# Assessing species biomass contributions in microbial communities via metaproteomics

**Authors:**

Manuel Kleiner[1,2]*[§], Erin Thorson[1]*, Christine E. Sharp[1], Xiaoli Dong[1], Dan Liu[1], Carmen Li[3] and Marc Strous[1]


[1]Department of Geoscience, University of Calgary, Calgary, AB, Canada

[2]Department of Plant and Microbial Biology, North Carolina State University, Raleigh, North Carolina, USA

[3]Department of Biological Sciences, University of Calgary, Calgary, AB, Canada


*contributed equally

§corresponding author

Email address: manuel.kleiner@ucalgary.ca

Phone: +1-403-210-6617

Key words: metagenomics, metaproteomics, microbial ecology, microbiome, soda lake biomats, oral microbiome, amplicon sequencing, tag sequencing, microbiota

## Abstract

20  Assessment of microbial community composition is the cornerstone of microbial ecology. Microbial

21  community composition can be analyzed by quantifying cell numbers or by quantifying biomass for

22  individual populations. However, as cell volumes can differ by orders of magnitude, these two approaches

23  yield vastly different results. Methods for quantifying cell numbers are already available (e.g. fluorescence

24  *in situ* hybridization, 16S rRNA gene amplicon sequencing), yet methods for assessing community

25  composition in terms of biomass are lacking.

26  We developed metaproteomics based methods for assessing microbial community composition using

27  protein abundance as a measure for biomass contributions of individual populations. We optimized the

28  accuracy and sensitivity of the method using artificially assembled microbial communities and found that

29  it is less prone to some of the biases found in sequencing-based methods. We applied the method using

30  communities from two different environments, microbial mats from two alkaline soda lakes and saliva

31  from multiple individuals.

## Introduction

32

33 Microbial communities are ubiquitous in all environments on Earth that support life and they play crucial

34 roles in global biogeochemical cycles, plant and animal health, and biotechnological processes [1]. One of

35 the most basic and crucial parameters that microbial ecologists determine when studying these

36 communities is their taxonomic composition. Currently, all methods for assessing community composition

37 provide a direct or indirect measure of cell numbers per taxon. For example, fluorescence *in situ*

38 hybridization (FISH) provides direct cell counts [2], while metagenomics or 16S rRNA gene amplicon

39 sequencing provide a more indirect measure of cell numbers as they essentially measure gene or genome

40 copy numbers [3].

41 Cell numbers, however, are often not the best measure for a species' contribution to a community, because

42 microbes can differ by several orders of magnitude in biomass and activity. For example, the unicellular

43 eukaryote *Schizosaccharomyces pombe* has a cell volume and per cell proteinaceous biomass that is ~6000

44 fold higher than that of the bacterium *Mycoplasma pneumoniae* [4]. Therefore, the development of methods

45 for the assessment of biomass contributions of community members is critical. Recently, FISH based

46 methods for the estimation of biovolume fractions of community members have been developed [5],

47 however, these methods are limited to a few community members as a separate fluorescently-labeled

48 probe is needed for each taxon that investigators want to analyze. Currently, there are no methods

49 available to estimate the biomass contribution and activity of individual community members on a large

50 scale.

51 Metaproteomics is an umbrella term for methods for identifying and quantifying proteins in microbial

52 communities [6] and may represent a suitable approach for assessing the taxonomic composition of a

53 microbial community based on species biomass contributions. Since proteins contribute a large amount of

54 biomass in microbial cells e.g. 55% of *Escherichia coli* dry weight (BNID 104954)[7], proteinaceous

55 biomass can be a good estimator of biomass contributions. Additionally, since proteins are the molecules

56 that provide the biological activities to cells, metaproteomics may also provide estimates of activities. In

57 recent years, several studies have been published, including some from our laboratory, which used

58 metaproteomic data to quantify biomass contributions of community members [8-10]. However, methods for

59 biomass assessment with metaproteomics have not been thoroughly developed and validated, and several

60 challenges and questions have not been addressed. The major challenge is the so-called protein inference

61 problem of shotgun proteomic approaches [11]. In shotgun proteomics, which is the most widely used

62 proteomic approach, proteins are identified by matching mass spectrometry derived peptide sequences to

63 protein sequences. The protein inference problem describes the fact that often the same peptide sequence

64 can match to multiple different proteins, which can lead to ambiguous protein identifications. This

3

65    problem was originally noted for eukaryotes, which often have multiple, very similar isoforms of a protein

66    [11]; however, the problem can be much more severe in metaproteomics, because in metaproteomic analysis

67    there are tens to hundreds of species that all have protein sequences sharing peptides with sequences from

68    other species. The protein inference problem will thus lead to incorrect interpretations of taxonomic

69    composition of metaproteomes [12] if not properly addressed. In fact, the protein inference problem is so

70    pervasive that it has been advantageously used in metaproteomics for cross-strain and -species protein

71    identification by using protein sequences from organisms closely related to the ones in the analyzed

72    community [13, 14]. Other challenges and questions include: How much mass spectrometric data is needed to

73    accurately quantify species in a community? And how do potentially incomplete protein sequence

74    databases for protein identification affect the outcome of the quantification?

75    Here we address these challenges and questions to develop a simple and robust metaproteomics-based

76    workflow for assessing species biomass contributions in microbial communities. Furthermore, we provide

77    a large dataset of metaproteomic, metagenomic and 16S rRNA gene amplicon data from three types of

78    artificial microbial communities for future method development and testing.

4

## Results

80  Overall, our method for species biomass assessment is similar to a basic workflow for metaproteomic

81  protein identification and label-free quantification (Fig. 1). However, in contrast to protein and function

82  focused metaproteomics, the label-free quantification data (spectral counts or peptide intensities) are not

83  summed for individual proteins, but rather for individual species or higher level taxonomic groups.

84  Importantly, the quantification data is summed based on the taxonomic assignment of inferred proteins

85  and not based on the taxonomic assignment of peptide identifications, because as mentioned above

86  peptides are frequently associated with multiple proteins from different taxa. Additionally, we assume that

87  a well annotated protein sequence database, which matches the studied environment as closely as possible,

88  is used. This database could either be based on metagenomes derived from samples that match the

89  metaproteomic samples or for well-studied environments, such as the human microbiome, a

90  comprehensive, non-redundant set of sequences from public databases.

91  For this study, we used the Proteome Discoverer software (version 2.0, Thermo Scientific) and MaxQuant

92  for protein identification, inference and quantification [15]. However, the methods discussed here are not

93  platform dependent and can be implemented on many other platforms using the mock community data that

94  we provide in this study for optimization.

### Achieving high specificity on species level protein identification with minimal losses in sensitivity

97  Before starting the actual species quantification, we first addressed the above mentioned protein inference

98  problem. For this we used proteomes from pure culture organisms and simulated metagenomic databases

99  to test what kind of protein inference parameters can be used to eliminate unwanted cross-strain and -

100 species protein identifications (specificity), while still identifying a large number of proteins for

101 quantification (sensitivity) (Fig. 2). We tested a variety of protein inference methods for the following four

102 scenarios: the simulated metagenomic database contained the protein sequences of the analyzed organism

103 and the sequences of (a) a very closely related strain from the same species, (b) several closely related

104 species from the same genus, (c) several related species from closely related genera, (d) no other

105 representative from the same domain (analyzed organism for (d) is an archaeon).

106 Commonly used protein inference filters that filter protein identifications simply for a false discovery rate

107 (FDR) of 5% based on target-decoy database searches (SQ 5% FDR) fail to identify proteins from the

108 analyzed organisms with high specificity for all scenarios except scenario (d) (Fig. 2). The same remains

109 true when another commonly used criterion of requiring two unique peptides is added (2 UP ⊂ SQ 5%

110 FDR). Here it is important to note that different protein identification platforms implement "unique

111 peptides" differently. While, for example, a "unique peptide" in Proteome Discoverer and MaxQuant

112     refers to a peptide that is unique to a group of highly similar protein sequences (protein group), it can also

113     refer to a peptide that is unique to a single protein sequence in other protein identification platforms. This

114     shows that to obtain high specificity on a taxonomic level protein inference has to be done differently

115     from common practices.

116     We tested five additional protein inference and filtering strategies (Fig. 2) and found that there are

117     multiple strategies that result in high specificity down to the species level i.e. removing almost all cross-

118     species protein identifications, while at the same time maintaining a high sensitivity i.e. the number of

119     identified proteins for the target organism is only slightly reduced as compared to the less specific

120     approaches (Fig. 2). As expected, the approaches tested were unable to resolve cross-strain protein

121     identifications in scenario (a) because protein sequences from the two strains were nearly identical in

122     many cases. This suggests that it might be beneficial to remove highly similar sequences by sequence

123     clustering when creating metaproteomic databases. Such a clustering would reduce database size,

124     redundancy and the number of ambiguous strain level protein identifications, thus providing clearer

125     species level identifications.

126     Going forward, we used two protein inference strategies for this study. The first strategy relies on the

127     SEQUEST algorithm for peptide identification and the Fido method for protein inference [16](2 PU $\subset$ SQ

128     Fido). Fido is available as a standalone program (https://noble.gs.washington.edu/proj/fido/) and as an

129     advanced implementation with convolution trees in Proteome Discoverer (FidoCT) [17]. For this strategy,

130     only proteins that are identified by FidoCT with an FDR of 5% and have at least two protein unique

131     peptides, are considered. The second strategy (SQ Fido $\cap$ (MQ || 3 PU), only considers proteins as

132     confidently inferred if they are identified by both FidoCT (FDR of 5%) and MaxQuant (FDR of 1%, at

133     least one unique peptide). Additionally, proteins are considered as confidently inferred if they have at least

134     three protein unique peptides in the FidoCT result even if not identified by MaxQuant.

135     We are confident that many more strategies can be devised with the pure culture proteome data and the

136     simulated metagenomic database, which we provide through the PRIDE repository (PXD006118).

**Label-free quantification enables accurate measurement of relative species protein abundance (proteinaceous biomass contribution per species)**

139     We used three types of mock communities to test and validate the methods for quantifying species

140     biomass contribution in microbial communities. The three communities were assembled using 32 species

141     and strains of Archaea, Bacteria, Eukaryotes and Bacteriophages (Fig. 3a, Supplementary Tables 1 to 3).

