

# Bacterial diversification through geological time

Stilianos Louca<sup>1,2\*</sup>, Patrick M. Shih<sup>3,4,5</sup>, Matthew W. Pennell<sup>1,2</sup>, Woodward W. Fischer<sup>6</sup>,  
Laura Wegener Parfrey<sup>1,2,7</sup> and Michael Doebeli<sup>1,2,8</sup>

**Numerous studies have estimated plant and animal diversification dynamics; however, no comparable rigorous estimates exist for bacteria—the most ancient and widespread form of life on Earth. Here, we analyse phylogenies comprising up to 448,112 bacterial lineages to reconstruct global bacterial diversification dynamics. To handle such large phylogenies, we developed methods based on the statistical properties of infinitely large trees. We further analysed sequencing data from 60 environmental studies to determine the fraction of extant bacterial diversity missing from the phylogenies—a crucial parameter for estimating speciation and extinction rates. We estimate that there are about 1.4–1.9 million extant bacterial lineages when lineages are defined by 99% similarity in the 16S ribosomal RNA gene, and that bacterial diversity has been continuously increasing over the past 1 billion years (Gyr). Recent bacterial extinction rates are estimated at 0.03–0.05 per lineage per million years (lineage<sup>-1</sup> Myr<sup>-1</sup>), and are only slightly below estimated recent bacterial speciation rates. Most bacterial lineages ever to have inhabited this planet are estimated to be extinct. Our findings disprove the notion that bacteria are unlikely to go extinct, and provide a valuable perspective on the evolutionary history of a domain of life with a sparse and cryptic fossil record.**

For over 3.5 Gyr, the geochemical composition of our planet has been shaped by the evolution and diversification of bacteria<sup>1</sup>. Most prominently, the Great Oxygenation Event was caused by cyanobacteria roughly 2.35 Gyr ago and dramatically altered Earth's surface environments and the subsequent evolution of life<sup>2</sup>. Despite the prominent role of bacteria in ancient and modern biospheres, little is known about the dynamics by which their diversity evolved over Earth's history. For many eukaryotes, the fossil record provides estimates of past diversity<sup>3–5</sup>, revealing that extant global eukaryotic diversity only represents a small fraction of the total diversity that existed in the past. Analogous estimates for bacterial diversity are lacking, largely because their fossil record is extremely poor, and thus the clades that are known are those with extant representatives. Fortunately, past diversification dynamics also leave a footprint in molecular phylogenies of extant organisms<sup>6</sup>. Many approaches have been developed to infer past diversification dynamics from these patterns<sup>7–9</sup>. Despite these methodological advances, global bacterial diversification dynamics remain largely unresolved and much less studied than eukaryotic diversification. Previous studies only examined diversification within a single bacterial genus<sup>10–12</sup> or a single archaeal phylum<sup>13</sup>, or phylogenies covering only a small and biased portion of diversity (~12,000 cultured bacterial and archaeal species)<sup>14</sup>. Many of these studies do not report absolute speciation or extinction rates<sup>10,12,13</sup>. Importantly, no previous study properly accounted for the incomplete sampling of bacterial diversity represented in the phylogenies. Knowledge of the 'sampling fraction', in addition to any phylogenetic information, is critical for estimating speciation and extinction rates from phylogenies, even to an order of magnitude<sup>15</sup>. As the extant global bacterial diversity was so far largely unknown, previous studies either assumed that the number of catalogued species was exhaustive<sup>14</sup> (an inaccurate assumption<sup>16</sup>), used local (rather than global) diversity estimates, such as for a small quantity of soil<sup>14</sup>, or estimated the unknown sampling fraction

directly from the phylogeny without additional information (an impossible task<sup>15</sup>). Consequently, there exists no rigorous estimate of global bacterial speciation rates, extinction rates or total diversity over time, and this uncertainty has clouded our interpretation of bacterial evolution over Earth's history. It is commonly hypothesized that bacterial extinction may not even occur at significant rates<sup>11,17–20</sup>, partly due to their large population sizes and wide dispersal ranges<sup>18,19</sup>, while others hypothesized that animal extinctions could cause substantial host-associated bacterial extinctions<sup>21</sup>.

To address these questions, we examined bacterial phylogenies comprising up to hundreds of thousands of clades, using mathematical tools that we developed specifically for large phylogenies. To properly account for the fraction of undiscovered diversity in our methods, thus resolving a long-standing problem in bacterial phylogenetics, we independently estimated global bacterial diversity using massive DNA sequencing data from 60 studies in diverse environments across the world. To evaluate the robustness of our results, we used numerical simulations and examined several phylogenies constructed using alternative methods. Importantly, some of our phylogenies were constructed from environmental sequences retrieved using culture-independent methods, providing a less biased (and thus more suitable<sup>22</sup>) representation of bacterial diversity compared with previous studies<sup>10,14</sup>. We used our methods, as well as the independently estimated global bacterial diversity, to reconstruct global bacterial speciation, extinction and diversification (speciation minus extinction) rates over the past 1 Gyr.

We used two time-calibrated bacterial phylogenies ('timetrees') based on the 16S ribosomal RNA (rRNA) gene—a popular marker gene in microbial ecology and evolution (448,112 and 162,371 tips, respectively; see Supplementary Table 1 for an overview and the Methods for details). We also analysed cyanobacteria alone due to their great importance to Earth's evolution, using four 16S rRNA-based timetrees constructed with various methods (586, 6,308,

<sup>1</sup>Biodiversity Research Centre, University of British Columbia, Vancouver, British Columbia, Canada. <sup>2</sup>Department of Zoology, University of British Columbia, Vancouver, British Columbia, Canada. <sup>3</sup>Joint BioEnergy Institute, Emeryville, CA, USA. <sup>4</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. <sup>5</sup>Department of Plant Biology, University of California, Davis, Davis, CA, USA. <sup>6</sup>Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, CA, USA. <sup>7</sup>Department of Botany, University of British Columbia, Vancouver, British Columbia, Canada. <sup>8</sup>Department of Mathematics, University of British Columbia, Vancouver, British Columbia, Canada.

\*e-mail: [louca.research@gmail.com](mailto:louca.research@gmail.com)

6,302 and 1,579 tips, respectively). In all cases, tips in the trees represent operational taxonomic units (OTUs); that is, clusters in the 16S rRNA gene delineated at 99% similarity—a common microbial ‘species’ measure<sup>23,24</sup>. We stress that bacterial OTUs only provide an approximate ‘species’ analogue to sexually reproducing organisms, and hence ‘speciation’ rates reported here should a priori only be interpreted as branching frequencies in 16S rRNA sequence space.

### Estimating diversification dynamics from large timetrees

Our methods were derived from standard stochastic models for cladogenesis, in which extant lineages can split or go extinct randomly and independent of each other as time proceeds<sup>25</sup>. These models predict the total number of extant lineages (total diversity) at each time point, as well as the number of lineages represented in the final timetree comprising only extant and sampled taxa (lineages through time (LTT))<sup>7</sup>. Our methods can account for the effect of incomplete taxon sampling, as well as for speciation and extinction rates that vary over time. In contrast to most existing methods, our methods consider timetrees in the continuum limit of infinitely many lineages, which yields novel ways to extract information from timetrees (Supplementary Information section 1.3). Notably, given some LTT curve, one can estimate a quantity that is related to the diversification rate at each time point, and which we refer to as the pulled diversification rate (PDR):

$$r_p = \lambda - \mu - \frac{1}{\lambda} \frac{d\lambda}{dt} \quad (1)$$

where  $t$  is time,  $\lambda$  is the instantaneous speciation rate and  $\mu$  is the instantaneous extinction rate. The PDR partly resembles the diversification rate ( $r = \lambda - \mu$ ), but is modified (pulled) by the term  $\lambda^{-1}d\lambda/dt$ , which represents the relative rate of change of  $\lambda$  over time and is small when  $\lambda$  varies slowly. In contrast to the diversification rate, the PDR can be estimated ‘non-parametrically’ from the curvature and slope of the LTT curve at any point in time. This approach does not require fitting a specific parameterized model<sup>25</sup>, nor a priori assumptions on how  $\lambda$  and/or  $\mu$  vary over time, nor assumptions about whether the PDR or diversification rate was positive or negative. More precisely, in the continuum limit, the PDR can be calculated using the LTT for any time  $t$  using the formula:

$$r_p(t) = -\tilde{\nu}(t) - \frac{1}{\tilde{\nu}(t)} \frac{d\tilde{\nu}}{dt} \quad (2)$$

where  $-\tilde{\nu}(t) = (1/\tilde{N}(t))d\tilde{N}/dt$  is the relative slope of the LTT and  $\tilde{N}(t)$  is the value of the LTT at time  $t$ . For finite trees, equation (2) is only an estimate.

Similar to the PDR, one can also estimate ‘pulled’ versions of other important variables, including the pulled extinction rate (PER),

$$\mu_p = \mu + (\lambda_o - \lambda) + \frac{1}{\lambda} \frac{d\lambda}{dt} \quad (3)$$

and the pulled total diversity (PTD),

$$N_p = N \frac{\lambda_o}{\lambda} \quad (4)$$

(estimation formulas provided in Supplementary Information section 1.3). Here,  $\lambda_o$  refers to the most recent speciation rate (that is, as observed near the tips of the tree) and  $N$  is the total diversity at any given point in time. The PER and PTD are equal to the extinction rate  $\mu$  and the total diversity  $N$ , respectively, when  $\lambda$  is constant ( $\lambda = \lambda_o$ ). If  $\lambda$  varies slowly ( $\lambda^{-1}d\lambda/dt \ll \mu$ ), the recent  $\mu_p$  still resembles the recent extinction rate, although the difference

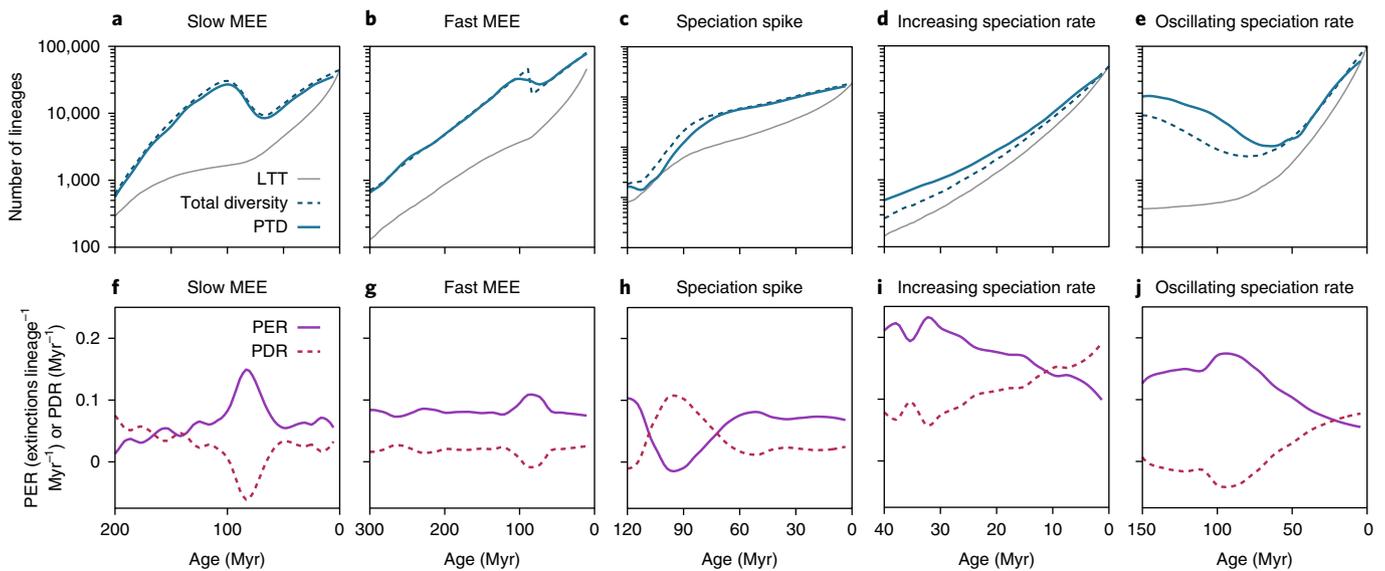
increases for older ages. Rapid variations in  $\lambda$  and/or  $\mu$  will usually lead to substantial variations in  $\mu_p$  and  $r_p$ .

In contrast to conventional maximum-likelihood or Bayesian methods<sup>26,27</sup> for estimating  $\lambda$ ,  $\mu$ ,  $r$  and  $N$ , the pulled variables  $\mu_p$ ,  $r_p$  and  $N_p$  can be estimated from the LTT for each past time point without any assumptions about how  $\lambda$  and  $\mu$  varied over time, and without fitting a specific parameterized model. Model fitting is the current de facto standard in phylogenetics-based reconstruction of diversification<sup>25,28</sup>, and is, in fact, included in the present study. However, it requires that a parameterized form be specified beforehand for  $\lambda$  and  $\mu$ ; for example, accounting for rate shifts at discrete time points, leading to well-known trade-offs between model realism and temporal resolution on the one hand versus model simplicity and confidence in parameter estimation on the other hand. The caveat is that  $\mu_p$ ,  $r_p$  and  $N_p$  are composite variables, and in general, solely knowing  $\mu_p$ ,  $r_p$  and  $N_p$  does not unambiguously determine the constituents  $\mu$ ,  $\lambda$ ,  $r$  and  $N$ . This limitation can be traced back to the fact that extinction partly erases a clade’s history<sup>29</sup> (further discussion in Supplementary Information section 5).

As we demonstrate here, pulled variables are a powerful tool for obtaining insight into past diversification dynamics and for testing model assumptions. Using timetrees simulated under realistic scenarios, we found that pulled variables can reveal past changes in diversification rates, such as those due to mass extinction events, oscillating speciation rates and short temporary spikes in the speciation rate, as well as diversity-dependent speciation and extinction rates (Fig. 1, Supplementary Figs. 1–6 and Supplementary Information section 2). In particular, our simulations revealed that changes in the speciation and/or extinction rate usually lead to similarly strong changes in  $\mu_p$  and that, reciprocally, a constant  $\mu_p$  over time is a strong indication that both  $\lambda$  and  $\mu$  were constant or varied only slowly over time (details in Supplementary Information section 4). Our simulations also revealed that the magnitude of the PDR is usually comparable to the magnitude of the diversification rate, and in fact, in all of our simulations, the two closely resembled each other. Furthermore, we found that  $N_p$  provides a quick way to roughly estimate past total diversities to order-of-magnitude accuracy (Fig. 1a–e), provided  $\lambda$  does not change drastically over time (that is, by orders of magnitude).

### Estimating extant global bacterial diversity

Estimating speciation and extinction rates and past total diversities from a timetree requires knowledge of the fraction of extant diversity represented in the tree<sup>15</sup>. Substantial uncertainty currently exists regarding the extent of extant bacterial diversity, with estimates ranging from a few million OTUs<sup>30</sup> to trillions of OTUs<sup>31</sup>. To better constrain extant bacterial diversity, we examined 165,422 bacterial OTUs recovered de novo from 16S rRNA sequences amplified from various environments, such as animal guts, the ocean, lakes and soils (60 distinct studies comprising 6,303 samples). De novo OTUs covered ~200 base pairs (bp) in the V4 region of the 16S rRNA gene—a region commonly targeted in microbial taxonomic surveys<sup>32</sup>. We calculated the overlap of these de novo OTUs with SILVA—one of the largest 16S rRNA sequence databases<sup>33</sup>—to estimate the fraction of extant bacterial and cyanobacterial OTUs covered by SILVA, as well as the total number of extant OTUs. Our approach is analogous to traditional mark–recapture approaches for estimating population sizes<sup>34</sup>, whereby the number of individuals found in a second survey (analogous to the number of OTUs in SILVA) is divided by the fraction of individuals marked in a first survey (analogous to our de novo dataset) that were recaptured by the second survey. We found that SILVA covers ~33% of de novo OTUs at 99% similarity. Based on the presence of 448,112 16S rRNA clusters in SILVA (that is, obtained by clustering SILVA’s full-length 16S rRNA sequences at 99% similarity), we estimate that there exist globally ~1.4 million bacterial full-length OTUs (overview in Supplementary Table 2).



**Fig. 1 | Non-parametric methods capture complex diversification scenarios.** **a–e**, LTT (grey continuous curves), total diversities (dashed curves) and non-parametrically estimated PTDs for trees simulated under various realistic scenarios, including a slow mass extinction event  $\sim 80$  Ma (**a**), a fast mass extinction event  $\sim 90$  Ma (**b**), a speciation spike  $\sim 90$  Ma (**c**), a gradually increasing diversity-dependent speciation rate (**d**) and an oscillating speciation rate (**e**). **f–j**, PERs (solid curves) and PDRs (dashed curves), estimated non-parametrically from the LTT in **a–e** under the same scenarios. In all scenarios, PER and diversification rates reveal changes in extinction and/or speciation rates, and PTDs approximately resemble the true total diversities (known in this case, since trees were generated via simulations). In **g**, the estimated PER and PDR are damped due to our noise filter, which blurs short ( $\sim 1$  Ma) fluctuations, although the mass extinction event's footprint is still clearly visible in the LTT and PTD (**b**). For more details and additional simulation examples, see Supplementary Information section 2 and Supplementary Figs. 1–6.

This estimate is robust (in order of magnitude) to variation in the methodology, sequencing depth and datasets used (all estimates are within 1.4–1.9 million; Supplementary Information section 3), and is comparable in order of magnitude to recent estimates by Schloss et al.<sup>30</sup>. Furthermore, based on the number of partial-length clusters in SILVA (that is, obtained by clustering the V4 region only), we estimate that there exist globally  $\sim 451,000$  bacterial partial-length (V4) OTUs.

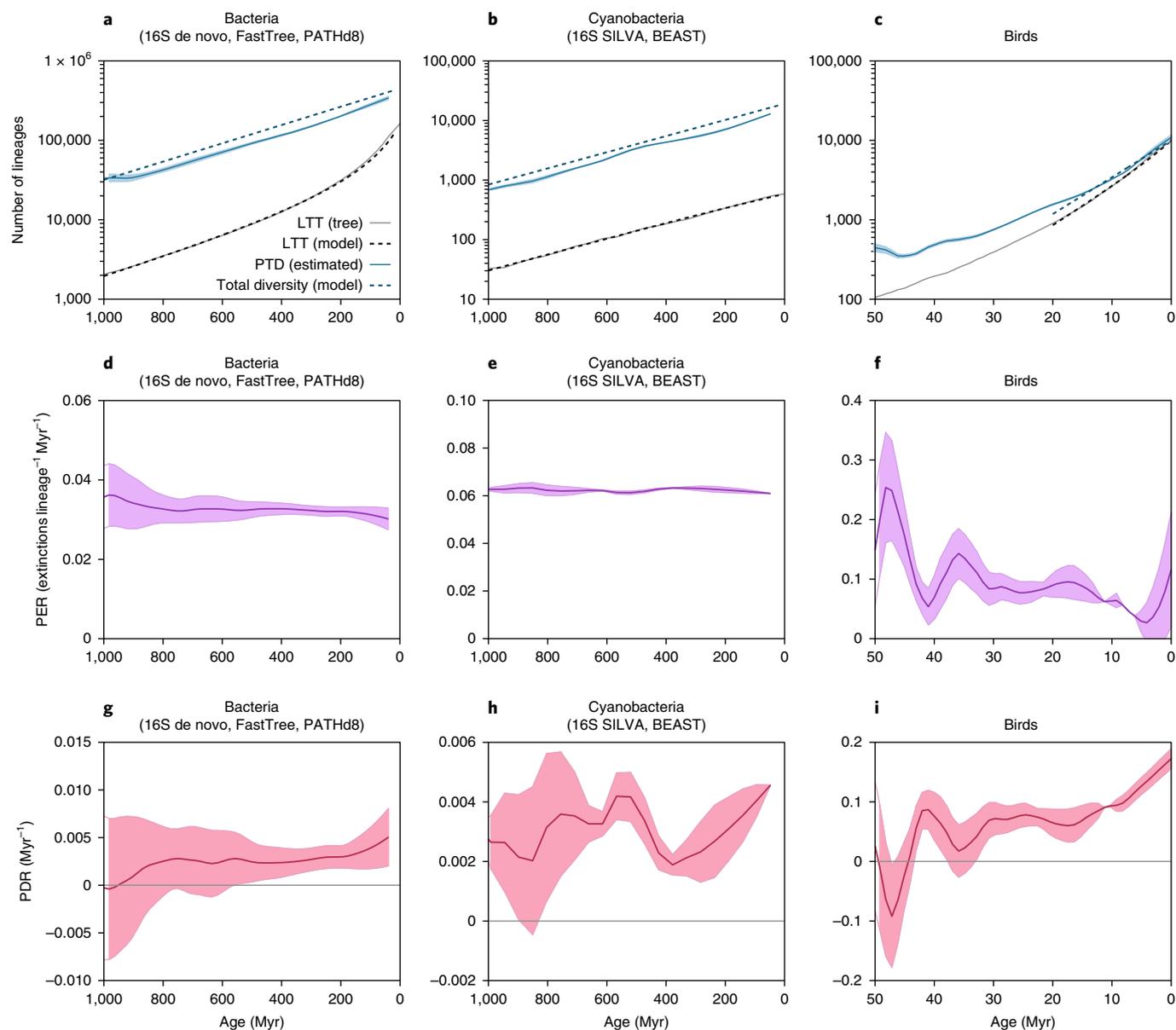
### Estimating bacterial speciation and extinction rates

To estimate recent bacterial speciation and extinction rates, we fitted parametric models to the LTT of the timetrees over a relatively short recent time interval (200 Myr). To gain further insight into past diversification dynamics and to scrutinize model assumptions, we also fitted models over a more extended time interval (1,000 Myr) and used non-parametric methods to estimate PDRs, PERs and PTDs. We found that simple models with constant speciation and extinction rates fitted bacterial and cyanobacterial LTTs well over the past 1 Gyr (mean relative deviation (MRD) below 5% in all cases; Fig. 2a,b and Supplementary Fig. 21a–d). This indicates that overall speciation and extinction rates were roughly constant over time. This conclusion is supported by our observation that fitted speciation and extinction rates change only moderately (by less than 35% in all cases; Supplementary Fig. 18) when models are fitted over the shorter time interval (200 Ma). Our conclusion is also consistent with our estimates that PERs were almost constant over time and almost identical to the extinction rates fitted in the models (less than 10% difference at any time point; Fig. 2d–f and Supplementary Fig. 18). As mentioned previously, a constant PER in and of itself is indicative of constant or only slowly varying speciation and extinction rates, because a rapidly and strongly varying speciation and/or extinction rate usually results in a varying PER (see simulations in Fig. 1 and discussion in Supplementary Information section 2). Based on the small variation observed in the PERs, speciation and extinction rates must have had relative rates of change below

$\sim 0.005 \text{ Myr}^{-1}$  (see explanation in Supplementary Information section 4). In comparison, estimated PERs for birds and vascular plants vary substantially over time (Fig. 2i and Supplementary Fig. 19c, using previously published timetrees<sup>35,36</sup>).

Our findings suggest that, during the past 1 Gyr, global bacterial speciation and extinction rates were not substantially affected during the mass extinction events seen in eukaryotic fossil records<sup>3–5</sup>. This conclusion does not support previous speculations that extinctions of plant- and animal-associated bacteria—resulting from extinction of their hosts—may contribute substantially to bacterial extinction rates<sup>21</sup>. The frequent existence of multiple ecotypes within single OTUs<sup>37,38</sup> may have facilitated bacterial lineage persistence during environmental perturbations and eukaryotic mass extinctions. Even if bacteria experienced extinctions at local scales because of environmental perturbations<sup>39</sup>, these extinctions may have been largely buffered at global scales due to wide dispersal ranges<sup>40</sup>. Our findings also suggest that overall bacterial speciation and extinction rates were not dramatically altered by eukaryotic radiation events, such as the radiation of animals  $\sim 600$  Myr ago (Ma) or the emergence of land plants  $\sim 465$  Ma. It is possible that diversification within individual bacterial clades may have been influenced by eukaryotic radiations and extinctions<sup>11,12</sup>, and that these cases are overshadowed when considering all bacteria together. We also cannot rule out slow effects on speciation and extinction rates (at time scales of billions of years), nor brief fluctuations (shorter than  $\sim 1$  Mya) with little effect on total diversity, both of which could be missed by our methods.

