

1 **Supplementary Materials**

2
3 **Quantifying taxonomy-function associations across**
4 **hierarchical scales of bacterial nitrogen cycling**

5 Mingli Jiang^a, Yanni Huang^a, Zhiming Wu^e, Qian Zhu^a, Qian Li^a, Kaihua Pan^a,
6 Mingliang Zhang^a, Liang Shi^d, Jiguo Qiu^a, Pengfa Li^c, Xin Yan^a, Yiyong Zhu^b, Qing
7 Hong^{a*}

8 ^a Department of Microbiology, College of Life Sciences, Nanjing Agricultural
9 University, Key Laboratory of Agricultural and Environmental Microbiology,
10 Ministry of Agriculture and Rural Affairs, Nanjing 210095, China.

11 ^b College of Resources and Environment Sciences, Nanjing Agricultural University,
12 Nanjing 210095, China.

13 ^c College of Resources and Environment, Fujian Agriculture and Forestry University,
14 Fuzhou 350002, China.

15 ^d College of Life Sciences, Nanjing Agricultural University, Nanjing 210095, China.

16 ^e Nanchang Key Laboratory of Microbial Resources Exploitation & Utilization from
17 Poyang Lake Wetland, College of Life Sciences, Jiangxi Normal University,
18 Nanchang 330022, China

19 **Corresponding author:** Qing Hong

20 **E-mail address:** hongqing@njau.edu.cn

21
22 **This PDF file includes:**

23 Supplementary Extended Methods EM1 to EM 7 (Page 2-6)

24 Supplementary Notes 1-4 (Page 7-20)

25 Supplementary Figures 1 to 7 (Page 21-25)

26 Supplementary Note Figures 8 to 12 (Page 26-30)

27 Supplementary Note Tables 1, 3, and 4 (Page 31-33)

28 Supplementary References (Page 34-37)

29

30 **Supplementary Extended Methods**

31 **Supplementary EM1: Functional gene assignment criteria and exclusion**

32 **rationale**

33 Given our objective to quantify taxonomy-function associations, high-confidence
34 functional assignments are essential to minimize false positives that would obscure
35 genuine taxonomic patterns. Our dual-requirement approach (core genes plus
36 auxiliary genes for most pathways) ensures detection of complete functional pathways
37 rather than fragmentary gene presence.

38 Functional potential for six nitrogen transformation pathways was determined
39 using pathway-specific core and auxiliary genes (Supplementary Table S1). Core
40 genes included: nitrate reductase (*nasA*) and nitrite reductase (*nirB*) for assimilatory
41 nitrate reduction (ANRA); periplasmic nitrate reductase (*napA*) or membrane-bound
42 nitrate reductase (*narG*) for dissimilatory nitrate reduction to nitrite (DNN); nitrite
43 reductases (*nirS* or *nirK*) for denitrification (DNF), with downstream genes (*cnorB*,
44 *qnor*, *nosZ*) determining gaseous products (Graf et al., 2014); nitrogenase reductase
45 (*nifH*) for nitrogen fixation (NF); ammonia monooxygenase (*amoA*) and
46 hydroxylamine oxidoreductase (*hao*) for nitrification (NIT); cytochrome c nitrite
47 reductase (*nrfA*) for dissimilatory nitrite reduction to ammonium (DNRA).

48 Alternative gene variants were excluded based on combined considerations of
49 prevalence and annotation reliability. Nitrite oxidoreductase (*nxr*) subunits share
50 substantial sequence similarity with nitrate reductase subunits, making reliable
51 differentiation challenging in homology-based annotation pipelines. NCycDB
52 database documentation explicitly notes that "*nxrAB* and *narGH* were difficult to
53 distinguish from each other due to their close genetic relationships" (Tu et al., 2019).
54 Additionally, *nxr* genes exhibit extremely narrow phylogenetic distribution, being
55 restricted to nitrite-oxidizing bacterial lineages including Nitrospira, Nitrobacter, and
56 related taxa (Daims et al., 2016; Lücker et al., 2010), combined with severe
57 underrepresentation in current genomic databases. These annotation reliability
58 concerns motivated *nxr* exclusion to prevent false positive inflation in nitrate

59 reduction pathway detection. Ferredoxin-nitrite reductase (*nirA*) was detected in 4.4%
60 of genomes but predominantly co-occurred with the canonical ANRA pathway (*nasA*
61 + *nirB*); among genomes containing both *nirA* and *nirB* (0.5%), 95.7% also contained
62 *nasA*, indicating functional redundancy. Only 0.5% of genomes contained *nirA*
63 without *nasA*, representing *nirA*-exclusive assimilatory reduction primarily in
64 cyanobacteria with distinct evolutionary origin (Guerrero and Lara, 1987). Alternative
65 nitrogenases (*vnfH*, *anfH*) exhibited very low prevalence (*vnfH*: 0.003%; *anfH*:
66 0.36%) compared to canonical molybdenum-dependent nitrogen fixation (*nifH*:
67 8.7%), representing specialized metal-limited systems rather than the widespread Mo-
68 dependent pathway (Kuypers et al., 2018). These exclusions ensure that functional
69 assignments reflect well-characterized pathways with adequate representation and
70 minimal annotation ambiguity.