142     Some of the bacterial strains were very closely related, but still distinguishable at the protein and

143     nucleotide sequence level. These included the *Rhizobium leguminosarum* and *Staphylococcus aureus*

6

144    strains. The three *Salmonella enterica* serotype typhimurium strains, however, only differed by a few

145    mutations or the presence of an additional plasmid. The UNEVEN mock community was designed to cover

146    a large range of species abundances both at the level of cell number and proteinaceous biomass to test for

147    the dynamic range and detection limits of the quantification methods (Fig. 3a). The EQUAL PROTEIN

148    AMOUNT and EQUAL CELL NUMBER mock communities contained either the same amount of protein for all

149    community members with varying cell numbers or the same number of cells for all members with varying

150    amounts of protein. Since the bacteriophages yield very little protein even if high particle numbers are

151    used we mixed them at a 10x lower ratio into the EQUAL PROTEIN AMOUNT community.

152    We tested three of the most commonly used label-free quantification methods for their accuracy in

153    measuring proteinaceous biomass contributions of individual species (Fig. 3b and 3c). These methods

154    included counting and summing of peptide-spectrum matches (PSMs), summing of peptide ion intensities

155    using only unique peptides (u intensities), and summing of peptide ion intensities using razor and unique

156    peptides as implemented in MaxQuant (r+u intensities) [15]. The input for these quantification methods were

157    two 8 hour long 1D-LC-MS/MS runs per sample (see methods).

158    All three methods produced a good representation of the diversity in the mock communities and detected

159    almost all species. The only exceptions were some of the bacteriophages and *N. ureae*, which were mixed

160    into the samples in low total protein amounts (Fig. 3b). As expected it was impossible to distinguish the

161    three *Salmonella enterica* strains and thus they are represented in Fig. 3b as one row. All three methods

162    performed similarly well when comparing the protein input amounts for the communities with the actual

163    measurements (Fig. 3c). In most cases the values for the measured % divided by the input % centered on

164    the expected value of 1, with the median values being very close to 1. Differences between the

165    quantification methods became apparent only for the UNEVEN community. Both peptide-intensity based

166    methods deviated strongly from the expectation and underestimated the abundance for many species. The

167    PSM based method was more robust for estimating abundances for the UNEVEN community which is

168    characterized by large differences in cell numbers and total protein amount between species.

169    **Metaproteomics is more accurate in assessing relative proteinaceous biomass**

170    **contributions as compared to metagenomics and amplicon sequencing methods**

171    We subjected subsamples of the above described mock communities to shotgun metagenomic sequencing

172    and 16S rRNA gene amplicon sequencing to test how well these commonly used methods for community

173    composition assessment estimate the proteinaceous biomass and cell number of species in communities in

174    comparison to the metaproteomic method presented here.

175    We sequenced 16S rRNA gene amplicons for four biological replicates of each community type yielding

176    an average of 5356 high-quality amplicon sequences per replicate (minimum 1686 and maximum 9986

7

177    sequences). The amplicon sequences were clustered into 21 operational taxonomic units (OTUs) using the

178    MetaAmp pipeline (version 1.3) [18]. Four of these OTUs were identified as Illumina in-run cross

179    contaminants from unrelated samples that were sequenced on the same lane. The remaining 17 OTUs

180    were taxonomically classified by MetaAmp at the genus level. A species level classification was not

181    possible because of the limited information content of the amplicon sequences. This meant, for example,

182    that there were three OTUs that were classified as *Pseudomonas*. Therefore, we had to assign the OTUs to

183    their respective species using BLASTn against the NCBI nr database and the prior knowledge about the

184    content of our mock communities. As expected none of the bacteriophages were detected by amplicon

185    sequencing due to the absence of a 16S rRNA gene in these phages (Fig. 3d). We also did not detect the

186    Archaeon *N. viennensis*, the eukaryotic green algae *Chl. reinhardtii* and six of the bacterial species by

187    amplicon sequencing. The primer pair that we used to generate the amplicons is optimized for the greatest

188    possible coverage of the bacterial domain [19], therefore it was not surprising that *N. viennensis* and *Chl.*

189    *reinhardtii* were not detected, although we successfully amplified at least the chloroplast sequence of

190    green algae using this primer pair in the past (data not shown). The failure to detect some of the bacteria in

191    all replicates is harder to explain. We have successfully generated amplicons from pure cultures of *N.*

192    *europaeae*, *N. ureae* and *N. multiformis* in the past with the primer pair used here (data not shown), thus

193    we have to assume that these species were not detected due to their low abundance in the UNEVEN

194    community samples or due to a primer bias leading to preferential amplification of the other bacterial

195    species. Such primer biases are a known problem for 16S rRNA gene amplicon sequencing [3, 20]. For the *R.*

196    *leg.* bv. viciae and *S. aureus* strains the amplicon sequences did not distinguish between each of the two

197    strains in the samples and thus only a minimum of one strain detection per species could be corroborated.

198    Metagenomic sequencing of 3 biological replicates of each community type yielded on average 33.5 M 75

199    bp reads (max. 37 M, min. 21 M). The same DNA was used for the metagenomic sequencing and the 16S

200    rRNA gene sequencing, however, only 3 of the 4 available biological replicates were metagenome

201    sequenced. For quantification, we mapped the metagenomic reads to the reference genomes and

202    assembled bins of the mock community members and normalized to the respective genome sizes.

203    All, except for one, organisms in the mock samples were detected by shotgun metagenomics, even

204    including the single-stranded DNA bacteriophage M13. As expected, the only organism not detected by

205    shotgun metagenomics was the single-stranded RNA bacteriophage F2, because the DNA extraction and

206    sequencing library preparation methods used effectively exclude RNA from being sequenced.

207    Surprisingly, the metagenomic sequencing yielded only a small number of reads for the green algae *Chl.*

208    *reinhardtii*, which was in no way representative of the input cell number for the mock communities (Fig.

209    3d). *Chl. reinhardtii* was much better represented in the metaproteomic data. One potential explanation for

210    the underrepresentation of *Chl. reinhardtii* in the sequencing data could be a bias of the DNA extraction

8

211    method used. The bead beating method used for DNA extraction, however, was quite rigorous. The

212    metagenomic data provided by far the best representation of the bacteriophages in the samples, with the

213    exception of the F2 phage, which was only detected in the metaproteomes.

214    Comparing all three methods, metaproteomics provided the most accurate estimates of proteinaceous

215    biomass for each species in the samples (Fig. 3e and 3f). The average x fold deviations of the measured

216    abundances from the expected abundance based on protein input were significantly lower for

217    metaproteomics as compared to metagenomic and amplicon sequencing (p-value <0.01, Supplementary

218    Table 4). Both the metagenomic and the amplicon based quantifications deviated from the actual values

219    when it came to assessing proteinaceous biomass. Particularly the metagenomic quantification produced

220    some extreme outliers (Fig. 3e, Supplementary Table 4).

221    All three methods performed badly, when it came to estimating the species cell numbers in the samples

222    (Fig. 3g), as they showed major deviations from the actual values in at least some of the mock

223    communities. Overall, metagenomic sequencing provided the estimates closest to the actual cell number

224    values, while the amplicon based quantification deviated the most from the actual numbers. The average x

225    fold deviations of the measured abundances from the expected abundance based on cell input were

226    significantly lower for metagenomics as compared to metaproteomics and amplicon sequencing (p-value

227    <0.01, Supplementary Table 4). The general overestimation of cell numbers by amplicon sequencing was

228    in part due to the fact that the amplicon sequencing failed to detect many of the species in the mock

229    communities driving up the relative abundances of the remaining ones.

230    Interestingly, the accuracy with which the three methods estimated the relative cell numbers in the mock

231    communities depended very much on the range of species abundances in them. All three methods

232    estimated the relative cell numbers quite well for the EQUAL CELL NUMBER community, but failed to

233    estimate them well for the EQUAL PROTEIN AMOUNT and UNEVEN communities, which represent a large

234    range of species abundances (Fig. 3a and 3g). This is likely due to the more inaccurate quantification of

235    low abundant strains/species that are close to the detection limit of the methods (see below and Fig. 4b).

**Low detection limit and high quantification accuracy with relatively little data**

237    To test the impact of the number of spectra acquired on the detection limit and dynamic range of species

238    proteinaceous biomass quantification, we ran five different LC-MS experimental setups for the four

239    biological replicates of the UNEVEN mock community (Fig. 4, Supplementary Table 5). These setups

240    provided varying numbers of $MS^2$ spectra for peptide identification. They included two basic 1D-LC-

241    MS/MS approaches of 260 min and 460 min run time. For each of these two approaches the amount of

242    data was doubled by running technical replicates. The fifth approach was a 2D-LC-MS/MS experiment in

9

243  which the sample was fractionated into 12 fractions using salt pulses on an SCX column followed by 120

244  min separations on a reverse phase column.

245  Each of the five approaches led to the detection of 27 out of the 30 distinguishable strains and species in

246  the community when the biological replicates were combined. We observed some small differences

247  between approaches in their detection sensitivity when looking at the data for individual biological

248  replicates. While we detected 25 to 26 species/strains (average 25.25) in the single 260 min runs, we

249  detected 26 to 27 (average 26.5) in the duplicate 460 min runs. From this follows that for the species

250  diversity and abundance distribution of the UNEVEN mock community a single 260 min (~130,000 $MS^2$

251  spectra) run provides a similar detection limit as compared to approaches that provide much more data

252  (e.g. 2x 460 min runs = ~390,000 $MS^2$ spectra). The detection limit for all five approaches was similar

253  and, interestingly, differed by organism group. The Archaeon *N. viennensis*, the Eukaryote *Chl.*

254  *reinhardtii* and all Bacteria were detected with all five approaches. The Bacterium *N. europaeae* was

255  mixed into the UNEVEN community with the lowest protein abundance of 0.08%, which suggests that at

256  least for Bacteria the detection limit is below 0.08%. Three out of the five bacteriophages in the

257  community were not detected by any of the approaches (Supplementary Table 6) even though they were

258  mixed into the community at protein abundances higher than that of *N. europaeae*, between 0.08 to 0.15%.

259  This is surprising, because these phages consist of only a few dominant proteins (e.g. capsid proteins),

260  which should enhance their detectability. Currently we do not have a good explanation for this result.

