study. Careful database construction (**Section 4.1.2**), combined with  
1496 robust FDR control during the search and advanced post-processing techniques, is critical  
1497 to ensure reliable peptide and protein identifications, thereby enabling meaningful biological  
1498 insights from metaproteomics datasets.

#### 1499 4.1.4 Protein inference

1500 Protein inference is a fundamental challenge in shotgun proteomics where the goal is to  
1501 determine the proteins present in a sample based on the peptides identified through tandem  
1502 mass spectrometry (Nesvizhskii and Aebersold 2005). This process is complicated by the  
1503 fact that peptides can often be mapped to multiple proteins or protein isoforms present in  
1504 the commonly large protein database. This is especially the case in complex samples such  
1505 as microbial communities where multiple species may contribute homologous proteins,  
1506 making it difficult to conclusively infer which proteins are actually present (Schallert et al.  
1507 2022).

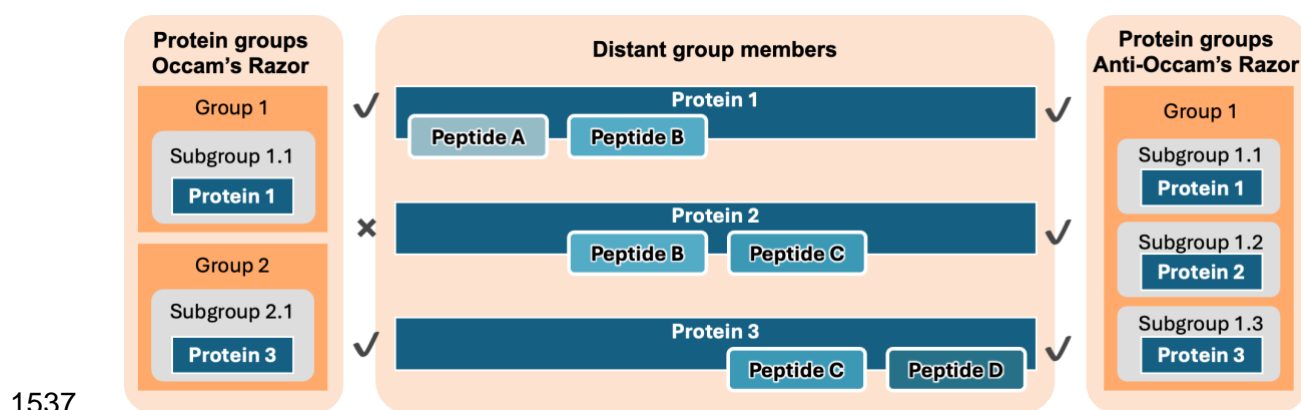
1508 To address this complexity, protein grouping is commonly used to generate a more  
1509 manageable list of identified protein (sub)groups for downstream analysis. However,  
1510 different methods for protein grouping exist, as depicted in **Figure 5**, and these are typically  
1511 performed by the search engine. It is essential to verify the default settings of the search  
1512 engine to understand which grouping approach it applies, and if needed, adjust it to align  
1513 with your research hypothesis. The two main approaches are Occam's razor and anti-  
1514 Occam's razor.

1515 Occam's razor is based on the principle of maximum parsimony, providing the smallest set  
1516 of proteins that can explain all observed peptides. However, this approach discards proteins

1517 not matched by a unique peptide, potentially losing their associated taxonomy and functions  
1518 that might be present in the sample. Occam's razor is particularly suited for simpler, single-  
1519 species samples or targeted proteomics experiments, where reducing complexity is key.

1520 In contrast, anti-Occam's razor adopts a more inclusive strategy, retaining all proteins that  
1521 can be mapped to at least one peptide, regardless of whether those peptides are shared  
1522 with other proteins. This approach is beneficial for complex metaproteomic samples, where  
1523 the goal is to capture as much protein diversity as possible. By being more inclusive, anti-  
1524 Occam's razor ensures that proteins from different species with minimal unique peptides  
1525 are not overlooked, providing a more comprehensive picture of the microbial community.  
1526 However, this inclusivity comes at the cost of increased complexity in the resulting protein  
1527 list.

1528 After choosing between Occam's and anti-Occam's razor principles, proteins can then be  
1529 grouped into protein groups or protein subgroups. Protein groups cluster proteins that share  
1530 at least one peptide, offering a broader overview of potential protein identifications. Protein  
1531 subgroups, on the other hand, are more specific and include proteins that share the exact  
1532 same set of peptides. For example, the anti-Occam's razor approach often benefits from  
1533 subgrouping to prevent excessively large and uninformative protein groups. In  
1534 metaproteomics, this approach helps disentangle the contributions of individual species,  
1535 even when closely related proteins share substantial sequence similarity (Schallert et al.  
1536 2022).



1537

1538 **Figure 5: Practical example of (sub)grouping approaches.** This grouping case deals with distant  
1539 group members, meaning that certain proteins in the group don't share a single peptide, in this case  
1540 protein 1 and 3. Applying the rule of parsimony separates the group in this specific case. In the anti-  
1541 Occam case, protein 2 remains in a separate subgroup.

1542 The choice of protein inference approach should align with the complexity of the sample  
1543 and the research objectives. For single-species or targeted studies, Occam's razor

1544 combined with protein grouping is advantageous for reducing false positives and simplifying  
1545 downstream analyses. This strategy was used, for example, in analyzing the SIHUMIx mock  
1546 community (Schäpe et al. 2019) as part of the CAMPI study (Van Den Bossche, Kunath, et  
1547 al. 2021). For complex, multi-species metaproteomic samples, anti-Occam's razor  
1548 combined with protein subgrouping is often preferred, as it maximizes protein diversity while  
1549 maintaining manageable group sizes. This inclusive approach was used for fecal sample  
1550 analysis in the CAMPI study (Van Den Bossche, Kunath, et al. 2021). Ultimately, the  
1551 selection of a protein inference method depends on the specific characteristics of the  
1552 sample and the research objectives. Researchers must balance the need for  
1553 comprehensive protein identification with the practical considerations of data complexity  
1554 and interpretability (Schallert et al. 2022).

#### 1555 4.1.5 Protein quantification

1556 Protein quantification is a central component of metaproteomics, offering valuable insights  
1557 into the functional dynamics of microbial communities. By quantifying proteins, researchers  
1558 can assess how microbes respond to environmental changes, revealing shifts in physiology  
1559 and metabolic processes. For example, changes in nutrient availability can trigger  
1560 significant alterations in protein expression within individual microbes (Caglar et al. 2017)  
1561 or entire microbial populations (Patnode et al. 2019). This section outlines the key concepts,  
1562 strategies, and challenges in metaproteomic quantification, focusing on label-free and  
1563 labeling-based approaches, as well as methods for downstream data analysis.

1564 Metaproteomics workflows typically rely on two main quantification strategies: label-free  
1565 quantification (LFQ) and labeling-based quantification. LFQ methods are widely used  
1566 because they do not require stable isotope labels, making them more suitable for diverse  
1567 and complex samples. Two common LFQ approaches are MS1 intensity-based  
1568 quantification and MS2 spectral counting. MS1 quantification measures precursor ion  
1569 intensities by calculating the area under the curve or apex intensity for each identified  
1570 peptide, with tools such as MaxQuant (Cox et al. 2014) or standalone alternatives like moFF  
1571 (Argentini et al. 2019) or FlashLFQ (Millikin et al. 2018). MS2 spectral counting, in contrast,  
1572 quantifies peptides based on the number of matched MS2 spectra. Although simpler to  
1573 implement, spectral counting typically has a narrower dynamic range and slightly lower  
1574 precision. Currently, there is limited validation to determine which of the two primary  
1575 quantification approaches—MS1 intensity-based quantification or MS2 spectral counting—  
1576 is more accurate for metaproteomics, or under which conditions one might outperform the  
1577 other. One study demonstrated that spectral counting provided a more accurate measure  
1578 of the proteinaceous biomass of members within a synthetic community compared to MS1

1579 intensities (Kleiner et al. 2017). Nonetheless, the prevailing consensus in the field suggests  
1580 that both methods are generally suitable for metaproteomic quantification, with their  
1581 applicability depending on the specific context and experimental goals.

1582 Labeling-based quantification approaches, while valuable in proteomics, are less commonly  
1583 used in metaproteomics due to the complexity of microbial communities. These methods,  
1584 including TMT and stable isotope labeling by amino acids in cell culture (SILAC), enable  
1585 absolute quantification and are particularly effective for controlled experimental designs  
1586 requiring precise comparisons across samples. However, applying these methods to  
1587 metaproteomics presents significant challenges. The diverse microbial populations and  
1588 high sample complexity of environmental or clinical samples make labeling-based  
1589 approaches less practical, favoring label-free strategies for most metaproteomics workflows.  
1590 Nevertheless, labeling remains a viable option for targeted studies with well-defined  
1591 microbial communities.

1592 Quantification in metaproteomics faces several challenges, particularly in aggregating  
1593 peptide-level data to infer protein abundances. This aggregation process is influenced by  
1594 the protein inference problem (Nesvizhskii and Aebersold 2005), which determines how  
1595 peptides are assigned to proteins or protein (sub)groups (see also **Section 4.1.4**). Most  
1596 software tools automatically assign peptides to proteins or protein groups, facilitating the  
1597 quantification process. Once protein abundance data is obtained, normalization and  
1598 transformation steps are crucial for meaningful statistical analysis. While various  
1599 normalization methods have been proposed for proteomic data (Bubis et al. 2017; Pavelka  
1600 et al. 2008; Välikangas, Suomi, and Elo 2018), the optimal approach for metaproteomics  
1601 remains an area of active research.

1602 One widely used normalization method, particularly for spectral count data, is the  
1603 normalized spectral abundance factor (NSAF) (Florens et al. 2006). This approach  
1604 compensates for biases introduced by protein length and sample variability. It involves  
1605 dividing a protein's PSM count by its amino acid length to account for protein size, followed  
1606 by normalizing against the total PSM count within the sample to reduce between-run batch  
1607 effects. NSAF is relatively simple to calculate, robust to missing values, and particularly  
1608 suited to the sparse data often encountered in metaproteomics. Further transformation,  
1609 such as log or square root normalization, is typically applied to meet the assumptions of  
1610 statistical tests.

1611 A key distinction between standard proteomics and metaproteomics is the need to account  
1612 for the diverse and complex nature of microbial communities. In metaproteomics, it may be

1613 advantageous to normalize protein abundances specifically for organisms or groups of  
1614 organisms within the community. This targeted normalization allows researchers to focus  
1615 on changes in gene expression and function within specific taxa, providing more granular  
1616 insights into microbial activity. The normalized spectral abundance factor per organism  
1617 (orgNSAF) normalization method has been proposed as a solution for this purpose, as it  
1618 enables normalization of protein abundances within defined taxonomic groups (Hinzke et  
1619 al. 2021; Mueller et al. 2010; Ponnudurai et al. 2020).

1620 A unique advantage of metaproteomic data is its ability to generate multiple datasets based  
1621 on the research question. These datasets generally involve summing the abundance of  
1622 constituent proteins into relevant categories. Broadly, there are three main categories: (1)  
1623 individual proteins or groups of proteins with similar sequences, which can offer insights  
1624 into the specific functionalities of individual organisms within the community; (2) categories  
1625 of biological functions assigned to proteins associated with the measured peptides,  
1626 enabling researchers to investigate shifts in overall community functions; and (3) taxonomic  
1627 categories, where protein abundances can be used to estimate the relative contributions of  
1628 different organisms within a microbial community.

1629 The accuracy of both functional and taxonomic quantification is heavily dependent on the  
1630 quality and completeness of protein annotations in the databases used. Functional  
1631 categories can range from highly specific annotations, such as biochemical reactions, to  
1632 broader descriptions of cellular processes like metabolism, gene expression, transport, or  
1633 replication. Similarly, taxonomic quantification can achieve high resolution, down to the  
1634 strain or species level (Brooks et al. 2015; Xiong et al. 2017), but this depends on the depth  
1635 and accuracy of protein annotations. In some cases, it is limited to higher taxonomic ranks  
1636 when annotations are incomplete or ambiguous (Blakeley-Ruiz et al. 2019). Metaproteomic  
1637 measurements, when processed correctly, can provide an accurate representation of the  
1638 relative proteinaceous biomass of microbial species within a community (Kleiner et al. 2017).  
1639 However, the specificity and accuracy of these measurements are closely tied to the  
1640 reliability of the annotations used for protein classification (Blakeley-Ruiz and Kleiner 2022;  
1641 Tanca et al. 2016).

1642 While these approaches enable the generation of robust datasets for understanding  
1643 microbial abundance and function, further validation is necessary to refine these  
1644 methodologies. Current quantification strategies in metaproteomics require additional  
1645 benchmarking to identify optimal or equivalent approaches for various types of studies.  
1646 Future research using mock communities with defined compositions and spike-in proteins

1647 will be crucial for systematically evaluating the accuracy, reproducibility, and reliability of  
1648 protein quantification methods in metaproteomics.

#### 1649 4.1.6 DIA data analysis

1650 The application of DIA-MS in metaproteomics, as discussed in **Section 3.7.2**, demands  
1651 tailored analytical workflows to manage the unique challenges posed by the complexity and  
1652 scale of microbial communities. Unlike DDA, which prioritizes peptide selection, DIA  
1653 generates complex spectra by fragmenting all ions within a predefined m/z range  
1654 simultaneously. This comprehensive approach requires advanced computational tools and  
1655 strategies to handle the resulting data.

1656 Extracting quantitative and identification data from DIA-MS involves specialized software,  
1657 such as Spectronaut (Bruderer et al. 2017), DIA-NN (Demichev et al. 2020), and  
1658 EncyclopeDIA (Searle et al. 2018). These tools rely heavily on pre-existing spectral libraries  
1659 to match experimental spectra to theoretical peptides. Such libraries are often generated  
1660 through prior DDA experiments or predicted from protein sequence databases. While  
1661 promising, library-free approaches that predict spectra directly from protein sequences  
1662 remain computationally intensive and impractical for complex metaproteomics samples  
1663 without additional data reduction strategies. One effective approach is using genome  
1664 sequencing to limit the database search space or performing a preliminary DDA step to  
1665 construct a targeted spectral library. These steps, although resource-intensive, are  
1666 essential for reducing ambiguity in protein and peptide identifications.

1667 Metaproteomics datasets amplify the inherent analytical challenges of DIA-MS due to their  
1668 immense scale, which frequently involves millions of proteins and peptides. This complexity  
1669 can lead to significant computational demands and requires extensive data processing  
1670 pipelines. Direct library-free DIA analysis for such datasets is virtually impossible with  
1671 current technology unless supplemental genome sequencing or DDA-based library  
1672 construction is performed. These preparatory steps add complexity but are critical for  
1673 optimizing DIA's utility in resolving the intricate dynamics of microbial communities.

1674 Recent advancements in MS, including DIA-PASEF (Gómez-Varela et al. 2023) and the  
1675 Orbitrap Astral analyzer (Dumas et al. 2024), have shown potential for enhancing the  
1676 application of DIA-MS in metaproteomics. These technologies allow for deeper proteome  
1677 coverage, improved sensitivity, and more accurate quantification. However, their integration  
1678 into workflows must be carefully aligned with the computational tools and spectral library  
1679 strategies mentioned above to fully exploit their capabilities.



1680 A recent benchmarking study has demonstrated the reproducibility and accuracy of DIA-  
1681 MS for metaproteomic workflows in comparison to DDA-MS methods (Rajczewski et al.  
1682 2024). Using mock communities of known taxonomic composition, DIA-MS consistently  
1683 identified and quantified more peptides and proteins across laboratories. Additionally, the  
1684 reproducibility of protein and peptide identifications was higher in DIA-MS workflows, which  
1685 also provided accurate quantification of both protein abundances and taxonomic groups.  
1686 These findings underscore the advantages of DIA-MS for metaproteomics, including its  
1687 capacity for deep sequencing, robust quantitation, and reproducibility across samples.  
1688 However, current studies also highlight the limitations of existing DIA tools when applied to  
1689 metaproteomic datasets, emphasizing the need for improvements in software capabilities  
1690 to handle the unique complexities of microbiome samples. These insights stress the  
1691 importance of optimizing library generation, computational tools, and workflows to fully  
1692 leverage the potential of DIA-MS for microbial community analysis.

1693 Although DIA-MS presents substantial benefits for reproducible and quantitative analysis,  
1694 its application in metaproteomics is still evolving and faces several technical and  
1695 computational challenges. Advances in mass spectrometry and bioinformatics hold promise  
1696 for addressing these hurdles, enabling deeper insights into microbial community dynamics.  
1697 Ongoing research is needed to refine workflows, optimize computational methods, and  
1698 explore the potential of library-free approaches to broaden its applicability in  
1699 metaproteomics.

## 1700 4.2 Taxonomic and functional analysis

1701 In metaproteomics, researchers aim to characterize microbial communities by determining  
1702 the organisms present (taxonomic analysis) and elucidating their physiological roles  
1703 (functional analysis). These analyses provide critical insights into the composition, diversity,  
1704 and ecological functions of microbial communities across diverse environments. The  
1705 accuracy of these assignments depends on the quality of peptide and protein identifications  
1706 (see **Section 4.1.1**) and is significantly influenced by the choice of database (see **Section**  
1707 **4.1.2**). Below, we describe the methodologies and tools available for taxonomic and  
1708 functional annotation in metaproteomics, emphasizing the importance of robust annotation  
1709 strategies and computational resources.

### 1710 4.2.1 Taxonomic analysis

1711 Taxonomic analysis in metaproteomics identifies the organisms present in a sample based  
1712 on their expressed proteins. This analysis provides insights into microbial community

1713 composition and diversity, linking proteins to their taxonomic origins. Taxonomic  
1714 assignment can be achieved using exact matching or homology-based searches against  
1715 comprehensive databases such as UniProtKB (The UniProt Consortium 2023) or NCBI NR  
1716 (O’Leary et al. 2016).

1717 While numerous metaproteomics-specific tools are available (described in **Section 4.2.4**),  
1718 researchers can also use tools originally developed for metagenomics, such as Centrifuge  
1719 (Kim et al. 2016) and Kraken 2 (Wood, Lu, and Langmead 2019). These tools match  
1720 peptides or proteins to known taxa, but their accuracy depends on the completeness of  
1721 publicly available genome databases. If organisms in the sample have not been previously  
1722 sequenced and deposited, taxonomic assignments may be incomplete or inaccurate.

1723 Alternatively, taxonomic assignments can leverage meta-omics databases derived from  
1724 metagenomic assemblies. Proteins are inherently tied to genomes, and clustering  
1725 metagenomic sequences into MAGs enables genome-centric taxonomy assignment. Tools  
1726 like GTDB-Tk (Chaumeil et al. 2020) use MAG taxonomy to assign taxa to proteins. For  
1727 proteins not linked to MAGs, tools such as CAT (von Meijenfeldt et al. 2019) can infer  
1728 taxonomy based on the context of all the genes in an assembled contig. Advances in long-  
1729 read sequencing are revolutionizing genome assembly from metagenomes, further  
1730 improving taxonomic assignments (Liu et al. 2022).

#### 1731 4.2.2 Functional analysis

1732 Functional analysis of metaproteomes reveals how microbial communities contribute to  
1733 environmental processes, human health, and disease. By measuring the abundance of  
1734 proteins involved in processes such as metabolism, transport, replication, and defense,  
1735 functional analysis provides a window into microbial community dynamics and their roles in  
1736 ecosystems.

1737 To describe microbial functions, various functional ontologies are used:

- 1738 • **Gene Ontology (GO):** Organizes annotations into three categories: molecular  
1739 functions, biological processes, and cellular components. GO terms are used to  
1740 describe what a gene product does (molecular function), the biological goals it helps  
1741 achieve (biological process), and where in the cell it acts (cellular component)  
1742 (The Gene Ontology Consortium 2019)
- 1743 • **Enzyme Commission (EC) numbers:** Categorizes enzymes by the chemical  
1744 reactions they catalyze, particularly useful in studies of enzymatic activity and the  
1745 role these enzymes play in metabolic pathways.

1746       • **Kyoto Encyclopedia of Genes and Genomes (KEGG)**: Maps proteins to  
1747       metabolic and signaling pathways, illustrating their interactions within larger  
1748       biological systems (Kanehisa and Goto 2000)

1749       There are also more specialized ontologies such as MEROPS (Rawlings et al. 2018) for  
1750       proteases and CAZy (Drula et al. 2022) for carbohydrate-active enzymes, including  
1751       glycoside hydrolases, offer enhanced specificity for analyzing distinct functional categories  
1752       within microbial communities.

1753       Functional annotations can rely on computational tools commonly used in metagenomics,  
1754       such as KoFamKOALA (Aramaki et al. 2020), InterProScan (Quevillon et al. 2005), and  
1755       eggNOG-mapper (Cantalapiedra et al. 2021). However, while these tools provide robust  
1756       frameworks for mapping protein functions, more tailored tools specifically designed for the  
1757       unique requirements of metaproteomics are available and discussed in **Section 4.2.4**.

#### 1758       4.2.3 Peptide-centric vs protein-centric approach

1759       In metaproteomics, taxonomic and functional analyses can be performed using either a  
1760       peptide-centric or protein-centric approach:

1761       • **Peptide-centric approach**: Peptides identified through MS are directly annotated  
1762       with taxa and functions based on their matches to *in silico* tryptic digests of known  
1763       protein sequences. This approach ensures that all potential protein matches are  
1764       retained during annotation, providing a broader view of possible taxa and  
1765       functions.

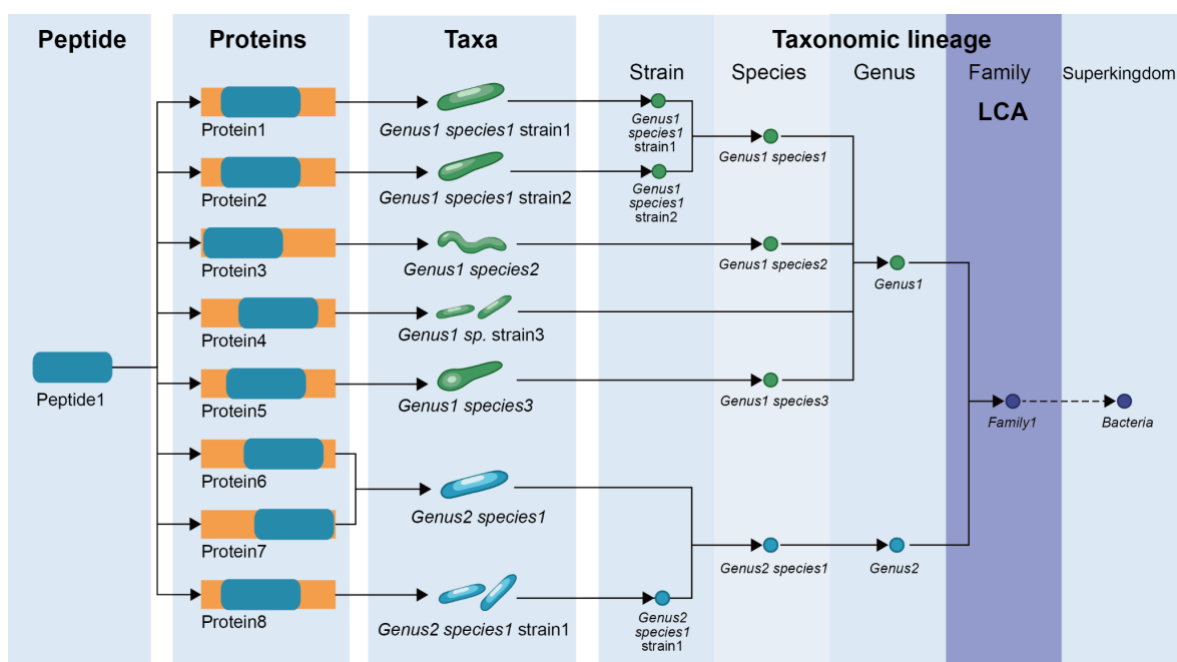
1766       • **Protein-centric approach**: Peptides are first mapped to their corresponding  
1767       proteins or protein (sub)groups, aggregating peptides that share common proteins.  
1768       This step addresses the protein inference problem, a challenge in assigning  
1769       peptides to proteins due to shared sequences among multiple proteins (see  
1770       **Section 4.1.4**).

