

1 **Special ISSUE:** "The de.NBI Network - Software Tools for Big Data Analysis in Life Sciences"

2 **TITLE: Challenges and perspectives of metaproteomic data analysis**

3 Robert Heyer<sup>1\*</sup>, Kay Schallert<sup>2</sup>, Roman Zoun<sup>3</sup>, Beatrice Becher<sup>4</sup>, Gunter Saake<sup>5</sup>, Dirk Benndorf<sup>6\*</sup>

4 **1.** Otto von Guericke University, Bioprocess Engineering, Universitätsplatz 2, 39106 Magdeburg,  
5 Germany; heyer@mpi-magdeburg.mpg.de (\*: [corresponding author](#))

6 **2.** Otto von Guericke University, Bioprocess Engineering, Universitätsplatz 2, 39106 Magdeburg,  
7 Germany; kay.schallert@ovgu.de

8 **3.** Otto von Guericke University, Institute for Technical and Business Information Systems,  
9 Universitätsplatz 2, 39106 Magdeburg, Germany; roman.zoun@ovgu.de

10 **4.** Otto von Guericke University, Bioprocess Engineering, Universitätsplatz 2, 39106 Magdeburg,  
11 Germany; beatrice.becher@st.ovgu.de

12 **5.** Otto von Guericke University, Institute for Technical and Business Information Systems,  
13 Universitätsplatz 2, 39106 Magdeburg, Germany; saake@iti.cs.uni-magdeburg.de

14 **6.** Otto von Guericke University, Bioprocess Engineering, Universitätsplatz 2, 39106 Magdeburg,  
15 Germany; Max Planck Institute for Dynamics of Complex Technical Systems, Bioprocess  
16 Engineering, Sandtorstraße 1, 39106 Magdeburg, Germany; benndorf@mpi-magdeburg.mpg.de  
17 (\*: [corresponding author](#))

18

19 **Abstract**

20 In nature microorganisms live in complex microbial communities. Comprehensive taxonomic and  
21 functional knowledge about microbial communities supports medical and technical application such as  
22 fecal diagnostics as well as operation of biogas plants or waste water treatment plants. Furthermore,  
23 microbial communities are crucial for the global carbon and nitrogen cycle in soil and in the ocean.  
24 Among the methods available for investigation of microbial communities, metaproteomics can  
25 approximate the activity of microorganisms by investigating the protein content of a sample. Although  
26 metaproteomics is a very powerful method, issues within the bioinformatic evaluation impede its  
27 success. In particular, construction of databases for protein identification, grouping of redundant  
28 proteins as well as taxonomic and functional annotation pose big challenges. Furthermore, growing  
29 amounts of data within a metaproteomics study require dedicated algorithms and software. This review  
30 summarizes recent metaproteomics software and addresses the introduced issues in detail.

31

32 **A. Highlights**

- 33 • **Metaproteomic studies profit from dedicated software tools**
- 34 • **Metagenomes and protein database constraints improve protein identification**
- 35 • **Grouping of proteins by shared peptides or sequence similarity reduce redundancy**
- 36 • **Several possibilities for taxonomic and functional classification of proteins exist**
- 37 • **Scalability of software and databases enables handling of big data amounts**

38

39 **B. Keywords**

- 40 • Bioinformatics
- 41 • Software
- 42 • Big data
- 43 • Environmental proteomics
- 44 • Microbial communities
- 45 • Mass spectrometry

46

47 **C. Content**

48 **1. Introduction**

49	<b>2.</b> <i>Status of proteomic software and latest trends</i>
50	<b>3.</b> <i>Software dedicated for metaproteomics</i>
51	<b>4.</b> <i>Construction of user databases for protein identification</i>
52	<b>5.</b> <i>Construction of user databases for protein identification: a use case</i>
53	<b>6.</b> <i>Protein inference problem and the grouping of proteins into “metaproteins”</i>
54	<b>7.</b> <i>Taxonomic and functional result evaluation</i>
55	<b>8.</b> <i>Quantitative data analysis in metaproteome studies</i>
56	<b>9.</b> <i>Strategies for storing and deployment of huge data</i>
57	<b>10.</b> <i>Future challenges, perspectives and demands</i>
58	<b>11.</b> <i>Conclusion</i>
59	

## 60 1. Introduction

61 Microorganisms represent 50–78% of Earth's total biomass ([Kallmeyer et al., 2012](#)) and occur in all  
62 environments. Some microorganisms produce biomass by photosynthesis whereas others act as  
63 composers and degrade dead biomass. Microbial species live in complex microbial communities in which  
64 they have to compete or cooperate with each other. Understanding the functioning of the microbial  
65 communities is important, because microbial communities in the human gut effect health ([Erickson et  
66 al., 2012](#); [Heintz-Buschart et al., 2016](#); [Kolmeder et al., 2016](#)) and several technical applications such as  
67 waste water treatment plants ([Püttker et al., 2015](#); [Wilmes et al., 2008](#)) and biogas plants ([Abram et al.,  
68 2011](#); [Hanreich et al., 2012](#)) rely on the metabolic activity of microbial communities.

69 Methods for the investigation of microbial communities target the microbial cells, their genes, their  
70 transcripts, their proteins and their metabolites ([Heyer et al., 2015](#)). Since proteins carry out most  
71 functions in cells, including catalysis of biochemical reactions, transport and cell structure, protein  
72 amounts correlate quite well with microbial activity ([Wilmes and Bond, 2006](#)). The investigation of all  
73 proteins from one species is called proteomics. In contrast metaproteomics is the study of proteins from  
74 multiple organisms. It was introduced by [Wilmes et al. \(2006+2004\)](#) and [Rodriguez-Valera \(2004\)](#). The  
75 typical metaproteomics workflow comprises protein extraction and purification, tryptic digestion into  
76 peptides, protein or peptide separation and tandem mass spectrometry (MS/MS) analysis. Proteins are  
77 identified by comparing experimental mass spectra and theoretical mass spectra predicted from  
78 comprehensive protein databases. For a detailed discussion about the metaproteomics workflow please  
79 refer to [Hettich et al. \(2013\)](#), [Becher et al. \(2013\)](#), [Heyer et al. \(2015\)](#), [Wöhlbrand et al. \(2013\)](#). Up to now  
80 most metaproteomics studies characterize the taxonomic and functional composition of complex  
81 microbial communities in their specific environment ([Abram et al., 2011](#); [Kan et al., 2005](#); [Ram et al.,  
82 2005](#); [Wilmes and Bond, 2006](#)). A few recent studies additionally correlated the taxonomic and  
83 functional composition with certain environmental/process parameters or diseases ([Erickson et al., 2012](#);  
84 [Heyer et al., 2016](#); [Kolmeder et al., 2016](#)). However, three issues within bioinformatic data evaluation  
85 hampered previous metaproteomics studies ([Muth et al., 2013](#)).

86 First, metaproteomes consist of up to 1,000 different species ([Schlüter et al., 2008](#)). Due to high  
87 complexity metaproteomics data analysis requires a greater computational effort, necessitating bigger  
88 hard drives, more memory, more processors and more efficient algorithms. A main issue is the database  
89 search against comprehensive protein databases. Whereas handling of small protein databases below 1  
90 GB is not critical, usage of the entire NCBI reference database requires extended computational time and  
91 may fail due to software or hardware limitations.

92 Second, identical peptides belonging to homologous proteins cause redundant protein identification  
93 ([Herbst et al., 2016](#)). As a result taxonomic and functional interpretation of results becomes ambiguous.  
94 A peptide may belong to the lactate dehydrogenase (1.1.1.27) of different members of the genus  
95 *Lactobacillus*, which ferment sugars to lactate. But it may also belong to some representatives of the  
96 order *Clostridiales fermenting* lactate to acetate ([Kohrs et al., 2014](#)).

97 Third, protein identification is difficult if the taxonomic composition is unknown or protein entries are  
98 missing from protein databases. For example the UniProt/TrEMBL database contains only proteins from

99 698,745 species (<http://www.ebi.ac.uk/uniprot/TrEMBLstats>, status 16.12.2016), but the number of  
100 microbial species on Earth is estimated to be up to one trillion (Locey and Lennon, 2016). Thereby,  
101 already small changes in the protein sequence between related microorganisms have a big impact on  
102 protein identification. One mutation in every tenth amino acid leads to completely different tryptic  
103 peptides which hinder the identification of any peptide for the investigated protein. ~~Since protein~~  
104 ~~identification relies on this sequence information,~~ Thus, researchers started to sequence metagenomes  
105 alongside metaproteomics studies (Ram et al., 2005; Tyson et al., 2004). Alternatively, they use  
106 metagenomes from similar samples for protein identification.

107 As a consequence of these issues, standard proteomics software is often insufficient for metaproteomics  
108 studies ~~missing the identification of unsequenced species or the comprehensive taxonomic and~~  
109 ~~functional description of microbial communities.~~ Thus, researchers favor special tools. Therefore, this  
110 review provides an overview about dedicated metaproteomics software and bioinformatic strategies.

111 ~~In addition to two previous reviews on bioinformatics in metaproteomics (Muth et al., 2013 +2016) we~~  
112 ~~present the impact of combining metagenomes on protein identification and address future hardware~~  
113 ~~requirements and handling of big data.~~

114 After a brief introduction to ~~metaproteomics studies and the state of proteomics software,~~ current  
115 metaproteomics software tools are ~~discussed.~~ Subsequently, this review ~~illuminates~~ the creation of  
116 protein databases for protein identification ~~investigating several biogas plant samples in a use case.~~ Then  
117 the grouping of redundant protein identifications, the evaluation of taxonomic and functional results ~~as~~  
118 ~~well as quantification in metaproteomics studies are discussed.~~ Finally, data storage and deployment  
119 solutions for big data as well as future challenges, perspectives and demand for metaproteomics  
120 software are considered.

## 121 ~~2. Brief introduction into the workflow of metaproteomic studies~~

122 ~~The following section briefly introduces standard metaproteomics workflows. Detailed discussions are~~  
123 ~~provided by Hettich et al. (2013), Becher et al. (2013) and Wöhlbrand et al. (2013). First, microbial cells~~  
124 ~~are lysed [Figure 1 A], using e.g. a ball mill or ultrasonic sound. Afterwards several centrifugation,~~  
125 ~~precipitation or extraction steps isolate the proteins from cell debris, DNA and the sample matrix. Protein~~  
126 ~~quantification assays determine the amount of extracted proteins. Of these the amido black assay~~  
127 ~~appears to be the most robust for samples containing impurities (Hanreich et al., 2013; Racusen, 1973).~~

128 ~~Subsequently, different fractionation steps reduce the sample complexity [Figure 1 C]. SDS PAGE~~  
129 ~~(Laemmli, 1970) and 2D PAGE (Klose, 1975; O'Farrell, 1975) separate proteins by their molecular weight~~  
130 ~~and isoelectric point. Alternatively, the protease trypsin cleaves proteins into peptides and one or two-~~  
131 ~~dimensional liquid chromatography (LC) separates the peptides according to their biochemical properties~~  
132 ~~[Figure 1 C]. Recently, researchers use a combination of both approaches. First, SDS PAGE separates the~~  
133 ~~proteins in several fractions, followed by tryptic digestion and reversed phase LC (Wilm et al., 1996).~~

134 ~~The LC is usually coupled online to the mass spectrometer (MS) [Figure 1 D]. A MS is a complex device~~  
135 ~~that determines the mass-to-charge ratio (m/z ratio) of each peptide as well as its quantity. After~~  
136 ~~ionization of peptides by the ion source, the mass analyzer separates peptide ions according to their~~

137 m/z-ratio before they are separately registered by the detector. Peptides eluting from the LC are  
138 continuously measured by the MS. Peptide ion intensities and their m/z ratios constitute the peptide ion  
139 spectrum or MS-spectrum for a given retention time.

140 The m/z-ratio of a peptide ion is quite specific for the masses of its amino acids, but not for the sequence  
141 of the amino acids. Thus, peptide ions are fragmented further to derive sequence information. This step  
142 is called tandem mass spectrometry (MS/MS) and results in the fragment ion spectrum or MS/MS-  
143 spectrum, which comprises the sequence of fragment ions. The m/z-difference between adjacent  
144 fragment ions represents a single amino acid. Consequently, a series of fragment ions reveals the  
145 peptide sequence. Precise m/z-ratio values specify the quality of MS-measurements. Recently, Orbitrap-  
146 MS deliver excellent accuracy ([Hu et al., 2005](#)).

147 Following the experiments, bioinformatic analyses are used to identify the proteins and help to evaluate  
148 protein significance. Database search algorithms for protein identification such as MS Amanda ([Dorfer et  
149 al., 2014](#)), MASCOT ([Perkins et al., 1999](#)) and XTANDEM ([Craig and Beavis, 2004](#)) [Figure 1 E] calculate the  
150 theoretical spectra for all proteins in a protein database [Figure 1 F] and compare these spectra against  
151 the measured MS/MS-spectrum.

152 As a result these algorithms provide a possible peptide for each MS/MS-spectrum, as well as the  
153 probability of the identification. In the next step, database search algorithms connect peptides to  
154 proteins. Some algorithms return only a single protein deemed best, while others return the entire list of  
155 proteins containing this peptide ([Muth et al., 2016](#)).

156 In addition to the probability of a spectrum identification, the false discovery rate (FDR) has evolved as  
157 the standard to evaluate identification quality ([Elias and Gygi, 2007](#)). The FDR can be calculated as the  
158 ratio of all spectra identified using a decoy database divided by the number of identified spectra  
159 using both, the original and the decoy database. The decoy database is a shuffled version of the original  
160 database ([Colaert et al., 2011](#)), which is supposed to contain false protein sequences, only.

161 In order to make the result evaluation more meaningful, the identified proteins are linked with their  
162 taxonomy and function [Figure 1 G, H]. Several systems are available to provide functional annotation of  
163 proteins, which will be discussed in detail later.

164 The next step in the metaproteome workflow is protein quantification. Different approaches for  
165 quantitative proteomics exist ([Vaudel et al., 2010](#)). In metaproteomics, protein amounts are often  
166 estimated by peptide count ([Ishihama et al., 2005](#)), spectral count ([Zybailov et al., 2007](#)) or peptide peak  
167 area ([Griffin et al., 2010](#)). To determine microbial activity and interaction, researchers can feed microbial  
168 communities with isotope labeled substrates. The incorporation of isotopes into proteins is measured via  
169 MS (Protein SIP ([Jehmlich et al., 2009](#); [Jehmlich et al., 2016](#))).

170 Finally the results of metaproteome studies are visualized in different ways, which were already  
171 reviewed by [Mehlan et al. \(2013\)](#) and [Oveland et al. \(2015\)](#) [Figure 1 I]. In summary, new visualization  
172 concepts for complex data improve data evaluation of metaproteomics studies. For example, Krona plots  
173 show the taxonomic profile for all taxonomic ranks simultaneously ([Ondov et al., 2011](#)). Voronoi  
174 treemaps highlight alterations of the protein expression sorted by protein functions ([Mehlan et al., 2013](#)).

175 Longterm storage and access of all MS files is archived by online repositories such as PRIDE (Vizcaino et  
176 al., 2016) [Figure 1 J].

