

Quantifying the taxonomic bias in enzymology

Chelsea J. Vickers¹ | Dean Fraga² | Wayne M. Patrick¹ 

¹Centre for Biodiscovery, School of Biological Sciences, Victoria University of Wellington, Wellington, New Zealand

²Department of Biology, The College of Wooster, Wooster, Ohio

Correspondence

Wayne M. Patrick, School of Biological Sciences, Victoria University of Wellington, PO Box 600, Wellington 6140, New Zealand.

Email: wayne.patrick@vuw.ac.nz

The copyright line for this article was changed on 19 April, 2021 after original online publication.

Abstract

The ongoing biotechnological revolution is rooted in our knowledge of enzymes. However, metagenomics is showing how little we know about Earth's enzyme repertoire. Deep sequencing has revolutionized our view of the tree of life. The genomes of newly-discovered organisms are replete with novel sequences, emphasizing the trove of enzyme structures and functions waiting to be explored by biochemists. Here, we sought to draw attention to the vastness of the “enzymatic dark matter” within the tree of life by placing enzymological knowledge in the context of phylogeny. We used kinetic parameters from the BRAunschweig ENzyme DATabase (BRENDA) as our proxy for enzymological knowledge. Mapping 12,677 BRENDA entries onto the phylogenetic tree revealed that 55% of these data were from eukaryotes, even though they are the least diverse part of the tree. At the next taxonomic level, only four of 18 archaeal phyla and 24 of 111 bacterial phyla are represented in the BRENDA dataset. One phylum, the Proteobacteria, accounts for over half of all bacterial entries. Similarly, the supergroup Amorphea, which includes animals and fungi, contains over half the data on eukaryotes. Many major taxonomic groups are notable for their complete absence from BRENDA, including the ultra-diverse bacterial Candidate Phyla Radiation. At the species level, five mammals (including human) contribute 15% of BRENDA entries. The taxonomic bias in enzymology is strong, but in the era of gene synthesis we now have the tools to address it. Doing so promises to enrich our biochemical understanding of life and uncover powerful new biocatalysts.

KEYWORDS

archaeal and bacterial phyla, enzymatic dark matter, eukaryotic supergroups, kinetic parameters, phylogenetics, tree of life

1 | INTRODUCTION

It has now been over 120 years since Eduard Buchner famously founded the field of enzymology by reporting that a cell-free extract of yeast could “institute fermentation of carbohydrates”.^{1,2} In the ensuing decades, enzymologists have learned an enormous amount about the structures, functions and mechanisms of thousands of enzymes.³ Increasingly, this fundamental knowledge is

allowing us to engineer or evolve enzymes into powerful biocatalysts for industry and synthetic biology.^{3,4}

Nevertheless, it is becoming apparent how little we actually know about the diversity of enzymes on Earth. Even the most biochemically well-characterized organisms – *Homo sapiens*, *Escherichia coli*, *Saccharomyces cerevisiae* and *Arabidopsis thaliana* – have had turnover numbers (k_{cat} values) reported for fewer than 10% of their enzymes.⁵ At the same time, metagenomics is

This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2021 The Authors. *Protein Science* published by Wiley Periodicals LLC on behalf of The Protein Society.

revolutionizing our view of the tree of life.^{6–9} Deep sequencing has proved the existence of previously unimagined bacterial and archaeal lifeforms, inhabiting ever more extreme environments. The genomes of these newly-discovered microorganisms demonstrate there is a stupendous diversity of protein sequences and enzyme functions waiting to be explored.^{10–13} Similarly, a combination of culturing and sequencing has massively expanded our knowledge of protist diversity, leading to substantial revisions of the eukaryotic part of the tree.^{6,14–16}

The goals of this study were to map our knowledge of enzymes in the context of phylogeny, and thence to draw attention to those parts of the tree of life that remain unexplored by biochemists. Our collective knowledge of enzymes takes many forms, including information on structure, function, kinetic parameters, expression level, folding, stability, mutability, promiscuity and evolution. As a proxy for all of this enzymological knowledge, we have chosen to focus on kinetic parameters for wild type enzymes catalyzing naturally occurring reactions. Previously, Davidi et al. filtered all of the entries in the BRAunschweig ENzyme DAtabase (BRENDA),¹⁷ as it was circa April 2017, for entries of this type. They used this filtered dataset to draw insightful conclusions about the chemical and physiological demands that shape enzyme evolution, while also acknowledging there has been a sampling bias in the enzymes studied to date.¹⁸ Here, we have used their dataset to examine this bias in more detail.

