

Insights into global planktonic diatom diversity: Comparisons between phylogenetically meaningful units that account for time

Teofil Nakov, Jeremy M. Beaulieu, and Andrew J. Alverson

*Department of Biological Sciences University of Arkansas 1 University of Arkansas,
SCEN 601 Fayetteville, AR 72701*

Abstract

Metabarcoding has offered unprecedented insights into microbial diversity. In many studies, short DNA sequences are binned into consecutively higher Linnaean ranks, and ranked groups (e.g., genera) are the units of biodiversity analyses. These analyses assume that Linnaean ranks are biologically meaningful and that identically ranked groups are comparable. We used a meta-barcode dataset for marine planktonic diatoms to illustrate the limits of this approach. We found that the 20 most abundant marine planktonic diatom genera ranged in age from 4 to 134 million years, indicating the non-equivalence of genera because some had more time to diversify than others. Still, species richness was only weakly correlated with genus age, highlighting variation in rates of speciation and/or extinction. Taxonomic classifications often do not reflect phylogeny, so genus-level analyses can include phylogenetically nested genera, further confounding rank-based analyses. These results underscore the indispensable role of phylogeny in understanding patterns of microbial diversity.

Keywords: diversification, metabarcoding, microbes, phylogeny

1 With potentially millions of species occupying all the world's aquatic and
2 terrestrial biomes, microbial species diversity is notoriously difficult to dis-
3 cover and catalog. Traditional approaches to species discovery are time and
4 labor intensive, and they miss species that cannot be cultivated in the lab [1].
5 The phylogenetic diversity of this undiscovered "microbial dark matter" is
6 often characterized through community DNA sequencing of barcode genes.
7 A typical workflow includes DNA extraction from an environmental sam-
8 ple, PCR amplification of a barcode fragment, and high-throughput DNA
9 sequencing of the amplicon. Sequencing reads are clustered into operational
10 taxonomic units (OTUs) that are subsequently binned into consecutively
11 higher taxonomic ranks, and these ranked groups, in turn, are often the
12 focus of biodiversity assessments [2].

13 Linnaean names and ranks are often taken to mean more than what they
14 are: arbitrary taxon delimitations disconnected from evolutionary history.
15 The treatment of named ranks as anything other than arbitrary implies that
16 identically ranked groups are somehow comparable, encouraging comparisons
17 of their ecology, biogeography, and species richness [3, 4, 5]. The only mean-
18 ingful comparisons involve groups with comparable evolutionary histories [6].
19 In this sense, monophyletic groups (clades) are more likely to be biologically
20 cohesive units, and they should have comparable species richness if they are
21 similar in age and have diversified at similar rates [7]. Comparison of mono-
22 phyletic groups, while accounting for time, provides a robust framework for
23 detecting clades with exceptional species richness and comparing their func-
24 tional, ecological, or biogeographic breadth [8].

25 The Tara Oceans Project sequenced 18S-V9 metabarcode fragments from

26 plankton samples to characterize microbial communities and species richness
27 across the world's oceans [9]. A total of 20 diatom genera accounted for nearly
28 99% of all diatom sequencing reads, and these genera were found to differ
29 in relative abundance, cell size, habitat preference, geographical distribution,
30 and species richness [2], however, it was unclear whether or not these patterns
31 deviated from expectations.

32 We focused our analyses on the genus-based patterns of species richness
33 and expected that older genera would be more species rich because they
34 have had a longer period of time to diversify [7]. We used a time-calibrated
35 phylogenetic tree of 1,151 diatoms [10] to calculate expected species richness
36 for the 20 most abundant marine diatom genera in the Tara Oceans survey
37 [2]. Given the crown age of diatoms [10], relative extinction (i.e., extinc-
38 tion/speciation) estimated from Cenozoic fossil species [11], and a minimum
39 approximation of total described and undescribed diatom diversity (30,000
40 species; [12]), we calculated a net diversification rate for diatoms (i.e., speci-
41 ation - extinction). We then used this rate to calculate the upper and lower
42 bounds of expected OTU richness for the 20 focal diatom genera [8].