## 177 **2. Status of proteomics software and latest trends**

178 For the comprehensive bioinformatic processing of MS data different software tools exist. These include  
179 software for peak picking in MS-spectra, software for protein identification via database search  
180 algorithms and tools for comparison of protein expression patterns. A comprehensive summary of all  
181 these software tools can be found in the OMIC tools database (<http://omictools.com/>, retrieved: 09-02-  
182 2017, (Henry et al., 2014)) and in several reviews (Cappadona et al., 2012; Gonzalez-Galarza et al., 2012).

183 Latest trends in proteomics software are the development of proteomics tool libraries such as OpenMS  
184 (Sturm et al., 2008), Compomics (Barsnes et al., 2011) or Trans-Proteomic Pipeline (Keller and  
185 Shteynberg, 2011). These libraries comprise software tools for each step of the processing workflow,  
186 ranging from data management to data analysis. Noteworthy are also webservices, such as ExPASy  
187 (Gasteiger et al., 2003), which provide a collection of small bioinformatic tools for biochemical analyses  
188 of proteins.

189 Repositories for MS-data such as PRIDE are used to enable long-term storage and to make published MS-  
190 data available to other researchers (Vizcaino et al., 2016). In this context general formats for exchange of  
191 MS results are necessary. **Current standard in the proteomics community are the mzIdentML format**  
192 **(Jones et al., 2012, mzTab format (Griss et al., 2014) and mzML format (Martens et al., 2011).**

193 Recent proteomics software combines several database search algorithms. For example, the SearchGUI  
194 tool (Vaudel et al., 2011) enables the parallel protein database search with eight different database  
195 search algorithms. Further developments are software tools for improved MS-operation and  
196 quantification. **Search items for these developments are “data independent acquisition” (Doerr, 2015),**  
197 **“multiple and single reaction monitoring” (Colangelo et al., 2013) as well as “absolute quantification”**  
198 **(Cappadona et al., 2012). However, a detailed discussion of these applications exceeds the scope of this**  
199 **review.**

200 Within the last years many powerful software tools were developed but their use was often restricted to  
201 a few scientific groups. Reasons were missing maintenance or availability after funding periods ended.  
202 Furthermore, many biological research groups lack bioinformatic skills to set up comprehensive software  
203 workflows or client-server architectures. In some cases even the conversion of data into the required  
204 input formats fail. In order to tackle these problems governments started to fund the collection,  
205 maintenance and support of research software tools. Examples are the Galaxy project  
206 (<https://usegalaxy.org/>, retrieved: 09-02-2017, (Afgan et al., 2016), ELIXIR ([https://www.elixir-](https://www.elixir-europe.org/)  
207 [europe.org/](https://www.elixir-europe.org/), retrieved: 09-02-2017, (Crosswell and Thornton, 2012)) or de.NBI (<https://www.denbi.de/>,  
208 retrieved: 09-02-2017).

## 209 **4. Software dedicated for metaproteomics**

210 To address the three issues specific to metaproteomics bioinformatic data evaluation, researchers  
211 started to develop special software tools and workflows [Table1, Figure 1]. These tools apply different  
212 concepts, which will be discussed later. Graph2Pep/Graph2Pro (Tang et al., 2016) and Compile

213 ([Chatterjee et al., 2016](#)) focus on tailoring protein databases for optimal protein identification. UniPept  
214 ([Mesuere et al., 2015](#)), ProPhane ([Schneider et al., 2011](#)), Megan CE ([Huson et al., 2016](#)) and Pipasic  
215 ([Penzlin et al., 2014](#)) enable taxonomic analysis, functional data evaluation and/or protein grouping.  
216 Additionally, several groups assembled comprehensive software workflows for metaproteomics, e.g.  
217 Galaxy-P ([Jagtap et al., 2015](#)), MetaPro-IQ ([Zhang et al., 2016](#)), MetaProteomeAnalyzer ([Muth et al.,  
218 2015a](#)) and others ([Heintz-Buschart et al., 2016](#); [May et al., 2016](#); [Tanca et al., 2013](#)). Among these  
219 workflows, the MPA is particularly user-friendly. It allows the user to control the entire bioinformatic  
220 workflow via an intuitive graphical user interface. Another noteworthy metaproteomics software tool is  
221 MetaProSIP ([Sachsenberg et al., 2015](#)). It supports the detection and quantification of isotope ratios for  
222 Protein-SIP experiments.

223 To ensure comparability of results between all these tools, standards for data exchange are crucial  
224 ([Timmins-Schiffman et al., 2017](#)). Consequentially, the Human Proteomics Standard Initiative is planning  
225 to extend the proteomics mzIdentML format in order to support metaproteomics data. Version 1.2.0 of  
226 the mzIdentML format ([Jones et al., 2012](#)) will support the representation of redundant protein groups  
227 (<http://www.psidev.info/mzidentml>, retrieved: 09-02-2017).

228 Another often neglected aspect is the reproducibility of results using different metaproteomics software  
229 tools. So far, only *Tanca et al. (2013)* tested their complete metaproteomics workflow for a defined  
230 mixed culture of nine different microorganisms. A comparison where multiple research groups evaluate  
231 an identical sample would also be desirable.

### 232 **5. Construction of user databases for protein identification**

233 Protein database selection affects the number of identified proteins as well as the identified taxonomies  
234 and identification increases. In consequence, the estimated FDR and thus, the threshold for accepting  
235 protein identifications are higher and may lead to the rejection of true protein identifications.

236 Optimal databases would only include proteins and posttranslational modifications present in the  
237 sample and detectable by MS. However, taxonomic composition and protein abundance are usually  
238 unknown for environmental samples. Furthermore, protein content between analyzed samples may  
239 differ significantly. Therefore, database selection is a challenging task ([Muth et al., 2015b](#); [Tanca et al.,  
240 2016](#)). This issue is further complicated by the adherence of the research community to the FDR concept  
241 ([Muth et al., 2015b](#)).

242 Originally *Elias et al. (2007)* established the FDR concept for comparable protein identification in pure  
243 culture proteomics. In particular, the FDR enables comparability between different mass spectrometers  
244 and database search algorithms. Subsequently, the proteomics community accepted the FDR calculation  
245 as the standard to control the quality of protein identifications. An FDR of 1% was defined as threshold  
246 ([Barnouin, 2011](#)). However, a condition for the successful estimation of the FDR is that the database fits  
247 well to the sample. This is not guaranteed for metaproteomics studies, resulting in inaccurate  
248 approximations of the FDR. Therefore, it would be desirable that the metaproteomics community revises  
249 the FDR concept questioning the decoy based approach. Instead protein identifications could be  
250 classified using machine learning approaches.

251 Principally researches have two options to construct their database for metaproteomics studies. The first  
252 strategy is to sequence the whole metagenome or metatranscriptome [Figure 2A] ([Ram et al., 2005](#);



253 [Tyson et al., 2004](#)) and to translate the genes to proteins by tools such as Transeq or Sixpack  
254 (<http://www.ebi.ac.uk/Tools/st/>, retrieved 07.06.2017). The second is to use comprehensive sequence  
255 databases [Figure 2\_1] and apply reasonable constraints. Recently, sequencing of metagenomes became  
256 affordable, due to high-throughput sequencing technologies such as Illumina sequencing ([Bentley et al.,](#)  
257 [2008](#); [Jünemann et al., 2014](#); [Jünemann et al., 2013](#)). However, several different processing states of  
258 metagenomes could be used as protein databases [Figure 2A]. After Illumina sequencing and quality  
259 control, metagenome data are present as reads. Reads are short fragments of about 150 base pairs,  
260 which can be translated into about 50 amino acids [Figure 2B]. Subsequently, the translated reads are  
261 assembled to contigs and redundant reads are removed [Figure 2C]. **Contigs may contain several genes.**  
262 In some high resolution metagenome studies, it is even possible to assemble the entire genome of single  
263 microorganisms ([Campanaro et al., 2016](#)). The disadvantage of reads and contigs is that all six reading  
264 frames are considered during the translation of DNA sequences into protein sequences. This multiplies  
265 the amount of data by six. Contigs may also contain several genes, which complicates the taxonomic and  
266 functional interpretation. Hence, genes are predicted from the contigs and non-coding DNA fragments  
267 are removed [Figure 2D]. Therefore, assembled metagenomes with gene predictions are the preferable  
268 databases for protein identification. **Sometimes it is even possible to reconstruct the whole genome of**  
269 **single microorganisms within the microbial community, which is called binning.**

270 Since these assembled metagenome protein databases match the actual sample, FDR estimation should  
271 be valid. However, the bioinformatic workflow to assemble metagenomes can also influence the protein  
272 identification ([Tanca et al., 2016](#)). For example, during metagenome assembly redundant reads where  
273 only one amino acid differs are sometimes condensed into a single read. This ignores protein isoforms  
274 and can lead to the loss of protein identifications. In contrast, a high number of translated reads in a  
275 database decrease protein identifications due to an increase in the FDR. In line with these problems,  
276 some authors experienced a higher number of protein identifications with read databases instead of  
277 contig databases ([Timmins-Schiffman et al., 2017](#)). Better protein identification was also observed by  
278 [Tang et al. 2016](#) ([Tang et al., 2016](#)) applying a graph-centric usage of reads as database.

279 The sequencing of metatranscriptomes is similar to metagenome sequencing [Figure 2A]. In principle  
280 only translation of RNA to DNA is required. Identification of proteins against metatranscriptomes is  
281 beneficial, since organisms only transcribe genes that are currently used ([Wilmes et al., 2015](#)).

282 Sequencing a metagenome or metatranscriptome for each sample is not always possible due to the high  
283 cost and effort **for the sequencing and the data processing**. Thus, researchers use metagenomes from  
284 similar samples or comprehensive databases such as UniProtKB/SwissProt, UniProtKB/TrEMBL ([UniProt,](#)  
285 [2015](#)), UniRef ([Suzek et al., 2007](#)), NCBI ([Coordinators, 2017](#)) or Ensemble ([Yates et al., 2016](#)) [Figure  
286 2\_1]. Database searches against complete comprehensive databases require long computation times and  
287 **decrease the number of identified proteins due to the overestimation of the FDR**. Reasonable constraints  
288 on these comprehensive databases are therefore necessary. For example [Jagtap et al. \(2013\)](#) proposed  
289 to search in two steps. Taxonomies or proteins identified in the first error-tolerant search are used to  
290 restrict the protein database for the second search [Figure 2\_2]. This obviously increases computation  
291 times, but reduces the FDR and the threshold for protein identifications. In the end more proteins are  
292 identified, but how well this approximates the real FDR remains unclear. Another option for reduction of

293 the FDR is to perform several searches against smaller sub databases and to merge their results  
294 afterwards (Muth et al., 2016; Tanca et al., 2016) [Figure 2\_3]. A more reasonable approach to constrain  
295 the protein database is taxonomic foreknowledge, because in some cases taxonomic composition of the  
296 sample is known (Tanca et al., 2016) [Figure 2\_4]. For example, sequencing of the 16S-rRNA gene  
297 provides a taxonomic profile. Nevertheless, performing pre-searches against all taxonomies can help to  
298 avoid excessive constraints on protein taxonomy during the actual searches.

299 A smart idea to decrease computational time for protein database searches was recently proposed by  
300 May et al. (2016). They searched against peptide databases instead of protein databases [Figure 2 E].  
301 This reduces the size of the search space due to the grouping of identical peptides from homologous  
302 proteins.

303 To summarize, all strategies to constrain protein databases carry some pitfalls and we would recommend  
304 researchers to try different approaches. Despite all these strategies for protein database construction,  
305 inaccurate FDR estimation hampers metaproteomics studies. Solutions other than the target-decoy  
306 approach are required to validate protein identifications across different MS and database search  
307 algorithms. A promising step towards this direction represent semi-supervised machine learning  
308 algorithms such as the software tools Percolator (Kall et al., 2007) or Nokoi (Gonnelli et al., 2015). They  
309 distinguish correct and incorrect peptide-to-spectrum matches using a classifier based on learning  
310 algorithms from real data.

### 311 **5. Construction of user databases for protein identification: A use case**

312 In order to visualize the impact of user databases a case study was conducted for a metaproteome  
313 analysis of three different biogas plant samples (BGP01, BGP02, BGP03). After phenol extraction, SDS-  
314 PAGE separation into ten fractions (Heyer et al., 2013) and LC-MS/MS measurement using an Orbitrap  
315 Elite (Heyer et al., 2016) different protein databases were tested [Figure 3]. First the samples were  
316 searched against the UniProtKB/SwissProt database. Second several metagenomes from biogas plants  
317 were tested (metagenome 1, metagenome 2, metagenome 4, metagenome 5 (Stolze et al., 2016),  
318 metagenome 6 (Schlüter et al., 2008). Of these metagenomes number 1 and 2 were prepared for BGP01  
319 resp. BGP02. A metagenome from a waste water treatment plant (WWTP) (Püttker et al., 2015) from was  
320 used as a negative control. Furthermore, the impacts of combining databases as well as of combining the  
321 results were evaluated.

322 The smallest numbers of identified metaproteins could be identified by the protein database search  
323 against the WWTP metagenome followed by the search against the UniProtKB/SwissProt database.  
324 Better results were obtained with the biogas plant metagenomes. Instead of 900 metaproteins for the  
325 protein database search against UniProtKB/SwissProt database about 2.000 metaproteins were  
326 identified using the biogas plant metagenomes. In some cases metagenomes appeared to be  
327 interchangeable, because metagenomes from other biogas plant samples showed equal or even better  
328 numbers of identified metaproteins as matching metagenomes, e.g. BGP02 and metagenome 2. This  
329 result questions whether the generation of a corresponding metagenome for each sample is always  
330 necessary. The combination of different metagenomes additionally increased the number of identified  
331 metaproteins to about 4.000 (combination metagenome 1+2+4+5+6). However, the number of  
332 additional metaprotein identifications decreased for each additional metagenome included in the

333 search. In contrast the combination of metagenome 5 and the poorly matching metagenome from a  
334 waste water treatment plant (WWTP) decreased the number of identified metaproteins showing that an  
335 increased size of the database led to an increased chance of false positive hits and an increased FDR. The  
336 highest number of identified metaproteins was obtained with the separate search against all  
337 metagenomes (metagenome 1;2;4;5;6) and subsequent combination of the results. Focusing on central  
338 metabolism and plotting the metaproteins into KEGG map 1200 clearly shows a higher coverage of  
339 pathways using the combined single searches (Figure 4). This strategy avoided the increase of the FDR  
340 due to the bigger database, but the statistical correctness of this approach is questionable. However, it  
341 circumvents the accumulation of redundant sequence data in a combined database contributing to  
342 increased database size and FDR. Therefore, the removal of redundancy using peptide based databases  
343 could be a strategy to combine databases without increasing the FDR. Furthermore, the fact that  
344 combined metagenomes outcompete single corresponding metagenomes points out that many  
345 metagenome sequences do not comprehensively represent the microbial communities.

#### 346 **6. Protein inference problem and the grouping of proteins into “metaproteins”**

347 Redundant identifications arising from homologous proteins share identical peptides and are therefore  
348 indistinguishable from each other. This hampers result evaluation and sample comparison within  
349 metaproteomic studies.