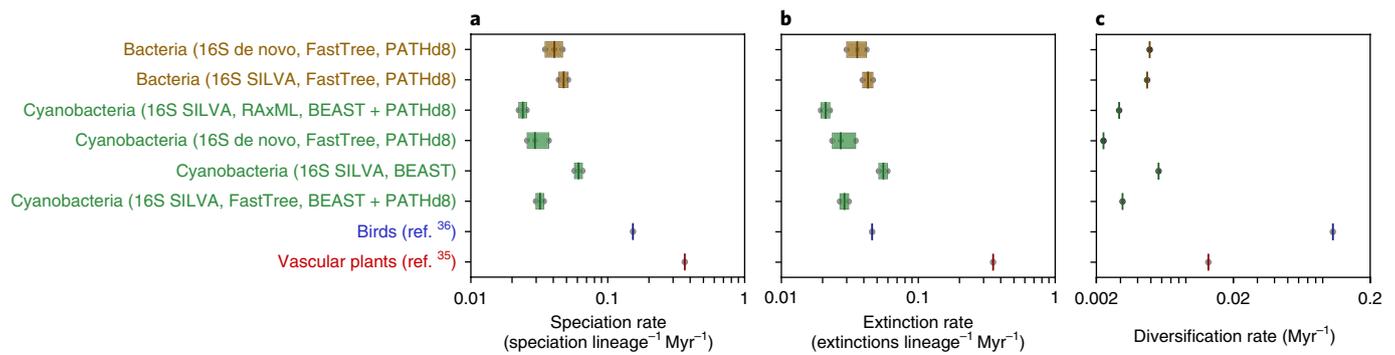
We emphasize that our results do not imply that speciation and extinction rates are homogeneous across clades ('clock-like'). For example, Marin et al.<sup>14</sup> found variable diversification rates across lineages of the Firmicutes bacterial phylum. Using simulations, we found that timetrees, in which speciation and extinction rates are evolving heritable traits and are thus not clock-like, can be fitted well by models with homogenous rates; in these cases, fitted speciation and extinction rates approximately correspond to the



**Fig. 2 | Bacterial, cyanobacterial and bird diversification dynamics through time.** **a–c**, LTT (grey solid curves) for bacteria (**a**), cyanobacteria (**b**) and birds (**c**), compared with speciation–extinction models fitted over the past 1 Gyr (grey dashed curves). Blue solid curves show non-parametrically estimated PTDs, while blue dashed curves show total diversities predicted by the fitted models. Note that each tree only comprises a subset of extant taxa; thus, the LTT does not coincide with total diversity at age 0 (the right-most point on the blue curve). The total diversity at age 0 is only an estimate, based on the fraction of de novo OTUs covered by SILVA (details in main text; overview in Supplementary Table 2). Also note the different time scales shown for birds (50 Myr) compared with bacteria and cyanobacteria (1,000 Myr). **d–f**, PERs (equation (13)), estimated non-parametrically from the same trees as in **a–c**. The roughly constant PERs in **d** and **e** are indicative of constant or only slowly varying speciation and extinction rates, consistent with the fitted models. **g–i**, PDRs (equation (1)), estimated non-parametrically from the same trees, respectively, as in **a–c**. In all panels, shading indicates standard errors of noise-filtered estimates. Summaries of timetree construction methods are indicated in the labels at the top. The bird tree was obtained from Jetz et al.<sup>36</sup>. For analogous figures using alternative timetrees, see Supplementary Fig. 21. For analogous figures for vascular plants, see Supplementary Fig. 19.

average speciation and extinction rates over all lineages (details in Supplementary Information section 4). This means that if bacterial speciation and extinction rates deviated from clock-like behaviour, our clock-like models would not necessarily be able to detect this deviation. Hence, despite our finding of roughly constant overall bacterial speciation and extinction rates over time, we cannot rule out potential differences between clades at any given time point. We also point out that our results only pertain to overall global bacterial diversification and do not distinguish between environments (for example, terrestrial versus marine).

Our fitted models suggest that recent overall extinction rates are 0.03–0.05 extinctions lineage<sup>-1</sup> Myr<sup>-1</sup> for bacteria and 0.02–0.06 extinctions lineage<sup>-1</sup> Myr<sup>-1</sup> for cyanobacteria (Fig. 3b and Supplementary Fig. 18b). These estimates are consistent with estimated recent PERs (Fig. 2d–f and Supplementary Figs. 21d–f and 22a). Our estimates are robust (to an order of magnitude) against variations in the dating of our timetrees (for example, 0.015–0.05 extinctions lineage<sup>-1</sup> Myr<sup>-1</sup> for bacteria; Supplementary Fig. 23). For comparison, using the same methods, we also estimated global extinction rates for vascular plants (~0.35 extinctions lineage<sup>-1</sup> Myr<sup>-1</sup>) and birds



**Fig. 3 | Estimated recent speciation, extinction and diversification rates. a–c.** Recent speciation rates (a), extinction rates (b) and diversification rates (c), estimated for various taxa and using various timetrees (one box per timetree). Estimates were obtained by fitting cladogenic models over a recent time interval of 200 Myr (for estimates over longer time intervals, see Supplementary Fig. 18). For each bacterial or cyanobacterial timetree, various alternative estimates of incomplete sampling fractions were used (Supplementary Tables 2–4); boxes span the results from all alternatives. Tree labels and boxes are coloured by taxon. Summaries of timetree sources and construction methods are indicated in brackets (see Methods for details).

( $\sim 0.05$  extinctions  $\text{lineage}^{-1} \text{Myr}^{-1}$ ), reproducing previous estimates for plants and animals (order of magnitude:  $\sim 0.1$  extinctions  $\text{lineage}^{-1} \text{Myr}^{-1}$ )<sup>41</sup>. We once more point out that bacterial OTUs are only approximately analogous to plant and animal species<sup>24,42</sup>; hence, comparisons between the two should be treated with caution.

We further found that bacterial speciation rates are only slightly above extinction rates—an observation commonly made for larger organisms<sup>29</sup>. Specifically, fitted bacterial diversification rates ( $\sim 0.004$ – $0.005 \text{Myr}^{-1}$ ) are much lower than fitted extinction and speciation rates (Fig. 3c). These values are consistent with similarly low estimated PDRs ( $\sim 0.003$ – $0.004 \text{Myr}^{-1}$ ; Fig. 2g and Supplementary Figs. 21i and 22b). Because bacterial extinction rates are so close to speciation rates, most bacterial lineages that ever existed are now extinct. Based on our fitted models, for each extant bacterial and cyanobacterial OTU there have been  $\sim 10$ – $14$  and  $\sim 5$ – $24$  extinctions over the past 1 Gyr, respectively. The conclusion that only a small fraction of bacterial lineages survived to the present resembles analogous observations for plants and animals<sup>29</sup>. Our finding of substantial bacterial extinction contrasts with reports that diversification within the *Aeromonas* bacterial genus is best explained without extinction<sup>11</sup>, although most analyses yielding zero extinction rates are arguably probably wrong<sup>29</sup>. Nevertheless, at this point, we cannot exclude the possibility that some younger clades (for example, genera) may exhibit much lower extinction rates than the bacterial average.

### Global bacterial diversity increases over time

According to all fitted models, bacterial and cyanobacterial diversification rates have been positive over the past 1 Gyr, suggesting an increase in the total diversity over time. Consistent with this, estimated PDRs are also mostly positive over the past 1 Gyr (Fig. 2g,h and Supplementary Fig. 21g–i), and PTDs ( $N_p$ ) increase roughly exponentially over time (Fig. 2a–c and Supplementary Fig. 21a–c). In principle, a positive PDR and an exponentially increasing  $N_p$  could be a mere result of a decreasing speciation rate over time (that is,  $\lambda > \lambda_0$  in equation (14)), rather than reflecting a truly increasing total diversity. This scenario seems unlikely because it would imply that  $\lambda$  (or the ratio  $\lambda/N$ , if  $N$  also varied) decreased substantially and approximately exponentially over the past 1 Gyr, and that  $\mu$  followed in a similar way (since  $\lambda - \mu \sim 0$ ). It is hard to imagine a simple scenario that would lead to these specific trajectories in  $\lambda$  and  $\mu$  (see discussion in Supplementary Information section 4). Instead, a simpler explanation for an exponentially increasing  $N_p$  is that  $\lambda$  and  $\mu$  were approximately constant and thus  $N$  increased fairly steadily over time. This interpretation is also supported by the fact that when plotted on a logarithmic axis over time, PTDs

exhibit a similar slope to the total diversities predicted by the fitted constant-rate models. A continuous increase in bacterial diversity has been observed previously in a smaller dataset<sup>14</sup>. Similarly, Gubry-Rangin et al.<sup>13</sup> found a stably high diversification rate within the Thaumarchaeota archaeal phylum over the past 400–700 Myr. A continuous increase in bacterial diversity is also comparable to the continuous increase of diversity observed in many eukaryotic taxa during the past 200 Myr<sup>4</sup>. However, we emphasize that diversification rates reconstructed here a priori only reflect 16S branching dynamics and may not reflect ecological diversification<sup>43</sup>.

### Conclusions

Our analysis sheds light on bacterial diversification over geological time. We found evidence that global bacterial diversity has mostly increased over the past 1 Gyr, with roughly constant or only slowly changing overall speciation and extinction rates when averaged over all clades. This conclusion has implications for how life unfolded over Earth's history, since bacteria are the most ancient and the most ubiquitous form of life on Earth<sup>44</sup>. We estimated that global bacterial extinction rates are only slightly below their speciation rates, and that only a small fraction of bacterial lineages that ever existed survived to the present. This has important implications for how we interpret records of ancient life. Some authors have interpreted morphological similarities between microfossils and extant bacterial taxa as signs of 'extreme evolutionary stasis' and absence of speciation and extinction<sup>17,20</sup>, while others even consider cyanobacteria to be living fossils that do not go extinct<sup>19</sup>. Our finding that lineage turnover is an important aspect of bacterial evolution suggests that it is possible that ancient microfossils belonged to extinct lineages, regardless of whether morphology was conserved or convergent<sup>45</sup>, although these extinct lineages may be stem lineages of extant groups. In a similar fashion, it is possible that some ancient molecular biomarkers, such as fossil lipids<sup>46</sup>, were produced by lineages that have gone extinct.

Our work extends various empirical palaeontological 'laws' of macrobial evolution<sup>29</sup> to bacteria—namely, that: extinction is an integral part of evolution; lineages are short-lived at geological time scales; the number of extinct species far exceeds the number of extant species; and speciation and extinction rates are typically similar in magnitude. Despite the high diversity of extant microorganisms, this diversity only represents a snapshot of the microbial diversity ever to have inhabited our planet.

### Methods

**Estimating the total number of extant OTUs.** To estimate the total number of extant bacterial OTUs, we used two alternative approaches. In the first approach,

a large random set of partial-length OTUs (99% identity in the V4 region of the 16S rRNA gene), recovered de novo from environmental samples, was compared with the SILVA SSU database (release 128)<sup>47</sup>, and the total number of extant OTUs was estimated based on the overlap between the de novo OTUs and SILVA. In the second approach, variable-length 16S rRNA sequences extracted from metagenome-assembled genomes (MAGs)<sup>48</sup> were compared with SILVA, and again the total number of extant OTUs was estimated based on the overlap between the MAG 16S rRNA sequences and SILVA. In both cases, non-redundant (NR99) full-length 16S rRNA sequences in the SILVA database were first clustered at 99% identity using uclust version 1.2.22 (ref. <sup>49</sup>) (options: `-usersort -nucleo`). Furthermore, to assess the potential phylogenetic bias of OTUs represented in SILVA, and how this bias affects the representation of older clades compared with random OTU sampling, we constructed and dated a phylogenetic tree of the de novo OTUs and counted the fraction of lineages in the tree over time that was represented in SILVA. Below, we describe these procedures in detail.

**Generating de novo OTUs.** We downloaded public raw Illumina reads of 16S rRNA gene amplicons (V4 region) from the European Nucleotide Archive (<https://www.ebi.ac.uk/ena>) for 6,303 samples from 60 studies across the globe, including animal guts, marine sediments and water columns, soils, bioreactors, lakes, and phytotelmata (henceforth referred to as the 'de novo dataset'; accession numbers provided in Supplementary File 1). We focused particular effort on the inclusion of soils (1,067 samples), which are thought to host a large fraction of Earth's bacterial diversity<sup>32</sup>. Any paired-end reads were merged using flash version 1.2.11 (ref. <sup>50</sup>) (options: `-min-overlap=20 -max-mismatch-density 0.25 -phred-offset=33 -allow-ouities`). Merged and single-end reads were trimmed and quality-filtered using vsearch version 2.4.3 (ref. <sup>51</sup>), keeping only reads at least 200 bp long after trimming (options: `-fastq_ascii 33 -fastq_minlen 200 -fastq_qmin 0 -fastq_maxee 1 -fastq_truncate 1 -fastq_maxee_rate 0.005 -fastq_stripleft 7`). Samples with more than 200,000 quality-filtered reads were rarefied down to 200,000 reads to reduce the computation time, by randomly picking reads without replacement. Rarefied reads were chimera-filtered de novo, separately for each sample, using vsearch (options: `--abskew 1.9 -mindiv 0.5 -minh 0.1`), yielding 265,640,818 quality-filtered and chimera-filtered reads in total, with a mean length of 262 nucleotides. Pooled reads from all samples were error-filtered and clustered de novo at 99% similarity using cd-hit-otu version 0.0.1 (ref. <sup>52</sup>), yielding 345,229 OTUs. OTUs were subsequently filtered anew for chimeras using vsearch (same options as before), yielding 216,707 OTUs. Lastly, any OTUs found in fewer than 2 samples were omitted to further reduce spurious OTUs, leaving us with 185,620 OTUs for downstream analysis.

We note that the above quality filters were chosen to be quite stringent in order to minimize the recovery of spurious OTUs (for example, stemming from sequencing errors or chimeras); however, this conservatism potentially came at the cost of also removing real OTUs. Minimizing the recovery of spurious OTUs is important for both a correct estimation of global bacterial diversity and improving the quality of the phylogenetic tree generated from the OTUs (see below). We emphasize that falsely omitting real OTUs does not affect our estimation of global bacterial diversity, as long as the inclusion or omission of an OTU is independent of the OTU's presence in SILVA. For similar reasons, while our omission of OTUs found in fewer than two samples may bias our census towards more cosmopolitan OTUs, it should not affect our estimation of global bacterial diversity. Indeed, when we also included OTUs found only in a single sample, our estimate of global bacterial diversity changed by less than 10%.

We point out that our mark-recapture-type approach may have underestimated global extant bacterial diversity if some bacterial OTUs are generally much more difficult to detect than others (for example, if they are only present in very specialized environments). In particular, extreme environments (such as hot springs) are under-represented in our de novo dataset, although these environments host relatively low diversity compared with other environments, such as soils. Future, more exhaustive sequencing studies, including a greater variety of environments, will undoubtedly improve estimates of extant bacterial diversity. An underestimation of extant bacterial diversity in the present study would mean that bacterial extinction rates are even higher than estimated here<sup>15</sup>. However, it would not affect our conclusions regarding the constancy of overall bacterial speciation and extinction rates over time, nor our conclusions regarding the continuous nearly exponential increase in bacterial diversity, unless the missed diversity was strongly phylogenetically biased and clustered within the tree—a scenario we view as unlikely.

**Fraction of de novo OTUs represented in SILVA.** De novo OTUs were taxonomically identified using a consensus approach based on the first 10 hits in SILVA at a similarity threshold of at least 70%. Specifically, OTUs were globally aligned to the SILVA non-redundant (NR99) SSU reference database using vsearch version 2.4.3 (ref. <sup>51</sup>) at a minimum similarity of 70% (options: `--id 0.7 --strand both --iddef 2 --maxaccepts 10 --uc_allhits`), while keeping track of the taxonomies provided by SILVA for each hit. For any given OTU, if at least one hit had a similarity of 100%, all hits with a similarity of 100% were considered candidates for forming a consensus taxonomy. Otherwise, if at least one hit had a similarity of  $\geq 70\%$ , all hits with a similarity of at least  $(s - 3\%)$  were used as candidates for a consensus taxonomy. For any candidate set of hits, the consensus taxonomy was defined as

the taxon at the lowest taxonomic level possible, containing all of the candidate hit taxonomies. If an OTU did not match any SILVA entry at or above 70% similarity, or did not form a consensus taxonomy even at the domain level, it was considered unidentified and was subsequently omitted. A total of 171,816 OTUs could be identified at some taxonomic level. OTUs identified as eukaryotes, chloroplasts and mitochondria were omitted from all subsequent analyses. Taxonomically identified OTUs were matched to the SILVA non-redundant (NR99) set using vsearch (options: `--iddef 2 --strand both`) at a similarity threshold of 99%. For any given focal taxon (for example, bacteria, archaea or cyanobacteria), we estimated the fraction ( $\rho$ ) of extant OTUs represented in SILVA as the fraction of taxonomically identified de novo OTUs that could be aligned to the clustered SILVA database (clustered at 99% identity) at a similarity of  $\geq 99\%$ .

**Total number of extant OTUs (based on overlap with SILVA).** Since de novo OTUs only cover a fraction of the 16S rRNA gene (~200 bp from the V4 hypervariable region), correctly estimating the number of extant partial-length (V4) OTUs requires knowledge of the number of V4 OTUs already contained in SILVA. Therefore, we extracted and clustered the part of 16S rRNA sequences in SILVA corresponding to the region covered by de novo OTUs. Specifically, we aligned de novo OTUs to SILVA using the QIIME script `align_seqs_pynast.py`<sup>53</sup>, and using a reduced set of the SILVA alignments (clustered at 90% similarity) as a template. We then identified the first nucleotide position in the OTU alignments that had a gap fraction below 0.9, and extracted the part of the NR99 SILVA alignments starting at that nucleotide position and extending 200 bp in the 5' → 3' direction (omitting gaps). Extracted partial SILVA sequences were then clustered at 99% identity using uclust version 1.2.22 (ref. <sup>49</sup>), yielding 161,070 bacterial and archaeal partial-length clusters.

The global number of extant V4 OTUs in the focal taxon was estimated as  $N_{V4}/\rho$ , where  $N_{V4}$  is the number of V4 clusters within the focal taxon in SILVA, and  $\rho$  is the previously estimated fraction of extant OTUs within the focal taxon represented in SILVA. To estimate the number of extant full-length OTUs, we multiplied the estimated number of extant V4 OTUs by the ratio  $N_{FL}/N_{V4}$ , where  $N_{FL}$  is the number of full-length clusters within the focal taxon in SILVA. Estimated fractions of V4 OTUs represented in SILVA at 99% similarity, as well as total numbers of extant V4 OTUs and full-length OTUs are listed in Supplementary Table 2. We note that one important assumption of the above estimation method is that the presence or absence of an OTU in the de novo dataset is independent of its presence or absence in SILVA. This assumption is probably approximately met, since OTUs in the de novo dataset were recovered without the use of any reference database and the de novo dataset covers a wide range of environments. We emphasize that this assumption does not imply that SILVA is phylogenetically unbiased (that is, we do not assume that OTUs co-occur in SILVA regardless of their phylogenetic relatedness). In fact, we detected substantial phylogenetic bias in SILVA when comparing the observed representation of deeply branching clades with the expectation under random unbiased sampling (Supplementary Fig. 7). This bias can explain why, despite the overall high fraction of extant OTUs already represented in SILVA, recent studies have discovered new deeply branching clades (for example, phyla) not represented in SILVA<sup>16,48</sup>.

**Fraction of MAG 16S rRNA sequences represented in SILVA.** To further verify our estimates of global extant bacterial and cyanobacterial OTUs using a separate method not constrained by potential primer bias, we also used 16S rRNA sequences from MAGs. Specifically, we downloaded 16S rRNA sequences for 2,853 MAGs<sup>48</sup> from [https://data.ace.uq.edu.au/public/misc\\_downloads/uba\\_genomes/](https://data.ace.uq.edu.au/public/misc_downloads/uba_genomes/) on 25 October 2017. Any sequences shorter than 500 bp were omitted, leaving us with 1,166 sequences for downstream analyses. Sequences were taxonomically identified and globally aligned to the clustered SILVA using the same approach as the de novo OTUs (see above). The global number of extant OTUs in each focal taxon was estimated as  $N_{FL}/\rho$ , where  $N_{FL}$  is the number of full-length clusters within the focal taxon in SILVA, and  $\rho$  is the fraction of MAG 16S rRNA sequences that could be aligned to the clustered SILVA at a similarity threshold of 99%. The results are shown in Supplementary Table 2. While this set of 16S rRNA sequences is much smaller than the de novo OTUs, and obtained from a much smaller set of samples, it can serve as a rough verification of the estimates obtained from de novo OTUs. For bacteria as well as cyanobacteria, estimates based on MAGs (Supplementary Table 3) are similar to estimates based on de novo OTUs (Supplementary Table 2). To assess the robustness of our estimated recent speciation, extinction and diversification rates (Fig. 3 and Supplementary Fig. 18), we performed our analyses based on the total number of extant OTUs estimated from de novo OTUs, as well as from MAGs (Supplementary Tables 2–4).

**Total number of extant OTUs (based on overlap with the Earth Microbiome Project (EMP)).** To obtain an additional independent estimate of the number of extant OTUs, we repeated our analysis by considering the overlap between our de novo OTUs and V4 OTUs recovered from a dataset published by the EMP<sup>32</sup>. Specifically, we downloaded raw reads generated by the EMP based on the run accession numbers provided on the EMP GitHub ([https://github.com/biocore/emp/blob/master/code/download-sequences/download\\_ebi\\_fasta.sh](https://github.com/biocore/emp/blob/master/code/download-sequences/download_ebi_fasta.sh)). Project number ERP010098 was omitted as the sequencing instrument was unspecified.

EMP reads were processed similarly to the de novo dataset described above, with the following differences: the minimum allowed read length (after quality filtering) was reduced to 100bp to accommodate the much shorter EMP reads, and the number of quality-filtered reads per sample was limited to 20,000. This yielded 400,528 chimera-filtered OTUs, representing 195,291,388 reads from 18,034 samples across 44 studies. OTUs were taxonomically identified as before, and any OTUs identified as eukaryotes, chloroplasts or mitochondria were omitted. OTUs found in fewer than 2 samples were also omitted, leaving us with 343,743 taxonomically identified OTUs. To calculate the fraction of EMP OTUs represented by our de novo OTUs, we matched the EMP OTUs against the de novo OTUs using vsearch (options: --iddef 2 --strand both) at a similarity threshold of 99%. The total number of extant V4 OTUs within a focal taxon was estimated as  $N_{\text{in}}/\rho$ , where  $N_{\text{in}}$  is the number of de novo OTUs within the focal taxon, and  $\rho$  is the fraction of EMP OTUs within the focal taxon that could be matched to a de novo OTU. The total number of extant full-length OTUs was estimated as before; that is, by multiplying the estimated number of extant V4 OTUs by the ratio  $N_{\text{FL}}/N_{\text{V4}}$ , where  $N_{\text{FL}}$  and  $N_{\text{V4}}$  are the number of full-length and V4 clusters, respectively, within the focal taxon in SILVA.

**Building a tree from de novo OTUs.** Representative sequences of taxonomically identified bacterial and archaeal de novo OTUs, excluding chloroplasts and mitochondria, and with a length of at least 200 nucleotides, were aligned using the QIIME script `parallel_align_seqs_pynast.py`<sup>53</sup>, and using a reduced set of the SILVA alignments (clustered at 90% similarity) as a template. A total of 171,510 OTUs could be successfully aligned. Alignments were reduced by first removing nucleotide positions with >95% gaps, and then removing the top 5% most entropic nucleotide positions. Taxonomic identities at the domain, phylum, class and order level from the preceding step were used to create split constraints for FastTree<sup>54</sup> by constraining each taxon to be on a single side of a split. Taxa with fewer than two OTUs were omitted from the constraints. A total of 603 constraints were defined. Using the alignments and the taxonomically generated constraints, we constructed a phylogenetic tree with FastTree (options: -spr 4 -gamma -no2nd -constraintWeight 100000). The phylogenetic tree was re-rooted so that bacteria and archaea were split at the root. The resulting tree was then dated with PATHd8 using the following dating anchors:

- GOE (secondary constraint). Most recent common ancestor (MRCA) of Oxyphotobacteria and Melainabacteria constrained by the Great Oxygenation Event<sup>55</sup> and based on a molecular clock analysis by Shih et al.<sup>36</sup> (table 3 therein).
  - Ri. MRCA of Rickettsiales constrained to before the earliest known appearance of mitochondria<sup>57</sup>.
  - CB. MRCA of Chlorobium and Bacteroidetes constrained to before the first known Chlorobium-specific biomarkers<sup>58</sup>.
  - Chr. MRCA of Chromatiaceae constrained to before the first known purple sulfur bacterial biomarkers<sup>46,58</sup>.
  - LUCA. MRCA of archaea and bacteria constrained to before the earliest known stromatolites and after the late heavy bombardment<sup>59,60</sup>.
- The dating anchors are summarized in Supplementary Table 5.

**Fraction of extant lineages represented in SILVA through time.** To assess the phylogenetic bias of SILVA, we calculated the fraction of lineages in the de novo tree represented in SILVA over time (that is, the fraction of discovered lineages (FDL); Supplementary Fig. 7a). Specifically, we extracted the de novo subtree comprising only the OTUs matched to SILVA at 99% similarity (see previous paragraph), counted the remaining LTT and divided those by the LTT in the full de novo tree. Note that this fraction loosely corresponds to the probability that a lineage at any age would be represented in SILVA, provided that it has not gone extinct. To examine how the FDL differed in the case of the null model of random independent sampling of OTUs, we replaced the OTUs represented in SILVA with a new, equally sized set of randomly chosen OTUs and recalculated the FDL. The variability of the outcome was assessed by repeating this procedure 100 times. To further examine, based both on SILVA and the null model, how the FDL through time depends on the fraction of OTUs discovered (sampling fraction), we subsampled SILVA down to various fractions (for example, 10 and 1%) and repeated the previous analysis (Supplementary Fig. 7b,c). We found that the FDL through time deviated substantially from the null model, indicating phylogenetically correlated representation of OTUs in SILVA, with some clades being over-represented or under-represented compared with random OTU sampling. For strong subsampling of SILVA (<1%; Supplementary Fig. 7c), this deviation diminished towards recent ages.

**Comparison with the 97% similarity threshold.** For comparison purposes, we repeated some of the above calculations for de novo OTUs clustered at 97% similarity. Specifically, we re-clustered de novo OTUs at a coarser similarity threshold of 97% using vsearch (options: --cluster\_fast --usersort -id 0.97 --iddef 2 --strand plus) and used the selected centroids as new representative sequences. Taxonomic identification was done as described above for de novo OTUs. Coverage by SILVA, as well as the total number of extant 97% OTUs, were

estimated similarly to above, the difference being that SILVA was clustered at 97% similarity and coverage by SILVA was calculated at 97% similarity. A timetree comprising 31,231 de novo 97% OTUs was constructed as above. We do not discuss these results in the main text, but provide them in Supplementary Table 6.

