


Avoiding quantification bias in metabarcoding: Application of a cell biovolume correction factor in diatom molecular biomonitoring

Valentin Vasselon¹  | Agnès Bouchez¹ | Frédéric Rimet¹ | Stéphan Jacquet¹ | Rosa Trobajo² | Méline Corniquel¹ | Kálmán Tapolczai¹ | Isabelle Domaizon¹

¹CARTELE, French National Institute for Agricultural Research (INRA), University of Savoie Mont Blanc, Thonon-les-Bains, France

²Aquatic Ecosystems, Institute for Food and Agricultural Research and Technology (IRTA), Catalunya, Spain

Correspondence

Valentin Vasselon
Email: valentin.vasselon@inra.fr

Funding information

Agence Française pour la Biodiversité (AFB)

Handling Editor: Andrew Mahon

Abstract

1. In recent years, remarkable progress has been made in developing environmental DNA metabarcoding. However, its ability to quantify species relative abundance remains uncertain, limiting its application for biomonitoring. In diatoms, although the *rbcl* gene appears to be a suitable barcode for diatoms, providing relevant qualitative data to describe taxonomic composition, improvement of species quantification is still required.
2. Here, we hypothesized that *rbcl* copy number is correlated with diatom cell biovolume (as previously described for the 18S gene) and that a correction factor (CF) based on cell biovolume should be applied to improve taxa quantification. We carried out a laboratory experiment using pure cultures of eight diatom species with contrasted cell biovolumes in order to (1) verify the relationship between *rbcl* copy numbers (estimated by qPCR) and diatom cell biovolumes and (2) define a potential CF. In order to evaluate CF efficiency, five mock communities were created by mixing different amounts of DNA from the eight species, and were sequenced using HTS and targeting the same *rbcl* barcode.
3. As expected, the correction of DNA reads proportions by the CF improved the congruence between morphological and molecular inventories. Final validation of the CF was obtained on environmental samples (metabarcoding data from 80 benthic biofilms) for which the application of CF allowed differences between molecular and morphological water quality indices to be reduced by 47%.
4. Overall, our results highlight the usefulness of applying a CF factor, which is effective in reducing over-estimation of high biovolume species, correcting quantitative biases in diatom metabarcoding studies and improving final water quality assessment.

KEYWORDS

benthic diatom, biovolume correction factor, freshwater ecosystems, gene copy number variation, quantitative metabarcoding

1 | INTRODUCTION

DNA metabarcoding allows species present in an environmental sample to be detected using a short DNA marker specific for a particular taxonomic group (Taberlet, Coissac, Hajibabaei, & Rieseberg, 2012). Combined with high-throughput sequencing (HTS), hundreds of samples can be analysed at the same time, offering an alternative to microscopy with higher resolution and accuracy, while being faster and cheaper (Stein, Martinez, Stiles, Miller, & Zakharov, 2014). This is particularly interesting for freshwater biomonitoring, in which thousands of river samples have to be analysed annually and management actions applied quickly (Keck, Vasselon, Tapolczai, Rimet, & Bouchez, 2017). The European Water Framework Directive (WFD, European Council, 2000) has implemented the use of benthic diatoms, among other biological indicators (fishes, macroinvertebrates, macrophytes and phytoplankton), for the assessment of aquatic ecosystem integrity. The different biotic diatom indices that have been developed are based on the relative abundances and the ecological values (sensitivity and tolerance to pollutants) of the species observed in rivers and lakes systems (e.g. Rimet, 2012). Different studies have already revealed the potential application of diatom metabarcoding in freshwater quality assessment (Apothéoz-Perret-Gentil et al., 2017; Kermarrec et al., 2014; Vasselon, Domaizon, Rimet, Kahlert, & Bouchez, 2017; Vasselon, Rimet, Tapolczai, & Bouchez, 2017; Visco et al., 2015). However, discrepancies between DNA metabarcoding and microscopy have been observed in species composition and relative abundance (Zimmermann, Glöckner, Jahn, Enke, & Gemeinholzer, 2015). This drawback is likely to affect the congruence between morphological and DNA metabarcoding quality index values and, *in fine*, the ecological assessment.

With respect to qualitative aspects, the incompleteness of the reference databases, the choice of the DNA marker and the efficiency of the PCR primers have been identified as important biases affecting species detection using DNA metabarcoding (Pawlowski, Lejzerowicz, Apothéoz-Perret-Gentil, Visco, & Esling, 2016). For benthic diatoms, the *rbcl* gene has proved to be an appropriate taxonomic marker for biomonitoring (Kermarrec et al., 2013, 2014; Vasselon, Domaizon, et al., 2017; Vasselon, Rimet, et al., 2017) and a well-curated barcode reference library is already available in open-access to assign species names to *rbcl* sequences (R-Syst::diatom, Rimet et al., 2016). However, no clear relationship has yet been demonstrated between the relative species abundances obtained by DNA metabarcoding with the *rbcl* barcode and those obtained by morphological observations (Rimet et al., 2014). As quantification of diatom species is required by the WFD for quality index calculation, more investigation is needed to understand and correct biases affecting diatom quantification based on HTS data.

Species quantification based on HTS data can be estimated from the number of DNA sequences (i.e. reads) assigned to each species, from which relative abundances can be calculated. Previous studies have documented a variety of problems that may affect the proportions of DNA reads obtained with HTS (Amend, Seifert, & Bruns, 2010; Deagle, Thomas, Shaffer, Trites, & Jarman, 2013; Pawlowski

et al., 2016; Tan et al., 2015; Thomas, Deagle, Eveson, Harsch, & Trites, 2016), including biological biases (e.g. gene copy number variation, tissue cell density, cell biovolume), technical biases (e.g. DNA extraction, PCR amplification) and biases linked to HTS itself (e.g. library construction, HTS technology used, bioinformatics treatments). Variation in gene copy number per cell constitutes a major bias known to affect the proportion of DNA read found for each species present in complex assemblages; this has been demonstrated for macroinvertebrates (Elbrecht, Peinert, & Leese, 2017), fish, amphibians (Evans et al., 2016), oligochaetes (Vivien, Lejzerowicz, & Pawlowski, 2016), foraminifera (Weber & Pawlowski, 2013) and microbial assemblages (Angly et al., 2014). However, to the best of our knowledge, no study has yet evaluated gene copy number variation bias on diatom metabarcoding quantification. While tissue cell density and species biomass are major biases likely to affect DNA metabarcoding quantification of multicellular organisms like macroinvertebrates (Elbrecht & Leese, 2015) or fish (Evans et al., 2016), diatoms are unicellular organisms for which gene copy number is mainly affected by the number of genomes and the number of gene copies per genome. This may be particularly true for non-nuclear markers like the chloroplast-encoded *rbcl* gene. Godhe et al. (2008) reported a clear correlation between the 18S gene copy number per cell with diatom cell length and biovolume, suggesting that the cell biovolume could be a proxy for the gene copy number. Keeping in mind that diatom biovolume varies from 10^1 to $10^9 \mu\text{m}^3$ (Snoeijs, Busse, & Potapova, 2002), gene copy number may vary greatly between the smallest and the biggest diatom species, affecting metabarcoding quantification.