1771       The peptide-centric approach typically considers all proteins that a peptide could originate  
1772       from, whereas protein-centric tools may discard information deemed redundant based on  
1773       the chosen protein (sub)grouping strategy. These different approaches may lead to  
1774       variations in the resulting annotations, and the debate over which method provides the most  
1775       accurate results remains an active topic in metaproteomics research (Van Den Bossche,  
1776       Kunath, et al. 2021).

#### 1777 4.2.4 Metaproteomics tools for taxonomic and functional analysis

1778 Various tools have been developed for taxonomic and functional analysis in  
1779 metaproteomics, each with distinct features and applications (Sajulga et al. 2020).

1780 Unipept is a powerful ecosystem of tools for the taxonomic and functional analysis of  
1781 metaproteomics samples, offering a command-line interface (CLI), a desktop application, a  
1782 web application, and an application programming interface (API) to accommodate diverse  
1783 user preferences and workflows. (Vande Moortele, Devlaminck, et al. 2024; Verschaffelt et  
1784 al. 2023; 2020). It follows a peptide-centric approach, assigning taxa and functions directly  
1785 to peptides by mapping them to the UniProtKB database. For taxonomic classification,  
1786 Unipept calculates the LCA by identifying the most specific, or lowest, shared taxonomic  
1787 rank among all taxa associated with a peptide's matched proteins (**Figure 6**). More details  
1788 on how the LCA is calculated can be found in a recent comprehensive tutorial (Van Den  
1789 Bossche, Verschaffelt, et al. 2024). Unipept also supports extensive functional analysis by  
1790 reporting functions based on the GO, EC, and InterPro classifications. For each peptide, it  
1791 aggregates all annotations associated with proteins matching the input peptide and counts  
1792 their occurrences. This information is displayed in a table within the web application.  
1793 Detailed tutorials and examples for using Unipept have been published (Mesuere et al.  
1794 2018; Van Den Bossche, Verschaffelt, et al. 2024), and the documentation available on the  
1795 website (<https://unipept.ugent.be/>) offers additional guidance to help users navigate the tool.



1796

1797 **Figure 6. Calculation of the Lowest Common Ancestor (LCA) for a tryptic peptide.** In  
1798 this figure, the hypothetical Peptide 1 is present in eight different proteins, which are

1799 associated with seven distinct organisms. The LCA for these organisms is identified as the  
1800 hypothetical Family 1. Figure adjusted from (Van Den Bossche, Verschaffelt, et al. 2024)

1801 The Peptonizer2000 is a novel metaproteomics pipeline for taxonomic inference that  
1802 models the errors and uncertainties introduced by a typical metaproteomics analysis  
1803 pipeline (Holstein et al. 2024). Indeed, the analysis of mass spectra is inherently challenging:  
1804 researchers need to match observed data to databases of protein sequences, where factors  
1805 such as database bias, ambiguous spectra, degenerate peptide sequences, and  
1806 interspecies sequence homology come into play. The Peptonizer2000 pipeline uses  
1807 Bayesian statistics to model peptide sequences, associated taxa, and the possible errors  
1808 and uncertainties introduced earlier as a graph. Then, subsequently, the Belief Propagation  
1809 algorithm is utilized on this graph to compute probability scores that indicate the potential  
1810 presence of a taxon in a sample under study.

1811 MetaLab (Cheng et al. 2017; 2023; Liao et al. 2018; L. Li et al. 2022) is an integrated  
1812 software platform that provides a streamlined pipeline for microbial identification,  
1813 quantification, and taxonomic profiling using mass spectrometry raw data. Employing a  
1814 hybrid approach, MetaLab combines information derived from both peptide-centric and  
1815 protein-centric metaproteomics analyses. MetaLab utilizes a precomputed index of the  
1816 UniProtKB for taxonomic classification of identified peptides and retrieves functional  
1817 annotations from the eggNOG database (Hernández-Plaza et al. 2023). The latest version  
1818 supports DDA and DIA workflows across various mass spectrometry platforms (Cheng et  
1819 al. 2024). Comprehensive resources on iMetaLab (L. Li et al. 2022) can be found on their  
1820 dedicated Wiki-page (<https://wiki.imetalab.ca/>).

1821 Prophane (Schiebenhoefer et al. 2020) is a software tool designed for taxonomic and  
1822 functional annotation of metaproteomes, offering interactive result visualization and an  
1823 intuitive web-based interface. It integrates data from various annotation databases,  
1824 including NCBI (Schoch et al. 2020), UniProtKB (The UniProt Consortium 2023), eggNOG  
1825 (Hernández-Plaza et al. 2023) or Pfam (Mistry et al. 2021). Unlike tools such as Unipept  
1826 and MetaLab, Prophane adopts a purely protein-centric approach for its analyses. The  
1827 software is accessible both as a Conda package (<https://anaconda.org/bioconda/prophane>)  
1828 and via a web service (<https://prophane.de/login>). Tutorials and example datasets are  
1829 available on the tool's website (<https://prophane.de/about/tutorial>).

1830 The MetaProteomeAnalyzer (MPA) (Muth, Behne, et al. 2015) is an open-source Java tool  
1831 designed for the taxonomic and functional analysis of metaproteomics data. MPA employs  
1832 both sequence-based and spectral-based approaches to identify organisms and functional

1833 pathways in a sample, enabling researchers to explore the metabolic activities of microbial  
1834 communities and their environmental interactions. The software supports multiple search  
1835 engines and incorporates features to reduce data redundancy by grouping protein hits into  
1836 so-called meta-proteins. MPA is available as a desktop application, and extensive tutorials,  
1837 documentation, and other resources are provided on its homepage ([www.mpa.ovgu.de](http://www.mpa.ovgu.de)).

## 1838 4.3 Downstream statistics

1839 A common question among researchers is how to determine the optimal approach for  
1840 downstream processing of metaproteomic data. Unfortunately, there is no universal  
1841 workflow that fits every scenario. This section aims to guide readers in constructing a  
1842 tailored decision tree for analyzing metaproteomic datasets. In earlier sections, we detailed  
1843 the generation of various metaproteomic data tables, including peptides, proteins,  
1844 taxonomy, and functional attributes. The next step involves uncovering the underlying  
1845 patterns and biological insights within these datasets through statistical analysis. Designing  
1846 a robust statistical analysis pipeline for metaproteomics requires researchers to make  
1847 several informed decisions, which are summarized in a “cheat sheet” in **Figure 7**.

### 1848 4.3.1 Identifying relevant scientific questions

1849 The foundation of any metaproteomics analysis begins with defining the key scientific  
1850 question(s) of the study. Metaproteomics allows us to address a variety of research  
1851 objectives. Below are some common examples of questions that can be explored (**Figure**  
1852 **7A**):

- 1853 i. **Cohort studies:** What differential features distinguish healthy individuals from  
1854 those with a disease? Are there potential biomarkers for specific conditions?
- 1855 ii. **Microbiome dynamics:** How does the microbiome vary over space and time?  
1856 Can beta diversity be observed at the functional ecological level? What is the  
1857 impact of specific environmental factors on the microbiome?
- 1858 iii. **Perturbation study:** How do microbial communities respond to external  
1859 perturbations at the taxonomic, functional, and ecological levels?
- 1860 iv. **Multi-omics study:** What (holistic) insights can be gained by integrating  
1861 metaproteomics with other omics approaches?

## 1862 4.3.2 Selecting appropriate levels of analytical insights

1863 Once the primary research questions are defined, the next step is to determine the level of  
1864 insights required to address these questions (**Figure 7B**). This involves selecting between  
1865 different analytical approaches tailored to the objectives of the study:

### 1866 i. **Feature-centric analysis:**

1867 Feature-based methods are the most commonly applied in metaproteomics. These  
1868 analyses focus on identifying differential features, which are quantifiable variables that  
1869 exhibit statistically significant differences between groups or conditions. Examples include  
1870 specific peptides, proteins, taxonomic groups, or annotated functions that vary significantly  
1871 under different experimental conditions.

1872 There are two key considerations that underpin feature-centric analysis: (i) the assumption  
1873 of standard statistical distributions, such as normality, to validate analytical methods, and  
1874 (ii) the treatment of features as independent variables, enabling the use of widely applied  
1875 statistical approaches like parametric or non-parametric tests.

1876

1877 By adhering to these principles, feature-centric analyses enable robust identification of  
1878 biologically meaningful differences across datasets.

### 1879 ii. **Community-centric analysis:**

1880 Unlike feature-centric analysis, community-centric analysis considers the dataset as a  
1881 reflection of a living ecological community. Here, proteins are viewed not as isolated  
1882 features but as components of interconnected networks, with functions linked through  
1883 evolutionary relationships and taxonomic origins. For example, proteins from different taxa  
1884 may exhibit functional redundancy, while ecological dynamics may influence functional and  
1885 taxonomic interactions.

1886 Due to these complex interactions, traditional statistical methods that assume feature  
1887 independence may not be suitable. To address these challenges, novel ecological  
1888 approaches have been developed in metaproteomics, inspired by advancements in  
1889 metagenomics.

1890 For example, metrics for functional redundancy utilize bipartite networks to link  
1891 taxonomic and functional attributes, serving as indicators of community health and stability  
1892 (Blakeley-Ruiz et al. 2019; L. Li et al. 2023). Similarly, PhyloFunc, integrates phylogenetic

1893 composition into functional beta diversity analysis by incorporating functional distances at  
1894 nodes of phylogenetic trees and applying a unifracs-like weighting scheme (Luman Wang et  
1895 al. 2024). This approach distinguishes whether functional changes result from  
1896 compensation among closely related species or shifts between distantly related taxa,  
1897 offering valuable insights into ecological dynamics.

### 1898 **iii. Cross-omics analysis**

1899 The metaproteome is not independent of other meta-omes; therefore, the integration of  
1900 multiple omics datasets is crucial for a deeper understanding of the systems ecology of  
1901 microbiomes. Different meta-omics approaches possess complementary strengths as they  
1902 collectively capture variations along the central dogma of molecular biology (DNA → RNA  
1903 → Protein), favoring a comprehensive understanding of biological processes and ecological  
1904 interactions within microbiomes.

1905 Despite the complementary nature of these datasets, most studies have traditionally  
1906 analyzed meta-omics using separate, stand-alone workflows. However, recent advances in  
1907 bioinformatics tools and platforms, such as Galaxy (Schiml et al. 2023) and MOSCA  
1908 (Sequeira et al. 2024), have facilitated the integration of these datasets, enabling more  
1909 seamless and coherent analysis. Cross-omics analysis can also provide an in-depth view  
1910 of the functional dynamics of community ecology.

1911 In a recent study, metagenomics and metaproteomics were paired to assess whether  
1912 certain proteins serve as niche proteins (proteins that contribute to the ecological role or  
1913 niche that a microbial community occupies within its environment) or play essential  
1914 metabolic roles within a community (T. Wang et al. 2024). To achieve this, genome- and  
1915 proteome-level functional redundancy within the community were compared simultaneously.  
1916 A larger discrepancy might indicate that certain genes are present but not expressed as  
1917 proteins, suggesting a more specialized or niche role. Smaller discrepancies might indicate  
1918 that the genes are actively translated into proteins, suggesting essential metabolic functions.

### 1919 **4.3.3 Data preprocessing strategies**

1920 After making the relevant decisions outlined in **Sections 4.3.1 and 4.3.2**, the first step in  
1921 downstream analysis is data preprocessing. Common preprocessing steps include data



1922 filtering, data transformation, data imputation, and data scaling (**Figure 7C**). However, there  
1923 is no universal approach for data preprocessing; the best strategy depends on the specific  
1924 research questions under investigation.

#### 1925 **i) Data transformation**

1926 Common data transformations used in proteomics and metaproteomics include logarithmic  
1927 transformations (e.g., log<sub>2</sub> or log<sub>10</sub>) and square root transformations. However, not all  
1928 scenarios are suitable for data transformation.

1929 **When to use data transformation:** Transformation is recommended when achieving near-  
1930 normality in the data is necessary. For feature-level analyses, log transformation of peak  
1931 intensities can make the data approximate a normal distribution. Normal distributions are  
1932 crucial for many commonly applied metaproteomic feature selection methods, such as  
1933 linear models, empirical Bayes, univariate t-tests, partial least squares discriminant analysis  
1934 (PLS-DA), and orthogonal partial least squares discriminant analysis (OPLS-DA). If the data  
1935 are not normally distributed, alternative non-parametric methods may be considered to  
1936 meet the assumptions of the chosen analysis.

1937 **When not to use data transformation:** Transformation should be avoided when reflecting  
1938 protein abundance. For example, volcano plots, often used for identifying differential  
1939 features, plot statistical significance (-log<sub>10</sub>(p-value)) against fold change (log<sub>2</sub> fold  
1940 change). While fold change values are log-transformed for visualization purposes, the  
1941 original fold change data should remain untransformed during statistical analyses or  
1942 comparisons. Additionally, in community-level analyses, log transformation can obscure  
1943 protein biomass information, which is essential for estimating taxonomic and functional  
1944 compositions. Protein intensities or PSM counts can serve as reliable measures of protein  
1945 biomass contributions by taxa (Kleiner et al. 2017). Therefore, composition-based analyses,  
1946 such as alpha and beta diversity or functional redundancy assessments, should use  
1947 untransformed data.

#### 1948 **ii) Data centering and scaling**

1949 In standard metaproteomics workflows, an equal amount of protein is typically extracted  
1950 from each sample, digested, and loaded into the mass spectrometer to ensure consistency  
1951 and comparability. However, in specific cases, metaproteomics may quantify overall protein  
1952 biomass responses based on the total protein biomass in a given system volume rather  
1953 than standardizing based on protein content (L. Li et al. 2020). In such cases, centering and

1954 scaling are not recommended. Alternative normalization techniques, such as total spectral  
1955 count normalization or median normalization, may be more appropriate for these scenarios.

### 1956 **iii) Data filtering**

1957 Filtering the dataset typically helps remove noise, irrelevant features, or outliers. The  
1958 application of data filtering should be tailored to the specific context of the study:

1959 **When to use stringent data filtering?** For feature-centric analysis. When identifying  
1960 biomarkers, it is essential to apply stringent data filtering. This involves setting a higher  
1961 threshold for the presence of proteins across samples to ensure that the identified  
1962 biomarkers are consistently found in the majority of subjects. By requiring proteins to be  
1963 present in a large percentage of samples (e.g., 70-90%), researchers can improve the  
1964 reliability and relevance of the identified biomarkers. This consistency is critical for  
1965 validating potential biomarkers, as it reduces the likelihood of identifying false positives.  
1966 Data filtering is also typically stringent for other types of feature-centric analysis to ensure  
1967 the validity of statistical hypotheses. However, the threshold and method of filtering (e.g.,  
1968 by the whole dataset or by group) must be properly applied to prevent over-filtering, which  
1969 could remove features that are truly missing in specific subgroups.

1970 **When is data filtering optional?** For community-centric analysis. While some level of  
1971 filtering is still beneficial to remove obvious noise, the thresholds can be less stringent  
1972 compared to feature-centric analysis. This allows for a more comprehensive view of  
1973 community dynamics. For example, unfiltered taxon-specific functional data can provide a  
1974 better review of the degree distribution of functions in a microbiome (L. Li et al. 2023).

### 1975 **iv) Data imputation**

1976 In a metaproteomic dataset, missingness often arises from two simultaneous mechanisms.  
1977 First, the diversity and sparse nature of the metaproteome lead to a significant proportion  
1978 of true missing proteins (missing not at random) (Plancade et al., 2022). Second, the  
1979 inherent depth limitation of current common metaproteomic techniques results in highly  
1980 sparse detection of low-abundance proteins across samples (missing at random) (Plancade  
1981 et al. 2022).

1982 Data imputation is the step that requires the most caution. Improper selection of the data  
1983 imputation approach can induce false positives. When a large proportion (e.g., >50%) of a  
1984 feature is missing, excessive imputation can lead to the creation of artificial values that do  
1985 not reflect the true biological scenario and, in some cases, can further lead to false positives.

1986 If the imputation method does not accurately reflect the nature of the missing data, it can  
1987 introduce bias, particularly if the data contains a mixture of both missingness mechanisms.  
1988 If features have been selected through a statistical test following data imputation, it is  
1989 recommended to always revisit the un-imputed data to double-check if the feature-level  
1990 difference is true before drawing solid conclusions.

1991 Alternatively, a univariate selection method has been which combines a test of association  
1992 between missingness and classes with a test for the difference in observed intensities  
1993 between classes. This method provides a robust alternative for handling missing data  
1994 without relying on imputation (Plancade et al. 2022).

1995 Notably, data imputation is essential for feature selection analysis, whereas for community-  
1996 level approaches, it is typically unnecessary, for reasons similar to those explained above.

#### 1997 4.3.4 Choosing data analysis methods

1998 After a thorough understanding and careful selection of preprocessing steps, the final step  
1999 in downstream data analysis is the selection of appropriate methods. This stage presents  
2000 significant opportunities for deriving diverse insights from the dataset and is often the most  
2001 engaging and time-consuming phase, allowing researchers to explore the data and uncover  
2002 meaningful biological or ecological patterns and conclusions. These strategies typically  
2003 include, but are not limited to:

- 2004 ● **Dimensionality reduction:** Dimensionality reduction methods are commonly used  
2005 to uncover underlying patterns or structures within the dataset and to assess  
2006 similarities between samples. Unsupervised methods such as Principal Component  
2007 Analysis (PCA), t-distributed Stochastic Neighbor Embedding (t-SNE), hierarchical  
2008 clustering, and k-means clustering are frequently applied. Supervised methods,  
2009 such as Partial Least Squares Discriminant Analysis (PLS-DA), are also widely  
2010 utilized. Dimensionality reduction is applicable not only to peptide, protein,  
2011 taxonomic, and functional tables but also at the MS1 level, especially when the  
2012 primary goal is to reveal patterns between samples (Simopoulos et al. 2022).
- 2013 ● **Enrichment analysis:** Enrichment analysis determines whether a subset of  
2014 selected features is significantly overrepresented compared to a background  
2015 database. While enrichment analysis can be implemented using programming  
2016 languages such as R, iMetaShiny (L. Li et al. 2022) offers interactive functionality  
2017 for taxonomic and functional enrichment analysis of protein IDs or COG IDs.

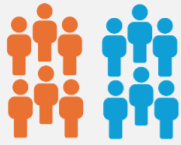
2018                    However, protein ID-based enrichment analysis is currently restricted to human gut  
2019                    metaproteome analysis using the IGC database.

- 2020                    ● **Feature Selection:** Several online tools, such as MetaFS (Tang et al. 2021),  
2021                    MetaQuantome (Easterly et al. 2019), MetaX (Qing Wu et al. 2024), iMetaShiny (L.  
2022                    Li et al. 2022), and stand-alone tools, such as Meta4P (Porcheddu et al. 2023) have  
2023                    been developed to facilitate feature-based metaproteomic data analysis without  
2024                    requiring extensive programming expertise.
- 2025                    ● **Pathway analysis:** Pathway analysis is typically employed to gain an overview of  
2026                    detected functions or to compare differentially expressed or enriched pathways  
2027                    across groups. The most commonly used tools for pathway analysis include KEGG  
2028                    mapper (Kanehisa and Goto 2000) and iPath (Letunic et al. 2008). More recently,  
2029                    PathwayPilot was developed to easily compare functions at the KEGG pathway  
2030                    level, either between selected taxa within a single sample or across different  
2031                    samples, by leveraging EC numbers to identify active enzymes as proxies for  
2032                    metabolites linked to KEGG maps, thereby facilitating investigations into functions  
2033                    associated with specific conditions while allowing targeted analysis of selected  
2034                    species (Vande Moortele, Verschaffelt, et al. 2024).
- 2035                    ● **Community analysis:** Beyond feature-driven analysis, community-level analysis  
2036                    focuses on viewing the entire metaproteome as a dynamic system. Such analyses  
2037                    may include inferring community composition, alpha diversity, beta diversity, and  
2038                    functional redundancy using metaproteomic data.

# Metaproteomics Down-stream Data Analysis “Cheat Sheet”

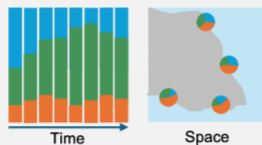
## A. Main domains of questions that metaproteomics down-stream analysis cares about

### Cohort Study



- Differential features?
- Biomarkers?
- Predictive model?

### Longitudinal Study



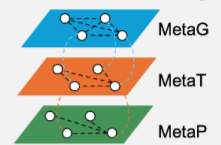
- Spatial-temporal variations?
- Beta-diversity?
- Effect of environment?

### Perturbation Study



- Taxonomic response?
- Pathway response?
- Ecosystem response?

### Multi-omic Study



- Biological process?
- Ecological insights?

## B. Identify desired insight levels to facilitate analysis strategy selection.

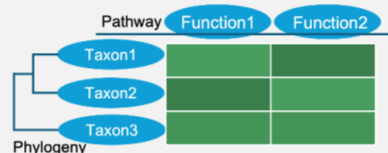
### Feature-centric analysis



Features are considered independent from each other

- Pros: wide statistical applicability
- Cons: usually ignores inter-dependency

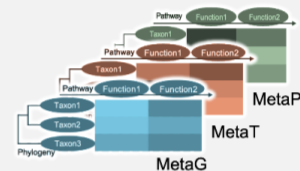
### Community-centric analysis



Interrelationship among features

- Pros: reflect ecological property
- Cons: less applicable methods to-date

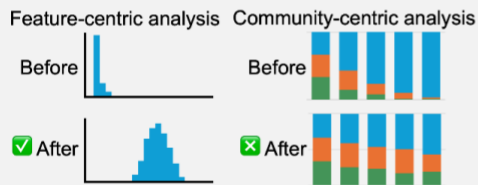
### Cross-omics analysis



Biological process along central dogma

## C. Proper choice of data pre-processing workflow

### Data transformation



- Pros: suitable for more statistical methods
- Cons: may mislead ecological interpretation

### Data centering and scaling



✓ Equal protein assumption

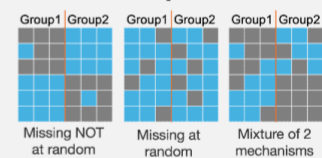
✗ For biomass evaluation

### Data filtering

✓ Usually stringent for feature-centric analysis

✗ Optional for community-centric analysis

### Data imputation

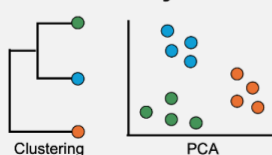


✗ Do not overly impute

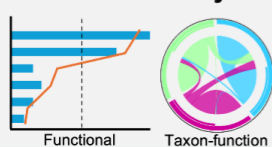
- Check results against original data
- Consider imputation-free methods

## D. Selection of data analysis method set

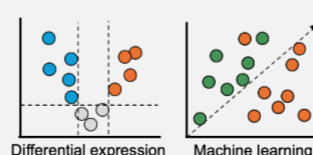
### Dimensionality reduction



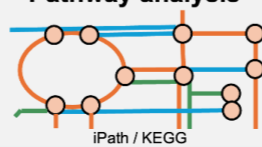
### Enrichment analysis



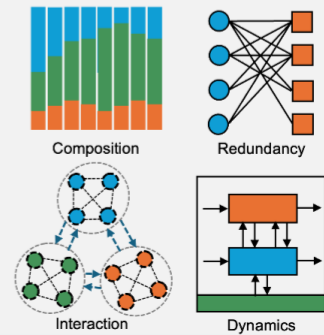
### Feature Selection



### Pathway analysis



### Community analysis



2039

2040 Figure 7. Metaproteomics down-stream data analysis “cheat sheet”

2041

## 2042 5. A collaborative effort: writing a comprehensive 2043 review with members of the Metaproteomics Initiative

2044 The Metaproteomics Initiative is an international community dedicated to advancing the  
2045 field of metaproteomics within microbiome research. Supported by HUPO and EuPA and  
2046 in collaboration with ELIXIR, this initiative serves as a central hub for researchers to  
2047 disseminate advancements, share methodologies, and establish standards across the  
2048 metaproteomics community.