261  Surprisingly, all approaches had a similar accuracy in terms of quantifying species abundances (Fig. 4a).

262  Our expectation was that an increased number of $MS^2$ spectra would increase the accuracy of the

263  abundance estimates. Our data suggests that with a 260 min run we already reached saturation in terms of

264  accuracy for the UNEVEN mock community type. Interestingly, all five approaches underestimated the

265  abundances of species/strains that are present in the samples in low amounts (Fig. 4a). If low-abundance

266  species (<0.5% in all approaches) are removed from the dataset resulting in 18 species remaining, then the

267  deviation of the measurement from the actual protein input amount becomes much smaller (Fig. 4b,

268  Supplementary Table 6). This suggests that, as with most other analytical methods, the accuracy of the

269  measurement is lower for quantities close to the detection limit and thus the proteinaceous biomass

270  estimates for low abundant species should be treated as less precise.

271  In summary, a single 260 min 1D-LC-MS/MS run on a QExactive Plus Mass Spectrometer provides

272  enough data to detect most species in a community that contains 30 distinguishable species and features a

273  range of proteinaceous biomass abundances of more than two orders of magnitude. The limit of detection

274  can be slightly lowered using longer peptide separations and by increasing the amount of data generated

275  per sample.

10

**276  Estimation of absolute biomass contribution is possible even with incomplete sequence**

**277  databases**

278   One potential drawback of metaproteomics based biomass quantification of species in a microbial

279   community is that proteomic protein identification relies on the availability of a protein sequence

280   database. Proteins can only be identified and quantified if protein sequences are present in the database

281   that have a high similarity to the actual proteins in a sample. Analogous to the primer bias based exclusion

282   or incorrect estimation of species abundances in 16S/18S rRNA gene amplicon sequencing [19, 21], the

283   incompleteness of the protein sequence database used for protein identification can lead to the exclusion

284   or incorrect estimation of species abundances. However, the metaproteomic data in theory allows

285   estimating how incomplete the sequence database used is based on the number of available mass spectra

286   and the known proportion of how many of these mass spectra lead to PSMs in a search with a mock

287   community for which all protein sequences are known. This should allow to correct the relative abundance

288   estimates to absolute estimates.

289   To test the influence of database incompleteness on quantification results and if the error in abundance

290   estimates resulting from it can be corrected for, we used two sequence databases of varying

291   incompleteness to quantify the species in the UNEVEN community. In the first incomplete database

292   (INCOMPLETE1) the protein sequences for *Pseudomonas denitrificans*, *Pseudomonas fluorescens* and

293   *Rhizobium leguminosarum* bv. viciae strain 3841 were removed leaving the sequences of the closely

294   related species/strains *Pseudomonas pseudoalcaligenes* and *Rhizobium leguminosarum* bv. viciae strain

295   VF39 in the database. In the second incomplete database (INCOMPLETE2) the remaining *Pseudomonas* and

296   *Rhizobium* sequences as well as the *Salmonella enterica* typhimurium LT2 sequences were removed.

297   As expected, the number of detected organisms dropped for the quantification with the incomplete

298   sequence databases (Fig. 5a). In the quantification with the INCOMPLETE1 database the number of PSMs

299   for the remaining *R. leg.* VF39 and *P. pseudoalcaligenes* increased and thus their relative abundance. This

300   increase in PSM number is due to the fact that in the absence of the protein sequences of the correct

301   species/strain some of the MS$^2$ spectra match to peptides from closely related species/strains. As expected,

302   for the very closely related *R. leguminosarum* strains a larger fraction of PSMs shifted from one strain to

303   the other as compared to the *Pseudomonas* species for which only a smaller fraction of PSMs shifted over.

304   The PSM number for most remaining organisms remained very similar across the database completeness

305   range with the exception of *E. coli*, which obtained a large number of additional PSMs from the closely

306   related *S. enterica* in the quantification based on the INCOMPLETE2 database. As expected, the drop in the

307   total number of PSMs led to an increase of relative organism abundance when more protein sequences

308   were removed from the database (Fig. 5a and 5b). We corrected these relative biomass estimates by

309   calculating the number of PSMs lost due to database incompleteness based on the known proportion of

11

310    $MS^2$ spectra to PSMs in the quantification with the complete database. The corrected relative abundance

311    estimates for the quantification with the INCOMPLETE2 database were in most cases very similar to the

312    quantification with the complete database (Fig. 5b). Therefore, the proteinaceous biomass abundances

313    adjusted for database incompleteness can be used as an approximation of absolute proteinaceous biomass

314    abundances.

315    **Case studies**

316    To demonstrate the power and application of the metaproteomics based methods for assessing species

317    biomass contributions in microbial communities, we applied the methods developed here to microbial

318    communities from two widely different environments.

319    For the first application example, we generated both metaproteomic data, as well as 16S rRNA gene

320    amplicon data from two phototrophic biomats from soda lakes in the Canadian Rocky Mountains (Fig.

321    6a). We summarized organism abundances at the phylum level. Even on this high taxonomic level major

322    differences between the lakes and the two approaches become apparent. While the 16S rRNA gene

323    amplicon data suggests that the lakes were rather similar in taxonomic composition on the phylum level,

324    the metaproteomes painted a very different picture. The metaproteomes indicate that the major

325    phototrophs between the lakes were different. Lake 1 was dominated by Cyanobacteria, whereas lake 2

326    was dominated by green algae. Additionally, we detected dsDNA viruses in lake 2, which despite the fact

327    that they contribute only a small amount of proteinaceous biomass could play an important ecological

328    role. Interestingly, some bacterial groups that made up a significant amount of the 16S rRNA gene

329    amplicons (e.g. Bacteroidetes/Chlorobi group) contributed only a minor amount based on the

330    metaproteomic data. Since, the cell lysis method used for both approaches was identical an extraction bias

331    is unlikely, suggesting that a primer bias may be responsible for the discrepancy.

332    For the second application example, we re-analyzed a recently published saliva metaproteome that

333    provided extensive insights into the diurnal and inter-individual variation of the oral microbiome [9]. Grassl

334    et al. provided two independent datasets in their study on the presence and abundance of specific taxa in

335    the oral microbiomes. The first dataset (Fig. 4 in the original publication) provides presence/absence

336    patterns of taxa based on unique peptide matches and cultivation results. The second dataset provides

337    quantification of taxa based on peptides identified by metaproteomics, however, without the specificity

338    increasing step of protein inference (Fig. 6 c in the original publication). Our results from the re-analysis

339    of the proteomic data corresponded well with the taxonomic presence and absence patterns inferred by

340    Grassl et al.. However, our metaproteomic quantification of the data showed very different abundance and

341    presence profiles for bacterial genera, as compared to the original metaproteomic analysis. We observed a

342    much larger inter-individual variation for organism abundances (Fig. 6b). Additionally, several genera that

12

343    were detected in the original analysis to be abundant in the samples were not detected at all or in much

344    lower abundances (e.g. *Enterococcus* and *Abiotropha*), while other genera were much higher in abundance

345    (e.g. *Veillonella*, *Actinomyce*s and *Rothia*) (Fig. 6b). Grassl et al. acknowledged in their study that the

346    quantification method they used could come "at the disadvantage that peptides shared by two genera could

347    lead to an overestimation of the taxon's abundance." Our analyses suggest that non-unique matching of

348    peptides between genera indeed led to the skew in the original quantification data. For example,

349    *Enterococcus* and *Abiotropha* share many peptides with *Streptococcus*, however, only streptococcal

350    proteins could be inferred confidently to be present in the samples. This demonstrates that using validated,

351    highly specific protein inference criteria for metaproteomic based species quantification is crucial and that

352    peptide identification without subsequent protein inference is not sufficient to achieve high enough

353    specificity for quantification.

354

## Discussion

The metaproteomics-based biomass assessment approach that we demonstrate here is not limited to a specific set of computational tools and parameters. To make the approach as broadly applicable as possible, we have chosen to use a route that should make it possible to transfer the approach to many other platforms. In this manuscript we highlight the most crucial considerations for developing concrete methods for metaproteomics-based biomass assessment (e.g. protein inference specificity) and supply a comprehensive dataset to transfer this approach to other computational or experimental platforms for proteomics. The provided pure-culture derived proteomes (PXD006118), for example, will allow investigators to determine parameters to achieve sufficient protein inference specificity, while the different mock community proteomes (PXD006118) will allow assessing parameters based on quantification accuracy and number of detected species.

As we demonstrate here, metaproteomics-based biomass assessment is a powerful approach that allows to accurately quantify the proteinaceous biomass of a large number of taxa in a community all at once. This approach complements existing high throughput approaches for determining community composition based on DNA sequencing, in that it provides a different, independent measure of community composition. Our case study on soda lake biomass nicely illustrates that sequencing-based methods and metaproteomics can provide very different pictures of a community. An added benefit of using metaproteomes for community composition analyses is that the proteomic information will also provide insights into which metabolic and physiological functions are expressed and play a major role in the community.

Recently, there has been a recurring interest in more quantitative methods for microbial ecology for the absolute quantification of community composition (e.g. cell counts per volume) [22]. Metaproteomics-based abundance estimates can be put into an absolute context by simple assays, for example, by measuring total protein content of a specified sample volume, wet weight or dry weight. The relative proteinaceous biomass abundances of community members can then be converted to absolute values after considering necessary corrections for database incompleteness (see Results).

There are several questions that go beyond the scope of this study that should be addressed in the future. First, is proteinaceous biomass an accurate representation of the total biomass of a species? We would argue that in many cases, proteinaceous biomass is a good estimate of total biomass. However, as always we expect exceptions, where proteinaceous biomass is not a good predictor of total biomass, which would for example be the case of microorganisms that store large amounts of carbon in form of polyhydroxyalkanoates or glycogen. Second, a likely much more difficult question to answer is, if and

14

387    under what circumstances proteinaceous biomass of a community member can be used as an

388    approximation of the biological activity of that community member?

## Online methods

### Assembly of mock communities

Cultures of 32 Archaea, Bacteria, Eukaryotes and Bacteriophages (Supplementary Table 1) were donated to us by very kind colleagues. Cells were washed using phosphate buffered saline pH 7.4 (Sigma-Aldrich) to remove the cultivation medium. Cell counts of washed cells were determined by microscopy using a Neubauer improved counting chamber. Cells were aliquoted and pelleted by centrifugation at 21,000 xg for 5 min to create cell aliquots with known cell number. Bacteriophages were purified by filtration and polyethylene glycol (PEG) precipitation as described in Kleiner et al. (2015) [23]. Phage titers were determined as particle forming units (PFUs) per ml using the soft-agar overlay method [24]. Liquid aliquots with known titer were made for all phages. Cell pellets and phage aliquots were stored at -80°C.