43 The 20 diatom genera ranged in age from 4134 million years (My). OTU
44 richness was only weakly correlated with clade age ($r=0.36$, 95% CI=-0.1–
45 0.7, $df=18$, $P=0.12$), with 12 of the 20 genera falling within expectations for
46 OTU diversity given their age (Figure 1). The most abundant and OTU-rich
47 genus, *Chaetoceros*, was also the oldest (Figure 1a). The birthdeath diversi-
48 fication model predicted that *Chaetoceros* diversity should range between 47
49 and 6567 species—the Tara Oceans dataset recovered 644 OTUs, consistent
50 with expectations for a clade of this age (Figure 1b). Some of the most diverse

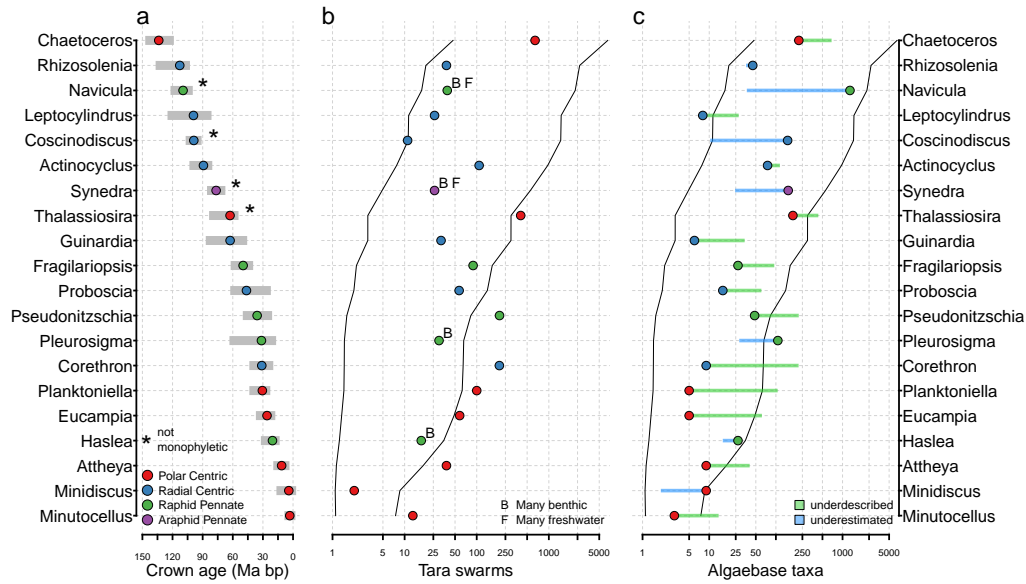


Figure 1: Age and estimated taxon richness of the 20 most abundant marine planktonic diatom genera identified by the Tara Oceans metabarcoding project [2]. Crown ages and uncertainty (grey bars) were estimated from 1000 bootstrap phylogenies [10] (a). Taxon richness was estimated from the number of OTU swarms in the Tara Oceans dataset (b) and the number of accepted species names in AlgaeBase [13] (c); lines delimit 95% confidence intervals of expected richness given the crown age of a clade, empirical extinction fraction, and diatom-wide estimate of the net diversification rate.

51 genera identified by metabarcoding (e.g., *Corethron* and *Pseudo-nitzschia*)
52 had OTU richness estimates that exceeded expectations. Assuming OTUs
53 correspond to species and that our estimates of clade age are not heavily bi-
54 ased, these genera have either exceptionally high speciation or low extinction
55 rates. Identifying the drivers of these patterns might offer new mechanistic
56 insights into phytoplankton diversification. Comparisons between OTU rich-
57 ness (Figure 1b) and accepted taxonomic names from AlgaeBase (Figure 1c)
58 showed expected discrepancies for lineages with substantial diversity in ben-
59 thic or freshwater habitats (e.g., *Navicula*; Figure 1b, B and F annotations;
60 Figure 1c, blue bars) and were also useful in highlighting clades that are
61 underdescribed at the species level (Figure 1c, green bars).