2049

2050 This Initiative aims to facilitate communication between experts and newcomers,  
2051 standardize practices, and accelerate developments in metaproteomic methodologies. Its  
2052 primary mission is to be the go-to resource for metaproteomics fundamentals,  
2053 advancements, and applications, fostering a collaborative network to drive forward  
2054 experimental and bioinformatic methodologies.

2055 The Metaproteomics Initiative supports on three pillars:

- 2056 1. **Communication & Collaboration:** This pillar focuses on sharing field  
2057 advancements, organizing benchmark studies like CAMPI, and hosting the  
2058 International Metaproteomics Symposium (IMS).
- 2059 2. **Education & Outreach:** The initiative educates the broader microbiome community  
2060 through accessible resources, including webinars and workshops, and facilitates  
2061 expert interactions.
- 2062 3. **Standardization:** Efforts are directed toward developing robust (meta)data  
2063 standards, promoting FAIR data principles to ensure accessible and reusable  
2064 research outputs.

2065 As part of our commitment to Education & Outreach, we created this review to make  
2066 metaproteomics accessible to a broad audience. To ensure a thorough and well-rounded  
2067 perspective, we first invited experts in various areas to draft individual sections. These  
2068 drafts were then reviewed internally, where initial feedback helped refine each section.  
2069 Once authors made adjustments, the document went through additional rounds, allowing  
2070 all contributors to share insights and address any remaining comments.

2071 In the next step, we brought in microbiome researchers who were new to metaproteomics  
2072 to review the manuscript, helping us ensure it was clear and approachable for those outside  
2073 the field. With their feedback integrated, all co-authors—including section authors and both

2074 expert and novice reviewers—had a final opportunity to review the work. This collaborative  
2075 approach allowed us to prepare a comprehensive, accessible resource, which we shared  
2076 as a preprint before journal submission.

## 2077 6. Conclusion

2078 This *Microbiologist's Guide to Metaproteomics* is designed for microbiome researchers  
2079 starting in metaproteomics, offering a practical introduction to reduce barriers to entry. It  
2080 covers the essentials of metaproteomics, including experimental design, sample  
2081 preparation, mass spectrometry data acquisition, peptide identification, protein inference,  
2082 taxonomic and functional analysis, and basic statistical methods. The guide provides the  
2083 foundational knowledge needed to apply metaproteomic technologies in microbiology and  
2084 microbiome studies.

2085 Metaproteomics is a rapidly evolving field with unresolved technical challenges and  
2086 unexplored areas. This guide focuses on foundational concepts rather than providing  
2087 exhaustive coverage. To address these challenges, the Metaproteomics Initiative launched  
2088 the "Critical Assessment of Metaproteome Investigations" (CAMPI) series, which facilitates  
2089 multi-laboratory collaborations to compare and improve workflows, including sample  
2090 preparation, mass spectrometry methods, and bioinformatics.

2091 Looking ahead, the next decade promises remarkable advancements in mass spectrometry,  
2092 with continually improving performance deepening the coverage of metaproteomic analysis.  
2093 These advancements, coupled with ongoing and future enhancements in wet-lab protocols,  
2094 strategies, and bioinformatic tools, will further propel the field. Collaborative efforts, such as  
2095 the CAMPI series of the Metaproteomics Initiative, underscores the power of cooperation  
2096 in driving metaproteomic progress. These developments, supported by input from  
2097 microbiome researchers, will help deepen our understanding of microbiomes and their  
2098 functions in diverse ecosystems.

## 2099 Author contributions

2100 This review is a collaborative effort led by Tim Van Den Bossche and Leyuan Li, and  
2101 overseen by the Scientific Committee of the Metaproteomics Initiative, who provided  
2102 overall guidance. Each section was contributed by nominated authors and internal  
2103 reviewers as follows: **Section 1: Why metaproteomics?** was written by Robert Hettich,  
2104 Jean Armengaud, Dirk Benndorf, and Paul Wilmes, and reviewed by Daniel Figeys.

2105 **Section 2: Basics of proteomics** was written by Zhibin Ning and Daniel Figeys, and  
2106 reviewed by Leyuan Li. **Section 3: Experimental methods in metaproteomics** includes  
2107 several subsections: **3.1 Experiment Design** was written by Lucia Grenga and Jean  
2108 Armengaud, reviewed by Céline Henry and Leyuan Li; **3.2 Sample collection,**  
2109 **preservation, and storage prior to preprocessing**, where **3.2.1 Sample collection and**  
2110 **preservation** was written by Sergio Uzzau and Alessandro Tanca, and **3.2.2 Storage**  
2111 **conditions to maintain sample integrity** was written by Lucia Grenga and Jean  
2112 Armengaud, both reviewed by Céline Henry and Leyuan Li; **3.3 Sample preprocessing**  
2113 was written by Lucia Grenga and Jean Armengaud, reviewed by Céline Henry; **3.4**  
2114 **Protein sample preparation: from extraction to digestion** was written by Nico  
2115 Jehmlich, reviewed by Xu Zhang and Céline Henry; **3.5 Separation and fractionation**  
2116 **techniques** was written by Xu Zhang and Marybeth Creskey, reviewed by Céline Henry;  
2117 **3.6 Automation** was written by Leyuan Li, reviewed by Sergio Uzzau and Alessandro  
2118 Tanca; **3.7 Mass spectrometry data acquisition methods** was written by Zhibin Ning  
2119 and Daniel Figeys, reviewed by Jean Armengaud and Céline Henry. **Section 4:**  
2120 **Computational analysis of metaproteomics data** includes several subsections: **4.1.1**  
2121 **Peptide identification with proteomics search engines** was written by Pratik Jagtap,  
2122 Subina Mehta, and Timothy Griffin, reviewed by Tanja Holstein and Kai Cheng; **4.1.2**  
2123 **Database construction or selection** was written by Paul Wilmes and Benoit Kunath,  
2124 reviewed by Jose Alfredo Blakely-Ruiz; **4.1.3 PSM FDR control**, by Tim Van Den  
2125 Bossche and Lennart Martens, reviewed by Tanja Holstein; **4.1.4 Protein inference** was  
2126 written by Tim Van Den Bossche, reviewed by Tanja Holstein; **4.1.5 Protein**  
2127 **quantification** was written by Jose Alfredo Blakely-Ruiz and Manuel Kleiner, reviewed by  
2128 Tanja Holstein and Kai Cheng; **4.1.6 DIA data analysis** was written by Pratik Jagtap,  
2129 reviewed by Tanja Holstein and Kai Cheng; **4.2: Taxonomic and functional Analysis**  
2130 was written by Pieter Verschaffelt and Bart Mesuere, reviewed by Tanja Holstein and Tim  
2131 Van Den Bossche; **4.3 Downstream statistics** was written by Leyuan Li, reviewed by  
2132 Tanja Holstein and Lucia Grenga. **Section 5: A collaborative effort: writing a**  
2133 **comprehensive review with members of the Metaproteomics Initiative** was written by  
2134 Tim Van Den Bossche, and reviewed by Leyuan Li. We invited Madita Brauer, Xuxa  
2135 Malliet, Jing Wang, Xin Zhang, Jong Kim to review the manuscript to ensure its  
2136 accessibility. All figures were artistically designed by Leyuan Li based on author drafts. To  
2137 homogenize the text, ensure consistency and avoid redundancy across sections, all  
2138 sections were rewritten by Tim Van Den Bossche. All authors commented and approved  
2139 the final version of the manuscript.  
2140



## Abbreviations

ABP	activity-based protein probing
ACN	acetonitrile
AI	artificial intelligence
API	application programming interface
AUC	area under the curve
BSHs	bile salt hydrolases
BSL	biosafety level
CLI	command-line interface
cRAP	common repository of adventitious proteins
CZE	capillary zone electrophoresis
DDA	data-dependent acquisition
DIA	data-independent acquisition
FASP	filter-aided sample preparation
FDR	false discovery rate
Galaxy-P	Galaxy for proteomics
GHs	glycoside hydrolases
HILIC	hydrophilic interaction liquid chromatography
HPLC	high-performance liquid chromatography
IMAC	Immobilized metal affinity chromatography
iST	in-stage tips
LC	liquid chromatography
LCA	lowest common ancestor
LFQ	label-free quantification
MAG	metagenome-assembled genome
MPA	MetaProteome Analyzer
MS	mass spectrometry

MuDPIT	multidimensional protein identification technology
NSAF	normalized spectral abundance factor
ORF	open-reading frame
PSMs	peptide-spectrum matches
PTM	post-translational modification
QC	quality control
RP	reverse phase
SCX	strong cation exchange
SILAC	stable isotope labeling by amino acids in cell culture
SP3	single-pot solid-phase-enhanced sample preparation
SPE	solid-phase extraction
TMT	tandem mass tags

2142

## 2143 Acknowledgements

2144 TV acknowledges funding from the Research Foundation Flanders (FWO) [grant  
 2145 1286824N]. JA acknowledges funding from the French National Agency for Research  
 2146 (Agence Nationale de la Recherche, grant ANR-20-CE34-0012) and Occitanie Région  
 2147 (grant 21023526-DeepMicro). JB acknowledges funding from the National Institute Of  
 2148 General Medical Sciences of the National Institutes of Health under Award Number  
 2149 R35GM138362. RH acknowledges funding from the United States Department of Energy,  
 2150 Biological and Environmental Research Program. PJ and TG acknowledge funding from  
 2151 the Minnesota Ovarian Cancer Alliance, the National Institutes of Health/National Cancer  
 2152 Institute [grants 5R01CA262153 (A.P.N.S.), 1R21CA267707], and The National Institutes  
 2153 of Health/National Cancer Institute [grant P30CA077598]. MK acknowledges funding from  
 2154 the National Institute Of General Medical Sciences of the National Institutes of Health [grant  
 2155 R35GM138362]. BK acknowledges funding from the FNR  
 2156 INTERMOBILITY/2022/BM/16965254. LM acknowledges funding from the Research  
 2157 Foundation Flanders (FWO) [grants G028821N, G010023N]. AT acknowledges funding  
 2158 from the Next Generation EU [grant PNRR-MAD-2022-12376416]. SU acknowledges  
 2159 funding from the Next Generation EU [grant PNRR-MAD-2022-12376416]. PV  
 2160 acknowledges funding from the Ghent University (BOF) [grant BOF/01P10623]. JW, XZ, LL

2161 acknowledges funding from the State Key Laboratory of Medical Proteomics, National  
2162 Center for Protein Sciences (Beijing), China. PW acknowledges funding from the European  
2163 Research Council (ERC) under the European Union's Horizon 2020 research and  
2164 innovation program [grant 863664].

## 2165 Conflicts of Interest

2166 DF is a Co-founder of MedBiome inc.

## 2167 References

- 2168 Aakko, Juhani, Sami Pietilä, Tomi Suomi, Mehrad Mahmoudian, Raine Toivonen, Petri  
2169 Kouvonen, Anne Rokka, Arno Hänninen, and Laura L. Elo. 2020. 'Data-  
2170 Independent Acquisition Mass Spectrometry in Metaproteomics of Gut  
2171 Microbiota—Implementation and Computational Analysis'. *Journal of Proteome  
2172 Research* 19 (1): 432–36. <https://doi.org/10.1021/acs.jproteome.9b00606>.
- 2173 Abdill, Richard J., Elizabeth M. Adamowicz, and Ran Blehman. 2022. 'Public Human  
2174 Microbiome Data Are Dominated by Highly Developed Countries'. *PLoS Biology*  
2175 20 (2): e3001536. <https://doi.org/10.1371/journal.pbio.3001536>.
- 2176 Almeida, Alexandre, Stephen Nayfach, Miguel Boland, Francesco Strozzi, Martin  
2177 Beracochea, Zhou Jason Shi, Katherine S. Pollard, et al. 2021. 'A Unified Catalog  
2178 of 204,938 Reference Genomes from the Human Gut Microbiome'. *Nature  
2179 Biotechnology* 39 (1): 105–14. <https://doi.org/10.1038/s41587-020-0603-3>.
- 2180 Andersen, Thea Os, Benoit J. Kunath, Live H. Hagen, Magnus Ø. Arntzen, and Phillip B.  
2181 Pope. 2021. 'Rumen Metaproteomics: Closer to Linking Rumen Microbial Function  
2182 to Animal Productivity Traits'. *Methods, Methods to face the challenges of  
2183 ruminant phenotyping*, 186 (February):42–51.  
2184 <https://doi.org/10.1016/j.ymeth.2020.07.011>.
- 2185 Aramaki, Takuya, Romain Blanc-Mathieu, Hisashi Endo, Koichi Ohkubo, Minoru  
2186 Kanehisa, Susumu Goto, and Hiroyuki Ogata. 2020. 'KofamKOALA: KEGG  
2187 Ortholog Assignment Based on Profile HMM and Adaptive Score Threshold'.  
2188 *Bioinformatics* 36 (7): 2251–52. <https://doi.org/10.1093/bioinformatics/btz859>.
- 2189 Argentini, Andrea, An Staes, Björn Grüning, Subina Mehta, Caleb Easterly, Timothy J.  
2190 Griffin, Pratik Jagtap, Francis Impens, and Lennart Martens. 2019. 'Update on the  
2191 moFF Algorithm for Label-Free Quantitative Proteomics'. *Journal of Proteome  
2192 Research* 18 (2): 728–31. <https://doi.org/10.1021/acs.jproteome.8b00708>.
- 2193 Arıkan, Muzaffer, Tuğçe Kahraman Demir, Zeynep Yıldız, Özkan Ufuk Nalbantoğlu, Nur  
2194 Damla Korkmaz, Nesrin H. Yılmaz, Aysu Şen, et al. 2023. 'Metaproteogenomic  
2195 Analysis of Saliva Samples from Parkinson's Disease Patients with Cognitive  
2196 Impairment'. *Npj Biofilms and Microbiomes* 9 (1): 1–10.  
2197 <https://doi.org/10.1038/s41522-023-00452-x>.
- 2198 Armengaud, Jean. 2023. 'Metaproteomics to Understand How Microbiota Function: The  
2199 Crystal Ball Predicts a Promising Future'. *Environmental Microbiology* 25 (1): 115–  
2200 25. <https://doi.org/10.1111/1462-2920.16238>.
- 2201 Armengaud, Jean, Joseph A Christie-Oleza, Gérémy Clair, Véronique Malard, and  
2202 Catherine Duport. 2012. 'Exoproteomics: Exploring the World around Biological  
2203 Systems'. *Expert Review of Proteomics* 9 (5): 561–75.  
2204 <https://doi.org/10.1586/epr.12.52>.
- 2205 Audain, Enriqe, Julian Uszkoreit, Timo Sachsenberg, Julianus Pfeuffer, Xiao Liang,

2206 Henning Hermjakob, Aniel Sanchez, et al. 2017. 'In-Depth Analysis of Protein  
2207 Inference Algorithms Using Multiple Search Engines and Well-Defined Metrics'.  
2208 *Journal of Proteomics* 150 (January):170–82.  
2209 <https://doi.org/10.1016/j.jprot.2016.08.002>.

2210 Aylward, Frank O., Kristin E. Burnum, Jarrod J. Scott, Garret Suen, Susannah G. Tringe,  
2211 Sandra M. Adams, Kerrie W. Barry, et al. 2012. 'Metagenomic and Metaproteomic  
2212 Insights into Bacterial Communities in Leaf-Cutter Ant Fungus Gardens'. *The*  
2213 *ISME Journal* 6 (9): 1688–1701. <https://doi.org/10.1038/ismej.2012.10>.

2214 Baldrian, Petr. 2017. 'Microbial Activity and the Dynamics of Ecosystem Processes in  
2215 Forest Soils'. *Current Opinion in Microbiology* 37 (June):128–34.  
2216 <https://doi.org/10.1016/j.mib.2017.06.008>.

2217 Bandick, Anna K., and Richard P. Dick. 1999. 'Field Management Effects on Soil Enzyme  
2218 Activities'. *Soil Biology and Biochemistry* 31 (11): 1471–79.  
2219 [https://doi.org/10.1016/S0038-0717\(99\)00051-6](https://doi.org/10.1016/S0038-0717(99)00051-6).

2220 Bankvall, Maria, Miguel Carda-Diéguez, Alex Mira, Anders Karlsson, Bengt Hasséus,  
2221 Roger Karlsson, and Jairo Robledo-Sierra. 2023. 'Metataxonomic and  
2222 Metaproteomic Profiling of the Oral Microbiome in Oral Lichen Planus - a Pilot  
2223 Study'. *Journal of Oral Microbiology* 15 (1): 2161726.  
2224 <https://doi.org/10.1080/20002297.2022.2161726>.

2225 Barglow, Katherine T., and Benjamin F. Cravatt. 2007. 'Activity-Based Protein Profiling for  
2226 the Functional Annotation of Enzymes'. *Nature Methods* 4 (10): 822–27.  
2227 <https://doi.org/10.1038/nmeth1092>.

2228 Bastida, F., T. Hernández, and C. García. 2014. 'Metaproteomics of Soils from Semiarid  
2229 Environment: Functional and Phylogenetic Information Obtained with Different  
2230 Protein Extraction Methods'. *Journal of Proteomics* 101 (April):31–42.  
2231 <https://doi.org/10.1016/j.jprot.2014.02.006>.

2232 Benndorf, Dirk, Gerd U. Balcke, Hauke Harms, and Martin von Bergen. 2007. 'Functional  
2233 Metaproteome Analysis of Protein Extracts from Contaminated Soil and  
2234 Groundwater'. *The ISME Journal* 1 (3): 224–34.  
2235 <https://doi.org/10.1038/ismej.2007.39>.

2236 Benndorf, Dirk, Carsten Vogt, Nico Jehmlich, Yvonne Schmidt, Henrik Thomas, Gary  
2237 Woffendin, Andrej Shevchenko, Hans-Hermann Richnow, and Martin von Bergen.  
2238 2009. 'Improving Protein Extraction and Separation Methods for Investigating the  
2239 Metaproteome of Anaerobic Benzene Communities within Sediments'.  
2240 *Biodegradation* 20 (6): 737–50. <https://doi.org/10.1007/s10532-009-9261-3>.

2241 Berard, Alicia R., Douglas K. Brubaker, Kenzie Birse, Alana Lamont, Romel D.  
2242 Mackelprang, Laura Noël-Romas, Michelle Perner, et al. 2023. 'Vaginal Epithelial  
2243 Dysfunction Is Mediated by the Microbiome, Metabolome, and mTOR Signaling'.  
2244 *Cell Reports* 42 (5): 112474. <https://doi.org/10.1016/j.celrep.2023.112474>.

2245 Beresford-Jones, Benjamin S., Samuel C. Forster, Mark D. Stares, George Notley, Elisa  
2246 Viciani, Hilary P. Browne, Daniel J. Boehmler, et al. 2022. 'The Mouse  
2247 Gastrointestinal Bacteria Catalogue Enables Translation between the Mouse and  
2248 Human Gut Microbiotas via Functional Mapping'. *Cell Host & Microbe* 30 (1): 124-  
2249 138.e8. <https://doi.org/10.1016/j.chom.2021.12.003>.

2250 Berg, Gabriele, Daria Rybakova, Doreen Fischer, Tomislav Cernava, Marie-Christine  
2251 Champomier Vergès, Trevor Charles, Xiaoyulong Chen, et al. 2020. 'Microbiome  
2252 Definition Re-Visited: Old Concepts and New Challenges'. *Microbiome* 8 (1): 103.  
2253 <https://doi.org/10.1186/s40168-020-00875-0>.

2254 Beyter, Doruk, Miin S. Lin, Yanbao Yu, Rembert Pieper, and Vineet Bafna. 2018.  
2255 'ProteoStorm: An Ultrafast Metaproteomics Database Search Framework'. *Cell*  
2256 *Systems* 7 (4): 463-467.e6. <https://doi.org/10.1016/j.cels.2018.08.009>.

2257 Bielow, Chris, Nils Hoffmann, David Jimenez-Morales, Tim Van Den Bossche, Juan  
2258 Antonio Vizcaíno, David L. Tabb, Wout Bittremieux, and Mathias Walzer. 2024.  
2259 'Communicating Mass Spectrometry Quality Information in mzQC with Python, R,  
2260 and Java'. *Journal of the American Society for Mass Spectrometry* 35 (8): 1875–

2261 82. <https://doi.org/10.1021/jasms.4c00174>.

2262 Bielow, Chris, Guido Mastrobuoni, and Stefan Kempa. 2016. 'Proteomics Quality Control:  
2263 Quality Control Software for MaxQuant Results'. *Journal of Proteome Research* 15  
2264 (3): 777–87. <https://doi.org/10.1021/acs.jproteome.5b00780>.

2265 Bihani, Surbhi, Aryan Gupta, Subina Mehta, Andrew T. Rajczewski, James Johnson,  
2266 Dhanush Borishetty, Timothy J. Griffin, Sanjeeva Srivastava, and Pratik D. Jagtap.  
2267 2023. 'Metaproteomic Analysis of Nasopharyngeal Swab Samples to Identify  
2268 Microbial Peptides in COVID-19 Patients'. *Journal of Proteome Research* 22 (8):  
2269 2608–19. <https://doi.org/10.1021/acs.jproteome.3c00040>.

2270 Birse, Kenzie D., Kateryna Kratzer, Christina Farr Zuend, Sarah Mutch, Laura Noël-  
2271 Romas, Alana Lamont, Max Abou, et al. 2020. 'The Neovaginal Microbiome of  
2272 Transgender Women Post-Gender Reassignment Surgery'. *Microbiome* 8 (1): 61.  
2273 <https://doi.org/10.1186/s40168-020-00804-1>.

2274 Blakeley-Ruiz, J. Alfredo, Alison R. Erickson, Brandi L. Cantarel, Weili Xiong, Rachel  
2275 Adams, Janet K. Jansson, Claire M. Fraser, and Robert L. Hettich. 2019.  
2276 'Metaproteomics Reveals Persistent and Phylum-Redundant Metabolic Functional  
2277 Stability in Adult Human Gut Microbiomes of Crohn's Remission Patients despite  
2278 Temporal Variations in Microbial Taxa, Genomes, and Proteomes'. *Microbiome* 7  
2279 (1): 18. <https://doi.org/10.1186/s40168-019-0631-8>.

2280 Blakeley-Ruiz, J. Alfredo, and Manuel Kleiner. 2022. 'Considerations for Constructing a  
2281 Protein Sequence Database for Metaproteomics'. *Computational and Structural  
2282 Biotechnology Journal* 20 (January):937–52.  
2283 <https://doi.org/10.1016/j.csbj.2022.01.018>.