We quantified the protein content of cell and phage aliquots for each strain using duplicate aliquots. For this, 300-600 ul SDT-lysis buffer (4% (w/v) sodium dodecyl sulfate (SDS), 100 mM Tris-HCl pH 7.6) were added to each pellet according to pellet size. The pellets in SDT-lysis buffer were vortexed and transferred to lysing matrix tubes (Matrix A, MP Biomedicals, Santa Ana, CA, USA) and lysed using a Bead Ruptor 24 (Omni International, https://www.omni-inc.com/) at 6 m/s for 45 seconds. The samples were heated for 10 minutes to 95°C and then centrifuged for 10 minutes at 21,000g. Dilutions of each sample were prepared and sample protein amounts were quantified using the Pierce bicinchoninic acid (BCA) assay (Thermo Scientific Pierce).

We assembled three types of mock communities by resuspending the frozen cell pellets of each microorganism in 150 ul ultrapure water and then combining varying amounts of each organism. The composition of each mock community type is detailed in Supplementary Tables 1 to 3. Four biological replicates of each mock community type were made and each replicate was divided into 20 aliquots. The UNEVEN mock community was designed to cover a large range of species abundances both on the level of cell number and proteinaceous biomass to test for the dynamic range and detection limits of the quantification methods (Fig. 3a). The EQUAL PROTEIN AMOUNT and EQUAL CELL NUMBER mock communities contained either the same amount of protein for all community members with varying cell numbers or the same number of cells for all members with varying amounts of protein. Since the bacteriophages yield very little protein even if high particle numbers are used, we mixed them at a 10x lower ratio into the EQUAL PROTEIN AMOUNT community.

### Sampling of soda lake biomats

Benthic microbial mats were sampled from two soda lakes located on the Cariboo Plateau, British Columbia, in June 2014 for 16S rRNA gene amplicon sequencing and metaproteomics and in May 2015

16

421    for metagenomic sequencing. Lake1 herein refers to Goodenough Lake (51°19'47.64"N 121°38'28.90"W)

422    and Lake2 refers to Last Chance Lake (51°19'39.3" N 121°37'59.3"W). Collected microbial mats from

423    each lake were pooled and immediately placed on ice in the field and frozen at -80 °C within two days of

424    sampling for DNA extraction.

### DNA extraction

426    For the mock community samples DNA was extracted from one aliquot of each of the four biological

427    replicates of each community type using the FastDNA Spin Kit (MP Biomedicals, Santa Ana, CA, USA)

428    according to the manufacturer's protocol with small modifications. Following addition of CLS-TC to each

429    aliquot, samples were homogenized in lysing matrix tubes (MP Biomedicals FastDNA Spin Kit, tube A)

430    for 45 seconds at 6 m/s using a Bead Ruptor 24 (OMNI). In addition, the DNA elution steps was repeated

431    twice. DNA concentrations were measured using a NanoDrop 2000 spectrophotometer (Thermo

432    Scientific).

433    DNA was extracted from the 2014 and 2015 Lake1 and Lake2 samples using the FastDNA Extraction Kit

434    for Soil (MP Biomedicals) with 10 minute centrifugation times for the spin filter steps and an additional

435    purification using 5.5 M guanidine thiocyanate as described in Sharp et al. (2017) [18].

### 16S rRNA gene amplicon sequencing and analysis of mock communities and soda lake biomats

438    DNA from all mock community samples and the 2014 soda lake biomats from Lake1 and Lake2 was used

439    for 16S rRNA gene amplicon libraries preparation as described in Sharp et al. (2017) [18]. We used the S-D-

440    Bact-0341-a-S-17 (also known as b341, 5'-

441    TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGAGGCAGCAG-3') [25] and S-D-

442    Bact-0785-a-A-21 (also known as Bakt_805R, 5'-

443    GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGACTACHVGGGTATCTAATCC-3') [26]

444    primers with added Illumina overhang adapters for the amplification of the HV regions 3-4 resulting in

445    427 bp amplicons (excluding the primers). Based on the evaluation by Klindworth et al. (2013) [19] this

446    primer pair yields a large coverage of the domain Bacteria. Libraries were pooled and normalized for

447    sequencing on the Illumina MiSeq Sequencer (San Diego, CA) using the 2 x 300 bp MiSeq Reagent Kit

448    v3. The resulting amplicon sequences were analyzed with MetaAmp [18]. Operational taxonomic units

449    (OTUs) were identified with a threshold of 97 % sequence similarity.

### Metagenomic sequencing of mock communities

451    Shotgun metagenomic sequencing (2x75 bp) of 3 replicates of each mock community type was performed

452    using the Illumina NextSeq 500 sequencer. The NEBNext Ultra II DNA Library Prep Kit (New England

453    Biolabs) was used for library preparation. Ten to nineteen million paired-end reads were generated for

454    each sample. We confirmed the library content using PhyloFlash (https://github.com/HRGV/phyloFlash)

455    and the quality of the data using FastQC (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/). We

456    used BBsplit from the BBmap package (version 35.85, http://sourceforge.net/projects/bbmap/) to map raw

457    reads against the mock community reference genomes to quantify the read coverage for each organism.

458    The reference genomes that were used are listed in Supplementary Table 7. Read mapping statistics for

459    each reference genome were generated using BBsplit's default parameters and by setting the 'refstats'

460    parameter. Relative read abundances for each organism were normalized to their genome sizes.

461    No reference genomes were available for AK199 and *Chromobacterium violaceum* CV026 in public

462    databases. Therefore we generated genomes for these two strains from the metagenomes using an iterative

463    assembly and binning strategy. All read files were trimmed for quality and adapters using BBduk from the

464    BBmap package (http://sourceforge.net/projects/bbmap/). The trimmed reads for the UNEVEN samples

465    were concatenated and assembled with metaSPAdes (version 3.8.1) [27]. The assembly quality was checked

466    by running metaQUAST (version 4.1) with the mock community reference genome set [28]. Metawatt

467    (version 3.5.2) was then used to create bins for AK199 and *C. violaceum* CV026 using default settings [29].

468    The bins were checked with metaQUAST to ensure that none of the included contigs aligned with any of

469    the other reference genomes for the mock community. The trimmed reads from all samples were

470    concatenated and BBmap was used to retrieve reads mapping to the AK199 and *C. violaceum* bins.

471    SPAdes (version 3.8.1) was used to assemble the mapped reads for AK199 and *C. violaceum* [30]. The

472    assembly quality was checked with metaQUAST, QUAST, and CheckM [31]. The AK199 genome was of

473    sufficient quality after this assembly round. The *C. violoceum* assembly was further improved by two

474    more rounds of read mapping and assembly. The AK199 and *C. violaceum* genomes were annotated using

475    the RAST server [32] and annotated protein sequences were retrieved for the construction of the protein

476    identification database.

### Soda lake biomat metagenomes, sequencing, assembly and annotation

478    DNA (250 ng) from the 2015 soda lake biomats from Lake1 and Lake2 was randomly sheared to a

479    fragment size of approximately 300 bp using a S2 focused-ultrasonicator (Covaris, Woburn, MA). The

480    fragmented DNA was then converted into an Illumina compatible sequencing library using the NEBNext

481    Ultra DNA Library Prep Kit according to the vendor's standard protocol. This included a size selection

482    step with SPRIselect magnetic beads and PCR enrichment (8 cycles) with NEBNext Multiplex Oligos for

483    Illumina. The libraries were measured using qPCR and the Kapa Library Quant Kit for Illumina and then

484    pooled in equal amounts for sequencing. A 1.8 pM solution was then sequenced on an Illumina NextSeq

485    500 sequencer using a 300 cycle (2x150 bp) high-output sequencing kit as per the Illumina protocol in the

486    Center for Health Genomics and Informatics in the Cumming School of Medicine, University of Calgary.

18

487  All raw Illumina reads were passed through an in-house Illumina read quality control program that filters

488  out known Illumina sequencing and library preparation artifacts. Specifically, all reads were removed that:

489  (i) matched the spike-in PhiX sequence; (ii) were shorter than 30 bp after clipping off the partial primer,

490  adapters, and the low-quality ranging at the ends; or (iii) were of low complexity. Reads that passed the

491  quality control stage were assembled into contigs using MEGAHIT v1.0.3 with options "--k-list

492  51,77,99,127 --min-count 2 –min-contig-len 500" [33]. The assembled contigs were merged into scaffolds

493  based on paired-end information using the SOAP v2.04 package [34]. The GapCloser v1.12 package was

494  applied to further close the gaps between contigs in scaffolds. All the scaffolds longer than 500 bp after

495  GapCloser post-processing were run through Prodigal v2.6.1 to identify coding sequences [35]. The coding

496  sequences (>= 60 aa) were annotated using DIAMOND [36] with options "-k 1 --seg no" to search against a

497  protein sequence reference database generated by GenomeDatabase

498  (https://sourceforge.net/projects/genomedatabase/) and the eggNOG database [37]. Best-hit matches

499  were filtered by query coverage >= 70% and percent identity >= 30%. Taxonomic assignments for protein

500  sequences were made on the basis of the filtered best-hit matches. The taxonomically annotated protein

501  sequences were then used to generate the protein identification database, by combining them with protein

502  sequences from several eukaryotic genomes and transcriptomes, which were chosen based on the results

503  from a 18S rRNA amplicon library. CD-HIT was used to remove redundant sequences from the database

504  using an identity threshold of 95% [38]. The cRAP protein sequence database (http://www.thegpm.org/crap/)

505  containing protein sequences of common laboratory contaminants was appended to the database. The final

506  database contained 4,171,024 protein sequences and is available from the PRIDE repository

507  (PXD006343).

508  **Protein extraction, quantification and peptide preparation**

509  Samples were lysed in SDT-lysis buffer with 0.1 M DTT. SDT-lysis buffer was added in a 1:10

510  sample/buffer ratio to the sample pellets. Cells were disrupted in lysing matrix tubes A (MP Biomedicals)

511  for 45 seconds at 6 m/s using the OMNI Bead Ruptor 24 and subsequently incubated at 95° C for 10

512  minutes followed by pelleting of debris for 5 min at 21,000 x g. We prepared tryptic digests following the

513  filter-aided sample preparation (FASP) protocol described by Wisniewski et al. (2009) [39]. In brief, 30 µl of

514  the cleared lysate were mixed with 200 µl of UA solution (8 M urea in 0.1 M Tris/HCl pH 8.5) in a 10

515  kDa MWCO 500 µl centrifugal filter unit (VWR International) and centrifuged at 14,000 x g for 40 min.