71 **Supplementary EM2: Sensitivity analysis for pathway annotation strategies**

72 To assess whether observed taxonomy-function associations reflect genuine
73 biological patterns or artifacts of pathway definition choices, we tested three
74 annotation strategies representing different stringency levels. All 73,472 genomes
75 were re-annotated under each strategy and genus-level Information Gain was
76 recalculated for all six nitrogen cycling pathways using identical computational
77 approaches. Spearman rank correlation quantified pattern consistency across
78 strategies.

79 Strict strategy employed annotation criteria detailed in Supplementary EM1,
80 requiring all pathway-specific core genes plus at least one auxiliary gene for most
81 pathways (Supplementary Table S1). Denitrification was assigned based solely on
82 nitrite reductase (*nirS* or *nirK*) presence, as this core enzyme is sufficient to indicate
83 denitrification capability, with downstream genes (*cnorB*, *qnor*, *nosZ*) determining
84 terminal products rather than pathway presence. Moderate strategy required only core
85 genes without auxiliary requirements: ANRA (*nasA* AND *nirB*), DNN (*napA* OR
86 *narG*), DNF (*nirS* OR *nirK*), DNRA (*nrfA*), NF (*nifH*), NIT (*amoA* AND *hao*).

87 Relaxed strategy assigned pathway presence if any gene listed in Supplementary
88 Table S1 for that pathway was detected.

89 Strict and Moderate strategies yielded highly correlated genus-level Information
90 Gain values (Spearman $\rho = 0.943$, $P = 0.005$) with identical pathway rankings
91 preserved. Relaxed strategy showed reduced correlation with Strict ($\rho = 0.657$, $P =$
92 0.156), attributable primarily to ANRA, whose positive genome proportion increased
93 from 15.2% under Strict criteria to 21.5% under Relaxed criteria—a larger expansion
94 than observed for any other pathway. The higher ANRA Information Gain under
95 Relaxed criteria therefore indicates that Strict criteria provide conservative rather than
96 inflated association estimates for ANRA. Pathway-specific results and correlation
97 analyses are provided in Supplementary Figure S4.

98 **Supplementary EM3: Environmental classification procedures and** 99 **methodological considerations**

100 Environmental sources were subdivided into 22 detailed subtypes based on
101 refined metadata extraction. Within the Aquatic category, genomes were classified
102 into Marine, Freshwater, Estuarine, and other aquatic subcategories. Within the
103 Terrestrial category, genomes were classified into Agricultural soils, Forest soils,
104 Grassland soils, and other terrestrial subtypes. The Other category comprised 15.1%
105 of genomes with ambiguous environmental descriptors, which were retained in global
106 analyses but excluded from environment-specific comparisons.

107 **Supplementary EM4: Environmental modulation analysis**

108 Environmental modulation was quantified using adaptation indices calculated as
109 Euclidean distances between environment-specific and global functional profiles for
110 classes present in multiple environments (≥ 10 genomes per environment). Function-
111 environment associations were assessed using Fisher's exact tests with Benjamini-
112 Hochberg correction (Benjamini and Hochberg, 1995).

113 **Supplementary EM5: Phylogenetic reconstruction, signal analysis, and** 114 **functional heterogeneity quantification**

115 Genus-level phylogenetic trees were constructed using 16 ribosomal proteins
116 (Parks et al., 2018). Representative genomes were selected from each genus based on
117 highest completeness and lowest contamination. Ribosomal protein sequences were
118 extracted using DIAMOND searches. Single-copy orthologs were identified and
119 aligned using MAFFT (Kato and Standley, 2013). Conserved regions were extracted
120 using TrimAl with gap threshold 0.2 and conservation threshold 50% (Capella-
121 Gutiérrez et al., 2009). Maximum-likelihood phylogenetic trees were constructed
122 using IQ-TREE with automated model selection and 1000 bootstrap replicates
123 (Nguyen et al., 2015).

124 Phylogenetic signal analysis employed Moran's I statistics. For binary ecological
125 strategy variables, phylogenetic signal was calculated using Abouheif's proximity
126 matrix method in adephylo (Jombart et al., 2010). For continuous variables, standard
127 Moran's I was computed using phylogenetic distance matrices with row-standardized
128 weights in ape (Paradis and Schliep, 2019). Statistical significance was assessed
129 through permutation testing (n=1000) with false discovery rate correction.

130 Environmental associations were analyzed using Fisher's exact tests. Functional
131 heterogeneity within genera was quantified using five metrics: Shannon Functional
132 Combination Diversity (adapted from Shannon entropy (Shannon, 1948)), Coefficient
133 of Variation of Functional Frequencies (CV_FFV), Functional Dominance,
134 Combination Richness, and Functional Evenness. Shannon FCD was calculated as H
135 $= -\sum(\pi_i \times \ln(\pi_i))$, where π_i represents the proportion of genomes exhibiting each
136 unique combination of six nitrogen cycling functions within a genus. Differences
137 among ecological strategy types were assessed using Kruskal-Wallis rank-sum tests
138 (Kruskal and Wallis, 1952) followed by pairwise Dunn's tests (Dunn, 1964) with
139 Holm correction (Holm, 1979). Effect sizes were quantified using epsilon-squared (ϵ^2)
140 (Tomczak and Tomczak, 2014).