2284 Blakeley-Ruiz, J Alfredo, Carlee S McClintock, Him K Shrestha, Suresh Poudel, Zamin K  
2285 Yang, Richard J Giannone, James J Choo, et al. 2022. 'Morphine and High-Fat  
2286 Diet Differentially Alter the Gut Microbiota Composition and Metabolic Function in  
2287 Lean versus Obese Mice'. *ISME Communications* 2 (1): 66.  
2288 <https://doi.org/10.1038/s43705-022-00131-6>.

2289 Blank, Clemens, Caleb Easterly, Bjoern Gruening, James Johnson, Carolin A. Kolmeder,  
2290 Praveen Kumar, Damon May, et al. 2018. 'Disseminating Metaproteomic  
2291 Informatics Capabilities and Knowledge Using the Galaxy-P Framework'.  
2292 *Proteomes* 6 (1): 7. <https://doi.org/10.3390/proteomes6010007>.

2293 Brooks, Brandon, Ryan S. Mueller, Jacque C. Young, Michael J. Morowitz, Robert L.  
2294 Hettich, and Jillian F. Banfield. 2015. 'Strain-Resolved Microbial Community  
2295 Proteomics Reveals Simultaneous Aerobic and Anaerobic Function during  
2296 Gastrointestinal Tract Colonization of a Preterm Infant'. *Frontiers in Microbiology* 6  
2297 (July). <https://doi.org/10.3389/fmicb.2015.00654>.

2298 Bruderer, Roland, Oliver M. Bernhardt, Tejas Gandhi, Yue Xuan, Julia Sondermann,  
2299 Manuela Schmidt, David Gomez-Varela, and Lukas Reiter. 2017. 'Optimization of  
2300 Experimental Parameters in Data-Independent Mass Spectrometry Significantly  
2301 Increases Depth and Reproducibility of Results \*'. *Molecular & Cellular Proteomics*  
2302 16 (12): 2296–2309. <https://doi.org/10.1074/mcp.RA117.000314>.

2303 Bubis, Julia A., Lev I. Levitsky, Mark V. Ivanov, Irina A. Tarasova, and Mikhail V.  
2304 Gorshkov. 2017. 'Comparative Evaluation of Label-Free Quantification Methods  
2305 for Shotgun Proteomics'. *Rapid Communications in Mass Spectrometry* 31 (7):  
2306 606–12. <https://doi.org/10.1002/rcm.7829>.

2307 Burns, Andrew P., Ya-Qin Zhang, Tuan Xu, Zhengxi Wei, Qin Yao, Yuhong Fang, Valeriu  
2308 Cebotaru, et al. 2021. 'A Universal and High-Throughput Proteomics Sample  
2309 Preparation Platform'. *Analytical Chemistry* 93 (24): 8423–31.  
2310 <https://doi.org/10.1021/acs.analchem.1c00265>.

2311 Burnum-Johnson, Kristin E., Jennifer E. Kyle, Amie J. Einfeld, Cameron P. Casey, Kelly  
2312 G. Stratton, Juan F. Gonzalez, Fabien Habyarimana, et al. 2017. 'MPLEx: A  
2313 Method for Simultaneous Pathogen Inactivation and Extraction of Samples for  
2314 Multi-Omics Profiling'. *Analyst* 142 (3): 442–48.  
2315 <https://doi.org/10.1039/C6AN02486F>.

- 2316 C. Silva, Ana S, Robbin Bouwmeester, Lennart Martens, and Sven Degroeve. 2019.  
2317 'Accurate Peptide Fragmentation Predictions Allow Data Driven Approaches to  
2318 Replace and Improve upon Proteomics Search Engine Scoring Functions'.  
2319 *Bioinformatics* 35 (24): 5243–48. <https://doi.org/10.1093/bioinformatics/btz383>.
- 2320 Caglar, Mehmet U., John R. Houser, Craig S. Barnhart, Daniel R. Boutz, Sean M. Carroll,  
2321 Aurko Dasgupta, Walter F. Lenoir, et al. 2017. 'The E. Coli Molecular Phenotype  
2322 under Different Growth Conditions'. *Scientific Reports* 7 (1): 45303.  
2323 <https://doi.org/10.1038/srep45303>.
- 2324 Cai, Xue, Zhangzhi Xue, Chunlong Wu, Rui Sun, Liuji Qian, Liang Yue, Weigang Ge, et  
2325 al. 2022. 'High-Throughput Proteomic Sample Preparation Using Pressure Cycling  
2326 Technology'. *Nature Protocols* 17 (10): 2307–25. <https://doi.org/10.1038/s41596-022-00727-1>.
- 2328 Cantalapiedra, Carlos P, Ana Hernández-Plaza, Ivica Letunic, Peer Bork, and Jaime  
2329 Huerta-Cepas. 2021. 'eggNOG-Mapper v2: Functional Annotation, Orthology  
2330 Assignments, and Domain Prediction at the Metagenomic Scale'. *Molecular  
2331 Biology and Evolution* 38 (12): 5825–29. <https://doi.org/10.1093/molbev/msab293>.
- 2332 Charlier, Philippe, Virginie Bourdin, Didier N'Dah, Mélodie Kielbasa, Olivier Pible, and  
2333 Jean Armengaud. 2024. 'Metaproteomic Analysis of King Ghezo Tomb Wall  
2334 (Abomey, Benin) Confirms 19th Century Voodoo Sacrifices'. *Proteomics* 24 (16):  
2335 e2400048. <https://doi.org/10.1002/pmic.202400048>.
- 2336 Chaumeil, Pierre-Alain, Aaron J Mussig, Philip Hugenholtz, and Donovan H Parks. 2020.  
2337 'GTDB-Tk: A Toolkit to Classify Genomes with the Genome Taxonomy Database'.  
2338 *Bioinformatics* 36 (6): 1925–27. <https://doi.org/10.1093/bioinformatics/btz848>.
- 2339 Chen, Jiahui, Yingying Sun, Jie Li, Mengge Lyu, Li Yuan, Jiancheng Sun, Shangqi Chen,  
2340 et al. 2024. 'In-Depth Metaproteomics Analysis of Tongue Coating for Gastric  
2341 Cancer: A Multicenter Diagnostic Research Study'. *Microbiome* 12 (1): 6.  
2342 <https://doi.org/10.1186/s40168-023-01730-8>.
- 2343 Cheng, Kai, Zhibin Ning, Leyuan Li, Xu Zhang, Joeselle M. Serrana, Janice Mayne, and  
2344 Daniel Figeys. 2023. 'MetaLab-MAG: A Metaproteomic Data Analysis Platform for  
2345 Genome-Level Characterization of Microbiomes from the Metagenome-Assembled  
2346 Genomes Database'. *Journal of Proteome Research* 22 (2): 387–98.  
2347 <https://doi.org/10.1021/acs.jproteome.2c00554>.
- 2348 Cheng, Kai, Zhibin Ning, Xu Zhang, Haonan Duan, Janice Mayne, and Daniel Figeys.  
2349 2024. 'MetaLab Platform Enables Comprehensive DDA and DIA Metaproteomics  
2350 Analysis'. *bioRxiv*, January, 2024.09.27.615406.  
2351 <https://doi.org/10.1101/2024.09.27.615406>.
- 2352 Cheng, Kai, Zhibin Ning, Xu Zhang, Leyuan Li, Bo Liao, Janice Mayne, Alain Stintzi, and  
2353 Daniel Figeys. 2017. 'MetaLab: An Automated Pipeline for Metaproteomic Data  
2354 Analysis'. *Microbiome* 5 (1): 157. <https://doi.org/10.1186/s40168-017-0375-2>.
- 2355 Cheng, Kai, Zhibin Ning, Xu Zhang, Janice Mayne, and Daniel Figeys. 2018. 'Separation  
2356 and Characterization of Human Microbiomes by Metaproteomics'. *TrAC Trends in  
2357 Analytical Chemistry* 108 (November):221–30.  
2358 <https://doi.org/10.1016/j.trac.2018.09.006>.
- 2359 Chourey, Karuna, Janet Jansson, Nathan VerBerkmoes, Manesh Shah, Krystle L.  
2360 Chavarria, Lauren M. Tom, Eoin L. Brodie, and Robert L. Hettich. 2010. 'Direct  
2361 Cellular Lysis/Protein Extraction Protocol for Soil Metaproteomics'. *Journal of  
2362 Proteome Research* 9 (12): 6615–22. <https://doi.org/10.1021/pr100787q>.
- 2363 Claeys, Tine, Tim Van Den Bossche, Yasset Perez-Riverol, Kris Gevaert, Juan Antonio  
2364 Vizcaíno, and Lennart Martens. 2023. 'lesSDRF Is More: Maximizing the Value of  
2365 Proteomics Data through Streamlined Metadata Annotation'. *Nature  
2366 Communications* 14 (1): 6743. <https://doi.org/10.1038/s41467-023-42543-5>.
- 2367 Combe, Colin W., Lars Kolbowski, Lutz Fischer, Ville Koskinen, Joshua Klein, Alexander  
2368 Leitner, Andrew R. Jones, Juan Antonio Vizcaíno, and Juri Rappsilber. 2024.  
2369 'mzIdentML 1.3.0 – Essential Progress on the Support of Crosslinking and Other  
2370 Identifications Based on Multiple Spectra'. *PROTEOMICS* 24 (17): 2300385.

- 2371 <https://doi.org/10.1002/pmic.202300385>.
- 2372 Cortay, Jean-Claude, Corinne Rieul, Françoise Bleicher, Mustapha Dadssi, and Alain J.  
2373 Cozzone. 1988. 'Evidence of Protein Kinase Activity and Characterization of  
2374 Substrate Proteins in Escherichia Coli'. In *Advances in Post-Translational*  
2375 *Modifications of Proteins and Aging*, edited by Vincenzo Zappia, Patrizia Galletti,  
2376 Raffaele Porta, and Finn Wold, 467–74. Boston, MA: Springer US.  
2377 [https://doi.org/10.1007/978-1-4684-9042-8\\_39](https://doi.org/10.1007/978-1-4684-9042-8_39).
- 2378 Cox, Jürgen, Marco Y. Hein, Christian A. Luber, Igor Paron, Nagarjuna Nagaraj, and  
2379 Matthias Mann. 2014. 'Accurate Proteome-Wide Label-Free Quantification by  
2380 Delayed Normalization and Maximal Peptide Ratio Extraction, Termed MaxLFQ \*'.  
2381 *Molecular & Cellular Proteomics* 13 (9): 2513–26.  
2382 <https://doi.org/10.1074/mcp.M113.031591>.
- 2383 Cox, Jürgen, and Matthias Mann. 2008. 'MaxQuant Enables High Peptide Identification  
2384 Rates, Individualized p.p.b.-Range Mass Accuracies and Proteome-Wide Protein  
2385 Quantification'. *Nature Biotechnology* 26 (12): 1367–72.  
2386 <https://doi.org/10.1038/nbt.1511>.
- 2387 Cox, Jürgen, Nadin Neuhauser, Annette Michalski, Richard A. Scheltema, Jesper V.  
2388 Olsen, and Matthias Mann. 2011. 'Andromeda: A Peptide Search Engine  
2389 Integrated into the MaxQuant Environment'. *Journal of Proteome Research* 10 (4):  
2390 1794–1805. <https://doi.org/10.1021/pr101065j>.
- 2391 Craig, Robertson, and Ronald C. Beavis. 2004. 'TANDEM: Matching Proteins with  
2392 Tandem Mass Spectra'. *Bioinformatics* 20 (9): 1466–67.  
2393 <https://doi.org/10.1093/bioinformatics/bth092>.
- 2394 Cravatt, Benjamin F., Aaron T. Wright, and John W. Kozarich. 2008. 'Activity-Based  
2395 Protein Profiling: From Enzyme Chemistry to Proteomic Chemistry'. *Annual*  
2396 *Review of Biochemistry* 77 (Volume 77, 2008): 383–414.  
2397 <https://doi.org/10.1146/annurev.biochem.75.101304.124125>.
- 2398 Creskey, Marybeth, Leyuan Li, Zhibin Ning, Emily EF Fekete, Janice Mayne, Krystal  
2399 Walker, Anna Ampaw, Robert Ben, Xu Zhang, and Daniel Figeys. 2023. 'An  
2400 Economic and Robust TMT Labeling Approach for High Throughput Proteomic  
2401 and Metaproteomic Analysis'. *PROTEOMICS* 23 (21–22): 2200116.  
2402 <https://doi.org/10.1002/pmic.202200116>.
- 2403 Dai, Chengxin, Anja Füllgrabe, Julianus Pfeuffer, Elizaveta M. Solovyeva, Jingwen Deng,  
2404 Pablo Moreno, Selvakumar Kamatchinathan, et al. 2021. 'A Proteomics Sample  
2405 Metadata Representation for Multiomics Integration and Big Data Analysis'. *Nature*  
2406 *Communications* 12 (1): 5854. <https://doi.org/10.1038/s41467-021-26111-3>.
- 2407 Degroeve, Sven, Davy Maddelein, and Lennart Martens. 2015. 'MS2PIP Prediction  
2408 Server: Compute and Visualize MS2 Peak Intensity Predictions for CID and HCD  
2409 Fragmentation'. *Nucleic Acids Research* 43 (W1): W326–30.  
2410 <https://doi.org/10.1093/nar/gkv542>.
- 2411 Delgado-Diaz, David Jose, Brianna Jesaveluk, Joshua A. Hayward, David Tyssen,  
2412 Arghavan Alisoltani, Matthys Potgieter, Liam Bell, et al. 2022. 'Lactic Acid from  
2413 Vaginal Microbiota Enhances Cervicovaginal Epithelial Barrier Integrity by  
2414 Promoting Tight Junction Protein Expression'. *Microbiome* 10 (1): 141.  
2415 <https://doi.org/10.1186/s40168-022-01337-5>.
- 2416 Delogu, F., B. J. Kunath, P. N. Evans, M. Ø. Arntzen, T. R. Hvidsten, and P. B. Pope.  
2417 2020. 'Integration of Absolute Multi-Omics Reveals Dynamic Protein-to-RNA  
2418 Ratios and Metabolic Interplay within Mixed-Domain Microbiomes'. *Nature*  
2419 *Communications* 11 (1): 4708. <https://doi.org/10.1038/s41467-020-18543-0>.
- 2420 Delogu, Francesco, Benoit J. Kunath, Pedro M. Queirós, Rashi Halder, Laura A. Lebrun,  
2421 Phillip B. Pope, Patrick May, Stefanie Widder, Emilie E. L. Muller, and Paul  
2422 Wilmes. 2024. 'Forecasting the Dynamics of a Complex Microbial Community  
2423 Using Integrated Meta-Omics'. *Nature Ecology & Evolution* 8 (1): 32–44.  
2424 <https://doi.org/10.1038/s41559-023-02241-3>.
- 2425 Demichev, Vadim, Christoph B. Messner, Spyros I. Vernardis, Kathryn S. Lilley, and

2426 Markus Ralser. 2020. 'DIA-NN: Neural Networks and Interference Correction  
 2427 Enable Deep Proteome Coverage in High Throughput'. *Nature Methods* 17 (1):  
 2428 41–44. <https://doi.org/10.1038/s41592-019-0638-x>.

2429 Demichev, Vadim, Lukasz Szyrwił, Fengchao Yu, Guo Ci Teo, George Rosenberger,  
 2430 Agathe Niewienda, Daniela Ludwig, et al. 2022. 'Dia-PASEF Data Analysis Using  
 2431 FragPipe and DIA-NN for Deep Proteomics of Low Sample Amounts'. *Nature*  
 2432 *Communications* 13 (1): 3944. <https://doi.org/10.1038/s41467-022-31492-0>.

2433 Deusch, Simon, Amélia Camarinha-Silva, Jürgen Conrad, Uwe Beifuss, Markus  
 2434 Rodehutsord, and Jana Seifert. 2017. 'A Structural and Functional Elucidation of  
 2435 the Rumen Microbiome Influenced by Various Diets and Microenvironments'.  
 2436 *Frontiers in Microbiology* 8 (August). <https://doi.org/10.3389/fmicb.2017.01605>.

2437 Deutsch, Eric W, Nuno Bandeira, Yasset Perez-Riverol, Vagisha Sharma, Jeremy J  
 2438 Carver, Luis Mendoza, Deepti J Kundu, et al. 2023. 'The ProteomeXchange  
 2439 Consortium at 10 Years: 2023 Update'. *Nucleic Acids Research* 51 (D1): D1539–  
 2440 48. <https://doi.org/10.1093/nar/gkac1040>.

2441 Deutsch, Eric W., Yasset Perez-Riverol, Jeremy Carver, Shin Kawano, Luis Mendoza,  
 2442 Tim Van Den Bossche, Ralf Gabriels, et al. 2021. 'Universal Spectrum Identifier for  
 2443 Mass Spectra'. *Nature Methods* 18 (7): 768–70. <https://doi.org/10.1038/s41592-021-01184-6>.

2445 Deutsch, Eric W., Juan Antonio Vizcaíno, Andrew R. Jones, Pierre-Alain Binz, Henry Lam,  
 2446 Joshua Klein, Wout Bittremieux, et al. 2023. 'Proteomics Standards Initiative at  
 2447 Twenty Years: Current Activities and Future Work'. *Journal of Proteome Research*  
 2448 22 (2): 287–301. <https://doi.org/10.1021/acs.jproteome.2c00637>.

2449 Dhabaria, Avantika, Paolo Cifani, Casie Reed, Hanno Steen, and Alex Kentsis. 2015. 'A  
 2450 High-Efficiency Cellular Extraction System for Biological Proteomics'. *Journal of*  
 2451 *Proteome Research* 14 (8): 3403–8.  
 2452 <https://doi.org/10.1021/acs.jproteome.5b00547>.

2453 Diz, Angel P., and Paula Sánchez-Marín. 2021. 'A Primer and Guidelines for Shotgun  
 2454 Proteomic Shotgun Proteomics Analysis in Non-Model Organisms'. In *Shotgun*  
 2455 *Proteomics: Methods and Protocols*, edited by Mónica Carrera and Jesús Mateos,  
 2456 77–102. New York, NY: Springer US. [https://doi.org/10.1007/978-1-0716-1178-4\\_6](https://doi.org/10.1007/978-1-0716-1178-4_6).

2458 Do, Katherine, Subina Mehta, Reid Wagner, Dechen Bhuming, Andrew T. Rajczewski,  
 2459 Amy P. N. Skubitz, James E. Johnson, Timothy J. Griffin, and Pratik D. Jagtap.  
 2460 2024. 'A Novel Clinical Metaproteomics Workflow Enables Bioinformatic Analysis  
 2461 of Host-Microbe Dynamics in Disease'. *mSphere* 9 (6): e00793-23.  
 2462 <https://doi.org/10.1128/msphere.00793-23>.

2463 Drula, Elodie, Marie-Line Garron, Suzan Dogan, Vincent Lombard, Bernard Henrissat,  
 2464 and Nicolas Terrapon. 2022. 'The Carbohydrate-Active Enzyme Database:  
 2465 Functions and Literature'. *Nucleic Acids Research* 50 (D1): D571–77.  
 2466 <https://doi.org/10.1093/nar/gkab1045>.

2467 Duan, Haonan, Kai Cheng, Zhibin Ning, Leyuan Li, Janice Mayne, Zhongzhi Sun, and  
 2468 Daniel Figeys. 2022. 'Assessing the Dark Field of Metaproteome'. *Analytical*  
 2469 *Chemistry* 94 (45): 15648–54. <https://doi.org/10.1021/acs.analchem.2c02452>.

2470 Duan, Haonan, Xu Zhang, and Daniel Figeys. 2023. 'An Emerging Field: Post-  
 2471 Translational Modification in Microbiome'. *PROTEOMICS* 23 (3–4): 2100389.  
 2472 <https://doi.org/10.1002/pmic.202100389>.

2473 Duchovni, Lirit, Genrieta Shmunis, and Lior Lobel. 2024. 'Posttranslational Modifications:  
 2474 An Emerging Functional Layer of Diet-Host-Microbe Interactions'. *mBio* 0 (0):  
 2475 e02387-24. <https://doi.org/10.1128/mbio.02387-24>.

2476 Dumas, Thibaut, Roxana Martinez Pinna, Clément Lozano, Sonja Radau, Olivier Pible,  
 2477 Lucia Grenga, and Jean Armengaud. 2024. 'The Astounding Exhaustiveness and  
 2478 Speed of the Astral Mass Analyzer for Highly Complex Samples Is a Quantum  
 2479 Leap in the Functional Analysis of Microbiomes'. *Microbiome* 12 (1): 46.  
 2480 <https://doi.org/10.1186/s40168-024-01766-4>.



2481 Easterly, Caleb W., Ray Sajulga, Subina Mehta, James Johnson, Praveen Kumar, Shane  
2482 Hubler, Bart Mesuere, Joel Rudney, Timothy J. Griffin, and Pratik D. Jagtap. 2019.  
2483 'metaQuantome: An Integrated, Quantitative Metaproteomics Approach Reveals  
2484 Connections Between Taxonomy and Protein Function in Complex Microbiomes \*'.  
2485 *Molecular & Cellular Proteomics* 18 (8): S82–91.  
2486 <https://doi.org/10.1074/mcp.RA118.001240>.  
2487 Elias, Joshua E., and Steven P. Gygi. 2007. 'Target-Decoy Search Strategy for Increased  
2488 Confidence in Large-Scale Protein Identifications by Mass Spectrometry'. *Nature*  
2489 *Methods* 4 (3): 207–14. <https://doi.org/10.1038/nmeth1019>.  
2490 Eng, Jimmy K., Tahmina A. Jahan, and Michael R. Hoopmann. 2013. 'Comet: An Open-  
2491 Source MS/MS Sequence Database Search Tool'. *PROTEOMICS* 13 (1): 22–24.  
2492 <https://doi.org/10.1002/pmic.201200439>.  
2493 Ezzeldin, Shahd, Aya El-Wazir, Shymaa Enany, Abdelrahman Muhammad, Dina Johar,  
2494 Aya Osama, Eman Ahmed, Hassan Shikshaky, and Sameh Magdeldin. 2019.  
2495 'Current Understanding of Human Metaproteome Association and Modulation'.  
2496 *Journal of Proteome Research* 18 (10): 3539–54.  
2497 <https://doi.org/10.1021/acs.jproteome.9b00301>.  
2498 Ferdous, Tahsin, Lai Jiang, Irina Dinu, Julie Groizeleau, Anita L. Kozyrskyj, Celia M. T.  
2499 Greenwood, and Marie-Claire Arrieta. 2022. 'The Rise to Power of the  
2500 Microbiome: Power and Sample Size Calculation for Microbiome Studies'.  
2501 *Mucosal Immunology* 15 (6): 1060–70. [https://doi.org/10.1038/s41385-022-00548-](https://doi.org/10.1038/s41385-022-00548-1)  
2502 [1](https://doi.org/10.1038/s41385-022-00548-1).  
2503 Fernández-Costa, Carolina, Salvador Martínez-Bartolomé, Daniel B. McClatchy, Anthony  
2504 J. Saviola, Nam-Kyung Yu, and John R. III Yates. 2020. 'Impact of the  
2505 Identification Strategy on the Reproducibility of the DDA and DIA Results'. *Journal*  
2506 *of Proteome Research* 19 (8): 3153–61.  
2507 <https://doi.org/10.1021/acs.jproteome.0c00153>.  
2508 Florens, Laurence, Michael J. Carozza, Selene K. Swanson, Marjorie Fournier, Michael K.  
2509 Coleman, Jerry L. Workman, and Michael P. Washburn. 2006. 'Analyzing  
2510 Chromatin Remodeling Complexes Using Shotgun Proteomics and Normalized  
2511 Spectral Abundance Factors'. *Methods, Chromatin and Transcriptional Regulation*,  
2512 40 (4): 303–11. <https://doi.org/10.1016/j.ymeth.2006.07.028>.  
2513 Frank, Ari, and Pavel Pevzner. 2005. 'PepNovo: De Novo Peptide Sequencing via  
2514 Probabilistic Network Modeling'. *Analytical Chemistry* 77 (4): 964–73.  
2515 <https://doi.org/10.1021/ac048788h>.  
2516 Fu, Qin, Christopher I. Murray, Oleg A. Karpov, and Jennifer E. Van Eyk. 2023.  
2517 'Automated Proteomic Sample Preparation: The Key Component for High  
2518 Throughput and Quantitative Mass Spectrometry Analysis'. *Mass Spectrometry*  
2519 *Reviews* 42 (2): e21750. <https://doi.org/10.1002/mas.21750>.  
2520 Geer, Lewis Y., Sanford P. Markey, Jeffrey A. Kowalak, Lukas Wagner, Ming Xu, Dawn  
2521 M. Maynard, Xiaoyu Yang, Wenyao Shi, and Stephen H. Bryant. 2004. 'Open  
2522 Mass Spectrometry Search Algorithm'. *Journal of Proteome Research* 3 (5): 958–  
2523 64. <https://doi.org/10.1021/pr0499491>.  
2524 The Gene Ontology Consortium. 2019. 'The Gene Ontology Resource: 20 Years and Still  
2525 GOing Strong'. *Nucleic Acids Research* 47 (D1): D330–38.  
2526 <https://doi.org/10.1093/nar/gky1055>.  
2527 Gessulat, Siegfried, Tobias Schmidt, Daniel Paul Zolg, Patroklos Samaras, Karsten  
2528 Schnatbaum, Johannes Zerweck, Tobias Knaute, et al. 2019. 'Prosit: Proteome-  
2529 Wide Prediction of Peptide Tandem Mass Spectra by Deep Learning'. *Nature*  
2530 *Methods* 16 (6): 509–18. <https://doi.org/10.1038/s41592-019-0426-7>.  
2531 Giagnoni, L., F. Magherini, L. Landi, S. Taghavi, A. Modesti, L. Bini, P. Nannipieri, D. Van  
2532 der Ielie, and G. Renella. 2011. 'Extraction of Microbial Proteome from Soil:  
2533 Potential and Limitations Assessed through a Model Study'. *European Journal of*  
2534 *Soil Science* 62 (1): 74–81. <https://doi.org/10.1111/j.1365-2389.2010.01322.x>.  
2535 Gómez-Varela, David, Feng Xian, Sabrina Grundtner, Julia Regina Sondermann,

2536 Giacomo Carta, and Manuela Schmidt. 2023. 'Increasing Taxonomic and  
2537 Functional Characterization of Host-Microbiome Interactions by DIA-PASEF  
2538 Metaproteomics'. *Frontiers in Microbiology* 14 (October).  
2539 <https://doi.org/10.3389/fmicb.2023.1258703>.