For all the reasons mentioned above, we hypothesized that a quantification correction factor (CF) based on diatom cell biovolume should be necessary to correct DNA read proportions to provide species quantification more comparable to microscopical counts. In order to confirm this hypothesis, we firstly conducted experiments on eight pure diatom cultures to examine whether variation in *rbcl* gene copy number per cell correlates with morphological characteristics (e.g. biovolume, cell length), from which a CF might be calculated. Secondly, the efficiency of the proposed CF was tested on (1) mock communities made by mixing known proportions of the eight diatom species cultures and (2) environmental diatom assemblages from rivers previously sequenced (Vasselon, Rimet, et al., 2017) and for which data are available online (Vasselon, Rimet, et al., 2017 dataset, <https://doi.org/10.5281/zenodo.400160>). Last, the capacity of the CF to improve the ecological assessment of rivers was tested by comparing water quality index values calculated from molecular data with corrected abundances to those calculated from classical morphological abundances.

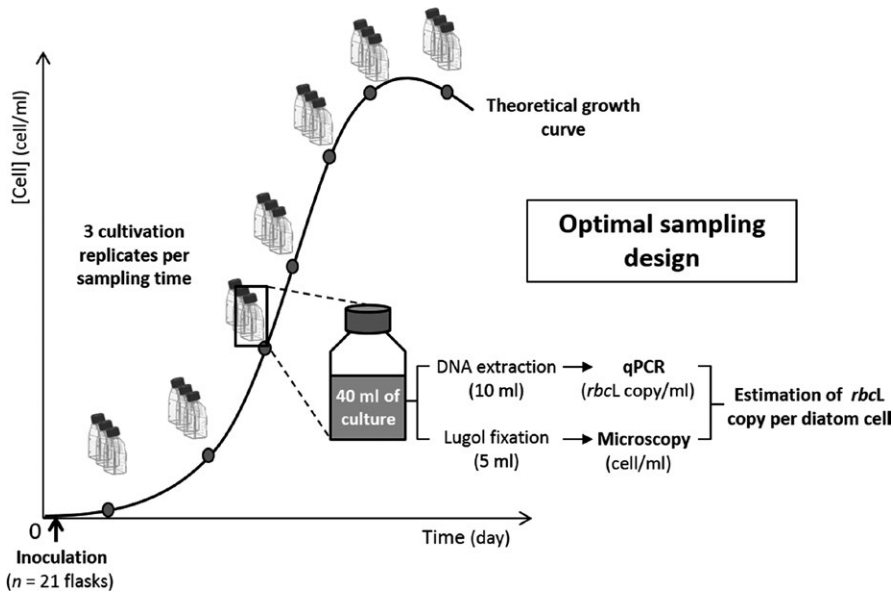
2 | MATERIALS AND METHODS

2.1 | Evaluation of the quantification bias and development of a quantification correction factor (CF)

To evaluate whether the *rbcl* copy number per cell varies between diatom species, strains from eight freshwater diatom

TABLE 1 Characteristics of the eight diatom species selected in the Thonon Culture Collection (TCC) and used in this study

Species	TCC code	Chloroplast (nb./cell)	Length (μm)	Width (μm)	Thickness (μm)	Biovolume (μm^3)
<i>Achnanthydium minutissimum</i> (Kützing) Czarnecki	TCC667	1	7.1	3.2	2.5	45
<i>Nitzschia palea</i> (Kützing) W.Smith	TCC139-1	2	22.7	4.0	4.0	183
<i>Ulnaria ulna</i> (Nitzsch) Compère	TCC670	2	54.6	7.9	9.5	4,087
<i>Pinnularia viridiformis</i> (Nitzsch) Ehrenberg	TCC890	2	51.4	14.3	17.8	10,282
<i>Diatoma tenuis</i> Kützing	TCC861	≈ 8	42.4	4.8	4.8	769
<i>Nitzschia inconspicua</i> Grunow	TCC488	2	8.1	4.3	3.6	98
<i>Fragilaria perminuta</i> (Grunow) Lange-Bertalot	TCC753	2	11.1	4.2	3.7	135
<i>Cyclotella meneghiniana</i> Kützing	TCC690	≈ 20	12.1		4.7	539

**FIGURE 1** Experimental design applied to the eight diatom species. After the inoculation of 21 flasks containing 40 ml of DV media, diatom culture growth was followed at seven sampling times (from T0 to T6) and analysis was performed in triplicate (3 flasks per sampling time)

species were selected from the Thonon Culture Collection (TCC; http://www6.inra.fr/carrtel-collection_eng/) (Table 1). The eight species were chosen for their contrasted morphological (size and cell biovolume), cytological (e.g. chloroplast number) and phylogenetic characteristics (Table 1). Cell dimensions (width, length, thickness) of the eight diatom species were measured under light microscopy ($1,000\times$ magnification) using a minimum of 10 specimens per species. Then, appropriate geometrical models were applied to calculate their cell biovolume (Sun & Liu, 2003) (Table 1). The eight diatom cultures were cultivated in triplicate in 40 ml sterile DV medium (Rimet et al., 2014) using 50 ml Nunc™ EasYFlasks™ (Thermo Fisher Scientific, Waltham, Massachusetts). Flasks were placed on a rotating platter (4 rpm) in a controlled thermostatic room ($21 \pm 2^\circ\text{C}$, 14 hr light/10 hr dark cycle, light intensity of c. $100 \mu\text{mol quanta m}^{-2} \text{s}^{-1}$). Flasks were inoculated in order to reach a concentration of ≈ 100 cells/ml at the beginning of the experiment for each species, except for *Ulnaria ulna* for which a concentration of $\approx 1,000$ cells/ml was used (due to its low growth rate). The growth of the eight diatom cultures was followed during 40 days, except for *Pinnularia*

viridiformis for which the survey lasted 73 days, due to its low growth rate. Cell concentrations, proportions of live/dead cells and *rbcL* gene copy concentrations per ml of media were measured for each culture at seven sampling times (referred to as T0 to T6) (Figure 1).

Diatom cell concentrations and proportions of live/dead cells were obtained by counting at least 400 specimens using inverted microscopy ($\times 1,000$ magnification) and the standard Utermöhl technique (European Committee for Standardization (CEN) 2006) (Figure 1). The proportion of live/dead cells was estimated by considering cells without visible intracellular contents as dead. Only living cells were taken into account to calculate the diatom cell concentration per ml of media. Flow cytometry using Sytox-Green was also used to confirm the microscopical data (not shown).