516  200 µl of UA solution were added again and centrifugal filter spun at 14,000 x g for 40 min. 100 µl of

517  IAA solution (0.05 M iodoacetamide in UA solution) were added to the filter and incubated at 22° C for

518  20 min. The IAA solution was removed by centrifugation and the filter was washed three times by adding

519  100 µl of UA solution and then centrifuging. The buffer on the filter was then changed to ABC (50 mM

520  Ammonium Bicarbonate), by washing the filter three times with 100 µl of ABC. 1 to 2 µg of MS grade

521   trypsin (Thermo Scientific Pierce, Rockford, IL, USA) in 40 µl of ABC was added to the filter and the

522   filters were incubated overnight in a wet chamber at 37° C. The next day, peptides were eluted by

523   centrifugation at 14,000 x g for 20 min, followed by addition of 50 µl of 0.5 M NaCl and further

524   centrifugation. Peptides were desalted using Sep-Pak C18 Plus Light Cartridges (Waters, Milford, MA,

525   USA) or C18 spin columns (Thermo Scientific Pierce, Rockford, IL, USA) according to the

526   manufacturer's instructions. Approximate peptide concentrations were determined using the Pierce Micro

527   BCA assay (Thermo Scientific Pierce, Rockford, IL, USA) following the manufacturer's instructions.

### 1D-LC-MS/MS and 2D-LC-MS/MS

529   The four biological replicates of each mock community type were analyzed using a block-randomized

530   design as outlined by Oberg and Vitek (2009) [40] using several LC-MS/MS methods. Two wash runs with

531   100% eluent B (80% acetonitrile, 0.1% formic acid) and one blank run were done between samples to

532   reduce carry over. For the 1D-LC-MS/MS mock community runs, 2 µg of peptide were loaded onto a 5

533   mm, 300 µm ID C18 Acclaim® PepMap100 pre-column (Thermo Fisher Scientific) using an UltiMate[TM]

534   3000 RSLCnano Liquid Chromatograph (Thermo Fisher Scientific) with loading solvent A (2%

535   acetonitrile, 0.05% TFA), eluent A (0.1% formic acid in water) and eluent B. After loading, the pre-

536   column was switched in line with a 50 cm x 75 µm analytical EASY-Spray column packed with PepMap

537   RSLC C18, 2µm material (Thermo Fisher Scientific), which was heated to 45° C. The analytical column

538   was connected via an Easy-Spray source to a Q Exactive Plus hybrid quadrupole-Orbitrap mass

539   spectrometer (Thermo Fisher Scientific). Peptides were separated on the analytical column at a flow rate

540   of 225 nl/min and mass spectra acquired in the Orbitrap as described by Petersen et al. (2016) [41]. A 260

541   min (from 2% B to 31% B in 200 min, in 40 min up to 50% B, 20 min at 99% B) and a 460 min gradient

542   (from 2% B to 31% B in 363 min, in 70 min up to 50% B, 27 min at 99% B) were used for 1D-LC. For the

543   2D-LC-MS/MS runs, 11 µg of peptide were loaded onto a 10 cm, 300 µm ID Poros 10 S SCX column

544   (Thermo Fisher Scientific) using the UltiMate[TM] 3000 RSLCnano LC with loading solvent B (2%

545   acetonitrile, 0.5% formic acid). Peptides were eluted from the SCX column onto the C18 pre-column

546   using 20 µl injection of salt plugs from the autosampler with increasing concentrations (12 salt plugs, 0 to

547   2000 mM NaCl). After each salt plug injection the pre-column was switched in line with the 50 cm x 75

548   µm analytical EASY-Spray column and peptides separated using a 120 minute gradient (from 2% B to

549   31% B in 82 min, in 10 min up to 50% B, 9 min at 99% B, 19 min at 2% B). Data acquisition in the Q

550   Exactive Plus was done as described by Petersen et al. (2016) [41].

551   The two soda lake samples were analyzed in technical quadruplicates by 1D-LC-MS/MS (1x 260 min and

552   3x 460 min runs for each). Two blank runs were done between samples to reduce carry over. For each 260

553   min run ~1 µg of peptide and for each 460 min run 2-4 µg of peptide were loaded onto a 2 cm, 75 µm ID

554   C18 Acclaim® PepMap 100 pre-column (Thermo Fisher Scientific) using an EASY-nLC 1000 Liquid

20

555 Chromatograph (Thermo Fisher Scientific) with eluent A (0.2% formic acid, 5% acetonitrile) and eluent B

556 (0.2% formic acid in acetonitrile). The pre-column was connected to a 50 cm x 75 µm analytical EASY-

557 Spray column packed with PepMap RSLC C18, 2µm material (Thermo Fisher Scientific), which was

558 heated to 35° C via the integrated heating module. The analytical column was connected via an Easy-

559 Spray source to a Q Exactive Plus. Peptides were separated on the analytical column at a flow rate of 225

560 nl/min using either a 260 min (from 0% to 20% B in 200 min, in 40 min to 35% B, ending with 20 min at

561 100% B) or a 460 min gradient (from 0% to 20% B in 354 min, in 71 min to 35% B, ending with 35 min

562 at 100% B). Eluting peptides were ionized with electrospray ionization and analyzed in the Q Exactive

563 Plus as described by Petersen et al. (2016) [41].

### Protein identification and quantification

565 For protein identification of the mock community samples a database was created using all protein

566 sequences from the reference genomes of the organisms used in the mock communities (Supplementary

567 Table 7). The cRAP protein sequence database (http://www.thegpm.org/crap/) containing protein

568 sequences of common laboratory contaminants was appended to the database. The final database

569 contained 123,100 protein sequences and is available from the PRIDE repository (PXD006118). For

570 protein identification of the soda lake mats we used the database described above.  For protein

571 identification of the human saliva metaproteomes we used the same public databases as described in

572 Grassl et al. [9] as a starting point. Namely the protein sequences from the human oral microbiome database

573 [42] and the human reference protein sequences from Uniprot (UP000005640). CD-HIT was used to remove

574 redundant sequences from the database using an identity threshold of 95% [38]. The saliva metaproteome

575 database contained 914,388 protein sequences and is available from the PRIDE repository (PXD006366).

576 For peptide identification and protein inference the MS/MS spectra were searched against the databases

577 using the Sequest HT node in Proteome Discoverer version 2.0.0.802 (Thermo Fisher Scientific) or the

578 MaxQuant software version 1.5.5.1 [15].

### Data availability

580 The mass spectrometry metaproteomics data and protein sequence databases have been deposited to the

581 ProteomeXchange Consortium via the PRIDE [43] partner repository with the dataset identifier PXD006118

582 for the pure culture and mock community data, with dataset identifier PXD006343 for the soda lake

583 biomats, and with dataset identifier PXD006366 for the re-analyses of the saliva metaproteomes by Grassl

584 et al. [9]. **Public release of the PRIDE projects will be requested as soon as a citable pre-print is online,**

585 **so it might take a few days for these identifiers to show up in the database.** A detailed overview of the

586 pure culture and mock community metaproteomic data for method development can be found in

587 Supplementary Table 5.

588    The sequencing data for the mock community metagenomes and 16S rRNA gene amplicons is available

589    from the European Nucleotide Archive with study accession number PRJEB19901.

590    The 16S rRNA gene amplicon sequencing data has been submitted to the NCBI short read archive (SRA)

591    with the following accession numbers SRR5291562 (Lake1) and SRR5291553 (Lake2).

**Author contributions**

M. K., conceived study, obtained and created bacterial stocks for mock communities, mock community experiments and mass spectrometry, data analysis, wrote the paper with input from all co-authors; E. T., performed mock community experiments and data analysis, wrote parts of the methods and revised the manuscript; C. S., obtained and created bacterial stocks for mock communities, soda lake sampling and sequencing data generation, wrote parts of the methods and revised manuscript; D. L., proteomics support and experiments, DNA extractions and 16S rRNA gene amplicon library preparation; X. D., assembled and annotated soda lake metagenomes, developed MetaAmp software; C. L., 16S rRNA amplicon library preparation and sequencing on MiSeq; M. S., conceived study, revised manuscript.