2540 Gonzalez, Carlos G., Hannah C. Wastyk, Madeline Topf, Christopher D. Gardner, Justin  
2541 L. Sonnenburg, and Joshua E. Elias. 2020. 'High-Throughput Stool  
2542 Metaproteomics: Method and Application to Human Specimens'. *mSystems* 5 (3):  
2543 10.1128/msystems.00200-20. <https://doi.org/10.1128/msystems.00200-20>.

2544 Graf, Alexander C., Johanna Striesow, Jan Pané-Farré, Thomas Sura, Martina Wurster,  
2545 Michael Lalk, Dietmar H. Pieper, Dörte Becher, Barbara C. Kahl, and Katharina  
2546 Riedel. 2021. 'An Innovative Protocol for Metaproteomic Analyses of Microbial  
2547 Pathogens in Cystic Fibrosis Sputum'. *Frontiers in Cellular and Infection  
2548 Microbiology* 11 (August). <https://doi.org/10.3389/fcimb.2021.724569>.

2549 Grenga, Lucia, Olivier Pible, Guylaine Miotello, Karen Culotta, Sylvie Ruat, Marie-Anne  
2550 Roncato, Fabienne Gas, et al. 2022. 'Taxonomical and Functional Changes in  
2551 COVID-19 Faecal Microbiome Could Be Related to SARS-CoV-2 Faecal Load'.  
2552 *Environmental Microbiology* 24 (9): 4299–4316. <https://doi.org/10.1111/1462-2920.16028>.

2554 Gu, Kongxin, Haruka Kumabe, Takumi Yamamoto, Naoto Tashiro, Takeshi Masuda,  
2555 Shingo Ito, and Sumio Ohtsuki. 2024. 'Improving Proteomic Identification Using  
2556 Narrow Isolation Windows with Zeno SWATH Data-Independent Acquisition'.  
2557 *Journal of Proteome Research* 23 (8): 3484–95.  
2558 <https://doi.org/10.1021/acs.jproteome.4c00149>.

2559 Guo, Xuan, Zhou Li, Qiuming Yao, Ryan S Mueller, Jimmy K Eng, David L Tabb, William  
2560 Judson Hervey IV, and Chongle Pan. 2018. 'Sipros Ensemble Improves Database  
2561 Searching and Filtering for Complex Metaproteomics'. *Bioinformatics* 34 (5): 795–  
2562 802. <https://doi.org/10.1093/bioinformatics/btx601>.

2563 Gurbich, Tatiana A., Alexandre Almeida, Martin Beracochea, Tony Burdett, Josephine  
2564 Burgin, Guy Cochrane, Shriya Raj, et al. 2023. 'MGnify Genomes: A Resource for  
2565 Biome-Specific Microbial Genome Catalogues'. *Journal of Molecular Biology* 435  
2566 (14): 168016. <https://doi.org/10.1016/j.jmb.2023.168016>.

2567 Haange, Sven-Bastiaan, Nico Jehmlich, Maximilian Hoffmann, Klaus Weber, Jörg  
2568 Lehmann, Martin von Bergen, and Ulla Slanina. 2019. 'Disease Development Is  
2569 Accompanied by Changes in Bacterial Protein Abundance and Functions in a  
2570 Refined Model of Dextran Sulfate Sodium (DSS)-Induced Colitis'. *Journal of  
2571 Proteome Research* 18 (4): 1774–86.  
2572 <https://doi.org/10.1021/acs.jproteome.8b00974>.

2573 Habeck, Tanja, Kyle A. Brown, Benjamin Des Soye, Carter Lantz, Mowei Zhou, Novera  
2574 Alam, Md Amin Hossain, et al. 2024. 'Top-down Mass Spectrometry of Native  
2575 Proteoforms and Their Complexes: A Community Study'. *Nature Methods*, May,  
2576 1–9. <https://doi.org/10.1038/s41592-024-02279-6>.

2577 Han, Lin, and Pamela V. Chang. 2023. 'Activity-Based Protein Profiling in Microbes and  
2578 the Gut Microbiome'. *Current Opinion in Chemical Biology* 76 (October):102351.  
2579 <https://doi.org/10.1016/j.cbpa.2023.102351>.

2580 Hansmeier, N., S Sharma, and TC Chao. 2022. 'Protein Purification and Digestion  
2581 Methods for Bacterial Proteomic Analyses.' In *Geddes-McAlister, J (Eds).  
2582 Proteomics in Systems Biology. Methods in Molecular Biology*. Vol. 2456. New  
2583 York, NY: Humana.

2584 Hao, Chunlin, Joshua E. Elias, Patrick K. H. Lee, and Henry Lam. 2023. 'metaSpectraST:  
2585 An Unsupervised and Database-Independent Analysis Workflow for  
2586 Metaproteomic MS/MS Data Using Spectrum Clustering'. *Microbiome* 11 (1): 176.  
2587 <https://doi.org/10.1186/s40168-023-01602-1>.

2588 Heintz-Buschart, Anna, and Paul Wilmes. 2018. 'Human Gut Microbiome: Function  
2589 Matters'. *Trends in Microbiology* 26 (7): 563–74.  
2590 <https://doi.org/10.1016/j.tim.2017.11.002>.

- 2591 Henry, Céline, Ariane Bassignani, Magali Berland, Olivier Langella, Harry Sokol, and  
2592 Catherine Juste. 2022. 'Modern Metaproteomics: A Unique Tool to Characterize  
2593 the Active Microbiome in Health and Diseases, and Pave the Road towards New  
2594 Biomarkers—Example of Crohn's Disease and Ulcerative Colitis Flare-Ups'. *Cells*  
2595 11 (8): 1340. <https://doi.org/10.3390/cells11081340>.
- 2596 Herbst, Florian-Alexander, Vanessa Lünsmann, Henrik Kjeldal, Nico Jehmlich, Andreas  
2597 Tholey, Martin von Bergen, Jeppe Lund Nielsen, Robert L. Hettich, Jana Seifert,  
2598 and Per Halkjær Nielsen. 2016. 'Enhancing Metaproteomics—The Value of  
2599 Models and Defined Environmental Microbial Systems'. *PROTEOMICS* 16 (5):  
2600 783–98. <https://doi.org/10.1002/pmic.201500305>.
- 2601 Hernández-Plaza, Ana, Damian Szklarczyk, Jorge Botas, Carlos P Cantalapiedra,  
2602 Joaquín Giner-Lamia, Daniel R Mende, Rebecca Kirsch, et al. 2023. 'eggNOG 6.0:  
2603 Enabling Comparative Genomics across 12 535 Organisms'. *Nucleic Acids*  
2604 *Research* 51 (D1): D389–94. <https://doi.org/10.1093/nar/gkac1022>.
- 2605 Hettich, Robert L., Chongle Pan, Karuna Chourey, and Richard J. Giannone. 2013.  
2606 'Metaproteomics: Harnessing the Power of High Performance Mass Spectrometry  
2607 to Identify the Suite of Proteins That Control Metabolic Activities in Microbial  
2608 Communities'. *Analytical Chemistry* 85 (9): 4203–14.  
2609 <https://doi.org/10.1021/ac303053e>.
- 2610 Heyer, Robert, Patrick Hellwig, Irena Maus, Daniel Walke, Andreas Schlüter, Julia Hassa,  
2611 Alexander Sczyrba, et al. 2024. 'Breakdown of Hardly Degradable Carbohydrates  
2612 (Lignocellulose) in a Two-Stage Anaerobic Digestion Plant Is Favored in the Main  
2613 Fermenter'. *Water Research* 250 (February):121020.  
2614 <https://doi.org/10.1016/j.watres.2023.121020>.
- 2615 Heyer, Robert, Kay Schallert, Anja Büdel, Roman Zoun, Sebastian Dorl, Alexander  
2616 Behne, Fabian Kohrs, et al. 2019. 'A Robust and Universal Metaproteomics  
2617 Workflow for Research Studies and Routine Diagnostics Within 24 h Using Phenol  
2618 Extraction, FASP Digest, and the MetaProteomeAnalyzer'. *Frontiers in*  
2619 *Microbiology* 10 (August). <https://doi.org/10.3389/fmicb.2019.01883>.
- 2620 Hiltemann, Saskia, Helena Rasche, Simon Gladman, Hans-Rudolf Hotz, Delphine  
2621 Larivière, Daniel Blankenberg, Pratik D. Jagtap, et al. 2023. 'Galaxy Training: A  
2622 Powerful Framework for Teaching!' *PLOS Computational Biology* 19 (1):  
2623 e1010752. <https://doi.org/10.1371/journal.pcbi.1010752>.
- 2624 Hinzke, Tjorven, Manuel Kleiner, and Stephanie Markert. 2018. 'Centrifugation-Based  
2625 Enrichment of Bacterial Cell Populations for Metaproteomic Studies on Bacteria–  
2626 Invertebrate Symbioses'. In *Microbial Proteomics: Methods and Protocols*, edited  
2627 by Dörte Becher, 319–34. New York, NY: Springer. [https://doi.org/10.1007/978-1-4939-8695-8\\_22](https://doi.org/10.1007/978-1-4939-8695-8_22).
- 2629 Hinzke, Tjorven, Manuel Kleiner, Mareike Meister, Rabea Schlüter, Christian Hentschker,  
2630 Jan Pané-Farré, Petra Hildebrandt, et al. 2021. 'Bacterial Symbiont  
2631 Subpopulations Have Different Roles in a Deep-Sea Symbiosis'. Edited by Gisela  
2632 Storz and Victoria Orphan. *eLife* 10 (January):e58371.  
2633 <https://doi.org/10.7554/eLife.58371>.
- 2634 Hinzke, Tjorven, Angela Kouris, Rebecca-Ayme Hughes, Marc Strous, and Manuel  
2635 Kleiner. 2019. 'More Is Not Always Better: Evaluation of 1D and 2D-LC-MS/MS  
2636 Methods for Metaproteomics'. *Frontiers in Microbiology* 10 (February).  
2637 <https://doi.org/10.3389/fmicb.2019.00238>.
- 2638 Holstein, Tanja, Pieter Verschaffelt, Tim Van den Bossche, Lennart Martens, and Thilo  
2639 Muth. 2024. 'The Peptonizer2000: Graphical Model Based Taxonomic  
2640 Identifications of Metaproteomic Samples'. bioRxiv.  
2641 <https://doi.org/10.1101/2024.05.20.594958>.
- 2642 Hughes, Christopher S, Sophia Foehr, David A Garfield, Eileen E Furlong, Lars M  
2643 Steinmetz, and Jeroen Krijgsveld. 2014. 'Ultrasensitive Proteome Analysis Using  
2644 Paramagnetic Bead Technology'. *Molecular Systems Biology* 10 (10): 757.  
2645 <https://doi.org/10.15252/msb.20145625>.

2646 Hustoft, Hanne, Helle Malerod, Steven Wilson, Léon Reubsaet, Elsa Lundanes, and Tyge  
2647 Greibrokk. 2012. 'A Critical Review of Trypsin Digestion for LC-MS Based  
2648 Proteomics'. In . <https://doi.org/10.13140/2.1.2226.7846>.

2649 Ishikawa, Masaki, Ryo Konno, Daisuke Nakajima, Mari Gotoh, Keiko Fukasawa, Hironori  
2650 Sato, Ren Nakamura, Osamu Ohara, and Yusuke Kawashima. 2022. 'Optimization  
2651 of Ultrafast Proteomics Using an LC-Quadrupole-Orbitrap Mass Spectrometer with  
2652 Data-Independent Acquisition'. *Journal of Proteome Research* 21 (9): 2085–93.  
2653 <https://doi.org/10.1021/acs.jproteome.2c00121>.

2654 Jabbar, Karolina S., Brendan Dolan, Lisbeth Eklund, Catharina Wising, Anna Ermund,  
2655 Åsa Johansson, Hans Törnblom, Magnus Simren, and Gunnar C. Hansson. 2021.  
2656 'Association between Brachyspira and Irritable Bowel Syndrome with Diarrhoea'.  
2657 *Gut* 70 (6): 1117–29. <https://doi.org/10.1136/gutjnl-2020-321466>.

2658 Jagtap, Pratik D., Alan Blakely, Kevin Murray, Shaun Stewart, Joel Kooren, James E.  
2659 Johnson, Nelson L. Rhodus, Joel Rudney, and Timothy J. Griffin. 2015.  
2660 'Metaproteomic Analysis Using the Galaxy Framework'. *PROTEOMICS* 15 (20):  
2661 3553–65. <https://doi.org/10.1002/pmic.201500074>.

2662 Jagtap, Pratik, Jill Goslinga, Joel A. Kooren, Thomas McGowan, Matthew S. Wroblewski,  
2663 Sean L. Seymour, and Timothy J. Griffin. 2013. 'A Two-Step Database Search  
2664 Method Improves Sensitivity in Peptide Sequence Matches for Metaproteomics  
2665 and Proteogenomics Studies'. *PROTEOMICS* 13 (8): 1352–57.  
2666 <https://doi.org/10.1002/pmic.201200352>.

2667 James, P., M. Quadroni, E. Carafoli, and G. Gonnet. 1993. 'Protein Identification by Mass  
2668 Profile Fingerprinting'. *Biochemical and Biophysical Research Communications*  
2669 195 (1): 58–64. <https://doi.org/10.1006/bbrc.1993.2009>.

2670 Jarman, Kristin H., Natalie C. Heller, Sarah C. Jenson, Janine R. Hutchison, Brooke L.  
2671 Deatherage Kaiser, Samuel H. Payne, David S. Wunschel, and Eric D. Merkley.  
2672 2018. 'Proteomics Goes to Court: A Statistical Foundation for Forensic  
2673 Toxin/Organism Identification Using Bottom-Up Proteomics'. *Journal of Proteome  
2674 Research* 17 (9): 3075–85. <https://doi.org/10.1021/acs.jproteome.8b00212>.

2675 Jensen, Marlene, Juliane Wippler, and Manuel Kleiner. 2021. 'Evaluation of RNAlater as a  
2676 Field-Compatible Preservation Method for Metaproteomic Analyses of Bacterium-  
2677 Animal Symbioses'. *Microbiology Spectrum* 9 (2): e01429-21.  
2678 <https://doi.org/10.1128/Spectrum.01429-21>.

2679 Jersie-Christensen, Rosa R., Liam T. Lanigan, David Lyon, Meaghan Mackie, Daniel  
2680 Belstrøm, Christian D. Kelstrup, Anna K. Fotakis, et al. 2018. 'Quantitative  
2681 Metaproteomics of Medieval Dental Calculus Reveals Individual Oral Health  
2682 Status'. *Nature Communications* 9 (1): 4744. <https://doi.org/10.1038/s41467-018-07148-3>.

2683

2684 Jiang, Yuming, Devasahayam Arokia Balaya Rex, Dina Schuster, Benjamin A. Neely,  
2685 Germán L. Rosano, Norbert Volkmar, Amanda Momenzadeh, et al. 2024.  
2686 'Comprehensive Overview of Bottom-Up Proteomics Using Mass Spectrometry'.  
2687 *ACS Measurement Science Au* 4 (4): 338–417.  
2688 <https://doi.org/10.1021/acsmeasuresciau.3c00068>.

2689 Johnson, Jethro S., Daniel J. Spakowicz, Bo-Young Hong, Lauren M. Petersen, Patrick  
2690 Demkowicz, Lei Chen, Shana R. Leopold, et al. 2019. 'Evaluation of 16S rRNA  
2691 Gene Sequencing for Species and Strain-Level Microbiome Analysis'. *Nature  
2692 Communications* 10 (1): 5029. <https://doi.org/10.1038/s41467-019-13036-1>.

2693 Johnson, Richard S., Brian C. Searle, Brook L. Nunn, Jason M. Gilmore, Molly Phillips,  
2694 Chris T. Amemiya, Michelle Heck, and Michael J. MacCoss. 2020. 'Assessing  
2695 Protein Sequence Database Suitability Using De Novo Sequencing \*'. *Molecular &  
2696 Cellular Proteomics* 19 (1): 198–208. <https://doi.org/10.1074/mcp.TIR119.001752>.

2697 Justice, Nicholas B., Zhou Li, Yingfeng Wang, Susan E. Spaulding, Annika C. Mosier,  
2698 Robert L. Hettich, Chongle Pan, and Jillian F. Banfield. 2014. '(15)N- and (2)H  
2699 Proteomic Stable Isotope Probing Links Nitrogen Flow to Archaeal Heterotrophic  
2700 Activity'. *Environmental Microbiology* 16 (10): 3224–37.

2701 <https://doi.org/10.1111/1462-2920.12488>.

2702 Käll, Lukas, Jesse D. Canterbury, Jason Weston, William Stafford Noble, and Michael J.  
2703 MacCoss. 2007. 'Semi-Supervised Learning for Peptide Identification from  
2704 Shotgun Proteomics Datasets'. *Nature Methods* 4 (11): 923–25.  
2705 <https://doi.org/10.1038/nmeth1113>.

2706 Käll, Lukas, John D. Storey, Michael J. MacCoss, and William Stafford Noble. 2008.  
2707 'Posterior Error Probabilities and False Discovery Rates: Two Sides of the Same  
2708 Coin'. *Journal of Proteome Research* 7 (1): 40–44.  
2709 <https://doi.org/10.1021/pr700739d>.

2710 Kanehisa, Minoru, and Susumu Goto. 2000. 'KEGG: Kyoto Encyclopedia of Genes and  
2711 Genomes'. *Nucleic Acids Research* 28 (1): 27–30.  
2712 <https://doi.org/10.1093/nar/28.1.27>.

2713 Keiblinger, Katharina M., Stephan Fuchs, Sophie Zechmeister-Boltenstern, and Katharina  
2714 Riedel. 2016. 'Soil and Leaf Litter Metaproteomics—a Brief Guideline from  
2715 Sampling to Understanding'. *FEMS Microbiology Ecology* 92 (11): fiw180.  
2716 <https://doi.org/10.1093/femsec/fiw180>.

2717 Keiblinger, Katharina M., Inés C. Wilhartitz, Thomas Schneider, Bernd Roschitzki,  
2718 Emanuel Schmid, Leo Eberl, Kathrin Riedel, and Sophie Zechmeister-Boltenstern.  
2719 2012. 'Soil Metaproteomics – Comparative Evaluation of Protein Extraction  
2720 Protocols'. *Soil Biology and Biochemistry* 54 (November):14–24.  
2721 <https://doi.org/10.1016/j.soilbio.2012.05.014>.

2722 Kieser, Silas, Evgeny M. Zdobnov, and Mirko Trajkovski. 2022. 'Comprehensive Mouse  
2723 Microbiota Genome Catalog Reveals Major Difference to Its Human Counterpart'.  
2724 *PLoS Computational Biology* 18 (3): e1009947.  
2725 <https://doi.org/10.1371/journal.pcbi.1009947>.

2726 Kim, Daehwan, Li Song, Florian P. Breitwieser, and Steven L. Salzberg. 2016.  
2727 'Centrifuge: Rapid and Sensitive Classification of Metagenomic Sequences'.  
2728 *Genome Research* 26 (12): 1721–29. <https://doi.org/10.1101/gr.210641.116>.

2729 Kleikamp, Hugo B. C., Denis Grouzdev, Pim Schaasberg, Ramon van Valderen, Ramon  
2730 van der Zwaan, Roel van de Wijngaart, Yuemei Lin, et al. 2023. 'Metaproteomics,  
2731 Metagenomics and 16S rRNA Sequencing Provide Different Perspectives on the  
2732 Aerobic Granular Sludge Microbiome'. *Water Research* 246 (November):120700.  
2733 <https://doi.org/10.1016/j.watres.2023.120700>.

2734 Kleikamp, Hugo B. C., Mario Pronk, Claudia Tugui, Leonor Guedes da Silva, Ben Abbas,  
2735 Yue Mei Lin, Mark C. M. van Loosdrecht, and Martin Pabst. 2021. 'Database-  
2736 Independent de Novo Metaproteomics of Complex Microbial Communities'. *Cell*  
2737 *Systems* 12 (5): 375-383.e5. <https://doi.org/10.1016/j.cels.2021.04.003>.

2738 Kleiner, Manuel. 2019. 'Metaproteomics: Much More than Measuring Gene Expression in  
2739 Microbial Communities'. *mSystems* 4 (3): 10.1128/msystems.00115-19.  
2740 <https://doi.org/10.1128/msystems.00115-19>.

2741 Kleiner, Manuel, Angela Kouris, Marlene Violette, Grace D'Angelo, Yihua Liu, Abigail  
2742 Korenek, Nikola Tolić, et al. 2023. 'Ultra-Sensitive Isotope Probing to Quantify  
2743 Activity and Substrate Assimilation in Microbiomes'. *Microbiome* 11 (1): 24.  
2744 <https://doi.org/10.1186/s40168-022-01454-1>.