RbcL copy number per ml was estimated by qPCR. From each cultivation replicate, 10 ml of culture was centrifuged at $17,000 \times g$ for 30 min (Figure 1). Total DNA was extracted from the resulting pellet using a protocol based on GenElute™-LPA DNA precipitation (Sigma-Aldrich, St Louis, Missouri) as previously described (Vasselon, Domaizon, et al., 2017). Then, qPCR assays were performed for each

of the eight diatom species on DNA extracted at all seven sampling times and with each of the three replicates, using the QuantiTect SYBR Green PCR Kit (Life Technologies, Carlsbad, USA) and the Rotor-Gene Q (Qiagen, Hilden, Germany). A short 312 bp region of the *rbcl* gene (the same as was used for HTS sequencing) was targeted using primers used by (Vasselon, Rimet, et al., 2017) and described in Table S1. qPCR reactions were performed following the method used by Vasselon, Domaizon, et al. (2017), using a final volume of 25 μ l using mix preparation and reaction conditions as described in Table S1. A fluorescence threshold of 0.01 was used to allow comparison of qPCR assays, denoising and determination of the cycles' threshold (Ct). Data analysis was performed using the ROTOR-GENE Q Series software (version 2.3.1) and the *rbcl* copy per ml of media was determined.

Finally, the number of *rbcl* gene copies per diatom cell was calculated for the eight diatom species by dividing the *rbcl* concentration (qPCR data) by the living cell concentration (microscopy data). A Kruskal–Wallis test was performed using R (R Development Core Team, 2013) to determine if the *rbcl* gene copy number per diatom cell varied significantly between the eight diatom species. Then, we tested the level of correlation between the number of *rbcl* gene copies per diatom cell and several morphological characteristics of the diatom cells (Table 1). Variables that did not approximate normal distributions were log transformed. Pearson correlation coefficients were calculated between the gene copy number per cell and the diatom cell morphological characteristics. This correlation was represented by a linear model.

2.2 | Validation in the quantification CF to mock and environmental HTS data

2.2.1 | Mock communities

The calculated CF was applied to metabarcoding data obtained from controlled diatom mock communities. Five mock communities (M1 to M5) were created by mixing DNA extracted from each of the eight diatom species sampled during their exponential growth phase, and for which the correspondence between cell abundances (microscopy) and qPCR counts was known. For each of the five mock communities, the volume of DNA used for seven species was kept unchanged (1 μ l) and only the volume of DNA of *P. viridiformis* varied as followed: M1 = 0.2 μ l, M2 = 0.4 μ l, M3 = 0.8 μ l, M4 = 1.6 μ l, M5 = 3.2 μ l. This resulted in contrasted *rbcl* proportions of the eight species among the five mock communities. Then, HTS sequencing of the *rbcl* 312 bp fragment was performed on three replicates of the five mock communities. The 15 corresponding libraries were prepared following the method described by Vasselon, Domaizon, et al. (2017) with the same primers and PCR reaction conditions as those used for *rbcl* qPCR (Table S1), changing only the cycle number to 30. Each library was diluted to 100 pm and all 15 were pooled together for one HTS run performed on the PGM Ion Torrent machine by the "Plateforme Génome Transcriptome" (PGTB, Bordeaux, France).

The sequencing platform provided a unique fastq file for each of the 15 libraries containing demultiplexed DNA reads without the sequencing adapters. Quality filtering of DNA reads was performed using

the MOTHUR software (Schloss et al., 2009) and bioinformatics process described previously (Vasselon, Domaizon, et al., 2017; Vasselon, Rimet, et al., 2017). Finally, a taxonomy was assigned to each DNA read with the "classify.seqs" command (Mothur) using default parameters with a confidence threshold of 85% and the R-Syst::diatom library (Rimet et al., 2016, version updated in January 2015 and available upon request) as a *rbcl* reference library. A molecular taxonomic list with the associated read numbers assigned to each of the eight diatom species was obtained for each of the five mock communities and used for subsequent analysis.

The quantification CF defined for the *rbcl* gene was then applied to the molecular taxonomic lists for the five mock communities by dividing the read number for each species by its corresponding CF. Both the uncorrected and corrected HTS relative abundances of species from the five mock communities were then compared to the relative abundances obtained using microscopy.

2.2.2 | Environmental diatom assemblages

To evaluate the efficiency of the CF to improve metabarcoding quantification from environmental samples, we used *rbcl* HTS data obtained from (Vasselon, Rimet, et al., 2017), corresponding to 80 benthic diatom samples collected from rivers in tropical island of Mayotte, Indian Ocean (Vasselon, Rimet, et al., 2017 dataset, <https://doi.org/10.5281/zenodo.400160>). A CF was calculated for each species (or genus when the species level was not reached) detected in molecular inventories of the rivers of Mayotte island using a generalized average of the morphological information (e.g. biovolume, length) available in the R-Syst::diatom library and applied to HTS data. Corrected molecular inventories were produced for all the 80 river samples using the CF. The impact of the CF on diatom taxa abundance rank in the molecular inventories was assessed by comparing original and corrected molecular diatom inventories. Then, the Specific Pollution-sensitivity Index (SPI) used for ecological assessment was calculated for each sample based on the corrected diatom molecular inventories using the OMNIDIA 5 software (Lecointe, Coste, & Prygiel, 1993; library 5.3 2015) and compared to the morphological SPI values for all river samples (Vasselon, Rimet, et al., 2017). Pearson correlation was used to evaluate the strength of correlations between original or corrected molecular SPI values and the morphological SPI values. Wilcoxon Signed Rank tests were conducted to determine whether the difference between the molecular and the morphological SPI (Δ SPI) varied significantly when using the original or the corrected molecular data for the molecular SPI calculation.

3 | RESULTS

3.1 | Variation in *rbcl* gene copy number between diatom species

Cell and *rbcl* gene concentrations were measured, by inverted microscopy and qPCR, respectively, for the eight diatom species at different cultivation stages corresponding to seven sampling points (T0 to T6).

Information has been summarized in Tables S2 and S3. As the eight diatom species reached the beginning of the stationary phase at the sampling time T2 (i.e. between 13 and 31 days of cultivation), only the [cell] and the [gene copy] values obtained for the T0, T1 and T2 sampling times were used for further analysis. The calculated mean values of the *rbcl* gene copy number per cell for each diatom species varied between 0.5 and 130 copies per cell (Figure 2). The Kruskal–Wallis test revealed that the *rbcl* copy number per cell was significantly different ($p < .001$) between the eight diatom species.

3.2 | Development of quantification CFs

The *rbcl* copy number per cell was highly correlated with cell biovolume ($r = .97$, $p < .001$), length ($r = .82$, $p < .001$), width ($r = .94$, $p < .001$) and thickness ($r = .96$, $p < .001$). The correlation between the *rbcl* copy number per cell and the cell biovolume followed a linear model (Figure 3). Assuming that this linear relation based on eight diatom species is applicable to all diatom species, the equation of this model allows calculation of an estimate of the relative *rbcl* copy number per cell as soon as the biovolume of the cell is known, and thus to define a CF specific to each species. Such quantification CFs were calculated for each of the eight diatom species of the mock communities (Table 2) and varied from 0.6 for *Achnantheidium minutissimum* to 78.5 for *P. viridiformis*. For each of the diatom taxa found in the environmental samples, CFs were also calculated using the biovolume information available for each taxa (from Rsystem::diatom library) (Table S4) and varied over a wider range, from 0.03 for *Fistulifera saprophila* to 649.8 for *Rhopalodia gibba*.