141 **Supplementary EM6: Molecular evolutionary analysis methods and quality**
142 **control**

143 For nucleotide diversity (π) calculations, gap positions and ambiguous
144 nucleotides were excluded from calculations. Software versions used: MAFFT v7.475
145 (Kato and Standley, 2013), boot package v1.3-28 (Canty and Ripley, 2017).

146 For dN/dS analyses, nucleotide sequences were aligned using MUSCLE v3.8.31
147 (Edgar, 2004), and alignments were refined using TrimAl v1.4 (Capella-Gutiérrez et
148 al., 2009) with moderate filtering parameters (-gt 0.7 -st 0.01 -cons 60) to remove
149 poorly aligned regions while preserving codon integrity. Codon alignment integrity
150 was ensured by removing internal stop codons, trimming alignments to multiples of
151 three nucleotides, and excluding alignment columns containing stop codons.
152 Phylogenetic trees for each genus-level alignment were constructed using IQ-TREE
153 v2.0.3 (Minh et al., 2020) under the GTR+G nucleotide substitution model.

154 The Single-Likelihood Ancestor Counting (SLAC) method implemented in
155 HyPhy v2.5.29 (Kosakovsky Pond and Frost, 2005) employs maximum likelihood to
156 estimate synonymous (dS) and non-synonymous (dN) substitution rates across the
157 phylogeny. Median ω values with interquartile ranges were reported to summarize
158 selection pressures across genera, as dN/dS distributions are typically right-skewed
159 due to occasional lineages under relaxed or positive selection.

160 **Supplementary EM7: Data visualization tools and parameters**

161 Data visualization was performed using ggplot2 v3.4.0 (Gómez-Rubio, 2017) for
162 scatter plots, bar charts, and box plots, with color palettes selected from viridis v0.6.2
163 (Garnier et al., 2021) to ensure accessibility for color-vision deficiencies. Heatmaps
164 were generated using pheatmap v1.0.12 (Kolde, 2019), with hierarchical clustering
165 (Euclidean distance, complete linkage) applied where appropriate. Correlation
166 matrices were visualized using corrplot v0.92 (Wei et al., 2017). Phylogenetic trees
167 with metadata annotations were created using ggtree v3.6.2 (Yu et al., 2017). Final
168 figure assembly and annotation were performed in Adobe Illustrator 2024.

169

170 **Supplementary Notes**

171 **Supplementary Note 1: Extended functional distribution data and genome** 172 **source patterns analysis**

173 Among functionally active genomes, single-pathway specialization dominated,
174 with the five most abundant single-function strategies closely reflecting the overall
175 functional hierarchy: DNN-only (6,646 genomes), ANRA-only (5,104), NF-only
176 (3,470), DNF-only (2,571), and DNRA-only (2,495), collectively representing 60.4%
177 of functionally active genomes (Figure S8). Multi-functional genomes exhibited
178 exponential frequency decline (dual-function: 13.4%, triple-function: 4.1%), with
179 higher-order combinations showing extreme rarity: quadruple-function genomes
180 comprised only 431 genomes (0.6%), quintuple-function genomes 27 genomes
181 (0.04%), and complete six-function genomes were virtually absent (Figure 1B). This
182 distribution follows an exponential decay pattern that supports metabolic cost-benefit
183 models where each additional pathway imposes increasing regulatory burdens (Lynch,
184 2007). Among dual-function combinations, the top five were DNF+DNN (2,920
185 genomes), ANRA+DNN (2,579), DNN+DNRA (1,853), ANRA+DNF (596), and
186 DNN+NF (490). Triple-function combinations were dominated by
187 ANRA+DNF+DNN (1,193 genomes, 39.4% of all triple-function genomes),
188 indicating modular integration of assimilatory and respiratory nitrogen cycling
189 strategies (Graf et al., 2014).

190 The unexpected ANRA-DNN pairing exhibits dual-level dominance in bacterial
191 nitrogen cycling: both pathways show high individual prevalence and strong
192 combinatorial association, establishing them as complementary metabolic strategies.
193 Organisms possessing both ANRA and DNN can flexibly allocate nitrate between
194 biosynthetic assimilation (supporting growth) and respiratory dissimilation
195 (supporting energy generation), potentially providing competitive advantages in
196 environments experiencing variable nitrogen availability and oxygen tension. In
197 contrast, DNN-DNRA represents a committed respiratory strategy that channels nitrite
198 exclusively toward ammonium production rather than allowing biosynthetic

199 flexibility. Recent studies have demonstrated that in facultative anaerobes,
200 assimilatory and dissimilatory nitrate reduction pathways can coexist with
201 coordinated expression even under aerobic conditions, allowing flexible allocation of
202 nitrate between biosynthetic and respiratory functions depending on nitrogen
203 availability (Ahn et al., 2025).