## Acknowledgements

## References

615

616   1.    Biteen, J.S. et al. Tools for the Microbiome: Nano and Beyond. *ACS Nano* **10**, 6-37 (2016).
617   2.    Amann, R. & Fuchs, B.M. Single-cell identification in microbial communities by improved
618         fluorescence in situ hybridization techniques. *Nat. Rev. Microbiol.* **6**, 339-348 (2008).
619   3.    Zhou, J. et al. High-throughput metagenomic technologies for complex microbial community
620         analysis: open and closed formats. *mBio* **6**, e02288-02214 (2015).
621   4.    Milo, R. What is the total number of protein molecules per cell volume? A call to rethink some
622         published values. *BioEssays* **35**, 1050-1055 (2013).
623   5.    Daims, H. in Cold Spring Harbor Protocols, Vol. 4 pdb.prot5253 (2009).
624   6.    Hettich, R.L., Sharma, R., Chourey, K. & Giannone, R.J. Microbial metaproteomics: identifying
625         the repertoire of proteins that microorganisms use to compete and cooperate in complex
626         environmental communities. *Curr. Opin. Microbiol.* **15**, 373-380 (2012).
627   7.    Milo, R., Jorgensen, P., Moran, U., Weber, G. & Springer, M. BioNumbers—the database of key
628         numbers in molecular and cell biology. *Nucleic Acids Res.* **38**, D750-D753 (2010).
629   8.    Hamann, E. et al. Environmental Breviatea harbour mutualistic *Arcobacter* epibionts. *Nature* **534**,
630         254-258 (2016).
631   9.    Grassl, N. et al. Ultra-deep and quantitative saliva proteome reveals dynamics of the oral
632         microbiome. *Genome Medicine* **8**, 44 (2016).
633   10.   Heyer, R. et al. Proteotyping of biogas plant microbiomes separates biogas plants according to
634         process temperature and reactor type. *Biotechnology for Biofuels* **9**, 155 (2016).
635   11.   Nesvizhskii, A.I. & Aebersold, R. Interpretation of shotgun proteomic data: The protein inference
636         problem. *Mol. Cell. Proteomics* **4**, 1419-1440 (2005).
637   12.   Timmins-Schiffman, E. et al. Critical decisions in metaproteomics: achieving high confidence
638         protein annotations in a sea of unknowns. *ISME J.* **11**, 309-314 (2017).
639   13.   Denef, V.J., Shah, M.B., VerBerkmoes, N.C., Hettich, R.L. & Banfield, J.F. Implications of
640         strain- and species-level sequence divergence for community and isolate shotgun proteomic
641         analysis. *J. Proteome Res.* **6**, 3152-3161 (2007).
642   14.   Kleiner, M. et al. Metaproteomics of a gutless marine worm and its symbiotic microbial
643         community reveal unusual pathways for carbon and energy use. *Proc. Natl. Acad. Sci. USA* **109**,
644         E1173-E1182 (2012).
645   15.   Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass spectrometry-
646         based shotgun proteomics. *Nat. Protoc.* **11**, 2301-2319 (2016).
647   16.   Serang, O., MacCoss, M.J. & Noble, W.S. Efficient marginalization to compute protein posterior
648         probabilities from shotgun mass spectrometry data. *J. Proteome Res.* **9**, 5346-5357 (2010).
649   17.   Serang, O. The probabilistic convolution tree: efficient exact bayesian inference for faster LC-
650         MS/MS protein inference. *PLoS ONE* **9**, e91507 (2014).
651   18.   Sharp, C.E. et al. Robust, high-productivity phototrophic carbon capture at high pH and alkalinity
652         using natural microbial communities. *Biotechnology for Biofuels* **10**, 84 (2017).
653   19.   Klindworth, A. et al. Evaluation of general 16S ribosomal RNA gene PCR primers for classical
654         and next-generation sequencing-based diversity studies. *Nucleic Acids Res.* **41**, e1 (2013).
655   20.   Parada, A.E., Needham, D.M. & Fuhrman, J.A. Every base matters: assessing small subunit rRNA
656         primers for marine microbiomes with mock communities, time series and global field samples.
657         *Environ. Microbiol.* **18**, 1403-1414 (2016).
658   21.   Tremblay, J. et al. Primer and platform effects on 16S rRNA tag sequencing. *Frontiers in*
659         *Microbiology* **6**, 771 (2015).
660   22.   Props, R. et al. Absolute quantification of microbial taxon abundances. *ISME J.* (2016).
661   23.   Kleiner, M., Hooper, L.V. & Duerkop, B.A. Evaluation of methods to purify virus-like particles
662         for metagenomic sequencing of intestinal viromes. *BMC Genomics* **16**, 7 (2015).
663   24.   Sambrook, J., Fritsch, E.F. & Maniatis, T. in Molecular Cloning, Vol. 1, Edn. second. (ed. C.
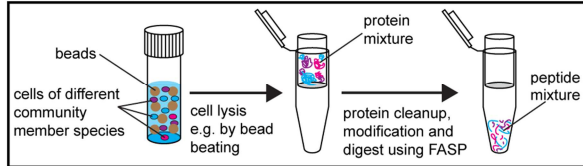664         Nolan) (Cold Spring Harbor Laboratory Press, 1989).

665   25.   Juck, D., Charles, T., Whyte, L.G. & Greer, C.W. Polyphasic microbial community analysis of
666         petroleum hydrocarbon-contaminated soils from two northern Canadian communities. *FEMS*
667         *Microbiol. Ecol.* **33**, 241-249 (2000).
668   26.   Herlemann, D.P.R. et al. Transitions in bacterial communities along the 2000 km salinity gradient
669         of the Baltic Sea. *ISME J.* **5**, 1571-1579 (2011).
670   27.   Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P.A. metaSPAdes: a new versatile
671         metagenomic assembler. *Genome Res.* **Published in Advance** (2017).
672   28.   Mikheenko, A., Saveliev, V. & Gurevich, A. MetaQUAST: evaluation of metagenome
673         assemblies. *Bioinformatics* **32**, 1088-1090 (2016).
674   29.   Strous, M., Kraft, B., Bisdorf, R. & Tegetmeyer, H. The binning of metagenomic contigs for
675         microbial physiology of mixed cultures. *Frontiers in Microbiology* **3** (2012).
676   30.   Bankevich, A. et al. SPAdes: a new genome assembly algorithm and its applications to single-cell
677         sequencing. *J. Comput. Biol.* **19**, 455-477 (2012).
678   31.   Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P. & Tyson, G.W. CheckM: assessing
679         the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome*
680         *Res.* **25**, 1043-1055 (2015).
681   32.   Aziz, R. et al. The RAST Server: Rapid Annotations using Subsystems Technology. *BMC*
682         *Genomics* **9**, 75 (2008).
683   33.   Li, D. et al. MEGAHIT v1.0: A fast and scalable metagenome assembler driven by advanced
684         methodologies and community practices. *Methods* **102**, 3-11 (2016).
685   34.   Luo, R. et al. SOAPdenovo2: an empirically improved memory-efficient short-read de novo
686         assembler. *Gigascience* **1**, 18 (2012).
687   35.   Hyatt, D., LoCascio, P.F., Hauser, L.J. & Uberbacher, E.C. Gene and translation initiation site
688         prediction in metagenomic sequences. *Bioinformatics* **28**, 2223-2230 (2012).
689   36.   Buchfink, B., Xie, C. & Huson, D.H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
690         *Meth.* **12**, 59-60 (2015).
691   37.   Huerta-Cepas, J. et al. eggNOG 4.5: a hierarchical orthology framework with improved functional
692         annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res.* **44**, D286-D293
693         (2016).
694   38.   Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or
695         nucleotide sequences. *Bioinformatics* **22**, 1658-1659 (2006).
696   39.   Wisniewski, J.R., Zougman, A., Nagaraj, N. & Mann, M. Universal sample preparation method
697         for proteome analysis. *Nat. Meth.* **6**, 359-362 (2009).
698   40.   Oberg, A.L. & Vitek, O. Statistical design of quantitative mass spectrometry-based proteomic
699         experiments. *J. Proteome Res.* **8**, 2144-2156 (2009).
700   41.   Petersen, J.M. et al. Chemosynthetic symbionts of marine invertebrate animals are capable of
701         nitrogen fixation. *Nature Microbiology* **2**, 16195 (2016).
702   42.   Chen, T. et al. The Human Oral Microbiome Database: a web accessible resource for investigating
703         oral microbe taxonomic and genomic information. *Database (Oxford)* **2010** (2010).
704   43.   Vizcaíno, J.A. et al. 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.*
705         **44**, D447-D456 (2016).
706   44.   Zaikova, E. et al. Microbial community dynamics in a seasonally anoxic fjord: Saanich Inlet,
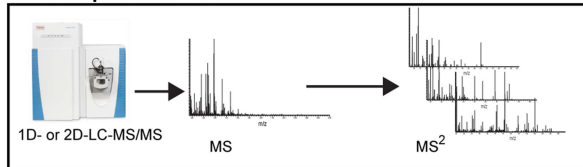707         British Columbia. *Environ. Microbiol.* **12**, 172-191 (2010).
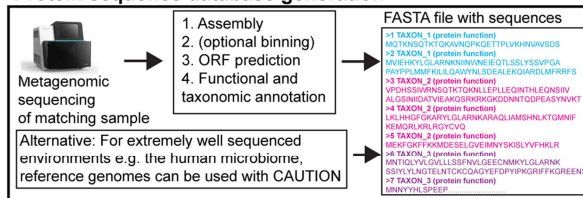708
709

26

710 **Figures**
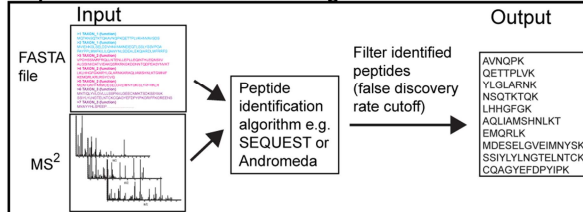
**Sample collection**

**Sample preparation**

beads

cells of different community member species

cell lysis e.g. by bead beating

protein mixture

protein cleanup, modification and digest using FASP

peptide mixture

**Data acquisition: LC-MS/MS**

1D- or 2D-LC-MS/MS

MS

MS$^2$

**Protein sequence database generation**

Metagenomic sequencing of matching sample

1. Assembly
2. (optional binning)
3. ORF prediction
4. Functional and taxonomic annotation

Alternative: For extremely well sequenced environments e.g. the human microbiome, reference genomes can be used with CAUTION

FASTA file with sequences

>1 TAXON_1 (protein function)
MQTKNSQTKTQKAVNQPKQETTPLVKHNVAVSDS
>2 TAXON_1 (protein function)
MVIEHKYLGLARNKNIINVNEIEQTLSSLYSSVPGA
PAYPPLMMFKILILQAWYNLSDEALEKQIARDLMFRRFS
>3 TAXON_2 (protein function)
VPDHSSIWRNSQTKTQKNLLEPLLEQINTHLEQNSIIV
ALGSINIIDATVIEAKQSRKRKGKDDNNTQDPEASYNVKT
>4 TAXON_2 (protein function)
LKLHHGFGKARYLGLARNKARAQLIAMSHNLKTGMNIF
KEMORLKRLRGYCVQ
>5 TAXON_2 (protein function)
MEKFGKFFKKMDESELGVEIMNYSKISLYVFHKLR
>6 TAXON_3 (protein function)
MNTIQLYVLGVLLLSSFNVLGEECNMKYLGLARNK
SSIYLYLNGTELNTCKCQAGYEFDPYIPKGRIFFKGREENS
>7 TAXON_3 (protein function)
MNNYYHLSPEEP

**Peptide identification and filtering**

Input

FASTA file

MS$^2$

Peptide identification algorithm e.g. SEQUEST or Andromeda

Filter identified peptides (false discovery rate cutoff)

Output

AVNQPK
QETTPLVK
YLGLARNK
NSQTKTQK
LHHGFGK
AQLIAMSHNLKT
EMQRLK
MDESELGVEIMNYSK
SSIYLYLNGTELNTCK
CQAGYEFDPYIPK

**Protein inference for specificity and sensitivity**

Peptides matched to protein sequences

>1 TAXON_1 (protein function)
MQTK NSQTKTQK AVNQPK QETTPLVK HNVAVSDS
>2 TAXON_1 (protein function)
MVIEHK YLGLARNK NIINWNEIEQTLSS LYSSVPGAPAYPPLMMFKILILQAWYNL SDEALEKQIARDLMFRRFS
>3 TAXON_2 (protein function)
VPDHSSIWR NSQTKTQK NLLEPLLEQI NTHLEQNSIIVALGSINIIDATVIEAKQSR KRKGKDDNNTQDPEASYNVKT
>4 TAXON_2 (protein function)
LK LHHGFGK AR YLGLARNK AR AQLIAMSHNLKT GMNIFK EMQRLK RLR GYCVQ
>5 TAXON_2 (protein function)
MEKFGKFFKK MDESELGVEIMNYSK ISLY VFHKLR
>6 TAXON_3 (protein function)
MNTIQLYVLGVLLLSSFNVLGEECNMK YLG LARNK SSIYLYLNGTELNTCK CQAGYEFD PYIPK GRIFFKGREENS

**Protein inference:** With simple filtering for false discovery rate and unique peptides

or

Using specialized inference algorithm e.g. Fido

Protein unique peptide

Non-unique peptide

Identified proteins

>1 TAXON_1 (protein function)
MQTK NSQTKTQK AVNQPK QETTPLVK HNVAVSDS
>4 TAXON_2 (protein function)
LK LHHGFGK AR YLGLARNK AR AQLIAMSHNLKT GMNIFK EMQRLK RLR GYCVQ
>5 TAXON_2 (protein function)
MEKFGKFFKK MDESELGVEIMNYSK ISLY VFHKLR
>6 TAXON_3 (protein function)
MNTIQLYVLGVLLLSSFNVLGEECNMK YLG LARNK SSIYLYLNGTELNTCK CQAGYEFD PYIPK GRIFFKGREENS
...
...
...