2 | RESULTS AND DISCUSSION

The dataset compiled by Davidi et al. comprises 12,721 entries (i.e., Table S1 from ref 18, which we accessed in mid-2020). Each entry contains at least one kinetic parameter (k_{cat} , K_M and/or k_{cat}/K_M) for an enzyme, catalyzing either a forward or backward reaction. Usefully, the authors also included the identity of the organism from which the enzyme was obtained. For our taxonomic analyses, we discarded the data on viral enzymes (40 entries), unclassified enzymes (three entries) and an entry on a synthetic ribozyme. This left a trimmed dataset of 12,677 entries. There is inevitably noise in this kind of global analysis, such as occasional inconsistencies between the BRENDA entry and the original paper (as discussed previously¹⁹). Nevertheless, our dataset serves as an informative snapshot of humankind's effort to study enzymes, as well as where there are gaps in our knowledge.

We began by using the NCBI Taxonomy database²⁰ to manually assign each BRENDA entry to its domain (Archaea, Bacteria or Eukaryota; Figure 1). In the dataset

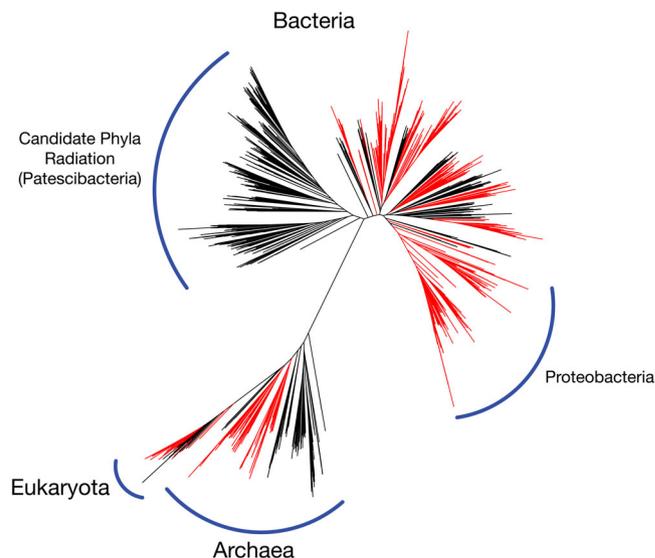


FIGURE 1 Low-resolution mapping of enzymological knowledge onto the tree of life. The phylogenetic tree was constructed by Castelle and Banfield using ribosomal protein sequences⁸ and highlights the relative levels of diversity in the three domains of life (Bacteria, Archaea and Eukaryota). We used the Interactive Tree of Life online tool (iTOL)⁴⁰ to color red the nodes where there is at least one entry in our BRENDA database of enzyme kinetic parameters. Major parts of the tree with no BRENDA entries are black. The coloration is at nodes approximately corresponding to phyla and aims to emphasize the broad gaps in enzymological knowledge

(Table S1), 901 entries (7%) were from the domain Archaea, 4,804 (38%) were from Bacteria and 6,972 (55%) were from Eukaryota. The eukaryotes comprise the least diverse part of the tree of life,^{8,21} and yet – unsurprisingly – over half our accumulated knowledge of enzyme kinetics comes from organisms in this domain. On the other hand, large taxonomic groups such as the bacterial Candidate Phyla Radiation (CPR) and multiple archaeal phyla are unexplored enzymatically (Figure 1).

While Figure 1 provided a low-resolution snapshot across the major phylogenetic groups, we also aimed to rigorously assign each BRENDA entry to one additional taxonomic rank. For Archaea and Bacteria, this was the phylum. There has been considerable recent interest in defining microbial relationships at the level of the phylum,^{7,8,22,23} and the Genome Taxonomy Database (GTDB; <https://gtdb.ecogenomic.org/>) has emerged as a comprehensive, genome-based resource for standardized taxonomic classification.²⁴ It also serves as a useful standard reference tool for phylogenetically aware enzymology. Therefore, we manually assigned each archaeal and bacterial BRENDA entry to its phylum as defined in release 95 of the GTDB (Table S1, column C). However, some of the GTDB classifications are unlikely to be