204 Major environmental categories showed distinct quantitative associations with
205 nitrogen cycling functions (Figure 1D, S8A). Aquatic isolates contributed 32.1% of
206 DNF genomes and 29.2% of DNN genomes despite comprising only 27.0% of total
207 genomes, consistent with the enrichment of anaerobic respiratory processes in
208 oxygen-depleted aquatic environments (Seitzinger et al., 2006). Terrestrial isolates
209 exhibited pronounced enrichment for assimilatory processes, contributing 27.9% of
210 ANRA genomes while representing just 14.1% of the dataset, consistent with nitrogen
211 limitation in soil ecosystems where assimilatory nitrogen uptake is essential for
212 microbial growth (Schimel and Bennett, 2004). Host-associated sources displayed a
213 distinctive functional signature: despite having the highest proportion of genomes
214 lacking detectable nitrogen cycling functions (72.4%), they exhibited the strongest
215 enrichment for nitrogen fixation among functionally active genomes, with NF
216 representing 6.4% of host-associated genomes—the highest proportion across all
217 environmental categories.

218 Fine-scale analysis across environmental subtypes revealed additional functional
219 specialization patterns (Figure S9B). Among aquatic sources, marine environments
220 showed stronger denitrification enrichment (DNF: 34.2% of genomes vs 27.0%
221 expected) than freshwater environments (DNF: 29.8% vs 27.0% expected), potentially
222 reflecting differences in organic carbon availability and oxygen depletion patterns.
223 Within terrestrial environments, forest soils showed the highest ANRA enrichment
224 (31.4% of ANRA genomes while representing 8.2% of the dataset) compared to
225 agricultural soils (24.1% vs 6.3% expected), consistent with different nitrogen
226 availability regimes in natural versus managed systems (Templer et al., 2012). Plant-
227 associated genomes within terrestrial sources demonstrated intermediate nitrogen

228 fixation patterns with rhizosphere isolates (7.8% NF) exceeding phyllosphere isolates
229 (3.1%). Host-associated environments revealed heterogeneous functional patterns
230 across host types: ruminant-associated genomes showed higher nitrogen fixation
231 frequencies (8.9% of genomes) than monogastric-associated genomes (4.2%). These
232 associations are interpreted as descriptive characterizations of functional enrichment
233 patterns across genome source categories. Observed patterns are consistent with
234 established biogeochemical principles, providing contextual support for the genome
235 annotation quality and environmental classification scheme used in this study.
236 Statistical analysis across 22 environmental subtypes with ≥ 100 genomes each
237 revealed 82 significant function-environment associations among 132 tested
238 combinations (62.1%, $q < 0.05$, Benjamini-Hochberg FDR correction); these patterns
239 document the co-distribution of bacterial functional potential and environmental
240 context within the current database, without implying causal environmental control of
241 functional organization.

242 **Supplementary Note 2: Extended class-level functional archetype analysis**

243 Bootstrap resampling ($n = 1,000$) yielded mean ARI = 0.553 ± 0.265 ; the wide
244 95% CI (0.003–1.000) is attributable to stochastic variation across numerically
245 unequal clusters ($n = 25, 34, 17$), with the dominant cluster structure remaining
246 reproducible, as indicated by 60.5% of resamples exceeding ARI = 0.5. Kruskal–
247 Wallis tests across the three multi-class archetypes confirmed significant inter-
248 archetype differences for all functional metrics (Supplementary Table S4), with η^2
249 ranging from 0.242 (ANRA) to 0.789 (participation rate). Cluster number validation
250 and individual class silhouette values are shown in Figure S3 and **Figure S11**,
251 respectively; archetype mean silhouettes were: Functionally Inactive, 0.625;
252 Functionally Moderate, 0.371; N-Retention Dominant, 0.155; Nitrification Specialist,
253 0 (by convention for single-member clusters in the R cluster package).

254 Among Functionally Inactive classes, Bacilli_A (948 genomes, PR = 0.040)
255 showed only residual DNRA (3.4%), while Coriobacteriia (945 genomes, PR = 0.099)
256 retained limited DNN (5.8%) and DNRA (4.9%), consistent with its host-associated

257 niche. Acidimicrobiia (203 genomes, PR = 0.074) showed modest DNN (6.4%). At the
258 archetype boundary, Fusobacteriia (92 genomes, PR = 0.141) and Rhodothermia (50
259 genomes, PR = 0.140) approach but do not cross the Jenks threshold. Thirteen classes
260 showed PR = 0 (including Vampiromicrobiia, Chlamydiia, Brachyspirae), indicating
261 complete genomic absence of nitrogen cycling genes.

