This is the author's final version of the contribution published as:

Bonk, F., Popp, D., Harms, H., Centler, F. (2018): PCR-based quantification of taxa-specific abundances in microbial communities: Quantifying and avoiding common pitfalls *J. Microbiol. Methods* **153**, 139 – 147

The publisher's version is available at:

http://dx.doi.org/10.1016/j.mimet.2018.09.015

1	PCR-based quantification of taxa-specific abundances in microbial communities:
2	Quantifying and avoiding common pitfalls
3	
4	Fabian Bonk [†] , Denny Popp ^{†,*} , Hauke Harms, Florian Centler
5	Department of Environmental Microbiology, UFZ - Helmholtz Centre for Environmental
6	Research, Permoserstraße 15, 04318 Leipzig, Germany
7	
8	[†] Authors contributed equally
9	
10	* Corresponding author: Denny Popp, Permoserstraße 15, 04318 Leipzig, Germany. Tel: +49
11	341 235 482247; e-mail address: denny.popp@ufz.de
12	
13	
14	
15	
16	
17	
18	
19	
20	

21 Abstract

22 The quantification of relative and absolute taxa-specific abundances in complex microbial 23 communities is crucial for understanding and modeling natural and engineered ecosystems. 24 Many errors inherent to this quantification are, though well-known, still insufficiently 25 addressed and can potentially lead to a completely different interpretation of experimental 26 results. This review provides a critical assessment of next generation sequencing (NGS) of 27 amplicons and quantitative real-time PCR for the quantification of relative and absolute taxa-28 specific genome abundances. Starting from DNA extraction, the following error sources were 29 considered: DNA extraction efficiency, PCR-associated bias, variance of strain-specific 16S 30 rRNA operon copy number per genome, and analysis of quantitative real-time PCR and NGS 31 data. Tools and methods for estimating and minimizing these errors are presented and 32 demonstrated on published data. In conclusion, amplicon sequencing and qPCR of 16S rRNA 33 genes are valuable tools to determine relative and absolute taxa-specific genome abundances, 34 but results can deviate by several orders of magnitudes from the true values if the reviewed 35 error sources are ignored. Many of these errors can be minimized in a cost-efficient manner 36 and large errors can be easily identified by plausibility checks as shown in this review. 37 Finally, the accurate conversion of genome abundances to cell numbers and microbial biomasses was pointed out as an important future research topic for the integration of PCR-38 39 based abundances into mathematical models.

40 Keywords: absolute abundance quantification; biomass quantification; 16S rRNA gene
41 amplicon sequencing; quantitative real-time PCR (qPCR); 16S rRNA gene copy number
42 variation; ploidy

43

44 **1. Introduction**

45 Microbial communities are the hidden champions in many natural and engineered ecosystems, 46 driving global elemental cycles, waste removal in waste water treatment plants, and methane 47 production in biogas plants to name a few examples. To understand and model these systems, 48 the accurate quantification of taxa-specific abundances is of crucial importance. Abundances 49 can hereby refer to gene abundances, genome abundances, cell numbers and biomasses. In 50 microbial ecology studies, abundances are often based on 16S rRNA gene copy numbers 51 while mechanistic mathematical models rather consider microbial abundances as biomass 52 measured as dry weight.

53 Various culture-independent molecular biology based techniques have been used to study 54 microbial communities such as metagenomics (Eloe-Fadrosh et al., 2016; Shakya et al., 55 2013), metaproteomics (Kleiner et al., 2017), flow cytometry (Lambrecht et al., 2017), fluorescence in-situ hybridization (FISH) (Nettmann et al., 2010), and amplicon sequencing 56 57 (Klassen et al., 2017). For amplicon sequencing, 16S rRNA genes are commonly used to 58 identify community composition, but also other genes have been targeted, for example the 59 mcrA gene to focus on methanogenic archaea (Steinberg and Regan, 2008). 16S rRNA gene 60 amplicon sequencing informs on relative taxa-specific gene abundances, which can be 61 converted to relative taxa-specific genome abundances using strain-specific 16S rRNA operon 62 copy number information. Relative gene and genome abundances are commonly used to 63 analyze the relationship between microbial community composition and environmental 64 parameters.

However, relative abundances are of limited use if total abundances are unknown (Props et al., 2017) (see Additional file 1, chapter A.1). Absolute taxa-specific genome abundances can be determined by quantitative real-time PCR (qPCR) for individual taxa (Yu et al., 2005) and for all taxa of a community at once by combining 16S rRNA gene amplicon sequencing with qPCR (Dannemiller et al., 2014). Both amplicon sequencing and qPCR can lead to substantial

errors of which many have been widely discussed in the literature. However, their impact is
still often neglected and not systematically corrected for when reporting on relative and
absolute genome abundances.

73 To guide practitioners, we critically evaluate potential errors of 16S rRNA gene amplicon 74 sequencing and qPCR in quantitative terms in this review. Moreover, tools and guidelines to 75 estimate and minimize these errors are provided. Errors discussed in this review are 76 associated with DNA extraction, PCR, and analysis of NGS and qPCR data. Furthermore, 77 methods are presented to estimate absolute abundances to check the plausibility of 78 experimental results. Finally, we addressed the errors associated with the conversion of 79 absolute taxa-specific genome abundances to taxa-specific cell numbers and biomasses to be 80 used in mathematical modelling. To illustrate the errors above, we mainly refer to examples 81 from anaerobic digestion which is driven by a complex microbial community. However, the 82 problems and solutions discussed in this study directly extend to any complex prokaryotic 83 community.

84

85 2. Errors associated with both amplicon sequencing and qPCR

86 2.1. Avoidance, quantification, and correction of DNA extraction efficiency associated
87 errors

Extracting DNA of a diverse community from a complex matrix is a challenging task. Lysis conditions have to be harsh enough to break up all types of cells but should not cause damage to the DNA. The strain-specific DNA extraction efficiency, i.e. the recovered amount of genomic DNA divided by the total amount of genomic DNA present in the sample for a certain strain, depends on the sample matrix, species' morphology, and the extraction method.

93 Inter-strain differences in DNA extraction efficiency, for example due to different cell walls 94 and membranes, lead to biased relative and absolute genome abundances. This error can be 95 minimized by spiking a microbial mock community of known and representative composition 96 to a sample and testing various DNA extraction methods (Willner et al., 2012). The mock 97 community should contain a variety of morphologies representative of the sample and E. coli 98 (if used to determine overall extraction efficiency, see below). The extraction method leading 99 to the least biased relative taxa-specific abundances of the mock community members should 100 be chosen.

101 After minimizing taxa-specific extraction biases, the overall extraction efficiency of 102 prokaryotic DNA needs to be determined. If DNA loss during extraction is neglected, 103 absolute genome abundances in the sample will be underestimated. Proposed standards for 104 efficiency estimation commonly rely on spiking a known number of E. coli cells to the 105 sample. A target gene in the genomic DNA or on a plasmid of E. coli is then quantified. The 106 overall extraction efficiency is calculated by the number of detected target genes divided by 107 the number of spiked target genes. If E. coli is already a member of the community to be 108 analyzed, its originally present target genes in the sample needs to be subtracted before 109 calculating the overall extraction efficiency. Using E. coli spikes to cattle manure, overall 110 extraction efficiencies were between 38% and 99.97% for different extraction methods 111 (Lebuhn et al., 2016). The extraction efficiency of 38% would result in an error on the 112 absolute genome abundances of 263%.