TABLE 1 The eight species with the most characterized enzymes, from each domain of life

Domain	No entries	% domain's total entries	Phylum ^a or supergroup ^b	Additional identifier ^c
Archaea				
<i>Sulfolobus solfataricus</i>	107	11.9%	Thermoproteota	Crenarchaeota
<i>Pyrococcus furiosus</i>	69	7.7%	Methanobacteriota	Euryarchaeota
<i>Archaeoglobus fulgidus</i>	67	7.4%	Halobacteriota	Euryarchaeota
<i>Methanocaldococcus jannaschii</i>	57	6.3%	Methanobacteriota	Euryarchaeota
<i>Sulfolobus tokodaii</i>	55	6.1%	Thermoproteota	Crenarchaeota
<i>Pyrococcus horikoshii</i>	46	5.1%	Methanobacteriota	Euryarchaeota
<i>Aeropyrum pernix</i>	38	4.2%	Thermoproteota	Crenarchaeota
<i>Thermococcus kodakarensis</i>	37	4.1%	Methanobacteriota	Euryarchaeota
Bacteria				
<i>Escherichia coli</i>	669	13.9%	Proteobacteria	Gammaproteobacteria
<i>Mycobacterium tuberculosis</i>	195	4.1%	Actinobacteriota	Actinobacteria
<i>Bacillus subtilis</i>	121	2.5%	Firmicutes	Bacilli
<i>Salmonella enterica</i>	112	2.3%	Proteobacteria	Gammaproteobacteria
<i>Pseudomonas putida</i>	107	2.2%	Proteobacteria	Gammaproteobacteria
<i>Pseudomonas aeruginosa</i>	103	2.1%	Proteobacteria	Gammaproteobacteria
Unspecified <i>Pseudomonas</i> sp.	101	2.1%	Proteobacteria	Gammaproteobacteria
<i>Thermus thermophilus</i>	84	1.7%	Deinococcota	Deinococcus-Thermus
Eukaryota				
<i>Homo sapiens</i>	766	11.0%	Amorphea	Metazoa (animals)
<i>Rattus norvegicus</i>	548	7.9%	Amorphea	Metazoa (animals)
<i>Saccharomyces cerevisiae</i>	298	4.3%	Amorphea	Fungi
<i>Bos taurus</i>	261	3.7%	Amorphea	Metazoa (animals)
<i>Arabidopsis thaliana</i>	242	3.5%	Archaeplastida	Viridiplantae (green plants)
<i>Mus musculus</i>	207	3.0%	Amorphea	Metazoa (animals)
<i>Sus scrofa</i>	165	2.4%	Amorphea	Metazoa (animals)
<i>Oryctolagus cuniculus</i>	115	1.6%	Amorphea	Metazoa (animals)

^aFor Archaea and Bacteria.^bFor Eukaryota.^cIdentifiers from the NCBI Taxonomy browser²⁰ that may or may not have standing in the taxonomy community, but remain commonly used by enzymologists.

more entries in BRENDA. At the other extreme, 682 species are represented by a single BRENDA entry. Over 88% of the species in our BRENDA dataset are represented by fewer than 10 entries and the trend is consistent across all three domains of life (Figure 4).

The skew towards a few very well-characterized model organisms is most evident in the eukaryotes, where over half of the total entries are from only 2.1% of species and the top eight eukaryotic species contribute 37.3% of all entries (Table 1). The mean number of entries per species is 7.3 but the median is only 2. Across

the entire tree of life, seven of the ten most well-characterized species are eukaryotes (Table 1).

Five of the ten most well-characterized species are not only eukaryotes, but mammals (human, rat, cattle, mouse and pig). Together, these five contribute 15% of all BRENDA entries. This observation is unlikely to surprise biochemists, most of whom have probably characterized enzymes from one or more of these species during their undergraduate training. For perspective, it is worth reflecting that there are ~12,000 species of mammals on Earth.²⁶ However, it has been estimated there are over