350 For pure culture proteomics *Niewjetzki et al. (2003)* proposed to use the least number of proteins to  
351 explain all peptides. But this neglects the presence of protein isoforms or proteins from unsequenced  
352 microorganisms ([Hettich et al., 2013](#)) often found in analyses of metaproteomics data. To solve this issue  
353 the metaproteomics community started to develop concepts for grouping of redundant protein  
354 identifications [Table 2]. The metaprotein concept, introduced by *Muth et al. (2015a)*, provides a good  
355 summary on protein grouping. Similar amino acid sequences (protein rules) or shared peptide  
356 identifications (peptide rules) constitute suitable criteria for grouping of homologous protein  
357 identifications into metaproteins. Conveniently, UniRef Clusters ([Lu et al., 2014](#); [Suzek et al., 2007](#)) and  
358 KEGG Ontologies ([Gotelli et al., 2012](#); [Kanehisa et al., 2016](#)) already classify most proteins on their  
359 sequence similarity. An easy retrieval of these classifications is enabled by the UniProtKB database,  
360 which is accessible through the UniProtJAPI library ([Patient et al., 2008](#)). Alternatively, proteins can be  
361 grouped when they share at least one identified peptide ([Kohrs et al., 2014](#); [Lu et al., 2014](#)) or an  
362 identical peptide set ([Keiblinger et al., 2012](#); [Kolmeder et al., 2012](#); [Schneider et al., 2011](#)). It should be  
363 noted that for peptide comparison, the isobaric amino acids leucine and isoleucine are not  
364 distinguishable from each other.

365 All these strategies reduce the redundancy of the protein identifications successfully. However, only  
366 grouping based on identified peptides considers different conservation levels of the protein sequences.  
367 Thus, it enables a better taxonomic classification. Unfortunately, sample comparison using the peptide  
368 rule requires the protein grouping across all samples. Furthermore, the grouping may change as soon as  
369 additional samples are added. In consequence, grouping according to sequence similarity, such as UniRef  
370 clusters, is better suited for sample comparisons ([Heyer et al., 2016](#); [Kohrs et al., 2017](#)).

371 In some instances it is desirable to consider the production of homologous proteins by different species.  
372 Homologous proteins often share peptides, which only differ in one or two amino acids. This indicates  
373 that these proteins should not be grouped together. To consider this bioinformatically, the Levenshtein  
374 distance ([Levenshtein, 1966](#)) between peptides of a protein group can be calculated ([Muth et al., 2015a](#)).  
375 Taxonomic foreknowledge is another option to improve metaprotein grouping. Protein groups can be  
376 restricted to certain phylogenetic affiliations, e.g. only proteins from the same genus.

377

### 378 **7. Taxonomic and functional result evaluation**

379 Comprehensive metaproteomics studies aim to describe the taxonomies and functions of complete  
380 microbial communities. In particular, the functions performed by each taxon should be elucidated.

381 Protein taxonomy [Table 3] is usually defined according to the NCBI Taxonomy ([Federhen, 2012](#)). It  
382 comprises the classification for all taxonomic levels into the phylogenetic tree starting from **species,**  
383 **genus and family via class, order and phylum to the kingdom and superkingdom levels.**

384 In contrast to pure culture proteomics, a large portion of identified peptides in metaproteomics may  
385 belong to several proteins from different species. Thus, the taxonomic value of an identified peptide is  
386 estimated using the lowest common ancestor (LCA) of the protein taxonomies where this peptide occurs.  
387 Protein taxonomy is then defined as the LCA of the peptide identifications ([Huson et al., 2011](#); [Jagtap et](#)  
388 [al., 2012](#)) or on the basis of unique peptides ([Karlsson et al., 2012](#); [Rooijers et al., 2011](#)). Certain taxa  
389 have a much larger number of unique peptides, which biases the taxonomic profile towards these taxa.  
390 In general, unique peptides are fairly uncommon, as the analyses by UniPept demonstrate ([Mesuere et](#)  
391 [al., 2015](#)). The LCA approach is imprecise as well, because peptide taxonomy is often assigned on the  
392 order level and not on the species level. To refine the taxonomy profile [Huson et al. \(2016\)](#) propose to  
393 weigh the identified peptides and their LCA taxonomy by the amount of unique peptides. Another  
394 approach to improve the precision of the taxonomic profile is to weigh identified peptides by their  
395 spectral count and their occurrence in reference proteomes ([Penzlin et al., 2014](#)). Still, evaluation and

396 comparison of taxonomic profiles is often challenging due to the high complexity of the data. This has led  
397 to several new approaches for data evaluation and visualization. The Krona plot ([Ondov et al., 2011](#))  
398 clearly visualizes the taxonomy profile of a sample over all taxonomic levels. Furthermore, calculating  
399 community indices such as richness and evenness can give a general overview about the taxonomic  
400 profile of different samples ([Heyer et al., 2016](#); [Marzorati et al., 2008](#)). In addition, specific interactions  
401 between single taxa can be examined by co-occurrence networks ([Heyer et al., 2016](#); [Huson et al., 2016](#);  
402 [Jenssen et al., 2001](#)).

403 Several approaches with varying degree of specificity exist to assign functions to proteins [Table 3]. The  
404 protein acetyl-coenzyme A synthetase (P27550) is selected as example. It belongs to the acetate  
405 catabolism, which is sufficient to classify this proteins function. In other cases however, it is necessary to  
406 know that this protein transfers a coenzyme or contributes to chemotaxis. Originally, researchers studied  
407 the function of proteins separately through biochemical assays. Later their results were compiled,  
408 standardized and stored in databases. Recently, the functions of proteins from new species are derived  
409 from sequence similarity to functionally classified proteins. Functional classification of proteins with

410 similar sequences is provided by databases such as KEGG ontology (KO) ([Kanehisa et al., 2016](#)), cluster of  
411 orthologous groups (COG) ([Tatusov et al., 2000](#)) and evolutionary genealogy of genes: non-supervised  
412 orthologous (eggNOG) ([Huerta-Cepas et al., 2016](#)).

413 Proteins of the same function possess differences in their amino acid sequence, but the sequences of  
414 their functional domains are highly conserved. Accordingly, the PFAM ([Finn et al., 2016](#)), the TIGRFAM  
415 database ([Haft et al., 2013](#)), the SMART database ([Letunic et al., 2015](#)) and the InterPro database ([Finn et  
416 al., 2017](#)) provide a functional classification based on similar functional domains. For example, acetyl-  
417 coenzyme A synthetase (P27550) possesses an acetyl-coenzyme A synthetase domain and an AMP-  
418 binding enzyme domain.

419 It is important to note that functional annotation of proteins can be divided into categories such as  
420 molecular function, biological process or ligand, which are organized hierarchically. This is achieved by  
421 gene ontologies (GO) ([Ashburner et al., 2000](#)) and UniProtKB keywords ([UniProt, 2015](#)). For acetyl-  
422 coenzyme A synthetase (P27550) the UniProtKB keyword of the category ligand is ATP-binding protein,  
423 which belongs to the group of nucleotide-binding proteins. Enzyme commission numbers (EC) are  
424 another functional characterization of proteins ([Bairoch, 2000](#)). They use a four digit number code to  
425 classify enzymes depending on the catalyzed biochemical reaction. The EC for acetyl-coenzyme A  
426 synthetase (P27550) is 6.2.1.1, where 6 classifies it as a ligase, 6.2 as forming carbon sulfur bonds, 6.2.1.  
427 as acid-thiol ligase and 6.2.1.1. as acetate Co A ligase.

428 Conveniently, access to this taxonomic and functional metainformation is already provided by well  
429 annotated databases, such as UniProtKB. The entire database is available via the UniProt webpage and  
430 can be accessed programmatically via connectors such as the UniProtJAPI ([Patient et al., 2008](#)).  
431 Metagenomes miss taxonomic and functional annotation. Therefore, metagenome sequences are  
432 annotated by BLAST ([Altschul et al., 1990](#)) to link them to sequences of annotated proteins. Contigs may  
433 contain several genes with different functions, which can lead to false annotations. Moreover, the best  
434 BLAST hit is not always the correct one ([Timmins-Schiffman et al., 2017](#)) and for searches with short  
435 sequences, such as peptides, parameters for the BLAST should be adapted (MS-BLAST ([Shevchenko et al.,  
436 2001](#))). Moreover, BLAST requires extensive computational time, which was addressed by development  
437 of the time-saving DIAMOND tool ([Buchfink et al., 2015](#)).

438 Another aim of metaproteomics studies is the analysis of certain metabolic pathways. Therefore,  
439 identified proteins can be visualized in the different metabolic and interaction pathways, using the  
440 pathway repositories MetaCyc ([Caspi et al., 2016](#)), KEGG pathways ([Kanehisa et al., 2016](#)) and Reactome  
441 ([Fabregat et al., 2016](#)). For KEGG pathways the web-based Interactive Pathways Explorer (iPath) ([Yamada  
442 et al., 2011](#)) provides an improved visualization and supports pathway analysis. Mapping of proteins to  
443 pathways is provided via the EC and KO numbers. Unfortunately, metabolic networks are incomplete,  
444 since many pathways are still unknown or specific for a minority of species. To overcome this limitation  
445 researcher started to create their own metabolic pathway maps. To achieve this, biochemical reactions,  
446 represented by EC numbers of identified proteins, were connected ([Tobalina et al., 2015](#)). A similar  
447 approach was chosen by *Roume et al. (2015)* aiming to identify key functions within a microbial  
448 community. Metabolic networks were modelled as a graph, where proteins (KO number) represented  
449 nodes and metabolites represented edges. Finally they defined key functions as nodes with high

450 neighborhood connectivity. In future, networks based on metaproteome data could be used to predict  
451 metabolic fluxes, using software tools such as the CellNetAnalyzer ([Klamt et al., 2007](#)).

#### 452 **8. Quantitative data analysis in metaproteome studies**

453 Protein quantification is crucial for comparative metaproteomics studies. Indeed different approaches  
454 for quantitative proteomics exist, e.g. isotopic chemical labelling of peptides ([Vaudel et al., 2010](#)). But  
455 due to interference of these approaches with contaminating compounds many metaproteomics studies  
456 simply rely on the estimation of protein amount by counting identified peptides or spectra and  
457 normalizing these results ([Ishihama et al., 2005](#)), ([Zybailov et al., 2007](#)). Depending on data-dependent  
458 selection of precursor ions and successful peptide identification these approaches are inaccurate and  
459 possess a small dynamic range [Tabb2009]. The quantification of the peptide peak intensity or area  
460 ([Griffin et al., 2010](#)) using tools such as Progenesis QI ([http://www.nonlinear.com/progenesis/qi-for-](http://www.nonlinear.com/progenesis/qi-for-proteomics/)  
461 [proteomics/](http://www.nonlinear.com/progenesis/qi-for-proteomics/)) or MaxQuant (Tyanova et al., 2016) is preferable. Alternatively, data-independent  
462 acquisition of MS/MS data (SWATH, MS<sup>E</sup>) combines peptide identification and quantification capturing  
463 all possible fragment information of all precursors for subsequent protein quantification from complex  
464 data (Bilbao et al., 2015). The most accurate quantification can be achieved by targeting only a single  
465 peptide (“single reaction monitoring”) or a limited selection of peptides of a certain protein (“single  
466 reaction monitoring”). For example, [Saito et al. \(2015\)](#) used this approach to quantify two nitrogen  
467 regulatory proteins for cyanobacterial taxa within microbial samples from the Central Pacific Ocean. The  
468 addition of isotopically labeled peptide for absolute quantification and the application of the Skyline  
469 software (MacLean et al., 2010) further improve this approach.

470 However, selection of peptides for targeted metaproteomics is more challenging than in pure culture  
471 proteomics, because a peptide may belong to multiple proteins from different taxa. Thus, the Unique  
472 Peptide Finder of the UniPept webservice ([Mesuere et al., 2016](#)) was developed to facilitate the selection  
473 of unique peptides for a certain taxa.

474

#### 475 **8. Strategies for storing and deployment of huge data**

476 Metaproteomics experiments comprise a massive amount of data including MS spectra, identified  
477 peptides and proteins as well as taxonomic and functional information. Our latest large-scale  
478 metaproteomics study produced about two Terabyte of data comprising roughly 15 million spectra and  
479 23,000 identified metaproteins (data not shown). Consequently, appropriate data storage using a  
480 **database management system (DBMS)** is beneficial. Key challenges for DBMS are high speed for writing  
481 and reading data as well as efficient data storage. Since MS acquisition and search algorithms are  
482 relatively slow, writing speed **has a negligible impact**. In contrast, reading speed **can be limiting**, because  
483 researches want to evaluate all data at once. Furthermore, lists of thousands of proteins are unfeasible  
484 when inspecting results. Instead, researchers favor meaningful summaries, comparisons and intuitive  
485 visualizations. But this requires demanding database queries.

486 Relational database management systems, which use the “Structured Query Language” (SQL), have been  
487 the norm to manage data in the past. In recent years, alternatives to SQL have gained popularity and are  
488 aggregated under the term NoSQL (“Not only SQL”). Relational database management systems store  
489 data in separate tables, which are connected via unique relations. NoSQL database management systems

490 use other concepts to store data like key-value associations (Berkeley DB  
491 ([http://www.oracle.com/technetwork/database/database-](http://www.oracle.com/technetwork/database/database-technologies/berkeleydb/overview/index.html)  
492 [technologies/berkeleydb/overview/index.html](http://www.oracle.com/technetwork/database/database-technologies/berkeleydb/overview/index.html), retrieved: 09-02-2017)), columns (Apache Cassandra  
493 (<http://cassandra.apache.org/>, retrieved: 09-02-2017)), documents (MongoDB  
494 (<https://www.mongodb.com>, retrieved: 09-02-2017)) or graphs (Neo4j, ([www.neo4j.com](http://www.neo4j.com), retrieved: 09-  
495 02-2017)).

496 NoSQL databases were motivated by the disadvantage in SQL databases to store all data in one place.  
497 In an analogy SQL databases can be imagined as a large building, which only a limited number of persons  
498 at a time can enter. An SQL query would be a person searching the building and collecting the  
499 information requested. If too many people search the building at a time, they will hinder each other and  
500 slow down the query process. NoSQL databases aim to address this issue of scalability. For instance, in  
501 our analogy Apache Cassandra creates a new identical building as soon as too many people try to enter.  
502 In consequence, NoSQL databases can handle more and more complex data requests. The disadvantage  
503 of NoSQL databases is reduced data consistency and large hard disc requirements due to multiple  
504 instances of the databases.

505 In sum NoSQL databases are highly beneficial for metaproteomics data. In line *Chatterlee et al. (2016)*  
506 already used MongoDB for storing sequence information and *Muth et al. (2015a)* Neo4j for flexible result  
507 queries. Additionally, *Measure et al. (2015)* are planning to use Berkeley databases to store the  
508 taxonomic value of each tryptic peptide.

509 Another trend of data storing and deployment which could be useful to increase the speed of data  
510 processing in metaproteomics is fast data (Braun et al., 2015 ). The fast data approach makes it possible  
511 to stream single spectra data to the cloud and process the data in real time for storing the results into  
512 the database. In other words, it parallelizes the data processing step and the measurement step to  
513 reduce experiment time. For example already the software MaxQuant Real-Time (Graumann et al., 2012)  
514 picks up this idea and processes the MS data in real time.