262 Within Functionally Moderate classes, Gammaproteobacteria (9,466 genomes,
263 PR = 0.710) showed balanced representation across DNN (48.6%), ANRA (31.0%),
264 DNF (26.2%), and DNRA (10.1%). Actinomycetes (6,752 genomes, PR = 0.600)
265 emphasised DNN (40.6%) and ANRA (33.2%), reflecting aerobic soil metabolism.
266 Bacteroidia (6,816 genomes, PR = 0.291) maintained 70.9% functional absence while
267 active genomes showed DNRA (13.0%) and DNF (9.3%). Alphaproteobacteria (6,094
268 genomes, PR = 0.550) exhibited distributed representation across DNN (29.9%),
269 ANRA (24.7%), and DNF (22.7%). Clostridia (5,690 genomes, PR = 0.209) showed
270 NF (18.2%) as its primary pathway despite low overall participation, while
271 Cyanobacteriia (523 genomes, PR = 0.317) showed NF (28.1%) as its dominant
272 function.

273 N-Retention Dominant classes are enriched for obligately or facultatively
274 anaerobic lineages. Desulfomicrobiia (175 genomes, PR = 0.886, NF = 54.3%,
275 DNRA = 66.3%, loss = 0), Desulfuromonadia (81 genomes, PR = 0.963, NF = 88.9%,
276 DNRA = 92.6%), and Desulfobulbia (75 genomes, PR = 0.987, NF = 72.0%,
277 DNRA = 72.0%) confirm NF + DNRA co-dominance as a shared constraint across
278 sulfate-reducing Desulfobacterota. Chlorobiia (86 genomes, PR = 0.977) showed near-
279 universal NF (97.7%) as its sole function, while Brocadiiia (22 genomes, PR = 0.955)
280 combined DNRA (54.5%), DNF (54.5%), and DNN (86.4%), consistent with the
281 anammox pathway (Strous et al., 1999). Syntrophia (27 genomes, PR = 0.556,
282 DNRA = 44.4%) represents the lower PR boundary of this archetype. NF and DNRA
283 frequencies were significantly higher in N-Retention Dominant than in Functionally
284 Moderate classes (Mann–Whitney U: $p = 0.002$ and $p < 0.001$, respectively), while
285 nitrogen loss was significantly lower ($p = 0.012$). The lower mean silhouette (0.155;

286 Figure S3) reflects a PR gradient within this archetype rather than misclassification,
287 as all 17 classes share the NF + DNRA dominance and low nitrogen loss profile that
288 define this archetype.

289 Nitrospira's silhouette of 0 (Figure S3; set to 0 by convention for single-member
290 clusters in the R cluster package) reflects its functional isolation: no comparable
291 functional profile exists among the remaining 76 classes across $k = 2-8$.

292 Genome source association analysis provides quantitative validation of the
293 intrinsic stability of taxonomy-based functional organization, examining whether
294 class-level functional archetypes remain consistent across genomes sampled from
295 different environmental categories. This analysis operates along two complementary
296 dimensions: strategy-level stability quantified by source association indices (Figure
297 S5), and functional frequency stability quantified by inter-class ranking consistency
298 across environments (Figure S10).

299 Source association indices ranged 0.003–0.813 (mean 0.234 ± 0.158). Major
300 generalist lineages exhibited the lowest indices (Gammaproteobacteria = 0.066,
301 Actinomycetes = 0.063, Negativicutes = 0.019), indicating that their functional
302 strategy profiles are intrinsic genomic properties insensitive to sampling provenance.
303 Classes with higher indices (e.g., Terriglobia = 0.313, Clostridia = 0.267) tend to be
304 functionally restricted specialists with low participation rates, for which sparse or
305 environmentally skewed sampling produces greater apparent strategy deviation. The
306 negative correlation between source association indices and class size ($\rho = -0.334$, p
307 < 0.001) confirms that this variation is itself taxonomically structured rather than
308 reflecting genuine environmental reprogramming of functional repertoires.

309 At the functional frequency level, inter-class rankings across six environmental
310 categories were highly conserved: Spearman correlations of class-level rankings
311 between environment pairs averaged $\rho = 0.978$ (range 0.937–1.000) for DNN, $\rho =$
312 0.947 (range 0.883–1.000) for ANRA, and $\rho = 0.925$ (range 0.709–1.000) for DNF.

313 Within individual classes, absolute frequencies showed limited variation:

314 Gammaproteobacteria DNN frequencies ranged 44–58% across environments ($CV =$

315 8.5%), representing stable functional identity, while higher CV values observed for
316 lower-prevalence functions (e.g., DNRA CV = 53.6%) reflect stochastic sampling
317 effects on rare capabilities rather than environmentally driven reprogramming. The 23
318 of 57 classes (40.4%) showing $CV > 0.3$ in participation rates (Figure S5) are
319 disproportionately low-prevalence specialists, consistent with this interpretation.
320 Taken together, these two dimensions of analysis confirm that class-level functional
321 archetypes represent stable intrinsic properties of bacterial lineages. Within-class
322 frequency variation across environments is real but occurs within a highly conserved
323 inter-class ranking structure, supporting the use of taxonomic identity as a reliable
324 organizational framework for nitrogen cycling functional inference across diverse
325 ecological contexts.