62 Metabarcoding identified *Thalassiosira* as one of the most abundant,
63 OTU-rich, and geographically widespread marine planktonic diatom gen-
64 era. A total of eight Thalassiosirales genera were detected in the Tara
65 Oceans project (*Cyclotella*, *Lauderia*, *Minidiscus*, *Planktoniella*, *Porosira*,
66 *Shionodiscus*, *Skeletonema*, and *Thalassiosira*), and these genera ranged in
67 age from 463 My (Figure 2). Thalassiosirales embodies many of the problems
68 with misappropriation of biological or evolutionary properties to taxa based
69 on their names [14]. The name *Thalassiosira* applies to a polyphyletic set
70 of species whose common ancestor dates to 63 My and gave rise to nearly
71 the full phylogenetic breadth Thalassiosirales diversity (Figure 2, diamond).
72 As a result, including *Thalassiosira* in genus-level analyses leads to highly
73 biased comparisons involving a genus that, in reality, corresponds roughly to
74 a taxonomic order (Figure 2). Moreover, four of the eight thalassiosiroid gen-
75 era detected by metabarcoding are nested within *Thalassiosira*, highlighting

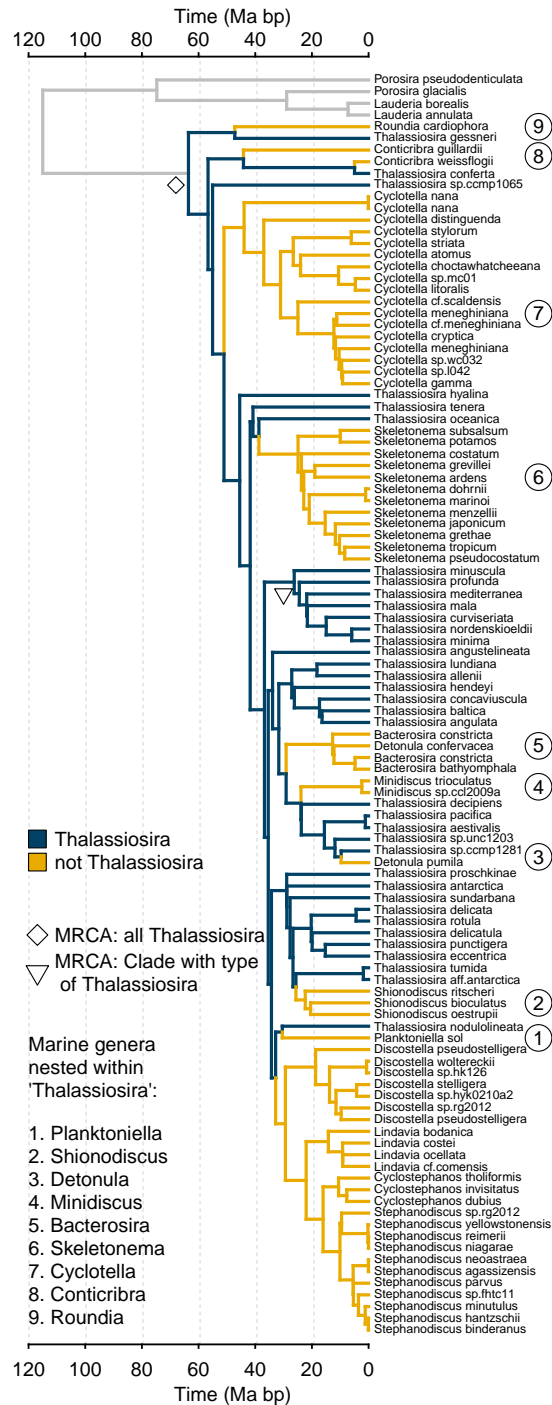


Figure 2: The genus *Thalassiosira* encompasses at least 10 marine planktonic diatom genera (including *Thalassiosira*) that range from 4-63 My in age. Topology and divergence times are based on [10].

76 a common source of non-independence in rank-based comparisons (Figure 2,
77 yellow branches). A more informative, phylogenetically based genus-level
78 classification may have revealed clade-specific habitat preferences or geo-
79 graphic distributions among the many distinct *Thalassiosira* lineages [15].