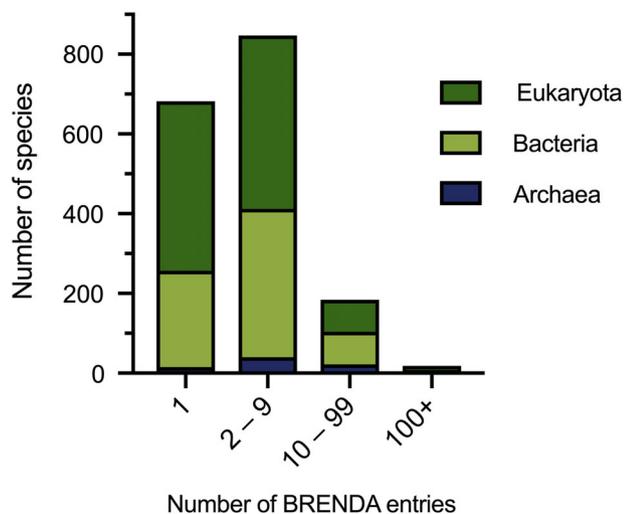


FIGURE 4 Counts of species according to how many entries they contribute to our BRENDA dataset. There are 18 species with more than 100 entries. Of these, one is archaeal, seven are bacterial and ten are eukaryotic

160 million species of animals (almost all arthropods and nematodes), and at least two billion species in total (mostly bacteria).²¹ Thus, a conservative estimate is that mammals comprise $\sim 0.0006\%$ of the species on Earth.

Phylogenetic distance increases the likelihood of discovering novel protein folds and functions.²⁷ This was illustrated in dramatic fashion when a single study reported 1,003 reference genomes selected solely because they maximized coverage of phylogenetic space, and which collectively increased the number of known protein sequence families by over 10%.¹¹ Sequence similarity networks and genome mining algorithms are also being used increasingly to identify enzymes with divergent sequences that expand our understanding, have biotechnological value, or both.^{13,28-32} For example, a recent study identified 33,000 sequences for the carbon-fixing enzyme, ribulose-1,5-bisphosphate carboxylase/oxygenase (RuBisCO), in genomic and metagenomic datasets.²⁸ A subset of 143 maximally diverse enzymes were characterized biochemically. RuBisCO has been studied for decades, from over 200 species, and yet five of these new-to-science enzymes had higher k_{cat} values than any previously characterized. Perhaps surprisingly, the fastest was not from a plant, alga or cyanobacterium, but instead from a soil-dwelling bacterium of the genus *Gallionella*.²⁸ It remains to be seen whether this fast RuBisCO might overcome the longstanding biotechnological challenge of improving carbon capture, and therefore growth, in crops.³³ There is clearly enormous potential for this approach of phylogenetically aware enzymology to realize the value of the deep sequencing revolution, both by bringing biochemical insights and by uncovering novel

biocatalysts with desirable properties for the biotechnology industry.

3 | CONCLUSIONS

Severe taxonomic biases have been quantified in other branches of biology. For example, it is widely acknowledged that the fields of ecology and conservation biology have been overly focused on a handful of charismatic, highly fundable animal species.^{26,34-36} In the case of biochemistry, it is equally clear we have a highly biased view of enzymatic diversity at all taxonomic levels. By quantifying that bias here, we have called attention to the vastness of the “enzymatic dark matter” within the tree of life.

In the era of gene synthesis and recombinant protein expression, it is possible to study any enzyme from any organism, regardless of whether that organism is culturable. Studying enzymes from maximally diverse, non-model organisms will broaden our understanding of the sequence-structure-function relationship and shed new light on the metabolic processes taking place throughout the biosphere.^{31,37,38} It will also uncover potent natural products and valuable new parts for synthetic biology.¹³ Perhaps most importantly, illuminating the enzymatic dark matter will undoubtedly inspire us with the “beautiful aspects of Nature’s chemistry”.³⁹

ACKNOWLEDGEMENTS

This work was supported via a grant from the Marsden Fund to WMP [grant number 18-VUW-050]. DF was supported by the Henry Luce III Fund for Distinguished Scholarship and the James and Jean McClung Endowment for International Research.

AUTHOR CONTRIBUTIONS

Chelsea J. Vickers: Conceptualization; data curation; formal analysis; investigation; visualization; writing-original draft; writing-review and editing. **Dean Fraga:** Conceptualization; data curation; formal analysis; funding acquisition; visualization; writing-review and editing. **Wayne Patrick:** Conceptualization; data curation; formal analysis; funding acquisition; investigation; visualization; writing-original draft; writing-review and editing.