515

## 516 **9. Future challenges, perspectives and demands**

517 Predictions about the future of metaproteomics software need to anticipate future applications for  
518 metaproteomics. Foreseeable trends are an increase in MS resolution and therefore more data that will  
519 be acquired. Since metaproteomics is still an emerging field, an increase in the number of research  
520 studies about complex microbial communities is expected. A great potential for the application of  
521 metaproteomics are process control in technical applications as well routine diagnostics of fecal samples.  
522 So far it is known that microbial communities in the human gut system are linked with autoimmune and  
523 allergic diseases, obesity, inflammatory bowel disease (IBD), and diabetes (Clemente et al., 2012).  
524 Consequently, the number of samples in clinical settings could rise to several thousand per day. Such an  
525 increase in sample numbers requires software tools that can handle huge data amounts. For routine  
526 diagnostics the total computation time may not exceed a few hours, so that a complete metaproteomics  
527 analysis may require less than one day. Another aspect is that software for medical applications has to  
528 conform to high quality standards and specific privacy regulations. Moreover, medical staff without a

529 special bioinformatic background should be able to operate such software tools. Although the routine  
530 usage of metaproteomics is still in question, the development may proceed quickly. For example, MALDI-  
531 MS based identification of microbial isolates became a standard procedure in clinical laboratories.

532 Strategies to facilitate software usage are to provide it via Docker (e.g. Bioconda  
533 <https://bioconda.github.io/>, retrieved: 09-02-2017) or web services to avoid problems with the  
534 installation and configuration of complex software frameworks. For example, developers of the MPA are  
535 planning to provide their software platform as web service within the de.NBI project. Most users with a  
536 medical or biological background would favor a graphical ready-to-use software tool. In contrast,  
537 bioinformaticians prefer modular software packages operated from the command line. The latter  
538 strategy enables flexible assembly of workflows and an easy improvement of single modules. The  
539 challenge for future development of metaproteomics software is to satisfy both sides.

540 Because metaproteomics is still a developing field, universal standards still have to be adopted by the  
541 community. Implementation of ring trials for metaproteomics data processing could further insights into  
542 the comparability of software tools, and enable the introduction of quality standards.

543 Further improvement requires the validation of protein identifications by the FDR estimation. In contrast  
544 to pure culture proteomics the estimated FDR is not always correct since the protein sequences for the  
545 investigated samples are often unknown. A solution might be the usage of semi-supervised machine  
546 learning algorithms such as the software tools Percolator or Nokoi (Gonnelli et al., 2015).

547 The use of protein databases could be standardized as well. While some researchers use comprehensive  
548 protein databases, others use diverse metagenomes, which differ in the processing state and origin. A  
549 solution might be the generation of non-redundant (May et al., 2016), fusion metagenomes for each  
550 type of microbial community. Thereby, this fusion metagenome should be assembled as far as possible.

551 Additionally, the binning of metagenomes may also improve the protein database quality. Proteins of the  
552 same function or metabolic pathway are often located adjacent on a contig or operon. Thus, they should  
553 feature equal expression patterns.

554 The key to handle the increased amount of data is the real-time processing of all arising MS data as well  
555 as the scalability of the software and the database. This means that the single computational steps  
556 operate in parallel and hardware resources can be allocated on demand, e.g. by cloud computing (Mell  
557 and Grance, 2010). To guarantee the long term maintenance and support for such systems, it is  
558 reasonable to follow the latest trends from the industry instead of developing own solutions. Suitable  
559 frameworks, among others, are Apache Spark (<http://spark.apache.org/>, retrieved: 09-02-2017) for  
560 analyzing data distributed in the cloud and OpenStack (<https://www.openstack.org/>, retrieved: 09-02-  
561 2017) to manage the instances running on the cloud.

562 Another strategy to decrease computation time is the smart deployment of hardware resources.  
563 Graphical processing units (GPU) can perform specific calculations in parallel. On the other hand central  
564 processing units (CPU) are suited for general tasks, but work serially. Identification of MS/MS spectra is a  
565 calculation that can be parallelized. In line, the protein database search algorithm X!Tandem was  
566 recently adopted to utilize a GPU (He and Li, 2015).



567 Beside adaptation of metaproteomics to bigger data volumes and the decrease of computation time,  
568 improved bioinformatic strategies are required to increase the number of identified spectra. State-of-  
569 the-art metaproteomics studies only achieve identification of 5-30% spectra. An estimated 30% of all  
570 spectra belong to solvent and background components (Griss et al., 2016). This means at least another  
571 30% spectra remain unidentified. Better metaproteomics software should contribute to overcome this  
572 issue. The generation of more suitable metagenomes for protein identification may increase the amount  
573 of identified spectra significantly. Inversely, assembly of metagenomes can be validated using peptides  
574 identified in metaproteomics studies (Nesvizhskii, 2014). ~~There are also alternatives to the generation of~~  
575 ~~metagenomes.~~

576 ~~Due to increased computational power and more precise MS it may become possible to search against a~~  
577 ~~database containing all theoretical peptides for a specific mass (Sadygov, 2015).~~ Spectral libraries  
578 represent another strategy to handle unidentified spectra (Lam et al., 2007). They could store and cluster  
579 spectra from any sample. Samples can be also compared based on their unidentified spectra. Interesting  
580 spectra can be annotated later using protein database search algorithms. Due to the drastic reduction of  
581 candidates, manual *de novo* sequencing is also possible (Frank and Pevzner, 2005). Function and  
582 taxonomy of *de novo* peptides can be derived by MS-BLAST search (Shevchenko et al., 2001). However,  
583 *de novo* sequencing of peptides is hampered by the short length of tryptic peptides which impede MS-  
584 BLAST identification. Better *de novo* and MS-BLAST results could be achieved by other proteases such as  
585 Lys-C (Jekel et al., 1983) or Arg-C, which result in longer peptides. ~~Due to increased computational power~~  
586 ~~and more precise MS it may become possible to search against a database containing all theoretical~~  
587 ~~peptides for a specific mass (Sadygov, 2015).~~ This would also solve problem with the database size  
588 ~~dependency of the FDR estimation.~~

589 Finally, metaproteomics software can benefit from the incorporation of data from other multi-omics  
590 techniques (Brink et al., 2016; Heintz-Buschart et al., 2016), e.g. metabolome data. For a detailed  
591 overview on multi-omics data processing, please refer to Franzosa et al. (2015) (Franzosa et al., 2015).

592 Due to increased computational power and more precise MS it may become possible to search against a  
593 database containing all theoretical peptides for a specific mass (Sadygov, 2015). Spectral libraries  
594 represent another strategy to handle unidentified spectra (Lam et al., 2007). They could store and cluster  
595 spectra from any sample. Samples can be also compared based on their unidentified spectra. Interesting  
596 spectra can be annotated later using protein database search algorithms. Due to the drastic reduction of  
597 candidates, manual *de novo* sequencing is also possible (Frank and Pevzner, 2005). Function and  
598 taxonomy of *de novo* peptides can be derived by MS-BLAST search (Shevchenko et al., 2001). However,  
599 *de novo* sequencing of peptides is hampered by the short length of tryptic peptides which impede MS-  
600 BLAST identification. Better *de novo* and MS-BLAST results could be achieved by other proteases such as  
601 Lys-C (Jekel et al., 1983) or Arg-C, which result in longer peptides.

602 Finally, metaproteomics software can benefit from the incorporation of data from other multi-omics  
603 techniques (Brink et al., 2016; Heintz-Buschart et al., 2016), e.g. metabolome data. For a detailed  
604 overview on multi-omics data processing, please refer to Franzosa et al. (2015) (Franzosa et al., 2015).

605



## 607 **10. Conclusions**

608 Metaproteomics represents a powerful tool for the taxonomic and functional characterization  
609 of complex microbial communities from environmental samples. In the future it has the  
610 potential to become a valuable tool for routine diagnostics, e.g. analysis of human feces.  
611 However, success of metaproteomics studies depends on dedicated software tools. These tools  
612 must be capable to handle big data, but also need to be useable by people with no background  
613 in bioinformatics. To achieve these goals, web services and software tools capable of parallel  
614 computing are reasonable (e.g. cloud computing). This would decrease computational costs and  
615 enables small laboratories to perform metaproteomics studies. Moreover, metaproteomics  
616 studies will benefit from software supporting the taxonomic and functional interpretation of  
617 results. Even if it is obvious, the close cooperation of bioinformaticians and biologists should  
618 also be considered during software development.

## 619 **Abbreviations**

620	CPU:	central processing unit
621	COG:	clusters of orthologous groups
622	DBMS:	database management system (DBMS)
623	de.NBI:	German Network for Bioinformatics Infrastructure
624	EC:	enzyme commission number
625	eggNOG:	evolutionary genealogy of genes: non-supervised orthologous
626	FDR:	false discovery rate
627	GPU:	graphical processing unit
628	GO:	gene ontologies
629	iPath:	Interactive Pathways Explorer
630	LC:	liquid chromatography
631	LCA:	lowest common ancestor
632	KO:	KEGG ontologies
633	MPA:	MetaProteomeAnalyzer
634	MS:	mass spectrometer
635	MS/MS:	tandem mass spectrometer

- 636 m/z-ratio: mass-to-charge ratio
- 637 NoSQL: not only SQL
- 638 SQL: structured query language
- 639

## 640 **5. Figures & tables**

641 **Figure 1:** Workflow for metaproteome analyses. Software tools specific for metaproteomics are  
642 highlighted in bold. Beside tools for single steps of the bioinformatic analysis also comprehensive  
643 software platforms are available (K).

644 **Figure 2:** Database construction for protein identification.

645 **Figure 3:** Impact of different metagenomes and their combination on the number of identified  
646 metaproteins.

647 **Figure 4:** This figure shows the identified metaproteins of sample BGP01 after protein database  
648 search against different databases mapped against the KEGG map 1200 (central carbon  
649 metabolism. Green: metaproteins identified by protein database search against  
650 UniProtKB/SwissProt; blue: metaproteins identified additionally by protein database search  
651 against the combined metagenomes (1+2+4+5+6); red: metaproteins identified additionally by  
652 protein database search against the single metagenomes (1;2;4;5;6).

653

654 **Table 1:** Overview about metaproteomic specific issues and appropriated software resp.  
655 bioinformatic strategies

656 **Table 2:** Strategies for grouping of redundant homologous proteins to metaproteins

657 **Table 3:** Strategies for taxonomic and functional annotation of proteins.

658

---

## 659 **6. Additional files**

660 Not applicable.

661

## 662 **Declarations**

663

## 664 **Acknowledgement**

665 Not applicable.

666

## 667 **Authors' contributions**

668 The manuscript was written by Robert Heyer (RH), Dirk Benndorf (DB), Kay Schallert (KS),  
669 Beatrice Becher (BB), Udo Reichl (UR) and Günther Saake (GS). All authors read and approved  
670 the final manuscript.

671

672 **Availability of data and material**

673 Not applicable.

674

675 **Ethics approval and consent to participate**

676 Not applicable.

677

678 **Consent for publication**

679 Not applicable.

680 **Competing interests**

681 The authors declare that they have no competing interest.

682

683 **Funding**

684 Funding: This work was supported by the Federal Ministry of Food, Agriculture and Consumer

685 Protection (BMELV) communicated by the Agency for Renewable Resources (FNR), grant no.

686 22404115 (“Biogas Measurement Program III”) and the de.NBI network (“MetaProtServ de-NBI-

687 039”).

688

## 689 8. Literature

- 690 Abram, F., Enright, A.M., O'Reilly, J., Botting, C.H., Collins, G., O'Flaherty, V., (2011) A  
691 metaproteomic approach gives functional insights into anaerobic digestion. *J Appl Microbiol* 110,  
692 1550-1560.
- 693 Afgan, E., Baker, D., van den Beek, M., Blankenberg, D., Bouvier, D., Cech, M., Chilton, J.,  
694 Clements, D., Coraor, N., Eberhard, C., Gruning, B., Guerler, A., Hillman-Jackson, J., Von Kuster,  
695 G., Rasche, E., Soranzo, N., Turaga, N., Taylor, J., Nekrutenko, A., Goecks, J., (2016) The Galaxy  
696 platform for accessible, reproducible and collaborative biomedical analyses: 2016 update.  
697 *Nucleic Acids Research* 44, W3-W10.
- 698 Altschul, S.F., Gish, W., Miller, W., Myers, E.W., Lipman, D.J., (1990) Basic local alignment search  
699 tool. *Journal of Molecular Biology* 215, 403-410.
- 700 Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K.,  
701 Dwight, S.S., Eppig, J.T., Harris, M.A., Hill, D.P., Issel-Tarver, L., Kasarskis, A., Lewis, S., Matese,  
702 J.C., Richardson, J.E., Ringwald, M., Rubin, G.M., Sherlock, G., (2000) Gene ontology: tool for the  
703 unification of biology. The Gene Ontology Consortium. *Nature Genetics* 25, 25-29.
- 704 Bairoch, A., (2000) The ENZYME database in 2000. *Nucleic Acids Research* 28, 304-305.
- 705 Barnouin, K., (2011) Guidelines for experimental design and data analysis of proteomic mass  
706 spectrometry-based experiments. *Amino Acids* 40, 259-260.
- 707 Barsnes, H., Vaudel, M., Colaert, N., Helsens, K., Sickmann, A., Berven, F.S., Martens, L., (2011)  
708 compomics-utilities: an open-source Java library for computational proteomics. *BMC*  
709 *Bioinformatics* 12, 70.
- 710 Bentley, D.R., Balasubramanian, S., Swerdlow, H.P., Smith, G.P., Milton, J., Brown, C.G., Hall, K.P.,  
711 Evers, D.J., Barnes, C.L., Bignell, H.R., Boutell, J.M., Bryant, J., Carter, R.J., Cheetham, R.K., Cox,  
712 A.J., Ellis, D.J., Flatbush, M.R., Gormley, N.A., Humphray, S.J., Irving, L.J., Karbelashvili, M.S., Kirk,  
713 S.M., Li, H., Liu, X.H., Maisinger, K.S., Murray, L.J., Obradovic, B., Ost, T., Parkinson, M.L., Pratt,  
714 M.R., Rasolonjatovo, I.M.J., Reed, M.T., Rigatti, R., Rodighiero, C., Ross, M.T., Sabot, A., Sankar,  
715 S.V., Scally, A., Schroth, G.P., Smith, M.E., Smith, V.P., Spiridou, A., Torrance, P.E., Tzonev, S.S.,  
716 Vermaas, E.H., Walter, K., Wu, X.L., Zhang, L., Alam, M.D., Anastasi, C., Aniebo, I.C., Bailey,  
717 D.M.D., Bancarz, I.R., Banerjee, S., Barbour, S.G., Baybayan, P.A., Benoit, V.A., Benson, K.F., Bevis,  
718 C., Black, P.J., Boodhun, A., Brennan, J.S., Bridgham, J.A., Brown, R.C., Brown, A.A., Buermann,  
719 D.H., Bundu, A.A., Burrows, J.C., Carter, N.P., Castillo, N., Catenazzi, M.C.E., Chang, S., Cooley,  
720 R.N., Crake, N.R., Dada, O.O., Diakoumakos, K.D., Dominguez-Fernandez, B., Earnshaw, D.J.,  
721 Egbujor, U.C., Elmore, D.W., Etchin, S.S., Ewan, M.R., Fedurco, M., Fraser, L.J., Fajardo, K.V.F.,  
722 Furey, W.S., George, D., Gietzen, K.J., Goddard, C.P., Golda, G.S., Granieri, P.A., Green, D.E.,  
723 Gustafson, D.L., Hansen, N.F., Harnish, K., Haudenschield, C.D., Heyer, N.I., Hims, M.M., Ho, J.T.,  
724 Horgan, A.M., Hoschler, K., Hurwitz, S., Ivanov, D.V., Johnson, M.Q., James, T., Jones, T.A.H.,  
725 Kang, G.D., Kerelska, T.H., Kersey, A.D., Khrebtukova, I., Kindwall, A.P., Kingsbury, Z., Kokko-  
726 Gonzales, P.I., Kumar, A., Laurent, M.A., Lawley, C.T., Lee, S.E., Lee, X., Liao, A.K., Loch, J.A., Lok,  
727 M., Luo, S.J., Mammen, R.M., Martin, J.W., McCauley, P.G., McNitt, P., Mehta, P., Moon, K.W.,  
728 Mullens, J.W., Newington, T., Ning, Z.M., Ng, B.L., Novo, S.M., O'Neill, M.J., Osborne, M.A.,