**Tree construction and dating.** To verify the robustness of our results, we examined a multitude of bacterial and cyanobacterial timetrees, constructed using various alternative methods, described below. Trees were constructed using either full-length 16S rRNA alignments from the SILVA reference database (release 128; non-redundant set)<sup>33</sup> or partial-length alignments of de novo clustered OTUs from public amplicon sequences from various environmental samples. Unless otherwise mentioned, tips in all 16S-based trees corresponded to OTUs clustered at 99% similarity. We note that bacterial OTUs were historically delineated using a similarity threshold of 97%. However, modern genomics revealed that taxa defined on this basis are usually underspecified, and that a greater similarity threshold (>99%) is required for distinguishing ecologically differentiated organisms<sup>25,24,42,61</sup>. In this study, we thus delineated OTUs at 99% similarity, but we also performed comparisons at 97% similarity (provided in Supplementary Figs. 8 and 9, and Supplementary Table 6). We stress that bacterial OTUs (even at the similarity threshold of 99%) only provide an approximate 'species' analogue, and any given OTU may still comprise multiple closely related strains with different genomic contents and ecological strategies<sup>38</sup>. Even formally named bacterial 'species' can display strong genomic and phenotypic strain diversity<sup>37</sup>. Hence, 'speciation' rates reported here probably represent a conservative estimate of the rate at which bacteria differentiate ecologically. Whether and how bacterial species can ever be reasonably defined remains an open question<sup>62</sup>; hence, the 16S rRNA gene remains a popular marker for cataloguing bacterial diversity and describing evolutionary relationships<sup>43</sup> in a well-defined and reproducible manner.

**Birds.** The bird tree was constructed and dated by Jetz et al.<sup>36</sup>, and was downloaded from the project's website (<http://litoria.eeb.yale.edu/bird-tree/archives/Stage2>) on 1 August 2017 ('Hackett\_backbone\_stage2\_tree\_0001.tre'). We assumed a sampling fraction ( $\rho$ ) of 1, since almost all bird species have probably already been discovered<sup>66</sup>.

**Vascular plants.** The vascular plants tree was constructed and dated by Zanne et al.<sup>35</sup>, and was downloaded from the Dryad Digital Repository (<http://datadryad.org/resource/doi:10.5061/dryad.63q27.2>). We assumed a sampling fraction of 0.724, according to estimates by Mora et al.<sup>67</sup> on the fraction of plant species discovered.

**Bacteria (16S SILVA, FastTree, PATHd8).** Non-redundant, full-length 16S rRNA gene alignments of 448,112 bacteria, 3 chloroplasts and one archaeon *Methanococcales* (for dating purposes), representing OTUs at 99% similarity, were extracted from SILVA. Alignments were reduced by first removing nucleotide positions with >95% gaps, and then removing the top 5% most entropic nucleotide positions. SILVA taxonomies provided for OTUs were used to define topological constraints for FastTree at the domain, phylum, class or order level, by constraining each taxon to be a monophyletic group in the new tree. Taxa with fewer than two OTUs were not constrained. A total of 625 constraints were defined. Using the reduced alignments, with the SILVA guide tree as a starting tree and using the taxonomic constraints, a new tree was generated with FastTree (options: -spr 4 -gamma -no2nd). The generated tree was re-rooted such that bacteria and archaea split at the root. The re-rooted tree was dated with PATHd8 using the anchors GOE, Chl, Ri, CB, Chr and LUCA, as listed in Supplementary Table 5. All archaea, chloroplasts and mitochondria were subsequently removed from the dated tree. An overview of the tree is shown in Supplementary Fig. 10.

**Bacteria (16S SILVA, 97%, FastTree, PATHd8).** This tree was created similarly to the previous bacterial tree (16S SILVA, FastTree, PATHd8), but with OTUs clustered at 97% similarity. An overview of the tree is shown in Supplementary Fig. 11.

**Bacteria (16S de novo, FastTree, PATHd8).** As described above, a dated phylogenetic tree comprising bacterial and archaeal partial-length OTUs (99% identity in the V4 region) was constructed from de novo clustered 16S rRNA gene amplicon sequences from a wide range of environments. From this 'de novo' tree, we extracted the subtree comprising those OTUs identified as bacterial. An overview of the tree is shown in Supplementary Fig. 12.

**Bacteria (16S de novo, 97%, FastTree, PATHd8).** This tree was created similarly to the previous bacterial tree (16S de novo, FastTree, PATHd8), but with OTUs clustered at 97% similarity. An overview of the tree is shown in Supplementary Fig. 13.

**Cyanobacteria (16S SILVA, FastTree, BEAST + PATHd8).** This tree was constructed with 16S rRNA sequences from SILVA, then dated using secondary constraints inferred from a previously dated multigenic cyanobacterial tree<sup>66</sup>, as follows. Non-redundant full-length 16S rRNA alignments of all non-chloroplast cyanobacteria and representative chloroplasts were extracted from SILVA. Alignments were pre-processed and used to construct a tree with FastTree, in the same way as described above for the bacteria (16S SILVA, FastTree, PATHd8). The generated tree was

re-rooted such that the root separates the Melainabacteria from the rest of the tree<sup>56</sup>. Next, a previously published dated multigene tree, including 60 cyanobacterial and 37 plastid taxa, was obtained from Shih et al.<sup>56</sup> (run T65). The multigene tree is based on full-length 16S rRNA gene sequences and 10 additional marker genes, and was dated using BEAST. To link tips in our 16S rRNA-based tree to tips in the multigene tree, the original SILVA 16S rRNA alignments were de-aligned (all gap characters removed) and mapped to the 16S rRNA sequences of the strains in the multigene tree, via global alignment using vsearch (options: --id 1.0 --iddef 2 --strand both --maxhits 1 --maxaccepts 1). A total of 60 tips could be mapped. From this mapping, and based on the divergence times of nodes in the multigene tree, secondary dating constraints were generated for the 16S rRNA tree using the congruency method by Eastman et al.<sup>68</sup>. This method was developed for generating very large timetrees based on a smaller previously dated 'reference' timetree (in our case, the multigene tree), by identifying nodes that are concordant in the target tree (in our case, the non-dated 16S rRNA-based tree) and the reference tree. Using the congruency method, which was performed with the R package *castor*<sup>69</sup>, a total of 17 concordant node pairs were identified. The divergence times of concordant nodes were then used as fixed constraints for dating the 16S rRNA-based tree using PATHd8. All chloroplasts were subsequently removed from the dated tree. This yielded a dated tree for 6,308 non-chloroplast cyanobacteria. An overview of the tree is shown in Supplementary Fig. 14.

**Cyanobacteria (16S SILVA, BEAST).** Non-redundant 16S rRNA alignments of non-chloroplast cyanobacteria and representative chloroplasts were subsampled randomly down to 586 OTUs (including 30 chloroplasts), and a tree was constructed and dated using BEAST, as follows. A log-normal relaxed molecular clock model was implemented using BEAST with the GTR+G substitution model on the 16S rRNA dataset. Chloroplast taxa from land plants were constrained to a normal distribution of  $477 \pm 70$  Ma based on Smith et al.<sup>70</sup>. A uniform prior on the base of crown group cyanobacteria ranged from 1,909–2,450 Ma<sup>56</sup>. The younger boundary of the age constraint is based on the conservatively younger age reported by Shih et al.<sup>56</sup>, while the older boundary represents a geological constraint on the origins of oxygenic photosynthesis and cyanobacteria based on the Great Oxygenation Event. A uniform prior was used to enable flexibility by allowing the MCMC search to agnostically converge on a date that best fit the data. Two separate MCMC chains were generated for 50 million generations, sampling every 10,000 generations, with the first 20 million generations discarded as burn-in. An overview of the tree is shown in Supplementary Fig. 15.

**Cyanobacteria (16S SILVA, RAXML, BEAST + PATHd8).** Non-redundant full-length 16S rRNA gene alignments of 6,302 cyanobacteria and 28 representative chloroplasts, representing OTUs at 99% similarity, were extracted from SILVA. Alignments were reduced by first removing nucleotide positions with >95% gaps, and then removing the top 5% most entropic nucleotide positions. Using the reduced alignments, and using the SILVA guide tree as a starting tree, a new tree was generated with RAXML Stamatakis2014RAXML (options: -m GTRCAT -p 34612 -f d). The generated tree was re-rooted such that the root separates the Melainabacteria from the rest of the tree<sup>56</sup>. The re-rooted tree was dated using PATHd8, based on secondary constraints extracted from a previously published multigene timetree<sup>56</sup>, as described above (cyanobacteria, '16S SILVA, FastTree, BEAST + PATHd8'). All chloroplasts were subsequently removed from the dated tree. An overview of the tree is shown in Supplementary Fig. 16.

**Cyanobacteria (16S de novo, FastTree, PATHd8).** This tree was extracted from the de novo tree (see Methods), similarly to the bacterial (16S de novo, FastTree, PATHd8) tree. An overview of the tree is shown in Supplementary Fig. 17.

Some bacterial trees were constructed and dated simultaneously using BEAST<sup>53</sup>, and some trees were first constructed using FastTree version 2.1.10 (ref. <sup>54</sup>) or RAXML version 8.2.9 (ref. <sup>64</sup>) and subsequently dated using PATHd8 version 1.0 (ref. <sup>65</sup>), depending on computational feasibility. Some trees were dated using primary dating anchors (summarized in Supplementary Table 5), while others were dated using secondary dating anchors extracted from timetrees previously constructed with BEAST. Note that, because PATHd8 (ref. <sup>65</sup>) requires at least one anchor with a fixed age, for timetrees dated using PATHd8 and primary anchors, the GOE anchor (split between Oxycyphobacteria and Melainabacteria) was fixed to an age of 2.55 Ga<sup>56</sup>. Below, we describe the source or construction method for each timetree in detail. An overview of all considered timetrees is provided in Supplementary Table 1.

Sampling fractions of trees ( $\rho$ ) were calculated by dividing the number of tips in each tree by the total number of extant full-length OTUs, extant partial-length (V4) OTUs or extant species (whichever was appropriate) in the corresponding taxon (as estimated in this study; overview in Supplementary Tables 2–4). Figure 2 and Supplementary Fig. 21 were generated using the estimates listed in Supplementary Table 2. The same approach for fitting speciation and extinction models, and estimating speciation and extinction rates (see Methods), was applied to all timetrees. The timetrees analysed in Fig. 2 are bacteria (16S de novo, FastTree, PATHd8), cyanobacteria (16S SILVA, BEAST) and vascular plants<sup>35</sup>. Analogous results are shown for additional timetrees in Supplementary Figs. 19 and 21.

Estimated recent speciation, extinction and diversification rates are summarized in Fig. 3 and Supplementary Fig. 18. As seen in Fig. 3 and Supplementary Fig. 18, all bacterial and cyanobacterial timetrees yielded similar estimates for speciation, extinction and diversification rates. This reproducibility underlines the robustness of our estimates, despite potential inaccuracies in tree construction due to short sequence alignments (in the case of the de novo dataset) and due to heuristic algorithms (for example, FastTree) used for some of the trees for computational feasibility. This result may not be surprising given that our estimates are entirely based on the LTT curve, which is a high-level summary statistic and which, for larger trees, may be rather invariant to uncertainties in tree topology.

We mention that very few geological anchors are currently available for dating bacterial phylogenies and, in principle, 16S rRNA nucleotide substitution rates could vary strongly between clades<sup>71</sup>. This variation could therefore introduce errors when translating phylogenetic distances to temporal distances. However, such errors are generally expected to further increase the deviations of the resulting trees from the simple cladogenic models fitted here. Hence, our conclusion that constant speciation and extinction rates are adequate models for global-scale bacterial diversification dynamics over geological time, as discussed in the main text, is actually conservative.

**Comparing tree topologies.** To quantify the variation in tree topologies obtained using the different tree construction methods, and to compare that variation with previously published trees, we proceeded as follows. In all cases, trees were either pruned to only include bacteria (excluding chloroplasts and mitochondria) or cyanobacteria (excluding chloroplasts), as appropriate. A 16S rRNA gene-based and manually curated guide tree (release 128, non-redundant set) was downloaded from the SILVA database<sup>35</sup>, and a tree of amplicon sequences previously published by the EMP (release 1, 'deblurred'; 150-bp sequences)<sup>35</sup> was obtained from [ftp://ftp.microbio.me/emp/release1/otu\\_info/deblur/emp150.5000\\_1000\\_rxb\\_placement\\_pruned75.tog.tre](ftp://ftp.microbio.me/emp/release1/otu_info/deblur/emp150.5000_1000_rxb_placement_pruned75.tog.tre). Each pair of trees (among our trees, the SILVA guide tree and the EMP tree) was compared after pruning trees to the set of tips shared by both trees, and using the Robinson–Foulds metric<sup>72</sup>. The Robinson–Foulds metric is widely used for comparing tree topologies, based on the number of tip clusters (sets of tips descending from internal nodes) that are unique to each tree. Robinson–Foulds distances were calculated using the R package *castor* version 1.3.3 (ref. <sup>69</sup>) (option: normalized=TRUE). To match tips across trees, all tips were renamed to SILVA sequence accessions whenever possible. Tips in the EMP tree, as well as tips in our de novo trees, were mapped to SILVA via global alignment at a similarity threshold of 99.5% using vsearch (options: --iddef 2 --strand both --maxhits 1 --maxaccepts 1)<sup>35</sup>; any unmapped tips were omitted from the comparisons. In cases where multiple tips matched the same SILVA entry, the nearest match was kept. The number of tips considered in each tree comparison, and calculated pairwise tree distances, are listed in Supplementary Table 7. To visualize relative distances between trees, we used multidimensional scaling ordination plots<sup>73</sup>, in which each tree is represented by a single point, and where the distance between any two points approximately corresponds to the original Robinson–Foulds distance (Supplementary Fig. 24). As can be seen in Supplementary Fig. 24, our trees fall within the typical range of variation of trees of this magnitude.

**Sensitivity analysis with respect to dating.** To assess the sensitivity of our rate estimates to uncertainties in tree dating, we analysed variants of our timetrees created by randomly varying dating constraints. Specifically, we created random variants of our bacterial and cyanobacterial PATHd8-dated timetrees for which we had originally used the primary dating anchors listed in Supplementary Table 5 (bacterial '16S SILVA, FastTree, PATHd8'; bacterial '16S de novo, FastTree, PATHd8' and cyanobacterial '16S de novo, FastTree, PATHd8'). For each random variant, we chose ages randomly and independently within the original age intervals of the dating anchor according to a triangular distribution, whose mode was set to the anchor node's calibrated age in the original timetree. These randomly drawn ages (one per anchor node) were then used as fixed dating constraints for dating the molecular phylogeny anew with PATHd8, thus obtaining a random variant. For each original timetree, we created ten random variants (see Supplementary Fig. 23 for example LTT calculated from these variants). For each timetree, we used its random variants to estimate speciation, extinction and diversification rates using the same methods as for the original tree. This yielded, for each timetree, a set of ten slightly different rate estimates (Supplementary Fig. 25), whose spread can be seen as a measure for estimation uncertainty due to errors in tree dating.

**Fitting models and estimating speciation, extinction and diversification rates.** Parameterized speciation and extinction models were fitted to the LTT of each considered timetree using the general approach described in Supplementary Information section 1.2. Models were fitted for three purposes: (1) to estimate recent speciation rates  $\lambda(\tau=0)$ ; (2) to assess whether a constant speciation and extinction rate provide an adequate description of the observed LTT within some sufficiently small age interval; and (3) to predict past total diversity,  $N(\tau)$ , using the fitted model within the considered age interval. Models were integrated backwards in time using the differential equation listed in Supplementary Information section 1.2 and with the initial condition  $N(0) = N(0)/\rho$ , where  $\rho$  is the

previously estimated sampling fraction (fraction of extant OTUs or species included in the tree) and  $\tilde{N}(0)$  is the number of tips in the tree. Model parameters were fitted by minimizing the MRD of the model's predicted LTT ( $\tilde{N}_m$ ) from the real LTT:

$$\text{MRD} = \frac{1}{n} \sum_{i=1}^n \frac{1}{\tilde{N}(\tau_i)} |\tilde{N}_m(\tau_i) - \tilde{N}(\tau_i)| \quad (5)$$

where  $\tau_1, \dots, \tau_n$  are discrete ages at which the model's predicted LTT is compared with the real LTT. This fitting objective, which is based on relative rather than absolute errors, was chosen so as to increase the importance of earlier time points in the tree, where the LTT can be orders of magnitude lower than at the tips. The  $\tau_i$  were chosen on a regular grid, comprising 100–200 points (depending on the size of the tree) and spanning the minimum and maximum ages considered for the tree. We avoided the most recent part of the LTT, where incomplete and phylogenetically biased taxon sampling can lead to deviations from the assumption of uncorrelated speciations, extinctions and discoveries, and where the choice of the OTU similarity threshold strongly affects branching frequencies (and thus the slope of the LTT). The last few time points (up to 20 Myr; overview in Supplementary Table 1) were therefore ignored during model fitting, in accordance with common practice<sup>74</sup>. For bacteria, the omitted age interval corresponds to ~1–2% divergence in the 16S rRNA gene<sup>75,76</sup>, and hence to one or two expected branching events in the timetrees. Age intervals considered for fitting are listed in Supplementary Table 1. We fitted each model by minimizing the MRD using the optimization function `stats::nlminb` in R. We repeated the optimization 1,000 times with random start values to avoid non-global local optima. This complete approach has been implemented in the R package `castor`<sup>69</sup>. All models assumed constant speciation and extinction rates,  $\lambda$  and  $\mu$ , which were fitted as described above. All fitted bacterial and cyanobacterial models achieved very good agreement with the observed LTT (MRD below 5% in all cases; overview in Supplementary Table 1).

Furthermore, to validate model assumptions and gain additional insight into past diversification dynamics, we estimated PERs ( $\mu_p$ ), PDRs ( $r_p$ ) and PTDs ( $N_p$ ) using the non-parametric methods described in Supplementary Information section 1.3. Before any non-parametric estimation of pulled variables (which involves derivatives of the LTT), the LTT was noise-filtered (smoothened) using a quadratic Savitzky–Golay filter. A noise filter is essential before estimating derivatives from the LTT, because finite-difference derivative estimators tend to amplify high-frequency noise in time series. Estimated PTDs, PERs and PDRs were also smoothened using local polynomial regression fitting (LOESS) to reduce noise, using the R function `msir::loess.sd` (`span` 0.2, `degree` 2)<sup>77</sup>. Standard errors of the smoothened estimates were calculated from the confidence intervals provided by `loess.sd`. As discussed in the main text, we found that  $\mu_p$  was almost constant over time for all prokaryotic trees examined (Fig. 2d–f), supporting the assumption of a roughly constant (or only slowly varying)  $\lambda$  and  $\mu$  made in the models (see discussion in Supplementary Information section 4).

Here, we have restricted our analyses to the most recent 1 Gyr because, for older time points, the smaller number of lineages in the tree, and thus the greater stochasticity and deviation from the continuum limit, lead to increased uncertainties in the estimated diversification dynamics (Supplementary Fig. 23). While future larger phylogenies will allow more accurate reconstruction of diversification dynamics for even more ancient times, we generally advise against using our non-parametric methods near the origin of a tree or clade of interest.

**Reporting Summary.** Further information on experimental design is available in the Nature Research Reporting Summary linked to this article.

**Code availability.** The R script for performing the diversification analyses on the timetrees, as well as the simulations discussed in Supplementary Information sections 2 and 4, is included as Supplementary File 3. The non-parametric methods introduced in the manuscript are implemented in the R package `castor`—a package for efficient phylogenetic analyses on very large trees<sup>69</sup> available on The Comprehensive R Archive Network (CRAN).

**Data availability.** Amplicon sequencing data used to recover de novo OTUs are publicly available under the accession numbers listed in Supplementary File 1. Accession numbers for sequencing data used from the EMP<sup>32</sup> are listed in Supplementary File 2. R code used in this study is provided as Supplementary File 3. Timetrees and undated phylogenetic trees constructed in this study are provided as Supplementary File 4. Taxonomic classifications of de novo OTUs are provided as Supplementary File 5.

Received: 1 February 2018; Accepted: 28 June 2018;  
Published online: 30 July 2018

## References

- Falkowski, P. G., Fenchel, T. & Delong, E. F. The microbial engines that drive Earth's biogeochemical cycles. *Science* **320**, 1034–1039 (2008).
- Fischer, W. W., Hemp, J. & Johnson, J. E. Evolution of oxygenic photosynthesis. *Annu. Rev. Earth Planet. Sci.* **44**, 647–683 (2016).
- Raup, D. M. & Sepkoski, J. J. Mass extinctions in the marine fossil record. *Science* **215**, 1501–1503 (1982).
- Signor, P. W. Biodiversity in geological time. *Am. Zool.* **34**, 23–32 (1994).
- McElwain, J. C. & Punyasena, S. W. Mass extinction events and the plant fossil record. *Trends Ecol. Evol.* **22**, 548–557 (2007).
- Nee, S., May, R. M. & Harvey, P. H. The reconstructed evolutionary process. *Phil. Trans. R. Soc. Lond. B* **344**, 305–311 (1994).
- Nee, S., Holmes, E. C., May, R. M. & Harvey, P. H. Extinction rates can be estimated from molecular phylogenies. *Phil. Trans. R. Soc. Lond. B* **344**, 77–82 (1994).
- Sanderson, M. J. & Donoghue, M. J. Reconstructing shifts in diversification rates on phylogenetic trees. *Trends Ecol. Evol.* **11**, 15–20 (1996).
- Morlon, H. Phylogenetic approaches for studying diversification. *Ecol. Lett.* **17**, 508–525 (2014).
- Morlon, H., Kems, B. D., Plotkin, J. B. & Brisson, D. Explosive radiation of a bacterial species group. *Evolution* **66**, 2577–2586 (2012).
- Lorén, J. G., Farfán, M. & Fusté, M. C. Molecular phylogenetics and temporal diversification in the genus *Aeromonas* based on the sequences of five housekeeping genes. *PLoS ONE* **9**, 1–15 (2014).
- Lebreton, F. et al. Tracing the Enterococci from paleozoic origins to the hospital. *Cell* **169**, 849–861 (2017).
- Gubry-Rangin, C. et al. Coupling of diversification and pH adaptation during the evolution of terrestrial Thaumarchaeota. *Proc. Natl Acad. Sci. USA* **112**, 9370–9375 (2015).
- Marin, J., Battistuzzi, F. U., Brown, A. C. & Hedges, S. B. The timetree of prokaryotes: new insights into their evolution and speciation. *Mol. Biol. Evol.* **34**, 437–446 (2017).
- Stadler, T. On incomplete sampling under birth–death models and connections to the sampling-based coalescent. *J. Theor. Biol.* **261**, 58–66 (2009).
- Hug, L. A. et al. A new view of the tree of life. *Nat. Microbiol.* **1**, 16048 (2016).
- Schopf, J. W. Disparate rates, differing fates: tempo and mode of evolution changed from the Precambrian to the Phanerozoic. *Proc. Natl Acad. Sci. USA* **91**, 6735–6742 (1994).
- Dykhuizen, D. E. Santa rosalia revisited: why are there so many species of bacteria? *Antonie Van Leeuwenhoek* **73**, 25–33 (1998).
- Butterfield, N. Macroevolution and macroecology through deep time. *Palaeontology* **50**, 41–55 (2007).
- Schopf, J. W. et al. Sulfur-cycling fossil bacteria from the 1.8-Ga duck creek formation provide promising evidence of evolution's null hypothesis. *Proc. Natl Acad. Sci. USA* **112**, 2087–2092 (2015).
- Weinbauer, M. G. & Rassoulzadegan, F. Extinction of microbes: evidence and potential consequences. *Endang. Species Res.* **3**, 205–215 (2007).
- Höhna, S., Stadler, T., Ronquist, F. & Britton, T. Inferring speciation and extinction rates under different sampling schemes. *Mol. Biol. Evol.* **28**, 2577–2589 (2011).
- Stackebrandt, E. & Ebers, J. Taxonomic parameters revisited: tarnished gold standards. *Microbiol. Today* **33**, 152–155 (2006).
- Edgar, R. C. Updating the 97% identity threshold for 16S ribosomal RNA OTUs. *Bioinformatics* **34**, 2371–2375 (2018).
- Sanmartín, I. & Meseguer, A. S. Extinction in phylogenetics and biogeography: from timetrees to patterns of biotic assemblage. *Front. Genet.* **7**, 35 (2016).
- Stadler, T. Mammalian phylogeny reveals recent diversification rate shifts. *Proc. Natl Acad. Sci. USA* **108**, 6187–6192 (2011).
- Silvestro, D., Schnitzler, J. & Zizka, G. A Bayesian framework to estimate diversification rates and their variation through time and space. *BMC Evol. Biol.* **11**, 311 (2011).
- Stadler, T. How can we improve accuracy of macroevolutionary rate estimates? *Syst. Biol.* **62**, 321–329 (2013).
- Marshall, C. R. Five palaeobiological laws needed to understand the evolution of the living biota. *Nat. Ecol. Evol.* **1**, 165 (2017).
- Schloss, P. D., Girard, R. A., Martin, T., Edwards, J. & Thrash, J. C. Status of the archaeal and bacterial census: an update. *mBio* **7**, e00201-16 (2016).
- Locey, K. J. & Lennon, J. T. Scaling laws predict global microbial diversity. *Proc. Natl Acad. Sci. USA* **113**, 5970–5975 (2016).
- Thompson, L. R. et al. A communal catalogue reveals Earth's multiscale microbial diversity. *Nature* **551**, 457–463 (2017).
- Glöckner, F. O. et al. 25 years of serving the community with ribosomal RNA gene reference databases and tools. *J. Biotechnol.* **261**, 169–176 (2017).
- Krebs, C. J. *Ecological Methodology* (Benjamin Cummings, San Francisco, CA, 1999).
- Zanne, A. E. et al. Three keys to the radiation of angiosperms into freezing environments. *Nature* **505**, 89–99 (2014).
- Jetz, W. et al. Global distribution and conservation of evolutionary distinctness in birds. *Curr. Biol.* **24**, 919–930 (2014).
- Welch, R. A. et al. Extensive mosaic structure revealed by the complete genome sequence of uropathogenic *Escherichia coli*. *Proc. Natl Acad. Sci. USA* **99**, 17020–17024 (2002).