**Quantification**

| Protein Accession | Description | Spectrum count | Peptide intensity |
|---|---|---|---|
| 1 | TAXON_1 … | 133 | 164772.37 |
| 21 | TAXON_1 … | 65 | 80527.85 |
| 309 | TAXON_1 … | 34 | 42122.26 |
| 4 | TAXON_2 … | 25 | 30972.25 |
| 5 | TAXON_2 … | 24 | 29733.36 |
| 411 | TAXON_2 … | 22 | 27255.58 |
| 2119 | TAXON_2 … | 23 | 28494.47 |
| 6 | TAXON_3 … | 54 | 66900.06 |
| 3334 | TAXON_3 … | 35 | 43361.15 |
| 236 | TAXON_3 … | 44 | 54511.16 |
| 123 | TAXON_3 … | 111 | 137516.79 |

**Sum** counts or intensities by taxon

100%
80%
60%
40%
20%
0%

TAXON_3
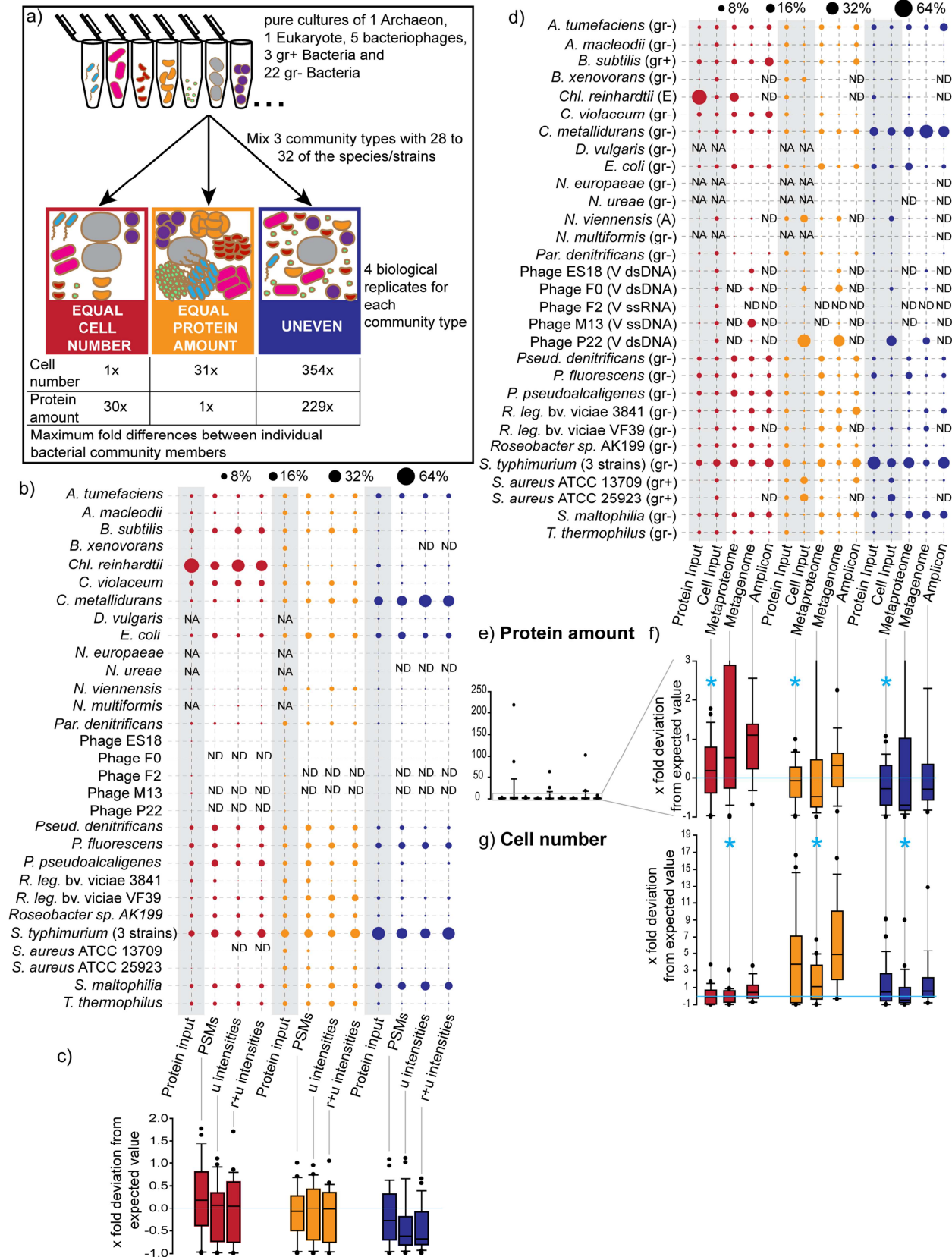TAXON_2
TAXON_1

711

712     **Figure 1: Workflow for assessing species biomass contributions using metaproteomics.**

29

**Figure 2: Specificity and sensitivity of protein identification using different protein inference strategies.** Proteins from different pure cultures were identified using a simulated metagenomic database that contained the protein sequences of 30 species and strains. For each pure culture, mass spectrometric data was acquired with a 260 minute long 1D-LC-MS/MS run (SEE METHODS). The database contained the sequences of the pure culture organism and sequences of organisms of various degrees of taxonomic relatedness. For the alphaproteobacterium *Rhizobium leguminosarum* bv. viciae 3841 the closest relative in the database was a highly similar strain of the same species (a), for the gammaproteobacterium *Pseudomonas denitrificans* the closest relatives in the database were other *Pseudomonas* species (b), for the betaproteobacerium *Nitrosospira multiformis* the closest relatives in the database were from the related betaproteobacterial genus *Nitrosomonas* (c), and for the archaeum *Nitrososphaera viennensis* the closest relatives in the database were bacteria (d). The same selection of protein inference strategies is shown for all cultures: **SQ 5% FDR**; SEQUEST search filtered based on search engine scores for 5% false discovery rate (FDR) using standard Target-Decoy strategy (implemented as Protein FDR Validator Node in Protein Discoverer). **2 UP ⊂ SQ 5% FDR;** same as previous, but only the subset of proteins identified with at least two unique peptides (UP) was considered. **SQ Fido;** Fido results filtered at 5% FDR based on protein q-value. **2 PU ⊂ SQ Fido;** same as previous, but only the subset of proteins identified with at least two protein unique (PU) peptides was considered. **SQ Fido ∩ MQ;** Only proteins considered that were identified both by Sequest-Fido (FDR 5%) and MaxQuant (1% protein FDR, at least 2 razor+unique peptides). **(2 PU ⊂ SQ Fido) ∩ MQ;** Same as previous, but for Sequest-Fido only the subset of proteins

733 identified with at least two protein unique (PU) peptides was considered. **SQ Fido ∩ (MQ || 3 PU);** Only

734 proteins considered that were identified both by Sequest-Fido and MaxQuant. Additionally, Sequest-Fido

735 identified proteins were retained even if they were not identified by MaxQuant if they had at least three

736 protein unique peptides

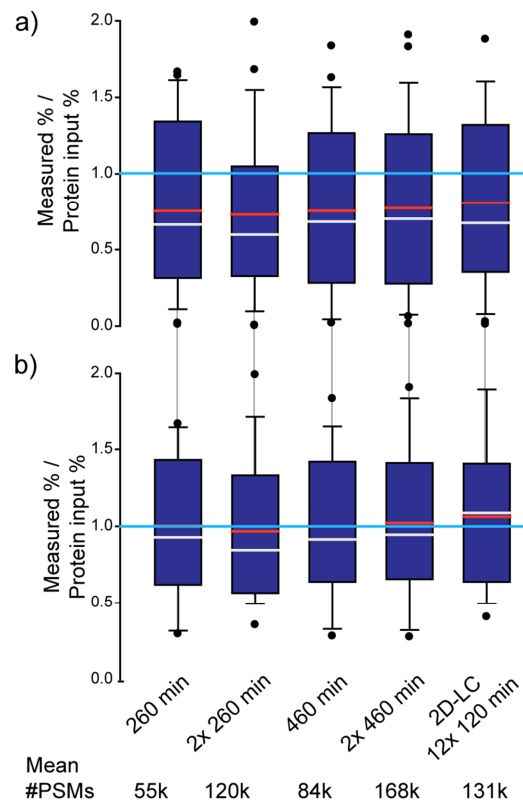738  **Figure 3: Quantification of mock community using three proteomic quantification methods and two**

739  **DNA sequencing based methods;** a) Illustration of mock community construction. 32 species and strains

740  were used for the construction of three distinct community types. b) comparison of three proteomic

741  quantification methods with the relative protein input amounts. Averages of four biological replicates per

742  community are shown (full data in Supplementary Table 8). Data from two 460 min long 1D-LC-MS/MS

743  runs were used per biological replicate. Three different quantification methods were used including sum of

744  peptide-spectrum matches (PSMs), sum of peptide ion intensities using only unique peptides (u

745  intensities), and sum of peptide ion intensities using razor and unique peptides (r+u intensities) as

746  implemented in MaxQuant. The bacteriophages were mixed at a 1:10 ratio into the "equal protein amount"

747  communities. Based on metagenomic sequencing we found that the *B. xenovorans* culture was

748  contaminated with *S. epidermidis* and thus the input protein amounts and cell numbers for *B. xenovorans*

749  were lower than calculated. We used three *Salmonella enterica* typhimurium strains in the mock

750  communities that differed only in a small number of genes and thus were de facto indistinguishable on the

751  proteomic and metagenomic level and thus the inputs for the three strains are reported as one species. The

752  bubble plot was generated with the bubble.pl script [44]. c) Box plots show the x fold deviation of the

753  amounts measured with the three proteomic quantification methods from the actual protein input amounts.