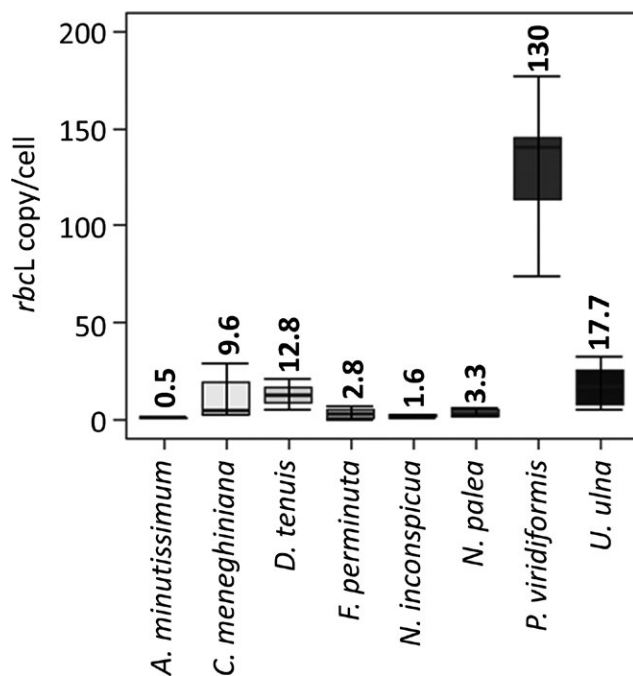


FIGURE 2 Estimation of the *rbcl* copy number per diatom cell for the eight diatom species. Mean values calculated using the gene and the diatom cell concentrations obtained, respectively, by qPCR and inverted microscopy at T0, T1 and T2 sampling points ($n = 9$)

3.3 | Application of CFs to mock and environmental HTS data

953,082 DNA reads were produced from the 15 libraries corresponding to the five DNA mock communities (3 replicates per mock). Following the bioinformatics quality filtering steps, 385,367 DNA reads were retained. A molecular taxonomic list was then created by removing DNA reads which remained unclassified (0.43% of the reads) or assigned to different taxa than the eight diatom species present in the mock communities (0.004% of the reads) (Table S5). The proportions of *P. viridiformis* reads in the five mock communities varied from 9% in M1 to 57% in M5 (Figure 4a) while observed cell proportions were lower; $\approx 0.03\%$ in M1 and 0.55% in M5 (Figure 4b). The application of the CF on DNA reads counts of the eight species changed their relative abundances in the five mock communities (Figure 4a). The rank of the eight species was also affected; for example in M5 the application of the CF changed the proportion of *P. viridiformis* from 57% to 4% and the proportion of *A. minutissimum* from 4% to 42%. The correspondence between morphological and molecular relative abundances was highly improved by applying the CF on the HTS data (Figure 4a,b).

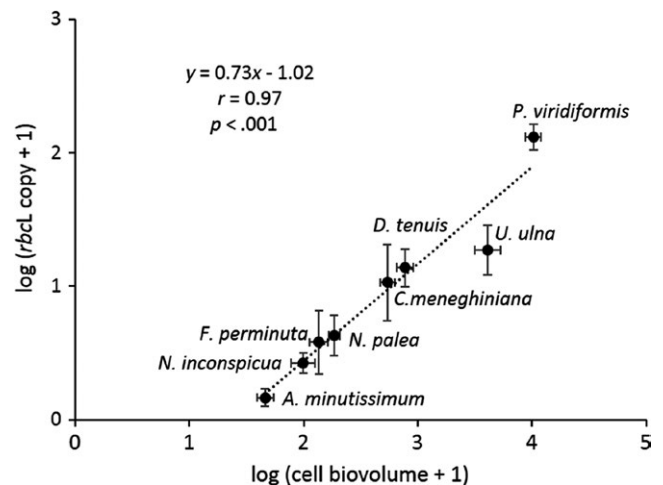
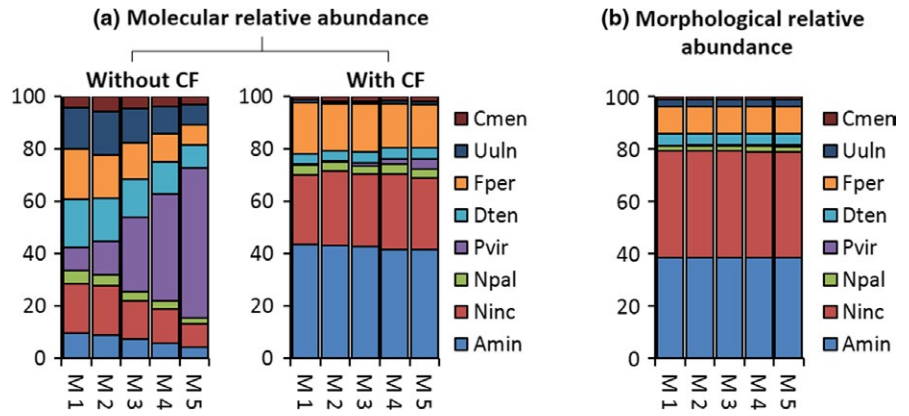


FIGURE 3 Correlation between the diatom cell biovolume and the *rbcl* gene copy number per cell after $\log(x + 1)$ transformation

TABLE 2 CF calculated for the eight diatom species using their respective cell biovolume (Table 1) and the linear equation between the *rbcl* copy number and the cell biovolume (Figure 3)

Species	Calculated CF
<i>Achnantheidium minutissimum</i>	0.6
<i>Nitzschia inconspicua</i>	1.7
<i>Nitzschia palea</i>	3.3
<i>Pinnularia viridiformis</i>	78.5
<i>Diatoma tenuis</i>	11.1
<i>Fragilaria perminuta</i>	2.4
<i>Ulnaria ulna</i>	39.6
<i>Cyclotella meneghiniana</i>	8.3

FIGURE 4 Relative abundances of the eight diatom species in the five DNA mock communities based (a) on mean of HTS DNA reads without (left) and with (right) correcting quantification using the biovolume correction factor and (b) on mean of morphological counts from inverted microscopy



From the 80 environmental samples previously sequenced (Vasselon, Rimet, et al., 2017), a molecular taxonomic list based on assigned DNA reads was produced including 23 families (75.1% of total reads assigned), 39 genera (72% of total reads assigned) and 66 diatom taxa, including taxa assigned at the genus and the species level, were used to calculate the SPI freshwater quality index. CFs calculated from cell biovolumes for those 84 taxa were then applied to correct the quantification of the environmental molecular inventories (Table S4). The proportions and ranks of the dominant taxa were affected by the application of the CFs (Figure 5). For example the application of CFs reduced the relative abundances of *Eunotia* and *Ulnaria* from 31.9% to 3.3% and 11.7% to 2.3%, respectively, making them more congruent with cell proportions observed with microscopy (3.1% for *Eunotia* and 0.4% for *Ulnaria*). The correlation between the morphological and the molecular SPI values for all river samples previously described ($r = .72$, $p < .001$) was slightly improved using SPI values based on inventories with corrected abundances ($r = .77$, $p < .001$). The application of the CF to correct the HTS quantification reduced significantly ($p < .001$) the differences between the molecular and morphological SPI values by 47% (Δ SPI reduced to 1.9 on average compared to 3.6 before correction, corresponding to 37.3% and 21.2% of error respectively) (Figure 6).