ORCID

Wayne M. Patrick  <https://orcid.org/0000-0002-2718-8053>

REFERENCES

1. Buchner E. Alkoholische Gärung ohne Hefezellen. Ber Dtsch Chem Ges. 1897;30:117-124.

2. Friedmann HC. Alcoholic fermentation without yeast cells. In: A. Cornish-Bowden, editor. *New Beer in an Old Bottle: Eduard Buchner and the Growth of Biochemical Knowledge*. València, Spain: Universitat de València, 1997; p. 25–31.
3. Heckmann CM, Paradisi F. Looking back: A short history of the discovery of enzymes and how they became powerful chemical tools. *Chem Cat Chem*. 2020;12:6082–6102.
4. Renata H, Wang ZJ, Arnold FH. Expanding the enzyme universe: Accessing non-natural reactions by mechanism-guided directed evolution. *Angew Chem Int Ed Engl*. 2015;54:3351–3367.
5. Davidi D, Milo R. Lessons on enzyme kinetics from quantitative proteomics. *Curr Opin Biotechnol*. 2017;46:81–89.
6. Burki F, Roger AJ, Brown MW, Simpson AGB. The new tree of eukaryotes. *Trends Ecol Evol*. 2020;35:43–55.
7. Zhu Q, Mai U, Pfeiffer W, et al. Phylogenomics of 10,575 genomes reveals evolutionary proximity between domains bacteria and archaea. *Nat Commun*. 2019;10:5477.
8. Castelle CJ, Banfield JF. Major new microbial groups expand diversity and alter our understanding of the tree of life. *Cell*. 2018;172:1181–1197.
9. Lloyd KG, Steen AD, Ladau J, Yin J, Crosby L. Phylogenetically novel uncultured microbial cells dominate earth microbiomes. *mSystems*. 2018;3:e00055–e00018.
10. Castelle CJ, Brown CT, Anantharaman K, Probst AJ, Huang RH, Banfield JF. Biosynthetic capacity, metabolic variety and unusual biology in the CPR and DPANN radiations. *Nat Rev Microbiol*. 2018;16:629–645.
11. Mukherjee S, Seshadri R, Varghese NJ, et al. 1,003 reference genomes of bacterial and archaeal isolates expand coverage of the tree of life. *Nat Biotechnol*. 2017;35:676–683.
12. León-Zayas R, Peoples L, Biddle JF, et al. The metabolic potential of the single cell genomes obtained from the challenger deep, Mariana trench within the candidate superphylum Parcubacteria (OD1). *Environ Microbiol*. 2017;19:2769–2784.
13. Gerlt JA. Genomic enzymology: Web tools for leveraging protein family sequence-function space and genome context to discover novel functions. *Biochemistry*. 2017;56:4293–4308.
14. Gawryluk RMR, Tikhonenkov DV, Hehenberger E, Husnik F, Mylnikov AP, Keeling PJ. Non-photosynthetic predators are sister to red algae. *Nature*. 2019;572:240–243.
15. Mahé F, de Vargas C, Bass D, et al. Parasites dominate hyperdiverse soil protist communities in neotropical rainforests. *Nat Ecol Evol*. 2017;1:91.
16. de Vargas C, Audic S, Henry N, et al. Eukaryotic plankton diversity in the sunlit ocean. *Science*. 2015;348:1261605.
17. Chang A, Jeske L, Ulbrich S, et al. BRENDA, the ELIXIR core data resource in 2021: New developments and updates. *Nucleic Acids Res*. 2021;49:D498–D508.
18. Davidi D, Longo LM, Jablonska J, Milo R, Tawfik DS. A Bird's-eye view of enzyme evolution: Chemical, physicochemical, and physiological considerations. *Chem Rev*. 2018;118:8786–8797.
19. Bar-Even A, Noor E, Savir Y, et al. The moderately efficient enzyme: Evolutionary and physicochemical trends shaping enzyme parameters. *Biochemistry*. 2011;50:4402–4410.
20. Schoch CL, Ciufo S, Domrachev M, et al. NCBI taxonomy: A comprehensive update on curation, resources and tools. *Database* 2020;2020:baaa062.
21. Larsen BB, Miller EC, Rhodes MK, Wiens JJ. Inordinate fondness multiplied and redistributed: The number of species on earth and the new pie of life. *Q Rev Biol*. 2017;92:229–265.