2745 Kleiner, Manuel, Erin Thorson, Christine E. Sharp, Xiaoli Dong, Dan Liu, Carmen Li, and  
2746 Marc Strous. 2017. 'Assessing Species Biomass Contributions in Microbial  
2747 Communities via Metaproteomics'. *Nature Communications* 8 (1): 1558.  
2748 <https://doi.org/10.1038/s41467-017-01544-x>.

2749 Kleiner, Manuel, Cecilia Wentrup, Christian Lott, Hanno Teeling, Silke Wetzel, Jacque  
2750 Young, Yun-Juan Chang, et al. 2012. 'Metaproteomics of a Gutless Marine Worm  
2751 and Its Symbiotic Microbial Community Reveal Unusual Pathways for Carbon and  
2752 Energy Use'. *Proceedings of the National Academy of Sciences* 109 (19).  
2753 <https://doi.org/10.1073/pnas.1121198109>.

2754 Kollipara, Laxmikanth, and René P. Zahedi. 2013. 'Protein Carbamylation: In Vivo  
2755 Modification or in Vitro Artefact?' *PROTEOMICS* 13 (6): 941–44.

2756 <https://doi.org/10.1002/pmic.201200452>.

2757 Kong, Andy T., Felipe V. Leprevost, Dmitry M. Avtonomov, Dattatreya Mellacheruvu, and  
2758 Alexey I. Nesvizhskii. 2017. 'MSFragger: Ultrafast and Comprehensive Peptide  
2759 Identification in Mass Spectrometry–Based Proteomics'. *Nature Methods* 14 (5):  
2760 513–20. <https://doi.org/10.1038/nmeth.4256>.

2761 Kong, Ling-Fen, Yan-Bin He, Zhang-Xian Xie, Xing Luo, Hao Zhang, Sheng-Hui Yi, Zhi-  
2762 Long Lin, et al. 2021. 'Illuminating Key Microbial Players and Metabolic Processes  
2763 Involved in the Remineralization of Particulate Organic Carbon in the Ocean's  
2764 Twilight Zone by Metaproteomics'. *Applied and Environmental Microbiology* 87  
2765 (20): e00986-21. <https://doi.org/10.1128/AEM.00986-21>.

2766 Kulak, Nils A., Garwin Pichler, Igor Paron, Nagarjuna Nagaraj, and Matthias Mann. 2014.  
2767 'Minimal, Encapsulated Proteomic-Sample Processing Applied to Copy-Number  
2768 Estimation in Eukaryotic Cells'. *Nature Methods* 11 (3): 319–24.  
2769 <https://doi.org/10.1038/nmeth.2834>.

2770 Kunath, B. J., O. Hickl, P. Queirós, C. Martin-Gallausiaux, L. A. Lebrun, R. Halder, C. C.  
2771 Laczny, et al. 2022. 'Alterations of Oral Microbiota and Impact on the Gut  
2772 Microbiome in Type 1 Diabetes Mellitus Revealed by Integrated Multi-Omic  
2773 Analyses'. *Microbiome* 10 (1): 243. <https://doi.org/10.1186/s40168-022-01435-4>.

2774 Kunath, Benoit J., Andreas Bremges, Aaron Weimann, Alice C. McHardy, and Phillip B.  
2775 Pope. 2017. 'Metagenomics and CAZyme Discovery'. *Methods in Molecular  
2776 Biology (Clifton, N.J.)* 1588:255–77. [https://doi.org/10.1007/978-1-4939-6899-  
2777 2\\_20](https://doi.org/10.1007/978-1-4939-6899-2_20).

2778 Kunath, Benoit J., Giusi Minniti, Morten Skaugen, Live H. Hagen, Gustav Vaaje-Kolstad,  
2779 Vincent G. H. Eijsink, Phillip B. Pope, and Magnus Ø Arntzen. 2019.  
2780 'Metaproteomics: Sample Preparation and Methodological Considerations.' In  
2781 *Emerging Sample Treatments in Proteomics (Ed. Capelo-Martínez, J.-L.)*, 187–  
2782 215. Cham: Springer International Publishing. [https://dokumen.pub/emerging-  
2783 sample-treatments-in-proteomics-1st-ed-978-3-030-12297-3-978-3-030-12298-  
2784 0.html](https://dokumen.pub/emerging-sample-treatments-in-proteomics-1st-ed-978-3-030-12297-3-978-3-030-12298-0.html).

2785 Lautenbacher, Ludwig, Kevin L. Yang, Tobias Kockmann, Christian Panse, Matthew  
2786 Chambers, Elias Kahl, Fengchao Yu, et al. 2024. 'Koina: Democratizing Machine  
2787 Learning for Proteomics Research'. bioRxiv.  
2788 <https://doi.org/10.1101/2024.06.01.596953>.

2789 Lazear, Michael R. 2023. 'Sage: An Open-Source Tool for Fast Proteomics Searching and  
2790 Quantification at Scale'. *Journal of Proteome Research* 22 (11): 3652–59.  
2791 <https://doi.org/10.1021/acs.jproteome.3c00486>.

2792 Lenčo, Juraj, Siddharth Jadeja, Denis K. Naplekov, Oleg V. Krokhin, Maria A. Khalikova,  
2793 Petr Chocholouš, Jiří Urban, Ken Broeckhoven, Lucie Nováková, and František  
2794 Švec. 2022. 'Reversed-Phase Liquid Chromatography of Peptides for Bottom-Up  
2795 Proteomics: A Tutorial'. *Journal of Proteome Research* 21 (12): 2846–92.  
2796 <https://doi.org/10.1021/acs.jproteome.2c00407>.

2797 Lesker, Till R., Abilash C. Durairaj, Eric J. C. Gálvez, Ilias Lagkouvardos, John F. Baines,  
2798 Thomas Clavel, Alexander Sczyrba, Alice C. McHardy, and Till Strowig. 2020. 'An  
2799 Integrated Metagenome Catalog Reveals New Insights into the Murine Gut  
2800 Microbiome'. *Cell Reports* 30 (9): 2909-2922.e6.  
2801 <https://doi.org/10.1016/j.celrep.2020.02.036>.

2802 Letunic, Ivica, Takuji Yamada, Minoru Kanehisa, and Peer Bork. 2008. 'iPath: Interactive  
2803 Exploration of Biochemical Pathways and Networks'. *Trends in Biochemical  
2804 Sciences* 33 (3): 101–3. <https://doi.org/10.1016/j.tibs.2008.01.001>.

2805 Levin, Yishai. 2011. 'The Role of Statistical Power Analysis in Quantitative Proteomics'.  
2806 *PROTEOMICS* 11 (12): 2565–67. <https://doi.org/10.1002/pmic.201100033>.

2807 Lezcano, María Ángeles, Laura Sánchez-García, Antonio Quesada, Daniel Carrizo,  
2808 Miguel Ángel Fernández-Martínez, Erika Cavalcante-Silva, and Víctor Parro. 2022.  
2809 'Comprehensive Metabolic and Taxonomic Reconstruction of an Ancient Microbial  
2810 Mat From the McMurdo Ice Shelf (Antarctica) by Integrating Genetic,

- 2811 Metaproteomic and Lipid Biomarker Analyses'. *Frontiers in Microbiology* 13 (April).  
 2812 <https://doi.org/10.3389/fmicb.2022.799360>.
- 2813 Li, Junhua, Huijue Jia, Xianghang Cai, Huanzi Zhong, Qiang Feng, Shinichi Sunagawa,  
 2814 Manimozhiyan Arumugam, et al. 2014. 'An Integrated Catalog of Reference Genes  
 2815 in the Human Gut Microbiome'. *Nature Biotechnology* 32 (8): 834–41.  
 2816 <https://doi.org/10.1038/nbt.2942>.
- 2817 Li, Leyuan, Janice Mayne, Adrian Beltran, Xu Zhang, Zhibin Ning, and Daniel Figeys.  
 2818 2024. 'RapidAIM 2.0: A High-Throughput Assay to Study Functional Response of  
 2819 Human Gut Microbiome to Xenobiotics'. *Microbiome Research Reports* 3 (2): N/A-  
 2820 N/A. <https://doi.org/10.20517/mrr.2023.57>.
- 2821 Li, Leyuan, Zhibin Ning, Kai Cheng, Xu Zhang, Caitlin M. A. Simopoulos, and Daniel  
 2822 Figeys. 2022. 'iMetaLab Suite: A One-Stop Toolset for Metaproteomics'. *iMeta* 1  
 2823 (2): e25. <https://doi.org/10.1002/imt2.25>.
- 2824 Li, Leyuan, Zhibin Ning, Xu Zhang, Janice Mayne, Kai Cheng, Alain Stintzi, and Daniel  
 2825 Figeys. 2020. 'RapidAIM: A Culture- and Metaproteomics-Based Rapid Assay of  
 2826 Individual Microbiome Responses to Drugs'. *Microbiome* 8 (1): 33.  
 2827 <https://doi.org/10.1186/s40168-020-00806-z>.
- 2828 Li, Leyuan, Tong Wang, Zhibin Ning, Xu Zhang, James Butcher, Joeselle M. Serrana,  
 2829 Caitlin M. A. Simopoulos, et al. 2023. 'Revealing Proteome-Level Functional  
 2830 Redundancy in the Human Gut Microbiome Using Ultra-Deep Metaproteomics'.  
 2831 *Nature Communications* 14 (1): 3428. [https://doi.org/10.1038/s41467-023-39149-](https://doi.org/10.1038/s41467-023-39149-2)  
 2832 2.
- 2833 Li, W., L. Jaroszewski, and A. Godzik. 2001. 'Clustering of Highly Homologous  
 2834 Sequences to Reduce the Size of Large Protein Databases'. *Bioinformatics*  
 2835 (*Oxford, England*) 17 (3): 282–83. <https://doi.org/10.1093/bioinformatics/17.3.282>.
- 2836 Li, Zhou, Yingfeng Wang, Qiuming Yao, Nicholas B. Justice, Tae-Hyuk Ahn, Dong Xu,  
 2837 Robert L. Hettich, Jillian F. Banfield, and Chongle Pan. 2014. 'Diverse and  
 2838 Divergent Protein Post-Translational Modifications in Two Growth Stages of a  
 2839 Natural Microbial Community'. *Nature Communications* 5 (1): 4405.  
 2840 <https://doi.org/10.1038/ncomms5405>.
- 2841 Li, Zhou, Qiuming Yao, Xuan Guo, Alexander Crits-Christoph, Melanie A. Mayes, William  
 2842 Judson Herve Iv, Sarah L. Lebeis, et al. 2019. 'Genome-Resolved Proteomic  
 2843 Stable Isotope Probing of Soil Microbial Communities Using <sup>13</sup>CO<sub>2</sub> and <sup>13</sup>C-  
 2844 Methanol'. *Frontiers in Microbiology* 10 (December):2706.  
 2845 <https://doi.org/10.3389/fmicb.2019.02706>.
- 2846 Liao, Bo, Zhibin Ning, Kai Cheng, Xu Zhang, Leyuan Li, Janice Mayne, and Daniel Figeys.  
 2847 2018. 'iMetaLab 1.0: A Web Platform for Metaproteomics Data Analysis'.  
 2848 *Bioinformatics* 34 (22): 3954–56. <https://doi.org/10.1093/bioinformatics/bty466>.
- 2849 Lin, Zongtao, Joanna Gongora, Xingyu Liu, Yixuan Xie, Chenfeng Zhao, Dongwen Lv,  
 2850 and Benjamin A. Garcia. 2023. 'Automation to Enable High-Throughput Chemical  
 2851 Proteomics'. *Journal of Proteome Research* 22 (12): 3676–82.  
 2852 <https://doi.org/10.1021/acs.jproteome.3c00467>.
- 2853 Liu, Lei, Yu Yang, Yu Deng, and Tong Zhang. 2022. 'Nanopore Long-Read-Only  
 2854 Metagenomics Enables Complete and High-Quality Genome Reconstruction from  
 2855 Mock and Complex Metagenomes'. *Microbiome* 10 (1): 209.  
 2856 <https://doi.org/10.1186/s40168-022-01415-8>.
- 2857 Long, Shuping, Yi Yang, Chengpin Shen, Yiwen Wang, Anmei Deng, Qin Qin, and Liang  
 2858 Qiao. 2020. 'Metaproteomics Characterizes Human Gut Microbiome Function in  
 2859 Colorectal Cancer'. *Npj Biofilms and Microbiomes* 6 (1): 1–10.  
 2860 <https://doi.org/10.1038/s41522-020-0123-4>.
- 2861 Low, Teck Yew, M. Aiman Mohtar, Pey Yee Lee, Nursyazwani Omar, Houjiang Zhou, and  
 2862 Mingliang Ye. 2021. 'Widening the Bottleneck of Phosphoproteomics: Evolving  
 2863 Strategies for Phosphopeptide Enrichment'. *Mass Spectrometry Reviews* 40 (4):  
 2864 309–33. <https://doi.org/10.1002/mas.21636>.
- 2865 Martens, Lennart, Matthew Chambers, Marc Sturm, Darren Kessner, Fredrik Levander,

- 2866 Jim Shofstahl, Wilfred H. Tang, et al. 2011. 'mzML—a Community Standard for  
2867 Mass Spectrometry Data \*'. *Molecular & Cellular Proteomics* 10 (1).  
2868 <https://doi.org/10.1074/mcp.R110.000133>.
- 2869 Matthiesen, Rune, and Jakob Bunkenborg. 2013. 'Introduction to Mass Spectrometry-  
2870 Based Proteomics'. In *Mass Spectrometry Data Analysis in Proteomics*, edited by  
2871 Rune Matthiesen, 1–45. Totowa, NJ: Humana Press. [https://doi.org/10.1007/978-1-62703-392-3\\_1](https://doi.org/10.1007/978-1-62703-392-3_1).
- 2873 Meijenfeldt, F. A. Bastiaan von, Ksenia Arkhipova, Diego D. Cambuy, Felipe H. Coutinho,  
2874 and Bas E. Dutilh. 2019. 'Robust Taxonomic Classification of Uncharted Microbial  
2875 Sequences and Bins with CAT and BAT'. *Genome Biology* 20 (1): 217.  
2876 <https://doi.org/10.1186/s13059-019-1817-x>.
- 2877 Mesuere, Bart, Felix Van der Jeugt, Toon Willems, Tom Naessens, Bart Devreese,  
2878 Lennart Martens, and Peter Dawyndt. 2018. 'High-Throughput Metaproteomics  
2879 Data Analysis with Unipept: A Tutorial'. *Journal of Proteomics*, Tutorials in  
2880 Bioinformatics for Biological Science, 171 (January):11–22.  
2881 <https://doi.org/10.1016/j.jprot.2017.05.022>.
- 2882 Meyer, Susann, Nicole Hüttig, Marianne Zenk, Udo Jäckel, and Dierk-Christoph Pöther.  
2883 2023. 'Bioaerosols in Swine Confinement Buildings: A Metaproteomic View'.  
2884 *Environmental Microbiology Reports* 15 (6): 684–97. <https://doi.org/10.1111/1758-2229.13208>.
- 2886 Millikin, Robert J., Stefan K. Solntsev, Michael R. Shortreed, and Lloyd M. Smith. 2018.  
2887 'Ultrafast Peptide Label-Free Quantification with FlashLFQ'. *Journal of Proteome  
2888 Research* 17 (1): 386–91. <https://doi.org/10.1021/acs.jproteome.7b00608>.
- 2889 Minniti, Giusi, Simen Rød Sandve, János Tamás Padra, Live Heldal Hagen, Sara Lindén,  
2890 Phillip B. Pope, Magnus Ø Arntzen, and Gustav Vaaje-Kolstad. 2019. 'The Farmed  
2891 Atlantic Salmon (*Salmo Salar*) Skin-Mucus Proteome and Its Nutrient Potential for  
2892 the Resident Bacterial Community'. *Genes* 10 (7): 515.  
2893 <https://doi.org/10.3390/genes10070515>.
- 2894 Mistry, Jaina, Sara Chuguransky, Lowri Williams, Matloob Qureshi, Gustavo A Salazar,  
2895 Erik L L Sonnhammer, Silvio C E Tosatto, et al. 2021. 'Pfam: The Protein Families  
2896 Database in 2021'. *Nucleic Acids Research* 49 (D1): D412–19.  
2897 <https://doi.org/10.1093/nar/gkaa913>.
- 2898 Mordant, Angie, and Manuel Kleiner. 2021. 'Evaluation of Sample Preservation and  
2899 Storage Methods for Metaproteomics Analysis of Intestinal Microbiomes'.  
2900 *Microbiology Spectrum* 9 (3): e01877-21. <https://doi.org/10.1128/Spectrum.01877-21>.
- 2902 Morris, Laura S., and Julian R. Marchesi. 2016. 'Assessing the Impact of Long Term  
2903 Frozen Storage of Faecal Samples on Protein Concentration and Protease  
2904 Activity'. *Journal of Microbiological Methods* 123 (April):31–38.  
2905 <https://doi.org/10.1016/j.mimet.2016.02.001>.
- 2906 Mueller, Ryan S, Vincent J Deneff, Linda H Kalnejais, K Blake Suttle, Brian C Thomas,  
2907 Paul Wilmes, Richard L Smith, et al. 2010. 'Ecological Distribution and Population  
2908 Physiology Defined by Proteomics in a Natural Microbial Community'. *Molecular  
2909 Systems Biology* 6 (1): 374. <https://doi.org/10.1038/msb.2010.30>.
- 2910 Muth, Thilo, Alexander Behne, Robert Heyer, Fabian Kohrs, Dirk Benndorf, Marcus  
2911 Hoffmann, Miro Lehtevä, Udo Reichl, Lennart Martens, and Erdmann Rapp. 2015.  
2912 'The MetaProteomeAnalyzer: A Powerful Open-Source Software Suite for  
2913 Metaproteomics Data Analysis and Interpretation'. *Journal of Proteome Research*  
2914 14 (3): 1557–65. <https://doi.org/10.1021/pr501246w>.
- 2915 Muth, Thilo, Carolin A. Kolmeder, Jarkko Salojärvi, Salla Keskitalo, Markku Varjosalo,  
2916 Froukje J. Verdam, Sander S. Rensen, et al. 2015. 'Navigating through  
2917 Metaproteomics Data: A Logbook of Database Searching'. *PROTEOMICS* 15 (20):  
2918 3439–53. <https://doi.org/10.1002/pmic.201400560>.
- 2919 Mysling, Simon, Giuseppe Palmisano, Peter Højrup, and Morten Thaysen-Andersen.  
2920 2010. 'Utilizing Ion-Pairing Hydrophilic Interaction Chromatography Solid Phase



- 2921 Extraction for Efficient Glycopeptide Enrichment in Glycoproteomics'. *Analytical*  
 2922 *Chemistry* 82 (13): 5598–5609. <https://doi.org/10.1021/ac100530w>.
- 2923 Nakayasu, Ernesto S., Marina Gritsenko, Paul D. Piehowski, Yuqian Gao, Daniel J. Orton,  
 2924 Athena A. Schepmoes, Thomas L. Fillmore, et al. 2021. 'Tutorial: Best Practices  
 2925 and Considerations for Mass-Spectrometry-Based Protein Biomarker Discovery  
 2926 and Validation'. *Nature Protocols* 16 (8): 3737–60. [https://doi.org/10.1038/s41596-](https://doi.org/10.1038/s41596-021-00566-6)  
 2927 021-00566-6.
- 2928 Narayanasamy, Shaman, Yohan Jarosz, Emilie E. L. Muller, Anna Heintz-Buschart, Malte  
 2929 Herold, Anne Kaysen, Cédric C. Laczny, Nicolás Pinel, Patrick May, and Paul  
 2930 Wilmes. 2016. 'IMP: A Pipeline for Reproducible Reference-Independent  
 2931 Integrated Metagenomic and Metatranscriptomic Analyses'. *Genome Biology* 17  
 2932 (1): 260. <https://doi.org/10.1186/s13059-016-1116-8>.
- 2933 Nebauer, Daniel J., Leanne A. Pearson, and Brett A. Neilan. 2024. 'Critical Steps in an  
 2934 Environmental Metaproteomics Workflow'. *Environmental Microbiology* 26 (5):  
 2935 e16637. <https://doi.org/10.1111/1462-2920.16637>.
- 2936 Nesvizhskii, Alexey I., and Ruedi Aebersold. 2005. 'Interpretation of Shotgun Proteomic  
 2937 Data'. *Molecular & Cellular Proteomics* 4 (10): 1419–40.  
 2938 <https://doi.org/10.1074/mcp.R500012-MCP200>.
- 2939 Nickerson, Jessica L., and Alan A. Doucette. 2020. 'Rapid and Quantitative Protein  
 2940 Precipitation for Proteome Analysis by Mass Spectrometry'. *Journal of Proteome*  
 2941 *Research* 19 (5): 2035–42. <https://doi.org/10.1021/acs.jproteome.9b00867>.
- 2942 Niu, Liangjie, Hang Zhang, Zhaokun Wu, Yibo Wang, Hui Liu, Xiaolin Wu, and Wei Wang.  
 2943 2018. 'Modified TCA/Acetone Precipitation of Plant Proteins for Proteomic  
 2944 Analysis'. *PLOS ONE* 13 (12): e0202238.  
 2945 <https://doi.org/10.1371/journal.pone.0202238>.
- 2946 Nowatzky, Yannek, Philipp Benner, Knut Reinert, and Thilo Muth. 2023. 'Mistle: Bringing  
 2947 Spectral Library Predictions to Metaproteomics with an Efficient Search Index'.  
 2948 *Bioinformatics* 39 (6): btad376. <https://doi.org/10.1093/bioinformatics/btad376>.
- 2949 Oberg, Ann L., and Olga Vitek. 2009. 'Statistical Design of Quantitative Mass  
 2950 Spectrometry-Based Proteomic Experiments'. *Journal of Proteome Research* 8  
 2951 (5): 2144–56. <https://doi.org/10.1021/pr8010099>.
- 2952 Odom, Aubrey R., Tyler Faits, Eduardo Castro-Nallar, Keith A. Crandall, and W. Evan  
 2953 Johnson. 2023. 'Metagenomic Profiling Pipelines Improve Taxonomic  
 2954 Classification for 16S Amplicon Sequencing Data'. *Scientific Reports* 13 (1):  
 2955 13957. <https://doi.org/10.1038/s41598-023-40799-x>.
- 2956 Okeke, Emmanuel Sunday, Richard Ekeng Ita, Egong John Egong, Lydia Etuk Udofia,  
 2957 Chiamaka Linda Mgbekidinma, and Otobong Donald Akan. 2021.  
 2958 'Metaproteomics Insights into Fermented Fish and Vegetable Products and  
 2959 Associated Microbes'. *Food Chemistry: Molecular Sciences* 3 (December):100045.  
 2960 <https://doi.org/10.1016/j.fochms.2021.100045>.
- 2961 O'Leary, Nuala A., Mathew W. Wright, J. Rodney Brister, Stacy Ciufo, Diana Haddad,  
 2962 Rich McVeigh, Bhanu Rajput, et al. 2016. 'Reference Sequence (RefSeq)  
 2963 Database at NCBI: Current Status, Taxonomic Expansion, and Functional  
 2964 Annotation'. *Nucleic Acids Research* 44 (D1): D733-745.  
 2965 <https://doi.org/10.1093/nar/gkv1189>.
- 2966 Olivella, Roger, Cristina Chiva, Marc Serret, Daniel Mancera, Luca Cozzuto, Antoni  
 2967 Hermoso, Eva Borràs, et al. 2021. 'QCloud2: An Improved Cloud-Based Quality-  
 2968 Control System for Mass-Spectrometry-Based Proteomics Laboratories'. *Journal*  
 2969 *of Proteome Research* 20 (4): 2010–13.  
 2970 <https://doi.org/10.1021/acs.jproteome.0c00853>.
- 2971 Pan, Haixia, Ruddy Wattiez, and David Gillan. 2024. 'Soil Metaproteomics for Microbial  
 2972 Community Profiling: Methodologies and Challenges'. *Current Microbiology* 81 (8):  
 2973 257. <https://doi.org/10.1007/s00284-024-03781-y>.
- 2974 Paramasivan, Selvam, Janna L. Morrison, Mitchell C. Lock, Jack R. T. Darby, Roberto A.  
 2975 Barrero, Paul C. Mills, and Pawel Sadowski. 2023. 'Automated Proteomics

2976 Workflows for High-Throughput Library Generation and Biomarker Detection Using  
2977 Data-Independent Acquisition'. *Journal of Proteome Research* 22 (6): 2018–29.  
2978 <https://doi.org/10.1021/acs.jproteome.3c00074>.