729 Osnowski, A., Ostadan, O., Paraschos, L.L., Pickering, L., Pike, A.C., Pike, A.C., Pinkard, D.C.,  
730 Pliskin, D.P., Podhasky, J., Quijano, V.J., Raczky, C., Rae, V.H., Rawlings, S.R., Rodriguez, A.C., Roe,  
731 P.M., Rogers, J., Bacigalupo, M.C.R., Romanov, N., Romieu, A., Roth, R.K., Rourke, N.J., Ruediger,  
732 S.T., Rusman, E., Sanches-Kuiper, R.M., Schenker, M.R., Seoane, J.M., Shaw, R.J., Shiver, M.K.,  
733 Short, S.W., Sizto, N.L., Sluis, J.P., Smith, M.A., Sohna, J.E.S., Spence, E.J., Stevens, K., Sutton, N.,  
734 Szajkowski, L., Tregidgo, C.L., Turcatti, G., vandeVondele, S., Verhovsky, Y., Virk, S.M., Wakelin, S.,  
735 Walcott, G.C., Wang, J.W., Worsley, G.J., Yan, J.Y., Yau, L., Zuerlein, M., Rogers, J., Mullikin, J.C.,  
736 Hurles, M.E., McCooke, N.J., West, J.S., Oaks, F.L., Lundberg, P.L., Klenerman, D., Durbin, R.,  
737 Smith, A.J., (2008) Accurate whole human genome sequencing using reversible terminator  
738 chemistry. *Nature* 456, 53-59.

739 Bilbao, A., Varesio, E., Luban, J., Strambio-De-Castillia, C., Hopfgartner, G., Muller, M., Lisacek, F.,  
740 (2015) Processing strategies and software solutions for data-independent acquisition in mass  
741 spectrometry. *Proteomics* 15, 964-980.

742 Braun, L., Etter, T., Gasparis, G., Kaufmann, M., Kossmann, D., Widmer, D., (2015 ) Analytics in  
743 Motion: High Performance Event-Processing AND Real-Time Analytics in the Same Database.  
744 Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data pp.  
745 251-264

746 Brink, B.G., Seidel, A., Kleinbolting, N., Nattkemper, T.W., Albaum, S.P., (2016) Omics Fusion - A  
747 Platform for Integrative Analysis of Omics Data. *Journal of Integrative Bioinformatics* 13, 296.

748 Buchfink, B., Xie, C., Huson, D.H., (2015) Fast and sensitive protein alignment using DIAMOND.  
749 *Nature Methods* 12, 59-60.

750 Campanaro, S., Treu, L., Kougias, P.G., De Francisci, D., Valle, G., Angelidaki, I., (2016)  
751 Metagenomic analysis and functional characterization of the biogas microbiome using high  
752 throughput shotgun sequencing and a novel binning strategy. *Biotechnology for Biofuels* 9.

753 Cappadona, S., Baker, P.R., Cutillas, P.R., Heck, A.J., van Breukelen, B., (2012) Current challenges  
754 in software solutions for mass spectrometry-based quantitative proteomics. *Amino Acids* 43,  
755 1087-1108.

756 Caspi, R., Billington, R., Ferrer, L., Foerster, H., Fulcher, C.A., Keseler, I.M., Kothari, A.,  
757 Krummenacker, M., Latendresse, M., Mueller, L.A., Ong, Q., Paley, S., Subhraveti, P., Weaver,  
758 D.S., Karp, P.D., (2016) The MetaCyc database of metabolic pathways and enzymes and the  
759 BioCyc collection of pathway/genome databases. *Nucleic Acids Research* 44, D471-480.

760 Chatterjee, S., Stupp, G.S., Park, S.K., Ducom, J.C., Yates, J.R., 3rd, Su, A.I., Wolan, D.W., (2016) A  
761 comprehensive and scalable database search system for metaproteomics. *BMC Genomics* 17,  
762 642.

763 Clemente, J.C., Ursell, L.K., Parfrey, L.W., Knight, R., (2012) The impact of the gut microbiota on  
764 human health: an integrative view. *Cell* 148, 1258-1270.

765 Colaert, N., Degroeve, S., Helsens, K., Martens, L., (2011) Analysis of the Resolution Limitations of  
766 Peptide Identification Algorithms. *Journal of Proteome Research* 10, 5555-5561.

767 Colangelo, C.M., Chung, L., Bruce, C., Cheung, K.H., (2013) Review of software tools for design  
768 and analysis of large scale MRM proteomic datasets. *Methods* 61, 287-298.



769 Coordinators, N.R., (2017) Database Resources of the National Center for Biotechnology  
770 Information. *Nucleic Acids Research* 45, D12-D17.

771 Craig, R., Beavis, R.C., (2004) TANDEM: matching proteins with tandem mass spectra.  
772 *Bioinformatics* 20, 1466-1467.

773 Crosswell, L.C., Thornton, J.M., (2012) ELIXIR: a distributed infrastructure for European biological  
774 data. *Trends Biotechnology* 30, 241-242.

775 Doerr, A., (2015) DIA mass spectrometry. *Nature Methods* 12, 35-35.

776 Dorfer, V., Pichler, P., Stranzl, T., Stadlmann, J., Taus, T., Winkler, S., Mechtler, K., (2014) MS  
777 Amanda, a universal identification algorithm optimized for high accuracy tandem mass spectra.  
778 *Journal of Proteome Research* 13, 3679-3684.

779 Elias, J.E., Gygi, S.P., (2007) Target-decoy search strategy for increased confidence in large-scale  
780 protein identifications by mass spectrometry. *Nature Methods* 4, 207-214.

781 Erickson, A.R., Cantarel, B.L., Lamendella, R., Darzi, Y., Mongodin, E.F., Pan, C., Shah, M.,  
782 Halfvarson, J., Tysk, C., Henrissat, B., Raes, J., Verberkmoes, N.C., Fraser, C.M., Hettich, R.L.,  
783 Jansson, J.K., (2012) Integrated metagenomics/metaproteomics reveals human host-microbiota  
784 signatures of Crohn's disease. *PLoS One* 7, e49138.

785 Fabregat, A., Sidiropoulos, K., Garapati, P., Gillespie, M., Hausmann, K., Haw, R., Jassal, B., Jupe,  
786 S., Korninger, F., McKay, S., Matthews, L., May, B., Milacic, M., Rothfels, K., Shamovsky, V.,  
787 Webber, M., Weiser, J., Williams, M., Wu, G., Stein, L., Hermjakob, H., D'Eustachio, P., (2016) The  
788 Reactome pathway Knowledgebase. *Nucleic Acids Research* 44, D481-487.

789 Federhen, S., (2012) The NCBI Taxonomy database. *Nucleic Acids Research* 40, D136-143.

790 Finn, R.D., Attwood, T.K., Babbitt, P.C., Bateman, A., Bork, P., Bridge, A.J., Chang, H.Y., Dosztanyi,  
791 Z., El-Gebali, S., Fraser, M., Gough, J., Haft, D., Holliday, G.L., Huang, H., Huang, X., Letunic, I.,  
792 Lopez, R., Lu, S., Marchler-Bauer, A., Mi, H., Mistry, J., Natale, D.A., Necci, M., Nuka, G., Orengo,  
793 C.A., Park, Y., Pesseat, S., Piovesan, D., Potter, S.C., Rawlings, N.D., Redaschi, N., Richardson, L.,  
794 Rivoire, C., Sangrador-Vegas, A., Sigrist, C., Sillitoe, I., Smithers, B., Squizzato, S., Sutton, G.,  
795 Thanki, N., Thomas, P.D., Tosatto, S.C., Wu, C.H., Xenarios, I., Yeh, L.S., Young, S.Y., Mitchell, A.L.,  
796 (2017) InterPro in 2017-beyond protein family and domain annotations. *Nucleic Acids Research*  
797 45, D190-D199.

798 Finn, R.D., Coggill, P., Eberhardt, R.Y., Eddy, S.R., Mistry, J., Mitchell, A.L., Potter, S.C., Punta, M.,  
799 Qureshi, M., Sangrador-Vegas, A., Salazar, G.A., Tate, J., Bateman, A., (2016) The Pfam protein  
800 families database: towards a more sustainable future. *Nucleic Acids Research* 44, D279-285.

801 Frank, A., Pevzner, P., (2005) PepNovo: de novo peptide sequencing via probabilistic network  
802 modeling. *Anal Chem* 77, 964-973.

803 Franzosa, E.A., Hsu, T., Sirota-Madi, A., Shafquat, A., Abu-Ali, G., Morgan, X.C., Huttenhower, C.,  
804 (2015) Sequencing and beyond: integrating molecular 'omics' for microbial community profiling.  
805 *Nature Reviews Microbiology* 13, 360-372.

806 Gasteiger, E., Gattiker, A., Hoogland, C., Ivanyi, I., Appel, R.D., Bairoch, A., (2003) ExPASy: The  
807 proteomics server for in-depth protein knowledge and analysis. *Nucleic Acids Research* 31, 3784-  
808 3788.

809 Gonnelli, G., Stock, M., Verwaeren, J., Maddelain, D., De Baets, B., Martens, L., Degroeve, S.,  
810 (2015) A decoy-free approach to the identification of peptides. *Journal of Proteome Research* 14,  
811 1792-1798.

812 Gonzalez-Galarza, F.F., Lawless, C., Hubbard, S.J., Fan, J., Bessant, C., Hermjakob, H., Jones, A.R.,  
813 (2012) A Critical Appraisal of Techniques, Software Packages, and Standards for Quantitative  
814 Proteomic Analysis. *Omics* 16, 431-442.

815 Gotelli, N.J., Ellison, A.M., Ballif, B.A., (2012) Environmental proteomics, biodiversity statistics  
816 and food-web structure. *Trends Ecology Evolution* 27, 436-442.

817 Graumann, J., Scheltema, R.A., Zhang, Y., Cox, J., Mann, M., (2012) A Framework for Intelligent  
818 Data Acquisition and Real-Time Database Searching for Shotgun Proteomics. *Molecular &  
819 Cellular Proteomics* 11.

820 Griffin, N.M., Yu, J., Long, F., Oh, P., Shore, S., Li, Y., Koziol, J.A., Schnitzer, J.E., (2010) Label-free,  
821 normalized quantification of complex mass spectrometry data for proteomic analysis. *Nature  
822 Biotechnology* 28, 83-89.

823 Griss, J., Jones, A.R., Sachsenberg, T., Walzer, M., Gatto, L., Hartler, J., Thallinger, G.G., Salek,  
824 R.M., Steinbeck, C., Neuhauser, N., Cox, J., Neumann, S., Fan, J., Reisinger, F., Xu, Q.W., Del Toro,  
825 N., Perez-Riverol, Y., Ghali, F., Bandeira, N., Xenarios, I., Kohlbacher, O., Vizcaino, J.A.,  
826 Hermjakob, H., (2014) The mzTab data exchange format: communicating mass-spectrometry-  
827 based proteomics and metabolomics experimental results to a wider audience. *Molecular &  
828 Cellular Proteomics* 13, 2765-2775.

829 Griss, J., Perez-Riverol, Y., Lewis, S., Tabb, D.L., Dianes, J.A., Del-Toro, N., Rurik, M., Walzer, M.W.,  
830 Kohlbacher, O., Hermjakob, H., Wang, R., Vizcaino, J.A., (2016) Recognizing millions of  
831 consistently unidentified spectra across hundreds of shotgun proteomics datasets. *Nature  
832 Methods* 13, 651-656.

833 Haft, D.H., Selengut, J.D., Richter, R.A., Harkins, D., Basu, M.K., Beck, E., (2013) TIGRFAMs and  
834 Genome Properties in 2013. *Nucleic Acids Research* 41, D387-395.

835 Hanreich, A., Heyer, R., Benndorf, D., Rapp, E., Pioch, M., Reichl, U., Klocke, M., (2012)  
836 Metaproteome analysis to determine the metabolically active part of a thermophilic microbial  
837 community producing biogas from agricultural biomass. *Canadian Journal of Microbiology* 58,  
838 917-922.

839 Hanreich, A., Schimpf, U., Zakrzewski, M., Schluter, A., Benndorf, D., Heyer, R., Rapp, E., Puhler,  
840 A., Reichl, U., Klocke, M., (2013) Metagenome and metaproteome analyses of microbial  
841 communities in mesophilic biogas-producing anaerobic batch fermentations indicate concerted  
842 plant carbohydrate degradation. *Systematic Applied Microbiology* 36, 330-338.

843 He, P., Li, K., (2015) MIC-Tandem: Parallel X! Tandem Using MIC on Tandem Mass Spectrometry  
844 Based Proteomics Data. *Cluster, Cloud and Grid Computing (CCGrid)*, 2015 15th IEEE/ACM  
845 International Symposium on, pp. 717-720.

846 Heintz-Buschart, A., May, P., Laczny, C.C., Lebrun, L.A., Bellora, C., Krishna, A., Wampach, L.,  
847 Schneider, J.G., Hogan, A., de Beaufort, C., Wilmes, P., (2016) Integrated multi-omics of the  
848 human gut microbiome in a case study of familial type 1 diabetes. *Nature Microbiology* 2, 16180.

849 Henry, V.J., Bandrowski, A.E., Pepin, A.S., Gonzalez, B.J., Desfeux, A., (2014) OMICtools: an  
850 informative directory for multi-omic data analysis. Database (Oxford) 2014.

851 Herbst, F.A., Lunsmann, V., Kjeldal, H., Jehmlich, N., Tholey, A., von Bergen, M., Nielsen, J.L.,  
852 Hettich, R.L., Seifert, J., Nielsen, P.H., (2016) Enhancing metaproteomics--The value of models  
853 and defined environmental microbial systems. *Proteomics* 16, 783-798.

854 Hettich, R.L., Pan, C.L., Chourey, K., Giannone, R.J., (2013) Metaproteomics: Harnessing the  
855 Power of High Performance Mass Spectrometry to Identify the Suite of Proteins That Control  
856 Metabolic Activities in Microbial Communities. *Anal Chem* 85, 4203-4214.