38. Shapiro, B. J. & Polz, M. F. Ordering microbial diversity into ecologically and genetically cohesive units. *Trends Microbiol.* **22**, 235–247 (2014).
39. Xie, S., Pancost, R. D., Yin, H., Wang, H. & Evershed, R. P. Two episodes of microbial change coupled with Permo/Triassic faunal mass extinction. *Nature* **434**, 494–497 (2005).
40. Gibbons, S. M. et al. Evidence for a persistent microbial seed bank throughout the global ocean. *Proc. Natl Acad. Sci. USA* **110**, 4651–4655 (2013).
41. Pimm, S. L. et al. The biodiversity of species and their rates of extinction, distribution, and protection. *Science* **344**, 1246752 (2014).
42. Kim, M., Oh, H.-S., Park, S.-C. & Chun, J. Towards a taxonomic coherence between average nucleotide identity and 16S rRNA gene sequence similarity for species demarcation of prokaryotes. *Int. J. Syst. Evol. Microbiol.* **64**, 346–351 (2014).
43. Straub, T. J. & Zhaxybayeva, O. A null model for microbial diversification. *Proc. Natl Acad. Sci. USA* **114**, E5414–E5423 (2017).
44. Whitman, W. B., Coleman, D. C. & Wiebe, W. J. Prokaryotes: the unseen majority. *Proc. Natl Acad. Sci. USA* **95**, 6578–6583 (1998).
45. Butterfield, N. J. Proterozoic photosynthesis—a critical review. *Palaeontology* **58**, 953–972 (2015).
46. Brocks, J. J. & Banfield, J. Unravelling ancient microbial history with community proteogenomics and lipid geochemistry. *Nat. Rev. Microbiol.* **7**, 601–609 (2009).
47. Quast, C. et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res.* **41**, D590–D596 (2013).
48. Parks, D. H. et al. Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nat. Microbiol.* **2**, 1533–1542 (2017).
49. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).
50. Magoc, T. & Salzberg, S. L. FLASH: fast length adjustment of short reads to improve genome assemblies. *Bioinformatics* **27**, 2957–2963 (2011).
51. Rognes, T., Flouri, T., Nichols, B., Quince, C. & Mahé, F. VSEARCH: a versatile open source tool for metagenomics. *PeerJ* **4**, e2584 (2016).
52. Li, W., Fu, L., Niu, B., Wu, S. & Wooley, J. Ultrafast clustering algorithms for metagenomic sequence analysis. *Brief. Bioinform.* **13**, 656–668 (2012).
53. Caporaso, J. G. et al. QIIME allows analysis of high-throughput community sequencing data. *Nat. Methods* **7**, 335–336 (2010).
54. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree: computing large minimum evolution trees with profiles instead of a distance matrix. *Mol. Biol. Evol.* **26**, 1641–1650 (2009).
55. Rasmussen, B., Fletcher, I. R., Brocks, J. J. & Kilburn, M. R. Reassessing the first appearance of eukaryotes and cyanobacteria. *Nature* **455**, 1101–1104 (2008).
56. Shih, P. M., Hemp, J., Ward, L. M., Matzke, N. J. & Fischer, W. W. Crown group Oxyphotobacteria postdate the rise of oxygen. *Geobiology* **15**, 19–29 (2017).
57. Parfrey, L. W., Lahr, D. J. G., Knoll, A. H. & Katz, L. A. Estimating the timing of early eukaryotic diversification with multigene molecular clocks. *Proc. Natl Acad. Sci. USA* **108**, 13624–13629 (2011).
58. Brocks, J. J. et al. Biomarker evidence for green and purple sulphur bacteria in a stratified palaeoproterozoic sea. *Nature* **437**, 866–870 (2005).
59. Walter, M., Buick, R. & Dunlop, J. Stromatolites 3,400–3,500 Myr old from the North Pole area, Western Australia. *Nature* **284**, 443–445 (1980).
60. Ryder, G., Koeberl, C. & Mojzsis, S. J. in *Origin of the Earth and Moon* (eds Canup, R. & Kevin Righter, K.) 475–492 (Univ. Arizona Press, Tucson, AZ, 2000).
61. Dykhuizen, D. Species numbers in bacteria. *Proc. Calif. Acad. Sci.* **56**, 62–71 (2005).
62. Fraser, C., Alm, E. J., Polz, M. F., Spratt, B. G. & Hanage, W. P. The bacterial species challenge: making sense of genetic and ecological diversity. *Science* **323**, 741–746 (2009).
63. Drummond, A. J. & Rambaut, A. BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evol. Biol.* **7**, 214 (2007).
64. Stamatakis, A. RAXML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
65. Britton, T., Anderson, C. L., Jacquet, D., Lundqvist, S. & Bremer, K. Estimating divergence times in large phylogenetic trees. *Syst. Biol.* **56**, 741–752 (2007).
66. May, R. M. How many species inhabit the earth? *Sci. Am.* **267**, 42–49 (1992).
67. Mora, C., Tittensor, D. P., Adl, S., Simpson, A. G. & Worm, B. How many species are there on Earth and in the ocean? *PLoS Biol.* **9**, e1001127 (2011).
68. Eastman, J. M., Harmon, L. J., Tank, D. C. & Paradis, E. Congruification: support for time scaling large phylogenetic trees. *Methods Ecol. Evol.* **4**, 688–691 (2013).
69. Louca, S. & Doebeli, M. Efficient comparative phylogenetics on large trees. *Bioinformatics* **34**, 1053–1055 (2017).
70. Smith, S. A., Beaulieu, J. M. & Donoghue, M. J. An uncorrelated relaxed-clock analysis suggests an earlier origin for flowering plants. *Proc. Natl Acad. Sci. USA* **107**, 5897–5902 (2010).
71. Kuo, C.-H. & Ochman, H. Inferring clocks when lacking rocks: the variable rates of molecular evolution in bacteria. *Biol. Direct* **4**, 35 (2009).
72. Day, W. H. E. Optimal algorithms for comparing trees with labeled leaves. *J. Classif.* **2**, 7–28 (1985).
73. Borg, I., Groenen, P. J. F. & Mair, P. *Applied Multidimensional Scaling* (Springer, Berlin, 2013).
74. Ricklefs, R. E. Estimating diversification rates from phylogenetic information. *Trends Ecol. Evol.* **22**, 601–610 (2007).
75. Ochman, H. & Wilson, A. Evolution in bacteria: evidence for a universal substitution rate in cellular genomes. *J. Mol. Evol.* **26**, 74–86 (1987).
76. Moran, N. A., Munson, M. A., Baumann, P. & Ishikawa, H. A molecular clock in endosymbiotic bacteria is calibrated using the insect hosts. *Proc. R. Soc. Lond. B* **253**, 167–171 (1993).
77. Scrucca, L. Model-based SIR for dimension reduction. *Comput. Stat. Data Anal.* **55**, 3010–3026 (2011).

### Acknowledgements

We thank D. H. Parks for providing the 16S rRNA sequences from MAGs<sup>48</sup>. S.L. was supported by an NSERC grant and a postdoctoral fellowship from the Biodiversity Research Centre, University of British Columbia. M.W.P., M.D. and L.W.P. were supported by NSERC Discovery Grants. P.M.S. was supported by The Branco Weiss Fellowship – Society in Science. W.W.F. acknowledges support from the Simons Collaboration on the Origins of Life and NASA Exobiology award number NNX16AJ57G.

### Author contributions

S.L., L.W.P. and M.D. conceived the project. S.L. developed the mathematical methods, performed the diversification analyses and wrote the first draft of the manuscript. P.M.S. performed the molecular clock analyses of the BEAST trees, provided the cyanobacterial multigene tree and contributed to the development of the project ideas. All authors helped to interpret the results, advised on methodological improvements and contributed to writing the manuscript.

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41559-018-0625-0>.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Correspondence and requests for materials** should be addressed to S.L.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see [Authors & Referees](#) and the [Editorial Policy Checklist](#).

### Statistical parameters

When statistical analyses are reported, confirm that the following items are present in the relevant location (e.g. figure legend, table legend, main text, or Methods section).

n/a Confirmed

- The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement
- An indication of whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided  
*Only common tests should be described solely by name; describe more complex techniques in the Methods section.*
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistics including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g.  $F$ ,  $t$ ,  $r$ ) with confidence intervals, effect sizes, degrees of freedom and  $P$  value noted  
*Give  $P$  values as exact values whenever suitable.*
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's  $d$ , Pearson's  $r$ ), indicating how they were calculated
- Clearly defined error bars  
*State explicitly what error bars represent (e.g. SD, SE, CI)*

*Our web collection on [statistics for biologists](#) may be useful.*

### Software and code

Policy information about [availability of computer code](#)

Data collection

No software was used for data collection. All data analyzed were already available on public repositories.

Data analysis

Any custom software used is either included as Supplemental material, or freely available at the CRAN R package repository (<https://cran.r-project.org>). Any 3rd party software used is publicly available and cited in the Methods.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors/reviewers upon request. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

### Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

Amplicon sequencing data used to recover de novo OTUs are publicly available under the accession numbers listed in Supplementary File 1. Accession numbers for sequencing data used from the Earth Microbiome Project are listed in Supplementary File 2. R code used in this study is provided as Supplementary File 3. Timetrees

as well as undated phylogenetic trees constructed in this study are provided as Supplementary File 4. Taxonomic classifications of de novo OTUs are provided as Supplementary File 5.

## Field-specific reporting

Please select the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences  Behavioural & social sciences  Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/authors/policies/ReportingSummary-flat.pdf](https://www.nature.com/authors/policies/ReportingSummary-flat.pdf)

## Ecological, evolutionary & environmental sciences study design

All studies must disclose on these points even when the disclosure is negative.

Study description	We reconstruct diversification dynamics of bacteria using publicly available 16S rRNA gene sequence data and novel phylogenetic methods.
Research sample	All analyses were based on existing, publicly available DNA sequence data. Accession numbers are provided as Supplementary Material.
Sampling strategy	Amplicon sequencing data, used to generate our de-novo trees, were chosen so as to represent as wide of an environmental range as possible, under the constraint that they must cover at least 200 bp of the 16S V4 region. For the SILVA trees, all OTUs available within the appropriate taxon were used. For the BEAST-computed trees, trees were sub-sampled randomly to enable computational feasibility. Any sub-sampling has been accounted for in our calculations.
Data collection	No novel data was collected.
Timing and spatial scale	No novel data was collected.
Data exclusions	No relevant data was explicitly excluded.
Reproducibility	All raw sequencing data used are available at public repositories under accession numbers provided as Supplemental material. SILVA 16S sequences are publicly available at the SILVA project website ( <a href="https://www.arb-silva.de">https://www.arb-silva.de</a> ). All trees created in this study, and computer code needed to analyze them, are provided as supplemental material. All other software used in our analyses are publicly available and cited in our Methods section.
Randomization	This study does not include experimental groups.
Blinding	This study does not include experimental groups.
Did the study involve field work?	<input type="checkbox"/> Yes <input checked="" type="checkbox"/> No

## Reporting for specific materials, systems and methods

### Materials & experimental systems

- n/a Involved in the study
- Unique biological materials
  - Antibodies
  - Eukaryotic cell lines
  - Palaeontology
  - Animals and other organisms
  - Human research participants

### Methods

- n/a Involved in the study
- ChIP-seq
  - Flow cytometry
  - MRI-based neuroimaging

In the format provided by the authors and unedited.

# Bacterial diversification through geological time

Stilianos Louca <sup>1,2\*</sup>, Patrick M. Shih<sup>3,4,5</sup>, Matthew W. Pennell<sup>1,2</sup>, Woodward W. Fischer<sup>6</sup>,  
Laura Wegener Parfrey <sup>1,2,7</sup> and Michael Doebeli<sup>1,2,8</sup>

---

<sup>1</sup>Biodiversity Research Centre, University of British Columbia, Vancouver, British Columbia, Canada. <sup>2</sup>Department of Zoology, University of British Columbia, Vancouver, British Columbia, Canada. <sup>3</sup>Joint BioEnergy Institute, Emeryville, CA, USA. <sup>4</sup>Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. <sup>5</sup>Department of Plant Biology, University of California, Davis, Davis, CA, USA. <sup>6</sup>Division of Geological and Planetary Sciences, California Institute of Technology, Pasadena, CA, USA. <sup>7</sup>Department of Botany, University of British Columbia, Vancouver, British Columbia, Canada. <sup>8</sup>Department of Mathematics, University of British Columbia, Vancouver, British Columbia, Canada.

\*e-mail: [louca.research@gmail.com](mailto:louca.research@gmail.com)

## S.1 Mathematical derivations

This section presents mathematical formulas for the behavior of speciation-extinction cladogenic models in the limit of infinitely large bifurcating timetrees. Using simulations, we found that our mathematical formulas are accurate for trees with more than 500 tips (see Methods and Supplement S.2), as is the case for all trees examined in this study. We note that for special cases, some of the presented formulas have been discussed previously by other authors<sup>1-3</sup>, however below we provide the full derivations for completeness. We then explain how these formulas can be used to reconstruct past diversification dynamics using trees comprising only extant OTUs. We will use the term “OTU” to refer to any group of closely related extant organisms represented as an individual tip in the tree.

### S.1.1 Cladogenesis in infinitely large trees

Let  $T$  be the time point at which taxa are sampled to construct the considered tree, that is, today. We will focus on the total number of extant lineages (“total diversity”, denoted  $N(t)$ ) at each time point  $t$  as well as the number of lineages represented in the considered tree, that is, with at least one extant discovered representative at time  $T$  (“lineages through time” or LTT, denoted  $\tilde{N}(t, T)$ ). We denote by  $\lambda(t)$  the (per-lineage) speciation rate and by  $\mu(t)$  the (per-lineage) extinction rate at any time  $t$ . We assume that speciations and extinctions occur randomly across the tree at Poissonian rates (i.e., with exponential uncorrelated waiting times), with all extant lineages being equally likely to speciate and equally likely to go extinct (but see our discussion on heterogenous rates in Supplement S.5). For any time  $t$  in the past ( $t \leq T$ ), denote by  $E(t, T)$  the probability that an extant lineage at time  $t$  is absent from the tree at time  $T$ , that is, has no extant discovered descendant at time  $T$ . In other words,  $E(t, T)$  is the probability that the entire clade, descending from a single lineage at time  $t$ , has gone extinct at time  $T$ . Thus,  $E(t, T)$  takes into account all possible speciation/extinction scenarios within that descending clade. For large trees the represented number of lineages,  $\tilde{N}(t, T)$  is related to the total number of lineages,  $N(t)$ , as follows:

$$\tilde{N}(t, T) = N(t) \cdot [1 - E(t, T)]. \quad (1)$$

Reciprocally, if one knows  $\tilde{N}(t, T)$  and  $E(t, T)$  at some time  $t$ , one can estimate the total number of lineages  $N(t)$ . In the next steps, we will derive differential equations that can be used to calculate  $E(t, T)$  and  $N(t)$  over time, provided certain information is available about speciation and extinction rates. For trees covering only a subset of total extant OTUs (incomplete OTU sampling) we assume random and phylogenetically uncorrelated sampling, that is, we assume that extant OTUs are included or excluded from the tree independently of one another. In this case, incomplete OTU sampling at time  $T$  can be mathematically treated as an instantaneous extinction event just prior to time  $T$  that affects all lineages independently and with equal probabilities<sup>4-6</sup>. We later assess the accuracy of the assumption of random OTU sampling, and the effects that non-random sampling may have on the representation of past lineages in the tree (see Methods).

Consider some focal clade of size  $n$  (i.e., comprising  $n$  extant lineages) at some time point  $t$ . The transition probability rate to any other size  $k$ , denoted  $\mathbb{Q}_{kn}$ , depends on the speciation and extinction rates at the time, as follows:

$$\begin{aligned} \mathbb{Q}_{kn}(t) &= n\lambda(t), & k = n + 1, \\ \mathbb{Q}_{kn}(t) &= n\mu(t), & k = n - 1, \\ \mathbb{Q}_{nn}(t) &= -\mathbb{Q}_{n+1,n}(t) - \mathbb{Q}_{n-1,n}(t), \\ \mathbb{Q}_{kn}(t) &= 0 & \text{for all other } k. \end{aligned} \quad (2)$$

Note that  $\mathbb{Q}_{kn}$  defines the transition rate matrix of a continuous-time Markov chain, whose value corresponds to the size of the focal clade over time. Hence, for a clade of size  $n$  at some time point  $t$ , the probability that the clade has size  $k$  at time  $T \geq t$  is given by the  $(k, n)$ -th entry of the transition matrix,  $\mathbb{M}(t, T)$ , defined as the solution to the ordinary differential equation (ODE):

$$\frac{d\mathbb{M}(t, T)}{dT} = \mathbb{Q}(T) \cdot \mathbb{M}(t, T), \quad \mathbb{M}(t, t) = \text{Id}, \quad (3)$$

where Id is the identity matrix. In particular, for a single extant lineage at time  $t$  (clade of size  $n = 1$ ), the probability of extinction  $E(t, T)$  is given by the  $(0, 1)$ -th entry (1st row, 2nd column) of the matrix  $\mathbb{M}(t, T)$ . The above ODE describes how  $\mathbb{M}(t, T)$  evolves forward in time (i.e. for increasing  $T$ ), however an analogous ODE also exists in the backward direction (i.e. for decreasing  $t$ ). Indeed, according to the Chapman-Kolmogorov equation for Markov chains<sup>7</sup> one has

$$\begin{aligned} \mathbb{M}(t, T) &= \mathbb{M}(t + \varepsilon, T) \cdot \mathbb{M}(t, t + \varepsilon) \\ &\stackrel{\text{Eq. (3)}}{=} \left[ \mathbb{M}(t, T) + \varepsilon \frac{d\mathbb{M}(t, T)}{dt} + \mathcal{O}(\varepsilon^2) \right] \cdot \left[ \text{Id} + \varepsilon \mathbb{Q}(t) + \mathcal{O}(\varepsilon^2) \right] \\ &= \mathbb{M}(t, T) + \varepsilon \mathbb{M}(t, T) \mathbb{Q}(t) + \varepsilon \frac{d\mathbb{M}(t, T)}{dt} + \mathcal{O}(\varepsilon^2), \end{aligned} \quad (4)$$

for any  $\varepsilon \rightarrow 0$ . Thus:

$$\frac{d\mathbb{M}(t, T)}{dt} = -\mathbb{M}(t, T) \cdot \mathbb{Q}(t). \quad (5)$$

In particular, since  $E(t, T) = \mathbb{M}_{01}(t, T)$ , one has:

$$\frac{dE(t, T)}{dt} = -[\mathbb{M}(t, T) \mathbb{Q}(t)]_{01} = -\underbrace{\mathbb{M}_{00}(t, T)}_1 \underbrace{\mathbb{Q}_{01}(t)}_\mu - \underbrace{\mathbb{M}_{01}(t, T)}_{E(t, T)} \underbrace{\mathbb{Q}_{11}(t)}_{-\lambda - \mu} - \underbrace{\mathbb{M}_{02}(t, T)}_{E(t, T)^2} \underbrace{\mathbb{Q}_{21}(t)}_\lambda, \quad (6)$$

that is,

$$\frac{dE(t, T)}{dt} = -\mu(t) + E(t, T)[\lambda(t) + \mu(t)] - E(t, T)^2 \lambda(t). \quad (7)$$

Note that we used the relationship  $\mathbb{M}_{02}(t, T) = E(t, T)^2$ , which holds under the assumption of this model that individual lineages survive or go extinct independently of one another.

In many cases (such as this study) it is useful to estimate total diversities in reverse time, that is, starting with the known current total number of lineages  $N(T)$  and moving backward in time to estimate  $N(t)$  for  $t < T$ . For that purpose, we re-parameterize all time functions in terms of the ‘‘age’’  $\tau = T - t$ , for example writing  $E(\tau)$  instead of  $E(T - \tau, T)$ . Written in terms of age, ODE (7) becomes

$$\frac{dE(\tau)}{d\tau} = \mu(\tau) - E(\tau) \cdot [\lambda(\tau) + \mu(\tau)] + E(\tau)^2 \lambda(\tau). \quad (8)$$

Hence,  $E(\tau)$  can be directly calculated backwards in time, assuming that  $E(0)$  is known and  $\lambda(\sigma)$  and  $\mu(\sigma)$  are known for all younger ages  $\sigma \leq \tau$ . The initial value  $E(0)$  corresponds to the probability of any given currently extant lineage being absent from the tree due to incomplete taxon sampling. For trees comprising all extant lineages,  $E(0)$  will be zero, but for trees only comprising a random fraction  $\rho$  (‘‘sampling fraction’’) of extant lineages,  $E(0)$  will be equal to  $1 - \rho$ . Solving Eq. (8) with initial condition  $E(0) = 0$  even if  $\rho \neq 1$ ,

will yield the probability of lineage extinction regardless of whether a lineage is sampled at time  $T$  or not.

For comparison with existing common models, we briefly consider the case where  $\lambda$  and  $\mu$  are constant over time. In that case, ODE (8) can be solved explicitly for any initial condition  $E(0) = 1 - \rho$ , yielding:

$$\begin{aligned} E(\tau) &= 1 - \frac{(\lambda - \mu)\rho}{\lambda\rho + [(1 - \rho)\lambda - \mu]e^{-\tau(\lambda - \mu)}}, \quad \text{if } \lambda \neq \mu, \\ E(\tau) &= 1 - \frac{\rho}{1 + \rho\tau\lambda}, \quad \text{if } \lambda = \mu. \end{aligned} \quad (9)$$

The solution in Eq. (9) has been previously derived by Nee *et al.*<sup>2</sup>. It is easy to see that  $E(\tau)$  converges to the value  $\mu/\lambda$  if  $\mu \leq \lambda$  and to the value 1 if  $\mu > \lambda$ , as  $\tau \rightarrow \infty$ . In particular, if  $\mu < \lambda$  (exponential diversification with constant speciation/extinction rates), one has  $\tilde{N}(\tau) \approx (1 - \mu/\lambda)N(\tau)$  for sufficiently large  $\tau$ . Hence, the ultimate probability of a lineage becoming extinct in the long run is given by  $\mu/\lambda$ , a well known result in paleobiology<sup>8</sup>. Further, in the distant past  $\tilde{N}$  will appear to increase at the same exponential rate as  $N$ , a result also known from previous studies<sup>6</sup>.

The ODE in Eq. (8) further allows us to derive an ODE for the represented number of lineages  $\tilde{N}(\tau)$  in the limit of infinitely large trees, as follows:

$$\begin{aligned} \frac{d\tilde{N}}{d\tau} &= \frac{d}{d\tau} [(1 - E)N] \\ &= \frac{dN}{d\tau}(1 - E) - N \frac{dE}{d\tau} \\ &\stackrel{\text{Eq. (8)}}{=} N \cdot (\mu - \lambda)(1 - E) - N \cdot [\mu - E(\lambda + \mu) + E^2\lambda] \\ &= \lambda N \cdot (2E - 1 - E^2), \end{aligned} \quad (10)$$

where we used the fact that  $dN/d\tau = (\mu - \lambda)N$ . Evaluating Eq. (10) at age  $\tau = 0$  yields:

$$\left. \frac{1}{\tilde{N}} \frac{d\tilde{N}}{d\tau} \right|_{\tau=0} = \lambda(0) \frac{N(0)}{\tilde{N}(0)} [2(1 - \rho) - 1 - (1 - \rho)^2] = -\rho\lambda(0). \quad (11)$$

In particular, in the case of complete sampling ( $\rho = 1$ ) at  $\tau = 0$  one has  $\tilde{N}^{-1}d\tilde{N}/d\tau = -\lambda$ , that is, at the tips of the tree the exponential growth rate of  $\tilde{N}$  is equal to the speciation rate. This is a well-known result in phylogenetics<sup>6</sup>, which is commonly explained by the fact that recently emerged lineages have not had the time yet to go extinct, making the tree's recent growth appear as if extinction rates were zero.

### S.1.2 Fitting speciation/extinction models to trees

If the absolute speciation and extinction rates  $B = \lambda N$  and  $D = \mu N$  are known, then integration of the ODE

$$\frac{dN}{d\tau} = D(\tau) - B(\tau) \quad (12)$$

allows calculation of the total diversity  $N$  over time. Further, knowledge of the sampling fraction  $\rho$  allows the calculation of  $E(\tau)$  by integrating ODE (8) backwards in time, and therefore of  $\tilde{N}(\tau)$  via Eq. (1). Today's total number of extant lineages,  $N(0) = \tilde{N}(0)/\rho$ , can be used as an initial condition for ODE (12). Thus, for a given tree, a given sampling fraction  $\rho$  and a model that specifies  $B$  and  $D$  as functions of  $\tau$  and/or

$N$ , the LTT predicted by the model,  $\tilde{N}_m(\tau)$ , can be compared to the LTT observed from the tree,  $\tilde{N}(\tau)$ . In principle, then, it is possible to fit model parameters by minimizing the deviation of the model's predicted  $\tilde{N}_m$  from the observed  $\tilde{N}$ .

In the simplest case the (per-lineage) speciation and extinction rates may be assumed to be constant over time, i.e.  $B = \lambda N$  and  $D = \mu N$  for constant parameters  $\lambda, \mu$ . In general, however, speciation and/or extinction rates could themselves depend on current diversities or depend explicitly on time, and thus additional model parameters may be needed to describe these dependencies<sup>5</sup>. As we explain in the main text, in this study we found that constant speciation/extinction rate models provide adequate descriptions of global bacterial and cyanobacterial diversification dynamics at the temporal resolutions and time intervals considered (mean relative deviation below 5% in all cases).

### S.1.3 Non-parametric estimations

Certain useful quantities may also be estimated non-parametrically from the LTT, that is, using the value and derivatives of the LTT at each time point, instead of fitting a parametric model to the entire LTT curve. The advantage of non-parametric estimations over parametric models is that they make fewer or no assumptions on how speciation and extinction rates vary over time, and they do not include a tradeoff between temporal resolution and model simplicity (e.g., in terms of the number of allowed rate shifts). Below, we describe novel non-parametric methods for extracting information from LTTs. In Supplement S.2 we demonstrate the power of these methods, using trees simulated under various scenarios.