Previous to DNA extraction, an additional step for removing extracellular DNA from dead cells can be performed using for example propidium monoazide (Emerson et al., 2017). While leading to more accurate results, the application of propidium monoazide to turbid samples remains a challenge in practice (Kirkegaard et al., 2017).

In conclusion, inter-species variations in DNA extraction efficiency can lead to substantially inaccurate relative genome abundances. This error cannot be corrected a posteriori. A mock community can be used to choose the best DNA extraction method for minimizing this error. Thereafter, the overall DNA extraction efficiency needs to be estimated by spiking a known amount of a standard, typically *E. coli*, to the sample prior to DNA extraction.

122 2.2. Avoiding PCR-associated biases

After extraction, the genomic DNA is used as a template for PCR amplification, either for amplicon sequencing of 16S rRNA genes or in qPCR. The PCR amplification step can be a major source of bias. These errors are associated with the PCR template properties (template concentration, GC content), primer choice (primer coverage and mismatch), polymerase choice, and the PCR protocol (annealing temperature and PCR cycle number).

128 PCR template properties

Diluting the DNA template concentration can positively affect the PCR efficiency if inhibitors are present in the sample. However, diluting the DNA template concentration for PCR for amplicon sequencing can exclude rare taxa from the detection and therefore results in a lower observed richness (Wu et al., 2010). The potential error of using too low DNA template concentrations can neither be estimated nor corrected for rare taxa and should therefore be avoided.

A high GC content of template DNA of a certain taxon can lead to an under-representation of this taxon (Pinto and Raskin, 2012) because high GC contents cause a less efficient initial denaturation during PCR leading to a lower amplification efficiency (Laursen et al., 2017). The resulting error of the genomic GC content on the relative gene abundances is hard to estimate as it is intertwined with other error sources and can be up to one order of magnitude.

140 Increasing the time for initial denaturation during PCR can reduce this error but cannot141 suppress it completely (Pinto and Raskin, 2012).

142 Primer coverage

143 Targeting complex microbial communities requires primer pairs covering all bacterial and 144 archaeal 16S rRNA gene sequences present in the sample. As indicated by several studies, 145 there are no universal primer pairs covering both domains equally well (Baker et al., 2003; 146 Bru et al., 2008; Takahashi et al., 2014). One suggested solution is the use of separate primer 147 pairs for each domain (Fischer et al., 2016). Furthermore, the use of several primer pairs for 148 bacteria is recommended because primer pairs covering all bacterial phyla are not available 149 (Klindworth et al., 2013) and might even not be feasible. An analysis of 500 bacterial 16S 150 rRNA gene sequences revealed that apart from a very limited region (position 788 to 798) the 151 absolute base conservation is restricted to four or less consecutive bases (Baker et al., 2003). 152 Hence, no primer having a suitable length can be designed to cover all bacterial and even less 153 so all prokaryotic 16S rRNA genes sequences.

For designing highly universal primer pairs, *in silico* coverages are calculated using tools like Silva TestPrime, see for example Klindworth et al. (2013). However, comparisons of predicted and measured coverages show high disagreements for some taxa (Claesson et al., 2010; Fischer et al., 2016; Thijs et al., 2017). Hence, *in silico* coverages cannot reliably predict *in situ* coverage and therefore cannot be applied to correct insufficiently covered taxa *a posteriori*.

The coverage of primers is determined by sequence similarity of primer and target sequence.
A single primer target mismatch can substantially decrease the PCR amplification efficiency
(up to 1000-fold), leading to a significant underrepresentation of the specific taxon (Bru et al.,
2008).

164 Many sequencing techniques have limited read lengths. That is why usually only a limited 165 number of variable regions of 16S rRNA genes are targeted, for example the V3/V4 regions 166 (Albertsen et al., 2015) or the V1-V3 regions (Cai et al., 2013). In a 16S rRNA gene database 167 study, targeting partial regions provided lower phylogenetic resolution than full-length 168 sequences (Kim et al., 2011). In contrast, another study reported similar community 169 compositions when targeting the V1 to V4 regions compared to full-length sequences 170 (Kraková et al., 2016). As a consequence of these conflicting results, there are no indisputable 171 recommendations on which variable regions to address to characterize complex microbial 172 communities. A recent method targeting the SSU rRNA molecules instead of their rRNA genes resulted in a million of full-length rRNA gene sequences free of the primer bias (Karst 173 174 et al., 2018). However, this method is laborious and complex and hence, not yet applicable as 175 standard method.

In conclusion, it is advisable (i) to use primer pairs which are specific for either bacteria or archaea and target each domain separately and (ii) to use two primer pairs for each domain targeting different variable regions.

179 *Choice of polymerase*

180 DNA polymerases used for elongation are characterized by their fidelity and proofreading 181 activity. Fidelity expressed as substitution error rate is in the range of 1 falsely incorporated 182 base per 2 kb or lower (Potapov and Ong, 2017). This means that this error can be neglected 183 considering the typical amplicon length of 300 to 500 bp. In theory, this error could be 184 corrected a posteriori as the substitution error profiles are known (Shagin et al., 2017). 185 However, such bioinformatics tools are not available yet. The proofreading activity of 186 polymerase can lead to primer editing (Gohl et al., 2016). If this happens early during PCR, 187 taxa which show a primer target mismatch will still be amplified efficiently, avoiding the potential 1000-fold bias due to the mismatch as described above. Hence, polymerases with ahigh fidelity and proofreading activity are desirable.

190 PCR protocol

191 Studies on the influence of PCR annealing temperature on relative abundances in microbial 192 communities show conflicting results. No significant effect of lowering the annealing 193 temperature was found for complex communities like chicken caecal samples (Sergeant et al., 194 2012) while an increased number of OTUs was observed in activated sludge (Albertsen et al., 195 2015). Its effect apparently depends on the sequence similarity between primer and template 196 DNA. In case of a perfect match, lowering the annealing temperature had no effect on the 197 PCR product ratios for a mock combination of two bacteria, while an exponential deviation 198 from the initial 1:1 ratio of two bacterial DNA templates was reported when there was one 199 mismatch (Sipos et al., 2007). In conclusion, high annealing temperatures should be avoided 200 as mismatches can lead to biased community composition results.

201 The effect of varying PCR cycle numbers on PCR products is also ambiguous. For activated 202 sludge, no effect was found (Albertsen et al., 2015). This contrasts another study reporting 203 increasing apparent richness when more PCR cycles were performed (Ahn et al., 2012). 204 However, this might be due to the formation of PCR artifacts like chimeric sequences which 205 is supposed to happen when PCR components become limiting. Chimeric sequences are 206 generated when the elongation is not completed within one PCR cycle and the DNA fragment 207 serves as primer for the next cycle binding to template DNA of another taxon present in the 208 sample. Chimeric sequences can be a substantial part of the PCR products as libraries with up 209 to 45% chimeras were found, which can be misinterpreted as additional taxa (Ashelford et al., 210 2006). The chimera formation cannot be quantified in complex communities and therefore, 211 the associated error cannot be corrected. In conclusion, it is beneficial to reduce chimera 212 formation to a minimum by restricting the number of PCR cycles and to filter at best 213 remaining chimeras by appropriate sequence analysis as discussed below.

214 Another type of PCR artifacts are heteroduplex molecules which arise from hybridization of 215 heterologous sequences and which are less frequent than chimeras (Oiu et al., 2001). They 216 will not influence the relative abundance obtained by common NGS methods like Illumina 217 sequencing and pyrosequencing as DNA is denatured before sequencing. This might be 218 different for other methods using dsDNA for sequencing like SMRT PacBio sequencing or 219 Oxford Nanopore sequencing. Formation of heteroduplexes can be reduced by reducing the 220 number of PCR cycles and avoiding high DNA template concentrations as mentioned above 221 (Thompson, 2002).