2979 Park, Sung Kyu Robin, Titus Jung, Peter S. Thuy-Boun, Ana Y. Wang, John R. III Yates,  
2980 and Dennis W. Wolan. 2019. 'ComPIL 2.0: An Updated Comprehensive  
2981 Metaproteomics Database'. *Journal of Proteome Research* 18 (2): 616–22.  
2982 <https://doi.org/10.1021/acs.jproteome.8b00722>.

2983 Pathak, Khyatiben V., Marissa I. McGilvrey, Charles K. Hu, Krystine Garcia-Mansfield,  
2984 Karen Lewandoski, Zahra Eftekhari, Yate-Ching Yuan, Frederic Zenhausern,  
2985 Emmanuel Menashi, and Patrick Pirrotte. 2020. 'Molecular Profiling of Innate  
2986 Immune Response Mechanisms in Ventilator-Associated Pneumonia'. *Molecular &  
2987 Cellular Proteomics* 19 (10): 1688–1705.  
2988 <https://doi.org/10.1074/mcp.RA120.002207>.

2989 Patnode, Michael L., Zachary W. Beller, Nathan D. Han, Jiye Cheng, Samantha L. Peters,  
2990 Nicolas Terrapon, Bernard Henrissat, et al. 2019. 'Interspecies Competition  
2991 Impacts Targeted Manipulation of Human Gut Bacteria by Fiber-Derived Glycans'.  
2992 *Cell* 179 (1): 59-73.e13. <https://doi.org/10.1016/j.cell.2019.08.011>.

2993 Pavelka, Norman, Marjorie L. Fournier, Selene K. Swanson, Mattia Pelizzola, Paola  
2994 Ricciardi-Castagnoli, Laurence Florens, and Michael P. Washburn. 2008.  
2995 'Statistical Similarities between Transcriptomics and Quantitative Shotgun  
2996 Proteomics Data \*'. *Molecular & Cellular Proteomics* 7 (4): 631–44.  
2997 <https://doi.org/10.1074/mcp.M700240-MCP200>.

2998 Perez-Riverol, Yasset, Chakradhar Bandla, Deepti J Kundu, Selvakumar Kamatchinathan,  
2999 Jingwen Bai, Suresh Hewapathirana, Nithu Sara John, et al. 2024. 'The PRIDE  
3000 Database at 20 Years: 2025 Update'. *Nucleic Acids Research*, November,  
3001 gkae1011. <https://doi.org/10.1093/nar/gkae1011>.

3002 Pettersen, Veronika Kuchařová, Luis Caetano Martha Antunes, Antoine Dufour, and  
3003 Marie-Claire Arrieta. 2022. 'Inferring Early-Life Host and Microbiome Functions by  
3004 Mass Spectrometry-Based Metaproteomics and Metabolomics'. *Computational  
3005 and Structural Biotechnology Journal* 20 (January):274–86.  
3006 <https://doi.org/10.1016/j.csbj.2021.12.012>.

3007 Pible, Olivier, Pauline Petit, Gérard Steinmetz, Corinne Rivasseau, and Jean Armengaud.  
3008 2023. 'Taxonomical Composition and Functional Analysis of Biofilms Sampled  
3009 from a Nuclear Storage Pool'. *Frontiers in Microbiology* 14 (April).  
3010 <https://doi.org/10.3389/fmicb.2023.1148976>.

3011 Pietilä, Sami, Tomi Suomi, and Laura L Elo. 2022. 'Introducing Untargeted Data-  
3012 Independent Acquisition for Metaproteomics of Complex Microbial Samples'. *ISME  
3013 Communications* 2 (1): 51. <https://doi.org/10.1038/s43705-022-00137-0>.

3014 Plancade, Sandra, Magali Berland, Mélisande Blein-Nicolas, Olivier Langella, Ariane  
3015 Bassignani, and Catherine Juste. 2022. 'A Combined Test for Feature Selection  
3016 on Sparse Metaproteomics Data—an Alternative to Missing Value Imputation'.  
3017 *PeerJ* 10 (June):e13525. <https://doi.org/10.7717/peerj.13525>.

3018 Ponnudurai, Ruby, Stefan E Heiden, Lizbeth Sayavedra, Tjorven Hinzke, Manuel Kleiner,  
3019 Christian Hentschker, Horst Felbeck, et al. 2020. 'Comparative Proteomics of  
3020 Related Symbiotic Mussel Species Reveals High Variability of Host–Symbiont  
3021 Interactions'. *The ISME Journal* 14 (2): 649–56. [https://doi.org/10.1038/s41396-  
3022 019-0517-6](https://doi.org/10.1038/s41396-019-0517-6).

3023 Porcheddu, Massimo, Marcello Abbondio, Laura De Diego, Sergio Uzzau, and Alessandro  
3024 Tanca. 2023. 'Meta4P: A User-Friendly Tool to Parse Label-Free Quantitative  
3025 Metaproteomic Data and Taxonomic/Functional Annotations'. *Journal of Proteome  
3026 Research* 22 (6): 2109–13. <https://doi.org/10.1021/acs.jproteome.2c00803>.

3027 Qian, Chen, and Robert L. Hettich. 2017. 'Optimized Extraction Method To Remove  
3028 Humic Acid Interferences from Soil Samples Prior to Microbial Proteome  
3029 Measurements'. *Journal of Proteome Research* 16 (7): 2537–46.  
3030 <https://doi.org/10.1021/acs.jproteome.7b00103>.

- 3031 Queirós, Pedro, Francesco Delogu, Oskar Hickl, Patrick May, and Paul Wilmes. 2021.  
 3032 'Mantis: Flexible and Consensus-Driven Genome Annotation'. *GigaScience* 10 (6):  
 3033 giab042. <https://doi.org/10.1093/gigascience/giab042>.
- 3034 Quevillon, E., V. Silventoinen, S. Pillai, N. Harte, N. Mulder, R. Apweiler, and R. Lopez.  
 3035 2005. 'InterProScan: Protein Domains Identifier'. *Nucleic Acids Research* 33  
 3036 (suppl\_2): W116–20. <https://doi.org/10.1093/nar/gki442>.
- 3037 Rabe, Alexander, Manuela Gesell Salazar, Stephan Michalik, Thomas Kocher, Harald  
 3038 Below, Uwe Völker, and Alexander Welk. 2022. 'Impact of Different Oral  
 3039 Treatments on the Composition of the Supragingival Plaque Microbiome'. *Journal*  
 3040 *of Oral Microbiology* 14 (1): 2138251.  
 3041 <https://doi.org/10.1080/20002297.2022.2138251>.
- 3042 Rajczewski, Andrew T., J. Alfredo Blakeley-Ruiz, Annaliese Meyer, Simina Vintila,  
 3043 Matthew R. McIlvin, Tim Van Den Bossche, Brian C. Searle, et al. 2024. 'Data-  
 3044 Independent Acquisition Mass Spectrometry as a Tool for Metaproteomics:  
 3045 Interlaboratory Comparison Using a Model Microbiome'. bioRxiv.  
 3046 <https://doi.org/10.1101/2024.09.18.613707>.
- 3047 Rawlings, Neil D, Alan J Barrett, Paul D Thomas, Xiaosong Huang, Alex Bateman, and  
 3048 Robert D Finn. 2018. 'The MEROPS Database of Proteolytic Enzymes, Their  
 3049 Substrates and Inhibitors in 2017 and a Comparison with Peptidases in the  
 3050 PANTHER Database'. *Nucleic Acids Research* 46 (D1): D624–32.  
 3051 <https://doi.org/10.1093/nar/gkx1134>.
- 3052 Rozanova, Svitlana, Julian Uszkoreit, Karin Schork, Bettina Serschnitzki, Martin  
 3053 Eisenacher, Lars Tönges, Katalin Barkovits-Boeddinghaus, and Katrin Marcus.  
 3054 2023. 'Quality Control—A Stepchild in Quantitative Proteomics: A Case Study for  
 3055 the Human CSF Proteome'. *Biomolecules* 13 (3): 491.  
 3056 <https://doi.org/10.3390/biom13030491>.
- 3057 Ruan, Wenhua, Chao Sun, Qikang Gao, and Neeraj Shrivastava. 2021. 'Metaproteomics  
 3058 Associated with Severe Early Childhood Caries Highlights the Differences in  
 3059 Salivary Proteins'. *Archives of Oral Biology* 131 (November):105220.  
 3060 <https://doi.org/10.1016/j.archoralbio.2021.105220>.
- 3061 Sachsenberg, Timo, Florian-Alexander Herbst, Martin Taubert, René Kermer, Nico  
 3062 Jehmlich, Martin von Bergen, Jana Seifert, and Oliver Kohlbacher. 2015.  
 3063 'MetaProSIP: Automated Inference of Stable Isotope Incorporation Rates in  
 3064 Proteins for Functional Metaproteomics'. *Journal of Proteome Research* 14 (2):  
 3065 619–27. <https://doi.org/10.1021/pr500245w>.
- 3066 Saito, Mak A., Vladimir V. Bulygin, Dawn M. Moran, Craig Taylor, and Chris Scholin.  
 3067 2011. 'Examination of Microbial Proteome Preservation Techniques Applicable to  
 3068 Autonomous Environmental Sample Collection'. *Frontiers in Microbiology* 2  
 3069 (November). <https://doi.org/10.3389/fmicb.2011.00215>.
- 3070 Sajulga, Ray, Caleb Easterly, Michael Riffle, Bart Mesuere, Thilo Muth, Subina Mehta,  
 3071 Praveen Kumar, et al. 2020. 'Survey of Metaproteomics Software Tools for  
 3072 Functional Microbiome Analysis'. *PLOS ONE* 15 (11): e0241503.  
 3073 <https://doi.org/10.1371/journal.pone.0241503>.
- 3074 Salvato, Fernanda, Robert L. Hettich, and Manuel Kleiner. 2021. 'Five Key Aspects of  
 3075 Metaproteomics as a Tool to Understand Functional Interactions in Host-  
 3076 Associated Microbiomes'. *PLOS Pathogens* 17 (2): e1009245.  
 3077 <https://doi.org/10.1371/journal.ppat.1009245>.
- 3078 Salvato, Fernanda, Simina Vintila, Omri M. Finkel, Jeffery L. Dangl, and Manuel Kleiner.  
 3079 2022. 'Evaluation of Protein Extraction Methods for Metaproteomic Analyses of  
 3080 Root-Associated Microbes'. *Molecular Plant-Microbe Interactions* 35 (11): 977–  
 3081 88. <https://doi.org/10.1094/MPMI-05-22-0116-TA>.
- 3082 Sapan, Christine V., and Roger L. Lundblad. 2015. 'Review of Methods for Determination  
 3083 of Total Protein and Peptide Concentration in Biological Samples'. *PROTEOMICS*  
 3084 – *Clinical Applications* 9 (3–4): 268–76. <https://doi.org/10.1002/prca.201400088>.
- 3085 Saunders, Jaclyn K., Matthew R. McIlvin, Chris L. Dupont, Drishti Kaul, Dawn M. Moran,

3086 Tristan Horner, Sarah M. Laperriere, et al. 2022. 'Microbial Functional Diversity  
3087 across Biogeochemical Provinces in the Central Pacific Ocean'. *Proceedings of*  
3088 *the National Academy of Sciences* 119 (37): e2200014119.  
3089 <https://doi.org/10.1073/pnas.2200014119>.

3090 Schallert, Kay, Pieter Verschaffelt, Bart Mesuere, Dirk Benndorf, Lennart Martens, and  
3091 Tim Van Den Bossche. 2022. 'Pout2Prot: An Efficient Tool to Create Protein  
3092 (Sub)Groups from Percolator Output Files'. *Journal of Proteome Research* 21 (4):  
3093 1175–80. <https://doi.org/10.1021/acs.jproteome.1c00685>.

3094 Schäpe, Stephanie Serena, Jannike Lea Krause, Beatrice Engelmann, Katarina Fritz-  
3095 Wallace, Florian Schattenberg, Zishu Liu, Susann Müller, et al. 2019. 'The  
3096 Simplified Human Intestinal Microbiota (SIHUMIx) Shows High Structural and  
3097 Functional Resistance against Changing Transit Times in In Vitro Bioreactors'.  
3098 *Microorganisms* 7 (12): 641. <https://doi.org/10.3390/microorganisms7120641>.

3099 Schiebenhoefer, Henning, Kay Schallert, Bernhard Y. Renard, Kathrin Trappe, Emanuel  
3100 Schmid, Dirk Benndorf, Katharina Riedel, Thilo Muth, and Stephan Fuchs. 2020.  
3101 'A Complete and Flexible Workflow for Metaproteomics Data Analysis Based on  
3102 MetaProteomeAnalyzer and Prophane'. *Nature Protocols* 15 (10): 3212–39.  
3103 <https://doi.org/10.1038/s41596-020-0368-7>.

3104 Schiebenhoefer, Henning, Tim Van Den Bossche, Stephan Fuchs, Bernhard Y. Renard,  
3105 Thilo Muth, and Lennart Martens. 2019. 'Challenges and Promise at the Interface  
3106 of Metaproteomics and Genomics: An Overview of Recent Progress in  
3107 Metaproteogenomic Data Analysis'. *Expert Review of Proteomics* 16 (5): 375–90.  
3108 <https://doi.org/10.1080/14789450.2019.1609944>.

3109 Schiml, Valerie C., Francesco Delogu, Praveen Kumar, Benoit Kunath, Bérénice Batut,  
3110 Subina Mehta, James E. Johnson, et al. 2023. 'Integrative Meta-Omics in Galaxy  
3111 and Beyond'. *Environmental Microbiome* 18 (1): 56.  
3112 <https://doi.org/10.1186/s40793-023-00514-9>.

3113 Schoch, Conrad L, Stacy Ciufu, Mikhail Domrachev, Carol L Hotton, Sivakumar Kannan,  
3114 Rogneda Khovanskaya, Detlef Leipe, et al. 2020. 'NCBI Taxonomy: A  
3115 Comprehensive Update on Curation, Resources and Tools'. *Database* 2020  
3116 (January):baaa062. <https://doi.org/10.1093/database/baaa062>.

3117 Searle, Brian C., Lindsay K. Pino, Jarrett D. Egertson, Ying S. Ting, Robert T. Lawrence,  
3118 Brendan X. MacLean, Judit Villén, and Michael J. MacCoss. 2018. 'Chromatogram  
3119 Libraries Improve Peptide Detection and Quantification by Data Independent  
3120 Acquisition Mass Spectrometry'. *Nature Communications* 9 (1): 5128.  
3121 <https://doi.org/10.1038/s41467-018-07454-w>.

3122 Searle, Brian C., Ariana E. Shannon, and Damien Beau Wilburn. 2023. 'Scribe: Next  
3123 Generation Library Searching for DDA Experiments'. *Journal of Proteome*  
3124 *Research* 22 (2): 482–90. <https://doi.org/10.1021/acs.jproteome.2c00672>.

3125 Sechi, Salvatore, and Brian T. Chait. 1998. 'Modification of Cysteine Residues by  
3126 Alkylation. A Tool in Peptide Mapping and Protein Identification'. *Analytical*  
3127 *Chemistry* 70 (24): 5150–58. <https://doi.org/10.1021/ac9806005>.

3128 Sequeira, João C., Vítor Pereira, M. Madalena Alves, M. Alcina Pereira, Miguel Rocha,  
3129 and Andreia F. Salvador. 2024. 'MOSCA 2.0: A Bioinformatics Framework for  
3130 Metagenomics, Metatranscriptomics and Metaproteomics Data Analysis and  
3131 Visualization'. *Molecular Ecology Resources* 24 (7): e13996.  
3132 <https://doi.org/10.1111/1755-0998.13996>.

3133 Shuken, Steven R. 2023. 'An Introduction to Mass Spectrometry-Based Proteomics'.  
3134 *Journal of Proteome Research* 22 (7): 2151–71.  
3135 <https://doi.org/10.1021/acs.jproteome.2c00838>.

3136 Simopoulos, Caitlin M. A., Zhibin Ning, Leyuan Li, Mona M. Khamis, Xu Zhang, Mathieu  
3137 Lavallée-Adam, and Daniel Figey. 2022. 'MetaProClust-MS1: An MS1 Profiling  
3138 Approach for Large-Scale Microbiome Screening'. *mSystems* 7 (4): e00381-22.  
3139 <https://doi.org/10.1128/msystems.00381-22>.

3140 Sinitcyn, Pavel, Jan Daniel Rudolph, and Jürgen Cox. 2018. 'Computational Methods for

3141 Understanding Mass Spectrometry–Based Shotgun Proteomics Data'. *Annual*  
3142 *Review of Biomedical Data Science* 1 (1): 207–34.  
3143 <https://doi.org/10.1146/annurev-biodatasci-080917-013516>.  
3144 Smyth, Patrick, Xu Zhang, Zhibin Ning, Janice Mayne, Jasmine Isabelle Moore, Krystal  
3145 Walker, Mathieu Lavallée-Adam, and Daniel Figeys. 2020. 'Studying the Temporal  
3146 Dynamics of the Gut Microbiota Using Metabolic Stable Isotope Labeling and  
3147 Metaproteomics'. *Analytical Chemistry* 92 (24): 15711–18.  
3148 <https://doi.org/10.1021/acs.analchem.0c02070>.  
3149 Speda, Jutta, Mikaela A. Johansson, Uno Carlsson, and Martin Karlsson. 2017.  
3150 'Assessment of Sample Preparation Methods for Metaproteomics of Extracellular  
3151 Proteins'. *Analytical Biochemistry* 516 (January):23–36.  
3152 <https://doi.org/10.1016/j.ab.2016.10.008>.  
3153 Starke, Robert, Nico Jehmlich, Trinidad Alfaro, Alice Dohnalkova, Petr Capek, Sheryl L.  
3154 Bell, and Kirsten S. Hofmockel. 2019. 'Incomplete Cell Disruption of Resistant  
3155 Microbes'. *Scientific Reports* 9 (1): 5618. [https://doi.org/10.1038/s41598-019-](https://doi.org/10.1038/s41598-019-42188-9)  
3156 [42188-9](https://doi.org/10.1038/s41598-019-42188-9).  
3157 Starke, Robert, Nico Jehmlich, and Felipe Bastida. 2019. 'Using Proteins to Study How  
3158 Microbes Contribute to Soil Ecosystem Services: The Current State and Future  
3159 Perspectives of Soil Metaproteomics'. *Journal of Proteomics*, 10 Year Anniversary  
3160 of Proteomics, 198 (April):50–58. <https://doi.org/10.1016/j.jprot.2018.11.011>.  
3161 Stewart, Robert D., Marc D. Auffret, Amanda Warr, Alan W. Walker, Rainer Roehe, and  
3162 Mick Watson. 2019. 'Compendium of 4,941 Rumen Metagenome-Assembled  
3163 Genomes for Rumen Microbiome Biology and Enzyme Discovery'. *Nature*  
3164 *Biotechnology* 37 (8): 953–61. <https://doi.org/10.1038/s41587-019-0202-3>.  
3165 Sun, Zhongzhi, Zhibin Ning, and Daniel Figeys. 2024. 'The Landscape and Perspectives  
3166 of the Human Gut Metaproteomics'. *Molecular & Cellular Proteomics* 23 (5):  
3167 100763. <https://doi.org/10.1016/j.mcpro.2024.100763>.  
3168 Tanca, Alessandro, Marcello Abbondio, Giovanni Fiorito, Giovanna Pira, Rosangela Sau,  
3169 Alessandra Manca, Maria Rosaria Muroli, et al. 2022. 'Metaproteomic Profile of  
3170 the Colonic Luminal Microbiota From Patients With Colon Cancer'. *Frontiers in*  
3171 *Microbiology* 13 (April). <https://doi.org/10.3389/fmicb.2022.869523>.  
3172 Tanca, Alessandro, Maria Antonietta Deledda, Laura De Diego, Marcello Abbondio, and  
3173 Sergio Uzzau. 2024. 'Benchmarking Low- and High-Throughput Protein Cleanup  
3174 and Digestion Methods for Human Fecal Metaproteomics'. *mSystems* 9 (7):  
3175 e00661-24. <https://doi.org/10.1128/msystems.00661-24>.  
3176 Tanca, Alessandro, Antonio Palomba, Cristina Fraumene, Daniela Pagnozzi, Valeria  
3177 Manghina, Massimo Deligios, Thilo Muth, et al. 2016. 'The Impact of Sequence  
3178 Database Choice on Metaproteomic Results in Gut Microbiota Studies'.  
3179 *Microbiome* 4 (1): 51. <https://doi.org/10.1186/s40168-016-0196-8>.  
3180 Tanca, Alessandro, Antonio Palomba, Salvatore Pisanu, Maria Filippa Addis, and Sergio  
3181 Uzzau. 2015. 'Enrichment or Depletion? The Impact of Stool Pretreatment on  
3182 Metaproteomic Characterization of the Human Gut Microbiota'. *PROTEOMICS* 15  
3183 (20): 3474–85. <https://doi.org/10.1002/pmic.201400573>.  
3184 Tanca, Alessandro, Antonio Palomba, Salvatore Pisanu, Massimo Deligios, Cristina  
3185 Fraumene, Valeria Manghina, Daniela Pagnozzi, Maria Filippa Addis, and Sergio  
3186 Uzzau. 2014. 'A Straightforward and Efficient Analytical Pipeline for Metaproteome  
3187 Characterization'. *Microbiome* 2 (1): 49. [https://doi.org/10.1186/s40168-014-0049-](https://doi.org/10.1186/s40168-014-0049-2)  
3188 [2](https://doi.org/10.1186/s40168-014-0049-2).  
3189 Tang, Jing, Minjie Mou, Yunxia Wang, Yongchao Luo, and Feng Zhu. 2021. 'MetaFS:  
3190 Performance Assessment of Biomarker Discovery in Metaproteomics'. *Briefings in*  
3191 *Bioinformatics* 22 (3): bbaa105. <https://doi.org/10.1093/bib/bbaa105>.  
3192 The Galaxy Community. 2024. 'The Galaxy Platform for Accessible, Reproducible, and  
3193 Collaborative Data Analyses: 2024 Update'. *Nucleic Acids Research* 52 (W1):  
3194 W83–94. <https://doi.org/10.1093/nar/gkae410>.  
3195 The, Matthew, Michael J. MacCoss, William S. Noble, and Lukas Käll. 2016. 'Fast and