We denote by  $\tilde{\nu} = (1/\tilde{N})d\tilde{N}/d\tau$  the relative slope of the LTT. Note that  $\tilde{\nu}$  is the apparent diversification rate that would be measured solely based on the LTT without consideration of extinctions or incomplete taxon sampling. As we show below,  $\tilde{\nu}$  and its slope at any time in the past can yield important information on speciation and extinction rates at that time. We note that examining  $\tilde{\nu}$  is analogous to examining branching frequencies in the tree over time, since  $\tilde{\nu}$  is proportional to the branching frequency per lineage. Dividing Eq. (10) by Eq. (1), yields the relationship

$$\tilde{\nu} = \frac{1}{\tilde{N}} \frac{d\tilde{N}}{d\tau} = \lambda \cdot (E - 1). \quad (13)$$

Solving Eq. (13) for  $E$  yields

$$E(\tau) = 1 + \frac{\tilde{\nu}(\tau)}{\lambda(\tau)}. \quad (14)$$

Inserting Eq. (14) into Eq. (1) yields an expression for the total diversity:

$$N(\tau) = \frac{\tilde{N}(\tau)}{1 - E(\tau)} = -\frac{\lambda(\tau)\tilde{N}(\tau)}{\tilde{\nu}(\tau)}. \quad (15)$$

Hence, if the speciation rate  $\lambda(\tau)$  is known at some age  $\tau$ , one can estimate  $E(\tau)$  and  $N(\tau)$  directly from the LTT using Eq. (14) and (15), respectively. Subsequently, one may estimate  $\mu(\tau)$  based on the relationship:

$$\frac{dN}{d\tau} = -\frac{dN}{dt} = -r \cdot N = -(\lambda - \mu) \cdot N, \quad (16)$$

where  $r = \lambda - \mu$  is the diversification rate. From Eq. (16) one obtains:

$$\mu = \lambda + \frac{1}{N} \frac{dN}{d\tau}. \quad (17)$$

We emphasize that the speciation rate  $\lambda(\tau)$  is usually unknown except for  $\tau = 0$ , where  $\lambda(\tau = 0) = -\tilde{\nu}(\tau = 0)/\rho$  according to Eq. (11). In the special case where  $\lambda$  is approximately constant over time,  $\lambda$  can be estimated as  $-\tilde{\nu}(0)/\rho$  or estimated by fitting a speciation-extinction model (as done in this study), and subsequently  $E$ ,  $N$  and  $\mu$  can be estimated from the LTT using Eqs. (14), (15) and (17). If  $\lambda$  varies over time, the erroneous assumption of a constant  $\lambda$ , equal to  $\lambda(0)$ , would introduce an error into the estimated  $E$ ,  $N$  and  $\mu$ . In that case, it may be useful to consider modified quantities, described below, that can be estimated from the LTT regardless of whether  $\lambda$  is constant or not.

### **Pulled total diversity**

The first modified quantity, which we refer to as “*pulled total diversity*”, is defined as:

$$N_p(\tau) = N(\tau) \cdot \frac{\lambda_o}{\lambda(\tau)}, \quad (18)$$

where  $\lambda_o = \lambda(0)$ . Observe that  $N_p$  is similar to the total diversity  $N$ , but differs from the latter by the factor  $\lambda(0)/\lambda(\tau)$ . From Eq. (15) it becomes clear that

$$N_p(\tau) = -\frac{\lambda_o \tilde{N}(\tau)}{\tilde{\nu}(\tau)}. \quad (19)$$

Since  $\lambda_o$  can be estimated from the LTT near the tips of the tree, the right hand side of Eq. (19) can be estimated directly from the LTT without any *a priori* assumptions on how  $\lambda$  may have varied over time. If  $\lambda$  is approximately constant over time, then  $N_p(\tau)$  will be approximately equal to  $N(\tau)$ . In particular, unless  $\lambda$  changed drastically (i.e., by orders of magnitude) over time,  $N_p$  provides a quick way to estimate past total diversities to order of magnitude accuracy.

### **Pulled extinction rate**

The second modified quantity, which we refer to as “*pulled extinction rate*”, is defined as:

$$\mu_p(\tau) = \mu(\tau) + [\lambda_o - \lambda(\tau)] - \frac{1}{\lambda} \frac{d\lambda}{d\tau}. \quad (20)$$

Observe that  $\mu_p$  is similar to the extinction rate  $\mu$ , but differs from the latter by the terms  $[\lambda_o - \lambda(\tau)]$  and  $\lambda^{-1}d\lambda/d\tau$ , both of which disappear when  $\lambda$  is constant over time. Eq. (20) can be rewritten as

$$\mu_p(\tau) = \lambda_o - r(t) - \frac{1}{\lambda} \frac{d\lambda}{d\tau}. \quad (21)$$

Hence, if diversity is close to speciation/extinction equilibrium ( $|r| \ll \lambda_o$ ), and  $\lambda$  changes only slowly, the pulled extinction rate is approximately equal to the recent extinction rate  $\lambda_o$ . Using Eqs. (15) and (17) in Eq.

(20) yields the formula:

$$\mu_p(\tau) = \lambda_o + \frac{d}{d\tau} \ln \left[ -\frac{\tilde{N}(\tau)}{\tilde{\nu}(\tau)} \right]. \quad (22)$$

Similarly to the pulled total diversity, the pulled extinction rate can be estimated from the LTT without any *a priori* assumptions on how  $\lambda$  and  $\mu$  may have varied over time.

### Pulled diversification rate

To derive the third modified quantity, we insert Eq. (15) into Eq. (17), thus obtaining the relationship:

$$\lambda - \mu + \frac{1}{\lambda} \frac{d\lambda}{d\tau} = -\frac{d}{d\tau} \ln \left[ -\frac{\tilde{N}(\tau)}{\tilde{\nu}(\tau)} \right] = -\tilde{\nu} + \frac{1}{\tilde{\nu}} \frac{d\tilde{\nu}}{d\tau}. \quad (23)$$

Observe that the right hand side of Eq. (23) can be readily calculated from the LTT without any assumption on  $\lambda$ ,  $\mu$  or  $\rho$ . The left hand side resembles the diversification rate ( $r = \lambda - \mu$ ), modified (“pulled”) by an additive correction term. We define the left hand side as the “*pulled diversification rate*” (PDR):

$$r_p = r + \frac{1}{\lambda} \frac{d\lambda}{d\tau}, \quad (24)$$

so that

$$r_p(\tau) = -\tilde{\nu} + \frac{1}{\tilde{\nu}} \frac{d\tilde{\nu}}{d\tau}. \quad (25)$$

We point out that the “pulled” variables ( $\mu_p$ ,  $r_p$ ,  $N_p$ ) and the recent speciation rate  $\lambda_o$  satisfy similar algebraic relationships as their “non-pulled” counterparts:

$$r_p = \lambda_o - \mu_p, \quad (26)$$

and

$$\frac{1}{N_p} \frac{dN_p}{d\tau} = -r_p. \quad (27)$$

Equation (25) can be used to estimate the pulled diversification rate from the LTT curve, without knowing  $\lambda$ ,  $\mu$  or  $\rho$ , without any assumptions on how  $\lambda$  or  $\mu$  may have varied over time, and without fitting any model parameters (see simulations in Supplementary Figs. 1–6 for examples). The downside is that, similarly to the pulled total diversity and the pulled extinction rate,  $r_p$  is a composite quantity that only resembles the diversification rate  $r$  when  $\lambda$  is roughly constant. If  $\lambda$  changes rapidly over time,  $r_p$  will differ substantially from  $r$ . Hence, solely knowing  $r_p$  does not *a priori* determine its constituents  $\mu$  and  $\lambda$ , nor  $r$ , unless either  $\lambda$  or  $\mu$  is estimated separately (or assumed to be constant).

## S.2 Evaluating non-parametric methods using simulations

To demonstrate the use of non-parametric estimation methods introduced in Supplement S.1.3, and to assess how the pulled extinction rates (PERs), pulled diversification rates (PDRs) and pulled total diversities (PTDs)

differ from their “non-pulled” analogs, we tested these methods on trees simulated under various scenarios (Figs. 1 and 1–6). Specifically, we simulated various trees with constant or non-constant speciation rates, and then examined the behavior of the PERs, PDRs and PTDs. In all scenarios, a random tree was generated according to the general model discussed in Supplement S.1.1, whereby lineages speciate or go extinct randomly at exponentially distributed time intervals and independently of one another. At any moment in time, speciation and extinction rates are assumed to be the same (homogenous) across all clades, but can vary over time. For an evaluation of how rate heterogeneity across clades could affect our results, see Supplement S.5. Following a simulation, the PER ( $\mu_p$ ) was estimated based on the generated tree’s LTT using formula (22), the PDR ( $r_p$ ) was estimated using formula (25) and the PTD ( $N_p$ ) was estimated using formula (19). Details of each simulation scenario are described below.

In the first examined scenario (Fig. 1), a tree was generated with constant speciation rate  $\lambda$ , and a density-dependent extinction rate ( $\mu \propto N^{0.2}$ ) with an additional temporary increase about 80 Ma ago, simulating a mass extinction event. Note that, since  $\lambda$  was constant, in this scenario the PER was equal to the extinction rate, the PDR was equal to the diversification rate and the PTD was equal to the total diversity. As seen in Fig. 1a, the total diversity (or PTD) is estimated very accurately except for the very early times near the root, and clearly reveals the effects of the mass extinction event. Similarly, the PDR, as estimated using Eq. (25), closely resembles the true PDR and clearly reflects the temporary drop in the diversification rate (Fig. 1b). The PER contains substantial noise and wide confidence intervals, although the mass extinction event is again clearly reproduced.

In the second scenario (Fig. 2), a tree was generated using a constant speciation rate and an extinction rate exhibiting a sharp temporary increase (mass extinction event, lasting only  $\sim 2$  Ma) about 90 Ma ago. Similarly to the first scenario, the spike in the extinction rate induces a spike in the PER (Fig. 2C) and in the PDR (Fig. 2b). Similarly, the tree’s LTT, the PTD and the estimated PTD all become distorted due to the mass extinction event. Due to the brevity of the event, the spikes in the estimated PER, PDR and PTD are severely damped by the noise-filter (smoothing) applied during estimation. That said, the tree’s LTT still reflects the mass extinction event (Fig. 2a). In particular, fitting a constant-rates cladogenic model, as done for the bacterial trees, reveals a clear deviation from the LTT near the mass extinction event (Fig. 2d).

In the third scenario (Fig. 3), a tree was generated using a strongly oscillating speciation rate ( $\pm 50\%$  around the mean, with period 200 million years) and a constant extinction rate. The estimated PDR clearly reflects the oscillatory nature of the diversification rate (Fig. 3b), although noticeable differences exist compared to the diversification rate due to the “pulling” term  $\lambda^{-1}d\lambda/d\tau$  discussed in Eq. (25). The estimated PER is strongly distorted by the oscillations in the true speciation rate (Fig. 3c). Interestingly, the estimated PTD only deviates moderately from the true total diversity, when compared on a logarithmic axis. This is because  $\lambda$  is always within the same order of magnitude as  $\lambda_o$ , and hence  $N_p$  also is always within the same order of magnitude as  $N$  (Eq. (18)).

In the fourth scenario (Fig. 4), a tree was generated using a constant extinction rate, and a diversity-dependent speciation rate that increased as the tree grew over time ( $\lambda \propto N^{0.1}$ ). As a result, near the tips  $\lambda$  was about twice as large as near the root. The diversity-dependence-driven increase in the speciation rate leads to a positive trend in the PDR (Fig. 4b), a negative trend in the PER (Fig. 4c), and a systematic difference between  $N_p$  and  $N$  that increases towards older ages (Fig. 4a).

In the fifth scenario (Fig. 5), a tree was generated using a constant extinction rate, and a speciation rate exhibiting a strong temporary increase (spike) about 90 Ma ago. The variation in the speciation rate causes a strong distortion in the PER, which may be confused with a temporary drop in extinction rate (Fig. 5c). The PDR also clearly reflects the temporary increase in diversification rate, and the difference between the two is only moderate (Fig. 5b). Similarly, the PTD approximately resembles the total diversity, and clearly

reproduces the sudden increase of diversity during the speciation spike (Fig. 5a).

In the sixth scenario (Fig. 6), a tree was generated using a constant extinction rate, and a speciation rate exhibiting a strong, sharp, temporary increase about 140 Ma ago (“speciation spike”, lasting only  $\sim 1$  Ma). Similarly to the previous scenario, the spike in the speciation rate induces a strong spike in the PER (Fig. 6c) and in the PDR (Fig. 6b). Similarly, the tree’s LTT, the PTD and the estimated PTD all become distorted due to the speciation spike. Due to the brevity of the speciation spike, the spikes in the estimated PER, PDR and PTD are severely damped by the noise-filter (smoothing) applied during estimation. That said, the tree’s LTT clearly reflects the past anomaly in the speciation rate (Fig. 6a). In particular, fitting a constant-rates cladogenic model, as done for the bacterial trees, does not accurately reproduce the tree’s LTT (Fig. 6d).

The above examples demonstrate three important points: First, for large trees generated under a constant speciation rate, non-parametric estimation methods can accurately reconstruct extinction rates, diversification rates and total diversities over time. Second, if the speciation rate varies over time (for example due to diversity-dependence), then this variation generally (i.e. if  $\mu$  is chosen arbitrarily) leads to similarly strong variations in the PER, even if the extinction rate is constant over time. It is generally hard (using solely the LTT) to determine whether variations in the PER and/or PDR are due to a varying speciation rate or due to a varying extinction rate. If, however, the estimated PER turns out to be constant over time, this is a strong indication that both  $\lambda$  and  $\mu$  were both constant or varied only slowly over time (see detailed explanation in Supplement S.4). Third, the magnitude of the PDR is generally comparable to the magnitude of the diversification rate, except perhaps during short isolated time intervals. In fact in all of the above simulations the PDR approximately resembled the diversification rate.

### S.3 Robustness of global extant bacterial diversity estimates

Here we discuss the robustness of our estimates of global extant bacterial OTU richness. Our estimate of 1.4 million bacterial OTUs (discussed in the main article) remains approximately the same ( $\sim 1.1$  million OTUs) when we omit *de novo* OTUs found in fewer than 10 samples (68,927 OTUs kept, of which  $\sim 42\%$  were covered by SILVA), suggesting that a potential bias towards more ubiquitous OTUs only weakly affects global bacterial diversity estimates. Our estimate also remains approximately the same when we strongly subsample reads (keeping  $\sim 13,000$  reads per sample on average, instead of  $\sim 40,000$  in the full dataset), yielding an estimated  $\sim 1.2$  million OTUs. The weak dependence on sequencing depth shows that our approach is not substantially affected by a detection bias towards more abundant organisms in each sample. The estimate is also similar ( $\sim 1.6$  million OTUs) when we consider 16S rRNA sequences from metagenome-assembled genomes<sup>9</sup>, instead of OTUs recovered from amplicon sequences, indicating that primer bias also only has a moderate influence. Finally, our estimate is similar ( $\sim 1.9$  million OTUs) when we consider the overlap between our *de novo* OTUs and OTUs that we recovered from another massive dataset generated by the Earth Microbiome Project (EMP)<sup>10</sup>, instead of an overlap between our *de novo* OTUs and SILVA (Supplementary Table 3). Our estimate is comparable to another recent estimate ( $\sim 1.5$  million unique 16S rRNA sequences)<sup>11</sup>, but is 6 orders of magnitude lower than the estimate by Locey *et al.*<sup>12</sup>, which was based on extrapolation of empirical scaling laws of local richness to global scales. The strong disagreement with Locey *et al.*’s estimates supports arguments that empirical scaling laws of local richness cannot be extrapolated to estimate global microbial richness<sup>13</sup>. Further, many of the extremely rare but diverse OTUs found in previous studies (termed the “rare biosphere”), such as in early releases of the EMP<sup>14</sup>, which reportedly recovered millions of rare V4-OTUs, may have been spurious due to methodology. Spurious OTUs are known to frequently occur in amplicon sequencing studies and can substantially inflate microbial richness estimates, notably due to sequencing errors and chimera formation during PCR<sup>15,16</sup>. Indeed, a more recent

release of the EMP that employed improved quality filters found only  $\sim 307,000$  unique V4 sequences across 96 studies<sup>10</sup>.

## S.4 Interpreting pulled variables of bacterial and cyanobacterial trees

### Pulled extinction rates

As discussed in the main text, here we found that estimated bacterial and cyanobacterial pulled extinction rates were almost constant over time. As we explain below, this strongly indicates that both  $\lambda$  and  $\mu$  were constant or varied only slowly over time. To see why this is the case, observe that for any time point the right hand side in Eq. (20) depends solely on the instantaneous speciation and extinction rates at that time point. Hence, if  $\mu$  and/or varied over time,  $\mu_p$  could only be constant if  $\mu$  instantaneously adjusted to (was determined by) the current value of  $\lambda$  and its relative rate of change, or vice versa. Specifically, for any given  $\mu$ , according to Eq. (21)  $\mu_p$  could only be (or estimated to be) approximately constant over time ( $|\mu_p(t) - \mu_p(T)| < \varepsilon$  for some small  $\varepsilon$ ) if  $\lambda(t)$  satisfied the differential equation

$$\frac{1}{\lambda} \frac{d\lambda}{dt} \approx r(t) - r_o, \quad (28)$$

where  $r_o$  is the recent diversification rate, and the error in Eq. (28) must be smaller than  $\varepsilon$ . For strongly varying  $r$ ,  $\lambda^{-1}d\lambda/dt$  must have closely and instantaneously followed  $r - r_o$ , however there is no known realistic mechanism that would cause  $\lambda^{-1}d\lambda/dt$  to behave in such way. In fact, for generic and independently chosen  $\mu(t)$  and  $\lambda(t)$ , one or both of which exhibits variations greatly exceeding  $\varepsilon$ , Eq. (28) would almost never be satisfied. The fact that for generic  $\mu$  and  $\lambda$  the pulled extinction rate is unlikely constant is also easily demonstrated using simulations (Figs. 1–6). It follows that  $r(t)$  must have been approximately constant ( $|r - r_o| \lesssim \varepsilon$ ) or small in magnitude ( $|r| \lesssim \varepsilon$ ), and thus  $\lambda^{-1}d\lambda/dt$  must also have been small in magnitude ( $\lesssim \varepsilon$ ). In the case of bacteria, this means that  $\lambda$  must have varied at billion-year time scales ( $|\lambda^{-1}d\lambda/dt| \sim 0.005 \text{ Ma}^{-1}$ ) or slower. We note that, since  $\lambda$  must have been constant or varied only slowly over time, according to Eq. (20) the recent pulled extinction rate is similar to the recent extinction rate ( $|\mu_p(T) - \mu(T)| \lesssim \varepsilon$ ). Towards older ages (e.g.,  $10^3 \text{ Ma}$  ago in the case of bacteria),  $\mu_p$  will remain close to  $\lambda_o$  even if  $\lambda$  and/or  $\mu$  changed substantially (Eq. 21), as observed in Figs. 2d,e.

### Pulled diversification rates

As discussed in the main text, all of our fitted models suggest low bacterial and cyanobacterial diversification rates ( $r \sim 0.002\text{--}0.005 \text{ Ma}^{-1}$ , Fig. 18c). Consistent with this observation, we also find low PDRs ( $|r_p| < 0.005 \text{ Ma}^{-1}$ , Figs. 2e,f and 21g–i). The fact that  $|r_p|$  is small over a prolonged time interval ( $|r_p| < \varepsilon$ , where  $\varepsilon \sim 0.005 \text{ Ma}^{-1}$ ) is indeed a strong indication that within that time interval the magnitude of the diversification rate ( $|r|$ ) was also not much larger than  $\varepsilon$ . To see why this is the case, note that at each time point  $t$  the right hand side of Eq. (24) only depends on the instantaneous diversification rate and the relative rate of change of  $\lambda$  at that time point. For generic  $\lambda$  and  $\mu$  with  $|r| \gg \varepsilon$  one would generally (and at most time points) also expect that  $|r_p| \gg \varepsilon$ , since otherwise  $\lambda^{-1}d\lambda/dt$  would have to instantaneously adjust to approximately cancel out  $r$ ; an unlikely coincidence. This logic is easily demonstrated using simulations, where  $r_p$  generally fluctuates at similar magnitudes as  $\mu$  and/or  $r$  (see simulation examples in Figs. 1–6).

## Pulled total diversities

In all examined bacterial and cyanobacterial trees, our fitted (constant-rate) models suggest a positive diversification rate and a roughly exponential increase of total diversity over the last 1 billion years, with bacterial and cyanobacterial diversity increasing roughly 10-fold (Figs. 2a,b and 21a–c). Consistent with this observation, bacterial and cyanobacterial pulled total diversities also exhibit an approximately exponential increase that closely resemble the predictions of the fitted models. An exponentially growing  $N_p$  suggests that  $N$  was itself growing approximately exponentially. To see why this is the case, recall that  $N_p = N\lambda/\lambda_o$  according to Eq. (18). Since  $N_p(\tau) \approx N_o e^{-r_o\tau}$ , where  $r_o$  is the diversification rate fitted by the model, one has

$$e^{r_o\tau} \approx \frac{\lambda(\tau)}{N(\tau)} \cdot \frac{N_o}{\lambda_o}. \quad (29)$$

The above equation can be satisfied in 3 alternative ways. First,  $\lambda$  could be approximately constant ( $\lambda \approx \lambda_o$ ), in which case  $N$  is itself growing exponentially ( $N \approx N_p$ ) subject to approximately constant speciation and extinction rates. Second,  $N(\tau)$  may be approximately constant, in which case  $\lambda$  (and thus  $\mu$ , since  $\lambda - \mu \approx 0$ ) must be approximately exponentially decreasing over time. It is, however, hard to find a reasonable explanation for why  $\lambda$  and  $\mu$  would decline exponentially at a roughly constant rate over the last 1 billion years, while  $N$  remains constant. Third, both  $N$  and  $\lambda$  vary substantially over time, in such a way that their ratio declines approximately exponentially. For generic varying  $\lambda(t)$  and  $N(t)$  this is highly unlikely, unless either  $\lambda$  is largely determined by  $N$  (and  $N$  varies such that  $\lambda(N(\tau))/N(\tau)$  is exponentially declining), or  $N$  is largely determined by  $\lambda$  (and  $\lambda$  varies such that  $\lambda(\tau)/N(\lambda(\tau))$  is exponentially declining). While the second and third scenario are mathematically perfectly valid, they are quite idiosyncratic and more complex than the first scenario, which solely requires that speciation and extinction rates be approximately constant over time. Further, both scenarios would not result in such a good fit of a constant-rates model as observed in this study (Figs. 2a,b and 21a–c). Hence, based on parsimony, and based on the fitted models, it is much more likely that indeed  $N$  grew approximately exponentially over time, with approximately constant speciation and extinction rates.

## S.5 Fitting homogenous-rate models to heterogenous-rate trees

As described in the main article, over geological time global microbial diversification appears to be well-described by simple branching models with homogenous (i.e. clade-independent) speciation and extinction rates (also known as “clock-like” models). At some sufficiently fine taxonomic scale, however, speciation and extinction rates may be clade-dependent, since speciation and extinction are coupled to ecological dynamics that vary across clades. The question thus arises how homogenous speciation/extinction rates, fitted in the aforementioned models, should be interpreted. To investigate this question, we simulated trees in which speciation and extinction rates were themselves evolving (and thus clade-dependent) traits, and then fitted models with homogenous speciation/extinction rates to the corresponding LTTs. In the simulated trees, the speciation rate  $\lambda$  and extinction rate  $\mu$  were modeled as independent Brownian motions evolving along the tree edges and constrained in a finite interval via reflection. Thus, at any moment in time the tree-wide absolute speciation (or extinction) rate was the sum of speciation (or extinction) rates of all extant tips.

Simulations were performed using the function `generate_tree_with_evolution_rates` in the R package `castor`<sup>17</sup>. Each simulated tree initially contained 10,000 tips, but was subsequently rarefied at some fraction  $\rho$  that was chosen randomly and uniformly between 0.1 and 1 to reflect incomplete sampling in our real trees. For each simulated tree, the maximum allowed speciation rate,  $\lambda_{\max}$ , was chosen randomly and

uniformly between 0 and 1 speciations per lineage per million years (S/LMa). The maximum allowed extinction rate,  $\mu_{\max}$ , was chosen randomly and uniformly between 0 and  $\lambda_{\max}$ . The minimum allowed speciation and extinction rates,  $\lambda_{\min}$  and  $\mu_{\min}$ , were set to 0. The diffusivity of the Brownian motion model for  $\lambda$ , denoted  $D_\lambda$ , was roughly chosen such that within a given expected number of speciations along a lineage (“rate memory”, denoted  $M$ ),  $\lambda$  would drift to a standard deviation equal to  $\frac{1}{2}(\lambda_{\max} + \lambda_{\min})$ , that is:

$$D_\lambda = \frac{1}{2M} (\lambda_{\min} + \lambda_{\max}) \cdot \left[ \frac{1}{2} (\lambda_{\max} + \lambda_{\min}) \right]^2. \quad (30)$$

The diffusivity  $D_\mu$  was chosen in a similar way:

$$D_\mu = \frac{1}{2M} (\lambda_{\min} + \lambda_{\max}) \cdot \left[ \frac{1}{2} (\mu_{\max} + \mu_{\min}) \right]^2. \quad (31)$$

Note that  $M$  can be interpreted as the approximate expected phylogenetic depth (number of consecutive branching points) at which  $\lambda$  and  $\mu$  are conserved. For example, for the parameters used in this study, a rate memory of  $M = 100$  corresponds to a diffusivity  $D_\lambda = 0.00125 \text{ (S/LMa)}^2/\text{Ma}$  and a phylogenetic conservatism of  $\sim 200 \text{ Ma}$  for  $\lambda$ .