222 2.3. Quantification and correction of errors associated with variations of strain-specific
 223 16S rRNA operon copy numbers per genome

An innate bias of using 16S rRNA genes for relative and absolute taxa-specific genome abundance quantification are variations of strain-specific 16S rRNA operon copy numbers per genome (Kembel et al., 2012), with reported ranges of 1-17 for bacteria and 1-4 for archaea (Stoddard et al., 2015). An unaccounted higher 16S rRNA operon copy number per genome of a specific taxon will appear as a higher relative and absolute genome abundance of this taxon. Hence, a correction for strain-specific 16S rRNA operon copy numbers per genome is necessary, for example by using the rrnDB database (Stoddard et al., 2015).

The correction for strain-specific 16S rRNA operon copy numbers per genome is incomplete as data is not available for all strains and as operational taxonomic units (OTUs) cannot be taxonomically assigned on strain level. If the strain-specific number is missing, either median or mean copy number for the respective higher taxonomic rank can be used. However, this leads to inaccuracies. 236 Here, we used three example data sets A (Klassen et al., 2017), B (Maus et al., 2017) and C 237 (Müller et al., 2016) to illustrate the effect of correcting for strain-specific 16S rRNA gene 238 operon copy numbers per genome on relative abundances (see Additional File 1, Chapter A.2, 239 for details). In our examples, the copy number correction changed the relative genome 240 abundances compared to the relative 16S rRNA gene abundances by 22% on average though 241 not uniformly for all taxa (Figure 2a). In a standard analysis of data set A without 16S rRNA 242 gene operon copy number correction, the phylum *Bacteroidetes* appeared to be most dominant 243 followed by Firmicutes and Chlorobi. In contrast, after 16S rRNA operon copy number 244 correction Chlorobi was the most dominant phylum, followed by Bacteroidetes and 245 Firmicutes. This was due to the low average 16S rRNA operon copy number per genome for 246 Chlorobi (2.0) compared to the other dominant phyla Bacteroidetes (3.7) and Firmicutes 247 (7.0). Hence, negligence of the strain-specific 16S rRNA operon copy numbers per genome 248 results in biased relative abundances and wrong dominance rankings.



Figure 2: Relative abundance of the 15 most abundant bacterial phyla for three example data sets A (Klassen et al., 2017), B (Maus et al., 2017) and C (Müller et al., 2016). a) Comparison of relative 16S rRNA gene abundance (the standard approach) with relative genome abundance after correction for strain-specific 16S rRNA operon copy number per genome

(corrected). b) Comparison of the relative 16S rRNA gene abundances after use of different strategies for selecting the representative sequence for each OTU. The standard approach uses the seed sequence which was used during the clustering process. Alternatively, the longest or the most abundant sequence per OTU was selected as the representative. Furthermore, an OTU-free approach was applied as implemented in the DADA2 pipeline. c) Comparison of relative 16S rRNA gene abundances obtained with different databases used for taxonomic assignment.

261

262 **3.** Errors solely associated with amplicon sequencing

263 3.1. Influence of sequencing technology on relative gene abundances

264 NGS platforms like 454 pyrosequencing, Illumina, IonTorrent, and PacBio employ different 265 principles for sequencing DNA (Goodwin et al., 2016). This in turn results in substantial 266 differences in length and quality of reads as well as sequencing depth, which, in turn, 267 influences the inferred community composition. As the sequence length is limited due to the 268 sequencing technology, only subsets of the variable regions of the 16S rRNA genes are 269 targeted. This is strongly linked to the choice of primers as discussed above. The read quality 270 heavily depends on the sequencing platform. For example, the frequently used MiSeq 271 Illumina sequencing platform is associated with a sequencing error rate of less than 1% 272 (Schirmer et al., 2016). The sequencing depths of NGS platforms depend on their read 273 numbers which usually range from thousands to millions of sequences. More reads generate 274 higher apparent richness as more rare taxa are detected (Claesson et al., 2010). Nonetheless, 275 there is a trade-off between the efficiency of rare taxa detection and artefactual taxa removal 276 (erroneous sequences) during bioinformatic analyses, see (Zhan and MacIsaac, 2015) for 277 further discussion. Apart from the rare taxa, the same taxa were detected by different sequencing platforms when the same primer pairs were applied (D'Amore et al., 2016;
Tremblay et al., 2015).

In conclusion, a sufficient sequencing depth is necessary to target rare taxa and sequencing error is considered only as a minor source of error. However due a limited read length, the sequencing platform restricts the primer choice which has a larger impact than the NGS platform itself (Hiergeist et al., 2016).

284 3.2. Avoidance of errors associated with NGS data analysis

Next to biases from wet-lab procedures, errors can be introduced during data processing. For the analysis of amplicon sequencing data, several pipelines are available. Different pipelines lead to considerably different relative 16S rRNA gene abundancies (Golob et al., 2017; Plummer and Twin, 2015; Werner et al., 2012). Here, we focused on general biases which can be introduced during analysis of amplicon sequencing data and which are not restricted to a specific analysis pipeline.

291 Filtering low quality and chimeric sequences

292 Quality control of raw sequences, i.e. filtering of low quality and chimeric sequences is 293 common to all pipelines. As chimeric sequences can make up a substantial portion of all 294 reads, their removal is of importance. Sequences can be filtered for chimeras using a reference 295 database (e.g., the Gold reference collection) containing 16S rRNA gene sequences of 296 cultivated bacteria and archaea. However, not all chimeric sequences are detected by this 297 approach (Ahn et al., 2012). Furthermore, this reference chimera detection will perform 298 poorly on microbial communities if they contain yet uncultivated organisms (Schloss et al., 2011). To increase chimera detection efficiency, de novo detection methods can be applied, 299 300 optionally in combination with reference-based detection (Schloss et al., 2011). However, as

still not all of the chimeras can be detected (Haas et al., 2011), it is important to reducechimera formation by adjusting the PCR conditions appropriately as described above.

In conclusion, chimera formation needs to be reduced by optimizing PCR conditions in the first place. A substantial part, but not all of the remaining chimeric sequences can be removed during data analysis. For microbial communities with a lot of undescribed microorganisms, *de novo* detection methods optionally combined with reference-based detection should be applied.

308 OTU clustering

Filtered sequences are clustered into OTUs according to sequence similarity. For clustering, a 16S rRNA gene sequence similarity threshold of 97% is commonly used to differentiate between species, though it represents an arbitrary threshold rather than being based on a commonly accepted species definition (Callahan et al., 2017). Hence, it may assign different species to the same OTU or one species to several OTUs due to sequence variation in multiple 16S rRNA operons.

315 After clustering, one representative sequence per OTU is selected and taxonomically 316 assigned. The effect of the choice of representative sequences for each OTU is illustrated in 317 Figure 2b. As standard in QIIME analysis, the centroid sequence which has been used for 318 defining the OTU is taken as the representative. Alternatively, the longest and the most 319 abundant sequence can be chosen. For the example data sets A (Klassen et al., 2017), B 320 (Maus et al., 2017) and C (Müller et al., 2016), the selection strategy has only a minor impact 321 on the obtained community composition. Relative 16S rRNA gene abundances change by 322 0.3% on average except for data set C if the most abundant sequences per OTU are selected. 323 Here, the average change of 5% is caused by a higher fraction of unclassified OTUs.