4 | DISCUSSION

Species quantification based on DNA metabarcoding is challenging for most of taxonomic groups as technical and biological biases affect DNA reads proportions. In order to limit those biases, several attempts were done to apply a CF on metabarcoding data, as shown for fishes (Thomas et al., 2016), bacteria and archaea (Angly et al., 2014) or oligochaetes (Vivien et al., 2016). For those studies, application of the CF, whether for correcting single (Angly et al., 2014) or multiple sources of quantification biases (Thomas et al., 2016), improved taxa quantification from metabarcoding data compare to morphological one. The result is generally a change in the ranks of the dominant taxa which affect directly the community structure and can lead to different ecological interpretations. For example the application of a CF on

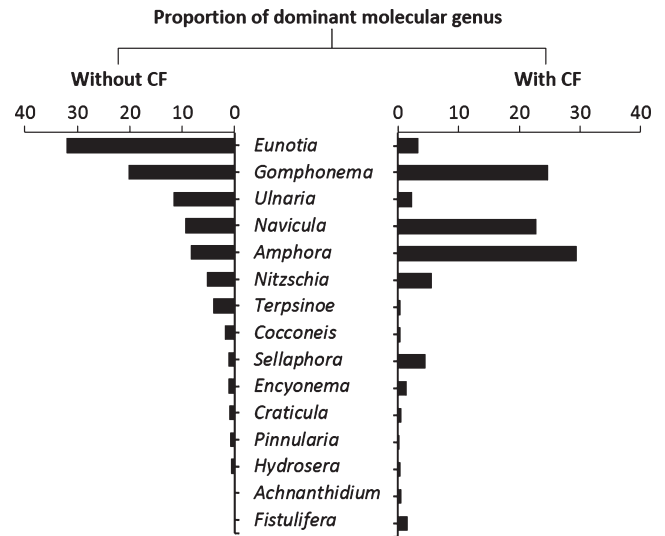


FIGURE 5 Dominant taxa (relative abundance > 0.5%) obtained in HTS Mayotte molecular inventories without (left) and with (right) application of the biovolume correction factor. All samples ($n = 80$) are considered

metabarcoding data obtained from aquatic oligochaetes samples improved the freshwater quality assessment based on molecular index calculation (Vivien et al., 2016). However, the development of CF can be challenging depending on the organism studied, as it requires finding a clear relationship between DNA reads and specimen proportions. This may be impossible due to accumulation of quantification biases (e.g. cell density, cell biomass, gene copy number). Nevertheless, the use of CF can be advantageous for organisms with a high variation in the DNA reads proportions between taxa (e.g. several log) and where a limited number of biases are involved like diatoms.

4.1 | Correlation between *rbcl* gene copy number and diatom cell biovolume: Impacts on HTS quantification

The copy number of the *rbcl* gene present in one diatom cell is affected by three parameters: (1) the number of chloroplasts per cell, (2) the number of genomes per chloroplast and (3) the number of

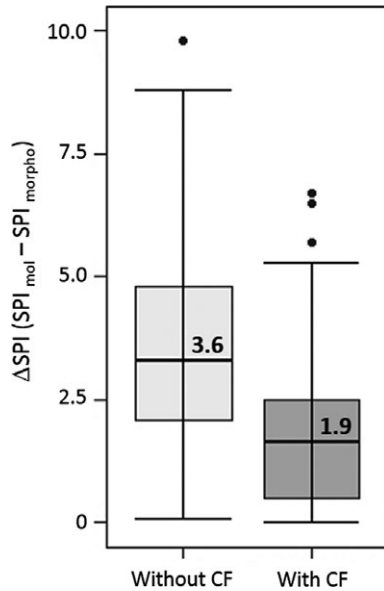


FIGURE 6 Distribution of the differences between the molecular and the morphological SPI (Δ SPI) for all Mayotte samples using original molecular SPI values (left) and new molecular SPI values based on molecular inventories corrected with the biovolume CF (right)

copies of the *rbcL* gene per chloroplast genome (Erslund, Aldrich, & Cattolico, 1981; Treusch et al., 2012). (1) For benthic diatoms, the chloroplast number per cell is quite stable inside a single genus with variations ranging from 1 to ≈ 8 chloroplast(s) per cell from a genus to another (Round, Crawford, & Mann, 1990), even if some centric genera may have tens of chloroplasts (e.g. *Melosira*, *Cyclotella*). (2) Regarding the chloroplast genome number per cell, higher plants can contain up to thousands of copies of chloroplast genome per cell (Bendich, 1987; Rauwolf, Golczyk, Greiner, & Herrmann, 2010) while unicellular algae generally exhibit a lower number of copies. For example *Olisthodiscus luteus* (Raphidophyceae), *Chlamydomonas reinhardtii* (Chlorophyceae), *Phaeodactylum tricornutum* (pennate diatom) and *Thalassiosira pseudonana* (centric diatom) contain, respectively, around 650, 80, 137 and 55 genome copies per cell (von Dassow, Petersen, Chepurnov, & Virginia Armbrust, 2008; Erslund et al., 1981; Gruber, 2008; Koop, Herz, Golds, & Nickelsen, 2007). (3) Finally, there is only one copy of the *rbcL* gene per chloroplast genome (e.g. Sabir et al., 2014), as in higher plants (Gutteridge & Gatenby, 1995).

Thus, the *rbcL* copy number may vary from tens to hundreds of copies per diatom cell. Our estimations are within this range with a maximum of 130 copies estimated for *P. viridiformis*. However, our method underestimates the *rbcL* gene copy number since 0.5 copy per cell was estimated for *A. minutissimum* (so implying that some cells have no *rbcL* copy). This may result from certain variability inherent to the estimation of gene copy number by qPCR and the quantification of cells by microscopical counts. Our results demonstrate, however, that the *rbcL* copy number varies significantly between the eight diatom species used in this study, according to the different diatom cell characteristics tested. In particular, we found a significant linear

relationship between the *rbcL* copy number and the cell biovolume. Although the size of the chloroplasts could not be estimated in this study, we assume that the increase in the cell biovolume is accompanied by an increase in the chloroplast biovolume (as shown by Okie, Smith, & Martin-Cereceda, 2016), inducing an increase in DNA quantity and chloroplast genome copies per chloroplast as shown by Rauwolf et al. (2010).

The correlation we found between the *rbcL* copy number and the diatom cell biovolume suggests that the relative abundance of diatom species with high cell biovolume is likely to be over-represented in metabarcoding data compared to microscopical counts. This is confirmed by the HTS data obtained for the mock communities, where diatom species with high cell biovolume are over-represented (e.g. *P. viridiformis*) and diatom species with low cell biovolume are under-represented (e.g. *A. minutissimum*). The relative abundance of *P. viridiformis* in the mock communities was negligible compared to other species, and doubling its proportion did not change its rank: the species remained the least abundant taxon within the morphological inventory. However, due to its high cell biovolume ($10^4 \mu\text{m}^3$) and relatively high *rbcL* copy number per cell, a marked over-representation of this species within the molecular inventory was observed. A CF was thus defined to correct these quantitative biases and was verified on mock communities and environmental samples.