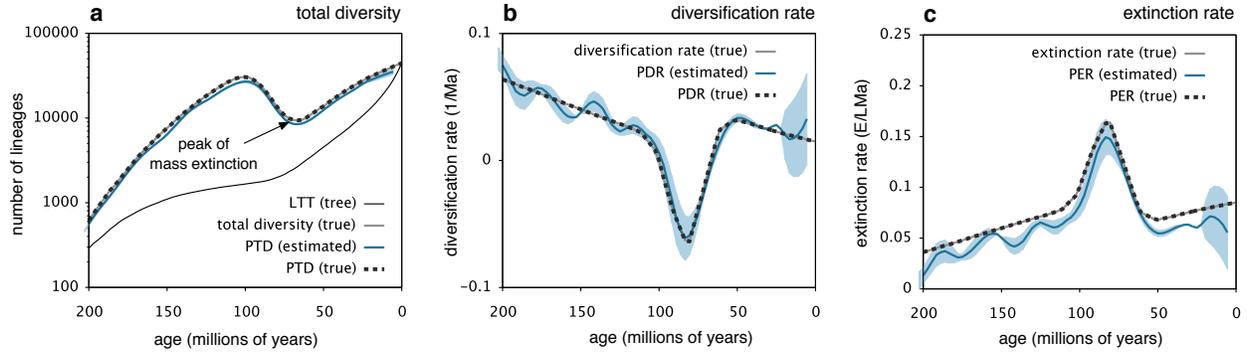
We investigated the suitability of homogenous rate models for representing speciation/extinction dynamics in these simulated trees, as well as the relationship between the fitted homogenous  $\lambda$ ,  $\mu$  and the probability distributions of the (clade-specific)  $\lambda$ ,  $\mu$  across tips of the simulated trees. During model fitting, the sampling fraction  $\rho$  applied to each simulated tree was assumed to be known. We repeated our investigation for various rate memories  $M$ , each time using 100 simulated trees.

In all cases, homogenous-rate models with constant  $\lambda$  and  $\mu$  fitted the simulated LTTs well (mean relative deviations below 10%). This suggests that timetrees, in which speciation and extinction rates are evolving heritable traits and thus not clock-like, can still be described by models with homogenous rates.

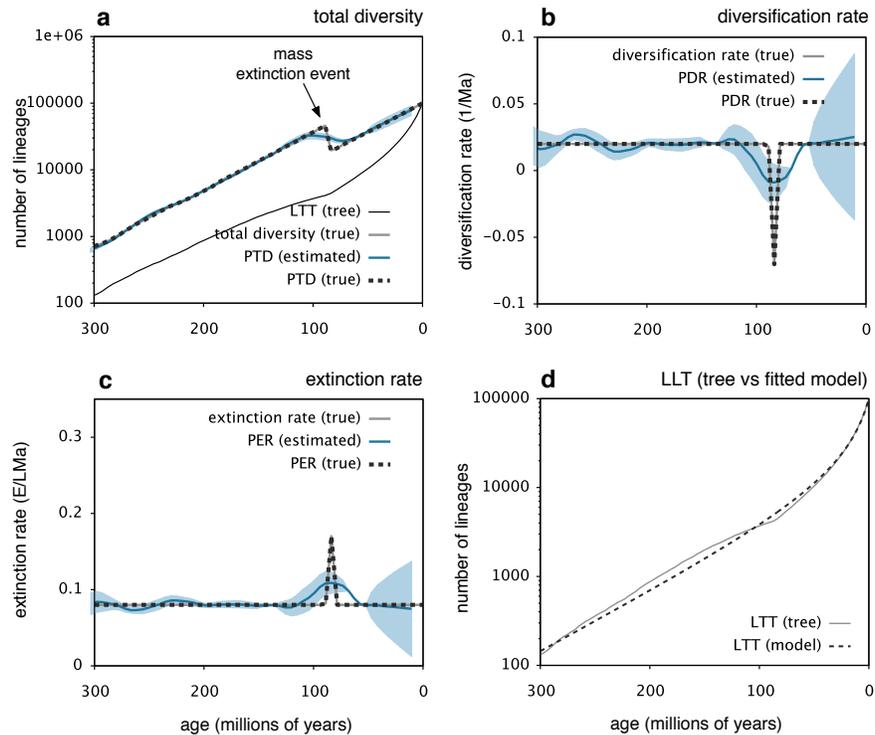
As shown in Supplementary Fig. 20, the  $\lambda$  and  $\mu$  estimated from homogenous-rate models approximately correspond to the average  $\lambda$  and  $\mu$  across extant tips in the tree. For  $\lambda$ , this approximation is quite accurate across all of our tested scenarios (mean  $R^2 > 0.95$ , Figs. 20a–c). For  $\mu$ , the accuracy of this approximation can vary substantially, depending on the rate of evolution of  $\lambda$  and  $\mu$  (i.e., their phylogenetic conservatism). This observation is consistent with previous findings that heterogenous diversification rates lead to inflated estimates of extinction rates when the latter are obtained from homogenous-rate models<sup>18</sup>. Our results show that the intensity of this bias is related to the rate of evolution of the traits determining speciation and extinction rates; the bias becomes stronger when  $\lambda$  and  $\mu$  exhibit strong phylogenetic conservatism. This relationship can be explained by the fact that our formulas have been derived under the assumption that the probabilities of two sister-lineages going extinct are independent ( $\mathbb{M}_{02}(t, T) = E(t, T)^2$  in Eq. 6). This assumption is only true if  $\lambda$  and  $\mu$  are homogenous, or if  $\lambda$  and  $\mu$  evolve sufficiently fast. That said, even for very strongly conserved  $\lambda$  and  $\mu$  (rate memory  $\sim 100$ ; Fig. 20f), the relative error is in the order of 80–100%, that is, even under such a scenario our methods are suitable for estimating average extinction rates to order of magnitude accuracy. In conclusion, for any given time point, rates estimated here should be interpreted as approximate current averages over the entire tree, but it is unknown whether rates were homogenous (i.e., equal to the average) or heterogenous (i.e., exhibiting a variance around the average) at that time point.

## S.6 Implications for reconstructing diversification from phylogenies

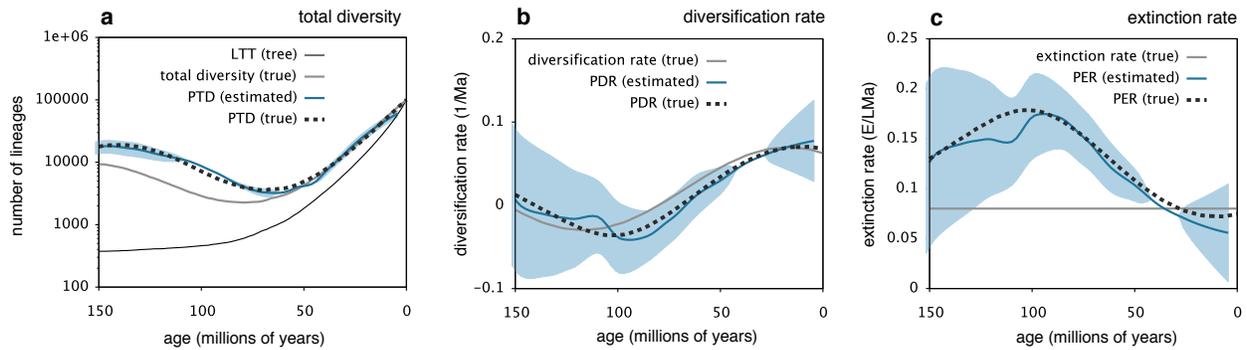
Debate exists over the information that can be extracted from molecular phylogenies in the absence of fossil data, with some authors even suggesting that extinctions cannot possibly be reconstructed from phylogenies alone<sup>8,18–20</sup>. Indeed, in principle similar LTTs (or branch length distributions) can be generated under different combinations of speciation and extinction rates varying over time in specific ways. In the most extreme case, a tree generated by a speciation-extinction process could in principle also be the result of a pure speciation process, with a particular idiosyncratically varying speciation rate over time generating the same LTT. This ambiguity is reflected in the PDR — a curve with similar information content as the LTT, representing an inseparable combination of the diversification rate and the relative rate of change of the speciation rate (Eq. 1). The problem is further amplified in some analyses that use simple (single-valued) summary statistics, such as the “ $\gamma$ -statistic”, to describe the overall shape of the LTT and then attempt to associate different values of the statistic with different diversification scenarios<sup>6,19,21</sup>. Importantly, molecular phylogenies completely lack information on the diversification dynamics of clades that are now entirely extinct. Hence, any reconstruction of past bacterial diversification assumes that the speciation/extinction rates within extinct clades were similar to (or distributed similarly to, if non-homogenous) those of extant clades. Despite the limitations of molecular phylogenies and reservations by some authors<sup>8,18</sup>, extinction does leave characteristic traces in phylogenies that tend to be markedly different from typical scenarios lacking extinction, as demonstrated here (e.g., Fig. 1) and elsewhere<sup>20</sup>. Molecular phylogenies thus contain information on past extinctions, albeit in a convoluted format, and this information can be extracted using methods such as ours.



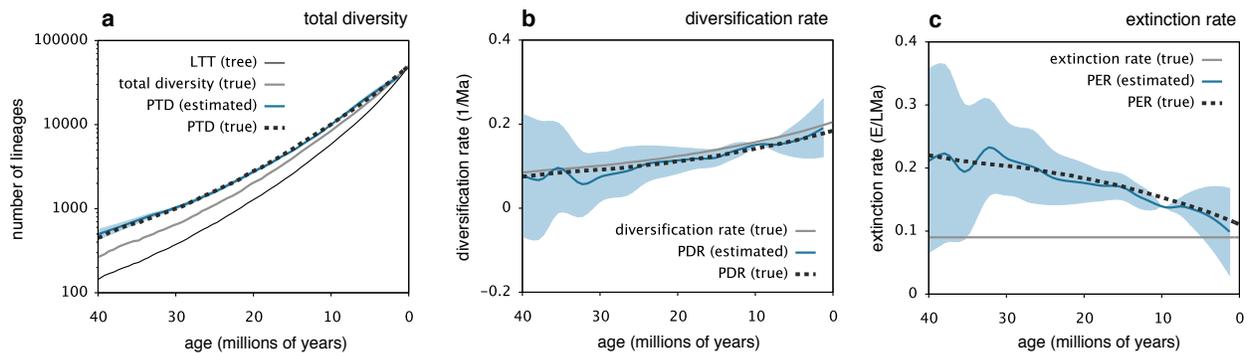
**Supplementary Figure 1: Demonstration of non-parametric estimation (mass extinction event).** Simulation and non-parametric analysis of a tree with constant speciation rate and variable extinction rate (including diversity dependence and a mass extinction event about 80 Ma ago). (a) Lineages through time (LTT) in the tree (black continuous curve), true total diversity (grey continuous curve), true pulled total diversity (PTD, dashed curve) and estimated PTD (blue continuous curve). The footprint of the mass extinction event is visible in the tree’s LTT. (b) Diversification rate (grey continuous curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue continuous curve, Eq. 25). (c) True extinction rate (grey continuous curve), true pulled extinction rate (PER, dashed curve) and estimated PER (blue continuous curve). Estimated quantities in A–C were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter. All “pulled” variables (PTD, PDR, PER) are equivalent to their non-pulled variants (total diversity, diversification rate, extinction rate) because the speciation rate is constant over time; consequently, curves showing true pulled variables (dashed lines) overlap completely with curves showing non-pulled variables (grey continuous lines).



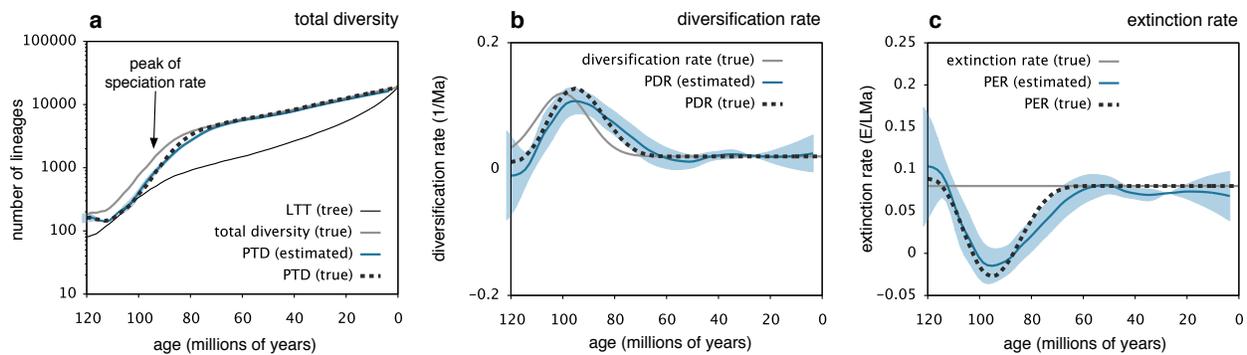
**Supplementary Figure 2: Demonstration of non-parametric estimation (sharp mass extinction event).** Simulation and non-parametric analysis of a tree with constant speciation rate and an extinction rate subject to a sharp temporary increase about 400 Ma ago (representing a sharp mass extinction event). (a) Lineages through time (LTT) in the tree (black continuous curve), true total diversity (grey continuous curve), true pulled total diversity (PTD, dashed curve) and estimated PTD (blue continuous curve). (b) Diversification rate (grey continuous curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue continuous curve, Eq. 25). (c) True extinction rate (grey continuous curve), true pulled extinction rate (PER, dashed curve) and estimated PER (blue continuous curve). (d) LTT in the tree (continuous curve), compared to the LTT predicted by a fitted constant-rates model (dashed curve). The deviation of the model’s prediction from the tree’s LTT is most clear as an inflection point at age  $\sim 400$  Ma. Estimated quantities in (a–c) were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter. The sharp spike in the extinction rate introduces a sharp spike in the PDR (b) and PER (c), which, however, is damped by the noise filter during estimation. All “pulled” variables (PTD, PDR, PER) are equivalent to their non-pulled variants (total diversity, diversification rate, extinction rate) because the speciation rate is constant over time; consequently, curves showing true pulled variables (dashed lines) overlap completely with curves showing non-pulled variables (grey continuous lines).



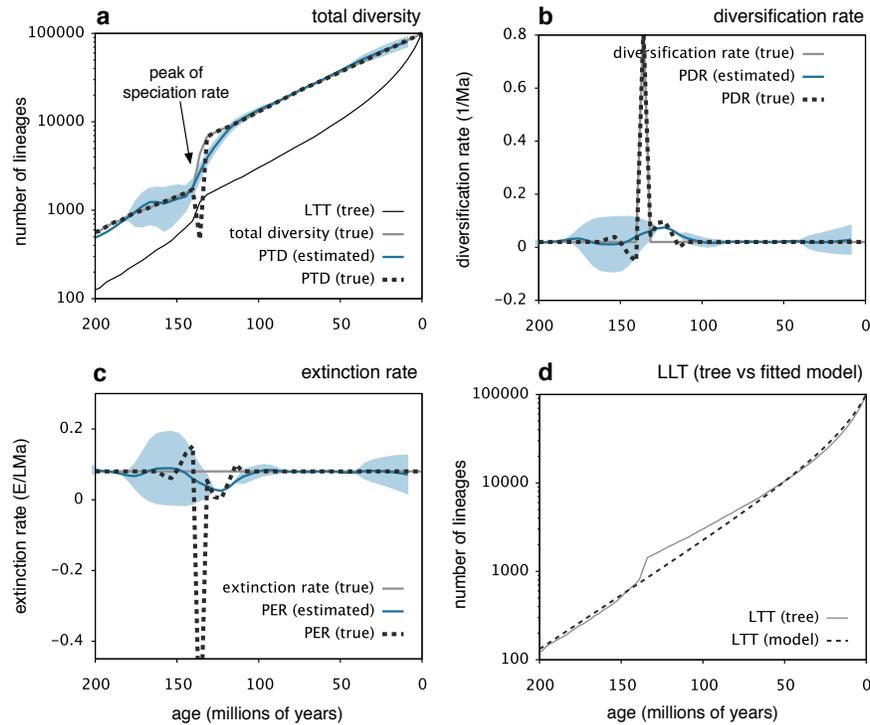
**Supplementary Figure 3: Demonstration of non-parametric estimation (oscillating speciation rate).** (a) Lineages through time for a simulated tree with strongly oscillating speciation rate ( $\pm 50\%$  around the mean) and variable extinction rate (including diversity dependence and a mass extinction event about 60 Ma in the past). Dashed curve: Lineages through time (LTT) in the tree. Grey continuous curve: True total diversity. Blue continuous curve: Pulled total diversity (PTD), estimated from the LTT. The fluctuations in the speciation rate lead to a deviation of the PTD from the total diversity. Nevertheless, the mass extinction event is clearly reflected in the PTD. (b) Diversification rate (thin curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue curve, Eq. 25), during the same simulation as in (a). (c) True extinction rate (grey curve) and pulled extinction rate (PER, blue curve) estimated from the LTT (Eq. 22). The fluctuations in the speciation rate clearly distort the pulled extinction rate. All estimated quantities were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter.



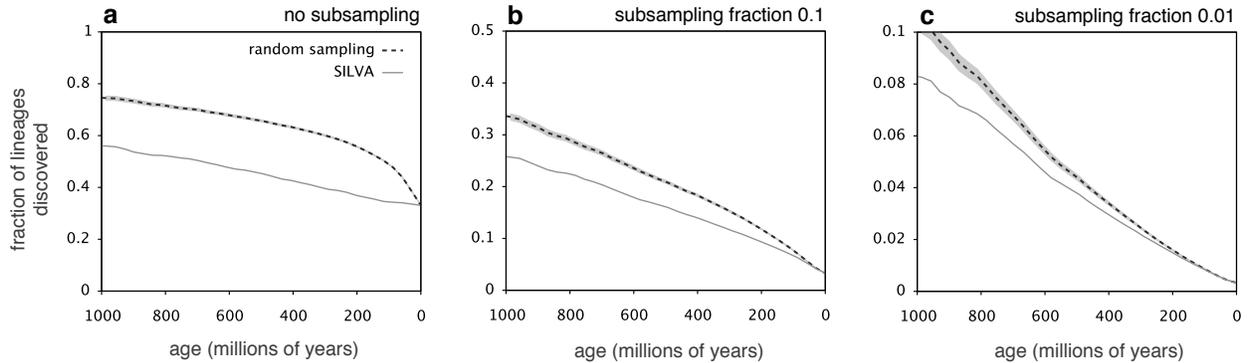
**Supplementary Figure 4: Demonstration of non-parametric estimation (diversity-dependent speciation rate).** (a) Lineages through time for a simulated tree with density-dependent speciation rate ( $\lambda \propto N^{0.1}$ ) and constant extinction rate. Dashed curve: Lineages through time (LTT) in the tree. Grey continuous curve: True total diversity. Blue continuous curve: Pulled total diversity (PTD), estimated from the LTT. (b) Diversification rate (thin curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue curve, Eq. 25), during the same simulation as in (a). (c) True extinction rate (grey curve) and pulled extinction rate (PER, blue curve) estimated from the LTT (Eq. 22). The density-dependence of the speciation rate introduces a trend in the pulled extinction rate. All estimated quantities were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter.



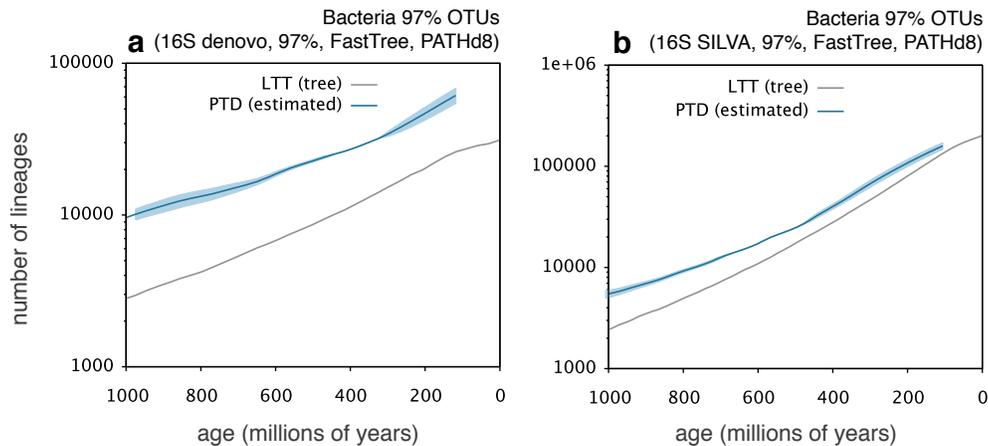
**Supplementary Figure 5: Demonstration of non-parametric estimation (speciation spike).** Simulation and non-parametric analysis of a tree with a temporarily increased speciation rate about 90 million years ago and a constant extinction rate. (a) Lineages through time (LTT) in the tree (black continuous curve), true total diversity (grey continuous curve), true pulled total diversity (PTD, dashed curve) and estimated PTD (blue continuous curve). (b) Diversification rate (grey continuous curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue continuous curve, Eq. 25), during the same simulation as in (a). (c) True extinction rate (grey continuous curve), true pulled extinction rate (PER, dashed curve) and estimated PER (blue continuous curve). The temporarily increased speciation rate introduces a similarly strong distortion in the PER. All estimated quantities were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter.



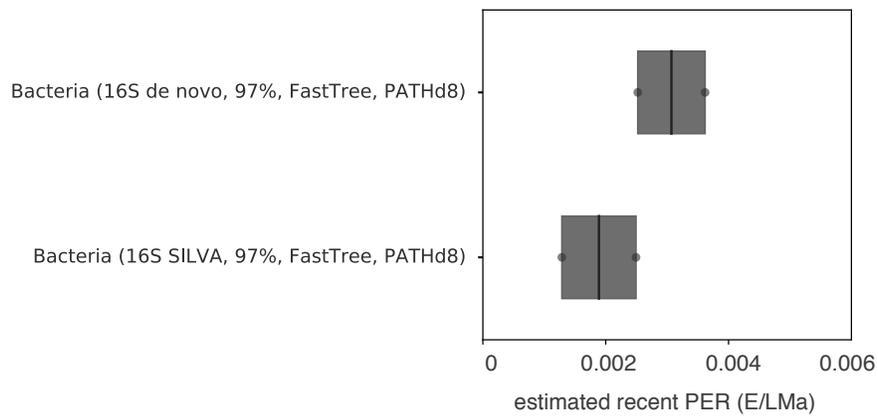
**Supplementary Figure 6: Demonstration of non-parametric estimation (brief speciation spike).** Simulation and non-parametric analysis of a tree with a briefly but strongly increased speciation rate about 140 million years ago and a constant extinction rate. (a) Lineages through time (LTT) in the tree (black continuous curve), true total diversity (grey continuous curve), true pulled total diversity (PTD, dashed curve) and estimated PTD (blue continuous curve). (b) Diversification rate (grey continuous curve), true pulled diversification rate (PDR, dashed curve) and estimated PDR (blue continuous curve, Eq. 25), during the same simulation as in (a). (c) True extinction rate (grey continuous curve), true pulled extinction rate (PER, dashed curve) and estimated PER (blue continuous curve). (d) LTT in the tree (continuous curve), compared to the LTT predicted by a fitted constant-rates model (dashed curve). The deviation between the tree's LTT and the model prediction is clearly visible. Estimated quantities in (a–c) were noise filtered similarly to the bacterial trees (Fig. 2); blue shades indicate standard errors of the noise filter. The brief spike in the speciation rate introduces a sharp spike in the PDR (b) and PER (c), which, however, is damped by the noise filter during estimation. Note that curves showing true pulled variables (dashed lines) overlap largely with curves showing non-pulled variables (grey continuous lines).



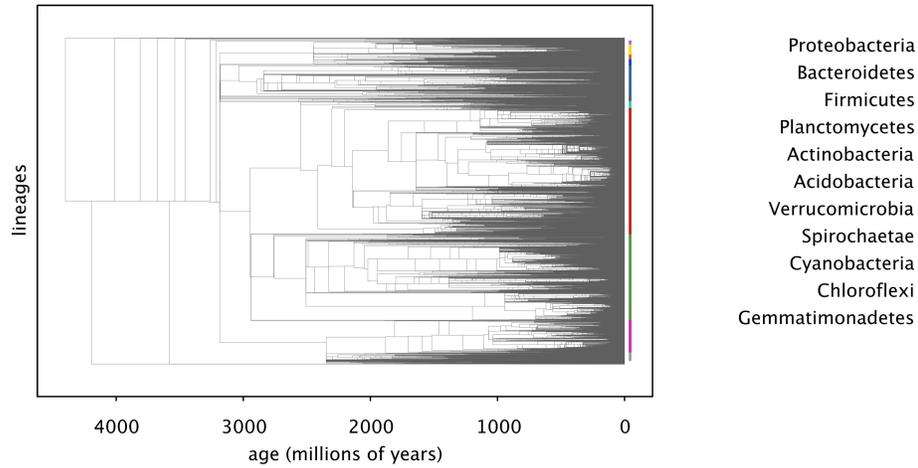
**Supplementary Figure 7: Fractions of discovered lineages over time (SILVA vs random OTU sampling).** (a) Fraction of bacterial lineages in the *de novo* dataset (162,371 *de-novo* OTUs at 99% 16S rRNA similarity), previously discovered (i.e. matched to SILVA at 99% similarity), as a function of lineage age (continuous line). The dashed line shows the expectation under the null model of random non-phylogenetically biased OTU discovery and the shading indicates the corresponding standard deviation. (b) Similar to a, but after subsampling SILVA down to a fraction of 10%. (c) Similar to a, but after subsampling SILVA down to a fraction of 1%. The deviation from the null model is due to a non-random (phylogenetically biased) representation of clades in SILVA.