In order to avoid arbitrary threshold setting for OTU clustering resulting in misassignments, recently developed analysis pipelines like DADA2 avoid OTU clustering altogether and instead account for sequence variation down to one nucleotide sequence differences (Callahan et al., 2016). Community compositions of the example data sets derived from the OTU clustering approach and exact sequence inference differed by 45% on average comparing the relative 16S rRNA genes abundances, see Figure 2b. A similar difference was reported for gut microbiome samples (Allali et al., 2017).

In conclusion, based on our example data sets, OTU clustering-free approaches should be used to avoid setting arbitrary sequence similarity thresholds. If OTU clustering is desired, selecting the most abundant sequence should be avoided. Instead, the centroid sequence of each OTU should be taken as the representative sequence.

335 Influence of 16S rRNA databases on taxonomic assignment

336 After selecting representative or inferring exact sequences, these are classified against a 337 taxonomic database. The choice of database has a strong influence on the observed 338 community composition for our example data sets A (Klassen et al., 2017), B (Maus et al., 339 2017) and C (Müller et al., 2016), see Figure 2c. Relative 16S rRNA gene abundances 340 changed on average by 9% when using SILVA or the SILVA-based MiDAS taxonomies 341 instead of the Greengenes taxonomy. Previous reports also showed that different databases for 342 taxonomic assignment give different community compositions (Werner et al., 2012). The 343 latest Greengenes version is from 2013 [25] and is problematic for the example data set, 344 because it lacks 16S rRNA sequences of microorganisms which were described more 345 recently, for example syntrophic acetate oxidizing bacteria. SILVA, RDP, and MiDAS are 346 more up to date. Using the RDP database results in a substantially higher number of 347 unclassified OTUs for the example data sets, and therefore does not seem to be 348 recommendable. The MiDAS database is built on the SILVA database and was amended for

taxa present in activated sludge, anaerobic digesters, and influent wastewater (McIlroy et al., 2015, 2017), making it the database of choice for the example data set and anaerobic digester samples in general. Similar to MiDAS, other dedicated databases exist, for example for human intestinal (Ritari et al., 2015), human oral (Chen et al., 2010), and bee intestinal (Newton and Roeselers, 2012) microbial communities.

354 4. Errors solely associated with qPCR data analysis

355 For qPCR, two reporter systems are commonly used, hybridization probes (also called TaqManTM probes) and intercalating dyes such as SYBR Green (Smith and Osborn, 2009). 356 357 Intercalating dyes bind non-specifically to all amplicons. Therefore, post-PCR melting curve 358 analyses need to confirm that only target genes were quantified and not non-specific PCR 359 products such as primer-dimers (Smith and Osborn, 2009). Hybridization probes are designed 360 to bind to a conserved site on the target gene, ensuring that only the target gene is quantified 361 (Smith and Osborn, 2009). Such a conserved site, however, might not exist for a target gene 362 present in various members of a mixed community with potentially unknown members. Here, 363 a probe might bind unequally to the 16S rRNA genes of all members and lead to biased 364 results.

For absolute gene quantification with qPCR, an external standard is necessary. From the 365 366 fluorescence of the sample compared to that of the standard, the 16S rRNA gene copy number 367 in the sample is determined. A qPCR standard contains a known number of 16S rRNA genes, 368 either from a single member or a mixture of members of the microbial community of interest. 369 The 16S rRNA genes can be in the form of a purified PCR product or a plasmid insert. 370 Plasmid standards can be linearized which gave more accurate results when quantifying 371 microalgae (Hou et al., 2010). In that study, the non-linearized standard led to an 372 overestimation of the target gene copy number by 777%. However, a study targeting two 373 bacterial and two archaeal species did not find a systematic overestimation by using circular plasmids, and similar absolute gene numbers were obtained using linearized, supercoiled or
nicked plasmids or a purified PCR product (Oldham and Duncan, 2012).

376 Ideally, a new standard representing the community composition is produced from each 377 sample individually. However this is resource intensive and therefore, it is common to use a 378 standard from a single species which neglects the possibility that the qPCR efficiency of the 379 standard could be atypical for the mixture of target genes of the community. However, for 380 example a 10% lower amplification efficiency of the sample can give rise to 275% 381 overestimated absolute gene copy numbers (Pérez et al., 2013). Lower amplification 382 efficiencies can also be caused by inhibitors in the sample matrix. The amplification 383 efficiencies of all individual samples and standards can be determined with for example the 384 LinRegPCRProgram as described in the literature (Brankatschk et al., 2012). If a difference in 385 amplification efficiency is detected, it is necessary to correct it for example by one-point-386 calibration (Brankatschk et al., 2012) or by repeating the analysis with diluted DNA template 387 concentrations. A convenient spreadsheet to apply the one-point-calibration method is 388 provided in Additional file 2. This spreadsheet can also be used to estimate the error of the 389 differences in efficiencies and other errors discussed above. The problem of amplification 390 efficiency differences between standard and sample in qPCR can be avoided by using digital PCR (dPCR) because it does not require a standard (Kim et al., 2015). 391

The absolute genome abundance of a single taxon can be obtained by qPCR with a specific primer pair. However, in a mixed culture other taxa than the targeted one can be additionally amplified leading to inaccurate quantification results. The use of hybridization probes can reduce this problem. In addition, amplicon sequencing using the same primer pair can be used to identify and correct the influence of unspecific amplification on the absolute genome abundance of a single taxon. Publication of qPCR data should follow the "minimum information for publication of quantitative real-time PCR experiments" (MIQE) (Bustin et al., 2009). In particular, the publication of the qPCR raw data is desirable to enable post-publication error estimates and corrections for amplification efficiency differences.

In conclusion, there are several error sources in qPCR that can sum up to several orders of
magnitude. Differences in amplification efficiency between standard and sample are often
neglected but can be corrected by one-point calibration. Either linearized plasmids, circular
(supercoiled or nicked) plasmids or PCR products can be used as a standard for prokaryotes.
A decision tree helping to avoid the above mentioned errors is provided in Figure 3.



408 Figure 3: Guideline for minimizing absolute 16S rRNA gene abundance quantification errors409 in qPCR

410

411 **5.** Error identification with plausibility checks based on environmental parameters

The quantification of absolute taxa-specific genome abundances can yield errors of up to several orders of magnitude. Plausibility checks are thus highly desirable to validate results and to avoid conclusions based on erroneous data. Several methods are available to estimate absolute biomass concentrations based on process parameters and environmental conditions.
Mechanistic mathematical models, such as the Anaerobic Digestion Model No.1 (Batstone et
al., 2002) can be used to predict microbial biomasses. However, such models require
extensive information as input. For systems with scarce information, black box approaches
are more suitable, in particular methods including thermodynamic considerations (Heijnen,
2013; Kleerebezem and Van Loosdrecht, 2010).

421 An example for a simple plausibility check is presented in Figure 4. Absolute archaeal 16S 422 rRNA gene copy numbers measured in anaerobic digesters (Lee et al., 2011; Nettmann et al., 423 2010) were compared with minimum gene copy numbers estimated by a black box approach. 424 These estimates were based on the assumption that the catabolism of any microbial cell is 425 limited by a maximum rate for electron transport (Heijnen, 2002). Given this assumption, the 426 minimum cell number required to produce the amount of methane measured in the digesters 427 was calculated (see Additional File 1, section A.4, for details). For the conversion of cell 428 numbers to gene copies, each cell was conservatively assumed to contain one archaeal 16S 429 rRNA gene copy.