326 **Supplementary Note 3: Extended genus-level ecological strategy analysis**

327 Specialist strategies maintain distinct ecological niches within structured
328 communities, while generalist strategies provide competitive advantages through
329 metabolic flexibility across variable nitrogen regimes. The moderate significance rate
330 (53.3%) observed across genome source-strategy combinations reflects the
331 complexity of these relationships, indicating that while systematic associations exist,
332 they are modulated by multiple ecological and evolutionary factors. Detailed
333 mechanistic analysis of environment-strategy associations reveals that aquatic
334 environments present variable redox conditions characteristic of fluctuating oxygen
335 availability, where metabolic flexibility in nitrogen transformations may provide
336 advantages (Mosley et al., 2022). In host-associated environments, genome
337 streamlining principles operate where microbes reduce investments in metabolically
338 expensive pathways when substrates are readily available through host provisioning
339 (Morris et al., 2012).

340 Moran's I analysis revealed the following detailed values: Multifunctional ($I =$
341 0.4394 , $p = 0.001$), Non-N-Cycling ($I = 0.4248$, $p = 0.001$), N-Retention ($I = 0.3217$,
342 $p = 0.001$), Intermediate ($I = 0.3075$, $p = 0.001$), and N-Loss ($I = 0.2067$, $p = 0.001$).
343 For comparison, genome source distribution patterns across five major categories

344 exhibited Moran's I values of 0.032, while nitrogen cycling function frequencies
345 displayed mean signals of 0.033, yielding an overall mean of 0.338 for ecological
346 strategies. The evolutionary basis for this disparity likely involves horizontal gene
347 transfer of small gene cassettes enabling individual gene mobility with relative ease
348 (Arnold et al., 2022), while complete multi-gene functional modules require
349 coordinated inheritance. This pattern aligns with observations that complex traits
350 encoded by multiple genes tend to be more phylogenetically conserved than simpler
351 traits (Isobe et al., 2019; Martiny et al., 2013). Phylogenetic clustering analysis
352 revealed the following distributions: Non-N-Cycling genera within Clostridia and
353 Bacilli (50.7% combined) and Bacteroidia (12.1%); Multifunctional genera within
354 Gammaproteobacteria (35.9%) and Alphaproteobacteria (27.4%), with additional
355 representation in Bacteroidia (12.9%).

356 Functional profile analysis across the five ecological strategy types revealed
357 characteristic pathway combinations within each strategy. N-Retention genera (n =
358 468, mean ANRA 15.3%, NF 18.9%, DNRA 16.0%) ranged from single-pathway
359 specialists to dual-retention lineages: *Geomonas* (23 genomes, NF 100%, DNRA
360 100%) and *Centipeda* (16 genomes, NF 100%, DNRA 100%) maintained maximum
361 frequencies across both retention pathways simultaneously, while *Campylobacter*
362 (251 genomes, DNRA 55.8%) and *Chlorobium* (68 genomes, NF 97.1%) exemplified
363 single-pathway dominance at high frequency. N-Loss genera (n = 77, mean DNF
364 40.9%, NIT 1.7%) were dominated by denitrification, with *Nitrosomonas* (84
365 genomes, DNF 92.9%, NIT 51.2%) and *Nitrospira* (24 genomes, DNF 87.5%, NIT
366 79.2%) carrying dual nitrification-denitrification capacity at high frequency.
367 Intermediate genera (n = 64, mean DNN 35.4%) showed wide internal variation
368 despite uniform absence of committed downstream pathways, with *Staphylococcus*
369 (181 genomes, DNN 84.5%) and *Proteus* (23 genomes, DNN 95.7%) at the high end
370 of DNN frequency. Multifunctional genera (n = 387, mean ANRA 30.7%, DNN
371 41.9%, DNF 34.9%, NF 9.5%, DNRA 9.1%) exhibited the broadest internal
372 heterogeneity, with 51 genera simultaneously exceeding 50% retention and 50% loss

373 scores. Among these, *Bradyrhizobium* (268 genomes, NF 63.4%, DNF 56.7%),
374 *Azospirillum* (61 genomes, NF 95.1%, DNF 72.1%), and *Methylobacter* (50 genomes,
375 ANRA 66.0%, NF 66.0%, DNF 80.0%) represented lineages with substantial
376 bidirectional nitrogen transformation capacity. As noted in the main text,
377 *Methylomonas* (31 genomes, NF 90.3%, DNF 93.5%) exemplifies functionally
378 paradoxical combinations in which nitrogen fixation and denitrification co-occur at
379 high frequency. Additional genera exhibiting this pattern included *Basfia* (13
380 genomes, DNRA 100%, DNF 100%), *Aromatoleum* (18 genomes, ANRA 66.7%, NF
381 55.6%, DNF 100%), and *Thauera* (25 genomes, ANRA 56.0%, DNF 100%),
382 collectively demonstrating that simultaneous high-frequency retention and loss
383 capacities are not restricted to a single lineage but distributed across phylogenetically
384 distinct genera within the Multifunctional strategy type.