857 Heyer, R., Kohrs, F., Benndorf, D., Rapp, E., Kausmann, R., Heiermann, M., Klocke, M., Reichl, U.,  
858 (2013) Metaproteome analysis of the microbial communities in agricultural biogas plants. *New*  
859 *Biotechnology* 30, 614-622.

860 Heyer, R., Benndorf, D., Kohrs, F., De Vrieze, J., Boon, N., Hoffmann, M., Rapp, E., Schluter, A.,  
861 Sczyrba, A., Reichl, U., (2016) Proteotyping of biogas plant microbiomes separates biogas plants  
862 according to process temperature and reactor type. *Biotechnology for Biofuels* 9, 155.

863 Heyer, R., Kohrs, F., Reichl, U., Benndorf, D., (2015) Metaproteomics of complex microbial  
864 communities in biogas plants. *Microbial Biotechnology* 8, 749-763.

865 Hu, Q., Noll, R.J., Li, H., Makarov, A., Hardman, M., Graham Cooks, R., (2005) The Orbitrap: a new  
866 mass spectrometer. *Journal of Mass Spectrometry* 40, 430-443.

867 Huerta-Cepas, J., Szklarczyk, D., Forslund, K., Cook, H., Heller, D., Walter, M.C., Rattei, T., Mende,  
868 D.R., Sunagawa, S., Kuhn, M., Jensen, L.J., von Mering, C., Bork, P., (2016) eggNOG 4.5: a  
869 hierarchical orthology framework with improved functional annotations for eukaryotic,  
870 prokaryotic and viral sequences. *Nucleic Acids Research* 44, D286-293.

871 Huson, D.H., Beier, S., Flade, I., Gorska, A., El-Hadidi, M., Mitra, S., Ruscheweyh, H.J., Tappu, R.,  
872 (2016) MEGAN Community Edition - Interactive Exploration and Analysis of Large-Scale  
873 Microbiome Sequencing Data. *PLoS Computational Biology* 12, e1004957.

874 Huson, D.H., Mitra, S., Ruscheweyh, H.J., Weber, N., Schuster, S.C., (2011) Integrative analysis of  
875 environmental sequences using MEGAN4. *Genome Res* 21, 1552-1560.

876 Ishihama, Y., Oda, Y., Tabata, T., Sato, T., Nagasu, T., Rappsilber, J., Mann, M., (2005)  
877 Exponentially modified protein abundance index (emPAI) for estimation of absolute protein  
878 amount in proteomics by the number of sequenced peptides per protein. *Molecular & Cellular*  
879 *Proteomics* 4, 1265-1272.

880 Jagtap, P., Goslinga, J., Kooren, J.A., McGowan, T., Wroblewski, M.S., Seymour, S.L., Griffin, T.J.,  
881 (2013) A two-step database search method improves sensitivity in peptide sequence matches for  
882 metaproteomics and proteogenomics studies. *Proteomics* 13, 1352-1357.

883 Jagtap, P., McGowan, T., Bandhakavi, S., Tu, Z.J., Seymour, S., Griffin, T.J., Rudney, J.D., (2012)  
884 Deep metaproteomic analysis of human salivary supernatant. *Proteomics* 12, 992-1001.

885 Jagtap, P.D., Blakely, A., Murray, K., Stewart, S., Kooren, J., Johnson, J.E., Rhodus, N.L., Rudney, J.,  
886 Griffin, T.J., (2015) Metaproteomic analysis using the Galaxy framework. *Proteomics* 15, 3553-  
887 3565.

888 Jehmlich, N., Schmidt, F., Taubert, M., Seifert, J., von Bergen, M., Richnow, H.H., Vogt, C., (2009)  
889 Comparison of methods for simultaneous identification of bacterial species and determination of  
890 metabolic activity by protein-based stable isotope probing (Protein-SIP) experiments. *Rapid*  
891 *Commun Mass Sp* 23, 1871-1878.

892 Jehmlich, N., Vogt, C., Lunsmann, V., Richnow, H.H., von Bergen, M., (2016) Protein-SIP in  
893 environmental studies. *Curr Opin Biotech* 41, 26-33.

894 Jekel, P.A., Weijer, W.J., Beintema, J.J., (1983) Use of endoproteinase Lys-C from  
895 *Lysobacter* in protein sequence analysis. *Analytical Biochemistry* 134, 347-354.

896 Jenssen, T.K., Laegreid, A., Komorowski, J., Hovig, E., (2001) A literature network of human genes  
897 for high-throughput analysis of gene expression. *Nature Genetics* 28, 21-28.

898 Jones, A.R., Eisenacher, M., Mayer, G., Kohlbacher, O., Siepen, J., Hubbard, S.J., Selley, J.N.,  
899 Searle, B.C., Shofstahl, J., Seymour, S.L., Julian, R., Binz, P.A., Deutsch, E.W., Hermjakob, H.,  
900 Reisinger, F., Griss, J., Vizcaino, J.A., Chambers, M., Pizarro, A., Creasy, D., (2012) The mzIdentML  
901 data standard for mass spectrometry-based proteomics results. *Molecular & Cellular Proteomics*  
902 11, M111 014381.

903 Junemann, S., Prior, K., Albersmeier, A., Albaum, S., Kalinowski, J., Goesmann, A., Stoye, J.,  
904 Harmsen, D., (2014) GABenchToB: a genome assembly benchmark tuned on bacteria and  
905 benchtop sequencers. *PLoS One* 9, e107014.

906 Junemann, S., Sedlazeck, F.J., Prior, K., Albersmeier, A., John, U., Kalinowski, J., Mellmann, A.,  
907 Goesmann, A., von Haeseler, A., Stoye, J., Harmsen, D., (2013) Updating benchtop sequencing  
908 performance comparison. *Nature Biotechnology* 31, 294-296.

909 Kall, L., Canterbury, J.D., Weston, J., Noble, W.S., MacCoss, M.J., (2007) Semi-supervised learning  
910 for peptide identification from shotgun proteomics datasets. *Nature Methods* 4, 923-925.

911 Kallmeyer, J., Pockalny, R., Adhikari, R.R., Smith, D.C., D'Hondt, S., (2012) Global distribution of  
912 microbial abundance and biomass in subseafloor sediment. *P Natl AcadSci USA* 109, 16213-  
913 16216.

914 Kan, J., Hanson, T.E., Ginter, J.M., Wang, K., Chen, F., (2005) Metaproteomic analysis of  
915 Chesapeake Bay microbial communities. *Saline Systems* 1, 7.

916 Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., Tanabe, M., (2016) KEGG as a reference  
917 resource for gene and protein annotation. *Nucleic Acids Research* 44, D457-D462.

918 Karlsson, R., Davidson, M., Svensson-Stadler, L., Karlsson, A., Olesen, K., Carlsohn, E., Moore,  
919 E.R.B., (2012) Strain-Level Typing and Identification of Bacteria Using Mass Spectrometry-Based  
920 Proteomics. *Journal of Proteome Research* 11, 2710-2720.

921 Keiblinger, K.M., Wilhartitz, I.C., Schneider, T., Roschitzki, B., Schmid, E., Eberl, L., Riedel, K.,  
922 Zechmeister-Boltenstern, S., (2012) Soil metaproteomics - Comparative evaluation of protein  
923 extraction protocols. *Soil Biol Biochem* 54, 14-24.

924 Keller, A., Shteynberg, D., (2011) Software pipeline and data analysis for MS/MS proteomics: the  
925 trans-proteomic pipeline. *Methods in Molecular Biology* 694, 169-189.

926 Klamt, S., Saez-Rodriguez, J., Gilles, E.D., (2007) Structural and functional analysis of cellular  
927 networks with CellNetAnalyzer. *BMC Systems Biology* 1, 2.

928 Klose, J., (1975) Protein mapping by combined isoelectric focusing and electrophoresis of mouse  
929 tissues. A novel approach to testing for induced point mutations in mammals. *Humangenetik* 26,  
930 231-243.

931 Kohrs, F., Heyer, R., Bissinger, T., Kottler, R., Schallert, K., Püttker, S., Behne, A., Rapp, E.,  
932 Benndorf, D., Reichl, U., (2017) Proteotyping of laboratory-scale biogas plants reveals multiple  
933 steady-states in community composition. *Anaerobe*.

934 Kohrs, F., Heyer, R., Magnussen, A., Benndorf, D., Muth, T., Behne, A., Rapp, E., Kausmann, R.,  
935 Heiermann, M., Klocke, M., Reichl, U., (2014) Sample prefractionation with liquid isoelectric  
936 focusing enables in depth microbial metaproteome analysis of mesophilic and thermophilic  
937 biogas plants. *Anaerobe* 29, 59-67.

938 Kolmeder, C.A., de Been, M., Nikkila, J., Ritamo, I., Matto, J., Valmu, L., Salojarvi, J., Palva, A.,  
939 Salonen, A., de Vos, W.M., (2012) Comparative metaproteomics and diversity analysis of human  
940 intestinal microbiota testifies for its temporal stability and expression of core functions. *PLoS*  
941 *One* 7, e29913.

942 Kolmeder, C.A., Salojarvi, J., Ritari, J., de Been, M., Raes, J., Falony, G., Vieira-Silva, S., Kekkonen,  
943 R.A., Corthals, G.L., Palva, A., Salonen, A., de Vos, W.M., (2016) Faecal Metaproteomic Analysis  
944 Reveals a Personalized and Stable Functional Microbiome and Limited Effects of a Probiotic  
945 Intervention in Adults. *PLoS One* 11, e0153294.

946 Laemmli, U.K., (1970) Cleavage of structural proteins during the assembly of the head of  
947 bacteriophage T4. *Nature* 227, 680-685.

948 Lam, H., Deutsch, E.W., Eddes, J.S., Eng, J.K., King, N., Stein, S.E., Aebersold, R., (2007)  
949 Development and validation of a spectral library searching method for peptide identification  
950 from MS/MS. *Proteomics* 7, 655-667.

951 Letunic, I., Doerks, T., Bork, P., (2015) SMART: recent updates, new developments and status in  
952 2015. *Nucleic Acids Research* 43, D257-260.

953 Levenshtein, V.I., (1966) Binary codes capable of correcting deletions, insertions and reversals.  
954 *Soviet Physics Doklady* 10, 707-710.

955 Locey, K.J., Lennon, J.T., (2016) Scaling laws predict global microbial diversity. *Proceedings of the*  
956 *National Academy of Sciences* 113, 5970-5975.

957 Lu, F., Bize, A., Guillot, A., Monnet, V., Madigou, C., Chapleur, O., Mazeas, L., He, P., Bouchez, T.,  
958 (2014) Metaproteomics of cellulose methanisation under thermophilic conditions reveals a  
959 surprisingly high proteolytic activity. *Isme J* 8, 88-102.

960 MacLean, B., Tomazela, D.M., Shulman, N., Chambers, M., Finney, G.L., Frewen, B., Kern, R.,  
961 Tabb, D.L., Liebler, D.C., MacCoss, M.J., (2010) Skyline: an open source document editor for  
962 creating and analyzing targeted proteomics experiments. *Bioinformatics* 26, 966-968.

963 Martens, L., Chambers, M., Sturm, M., Kessner, D., Levander, F., Shofstahl, J., Tang, W.H., Ropp,  
964 A., Neumann, S., Pizarro, A.D., Montecchi-Palazzi, L., Tasman, N., Coleman, M., Reisinger, F.,  
965 Souda, P., Hermjakob, H., Binz, P.A., Deutsch, E.W., (2011) mzML-a Community Standard for  
966 Mass Spectrometry Data. *Molecular & Cellular Proteomics* 10.

967 Marzorati, M., Wittebolle, L., Boon, N., Daffonchio, D., Verstraete, W., (2008) How to get more  
968 out of molecular fingerprints: practical tools for microbial ecology. *Environ Microbiol* 10, 1571-  
969 1581.

970 May, D.H., Timmins-Schiffman, E., Mikan, M.P., Harvey, H.R., Borenstein, E., Nunn, B.L., Noble,  
971 W.S., (2016) An Alignment-Free "Metapeptide" Strategy for Metaproteomic Characterization of  
972 Microbiome Samples Using Shotgun Metagenomic Sequencing. *Journal of Proteome Research*  
973 15, 2697-2705.

974 Mehlan, H., Schmidt, F., Weiss, S., Schuler, J., Fuchs, S., Riedel, K., Bernhardt, J., (2013) Data  
975 visualization in environmental proteomics. *Proteomics* 13, 2805-2821.

976 Mell, P., Grance, T., (2010) The NIST Definition of Cloud Computing. *CommunAcm* 53, 50-50.

977 Mesuere, B., Debyser, G., Aerts, M., Devreese, B., Vandamme, P., Dawyndt, P., (2015) The  
978 Unipept metaproteomics analysis pipeline. *Proteomics* 15, 1437-1442.

979 Mesuere, B., Van der Jeugt, F., Devreese, B., Vandamme, P., Dawyndt, P., (2016) The unique  
980 peptidome: Taxon-specific tryptic peptides as biomarkers for targeted metaproteomics.  
981 *Proteomics* 16, 2313-2318.

982 Muth, T., Behne, A., Heyer, R., Kohrs, F., Benndorf, D., Hoffmann, M., Lehteve, M., Reichl, U.,  
983 Martens, L., Rapp, E., (2015a) The MetaProteomeAnalyzer: a powerful open-source software  
984 suite for metaproteomics data analysis and interpretation. *Journal of Proteome Research* 14,  
985 1557-1565.

986 Muth, T., Benndorf, D., Reichl, U., Rapp, E., Martens, L., (2013) Searching for a needle in a stack  
987 of needles: challenges in metaproteomics data analysis. *Molecular BioSystems* 9, 578-585.

988 Muth, T., Kolmeder, C.A., Salojarvi, J., Keskitalo, S., Varjosalo, M., Verdam, F.J., Rensen, S.S.,  
989 Reichl, U., de Vos, W.M., Rapp, E., Martens, L., (2015b) Navigating through metaproteomics  
990 data: A logbook of database searching. *Proteomics* 15, 3439-3453.

991 Muth, T., Renard, B.Y., Martens, L., (2016) Metaproteomic data analysis at a glance: advances in  
992 computational microbial community proteomics. *Expert Rev Proteomic* 13, 757-769.

993 Nesvizhskii, A.I., (2014) Proteogenomics: concepts, applications and computational strategies.  
994 *Nature Methods* 11, 1114-1125.

995 Nesvizhskii, A.I., Keller, A., Kolker, E., Aebersold, R., (2003) A statistical model for identifying  
996 proteins by tandem mass spectrometry. *Anal Chem* 75, 4646-4658.

997 O'Farrell, P.H., (1975) High resolution two-dimensional electrophoresis of proteins. *The Journal*  
998 *of Biological Chemistry* 250, 4007-4021.

999 Ondov, B.D., Bergman, N.H., Phillippy, A.M., (2011) Interactive metagenomic visualization in a  
1000 Web browser. *BMC Bioinformatics* 12, 385.

1001 Oveland, E., Muth, T., Rapp, E., Martens, L., Berven, F.S., Barsnes, H., (2015) Viewing the  
1002 proteome: How to visualize proteomics data? *Proteomics* 15, 1341-1355.