4.2 | Current potential and limits of the quantification CF

The use of the same *rbcL* primers for the qPCR assays and the HTS enabled us to generate a specific CF well suited to correct *rbcL* metabarcoding quantifications. Its application to the HTS data of the mock communities allowed us to obtain comparable species proportions in morphological and molecular based approaches of mock communities. This was also confirmed with the Mayotte river samples, for which the quantification CF resulted in a better congruence between DNA reads and cells proportions, reducing the over-representation of high biovolume *Eunotia* and *Ulnaria* species. Furthermore, SPI calculation based on corrected metabarcoding data gives SPI values more comparable to SPI values obtained from morphological data, suggesting that it may be possible to replace morphological by molecular monitoring for the ecological assessment of Mayotte rivers. In the same way, (Vivien et al., 2016) have shown that application of a CF to correct DNA reads proportions allows a more accurate estimation of oligochaete proportions, improving quality index calculation and quality assessment of watercourse sediments. Our results confirm that water quality index based on diatom metabarcoding and DNA read proportions are directly affected by gene copy number variation, and show the potential value of integrating CFs into molecular SPI calculation. However, as the biovolume-copy number relationship was based on only eight diatom species and the efficiency of the resulting CFs validated on only one HTS dataset, further experiments including more species and larger datasets will be required to develop and fully validate CFs for use in molecular biomonitoring.

The CF developed in this study assumes that gene copy number is constant in each taxon. However, gene copy number may vary with the physiological status of the cell and stage of the life cycle, since in most diatoms cell volume decreases during the vegetative phase. The physiological status varies with cell cycle progression; additionally several factors may affect the physiological status of diatoms like changes in environmental conditions (e.g. nutrient availability, pollutants, temperature, etc.) (Pandey et al., 2017). Altered physiological status of a given population is generally characterized by a higher proportion of damaged cells. The compromised/damaged cells are characterized by alteration of membrane integrity, degradation of the photosynthetic pigments or fragmentation of genomic DNA (Zetsche & Meysman, 2012; Znachor, Rychtecký, Nedoma, & Visocká, 2015). Variations in DNA integrity and chloroplast physiology between cells of a given population can impact directly the *rbcl* gene copy number per cell and thus DNA metabarcoding quantification. (Eberhard, Drapier, & Wollman, 2002) showed that chloroplast genome copy number is reduced when the green alga *Chlamydomonas reinhardtii* is cultivated under phototrophic conditions compared to cultivation in mixotrophic conditions. Limitation by mineral nutrients may also have an impact; for instance iron limitation can reduce the number of the chloroplast per cell (from 4 to 2) and their size in the marine diatom *Thalassiosira oceanica* (Hustedt) Hasle et Heimdal (Lommer et al., 2012). Variation in the cell physiological state was not taken into account in developing CFs for diatom metabarcoding. However, during our experiments we discriminated live and dead cells; we observed that their respective proportions did not affect significantly the correlation between the gene copy number per cell and the cell biovolume (Figure S1). Further experiments should be performed to evaluate the impact on the final CFs of *rbcl* gene copy number variation linked to physiological status.

The biovolume of each diatom species is required to apply the CF and hence correct the quantification in metabarcoding datasets. Several reference databases provide biovolume information for a lot diatom species (e.g. Rimet et al., 2016), but they do not generally account for biovolume variability, which is a complicating factor in diatoms because of the peculiarities of the life cycle. Diatom cell size within a population is not constant due to the method of vegetative reproduction, which leads to a progressive cell size reduction in the population (Crawford, 1981), followed by restoration of cell size via a sexual event. For this reason, different cell sizes can be observed in the same diatom population, either in pure cultures (e.g. in the marine diatom *Thalassiosira weissflogii* Grunow: Armbrust & Chisholm, 1992) or in environmental populations (e.g. the freshwater species *Sellaphora pupula* (Kützting) Mereschk.: (Mann, Chepurnov, & Droop, 1999)). However, although the range of cell sizes within a given diatom population may vary by a factor of 2 to 5 in the environment (Hense & Beckmann, 2015), natural populations usually have a rather narrow range of sizes and larger cells form a negligible fraction of the population (Mann, 2011). Furthermore, the distribution of cell size within environmental populations is often close to being normal (Mann et al., 1999; Spaulding, Jewson, Bixby, Nelson, & McKnight, 2012). The balance between small and big individuals in the same population will

therefore limit errors associated with the use of a mean biovolume. Hence, we propose to use the mean of biovolume to calculate CFs; without considering other potential HTS quantification biases, its application to DNA reads of environmental material should allow a good correction of their proportions.

ACKNOWLEDGEMENTS

The authors declare no conflict of interest. Funding provided by the French National Agency for Water and Aquatic Environments (ONEMA-AFB) and supported by the European COST action DNAqua-Net (CA 15219). A special thanks to David G. Mann for the constructive discussions that helped to improve the manuscript.

AUTHORS' CONTRIBUTIONS

V.V., A.B., F.R., S.J., M.C., K.T. and I.D. contributed to the study designed. V.V., M.C. and S.J. conducted the laboratory work. V.V. analysed the data and wrote the manuscript. All the authors contributed to the discussions and to manuscript editing.

DATA ACCESSIBILITY

All PGM raw sequence data are available for the 15 libraries, corresponding to the five DNA mock communities with three replicates, on the Zenodo repository website (<https://doi.org/10.5281/zenodo.807178>).

ORCID

Valentin Vasselon  <http://orcid.org/0000-0001-5038-7918>

REFERENCES

- Amend, A. S., Seifert, K. A., & Bruns, T. D. (2010). Quantifying microbial communities with 454 pyrosequencing: Does read abundance count? *Molecular Ecology*, *19*, 5555–5565. <https://doi.org/10.1111/j.1365-294X.2010.04898.x>
- Angly, F. E., Dennis, P. G., Skarshewski, A., Vanwongerghem, I., Hugenholtz, P., & Tyson, G. W. (2014). CopyRighter: A rapid tool for improving the accuracy of microbial community profiles through lineage-specific gene copy number correction. *Microbiome*, *2*, 11. <https://doi.org/10.1186/2049-2618-2-11>
- Apothéloz-Perret-Gentil, L., Cordonier, A., Straub, F., Iseli, J., Esling, P., & Pawlowski, J. (2017). Taxonomy-free molecular diatom index for high-throughput eDNA biomonitoring. *Molecular Ecology Resources*, *17*, 1231–1242. <https://doi.org/10.1111/1755-0998.12668>
- Armbrust, E. V., & Chisholm, S. W. (1992). Patterns of cell size change in a marine centric diatom: Variability evolving from clonal isolates. *Journal of Phycology*, *28*, 146–156. <https://doi.org/10.1111/j.0022-3646.1992.00146.x>
- Bendich, A. J. (1987). Why do chloroplasts and mitochondria contain so many copies of their genome? *BioEssays*, *6*, 279–282. [https://doi.org/10.1002/\(ISSN\)1521-1878](https://doi.org/10.1002/(ISSN)1521-1878)
- Council, E. (2000). *Directive 2000/60/EC of the European parliament and of the council of 23 October 2000 establishing a framework for community*