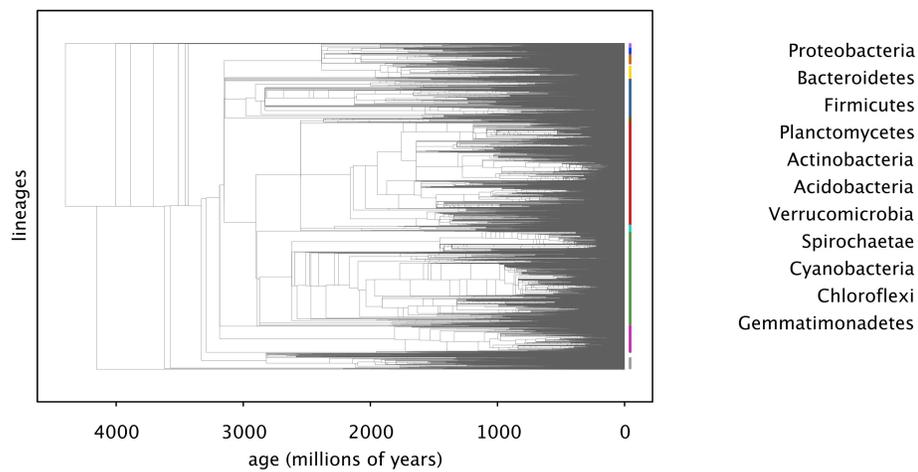
**Supplementary Figure 8: Bacterial diversities over time (based on 97%-OTUs).** Lineages through time (LTT, grey curves), compared to estimated pulled total diversities (PTD, blue curves) of bacteria, based on timetrees of 16S rRNA OTUs delineated at 97% similarity. (a) Based on partial-length (V4) *de novo* 97%-OTUs. (b) Based on full-length 97%-OTUs from SILVA. Summaries of timetree construction methods are indicated in brackets.



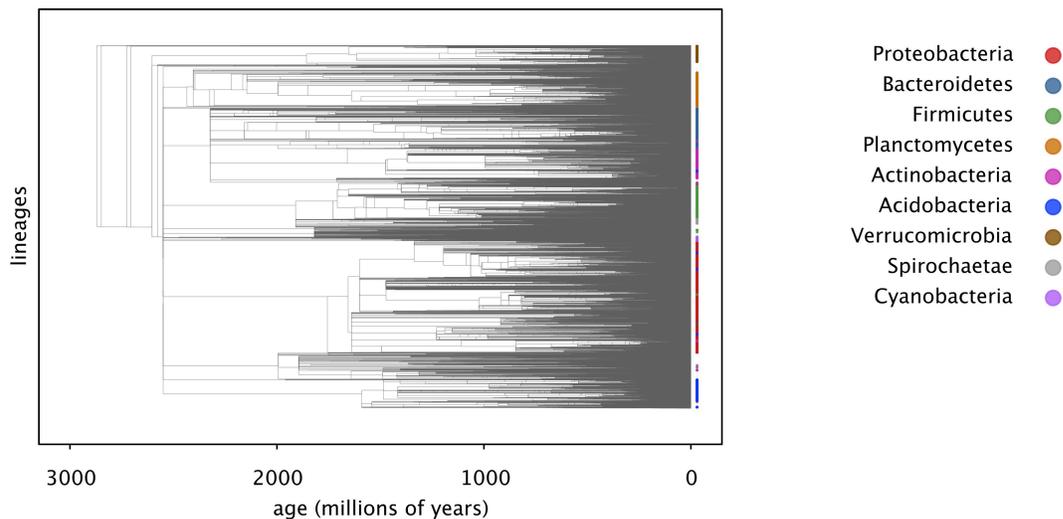
**Supplementary Figure 9: Recent pulled extinction rates (bacterial 97%-OTUs).** Non-parametrically estimated recent pulled extinction rates of bacteria, based on timetrees of 16S rRNA OTUs delineated at 97% similarity. Each box corresponds to a different tree, and comprises results obtained by assuming various numbers of total extant OTUs (Supplementary Table 6). Summaries of timetree construction methods are indicated in brackets.



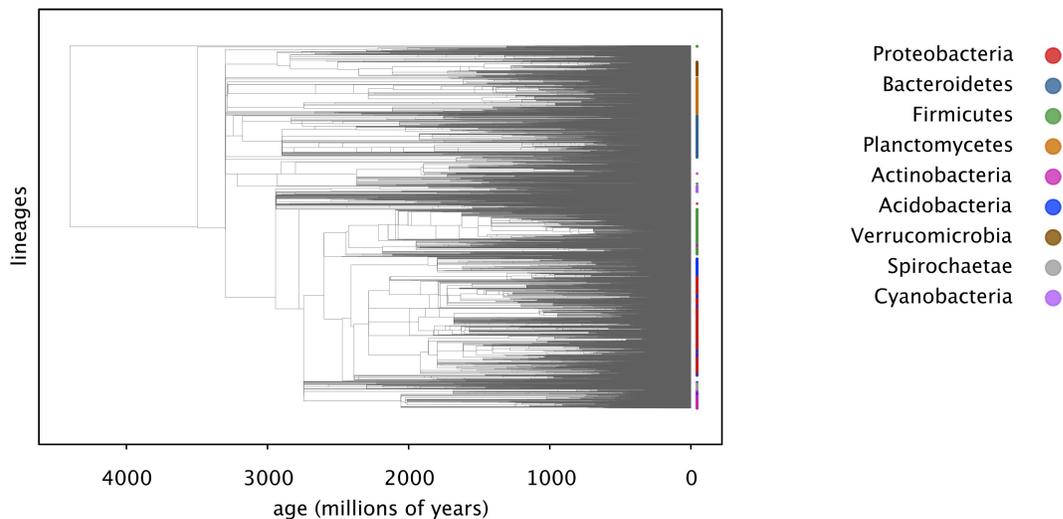
**Supplementary Figure 10:** Bacterial timetree (16S SILVA, FastTree, PATHd8), used for the diversification analysis. Major phyla (based on SILVA 128)<sup>22</sup> are indicated as colored segments.



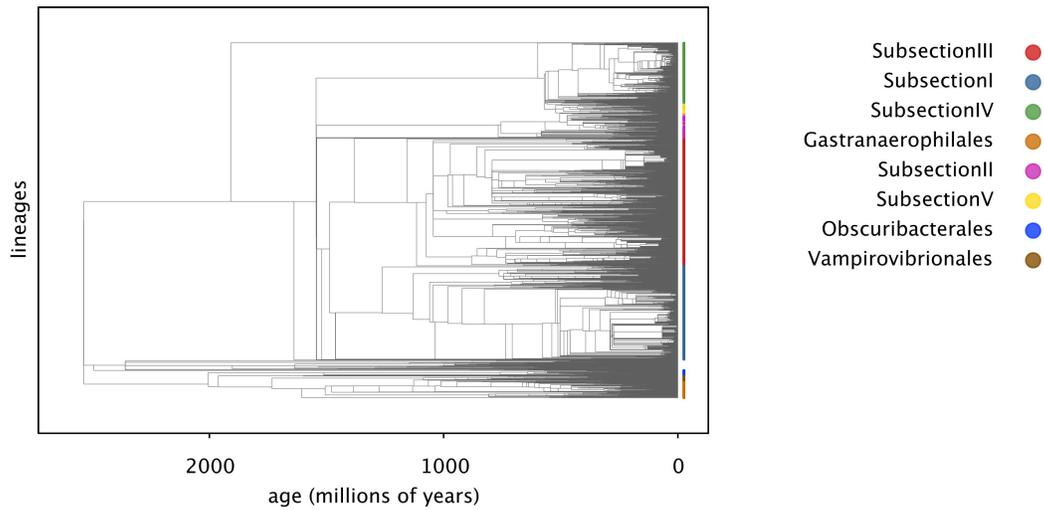
**Supplementary Figure 11:** Bacterial timetree (16S SILVA, 97%, FastTree, PATHd8), used for the diversification analysis. Major phyla (based on SILVA 128)<sup>22</sup> are indicated as colored segments.



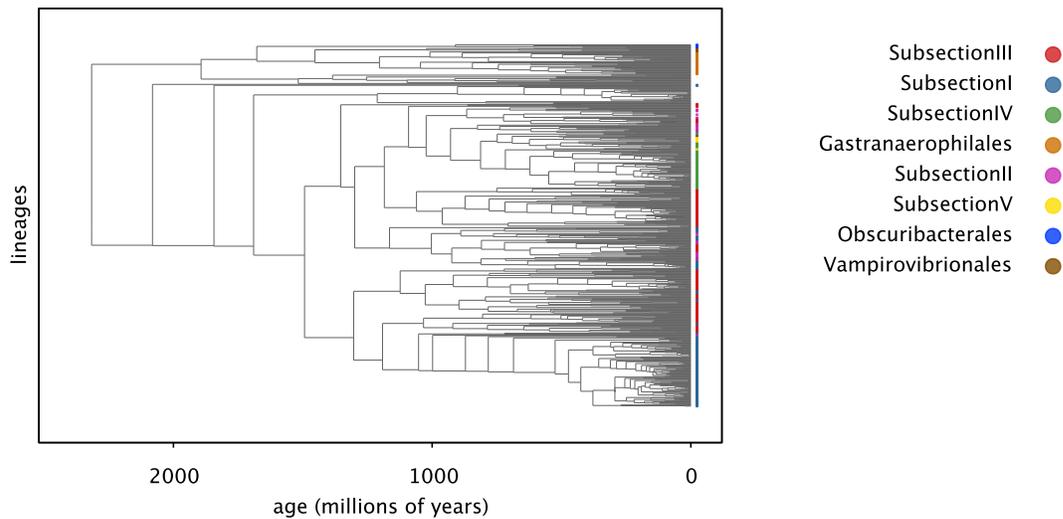
**Supplementary Figure 12:** Bacterial timetree (16S V4 de novo, FastTree, PATHd8), used for the diversification analysis. Major phyla (based on SILVA 128)<sup>22</sup> are indicated as colored segments. Note that some tips belong to unidentified phyla, and are thus not highlighted in color. Also note that the phylogeny is not absolutely congruent with taxonomic assignments.



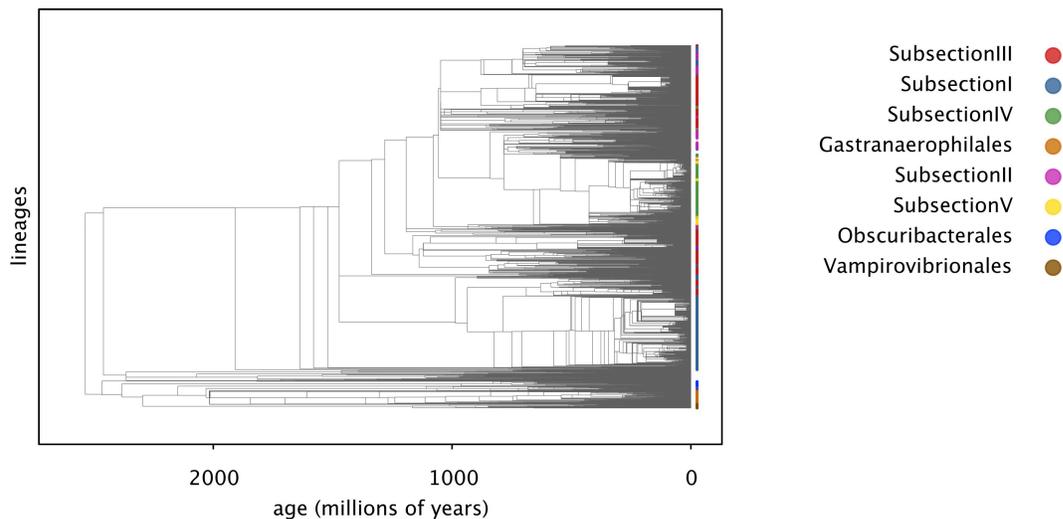
**Supplementary Figure 13:** Bacterial timetree (16S de novo, 97%, FastTree, PATHd8), used for the diversification analysis. Major phyla (based on SILVA 128)<sup>22</sup> are indicated as colored segments. Note that some tips belong to unidentified phyla, and are thus not highlighted in color. Also note that the tree is not absolutely congruent with taxonomic assignments.



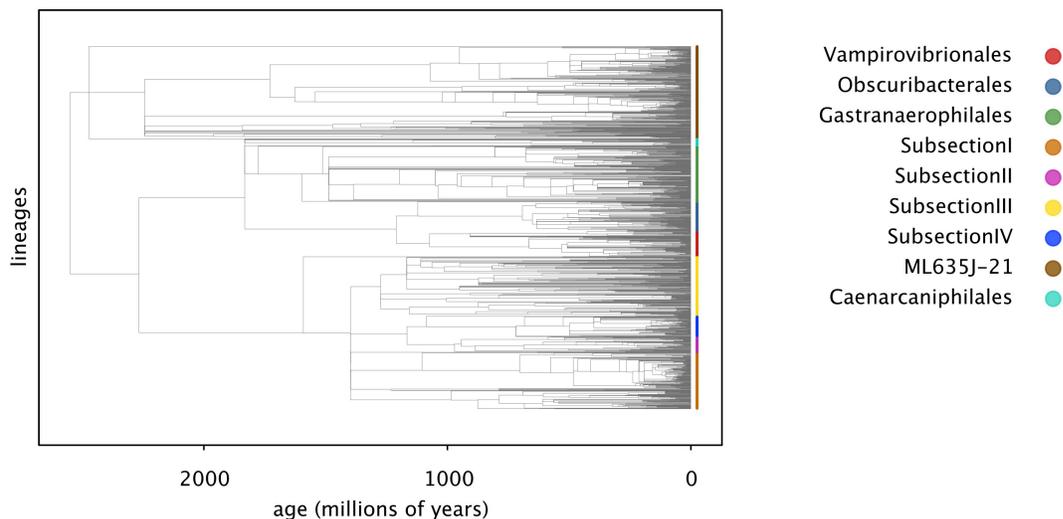
**Supplementary Figure 14:** Cyanobacterial timetree (16S SILVA, FastTree, BEAST+PATHd8), used for the diversification analysis. Major classes or orders (based on SILVA 128)<sup>22</sup> are indicated as colored segments. Note that some tips belong to unidentified classes or orders, and are thus not highlighted in color.



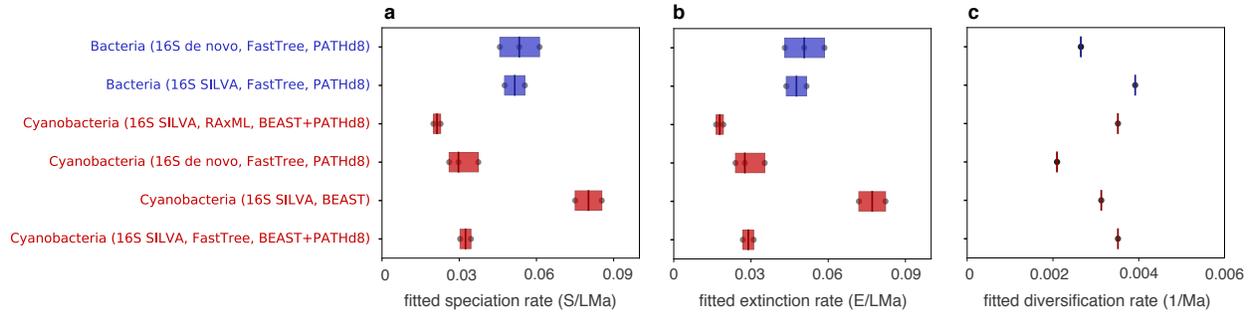
**Supplementary Figure 15:** Cyanobacterial timetree (16S SILVA, BEAST), used for the diversification analysis. Major classes or orders (based on SILVA 128)<sup>22</sup> are indicated as colored segments. Note that some tips belong to unidentified classes or orders, and are thus not highlighted in color. Also note that the tree is not absolutely congruent with taxonomic assignments.



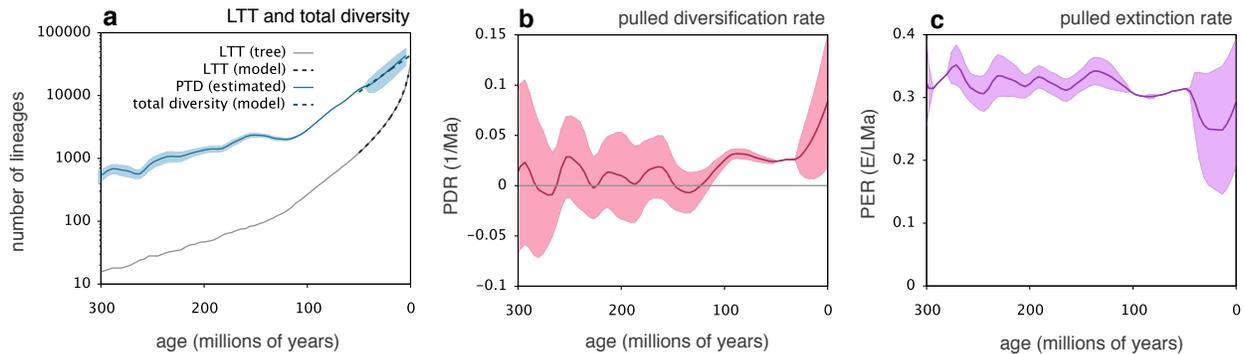
**Supplementary Figure 16:** Cyanobacterial timetree (16S SILVA, RAXML, BEAST+PATHd8), used for the diversification analysis. Major classes or orders (based on SILVA 128)<sup>22</sup> are indicated as colored segments. Note that some tips belong to unidentified classes or orders, and are thus not highlighted in color. Also note that the tree is not absolutely congruent with taxonomic assignments.



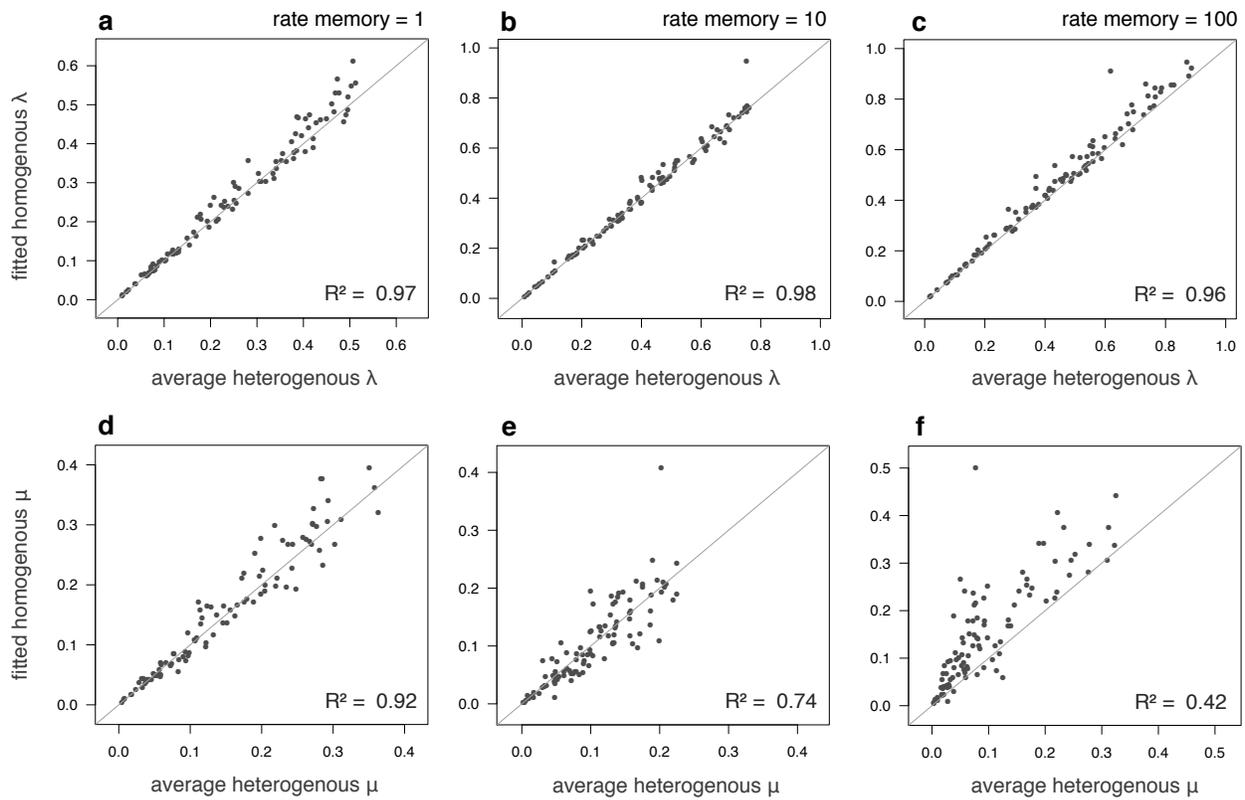
**Supplementary Figure 17:** Cyanobacterial timetree (16S de novo, FastTree, PATHd8), used for the diversification analysis. Major classes or orders (based on SILVA 128)<sup>22</sup> are indicated as colored segments.



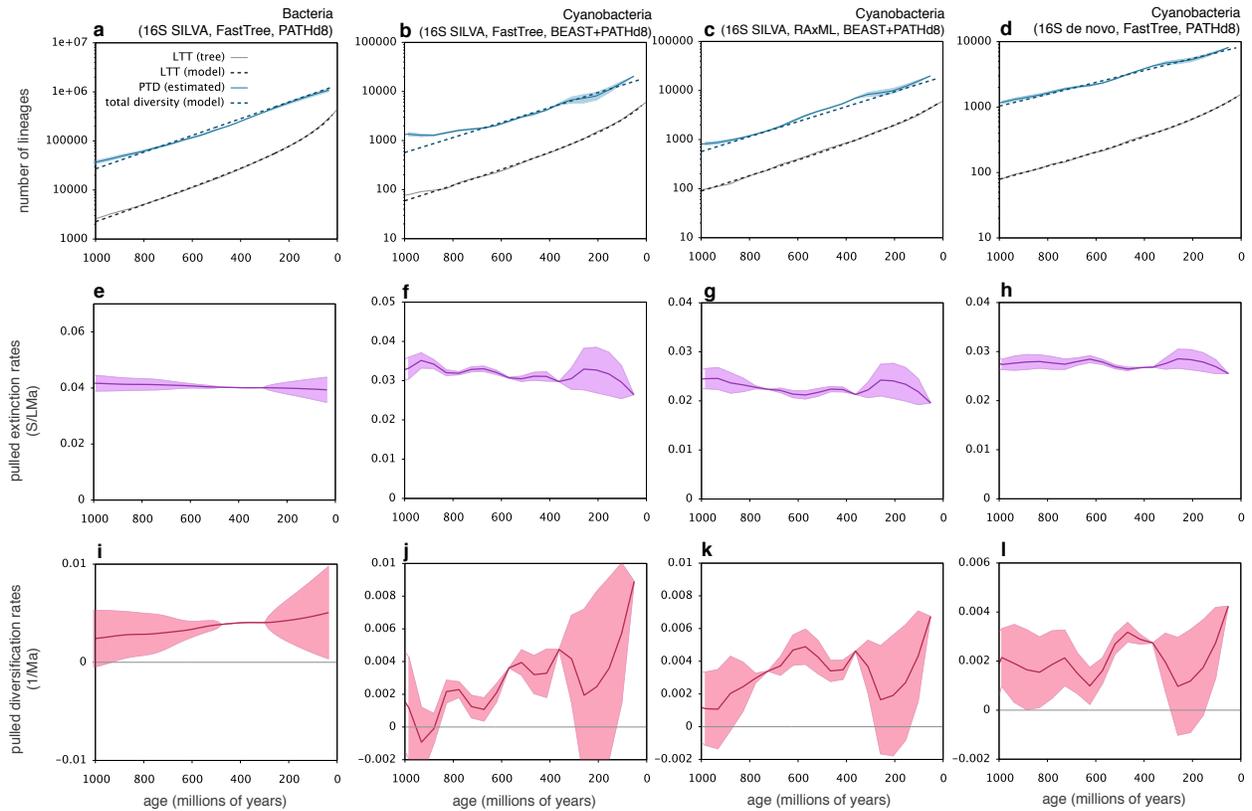
**Supplementary Figure 18: Estimated long-term speciation, extinction and diversification rates.** Speciation rates (a), extinction rates (b) and diversification rates (c), estimated for various taxa and using various dated timetrees (one box per timetree). Estimates are obtained by fitting cladogenic models over the past 1 billion years (for estimates over shorter time intervals see Fig. 3 in the main article). Each box comprises results obtained by assuming various numbers of total extant OTUs (Supplementary Tables 2, 4 and 3). Tree labels and boxes are colored by taxon. Summaries of timetree sources or construction methods are indicated in brackets (see Methods for details).



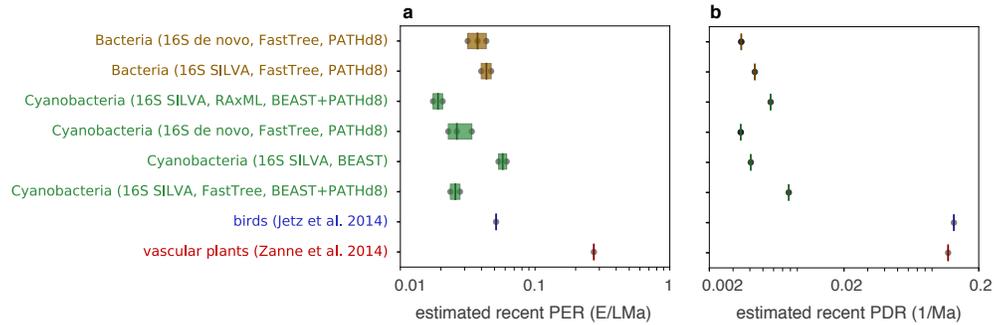
**Supplementary Figure 19: Diversification of vascular plants over time.** (a): Lineages through time for vascular plants (grey continuous line), compared to a speciation-extinction model fitted to a recent age interval (grey dashed lines). The blue continuous curve shows non-parametrically estimated pulled total diversities (PTDs), the blue dashed curve shows total diversities predicted by the fitted model. (b) Pulled extinction rates over time, estimated non-parametrically. (c): Pulled diversification rates, estimated non-parametrically. The tree was taken from Zanne *et al.*<sup>23</sup>.