1003 Patient, S., Wieser, D., Kleen, M., Kretschmann, E., Jesus Martin, M., Apweiler, R., (2008)  
1004 UniProtJAPI: a remote API for accessing UniProt data. *Bioinformatics* 24, 1321-1322.

1005 Penzlin, A., Lindner, M.S., Doellinger, J., Dabrowski, P.W., Nitsche, A., Renard, B.Y., (2014)  
1006 Pipasic: similarity and expression correction for strain-level identification and quantification in  
1007 metaproteomics. *Bioinformatics* 30, i149-156.

1008 Perkins, D.N., Pappin, D.J.C., Creasy, D.M., Cottrell, J.S., (1999) Probability-based protein  
1009 identification by searching sequence databases using mass spectrometry data. *Electrophoresis*  
1010 20, 3551-3567.

1011 Püttker, S., Kohrs, F., Benndorf, D., Heyer, R., Rapp, E., Reichl, U., (2015) Metaproteomics of  
1012 activated sludge from a wastewater treatment plant - a pilot study. *Proteomics* 15, 3596-3601.

1013 Racusen, D., (1973) Stoichiometry of the amido black reaction with proteins. *Analytical*  
1014 *Biochemistry* 52, 96-101.

1015 Rodriguez-Valera, F., (2004) Environmental genomics, the big picture? *FEMS Microbiology*  
1016 *Letters* 231, 153-158.

1017 Ram, R.J., Verberkmoes, N.C., Thelen, M.P., Tyson, G.W., Baker, B.J., Blake, R.C., 2nd, Shah, M.,  
1018 Hettich, R.L., Banfield, J.F., (2005) Community proteomics of a natural microbial biofilm. *Science*  
1019 308, 1915-1920.

1020 Rooijers, K., Kolmeder, C., Juste, C., Dore, J., de Been, M., Boeren, S., Galan, P., Beauvallet, C., de  
1021 Vos, W.M., Schaap, P.J., (2011) An iterative workflow for mining the human intestinal  
1022 metaproteome. *BMC Genomics* 12, 6.

1023 Roume, H., Heintz-Buschart, A., Muller, E.E.L., May, P., Satagopam, V.P., Laczny, C.e., dric C.,  
1024 Narayanasamy, S., Lebrun, L.A., Hoopmann, M.R., Schupp, J.M., others, (2015) Comparative  
1025 integrated omics: identification of key functionalities in microbial community-wide metabolic  
1026 networks. *npj Biofilms and Microbiomes* 1.

1027 Sachsenberg, T., Herbst, F.A., Taubert, M., Kermer, R., Jehmlich, N., von Bergen, M., Seifert, J.,  
1028 Kohlbacher, O., (2015) MetaProSIP: automated inference of stable isotope incorporation rates in  
1029 proteins for functional metaproteomics. *Journal of Proteome Research* 14, 619-627.

1030 Sadygov, R.G., (2015) Using SFQUFST with Theoretically Complete Sequence Databases. *J Am Soc*  
1031 *Mass Spectr* 26, 1858-1864.

1032 Saito, M.A., Dorsk, A., Post, A.F., Mcllvin, M.R., Rappe, M.S., DiTullio, G.R., Moran, D.M., (2015)  
1033 Needles in the blue sea: Sub-species specificity in targeted protein biomarker analyses within the  
1034 vast oceanic microbial metaproteome. *Proteomics* 15, 3521-3531.

1035 Schlüter, A., Bekel, T., Diaz, N.N., Dondrup, M., Eichenlaub, R., Gartemann, K.H., Krahn, I., Krause,  
1036 L., Kromeke, H., Kruse, O., Mussnug, J.H., Neuweger, H., Niehaus, K., Pühler, A., Runte, K.J.,  
1037 Szczepanowski, R., Tauch, A., Tilker, A., Viehover, P., Goesmann, A., (2008) The metagenome of a  
1038 biogas-producing microbial community of a production-scale biogas plant fermenter analysed by  
1039 the 454-pyrosequencing technology. *Journal of Biotechnology* 136, 77-90.

1040 Schneider, T., Schmid, E., de Castro, J.V., Jr., Cardinale, M., Eberl, L., Grube, M., Berg, G., Riedel,  
1041 K., (2011) Structure and function of the symbiosis partners of the lung lichen (*Lobariapulmonaria*  
1042 *L. Hoffm.*) analyzed by metaproteomics. *Proteomics* 11, 2752-2756.

1043 Seifert, J., Herbst, F.A., Nielsen, P.H., Planes, F.J., Jehmlich, N., Ferrer, M., von Bergen, M., (2013)  
1044 Bioinformatic progress and applications in metaproteogenomics for bridging the gap between

1045 genomic sequences and metabolic functions in microbial communities. *Proteomics* 13, 2786-  
1046 2804.

1047 Shevchenko, A., Sunyaev, S., Loboda, A., Shevchenko, A., Bork, P., Ens, W., Standing, K.G., (2001)  
1048 Charting the proteomes of organisms with unsequenced genomes by MALDI-quadrupole time-  
1049 of-flight mass spectrometry and BLAST homology searching. *Anal Chem* 73, 1917-1926.

1050 Stolze, Y., Bremges, A., Rummig, M., Henke, C., Maus, I., Puhler, A., Sczyrba, A., Schluter, A.,  
1051 (2016) Identification and genome reconstruction of abundant distinct taxa in microbiomes from  
1052 one thermophilic and three mesophilic production-scale biogas plants. *Biotechnology for*  
1053 *Biofuels* 9.

1054 Sturm, M., Bertsch, A., Gropl, C., Hildebrandt, A., Hussong, R., Lange, E., Pfeifer, N., Schulz-  
1055 Trieglaff, O., Zerck, A., Reinert, K., Kohlbacher, O., (2008) OpenMS - an open-source software  
1056 framework for mass spectrometry. *BMC Bioinformatics* 9, 163.

1057 Suzek, B.E., Huang, H., McGarvey, P., Mazumder, R., Wu, C.H., (2007) UniRef: comprehensive and  
1058 non-redundant UniProt reference clusters. *Bioinformatics* 23, 1282-1288.

1059 Tanca, A., Palomba, A., Deligios, M., Cubeddu, T., Fraumene, C., Biossa, G., Pagnozzi, D., Addis,  
1060 M.F., Uzzau, S., (2013) Evaluating the impact of different sequence databases on metaproteome  
1061 analysis: insights from a lab-assembled microbial mixture. *PLoS One* 8, e82981.

1062 Tanca, A., Palomba, A., Fraumene, C., Pagnozzi, D., Manghina, V., Deligios, M., Muth, T., Rapp, E.,  
1063 Martens, L., Addis, M.F., Uzzau, S., (2016) The impact of sequence database choice on  
1064 metaproteomic results in gut microbiota studies. *Microbiome* 4.

1065 Tang, H., Li, S., Ye, Y., (2016) A Graph-Centric Approach for Metagenome-Guided Peptide and  
1066 Protein Identification in Metaproteomics. *PLOS Computational Biology* 12, e1005224.

1067 Tatusov, R.L., Galperin, M.Y., Natale, D.A., Koonin, E.V., (2000) The COG database: a tool for  
1068 genome-scale analysis of protein functions and evolution. *Nucleic Acids Research* 28, 33-36.

1069 Timmins-Schiffman, E., May, D.H., Mikan, M., Riffle, M., Frazar, C., Harvey, H.R., Noble, W.S.,  
1070 Nunn, B.L., (2017) Critical decisions in metaproteomics: achieving high confidence protein  
1071 annotations in a sea of unknowns. *ISME J* 11, 309-314.

1072 Tobalina, L., Bargiela, R., Pey, J., Herbst, F.A., Lores, I., Rojo, D., Barbas, C., Pelaez, A.I., Sanchez,  
1073 J., von Bergen, M., Seifert, J., Ferrer, M., Planes, F.J., (2015) Context-specific metabolic network  
1074 reconstruction of a naphthalene-degrading bacterial community guided by metaproteomic data.  
1075 *Bioinformatics* 31, 1771-1779.

1076 Tyanova, S., Temu, T., Cox, J., (2016) The MaxQuant computational platform for mass  
1077 spectrometry-based shotgun proteomics. *Nature Protocols* 11, 2301-2319.

1078 Tyson, G.W., Chapman, J., Hugenholtz, P., Allen, E.E., Ram, R.J., Richardson, P.M., Solovyev, V.V.,  
1079 Rubin, E.M., Rokhsar, D.S., Banfield, J.F., (2004) Community structure and metabolism through  
1080 reconstruction of microbial genomes from the environment. *Nature* 428, 37-43.

1081 UniProt, C., (2015) UniProt: a hub for protein information. *Nucleic Acids Research* 43, D204-212.

1082 Vaudel, M., Barsnes, H., Berven, F.S., Sickmann, A., Martens, L., (2011) SearchGUI: An open-  
1083 source graphical user interface for simultaneous OMSSA and X!Tandem searches. *Proteomics* 11,  
1084 996-999.



1085 Vaudel, M., Sickmann, A., Martens, L., (2010) Peptide and protein quantification: a map of the  
1086 minefield. *Proteomics* 10, 650-670.

1087 Vizcaino, J.A., Csordas, A., Del-Toro, N., Dianes, J.A., Griss, J., Lavidas, I., Mayer, G., Perez-Riverol,  
1088 Y., Reisinger, F., Ternent, T., Xu, Q.W., Wang, R., Hermjakob, H., (2016) 2016 update of the PRIDE  
1089 database and its related tools. *Nucleic Acids Research* 44, 11033.

1090 Wilm, M., Shevchenko, A., Houthaeve, T., Breit, S., Schweigerer, L., Fotsis, T., Mann, M., (1996)  
1091 Femtomole sequencing of proteins from polyacrylamide gels by nano-electrospray mass  
1092 spectrometry. *Nature* 379, 466-469.

1093 Wilmes, P., Andersson, A.F., Lefsrud, M.G., Wexler, M., Shah, M., Zhang, B., Hettich, R.L., Bond,  
1094 P.L., VerBerkmoes, N.C., Banfield, J.F., (2008) Community proteogenomics highlights microbial  
1095 strain-variant protein expression within activated sludge performing enhanced biological  
1096 phosphorus removal. *Isme J* 2, 853-864.

1097 Wilmes, P., Bond, P.L., (2004) The application of two-dimensional polyacrylamide gel  
1098 electrophoresis and downstream analyses to a mixed community of prokaryotic microorganisms.  
1099 *Environ Microbiol* 6, 911-920.

1100 Wilmes, P., Bond, P.L., (2006) Metaproteomics: studying functional gene expression in microbial  
1101 ecosystems. *Trends Microbiol* 14, 92-97.

1102 Wilmes, P., Heintz-Buschart, A., Bond, P.L., (2015) A decade of metaproteomics: Where we stand  
1103 and what the future holds. *Proteomics* 15, 3409-3417.

1104 Wohlbrand, L., Trautwein, K., Rabus, R., (2013) Proteomic tools for environmental microbiology--  
1105 a roadmap from sample preparation to protein identification and quantification. *Proteomics* 13,  
1106 2700-2730.

1107 Yamada, T., Letunic, I., Okuda, S., Kanehisa, M., Bork, P., (2011) iPath2.0: interactive pathway  
1108 explorer. *Nucleic Acids Research* 39, W412-415.

1109 Yates, A., Akanni, W., Amode, M.R., Barrell, D., Billis, K., Carvalho-Silva, D., Cummins, C.,  
1110 Clapham, P., Fitzgerald, S., Gil, L., Giron, C.G., Gordon, L., Hourlier, T., Hunt, S.E., Janacek, S.H.,  
1111 Johnson, N., Juettemann, T., Keenan, S., Lavidas, I., Martin, F.J., Maurel, T., McLaren, W.,  
1112 Murphy, D.N., Nag, R., Nuhn, M., Parker, A., Patricio, M., Pignatelli, M., Rahtz, M., Riat, H.S.,  
1113 Sheppard, D., Taylor, K., Thormann, A., Vullo, A., Wilder, S.P., Zadissa, A., Birney, E., Harrow, J.,  
1114 Muffato, M., Perry, E., Ruffier, M., Spudich, G., Trevanion, S.J., Cunningham, F., Aken, B.L.,  
1115 Zerbino, D.R., Flicek, P., (2016) Ensembl 2016. *Nucleic Acids Research* 44, D710-716.

1116 Zhang, X., Ning, Z., Mayne, J., Moore, J.I., Li, J., Butcher, J., Deeke, S.A., Chen, R., Chiang, C.K.,  
1117 Wen, M., Mack, D., Stintzi, A., Figeys, D., (2016) MetaPro-IQ: a universal metaproteomic  
1118 approach to studying human and mouse gut microbiota. *Microbiome* 4, 31.

1119 Zybilov, B.L., Florens, L., Washburn, M.P., (2007) Quantitative shotgun proteomics using a  
1120 protease with broad specificity and normalized spectral abundance factors. *Molecular*  
1121 *Biosystems* 3, 354-360.