- action in the field of water policy. Brussels: Office for Official Publications of the European Communities.
- Crawford, R. M. (1981). The siliceous components of the diatom cell wall and their morphological variation. In T. L. Simpson & B. E. Volcani (Eds.), *Silicon and siliceous structures in biological systems* (pp. 129–156). New York, NY: Springer New York. <https://doi.org/10.1007/978-1-4612-5944-2>
- von Dassow, P., Petersen, T. W., Chepurinov, V. A., & Virginia Armbrust, E. (2008). Inter- and intraspecific relationships between nuclear DNA content and cell size in selected members of the centric diatom genus *Thalassiosira* (Bacillariophyceae). *Journal of Phycology*, 44, 335–349. <https://doi.org/10.1111/j.1529-8817.2008.00476.x>
- Deagle, B. E., Thomas, A. C., Shaffer, A. K., Trites, A. W., & Jarman, S. N. (2013). Quantifying sequence proportions in a DNA-based diet study using Ion Torrent amplicon sequencing: Which counts count? *Molecular Ecology Resources*, 13, 620–633. <https://doi.org/10.1111/1755-0998.12103>
- Eberhard, S., Drapier, D., & Wollman, F.-A. (2002). Searching limiting steps in the expression of chloroplast-encoded proteins: Relations between gene copy number, transcription, transcript abundance and translation rate in the chloroplast of *Chlamydomonas reinhardtii*. *The Plant Journal*, 31, 149–160. <https://doi.org/10.1046/j.1365-313X.2002.01340.x>
- Elbrecht, V., & Leese, F. (2015). Can DNA-based ecosystem assessments quantify species abundance? Testing primer bias and biomass – Sequence relationships with an innovative metabarcoding protocol. *PLoS ONE*, 10, e0130324. <https://doi.org/10.1371/journal.pone.0130324>
- Elbrecht, V., Peinert, B., & Leese, F. (2017). Sorting things out: Assessing effects of unequal specimen biomass on DNA metabarcoding. *Ecology and Evolution*, 7, 6918–6926. <https://doi.org/10.1002/ece3.3192>
- Ersland, D. R., Aldrich, J., & Cattolico, R. A. (1981). Kinetic complexity, homogeneity, and copy number of chloroplast DNA from the marine alga *Olisthodiscus luteus*. *Plant Physiology*, 68, 1468–1473. <https://doi.org/10.1104/pp.68.6.1468>
- European Committee for Standardization (CEN). (2006). EN 15204 – Water quality – Guidance standard on the enumeration of phytoplankton using inverted microscopy (Utermöhl technique). *European Standard*, 1–42.
- Evans, N. T., Olds, B. P., Renshaw, M. A., Turner, C. R., Li, Y., Jerde, C. L., ... Lodge, D. M. (2016). Quantification of mesocosm fish and amphibian species diversity via environmental DNA metabarcoding. *Molecular Ecology Resources*, 16, 29–41. <https://doi.org/10.1111/1755-0998.12433>
- Godhe, A., Asplund, M. E., Härnström, K., Saravanan, V., Tyagi, A., & Karunasagar, I. (2008). Quantification of diatom and dinoflagellate biomasses in coastal marine seawater samples by real-time PCR. *Applied and Environmental Microbiology*, 74, 7174–7182. <https://doi.org/10.1128/AEM.01298-08>
- Gruber, A. (2008). Molecular characterisation of diatom plastids (PhD thesis). University of Konstanz.
- Gutteridge, S., & Gatenby, A. (1995). Rubisco synthesis, assembly, mechanism, and regulation. *The Plant Cell Online*, 7, 809–819. <https://doi.org/10.1105/tpc.7.7.809>
- Hense, I., & Beckmann, A. (2015). A theoretical investigation of the diatom cell size reduction–restitution cycle. *Ecological Modelling*, 317, 66–82. <https://doi.org/10.1016/j.ecolmodel.2015.09.003>
- Keck, F., Vasselton, V., Tapolczai, K., Rimet, F., & Bouchez, A. (2017). Freshwater biomonitoring in the information age. *Frontiers in Ecology and the Environment*, 15, 266–274. <https://doi.org/10.1002/fee.1490>
- Kermarrec, L., Franc, A., Rimet, F., Chaumeil, P., Frigerio, J.-M., Humbert, J., & Bouchez, A. (2014). A next-generation sequencing approach to river biomonitoring using benthic diatoms. *Freshwater Science*, 33, 349–363. <https://doi.org/10.1086/675079>
- Kermarrec, L., Franc, A., Rimet, F., Chaumeil, P., Humbert, J. F., & Bouchez, A. (2013). Next-generation sequencing to inventory taxonomic diversity in eukaryotic communities: A test for freshwater diatoms. *Molecular Ecology Resources*, 13, 607–619. <https://doi.org/10.1111/1755-0998.12105>
- Koop, H.-U., Herz, S., Golds, T. J., & Nickelsen, J. (2007). The genetic transformation of plastids. In R. Bock (Ed.) *Cell and molecular biology of plastids. Topics in current genetics* (Vol. 19). Berlin, Heidelberg: Springer.
- Lecointe, C., Coste, M., & Prygiel, J. (1993). 'OMNIDIA': Software for taxonomy, calculation of diatom indices and inventories management. *Hydrobiologia*, 269–270, 509–513. <https://doi.org/10.1007/BF00028048>
- Lommer, M., Specht, M., Roy, A.-S., Kraemer, L., Andreson, R., Gutowska, M. A., ... LaRoche, J. (2012). Genome and low-iron response of an oceanic diatom adapted to chronic iron limitation. *Genome biology*, 13, R66. <https://doi.org/10.1186/gb-2012-13-7-r66>
- Mann, D. G. (2011). Size and sex. In E. J. Seckbach, & J. P. Kocielek (Eds.), *The diatom world* (pp. 145–166). Dordrecht: Springer. <https://doi.org/10.1007/978-94-007-1327-7>
- Mann, D. G., Chepurinov, V. A., & Droop, S. J. M. (1999). Sexuality, incompatibility, size variation, and preferential polyandry in natural populations and clones of *Sellaphora pupula* (Bacillariophyceae). *Journal of Phycology*, 35, 152–170. <https://doi.org/10.1046/j.1529-8817.1999.3510152.x>
- Okie, J. G., Smith, V. H., & Martin-Cereceda, M. (2016). Major evolutionary transitions of life, metabolic scaling and the number and size of mitochondria and chloroplasts. *Proceedings Biological Sciences*, 283, 20160611. <https://doi.org/10.1098/rspb.2016.0611>
- Pandey, L. K., Bergey, E. A., Lyu, J., Park, J., Choi, S., Lee, H., ... Han, T. (2017). The use of diatoms in ecotoxicology and bioassessment: Insights, advances and challenges. *Water Research*, 118, 39–58. <https://doi.org/10.1016/j.watres.2017.01.062>
- Pawlowski, J., Lejzerowicz, F., Apotheloz-Perret-Gentil, L., Visco, J., & Esling, P. (2016). Protist metabarcoding and environmental biomonitoring: Time for change. *European Journal of Protistology*, 55, 12–25. <https://doi.org/10.1016/j.ejop.2016.02.003>
- R Development core team (2013). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Rauwolf, U., Golczyk, H., Greiner, S., & Herrmann, R. G. (2010). Variable amounts of DNA related to the size of chloroplasts III. Biochemical determinations of DNA amounts per organelle. *Molecular Genetics and Genomics*, 283, 35–47. <https://doi.org/10.1007/s00438-009-0491-1>
- Rimet, F. (2012). Recent views on river pollution and diatoms. *Hydrobiologia*, 683, 1–24. <https://doi.org/10.1007/s10750-011-0949-0>
- Rimet, F., Chaumeil, P., Keck, F., Kermarrec, L., Vasselton, V., Kahlert, M., ... Bouchez, A. (2016). R-Syst::diatom: An open-access and curated barcode database for diatoms and freshwater monitoring. *Database*, 2016, baw016. <https://doi.org/10.1093/database/baw016>
- Rimet, F., Trobajo, R., Mann, D. G., Kermarrec, L., Franc, A., Domaizon, I., & Bouchez, A. (2014). When is sampling complete? the effects of geographical range and marker choice on perceived diversity in *Nitzschia palea* (Bacillariophyta). *Protist*, 165, 245–259. <https://doi.org/10.1016/j.protis.2014.03.005>
- Round, F. E., Crawford, R. M., & Mann, D. G. (1990). *Diatoms: Biology and morphology of the genera*. London: Cambridge University Press.
- Sabir, J. S. M., Yu, M., Ashworth, M. P., Baeshen, N. A., Baeshen, M. N., Bahieldin, A., ... Jansen, R. K. (2014). Conserved gene order and expanded inverted repeats characterize plastid genomes of Thalassiosirales. *PLoS ONE*, 9, e107854. <https://doi.org/10.1371/journal.pone.0107854>
- Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., ... Weber, C. F. (2009). Introducing MOTHUR: Open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Applied and Environmental Microbiology*, 75, 7537–7541. <https://doi.org/10.1128/AEM.01541-09>
- Snoeijs, P., Busse, S., & Potapova, M. (2002). The importance of diatom cell size in community analysis. *Journal of Phycology*, 38, 265–281. <https://doi.org/10.1046/j.1529-8817.2002.01105.x>