430 The archaeal 16S rRNA gene copy numbers measured by Nettmann et al. (2010) lie well 431 above the estimated minimum copy numbers. However, the copy numbers measured by Lee et 432 al. (2010) are orders of magnitude below the estimated minimum. This indicates that their 433 results likely underestimate absolute quantities in the sampled digesters. DNA extraction 434 efficiencies were not considered in that study, which might be a reason for the implausibly 435 low copy numbers. This example illustrates how even simple plausibility checks can be used 436 to identify implausible quantification results. Such plausibility checks are not restricted to 437 methanogenic environments. Kleerebezem and Van Loosdrecht (2010) for example provided 438 biomass yield estimates for 61 organic compounds with either oxygen, nitrate, sulfate and 439 carbon dioxide as electron acceptors.



Figure 4: Example for checking the plausibility of absolute 16S rRNA gene quantification.
Comparison of experimentally derived archaeal 16S rRNA gene copies in anaerobic digesters
(Lee et al., 2011; Nettmann et al., 2010) with estimated minimum gene copy number based on
the maximum electron transfer rate (Heijnen, 2002) and the methane production rates
measured in the digesters (normalized to the digesters working volume).

447 6. Errors associated with converting genome abundances to cell numbers and448 biomasses

Ideally, taxa-specific genome abundances could be converted to more tangible cell numbers. However, this requires taxa-specific information on the ploidy, i.e. the number of genome copies per cell. Prokaryotes have historically been considered as monoploid (1 genome copy per cell) (Pecoraro et al., 2011). However, several recent studies have found oligoploid (<10 genome copies per cell) and polyploid (>10 genome copies per cell) archaea and bacteria and it appears that monoploid prokaryotes are rather the exception than the rule (Soppa, 2014). 455 A correction for ploidy is difficult. The ploidy can differ even within one genus, for example 456 within Desulfovibrio and within Neisseria (Pecoraro et al., 2011). Furthermore, the ploidy of a 457 species does not only vary between one and two during the cell cycle but can differ greatly 458 between different growth phases, for example from 3-15 genome copies per cell in the 459 exponential phase to 2-4 genome copies in the stationary phase for Methanocaldococcus 460 *jannaschii* (Pecoraro et al., 2011). Taxa-specific ploidy has been determined experimentally 461 for pure cultures by combining qPCR with cell counting (Pecoraro et al., 2011). For complex 462 communities, qPCR could be combined with fluorescence-activated cell sorting in the future.

463 Deriving taxa-specific cell numbers from complex communities was recently suggested by 464 combining absolute cell numbers derived from flow cytometry with taxa-specific 16S rRNA 465 genome abundances (Props et al., 2017). However, this method requires the same ploidy of all 466 taxa which is unlikely given the high variance of ploidy found in prokaryotes (Pecoraro et al., 467 2011).

468 As mentioned above, mechanistic models often consider microorganisms not as cells but as 469 biomass, but it is difficult to determine experimentally the taxa-specific biomasses of species-470 rich complex microbial communities. Nevertheless, the distinction between cell number and 471 biomass must not be ignored, because the prokaryotic cell masses can vary over several orders 472 of magnitude (Loferer-Krößbacher et al., 1998). Taxa-specific biomasses can be inferred from 473 cell volumes determined by FISH combined with digital image analysis (Daims, 2009). In 474 addition to FISH, metaproteomics has recently been suggested as a method for taxa-specific 475 biomass quantification (Kleiner et al., 2017).

In conclusion, taxa-specific average genome copies and biomasses per cell can vary
substantially between taxa and neglecting this fact can lead to errors of up to 10,000%.
Consequently, the accurate conversion between genome abundances and cell numbers as well
as biomasses remains a challenge.

481 **7.** Conclusions

482 Relative and absolute taxa-specific genome abundances are important parameters for studying 483 microbial community dynamics, but their quantification with PCR based approaches has a 484 number of potential errors that can reach several orders of magnitudes. These errors and 485 suggested measures to avoid or reduce them are summarized in Table 1. Many errors can 486 already be reduced by proper data analysis. Others require additional experimental effort, 487 such as spiking of known microorganisms to estimate DNA extraction efficiencies. Using 488 different primer pairs for bacteria and archaea is essential for accurate analyses but adds 489 substantial experimental effort. The accurate conversion of taxa-specific genome abundances 490 to cell numbers and biomasses is important for their use in mathematical models but remains a 491 challenge. Flow cytometry, FISH and metaproteomics might bring valuable culture-492 independent contributions in the future to solve this problem.

495 Additional files

Additional file 1: Further information on the studied errors and concepts, .pdf,
Details on the example NGS data sets and the data analysis. Details on the minimum
gene number estimate. Illustration of the difference in using relative and absolute cell
numbers for interpreting the influence of distinct taxa on the process performance.
Additional information on correcting strain-specific 16S rRNA operon copy numbers.
Additional file 2: Tool for qPCR data analysis und error estimation, .xlsx, Excel

sheet to simplify the analysis of qPCR data for absolute quantification and errorestimation.

505 **Funding**

- 506 This work was funded by the German Federal Ministry of Education and Research (e:Bio
- 507 project "McBiogas", FKZ 031A317).

508

509 **Conflict of interest**

510 The authors declare they have no conflict of interest.

511 **References**

- Ahn, J.-H., Kim, B.-Y., Song, J., and Weon, H.-Y. (2012). Effects of PCR cycle number and
 DNA polymerase type on the 16S rRNA gene pyrosequencing analysis of bacterial
 communities. J. Microbiol. 50, 1071–1074. doi:10.1007/s12275-012-2642-z.
- Albertsen, M., Karst, S. M., Ziegler, A. S., Kirkegaard, R. H., and Nielsen, P. H. (2015). Back
 to basics the influence of DNA extraction and primer choice on phylogenetic analysis
 of activated sludge communities. *PLoS One* 10, e0132783.
 doi:10.1371/journal.pone.0132783.
- Allali, I., Arnold, J. W., Roach, J., Cadenas, M. B., Butz, N., Hassan, H. M., et al. (2017). A
 comparison of sequencing platforms and bioinformatics pipelines for compositional
 analysis of the gut microbiome. *BMC Microbiol*. 17, 194. doi:10.1186/s12866-017-11018.
- Ashelford, K. E., Chuzhanova, N. a., Fry, J. C., Jones, A. J., and Weightman, A. J. (2006).
 New screening software shows that most recent large 16S rRNA gene clone libraries
 contain chimeras. *Appl. Environ. Microbiol.* 72, 5734–5741. doi:10.1128/AEM.0055606.
- Baker, G. C., Smith, J. J., and Cowan, D. A. (2003). Review and re-analysis of domainspecific 16S primers. J. Microbiol. Methods 55, 541–555.
 doi:10.1016/j.mimet.2003.08.009.
- Batstone, D. J., Keller, J., Angelidaki, I., Kalyuzhnyi, S. V., Pavlostathis, S. G., Rozzi, A., et
 al. (2002). *Anaerobic Digestion Model No.1*. (*ADM1*). London: IWA Publishing.
- Brankatschk, R., Bodenhausen, N., Zeyer, J., and Bürgmann, H. (2012). Simple absolute
 quantification method correcting for quantitative PCR efficiency variations for microbial
 community samples. *Appl. Environ. Microbiol.* 78, 4481–4489.
 doi:10.1128/AEM.07878-11.
- Bru, D., Martin-Laurent, F., and Philippot, L. (2008). Quantification of the detrimental effect
 of a single primer-template mismatch by real-Time PCR using the 16S rRNA gene as an
 example. *Appl. Environ. Microbiol.* 74, 1660–1663. doi:10.1128/AEM.02403-07.
- Bustin, S. A., Benes, V., Garson, J. A., Hellemans, J., Huggett, J., Kubista, M., et al. (2009).
 The MIQE guidelines: minimum information for publication of quantitative real-time
 PCR experiments. *Clin. Chem.* 55, 611–622. doi:10.1373/clinchem.2008.112797.
- 542 Cai, L., Ye, L., Tong, A. H. Y., Lok, S., and Zhang, T. (2013). Biased diversity metrics
 543 revealed by bacterial 16S pyrotags derived from different primer sets. *PLoS One* 8, e53649. doi:10.1371/journal.pone.0053649.
- Callahan, B. J., McMurdie, P. J., and Holmes, S. P. (2017). Exact sequence variants should
 replace operational taxonomic units in marker-gene data analysis. *ISME J.* 11, 2639–2643. doi:10.1038/ismej.2017.119.
- Callahan, B. J., McMurdie, P. J., Rosen, M. J., Han, A. W., Johnson, A. J. A., and Holmes, S.
 P. (2016). DADA2: High-resolution sample inference from Illumina amplicon data. *Nat. Methods* 13, 581–583. doi:10.1038/nmeth.3869.
- Chen, T., Yu, W.-H., Izard, J., Baranova, O. V., Lakshmanan, A., and Dewhirst, F. E. (2010).
 The Human Oral Microbiome Database: a web accessible resource for investigating oral
 microbe taxonomic and genomic information. *Database* 2010, baq013-baq013.
 doi:10.1093/database/baq013.