754  The box indicates the $1^{st}$ and $3^{rd}$ quartile, the line indicates the median and the whiskers indicate the $10^{th}$

755  and $90^{th}$ percentile. Outliers are indicated as individual points. If measurement and input were equal then

756  all values would be exactly 0 (indicated by bright blue line). Zeros (species that were not detected i.e.

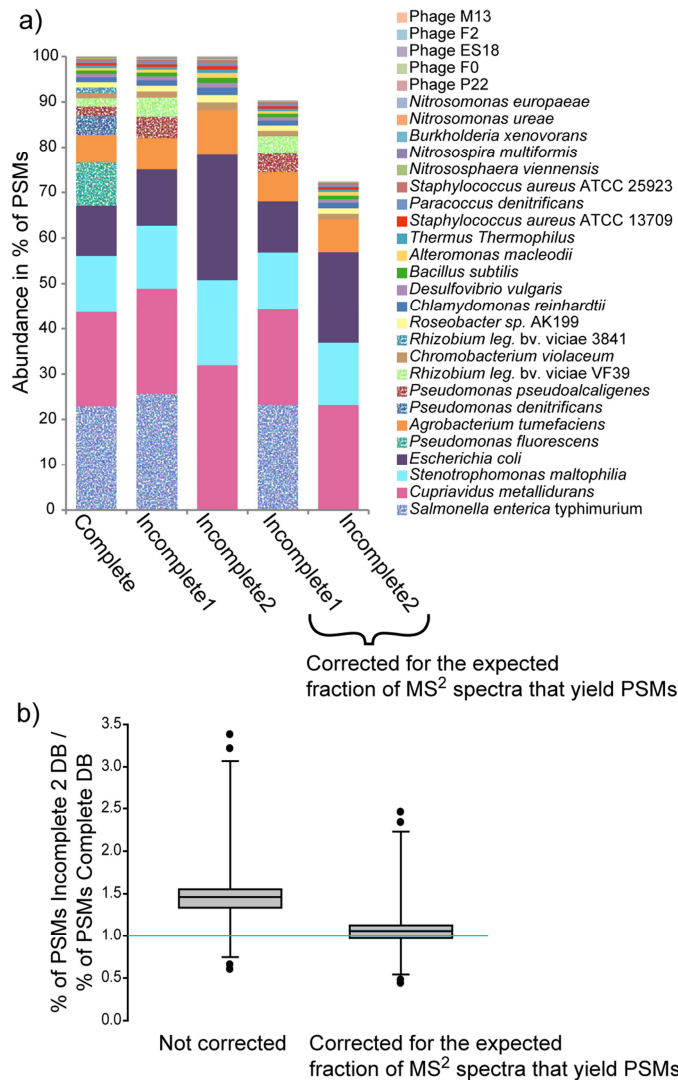757  'ND' in panel b) were removed before plotting.

758  d) Comparison of metaproteomic, shotgun metagenomic and 16S rRNA gene amplicon based

759  quantification of the mock communities with the input protein amounts and cell numbers. For 16S rRNA

760  gene amplicons and metaproteomes averages of four biological replicates per community type are shown.

761  For the shotgun metagenomic data averages of three biological replicates are shown. For the

762  metaproteomes the PSM based quantification is shown (see panel b). (gr- or gr+) gram positive or

763  negative Bacterium, (A) Archaeum, (E) Eukaryote, (V dsDNA, ssDNA or ssRNA) virus specifying

764  nucleic acid type of genome.

765  e-g) Box plots show the x fold deviation of the species abundance quantification with metaproteomics,

766  metagenomics and 16S rRNA gene amplicon sequencing from the actual input amounts for protein e) and

767  f) and cell number g). f) is an enlargement of the lower part of e). If measured and input species

768  abundance were equal, then all values would be exactly 0 (indicated by bright blue line). Zeros (species

769  that were not detected i.e. 'ND' in panel d) were removed before plotting. For each community type and

770  method the method with the significantly lowest x fold deviation (p-value <0.01) is indicated with a bright

771  blue '*' (see Supplementary Table 4 for details on statistics).

33

772    NA: species not added to this mock community; ND: Not detected with this method.
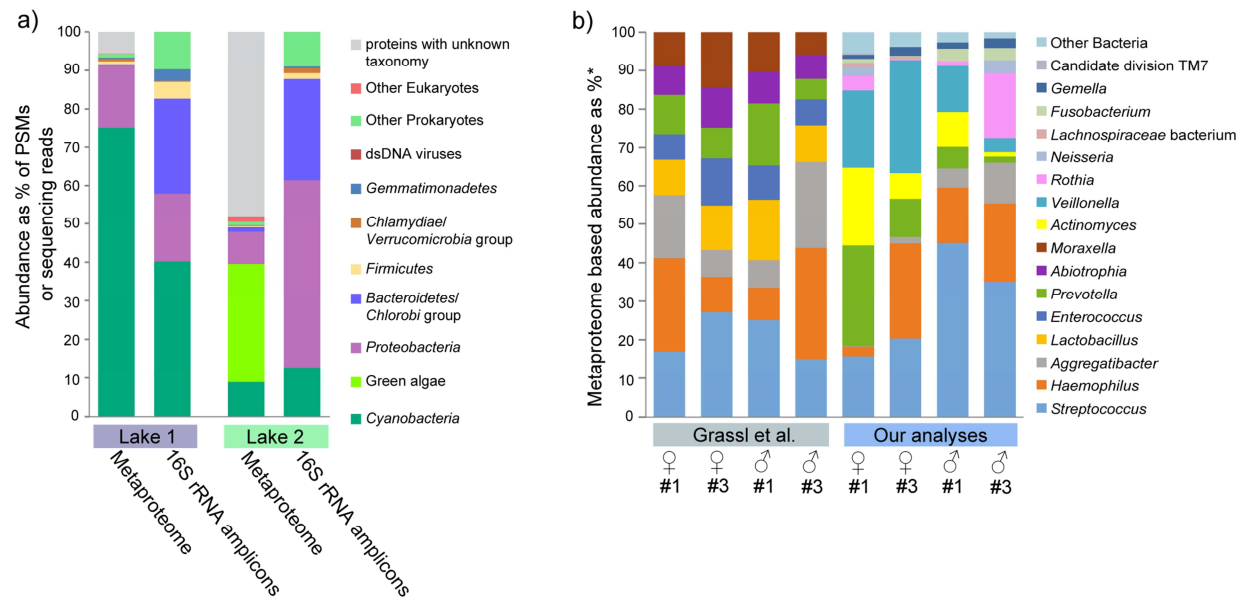
**Figure 4: Comparison of quantification accuracy for the UNEVEN mock community depending on LC method, gradient length and technical replication.** Four biological replicates were run for each method and the number of peptide-spectrum matches (PSMs) was averaged per organism. The numbers below the plots give the average total number of PSMs generated for each of the methods. Box plots show the deviation of the amounts measured by summing of PSMs from the actual protein input amounts. The box indicates the 1st and 3rd quartile, the grey line indicates the median, the red line the average and the whiskers indicate the 10th and 90th percentile. If measurement and input were equal, then all values would be exactly 1 (indicated by bright blue line). The first four methods were 1D-LC-MS/MS runs of the given length. The '2x' indicates if technical replicates were run. The fifth method was 2D-LC-MS/MS runs with each of the 12 fractions measured for 120 minutes (detailed data for this Figure is in Supplementary Table 6). In a) the deviation values for all 27 detected strains and species are shown. In b) the deviation values are only shown for the 18 strains and species that had an abundance >0.5% based on at least one of the methods.

35

787

**Figure 5: Effect of using incomplete protein sequence databases on quantification results.** Species in the UNEVEN community were quantified using the complete protein sequence database containing the protein sequences of all species in the community for protein identification and two sequence databases of varying incompleteness. In the first incomplete database (INCOMPLETE1) the protein sequences for *Pseudomonas denitrificans*, *Pseudomonas fluorescens* and *Rhizobium leguminosarum* bv. viciae (strain 3841) were removed leaving the sequences of the closely related species/strains *Pseudomonas pseudoalcaligenes* and *Rhizobium leguminosarum* bv. viciae (strain VF39) in the database. In the second incomplete database (INCOMPLETE2) the remaining *Pseudomonas* and *Rhizobium* sequences as well as the *Salmonella enterica* typhimurium LT2 sequences were removed. In a) the average quantification results for the four UNEVEN community biological replicates are shown. For the 4th and 5th bar the quantification results were corrected by considering the percentage of PSMs lost due to database incompleteness, i.e. based on searches with the complete database the expected fraction of MS$^2$ spectra that yielded PSMs was

800     known and thus we could calculate the difference between the expected number of PSMs and actual

801     number of PSMs in the quantification with the incomplete databases. In b) the comparison of the

802     quantification results with the complete and the INCOMPLETE2 databases are shown before and after

803     correction of the INCOMPLETE2 quantification results. If the quantification results were in perfect

804     agreement then all values would be 1 (indicated by the bright blue line). The box indicates the 1$^{st}$ and 3$^{rd}$

805     quartile, the line indicates the median and the whiskers indicate the 10$^{th}$ and 90$^{th}$ percentile.

806

**Figure 6: Two application examples using metaproteomics of estimating species biomass in communities.** a) Comparison of phylum level quantification of two soda lake biomats using metaproteomics and 16S rRNA gene amplicon sequencing. b) Re-analysis of a published saliva metaproteome dataset by Grassl et al. (2016) [9] using our method and comparison with the original analysis. *For the Grassl et al. analyses the abundances are given in % of the summed MS intensities of the 10 most abundant peptides per genus across all samples. For our analyses the abundances are given as % of all PSMs from proteins inferred by FidoCT with an FDR of 5% and at least 2 protein unique peptides.