385 These genus-level functional profiles provide reference baselines for community-
386 scale nitrogen cycling inference and targeted cultivation. For 16S rRNA gene
387 amplicon studies, community-level functional potential is estimated through
388 abundance-weighted summation across all detected genera using values from
389 Supplementary Table S2. For example, if 16S sequencing reveals *Pseudomonas* at
390 15% relative abundance, querying Supplementary Table S2 shows this genus (based
391 on 1,020 representative genomes) exhibits DNN at 39.1% frequency, ANRA at
392 35.9%, and DNF at 28.9%, yielding genus-specific contributions to community
393 functional potential of 5.9% for DNN, 5.4% for ANRA, and 4.3% for DNF. Summing
394 such contributions across all detected genera provides comprehensive community-
395 level functional estimates for each nitrogen transformation pathway. For targeted
396 cultivation, genera combining substantial genome representation with elevated
397 functional frequencies serve as efficient screening targets: *Neisseria* (105 genomes,
398 DNF 92.4%) and *Nitrosomonas* for denitrification isolation; *Campylobacter* and
399 *Citrobacter* (41 genomes, DNRA 97.6%) for DNRA-encoding bacteria; *Nostoc* (68
400 genomes, NF 100%) and *Chlorobium* for nitrogen fixation. These genomic potential
401 estimates require validation through activity-based measurements including

402 quantitative PCR of functional genes, metatranscriptomics, or biogeochemical rate
403 determinations, as environmental modulation, horizontal gene transfer, and post-
404 transcriptional regulation introduce uncertainty in translating gene presence to
405 metabolic activity.

406 Functional heterogeneity analysis across ecological strategy types revealed
407 systematic differences in intra-genus functional diversity patterns (Figure S12A).
408 Kruskal-Wallis tests demonstrated highly significant variation across all five
409 heterogeneity indices: Shannon Functional Combination Diversity ($\chi^2 = 913.96$, $p <$
410 0.001), Coefficient of Variation of Functional Frequencies ($\chi^2 = 971.11$, $p <$
411 0.001), Functional Dominance ($\chi^2 = 828.51$, $p <$ 0.001), Combination Richness ($\chi^2 = 977.75$,
412 $p <$ 0.001), and Functional Evenness ($\chi^2 = 677.22$, $p <$ 0.001). Among nitrogen-
413 cycling genera, Multifunctional genera exhibited the highest functional heterogeneity
414 levels (mean Shannon FCD = 1.380), which likely reflects the complexity of
415 coordinating both nitrogen retention and loss capabilities within genera. Retention
416 genera showed intermediate heterogeneity (mean Shannon FCD = 0.732), while
417 Intermediate and N-Loss genera displayed more constrained functional variation
418 patterns. These differences indicate that coordination of multiple nitrogen
419 transformation pathways within genera creates greater functional complexity than
420 specialized single-pathway strategies (Shade and Gilbert, 2015).

421 Genome source distribution breadth analysis demonstrated a positive relationship
422 between genus-level genome source distribution diversity and functional
423 heterogeneity across ecological strategy types (Figure S12B). Integration of genome
424 source plasticity and functional specialization patterns across qualified genera
425 revealed four distinct categories through standardized axis analysis (Figure S13).
426 Quadrant 2 (Phylogenetic Specialists) contained the largest proportion of genera (478
427 genera, 37.3%), characterized by low genome source plasticity and high functional
428 specialization, dominated by Non-N-Cycling genera. Quadrant 3 (Phylogenetic
429 Generalists) included 413 genera (32.2%) with low genome source plasticity and low
430 functional specialization, primarily comprising Multifunctional and Retention genera.

431 Quadrant 4 (Genome Source Generalists) encompassed 227 genera (17.7%) exhibiting
432 high genome source plasticity and low functional specialization, with Multifunctional
433 genera representing the majority. Quadrant 1 (Genome Source Specialists) contained
434 163 genera (12.7%) displaying high genome source plasticity and high functional
435 specialization.

436 **Supplementary Note 4: Extended molecular evolutionary analysis**

437 Class-level analysis of *nrfA* revealed systematic phylogenetic structuring in
438 contrast to its environmental stability: Alphaproteobacteria exhibited the lowest
439 diversity ($\pi = 0.390$), while Clostridia and Coriobacteriia displayed elevated values ($\pi =$
440 $= 0.447-0.463$), and Gammaproteobacteria showed intermediate-to-high values ($\pi =$
441 0.453) (Figure 5C). The combination of environmental stability and phylogenetic
442 structuring indicates that the bimodal dN/dS distribution reflects lineage-specific
443 rather than environment-driven evolutionary dynamics. The nine genera exhibiting
444 *nrfA* dN/dS ≥ 1.0 span seven phylogenetically distinct bacterial classes:

445 Gammaproteobacteria (*Shewanella*, *Actinobacillus*), Actinomycetes (*Arachnia*,
446 *Actinomyces*), Myxococcia (*Myxococcus*), Bacteroidia (*Parabacteroides*),
447 Blastocatellia (*OLB17*), Bacilli_A (*Bulleidia*), and Coriobacteriia (*Adlercreutzia*),
448 with ω values ranging from 1.591 to 4.190. Despite exhibiting relaxed or positive
449 selection at the sequence level, these nine genera maintain high DNRA functional
450 frequencies (mean 70.1%, range 40.3-95.8% across genera), with elevated *nrfA* gene
451 presence rates observed in sequenced genomes (e.g., *Shewanella* 94%, *Actinobacillus*
452 96%, *Myxococcus* 82%). This functional retention distinguishes the high- ω pattern
453 from simple pathway degradation, where declining functional prevalence would be
454 expected to accompany relaxed selection. Given that DNRA is deployed across
455 diverse regulatory contexts (Durand and Guillier, 2021) and electron donor regimes
456 (Liu et al., 2021; Yoon et al., 2015) in different bacterial lineages, the
457 phylogenetically structured selection patterns likely reflect lineage-specific
458 optimization of *nrfA* for integration within distinct metabolic networks, where

459 regulatory and energetic contexts of each bacterial clade necessitate particular
460 sequence configurations while preserving the fundamental nitrite reduction chemistry.