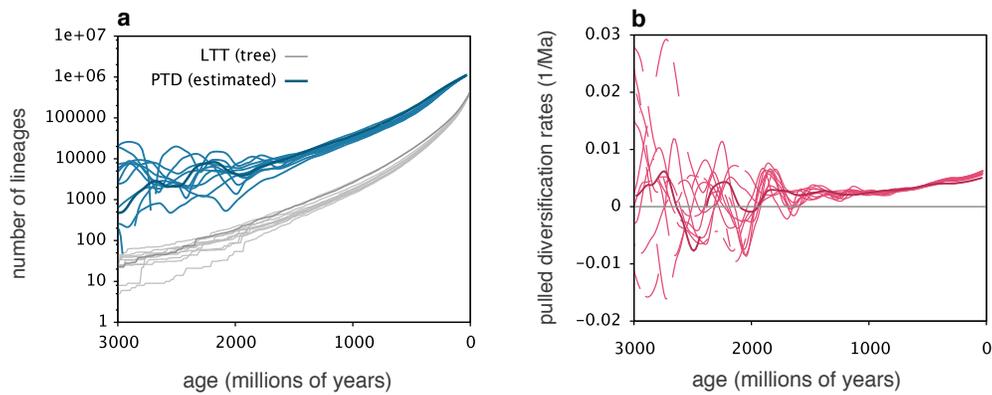
**Supplementary Figure 20: Fitting homogenous-rate models to heterogenous-rate trees.** (a–c) Average speciation rates across extant tips of simulated trees with evolving  $\lambda$  (horizontal axes), compared to fitted homogenous speciation rates (vertical axis). One point per simulated tree. Rate memories of speciation rates were 1 in (a), 10 in (b) and 100 in (c); recall that a lower rate memory corresponds to a faster evolving  $\lambda$  and  $\mu$ . Diagonals are shown for reference. Fractions of variance in the vertical axis explained by the horizontal axis ( $R^2$ ) are written in each figure. (d–f) Analogous to (a–c), but showing extinction rates. Methodological details provided in Supplement S.5.



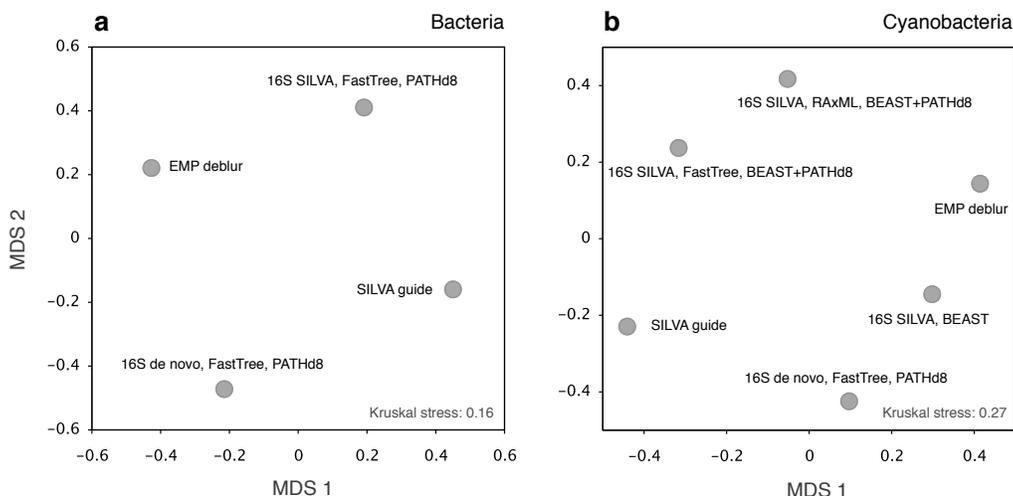
**Supplementary Figure 21: Bacterial and cyanobacterial diversification over time (alternative trees).** (a–d): Lineages through time (LTT) of a bacterial timetree (a) and three alternative cyanobacterial timetrees (b–d), compared to predictions by speciation-extinction models fitted over the last 1 billion years (grey dashed curves). Blue continuous curves show non-parametrically estimated pulled total diversities (PTDs), blue dashed curves show total diversities predicted by the fitted models. Note that each tree only comprises a subset of extant species, and thus the true extant diversity (right-most point on blue curves) is only an estimate (overview in Table 2). (e–h) Pulled extinction rate (PER) over time, estimated non-parametrically from the trees used in (a–d) (Supplement S.1.3). (i–l): Pulled diversification rate (PDR), estimated non-parametrically from the same trees as in (a–d). In all figures, shades indicate standard errors of noise-filtered estimates. Summaries of tree construction methods are indicated in brackets (see Methods for details).



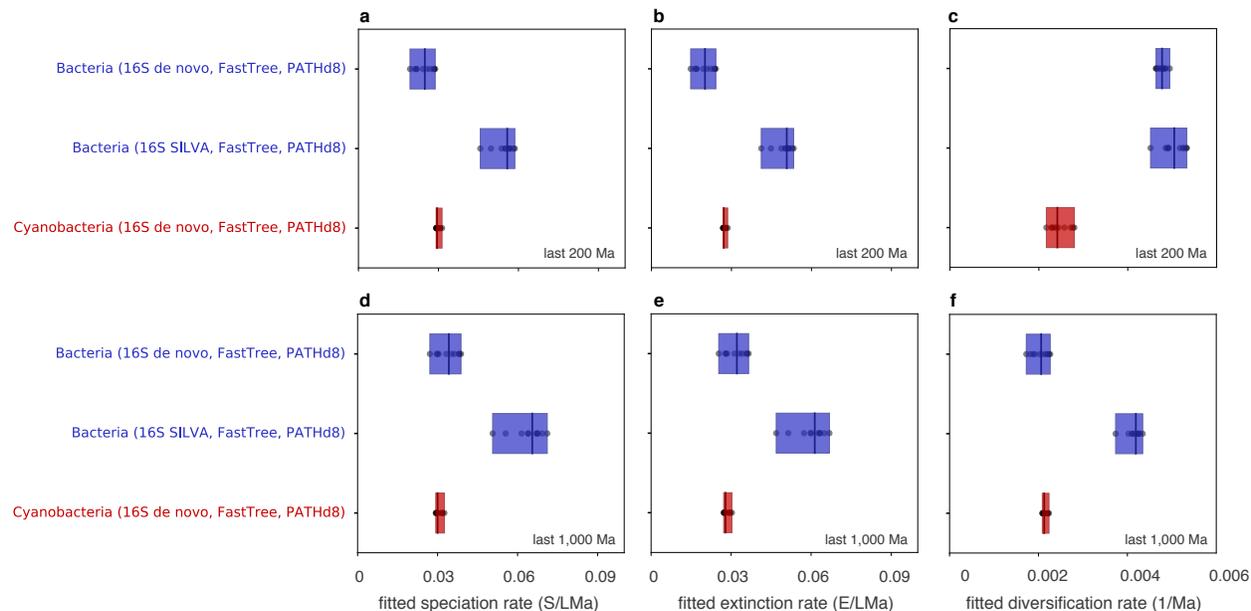
**Supplementary Figure 22: Estimated recent pulled extinction and diversification rates.** Recent pulled extinction rates (a) and recent pulled diversification rates (b), estimated non-parametrically for various taxa and using various dated timetrees (one box per timetree). Each box corresponds to a different tree, and comprises results obtained by assuming various numbers of extant bacterial or cyanobacterial OTUs (Supplementary Tables 2, 4 and 3). Tree labels and boxes are colored by taxon. Summaries of timetree sources or construction methods are indicated in brackets (see Methods for details).



**Supplementary Figure 23: Bacterial diversities over time (sensitivity to dating).** (a) Lineages through time (LTT, grey curves), compared to estimated pulled total diversities (PTD, blue curves) of bacteria, based on timetrees constructed from 16S sequences in SILVA (timetree “16S SILVA, FastTree, PATHd8” and random variants). Individual curves correspond to slightly different timetrees, dated using randomly varied dating constraints for purposes of sensitivity analysis (see Methods for details). (b) Estimated pulled diversification rates (PDR) for the same trees as in (a). In both figures, darker curves correspond to the point-estimate timetree (“16S SILVA, FastTree, PATHd8”) discussed in the main article. Observe that the variability between tree variants as well as spurious fluctuations, and thus the uncertainty in the estimated PTD and PDR, increases drastically for ages older than 1.5 billion years.



**Supplementary Figure 24: Comparison of tree topologies.** Multidimensional scaling plot showing pairwise distances between bacterial (a) or cyanobacterial (b) timetrees investigated in the main article (one point per tree). Points closer to each other indicate a greater similarity. Tree distances are based on the Robinson-Foulds metric<sup>24</sup>, which measures the difference in tree topology when restricted to tips common to both trees being compared. Tree summaries are indicated next to each point (see Methods for details). “EMP deblur” is taken from<sup>10</sup>, and “SILVA guide” refers to the SILVA SSU guide tree (release 128)<sup>22</sup>. For a list of pairwise tree distances, see Supplementary Table 7.



**Supplementary Figure 25: Sensitivity of rate estimates to dating constraints.** Speciation rates (a,d), extinction rates (b,e) and diversification rates (c,f), estimated for Bacteria or Cyanobacteria, using various dated trees (one box per tree). Estimates are obtained by fitting cladogenic models over the last 200 million years (top row) or the last 1 billion years (bottom row). Each box comprises results obtained by varying dating anchors within their uncertainty intervals (Supplementary Table 5, details in Methods). Tree labels and boxes are colored by taxon. Summaries of tree sources or construction methods are indicated in brackets. Only timetrees dated using PATHd8 and based on primary constraints (Table 5) are shown.

**Supplementary Table 1: Overview of timetrees.** Overview of timetrees analyzed in the main article, including source or construction method, tree size (number of tips), the age interval considered for estimating recent speciation/extinction rates (via model fitting, Supplement S.1.2), fitted recent speciation/extinction rates ( $\lambda$  and  $\mu$ ) and mean relative deviation (MRD) of models fitted over 1 billion years. Each tree's sampling fraction ( $\rho$ ) was set to its size divided by the corresponding (independently estimated) total number of extant OTUs, as listed in Table 2.

<b>tree</b>	<b>size (tips)</b>	<b>ages (Ma)</b>	$\lambda$ <b>(S/LMa)</b>	$\mu$ <b>(E/LMa)</b>	<b>MRD</b>
Bacteria (16S SILVA, FastTree, PATHd8)	448,112	20–200	0.0440	0.0393	0.030
Bacteria (16S de novo, FastTree, PATHd8)	162,371	20–200	0.0350	0.0300	0.025
Cyanobacteria (16S SILVA, BEAST)	586	20–200	0.0651	0.0594	0.024
Cyanobacteria (16S de novo, FastTree, PATHd8)	1,579	20–200	0.0294	0.0271	0.020
Cyanobacteria (16S SILVA, FastTree, BEAST+PATHd8)	6,308	20–200	0.0340	0.0309	0.046
Cyanobacteria (16S SILVA, RAxML, BEAST+PATHd8)	6,302	20–200	0.0255	0.0226	0.035
birds <sup>25</sup>	9,993	0–20	0.172	0.0828	–
vascular plants <sup>23</sup>	31,749	5–50	0.332	0.306	–

**Supplementary Table 2:** Column 2: Fraction of *de novo* bacterial OTUs (99% similarity in the 16S rRNA gene) that were found to be represented in the SILVA v128 database (at 99% similarity). Columns 3 & 4: Number of full-length (FL) clusters and partial-length (V4) clusters, respectively, in SILVA (16S rRNA sequences clustered at 99% similarity). Columns 5 & 6: Global number of extant full-length and partial-length OTUs, respectively, based on the fraction of *de novo* OTUs represented in SILVA (columns 3 & 4 divided by column 2). See Methods for details.

<b>taxon</b>	<b>fraction <i>de novo</i> OTUs</b>	<b>FL clusters in SILVA</b>	<b>V4 clusters in SILVA</b>	<b>estimated extant FL</b>	<b>estimated extant V4</b>
Bacteria	0.329	448,112	148,471	1,360,260	450,691
Cyanobacteria	0.330	6,308	2,770	19,127	8,399
Firmucutes	0.442	115,864	34,275	264,416	77,628
Bacteroidetes	0.377	45,351	16,896	120,265	44,806
Proteobacteria	0.367	174,646	47,027	476,125	128,206
Acidobacteria	0.468	12,676	5,448	27,102	11,648
Spirochaetes	0.177	3,683	2,103	20,793	11,872
Planctomycetes	0.181	8,038	5,659	44,209	31,125

**Supplementary Table 3:** Columns 2 & 3: Number of *de novo* OTUs and number of OTUs generated from the EMP dataset, respectively. Column 4: Overlap between *de novo* OTUs and EMP OTUs (fraction of EMP OTUs matched to *de novo* OTUs at 99% identity). Column 5: Global number of extant partial-length (V4) OTUs, estimated based on the overlap between *de novo* OTUs and EMP OTUs (column 2 divided by column 4). Column 6: Estimated global number of extant full-length (FL) OTUs, based on the previously estimated number of extant V4-OTUs and the ratio of full-length over V4 OTUs in SILVA (Table 2).

taxon	# <i>de novo</i> OTUs	# EMP OTUs	overlap <i>de novo</i> vs EMP	extant V4 OTUs	extant FL OTUs
Bacteria	165,422	333,524	0.270	612,674	1,849,358
Cyanobacteria	1,598	5,547	0.149	10,725	24,424

**Supplementary Table 4:** Column 2: Fraction of metagenome-assembled-genome (MAG) 16S rRNA sequences that were found to be represented in the SILVA v128 database (at 99% similarity). Column 3: Number of full-length 16S rRNA sequence clusters in SILVA (at 99% similarity) within each considered taxon. Column 4: Global number of extant full-length OTUs, estimated based on the fraction of MAGs represented in SILVA (column 3 divided by column 2). See Methods for details.

taxon	fraction MAGs	FL clusters in SILVA	estimated extant FL
Bacteria	0.283	448,112	1,583,740
Cyanobacteria	0.375	6,308	16,821

**Supplementary Table 5: Prokaryotic dating anchors.** Anchors used to date bacterial trees and the *de novo* tree. Each anchor was defined as the most recent common ancestor (MRCA) of one or more taxa. For BEAST-calibrated trees, all age intervals had a uniform prior. For PATHd8, no prior was specified. Note that some trees only contained a subset of these anchors. Also note that, because PATHd8<sup>26</sup> requires at least one anchor with fixed age, for PATHd8-dated trees the GOE anchor was fixed to an age of 2.55 Ga<sup>27</sup>.

ID	MRCA	range (Ga)	description
GOE	Oxyphotobacteria, Melainabacteria	2.238–2.63	Great Oxygenation Event (GOE) and mol. clock analysis <sup>27,28</sup>
Chl	Chloroplasts	1.047–4.4	Rhodophyte (red algae), Bangiomorpha <sup>29</sup>
Ri	Rickettsiales	1.6–4.4	ancestor of mitochondrion <sup>30</sup>
CB	Chlorobium, Bacteroidetes	1.64–4.4	Chlorobium-specific biomarkers <sup>31</sup>
Chr	Chromatiaceae	1.64–4.4	purple sulfur bacteria (gammaproteobacteria) biomarkers <sup>31</sup>
LUCA	Archaea, Bacteria	3.5–4.4	stromatolites < LUCA < detrital zircons <sup>32,33</sup>

**Supplementary Table 6:** Fraction of bacterial *de novo* 97%-OTUs (16S rRNA sequence clusters at 97% similarity) and metagenome-assembled-genome (MAG) 16S rRNA sequence sequences that were found to be represented in the SILVA v128 database (at 97% similarity). Columns 4 & 5: Number of full-length (FL) clusters and partial-length (V4) clusters, respectively, in SILVA (16S rRNA gene sequences clustered at 97% similarity). Columns 6 & 7: Global number of extant partial-length and full-length OTUs, respectively, estimated based on the fraction of *de novo* OTUs represented in SILVA (columns 4 & 5 divided by column 2). Column 8: Global number of extant full-length OTUs, estimated based on the fraction of MAGs represented in SILVA (column 4 divided by column 3).

taxon	fraction <i>de novo</i> OTUs represented	fraction MAGs represented	FL clusters in SILVA	V4 clusters in SILVA	extant FL based on <i>de novo</i>	extant V4 based on <i>de novo</i>	extant FL based on MAGs
Bacteria	0.496	0.382	201,586	84,397	406,423	170,155	527,712
Cyanobacteria	0.514	0.375	2,906	1,673	5,654	3,255	7,749

**Supplementary Table 7: Tree distances.** Pairwise topological distances between timetrees considered in the main article and previously published trees (Earth Microbiome Project deblurred 150 bp sequences<sup>10</sup> and SILVA 16S rRNA-based guide tree<sup>22</sup>), using the normalized Robinson-Foulds metric<sup>24</sup>. Also listed are the numbers of tips included in each comparison (i.e., common to both trees compared). See the Methods for details. For a visualization of tree distances see Supplementary Fig. 24.

tree 1	tree 2	distance	tips
Bacteria (16S <i>de novo</i> , FastTree, PATHd8)	EMP deblur	0.82	12,416
Bacteria (16S <i>de novo</i> , FastTree, PATHd8)	Bacteria (16S SILVA, FastTree, PATHd8)	0.82	24,169
Bacteria (16S <i>de novo</i> , FastTree, PATHd8)	SILVA guide	0.84	25,443
EMP deblur	Bacteria (16S SILVA, FastTree, PATHd8)	0.76	28,302
EMP deblur	SILVA guide	0.80	29,930
Bacteria (16S SILVA, FastTree, PATHd8)	SILVA guide	0.74	448,112
Cyan. (16S <i>de novo</i> , FastTree, PATHd8)	EMP deblur	0.73	153
Cyan. (16S <i>de novo</i> , FastTree, PATHd8)	Cyan. (16S SILVA, FastTree, BEAST+PATHd8)	0.71	300
Cyan. (16S <i>de novo</i> , FastTree, PATHd8)	Cyan. (16S SILVA, RAxML, BEAST+PATHd8)	0.71	300
Cyan. (16S <i>de novo</i> , FastTree, PATHd8)	SILVA guide	0.74	326
Cyan. (16S <i>de novo</i> , FastTree, PATHd8)	Cyan. (16S SILVA, BEAST)	0.44	26
EMP deblur	Cyan. (16S SILVA, FastTree, BEAST+PATHd8)	0.71	435
EMP deblur	Cyan. (16S SILVA, RAxML, BEAST+PATHd8)	0.70	435
EMP deblur	SILVA guide	0.74	473
EMP deblur	Cyan. (16S SILVA, BEAST)	0.40	36
Cyan. (16S SILVA, FastTree, BEAST+PATHd8)	Cyan. (16S SILVA, RAxML, BEAST+PATHd8)	0.51	6302
Cyan. (16S SILVA, FastTree, BEAST+PATHd8)	SILVA guide	0.68	6308
Cyan. (16S SILVA, FastTree, BEAST+PATHd8)	Cyan. (16S SILVA, BEAST)	0.68	483
Cyan. (16S SILVA, RAxML, BEAST+PATHd8)	SILVA guide	0.69	6302
Cyan. (16S SILVA, RAxML, BEAST+PATHd8)	Cyan. (16S SILVA, BEAST)	0.69	483
SILVA guide	Cyan. (16S SILVA, BEAST)	0.75	557

## S.7 Summary of supplementary files

Supplementary File 1: Accession numbers and summaries of amplicon sequencing data used to recover *de novo* OTUs.

Supplementary File 2: Accession numbers and summaries of sequencing data used from the Earth Microbiome Project<sup>10</sup>.

Supplementary File 3: R code used for analyzing diversification dynamics.

Supplementary File 4: Timetrees and undated phylogenetic trees constructed in this study, including required input files.

Supplementary File 5: Taxonomic classifications of *de novo* OTUs.

## References

- [1] Kendall, D. G. On some modes of population growth leading to ra fisher’s logarithmic series distribution. *Biometrika* **35**, 6–15 (1948).
- [2] Nee, S., May, R. M. & Harvey, P. H. The reconstructed evolutionary process. *Philosophical Transactions of the Royal Society of London B: Biological Sciences* **344**, 305–311 (1994).
- [3] Maddison, W. P., Midford, P. E., Otto, S. P. & Oakley, T. Estimating a binary character’s effect on speciation and extinction. *Systematic Biology* **56**, 701–710 (2007).
- [4] Sanderson, M. J. & Donoghue, M. J. Reconstructing shifts in diversification rates on phylogenetic trees. *Trends in Ecology & Evolution* **11**, 15–20 (1996).
- [5] Morlon, H. Phylogenetic approaches for studying diversification. *Ecology Letters* **17**, 508–525 (2014).
- [6] Sanmartín, I. & Meseguer, A. S. Extinction in phylogenetics and biogeography: From timetrees to patterns of biotic assemblage. *Frontiers in Genetics* **7**, 35 (2016).
- [7] Norris, J. *Markov Chains*. Cambridge Series in Statistical and Probabilistic Mathematics (Cambridge University Press, 1998).
- [8] Marshall, C. R. Five palaeobiological laws needed to understand the evolution of the living biota. *Nature Ecology & Evolution* **1**, 165 (2017).
- [9] Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life. *Nature Microbiology* (2017).
- [10] Thompson, L. R. *et al.* A communal catalogue reveals earth’s multiscale microbial diversity. *Nature* **551**, 457–463 (2017).
- [11] Schloss, P. D., Girard, R. A., Martin, T., Edwards, J. & Thrash, J. C. Status of the archaeal and bacterial census: an update. *mBio* **7** (2016).
- [12] Locey, K. J. & Lennon, J. T. Scaling laws predict global microbial diversity. *Proceedings of the National Academy of Sciences* **113**, 5970–5975 (2016).
- [13] Willis, A. Extrapolating abundance curves has no predictive power for estimating microbial biodiversity. *Proceedings of the National Academy of Sciences* **113**, E5096 (2016).
- [14] Gilbert, J. A., Jansson, J. K. & Knight, R. The Earth Microbiome project: successes and aspirations. *BMC Biology* **12**, 69 (2014).
- [15] Reeder, J. & Knight, R. The ‘rare biosphere’: a reality check. *Nat Meth* **6**, 636–637 (2009).
- [16] Edgar, R. C. Accuracy of microbial community diversity estimated by closed- and open-reference otus. *PeerJ* **5**, e3889 (2017).
- [17] Louca, S. & Doebeli, M. Efficient comparative phylogenetics on large trees. *Bioinformatics* **34**, 1053–1055 (2017).
- [18] Rabosky, D. L. Extinction rates should not be estimated from molecular phylogenies. *Evolution* **64**, 1816–1824 (2010).

- [19] Liow, L. H., Quental, T. B. & Marshall, C. R. When can decreasing diversification rates be detected with molecular phylogenies and the fossil record? *Systematic Biology* **59**, 646–659 (2010).
- [20] Beaulieu, J. M. & O’Meara, B. C. Extinction can be estimated from moderately sized molecular phylogenies. *Evolution* **69**, 1036–1043 (2015).
- [21] Aldous, D. J. Stochastic models and descriptive statistics for phylogenetic trees, from yule to today. *Statistical Science* **16**, 23–34 (2001).
- [22] Glöckner, F. O. *et al.* 25 years of serving the community with ribosomal rna gene reference databases and tools. *Journal of Biotechnology* **261** (2017).
- [23] Zanne, A. E. *et al.* Three keys to the radiation of angiosperms into freezing environments. *Nature* **505**, 89–99 (2014).
- [24] Day, W. H. E. Optimal algorithms for comparing trees with labeled leaves. *Journal of Classification* **2**, 7–28 (1985).
- [25] Jetz, W. *et al.* Global distribution and conservation of evolutionary distinctness in birds. *Current Biology* **24**, 919–930 (2014).
- [26] Britton, T., Anderson, C. L., Jacquet, D., Lundqvist, S. & Bremer, K. Estimating divergence times in large phylogenetic trees. *Systematic Biology* **56**, 741–752 (2007).
- [27] Shih, P. M., Hemp, J., Ward, L. M., Matzke, N. J. & Fischer, W. W. Crown group Oxyphotobacteria postdate the rise of oxygen. *Geobiology* **15**, 19–29 (2017).
- [28] Rasmussen, B., Fletcher, I. R., Brocks, J. J. & Kilburn, M. R. Reassessing the first appearance of eukaryotes and cyanobacteria. *Nature* **455**, 1101–1104 (2008).
- [29] Gibson, T. M. *et al.* Precise age of *Bangiomorpha pubescens* dates the origin of eukaryotic photosynthesis. *Geology* **46**, 135–183 (2017).
- [30] Parfrey, L. W., Lahr, D. J. G., Knoll, A. H. & Katz, L. A. Estimating the timing of early eukaryotic diversification with multigene molecular clocks. *Proceedings of the National Academy of Sciences* **108**, 13624–13629 (2011).
- [31] Brocks, J. J. *et al.* Biomarker evidence for green and purple sulphur bacteria in a stratified palaeoproterozoic sea. *Nature* **437**, 866–870 (2005).
- [32] Walter, M., Buick, R. & Dunlop, J. Stromatolites 3,400–3,500 Myr old from the North Pole area, Western Australia. *Nature* **284**, 443–445 (1980).
- [33] Wilde, S. A., Valley, J. W., Peck, W. H. & Graham, C. M. Evidence from detrital zircons for the existence of continental crust and oceans on the earth 4.4 gyr ago. *Nature* **409**, 175–178 (2001).