22. Baker BJ, De Anda V, Seitz KW, Dombrowski N, Santoro AE, Lloyd KG. Diversity, ecology and evolution of Archaea. *Nat Microbiol*. 2020;5:887–900.
23. Parks DH, Chuvochina M, Waite DW, et al. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol*. 2018;36:996–1004.
24. Parks DH, Chuvochina M, Chaumeil PA, Rinke C, Mussig AJ, Hugenholtz P. A complete domain-to-species taxonomy for bacteria and Archaea. *Nat Biotechnol*. 2020;38:1079–1086.
25. del Campo J, Sieracki ME, Molestina R, Keeling P, Massana R, Ruiz-Trillo I. The others: Our biased perspective of eukaryotic genomes. *Trends Ecol Evol*. 2014;29:252–259.
26. Troudet J, Grandcolas P, Blin A, Vignes-Lebbe R, Legendre F. Taxonomic bias in biodiversity data and societal preferences. *Sci Rep*. 2017;7:9132.
27. Kunin V, Cases I, Enright AJ, de Lorenzo V, Ouzounis CA. Myriads of protein families and still counting. *Genome Biol*. 2003;4:401.
28. Davidi D, Shamsoum M, Guo Z, et al. Highly active rubiscos discovered by systematic interrogation of natural sequence diversity. *EMBO J*. 2020;39:e104081.
29. Hon J, Borko S, Stourac J, et al. EnzymeMiner: Automated mining of soluble enzymes with diverse structures, catalytic properties and stabilities. *Nucleic Acids Res*. 2020;48:W104–W109.
30. Mak WS, Tran S, Marcheschi R, et al. Integrative genomic mining for enzyme function to enable engineering of a non-natural biosynthetic pathway. *Nat Commun*. 2015;6:10005.
31. Copp JN, Akiva E, Babbitt PC, Tokuriki N. Revealing unexplored sequence-function space using sequence similarity networks. *Biochemistry*. 2018;57:4651–4662.
32. Akiva E, Copp JN, Tokuriki N, Babbitt PC. Evolutionary and molecular foundations of multiple contemporary functions of the nitroreductase superfamily. *Proc Natl Acad Sci U S A*. 2017;114:E9549–E9558.
33. Ort DR, Merchant SS, Alric J, et al. Redesigning photosynthesis to sustainably meet global food and bioenergy demand. *Proc Natl Acad Sci U S A*. 2015;112:8529–8536.
34. dos Santos JW, Correia RA, Malhado ACM, et al. Drivers of taxonomic bias in conservation research: A global analysis of terrestrial mammals. *Anim Conserv*. 2020;23:679–688.
35. Clark JA, May RM. Taxonomic bias in conservation research. *Science*. 2002;297:191–192.
36. Hochkirch A, Samways MJ, Gerlach J, et al. A strategy for the next decade to address data deficiency in neglected biodiversity. *Conserv Biol*: PMID. 2020;32656858.
37. Ferla MP, Brewster JL, Hall KR, Evans GB, Patrick WM. Primal-like enzymes from bacteria with reduced genomes. *Mol Microbiol*. 2017;105:508–524.
38. Newton MS, Arcus VL, Gerth ML, Patrick WM. Enzyme evolution: Innovation is easy optimization is complicated. *Curr Opin Struct Biol*. 2018;48:110–116.
39. Tawfik DS, van der Donk WA. Editorial overview: Biocatalysis and biotransformation: esoteric niche enzymology. *Curr Opin Chem Biol*. 2016;31:v–vii.
40. Letunic I, Bork P. Interactive tree of life (iTOL) v4: Recent updates and new developments. *Nucleic Acids Res*. 2019;47:W256–W259.
41. Mandler K, Chen H, Parks DH, Lobb B, Hug LA, Doxey AC. AnnoTree: Visualization and exploration of a functionally

annotated microbial tree of life. *Nucleic Acids Res.* 2019;47:4442–4448.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.

How to cite this article: Vickers CJ, Fraga D, Patrick WM. Quantifying the taxonomic bias in enzymology. *Protein Science.* 2021;30:914–921.
<https://doi.org/10.1002/pro.4041>