- Spaulding, S. A., Jewson, D. H., Bixby, R. J., Nelson, H., & McKnight, D. M. (2012). Automated measurement of diatom size. *Limnology and Oceanography: Methods*, *10*, 882–890. <https://doi.org/10.4319/lom.2012.10.882>
- Stein, E. D., Martinez, M. C., Stiles, S., Miller, P. E., & Zakharov, E. V. (2014). Is DNA barcoding actually cheaper and faster than traditional morphological methods: Results from a survey of freshwater bioassessment efforts in the United States?. *PLoS ONE*, *9*, e95525. <https://doi.org/10.1371/journal.pone.0095525>
- Sun, J., & Liu, D. (2003). Geometric models for calculating cell biovolume and surface area for phytoplankton. *Journal of Plankton Research*, *25*, 1331–1346. <https://doi.org/10.1093/plankt/fbg096>
- Taberlet, P., Coissac, E., Hajibabaei, M., & Rieseberg, L. H. (2012). Environmental DNA. *Molecular Ecology*, *21*, 1789–1793. <https://doi.org/10.1111/j.1365-294X.2012.05542.x>
- Tan, B., Ng, C., Nshimiyimana, J. P., Loh, L. L., Gin, K. Y.-H., & Thompson, J. R. (2015). Next-generation sequencing (NGS) for assessment of microbial water quality: Current progress, challenges, and future opportunities. *Frontiers in microbiology*, *6*, 1027.
- Thomas, A. C., Deagle, B. E., Eveson, J. P., Harsch, C. H., & Trites, A. W. (2016). Quantitative DNA metabarcoding: Improved estimates of species proportional biomass using correction factors derived from control material. *Molecular Ecology Resources*, *16*, 714–726. <https://doi.org/10.1111/1755-0998.12490>
- Treusch, A. H., Demir-Hilton, E., Vergin, K. L., Worden, A. Z., Carlson, C. A., Donatz, M. G., ... Giovannoni, S. J. (2012). Phytoplankton distribution patterns in the northwestern Sargasso Sea revealed by small subunit rRNA genes from plastids. *The ISME Journal*, *6*, 481–492. <https://doi.org/10.1038/ismej.2011.117>
- Vasselon, V., Domaizon, I., Rimet, F., Kahlert, M., & Bouchez, A. (2017a). Application of high-throughput sequencing (HTS) metabarcoding to diatom biomonitoring: Do DNA extraction methods matter? *Freshwater Science*, *36*, 162–177. <https://doi.org/10.1086/690649>
- Vasselon, V., Rimet, F., Tapolczai, K., & Bouchez, A. (2017b). Assessing ecological status with diatoms DNA metabarcoding: Scaling-up on a WFD monitoring network (Mayotte Island, France). *Ecological Indicators*, *82*, 1–12. <https://doi.org/10.1016/j.ecolind.2017.06.024>
- Visco, J. A., Apothéloz-Perret-Gentil, L., Cordonier, A., Esling, P., Pillet, L., & Pawlowski, J. (2015). Environmental monitoring: Inferring the diatom index from next-generation sequencing data. *Environmental Science & Technology*, *49*, 7597–7605. <https://doi.org/10.1021/es506158m>
- Vivien, R., Lejzerowicz, F., & Pawlowski, J. (2016). Next-generation sequencing of aquatic oligochaetes: Comparison of experimental communities. *PLoS ONE*, *11*, 1–14.
- Weber, A. A.-T., & Pawlowski, J. (2013). Can abundance of protists be inferred from sequence data: A case study of foraminifera. *PLoS ONE*, *8*, e56739. <https://doi.org/10.1371/journal.pone.0056739>
- Zetsche, E.-M., & Meysman, F. J. R. (2012). Dead or alive? Viability assessment of micro- and mesoplankton. *Journal of Plankton Research*, *34*, 493–509. <https://doi.org/10.1093/plankt/fbs018>
- Zimmermann, J., Glöckner, G., Jahn, R., Enke, N., & Gemeinholzer, B. (2015). Metabarcoding vs. morphological identification to assess diatom diversity in environmental studies. *Molecular Ecology Resources*, *15*, 526–542. <https://doi.org/10.1111/1755-0998.12336>
- Znachor, P., Rychtecký, P., Nedoma, J., & Visocká, V. (2015). Factors affecting growth and viability of natural diatom populations in the meso-eutrophic Řimov Reservoir (Czech Republic). *Hydrobiologia*, *762*, 253–265. <https://doi.org/10.1007/s10750-015-2417-8>

SUPPORTING INFORMATION

Additional Supporting Information may be found online in the supporting information tab for this article.

How to cite this article: Vasselon V, Bouchez A, Rimet F, et al. Avoiding quantification bias in metabarcoding: Application of a cell biovolume correction factor in diatom molecular biomonitoring. *Methods Ecol Evol*. 2018;9:1060–1069. <https://doi.org/10.1111/2041-210X.12960>