80 The problems with rank-based comparisons, including as they relate to di-
81 atoms, are well known [14, 16, 15]. A frequently cited advantage of metabar-
82 coding is that it does not require taxonomic expertise. Still, the taxonomic
83 affiliations of metabarcode sequences often become the units of biodiversity
84 analyses. A working knowledge of phylogeny and systematics—which invari-
85 ably highlight the deficiencies of Linnaean classifications—can lead to more
86 meaningful analyses that explicitly incorporate phylogenetic history, ensuring
87 robust comparisons of biologically equivalent units that account for time.

88 [1] C. Rinke, P. Schwientek, A. Sczyrba, N. N. Ivanova, I. J. Anderson, J.-F.
89 Cheng, A. Darling, S. A. Malfatti, B. K. Swan, E. A. Gies, et al., Insights
90 into the phylogeny and coding potential of microbial dark matter (2013).

91 [2] S. Malviya, E. Scalco, S. Audic, F. Vincent, A. Veluchamy, J. Poulain,
92 P. Wincker, D. Iudicone, C. de Vargas, L. Bittner, et al., Insights into
93 global diatom distribution and diversity in the worlds ocean, Proceed-
94 ings of the National Academy of Sciences 113 (2016) E1516–E1525.

95 [3] P. D. Cantino, K. de Queiroz, et al., Phylocode: a phylogenetic code of
96 biological nomenclature, 2000.

97 [4] F. Pleijel, G. W. Rouse, Ceci n'est pas une pipe: names, clades and
98 phylogenetic nomenclature, Journal of Zoological Systematics and Evo-
99 lutionary Research 41 (2003) 162–174.

- 100 [5] P. Sundberg, F. Pleijel, Phylogenetic classification and the definition of
101 taxon names, *Zoologica scripta* 23 (1994) 19–25.
- 102 [6] P. H. Harvey, M. D. Pagel, et al., *The comparative method in evolu-*
103 *tionary biology*, volume 239, Oxford University Press Oxford, 1991.
- 104 [7] T. Stadler, D. L. Rabosky, R. E. Ricklefs, F. Bokma, On age and species
105 richness of higher taxa, *The American Naturalist* 184 (2014) 447–455.
- 106 [8] S. Magallón, M. J. Sanderson, Absolute diversification rates in an-
107 giosperm clades, *Evolution* 55 (2001) 1762–1780.
- 108 [9] S. Sunagawa, L. P. Coelho, S. Chaffron, J. R. Kultima, K. Labadie,
109 G. Salazar, B. Djahanschiri, G. Zeller, D. R. Mende, A. Alberti, et al.,
110 Structure and function of the global ocean microbiome, *Science* 348
111 (2015) 1261359.
- 112 [10] T. Nakov, J. M. Beaulieu, A. J. Alverson, A time-calibrated phylogeny
113 of diatoms (bacillariophyta), under review (2017).
- 114 [11] D. Lazarus, J. Barron, J. Renaudie, P. Diver, A. Türke, Cenozoic plank-
115 tonic marine diatom diversity and correlation to climate change, *PLoS*
116 *One* 9 (2014) e84857.
- 117 [12] D. G. Mann, P. Vanormelingen, An inordinate fondness? the num-
118 ber, distributions, and origins of diatom species, *Journal of eukaryotic*
119 *microbiology* 60 (2013) 414–420.
- 120 [13] M. Guiry, G. Guiry, *Algaebase*. national university of ireland, galway,
121 2017.

- 122 [14] A. J. Alverson, B. Beszteri, M. L. Julius, E. C. Theriot, The model ma-
123 rine diatom *thalassiosira pseudonana* likely descended from a freshwater
124 ancestor in the genus *cyclotella*, *BMC evolutionary biology* 11 (2011)
125 125.
- 126 [15] R. Wiese, J. Renaudie, D. B. Lazarus, Testing the accuracy of genus-
127 level data to predict species diversity in cenozoic marine diatoms, *Ge-*
128 *ology* 44 (2016) 1051–1054.
- 129 [16] J. P. Kociolek, Taxonomy and ecology: further considerations,
130 *Proceedings-California Academy of Sciences* 56 (2005) 99.