- Claesson, M. J., Wang, Q., O'Sullivan, O., Greene-Diniz, R., Cole, J. R., Ross, R. P., et al.
 (2010). Comparison of two next-generation sequencing technologies for resolving highly
 complex microbiota composition using tandem variable 16S rRNA gene regions. *Nucleic Acids Res.* 38, e200–e200. doi:10.1093/nar/gkq873.
- D'Amore, R., Ijaz, U. Z., Schirmer, M., Kenny, J. G., Gregory, R., Darby, A. C., et al. (2016).
 A comprehensive benchmarking study of protocols and sequencing platforms for 16S
 rRNA community profiling. *BMC Genomics* 17, 55. doi:10.1186/s12864-015-2194-9.
- Daims, H. (2009). Use of fluorescence in situ hybridization and the daime image analysis
 program for the cultivation-independent quantification of microorganisms in
 environmental and medical samples. *Cold Spring Harb. Protoc.* 4, 1–7.
 doi:10.1101/pdb.prot5253.
- Dannemiller, K. C., Lang-Yona, N., Yamamoto, N., Rudich, Y., and Peccia, J. (2014).
 Combining real-time PCR and next-generation DNA sequencing to provide quantitative
 comparisons of fungal aerosol populations. *Atmos. Environ.* 84, 113–121.
 doi:10.1016/j.atmosenv.2013.11.036.
- Eloe-Fadrosh, E. A., Ivanova, N. N., Woyke, T., and Kyrpides, N. C. (2016). Metagenomics
 uncovers gaps in amplicon-based detection of microbial diversity. *Nat. Microbiol.* 1,
 15032. doi:10.1038/nmicrobiol.2015.32.
- 573 Emerson, J. B., Adams, R. I., Román, C. M. B., Brooks, B., Coil, D. a., Dahlhausen, K., et al.
 574 (2017). Schrödinger's microbes: tools for distinguishing the living from the dead in
 575 microbial ecosystems. *Microbiome* 5, 86. doi:10.1186/s40168-017-0285-3.
- Fischer, M. A., Güllert, S., Neulinger, S. C., Streit, W. R., and Schmitz, R. A. (2016).
 Evaluation of 16S rRNA gene primer pairs for monitoring microbial community
 structures showed high reproducibility within and low comparability between datasets
 generated with multiple archaeal and bacterial primer pairs. *Front. Microbiol.* 7.
 doi:10.3389/fmicb.2016.01297.
- Gohl, D. M., Vangay, P., Garbe, J., MacLean, A., Hauge, A., Becker, A., et al. (2016).
 Systematic improvement of amplicon marker gene methods for increased accuracy in microbiome studies. *Nat. Biotechnol.* 34, 942–949. doi:10.1038/nbt.3601.
- Golob, J. L., Margolis, E., Hoffman, N. G., and Fredricks, D. N. (2017). Evaluating the
 accuracy of amplicon-based microbiome computational pipelines on simulated human
 gut microbial communities. *BMC Bioinformatics* 18, 283. doi:10.1186/s12859-0171690-0.
- Goodwin, S., McPherson, J. D., and McCombie, W. R. (2016). Coming of age: ten years of
 next-generation sequencing technologies. *Nat. Rev. Genet.* 17, 333–351.
 doi:10.1038/nrg.2016.49.
- Haas, B. J., Gevers, D., Earl, A. M., Feldgarden, M., Ward, D. V., Giannoukos, G., et al.
 (2011). Chimeric 16S rRNA sequence formation and detection in Sanger and 454pyrosequenced PCR amplicons. *Genome Res.* 21, 494–504. doi:10.1101/gr.112730.110.
- Heijnen, J. J. (2002). "Bioenergetics of Microbial GrowthMicrobial Growth," in *Encyclopedia of Bioprocess Technology*, eds. M. C. Flickinger and S. W. Drew (Hoboken, NJ, USA:
 John Wiley & Sons, Inc.), 267–291. doi:10.1002/0471250589.ebt026.
- Heijnen, J. J. (2013). "A thermodynamic approach to predict black box model parameters for
 microbial growth," in *Biothermodynamics* (EPFL Press), 443–472.
- Hiergeist, A., Reischl, U., and Gessner, A. (2016). Multicenter quality assessment of 16S
 ribosomal DNA-sequencing for microbiome analyses reveals high inter-center