Table 1

Issue	Solution/ bioinformatic strategie	Reference
<b>Grouping of redundant homologous proteins</b>	1. Flexible grouping to metaproteins based on protein, peptide and taxonomy similarity	MetaProteomeAnalyzer (Muth et al., 2015a)
<b>Database tailoring</b>	2. Grouping by shared peptide 1. Two step database search 2. Metapeptide database 3. A "Graph-Centric Approach"	Prophane (Schneider et al., 2011) (Jagtap et al., 2013) (May et al., 2016) Graph2Pep/ Graph2Prot (Zhang et al., 2016)
<b>Taxonomic and functional evaluation</b>	1. Calculate taxonomic value for each identified peptide (LCA) and visualize results 2. Calculate taxonomic value for peptides using peptide similarity estimation and expression level weighting 3. Taxonomic evaluation (LCA) and functional prediction using RPSBLAST or HMMER3 4. Taxonomic (LCA) and functional evaluation using ECs, KEGG Ontologies and KEGG Pathways. Unknown sequences can be annotated using Diamond. 5. Taxonomic (LCA) and functional evaluation using UniProt Keywords, ECs, KEGG Ontologies, KEGG Pathways. Unknown sequences can be annotated using BLAST.	UniPept (Mesuere et al., 2015) Pipasic (Penzlin et al., 2014) Prophane (Schneider et al., 2011) Megan CE (Huson et al., 2016) MPA (Muth et al., 2015a)
<b>Storing and deployment of big data</b>	1. Scalable set of sequence databases and specific database search algorithm	Compile and Blazmass (Chatterjee et al., 2016)
<b>Quantitation</b>	1. Detection and quantification of isotope ratios for Protein-SIP	MetaProSip (Sachsenberg et al., 2015)

**Table 2**

<b>Rule</b>	<b>Principle</b>	<b>Explanation</b>	<b>Reference</b>
<b>Protein rule</b>	1. UniRef-Cluster	Grouping of proteins when they have 50%, 90% or 100% sequence similarity. Protein clustering provided by UniRef Cluster [Suzek2007].	(Lu et al., 2014; Suzek et al., 2007)
	2. KEGG Ontologies	Grouping of proteins when they are similar to functional classified genes within KEGG Ontology [Mai 2005]. KEGG Ontologies are provide by UniProtKB databases [JAPI PAPER].	(Gotelli et al., 2012; Kanehisa et al., 2016)
<b>Peptide rule</b>	1. Shared peptide set	Group proteins when they share the same peptides.	(Keiblinger et al., 2012; Kolmeder et al., 2012; Schneider et al., 2011)
	2. One shared peptide	Group proteins when they have one identified peptide in common	(Kohrs et al., 2014; Lu et al., 2014)
	3. One shared peptide + Levenshtein, distance < 2	Group proteins when they share the same peptides, but not if they have two similar peptides with less than 2 point mutations differences. This tracks the production of one protein by different microorganisms.	(Muth et al., 2015a)
<b>Taxonomy rule</b>	1. Phylogenetic affiliation	Extends other rules by a certain phylogenetic affiliation.	(Muth et al., 2015a)

Table 3

Issue	Name/ principle	Explanation	Reference
Taxonomic classification	1. Lowest common ancestor	Define taxonomy as the lowest common ancestor into the phylogenetic tree.	(Huson et al., 2011; Jagtap et al., 2012)
	2. Weighted lowest common ancestor	Adjust the lowest common ancestor by unique identification for the single taxa.	(Huson et al., 2016)
	3. Peptide similarity estimation and expression level weighting	Weight taxonomy of identified peptides by their spectra abundance and their occurrence in a reference proteome.	(Penzlin et al., 2014)
	4. Unique peptides	Define taxonomy and taxonomy profiles only based on unique peptides.	(Rooijers et al., 2011; Karlsson et al., 2012)
Functional classification	1. KEGG Orthologies (KO)	Grouping of genes with same function by sequence similarity.	(Kanehisa et al., 2016)
	2. Cluster of orthologues genes (COG)	Grouping of genes with same function by sequence similarity.	(Tatusov et al., 2000)
	3. Evolutionary genealogy of genes: Non-supervised Orthologous Groups (eggNOG)	Extension off COG by non-supervised orthologous groups constructed from numerous organisms.	(Huerta-Cepas et al., 2016)
	4. PFAM	Database of conserved functional units, represented by a set of aligned sequences with their probabilistic representation (hidden Markov model).	(Finn et al., 2016)
	5 TIGRFAM	Database of conserved functional units, represented by a set of aligned sequences with their probabilistic representation (hidden Markov model). In contrast to PFAM TIGRFAM emphasize protein function and	(Haft et al., 2013)
	6. SMART	Functional domain database based on manually curated hidden Markov models.	(Letunic et al., 2015)
	7. InterPro	Functional analyses of protein sequences by classifying them into families and predicting the presence of domains and important sites. Signatures are provided by 14 different member databases (among others PFAM, TIGRFAMS, SMART).	(Finn et al., 2017)
Pathway mapping	8. Enzyme Comission number (EC)	Numerical classification scheme for enzymes, based on the chemical reactions they catalyze	(Bairoch, 2000)
	9. UniProt Keywords	Hierachical classification of protein functions.	(UniProt, 2015)
	10. Gene ontologies	Hierachical classification of protein functions.	(Ashburner et al., 2000)
	1. MetaCyc	Curated database of experimentally confirmed metabolic pathways.	(Caspi et al., 2016)
	2. KEGG pathways	Collection of manually drawn pathway maps representing knowledge on the molecular interaction and reaction networks.	(Kanehisa et al., 2016)
3. Reactome	Pathway database.	(Fabregat et al., 2016)	
4. Interactive Pathways Explorer (iPath)	Web-based tool for the visualization, analysis and customization of pathways maps.	(Yamada et al., 2011)	
5. CellNetAnalyzer	MATLAB toolbox providing computational methods and algorithms for exploring structural and functional properties of metabolic, signaling, and regulatory networks.	(Klamt et al., 2007)	

Figure 1

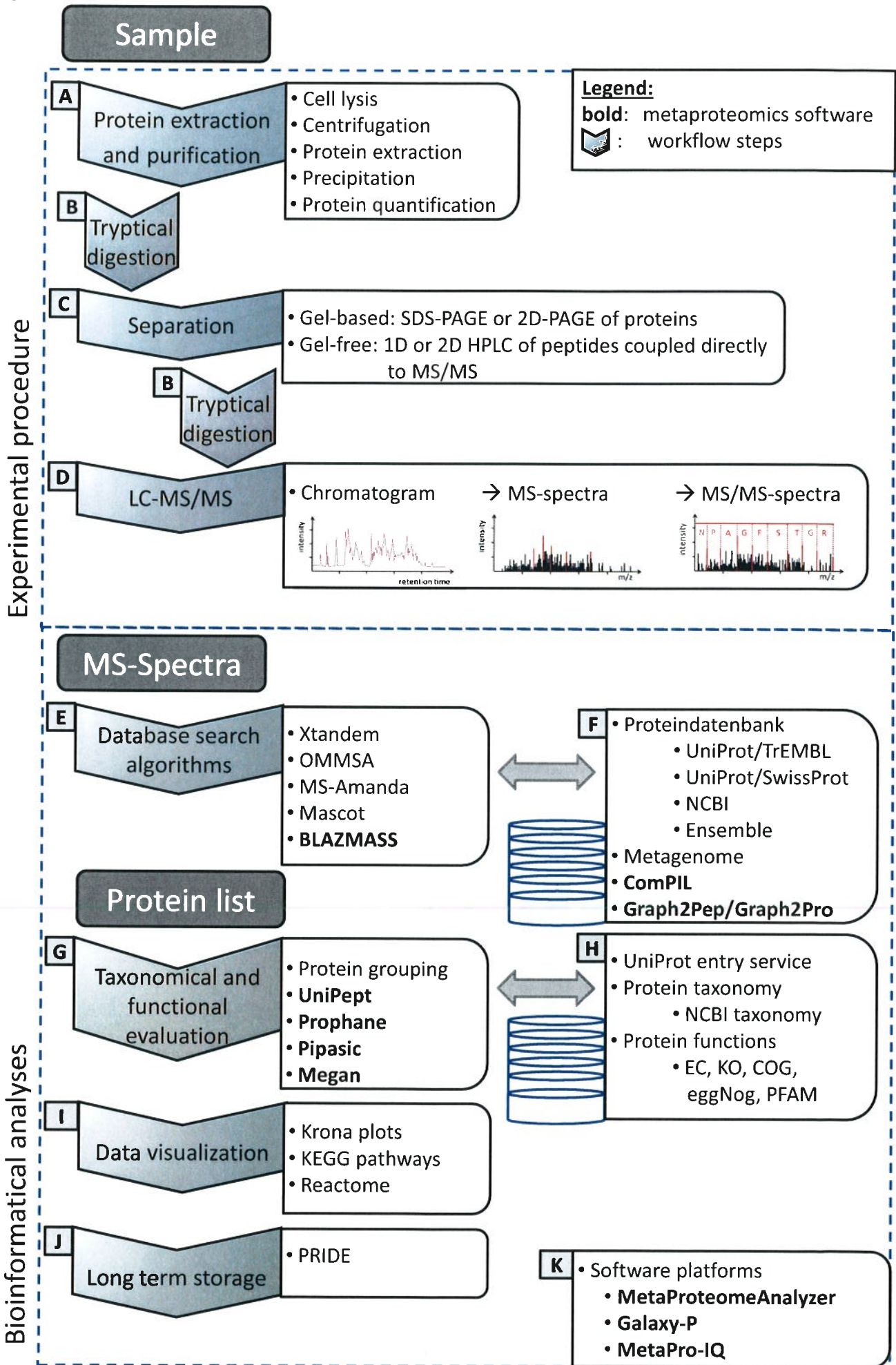


Figure 2

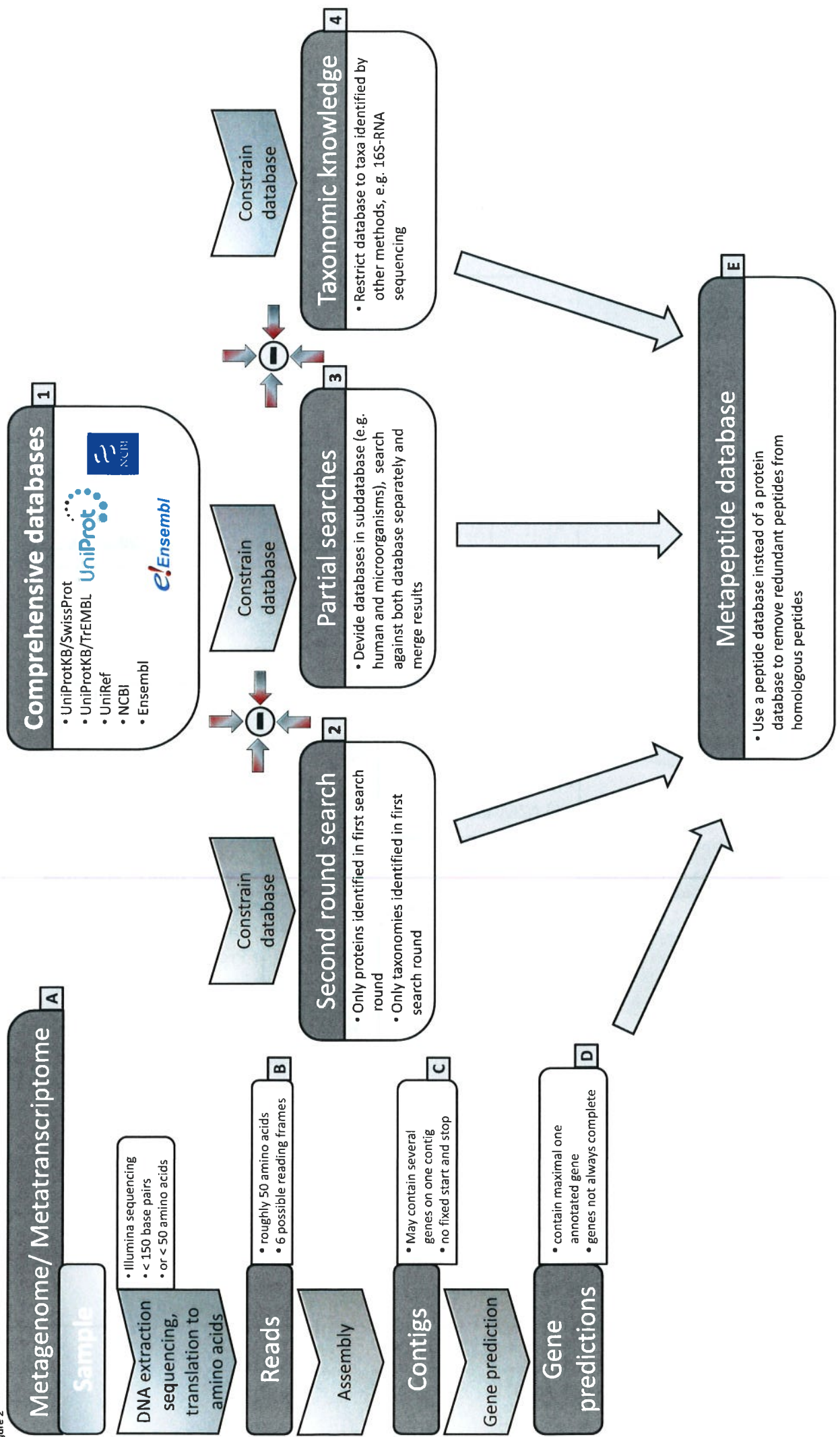
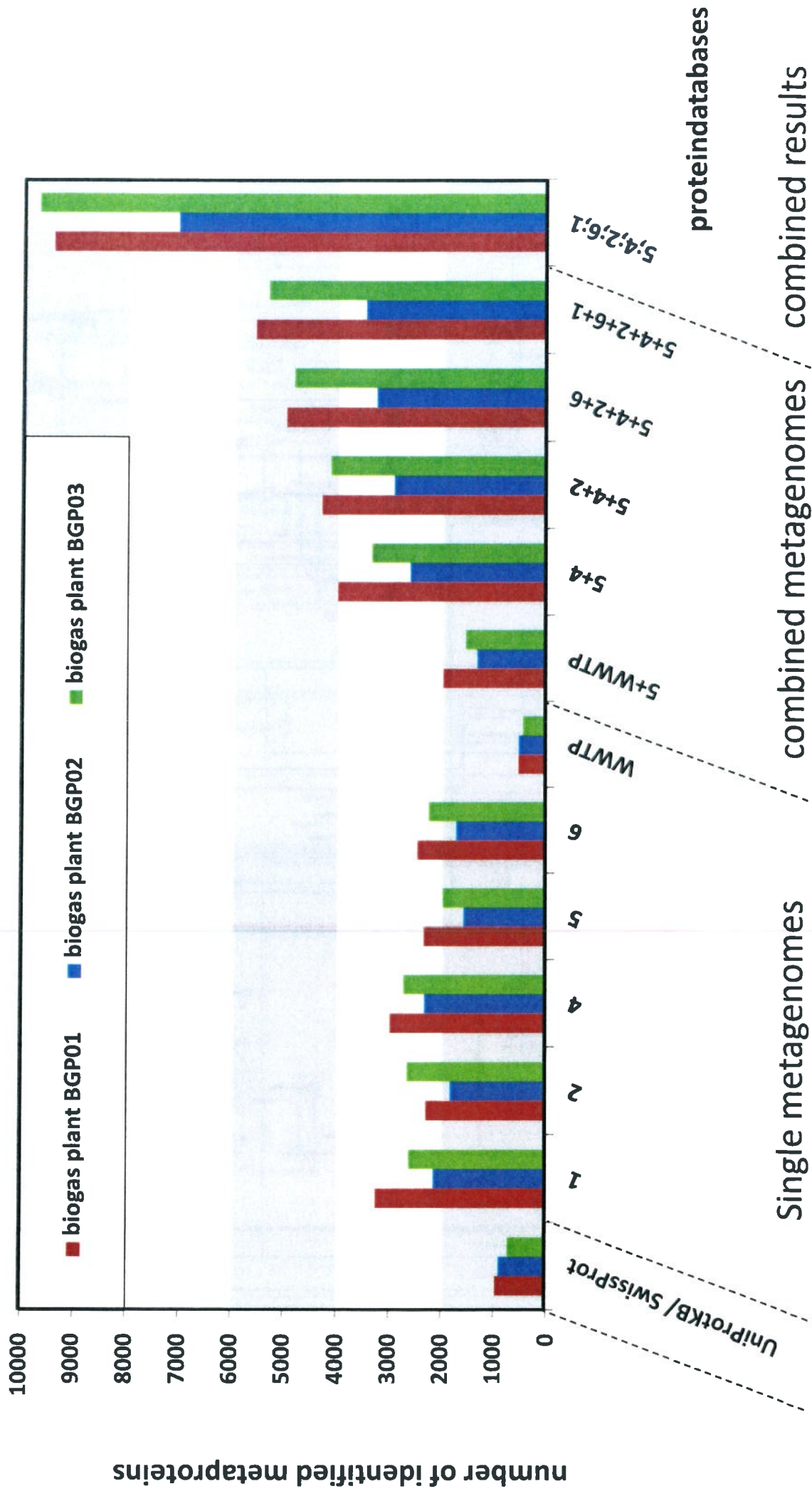


Figure 3



CARBON METABOLISM