3196 Accurate Protein False Discovery Rates on Large-Scale Proteomics Data Sets  
3197 with Percolator 3.0'. *Journal of the American Society for Mass Spectrometry* 27  
3198 (11): 1719–27. <https://doi.org/10.1007/s13361-016-1460-7>.  
3199 The UniProt Consortium. 2023. 'UniProt: The Universal Protein Knowledgebase in 2023'.  
3200 *Nucleic Acids Research* 51 (D1): D523–31. <https://doi.org/10.1093/nar/gkac1052>.  
3201 Välikangas, Tommi, Tomi Suomi, and Laura L. Elo. 2018. 'A Systematic Evaluation of  
3202 Normalization Methods in Quantitative Label-Free Proteomics'. *Briefings in*  
3203 *Bioinformatics* 19 (1): 1–11. <https://doi.org/10.1093/bib/bbw095>.  
3204 Van Den Bossche, Tim, Magnus Ø. Arntzen, Dörte Becher, Dirk Benndorf, Vincent G. H.  
3205 Eijnsink, Céline Henry, Pratik D. Jagtap, et al. 2021. 'The Metaproteomics Initiative:  
3206 A Coordinated Approach for Propelling the Functional Characterization of  
3207 Microbiomes'. *Microbiome* 9 (1): 243. [https://doi.org/10.1186/s40168-021-01176-](https://doi.org/10.1186/s40168-021-01176-w)  
3208 [w](https://doi.org/10.1186/s40168-021-01176-w).  
3209 Van Den Bossche, Tim, Denis Beslic, Sam van Puyenbroeck, Tomi Suomi, Tanja  
3210 Holstein, Lennart Martens, Laura L. Elo, and Thilo Muth. 2024. 'Metaproteomics  
3211 beyond Databases: Addressing the Challenges and Potentials of de Novo  
3212 Sequencing'. ChemRxiv. <https://doi.org/10.26434/chemrxiv-2024-4v6q0>.  
3213 Van Den Bossche, Tim, Benoit J. Kunath, Kay Schallert, Stephanie S. Schäpe, Paul E.  
3214 Abraham, Jean Armengaud, Magnus Ø Arntzen, et al. 2021. 'Critical Assessment  
3215 of MetaProteome Investigation (CAMPI): A Multi-Laboratory Comparison of  
3216 Established Workflows'. *Nature Communications* 12 (1): 7305.  
3217 <https://doi.org/10.1038/s41467-021-27542-8>.  
3218 Van Den Bossche, Tim, Pieter Verschaffelt, Kay Schallert, Harald Barsnes, Peter  
3219 Dawyndt, Dirk Benndorf, Bernhard Y. Renard, Bart Mesuere, Lennart Martens,  
3220 and Thilo Muth. 2020. 'Connecting MetaProteomeAnalyzer and PeptideShaker to  
3221 Unipept for Seamless End-to-End Metaproteomics Data Analysis'. *Journal of*  
3222 *Proteome Research* 19 (8): 3562–66.  
3223 <https://doi.org/10.1021/acs.jproteome.0c00136>.  
3224 Van Den Bossche, Tim, Pieter Verschaffelt, Tibo Vande Moortele, Peter Dawyndt, Lennart  
3225 Martens, and Bart Mesuere. 2024. 'Biodiversity Analysis of Metaproteomics  
3226 Samples with Unipept: A Comprehensive Tutorial'. In *Protein Bioinformatics*,  
3227 edited by Frédérique Lisacek, 183–215. New York, NY: Springer US.  
3228 [https://doi.org/10.1007/978-1-0716-4007-4\\_11](https://doi.org/10.1007/978-1-0716-4007-4_11).  
3229 Vande Moortele, Tibo, Bram Devlaminck, Simon Van de Vyver, Tim Van Den Bossche,  
3230 Lennart Martens, Peter Dawyndt, Bart Mesuere, and Pieter Verschaffelt. 2024.  
3231 'Unipept 6.0: Expanding Metaproteomics Analysis with Support for Missed  
3232 Cleavages, Semi-Tryptic and Non-Tryptic Peptides'. bioRxiv.  
3233 <https://doi.org/10.1101/2024.09.26.615136>.  
3234 Vande Moortele, Tibo, Pieter Verschaffelt, Qingyao Huang, Nadezhda T. Doncheva,  
3235 Tanja Holstein, Caroline Jachmann, Peter Dawyndt, Lennart Martens, Bart  
3236 Mesuere, and Tim Van Den Bossche. 2024. 'PathwayPilot: A User-Friendly Tool  
3237 for Visualizing and Navigating Metabolic Pathways'. bioRxiv.  
3238 <https://doi.org/10.1101/2024.06.21.599989>.  
3239 Vaudel, Marc, Harald Barsnes, Frode S. Berven, Albert Sickmann, and Lennart Martens.  
3240 2011. 'SearchGUI: An Open-Source Graphical User Interface for Simultaneous  
3241 OMSSA and X!Tandem Searches'. *PROTEOMICS* 11 (5): 996–99.  
3242 <https://doi.org/10.1002/pmic.201000595>.  
3243 Vaudel, Marc, Julia M. Burkhart, René P. Zahedi, Eystein Oveland, Frode S. Berven,  
3244 Albert Sickmann, Lennart Martens, and Harald Barsnes. 2015. 'PeptideShaker  
3245 Enables Reanalysis of MS-Derived Proteomics Data Sets'. *Nature Biotechnology*  
3246 33 (1): 22–24. <https://doi.org/10.1038/nbt.3109>.  
3247 Vaudel, Marc, A. Saskia Venne, Frode S. Berven, René P. Zahedi, Lennart Martens, and  
3248 Harald Barsnes. 2014. 'Shedding Light on Black Boxes in Protein Identification'.  
3249 *PROTEOMICS* 14 (9): 1001–5. <https://doi.org/10.1002/pmic.201300488>.  
3250 Verberkmoes, Nathan C., Alison L. Russell, Manesh Shah, Adam Godzik, Magnus

- 3251 Rosenquist, Jonas Halfvarson, Mark G. Lefsrud, et al. 2009. 'Shotgun  
3252 Metaproteomics of the Human Distal Gut Microbiota'. *The ISME Journal* 3 (2):  
3253 179–89. <https://doi.org/10.1038/ismej.2008.108>.
- 3254 Verheggen, Kenneth, Helge Ræder, Frode S. Berven, Lennart Martens, Harald Barsnes,  
3255 and Marc Vaudel. 2020. 'Anatomy and Evolution of Database Search Engines—a  
3256 Central Component of Mass Spectrometry Based Proteomic Workflows'. *Mass  
3257 Spectrometry Reviews* 39 (3): 292–306. <https://doi.org/10.1002/mas.21543>.
- 3258 Verschaffelt, Pieter, Alessandro Tanca, Marcello Abbondio, Tim Van Den Bossche, Tibo  
3259 Vande Moortele, Peter Dawyndt, Lennart Martens, and Bart Mesuere. 2023.  
3260 'Unipept Desktop 2.0: Construction of Targeted Reference Protein Databases for  
3261 Metaproteogenomics Analyses'. *Journal of Proteome Research* 22 (8): 2620–28.  
3262 <https://doi.org/10.1021/acs.jproteome.3c00091>.
- 3263 Verschaffelt, Pieter, Philippe Van Thienen, Tim Van Den Bossche, Felix Van der Jeugt,  
3264 Caroline De Tender, Lennart Martens, Peter Dawyndt, and Bart Mesuere. 2020.  
3265 'Unipept CLI 2.0: Adding Support for Visualizations and Functional Annotations'.  
3266 *Bioinformatics (Oxford, England)* 36 (14): 4220–21.  
3267 <https://doi.org/10.1093/bioinformatics/btaa553>.
- 3268 Vertommen, Annelies, Bart Panis, Rony Swennen, and Sebastien Christian Carpentier.  
3269 2010. 'Evaluation of Chloroform/Methanol Extraction to Facilitate the Study of  
3270 Membrane Proteins of Non-Model Plants'. *Planta* 231 (5): 1113–25.  
3271 <https://doi.org/10.1007/s00425-010-1121-1>.
- 3272 Waibel, Matthias, Kevin McDonnell, Maria Tuohy, Sally Shirran, Sylvia Synowsky, Barry  
3273 Thornton, Eric Paterson, Fiona Brennan, and Florence Abram. 2023. 'Assessing  
3274 the Impact of Interfering Organic Matter on Soil Metaproteomic Workflow'.  
3275 *European Journal of Soil Science* 74 (3): e13392.  
3276 <https://doi.org/10.1111/ejss.13392>.
- 3277 Wang, Angela, Emily E. F. Fekete, Marybeth Creskey, Kai Cheng, Zhibin Ning, Annabelle  
3278 Pfeifle, Xuguang Li, Daniel Figeys, and Xu Zhang. 2024. 'Assessing Fecal  
3279 Metaproteomics Workflow and Small Protein Recovery Using DDA and DIA  
3280 PASEF Mass Spectrometry'. *Microbiome Research Reports* 3 (3): N/A-N/A.  
3281 <https://doi.org/10.20517/mrr.2024.21>.
- 3282 Wang, Jiaqin, Xu Zhang, Leyuan Li, Zhibin Ning, Janice Mayne, Cian Schmitt-Ulms,  
3283 Krystal Walker, Kai Cheng, and Daniel Figeys. 2020. 'Differential Lysis Approach  
3284 Enables Selective Extraction of Taxon-Specific Proteins for Gut Metaproteomics'.  
3285 *Analytical Chemistry* 92 (7): 5379–86.  
3286 <https://doi.org/10.1021/acs.analchem.0c00062>.
- 3287 Wang, Le-heng, De-Quan Li, Yan Fu, Hai-Peng Wang, Jing-Fen Zhang, Zuo-Fei Yuan,  
3288 Rui-Xiang Sun, Rong Zeng, Si-Min He, and Wen Gao. 2007. 'pFind 2.0: A  
3289 Software Package for Peptide and Protein Identification via Tandem Mass  
3290 Spectrometry'. *Rapid Communications in Mass Spectrometry* 21 (18): 2985–91.  
3291 <https://doi.org/10.1002/rcm.3173>.
- 3292 Wang, Luman, Caitlin M. A. Simopoulos, Joeselle M. Serrana, Zhibin Ning, Boyan Sun,  
3293 Jinhui Yuan, Daniel Figeys, and Leyuan Li. 2024. 'PhyloFunc: Phylogeny-Informed  
3294 Functional Distance as a New Ecological Metric for Metaproteomic Data Analysis'.  
3295 bioRxiv. <https://doi.org/10.1101/2024.05.28.596184>.
- 3296 Wang, Songduo, Zenghu Zhang, Kaiguang Yang, Jiulong Zhao, Weijie Zhang, Zhiting  
3297 Wang, Zhen Liang, et al. 2024. 'SMMP: A Deep-Coverage Marine Metaproteome  
3298 Method for Microbial Community Analysis throughout the Water Column Using 1 L  
3299 of Seawater'. *Analytical Chemistry* 96 (29): 12030–39.  
3300 <https://doi.org/10.1021/acs.analchem.4c02079>.
- 3301 Wang, Tong, Leyuan Li, Daniel Figeys, and Yang-Yu Liu. 2024. 'Pairing Metagenomics  
3302 and Metaproteomics to Characterize Ecological Niches and Metabolic Essentiality  
3303 of Gut Microbiomes'. *ISME Communications* 4 (1): ycae063.  
3304 <https://doi.org/10.1093/ismeco/ycae063>.
- 3305 Wang, Yuqiu, Yanting Zhou, Xiao Xiao, Jing Zheng, and Hu Zhou. 2020. 'Metaproteomics:

3306 A Strategy to Study the Taxonomy and Functionality of the Gut Microbiota'.  
 3307 *Journal of Proteomics* 219 (May):103737.  
 3308 <https://doi.org/10.1016/j.jprot.2020.103737>.  
 3309 Washburn, Michael P., Dirk Wolters, and John R. Yates. 2001. 'Large-Scale Analysis of  
 3310 the Yeast Proteome by Multidimensional Protein Identification Technology'. *Nature*  
 3311 *Biotechnology* 19 (3): 242–47. <https://doi.org/10.1038/85686>.  
 3312 Wessel, D., and U. I. Flügge. 1984. 'A Method for the Quantitative Recovery of Protein in  
 3313 Dilute Solution in the Presence of Detergents and Lipids'. *Analytical Biochemistry*  
 3314 138 (1): 141–43. [https://doi.org/10.1016/0003-2697\(84\)90782-6](https://doi.org/10.1016/0003-2697(84)90782-6).  
 3315 Wilkins, Marc R., Christian Pasquali, Ron D. Appel, Keli Ou, Olivier Golaz, Jean-Charles  
 3316 Sanchez, Jun X. Yan, et al. 1996. 'From Proteins to Proteomes: Large Scale  
 3317 Protein Identification by Two-Dimensional Electrophoresis and Amino Acid  
 3318 Analysis'. *Bio/Technology* 14 (1): 61–65. <https://doi.org/10.1038/nbt0196-61>.  
 3319 Wilmes, Paul, Anna Heintz-Buschart, and Philip L. Bond. 2015. 'A Decade of  
 3320 Metaproteomics: Where We Stand and What the Future Holds'. *PROTEOMICS* 15  
 3321 (20): 3409–17. <https://doi.org/10.1002/pmic.201500183>.  
 3322 Wiśniewski, Jacek R., Alexandre Zougman, Nagarjuna Nagaraj, and Matthias Mann.  
 3323 2009. 'Universal Sample Preparation Method for Proteome Analysis'. *Nature*  
 3324 *Methods* 6 (5): 359–62. <https://doi.org/10.1038/nmeth.1322>.  
 3325 Wolf, Maximilian, Kay Schallert, Luca Knipper, Albert Sickmann, Alexander Sczyrba, Dirk  
 3326 Benndorf, and Robert Heyer. 2023. 'Advances in the Clinical Use of  
 3327 Metaproteomics'. *Expert Review of Proteomics* 20 (4–6): 71–86.  
 3328 <https://doi.org/10.1080/14789450.2023.2215440>.  
 3329 Wood, Derrick E., Jennifer Lu, and Ben Langmead. 2019. 'Improved Metagenomic  
 3330 Analysis with Kraken 2'. *Genome Biology* 20 (1): 257.  
 3331 <https://doi.org/10.1186/s13059-019-1891-0>.  
 3332 Wu, Enhui, Guanyang Xu, Dong Xie, and Liang Qiao. 2024. 'Data-Independent  
 3333 Acquisition in Metaproteomics'. *Expert Review of Proteomics* 21 (7–8): 271–80.  
 3334 <https://doi.org/10.1080/14789450.2024.2394190>.  
 3335 Wu, Qing, Zhibin Ning, Ailing Zhang, Xu Zhang, Zhongzhi Sun, and Daniel Figeys. 2024.  
 3336 'MetaX: A Peptide Centric Metaproteomic Data Analysis Platform Using  
 3337 Operational Taxa-Functions (OTF)'. *bioRxiv*.  
 3338 <https://doi.org/10.1101/2024.04.19.590315>.  
 3339 Wu, Qiong, Jiangnan Zheng, Xintong Sui, Changying Fu, Xiaozhen Cui, Bin Liao,  
 3340 Hongchao Ji, et al. 2024. 'High-Throughput Drug Target Discovery Using a Fully  
 3341 Automated Proteomics Sample Preparation Platform'. *Chemical Science* 15 (8):  
 3342 2833–47. <https://doi.org/10.1039/D3SC05937E>.  
 3343 Xiao, Xiaolian, Haidan Sun, Xiaoyan Liu, Zhengguang Guo, Shuxin Zheng, Jiyu Xu,  
 3344 Jiameng Sun, Ying Lan, Chen Shao, and Wei Sun. 2022. 'Qualitative and  
 3345 Quantitative Proteomic and Metaproteomic Analyses of Healthy Human Urine  
 3346 Sediment'. *PROTEOMICS – Clinical Applications* 16 (2): 2100007.  
 3347 <https://doi.org/10.1002/prca.202100007>.  
 3348 Xiao, Xiaolian, Xiaoping Xiao, Yaoran Liu, Haidan Sun, Xiaoyan Liu, Zhengguang Guo,  
 3349 Qian Li, and Wei Sun. 2023. 'Metaproteomics Characterizes the Human Gingival  
 3350 Crevicular Fluid Microbiome Function in Periodontitis'. *Journal of Proteome*  
 3351 *Research* 22 (7): 2411–20. <https://doi.org/10.1021/acs.jproteome.3c00143>.  
 3352 Xie, Fei, Wei Jin, Huazhe Si, Yuan Yuan, Ye Tao, Junhua Liu, Xiaoxu Wang, et al. 2021.  
 3353 'An Integrated Gene Catalog and over 10,000 Metagenome-Assembled Genomes  
 3354 from the Gastrointestinal Microbiome of Ruminants'. *Microbiome* 9 (1): 137.  
 3355 <https://doi.org/10.1186/s40168-021-01078-x>.  
 3356 Xiong, Weili, Christopher T. Brown, Michael J. Morowitz, Jillian F. Banfield, and Robert L.  
 3357 Hettich. 2017. 'Genome-Resolved Metaproteomic Characterization of Preterm  
 3358 Infant Gut Microbiota Development Reveals Species-Specific Metabolic Shifts and  
 3359 Variabilities during Early Life'. *Microbiome* 5 (1): 72.  
 3360 <https://doi.org/10.1186/s40168-017-0290-6>.



- 3361 Xiong, Weili, Richard J. Giannone, Michael J. Morowitz, Jillian F. Banfield, and Robert L.  
3362 Hettich. 2015. 'Development of an Enhanced Metaproteomic Approach for  
3363 Deepening the Microbiome Characterization of the Human Infant Gut'. *Journal of*  
3364 *Proteome Research* 14 (1): 133–41. <https://doi.org/10.1021/pr500936p>.
- 3365 Xu, Ping, Duc M. Duong, and Junmin Peng. 2009. 'Systematical Optimization of Reverse-  
3366 Phase Chromatography for Shotgun Proteomics'. *Journal of Proteome Research* 8  
3367 (8): 3944–50. <https://doi.org/10.1021/pr900251d>.
- 3368 Xue, Ming-Yuan, Yun-Yi Xie, Xin-Wei Zang, Yi-Fan Zhong, Xiao-Jiao Ma, Hui-Zeng Sun,  
3369 and Jian-Xin Liu. 2024. 'Deciphering Functional Groups of Rumen Microbiome and  
3370 Their Underlying Potentially Causal Relationships in Shaping Host Traits'. *iMeta* 3  
3371 (4): e225. <https://doi.org/10.1002/imt2.225>.
- 3372 Yang, Liang, Wenlai Fan, and Yan Xu. 2020. 'Metaproteomics Insights into Traditional  
3373 Fermented Foods and Beverages'. *Comprehensive Reviews in Food Science and*  
3374 *Food Safety* 19 (5): 2506–29. <https://doi.org/10.1111/1541-4337.12601>.
- 3375 Yang, Tingpeng, Tianze Ling, Boyan Sun, Zhendong Liang, Fan Xu, Xiansong Huang,  
3376 Linhai Xie, et al. 2024. 'Introducing  $\pi$ -HelixNovo for Practical Large-Scale de Novo  
3377 Peptide Sequencing'. *Briefings in Bioinformatics* 25 (2): bbae021.  
3378 <https://doi.org/10.1093/bib/bbae021>.
- 3379 Yilmaz, Melih, William E. Fondrie, Wout Bittremieux, Carlo F. Melendez, Rowan Nelson,  
3380 Varun Ananth, Sewoong Oh, and William Stafford Noble. 2024. 'Sequence-to-  
3381 Sequence Translation from Mass Spectra to Peptides with a Transformer Model'.  
3382 bioRxiv. <https://doi.org/10.1101/2023.01.03.522621>.
- 3383 Zhang, Weipeng, Jin Sun, Huiluo Cao, Renmao Tian, Lin Cai, Wei Ding, and Pei-Yuan  
3384 Qian. 2016. 'Post-Translational Modifications Are Enriched within Protein  
3385 Functional Groups Important to Bacterial Adaptation within a Deep-Sea  
3386 Hydrothermal Vent Environment'. *Microbiome* 4 (1): 49.  
3387 <https://doi.org/10.1186/s40168-016-0194-x>.
- 3388 Zhang, Xu, Wendong Chen, Zhibin Ning, Janice Mayne, David Mack, Alain Stintzi, Ruijun  
3389 Tian, and Daniel Figeys. 2017. 'Deep Metaproteomics Approach for the Study of  
3390 Human Microbiomes'. *Analytical Chemistry* 89 (17): 9407–15.  
3391 <https://doi.org/10.1021/acs.analchem.7b02224>.
- 3392 Zhang, Xu, Kai Cheng, Zhibin Ning, Janice Mayne, Krystal Walker, Hao Chi, Charles L.  
3393 Farnsworth, Kimberly Lee, and Daniel Figeys. 2021. 'Exploring the Microbiome-  
3394 Wide Lysine Acetylation, Succinylation, and Propionylation in Human Gut  
3395 Microbiota'. *Analytical Chemistry* 93 (17): 6594–98.  
3396 <https://doi.org/10.1021/acs.analchem.1c00962>.
- 3397 Zhang, Xu, and Daniel Figeys. 2019. 'Perspective and Guidelines for Metaproteomics in  
3398 Microbiome Studies'. *Journal of Proteome Research* 18 (6): 2370–80.  
3399 <https://doi.org/10.1021/acs.jproteome.9b00054>.
- 3400 Zhang, Xu, Leyuan Li, Janice Mayne, Zhibin Ning, Alain Stintzi, and Daniel Figeys. 2018.  
3401 'Assessing the Impact of Protein Extraction Methods for Human Gut  
3402 Metaproteomics'. *Journal of Proteomics, Proteomics in Infectious Diseases*, 180  
3403 (May):120–27. <https://doi.org/10.1016/j.jprot.2017.07.001>.
- 3404 Zhang, Xu, Zhibin Ning, Janice Mayne, Jasmine I. Moore, Jennifer Li, James Butcher,  
3405 Shelley Ann Deeke, et al. 2016. 'MetaPro-IQ: A Universal Metaproteomic  
3406 Approach to Studying Human and Mouse Gut Microbiota'. *Microbiome* 4 (1): 31.  
3407 <https://doi.org/10.1186/s40168-016-0176-z>.
- 3408 Zhang, Xu, Zhibin Ning, Janice Mayne, Yidai Yang, Shelley A. Deeke, Krystal Walker,  
3409 Charles L. Farnsworth, et al. 2020. 'Widespread Protein Lysine Acetylation in Gut  
3410 Microbiome and Its Alterations in Patients with Crohn's Disease'. *Nature*  
3411 *Communications* 11 (1): 4120. <https://doi.org/10.1038/s41467-020-17916-9>.
- 3412 Zhao, Jinzhi, Yi Yang, Liangqiang Chen, Jianxujie Zheng, Xibin Lv, Dandan Li, Ziyu Fang,  
3413 et al. 2023. 'Quantitative Metaproteomics Reveals Composition and Metabolism  
3414 Characteristics of Microbial Communities in Chinese Liquor Fermentation  
3415 Starters'. *Frontiers in Microbiology* 13 (January).

3416 <https://doi.org/10.3389/fmicb.2022.1098268>.  
3417 Zhao, Jinzhi, Yi Yang, Mengjing Teng, Jianxujie Zheng, Bing Wang, Vijini  
3418 Mallawaarachchi, Yu Lin, et al. 2023. 'Metaproteomics Profiling of the Microbial  
3419 Communities in Fermentation Starters (Daqu) during Multi-Round Production of  
3420 Chinese Liquor'. *Frontiers in Nutrition* 10 (June).  
3421 <https://doi.org/10.3389/fnut.2023.1139836>.  
3422 Zhao, Jinzhi, Yi Yang, Hua Xu, Jianxujie Zheng, Chengpin Shen, Tian Chen, Tao Wang,  
3423 et al. 2023. 'Data-Independent Acquisition Boosts Quantitative Metaproteomics for  
3424 Deep Characterization of Gut Microbiota'. *Npj Biofilms and Microbiomes* 9 (1): 1–  
3425 14. <https://doi.org/10.1038/s41522-023-00373-9>.