- 601 variability. Int. J. Med. Microbiol. 306, 334–342. doi:10.1016/j.ijmm.2016.03.005.
- Hou, Y., Zhang, H., Miranda, L., and Lin, S. (2010). Serious overestimation in quantitative
 PCR by circular (supercoiled) plasmid standard: microalgal pcna as the model gene. *PLoS One* 5, e9545. doi:10.1371/journal.pone.0009545.
- Kembel, S. W., Wu, M., Eisen, J. A., and Green, J. L. (2012). Incorporating 16S gene copy
 number information improves estimates of microbial diversity and abundance. *PLoS Comput. Biol.* 8, e1002743. doi:10.1371/journal.pcbi.1002743.
- Kim, M., Morrison, M., and Yu, Z. (2011). Evaluation of different partial 16S rRNA gene
 sequence regions for phylogenetic analysis of microbiomes. *J. Microbiol. Methods* 84,
 81–87. doi:10.1016/j.mimet.2010.10.020.
- Kim, T. G., Jeong, S.-Y., and Cho, K.-S. (2015). Development of droplet digital PCR assays
 for methanogenic taxa and examination of methanogen communities in full-scale
 anaerobic digesters. *Appl. Microbiol. Biotechnol.* 99, 445–458. doi:10.1007/s00253-0146007-x.
- Kirkegaard, R. H., McIlroy, S. J., Kristensen, J. M., Nierychlo, M., Karst, S. M., Dueholm, M.
 S., et al. (2017). The impact of immigration on microbial community composition in fullscale anaerobic digesters. *Sci. Rep.* 7, 9343. doi:10.1038/s41598-017-09303-0.
- Klassen, V., Blifernez-Klassen, O., Wibberg, D., Winkler, A., Kalinowski, J., Posten, C., et
 al. (2017). Highly efficient methane generation from untreated microalgae biomass. *Biotechnol. Biofuels* 10, 186. doi:10.1186/s13068-017-0871-4.
- Kleerebezem, R., and Van Loosdrecht, M. C. M. (2010). A generalized method for
 thermodynamic state analysis of environmental systems. *Crit. Rev. Environ. Sci. Technol.* 40, 1–54. doi:10.1080/10643380802000974.
- Kleiner, M., Thorson, E., Sharp, C. E., Dong, X., Liu, D., Li, C., et al. (2017). Assessing
 species biomass contributions in microbial communities via metaproteomics. *Nat. Commun.* 8, 1558. doi:10.1038/s41467-017-01544-x.
- Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M., et al. (2013).
 Evaluation of general 16S ribosomal RNA gene PCR primers for classical and nextgeneration sequencing-based diversity studies. *Nucleic Acids Res.* 41, e1–e1.
 doi:10.1093/nar/gks808.
- Kraková, L., Šoltys, K., Budiš, J., Grivalský, T., Ďuriš, F., Pangallo, D., et al. (2016).
 Investigation of bacterial and archaeal communities: novel protocols using modern sequencing by Illumina MiSeq and traditional DGGE-cloning. *Extremophiles* 20, 795– 808. doi:10.1007/s00792-016-0855-5.
- Lambrecht, J., Cichocki, N., Hübschmann, T., Koch, C., Harms, H., and Müller, S. (2017).
 Flow cytometric quantification, sorting and sequencing of methanogenic archaea based on F420 autofluorescence. *Microb. Cell Fact.* 16, 1–15. doi:10.1186/s12934-017-0793-7.
- Laursen, M. F., Dalgaard, M. D., and Bahl, M. I. (2017). Genomic GC-content affects the
 accuracy of 16S rRNA gene sequencing based microbial profiling due to PCR bias. *Front. Microbiol.* 8, 1–8. doi:10.3389/fmicb.2017.01934.
- Lebuhn, M., Derenkó, J., Rademacher, A., Helbig, S., Munk, B., Pechtl, A., et al. (2016).
 DNA and RNA extraction and quantitative real-time PCR-based assays for biogas
 biocenoses in an interlaboratory comparison. *Bioengineering* 3, 7.
 doi:10.3390/bioengineering3010007.
- Lee, I. S., Parameswaran, P., and Rittmann, B. E. (2011). Effects of solids retention time on methanogenesis in anaerobic digestion of thickened mixed sludge. *Bioresour. Technol.*

- 647 102, 10266–10272. doi:10.1016/j.biortech.2011.08.079.
- Loferer-Krößbacher, M., Klima, J., and Psenner, R. (1998). Determination of bacterial cell
 dry mass by transmission electron microscopy and densitometric image analysis. *Appl. Environ. Microbiol.* 64, 688–694.
- 651
- Maus, I., Kim, Y. S., Wibberg, D., Stolze, Y., Off, S., Antonczyk, S., et al. (2017). Biphasic
 study to characterize agricultural biogas plants by high-throughput 16S rRNA gene
 amplicon sequencing and microscopic analysis. *J. Microbiol. Biotechnol.* 27, 321–334.
 doi:10.4014/jmb.1605.05083.
- McIlroy, S. J., Kirkegaard, R. H., McIlroy, B., Nierychlo, M., Kristensen, J. M., Karst, S. M.,
 et al. (2017). MiDAS 2.0: an ecosystem-specific taxonomy and online database for the
 organisms of wastewater treatment systems expanded for anaerobic digester groups. *Database* 2017. doi:10.1093/database/bax016.
- McIlroy, S. J., Saunders, A. M., Albertsen, M., Nierychlo, M., McIlroy, B., Hansen, A. A., et
 al. (2015). MiDAS: the field guide to the microbes of activated sludge. *Database* 2015,
 bav062. doi:10.1093/database/bav062.
- Müller, B., Sun, L., Westerholm, M., and Schnürer, A. (2016). Bacterial community
 composition and fhs profiles of low- and high-ammonia biogas digesters reveal novel
 syntrophic acetate-oxidising bacteria. *Biotechnol. Biofuels* 9, 48. doi:10.1186/s13068016-0454-9.
- Nettmann, E., Bergmann, I., Pramschufer, S., Mundt, K., Plogsties, V., Herrmann, C., et al.
 (2010). Polyphasic analyses of methanogenic archaeal communities in agricultural
 biogas plants. *Appl. Environ. Microbiol.* 76, 2540–2548. doi:10.1128/AEM.01423-09.
- Newton, I. L., and Roeselers, G. (2012). The effect of training set on the classification of
 honey bee gut microbiota using the Naïve Bayesian Classifier. *BMC Microbiol.* 12, 221.
 doi:10.1186/1471-2180-12-221.
- Oldham, A. L., and Duncan, K. E. (2012). Similar gene estimates from circular and linear
 standards in quantitative PCR analyses using the prokaryotic 16S rRNA gene as a model. *PLoS One* 7, e51931. doi:10.1371/journal.pone.0051931.
- Pecoraro, V., Zerulla, K., Lange, C., and Soppa, J. (2011). Quantification of ploidy in
 proteobacteria revealed the existence of monoploid, (mero-)oligoploid and polyploid
 species. *PLoS One* 6. doi:10.1371/journal.pone.0016392.
- Pérez, L. M., Fittipaldi, M., Adrados, B., Morató, J., and Codony, F. (2013). Error estimation
 in environmental DNA targets quantification due to PCR efficiencies differences
 between real samples and standards. *Folia Microbiol. (Praha).* 58, 657–662.
 doi:10.1007/s12223-013-0255-5.
- Pinto, A. J., and Raskin, L. (2012). PCR biases distort bacterial and archaeal community
 structure in pyrosequencing datasets. *PLoS One* 7, e43093.
 doi:10.1371/journal.pone.0043093.
- Plummer, E., and Twin, J. (2015). A comparison of three bioinformatics pipelines for the
 analysis of preterm gut microbiota using 16S rRNA gene sequencing data. J. Proteomics *Bioinform.* 8. doi:10.4172/jpb.1000381.
- Potapov, V., and Ong, J. L. (2017). Examining sources of error in PCR by single-molecule
 sequencing. *PLoS One* 12, 1–19. doi:10.1371/journal.pone.0181128.
- 691 Props, R., Kerckhof, F.-M., Rubbens, P., De Vrieze, J., Hernandez Sanabria, E., Waegeman,

- W., et al. (2017). Absolute quantification of microbial taxon abundances. *ISME J.* 11,
 584–587. doi:10.1038/ismej.2016.117.
- Qiu, X., Wu, L., Huang, H., McDonel, P. E., Palumbo, a. V., Tiedje, J. M., et al. (2001).
 Evaluation of PCR-generated chimeras, mutations, and heteroduplexes with 16S rRNA
 gene-based cloning. *Appl. Environ. Microbiol.* 67, 880–887. doi:10.1128/AEM.67.2.880887.2001.
- Ritari, J., Salojärvi, J., Lahti, L., and de Vos, W. M. (2015). Improved taxonomic assignment
 of human intestinal 16S rRNA sequences by a dedicated reference database. *BMC Genomics* 16, 1056. doi:10.1186/s12864-015-2265-y.
- Schirmer, M., D'Amore, R., Ijaz, U. Z., Hall, N., and Quince, C. (2016). Illumina error
 profiles: resolving fine-scale variation in metagenomic sequencing data. *BMC Bioinformatics* 17, 1–15. doi:10.1186/s12859-016-0976-y.
- Schloss, P. D., Gevers, D., and Westcott, S. L. (2011). Reducing the effects of PCR
 amplification and sequencing artifacts on 16S rRNA-based studies. *PLoS One* 6, e27310.
 doi:10.1371/journal.pone.0027310.
- Sergeant, M. J., Constantinidou, C., Cogan, T., Penn, C. W. C. W., and Pallen, M. J. M. J.
 (2012). High-throughput sequencing of 16S rRNA gene amplicons: effects of extraction
 procedure, primer length and annealing temperature. *PLoS One* 7, e38094.
 doi:10.1371/journal.pone.0038094.
- Shagin, D. a., Shagina, I. a., Zaretsky, A. R., Barsova, E. V., Kelmanson, I. V., Lukyanov, S.,
 et al. (2017). A high-throughput assay for quantitative measurement of PCR errors. *Sci. Rep.* 7, 2718. doi:10.1038/s41598-017-02727-8.
- Shakya, M., Quince, C., Campbell, J. H., Yang, Z. K., Schadt, C. W., and Podar, M. (2013).
 Comparative metagenomic and rRNA microbial diversity characterization using archaeal
 and bacterial synthetic communities. *Environ. Microbiol.* 15, 1882–1899.
 doi:10.1111/1462-2920.12086.
- Sipos, R., Székely, A. J., Palatinszky, M., Révész, S., Márialigeti, K., and Nikolausz, M.
 (2007). Effect of primer mismatch, annealing temperature and PCR cycle number on 16S
 rRNA gene-targetting bacterial community analysis. *FEMS Microbiol. Ecol.* 60, 341–
 350. doi:10.1111/j.1574-6941.2007.00283.x.
- Smith, C. J., and Osborn, a. M. (2009). Advantages and limitations of quantitative PCR (QPCR)-based approaches in microbial ecology. *FEMS Microbiol. Ecol.* 67, 6–20.
 doi:10.1111/j.1574-6941.2008.00629.x.
- Soppa, J. (2014). Polyploidy in archaea and bacteria: About desiccation resistance, giant cell
 size, long-term survival, enforcement by a eukaryotic host and additional aspects. *J. Mol. Microbiol. Biotechnol.* 24, 409–419. doi:10.1159/000368855.
- Steinberg, L. M., and Regan, J. M. (2008). Phylogenetic comparison of the methanogenic communities from an acidic, oligotrophic fen and an anaerobic digester treating municipal wastewater sludge. *Appl. Environ. Microbiol.* 74, 6663–6671. doi:10.1128/AEM.00553-08.
- Stoddard, S. F., Smith, B. J., Hein, R., Roller, B. R. K., and Schmidt, T. M. (2015). rrnDB:
 improved tools for interpreting rRNA gene abundance in bacteria and archaea and a new
 foundation for future development. *Nucleic Acids Res.* 43, D593–D598.
 doi:10.1093/nar/gku1201.
- Takahashi, S., Tomita, J., Nishioka, K., Hisada, T., and Nishijima, M. (2014). Development
 of a prokaryotic universal primer for simultaneous analysis of bacteria and archaea using

- next-generation sequencing. *PLoS One* 9, e105592. doi:10.1371/journal.pone.0105592.
- Thijs, S., Op De Beeck, M., Beckers, B., Truyens, S., Stevens, V., Van Hamme, J. D., et al.
 (2017). Comparative evaluation of four bacteria-specific primer pairs for 16S rRNA gene
 surveys. *Front. Microbiol.* 8. doi:10.3389/fmicb.2017.00494.
- Thompson, J. R. (2002). Heteroduplexes in mixed-template amplifications: formation,
 consequence and elimination by "reconditioning PCR." *Nucleic Acids Res.* 30, 2083–
 2088. doi:10.1093/nar/30.9.2083.
- Tremblay, J., Singh, K., Fern, A., Kirton, E. S., He, S., Woyke, T., et al. (2015). Primer and
 platform effects on 16S rRNA tag sequencing. *Front. Microbiol.* 6.
 doi:10.3389/fmicb.2015.00771.
- Werner, J. J., Koren, O., Hugenholtz, P., DeSantis, T. Z., Walters, W. A., Caporaso, J. G., et
 al. (2012). Impact of training sets on classification of high-throughput bacterial 16s
 rRNA gene surveys. *ISME J.* 6, 94–103. doi:10.1038/ismej.2011.82.
- Willner, D., Daly, J., Whiley, D., Grimwood, K., Wainwright, C. E., and Hugenholtz, P.
 (2012). Comparison of DNA extraction methods for microbial community profiling with
 an application to pediatric bronchoalveolar lavage samples. *PLoS One* 7, e34605.
 doi:10.1371/journal.pone.0034605.
- Wu, J.-Y., Jiang, X.-T., Jiang, Y.-X., Lu, S.-Y., Zou, F., and Zhou, H.-W. (2010). Effects of
 polymerase, template dilution and cycle number on PCR based 16 S rRNA diversity
 analysis using the deep sequencing method. *BMC Microbiol.* 10, 255. doi:10.1186/14712180-10-255.
- Yu, Y., Lee, C., Kim, J., and Hwang, S. (2005). Group-specific primer and probe sets to
 detect methanogenic communities using quantitative real-time polymerase chain
 reaction. *Biotechnol. Bioeng.* 89, 670–679. doi:10.1002/bit.20347.
- Zhan, A., and MacIsaac, H. J. (2015). Rare biosphere exploration using high-throughput
 sequencing: research progress and perspectives. *Conserv. Genet.* 16, 513–522.
 doi:10.1007/s10592-014-0678-9.

766 Table 1: Estimation, avoidance and correction of errors in PCR-based taxa-specific absolute genome abundance quantification

Error type and source	Potential	Estimation of error	Avoidance of error a priori	Correction of error a posteriori	Additional effort for error		
	error size				avoidance/correction		
NGS and qPCR associated errors							
Strain-specific DNA extraction bias	N/A	Spike mock community	Optimize extraction protocol	Impractical	High		
Extraction efficiency	~1000 %	Spike standard (E. coli cells)	Optimize extraction protocol	Quantifying spiked standard	Medium		
DNA Template concentration PCR for	Low (only	Not possible (?)	Avoid too low concentrations	Impractical	Low		
amplicon sequencing	rare taxa)						
Primer coverage	~100,000%	Impractical	(1) Separate primer pairs for bacteria and archaea;	Impractical	High for (1)		
			(2) Additional primer pair for each domain covering different		Very high for (1) + (2)		
			variable region				
PCR	~1000%	Impractical	Reduce cycles to reduce PCR artifacts, use high fidelity	Identification of chimeras and	low		
			polymerase	their removal			
NGS platform	low	Impractical	Impractical	Impractical	N/A		
NGS data analysis							
OTU clustering	~10 %	N/A	Use OTU free clustering approach		Low (only computational)		
Taxonomic database choice	~10%	N/A	Use dedicated database		Low (only computational)		
Strain-specific 16S rRNA operon copy	~100%	N/A	N/A	Use amplicon sequencing data	Low (only computational)		
number per genome				and rrnDB database			
qPCR data analysis							
PCR efficiency standard vs sample	~100-1000%	Use amplification curves	Optimize PCR protocol + standard. Use dPCR instead.	Use one-point calibration method	Low (only computational)		