461 Phylogenetic analysis of *napA* conservation across 20 bacterial classes revealed
462 pronounced lineage-specific variation. Deinococci exhibited exceptionally stringent
463 purifying selection ($\pi = 0.116$)—the lowest value observed for any gene-class
464 combination—while Campylobacteria ($\pi = 0.313$) and Gammaproteobacteria ($\pi =$
465 0.331) showed substantially elevated diversity, representing a nearly 3-fold range in
466 sequence conservation (Figure 5C). Environmental deployment contexts revealed
467 additional variation: aquatic ($\pi = 0.399$) and engineered system ($\pi = 0.386$)
468 environments occupied intermediate positions between the terrestrial and extreme
469 environment extremes (Figure 5D). Unlike *narG*'s function during strict anaerobic
470 respiration, *napA* enables periplasmic nitrate scavenging under microoxic-to-oxic
471 conditions (Potter et al., 2001; Sparacino-Watkins et al., 2014), where maintenance of
472 electron transfer efficiency across fluctuating redox states likely imposes stringent
473 structural requirements. The correlation between *napA*'s high genus-level
474 conservation and its environmental responsiveness demonstrates that ecological
475 deployment contexts can impose selection pressures that modulate sequence evolution
476 independently of pathway-level taxonomy-function association strength.

477 Detailed structural analysis of denitrification genes revealed mechanistic bases
478 for within-pathway heterogeneity. The exceptional conservation of *cnorB* reflects
479 structural constraints imposed by the binuclear heme b_3 -FeB active site, where three
480 conserved histidine residues and a glutamate residue coordinate the metal centers
481 required for NO reduction (Murali et al., 2024; Schurig-Briccio et al., 2013). The
482 evolutionarily unrelated nitrite reductases *nirS* and *nirK* showed significant
483 differences reflecting their divergent structural architectures—cytochrome cd_1 (*nirS*)
484 versus copper-containing (*nirK*) metalloforms. The divergent environmental
485 sensitivities of *nirS* and *nirK* reflect their structural architectures and functional
486 autonomy differences: *nirK* can function independently while *nirS* requires additional
487 gene products for nitrite reduction (Ming et al., 2024), potentially accounting for

488 *nirS*'s greater responsiveness to environmental variation. These patterns reveal that
489 taxonomy-function association strength, while shaped by pathway-level associations,
490 is further modulated by gene-specific structural and catalytic properties that contribute
491 comparably to evolutionary rate variation.

492 Detailed analysis of nucleotide diversity patterns across 13 nitrogen cycling genes
493 revealed substantial variation in evolutionary constraints beyond those discussed
494 extensively in the main text (Figure 5). Genes not analyzed in detail in the main text
495 include nitrous oxide reductase *nosZ* ($\pi = 0.20 \pm 0.073$), ammonia monooxygenase
496 *amoA* ($\pi = 0.20 \pm 0.071$), quinol-dependent nitric oxide reductase *qnor* ($\pi = 0.24 \pm$
497 0.089), and hydroxylamine oxidoreductase *hao* ($\pi = 0.27 \pm 0.13$). Notably, *hao*
498 displayed the highest standard deviation among all genes analyzed, suggesting
499 substantial variation in evolutionary constraints across nitrifying lineages.

500 Comparative analysis of functional gene groups revealed that assimilatory genes (*nas*,
501 *nirB*; mean $\pi = 0.285$) exhibited significantly higher diversity than dissimilatory
502 respiratory genes (*napA*, *narG*, *nrfA*; mean $\pi = 0.207$), indicating stronger purifying
503 selection on respiratory pathway components. Within denitrification, evolutionary
504 rates varied systematically, with *cnorB* and *nosZ* showing the highest conservation (π
505 $= 0.18-0.20$), followed by *nirS* ($\pi = 0.22$), *qnor* ($\pi = 0.24$), and *nirK* ($\pi = 0.25$),
506 reflecting structural complexity hierarchies from binuclear metal centers to copper-
507 dependent catalytic sites.

508 Class-level molecular evolution analysis across 20 major bacterial classes
509 demonstrated systematic phylogenetic structuring of nitrogen cycling gene evolution
510 (Figure 5C). Beyond the range extremes noted in the main text (Cyanobacteriia mean
511 $\pi = 0.29$, Clostridia mean $\pi = 0.45$), intermediate conservation patterns were observed
512 across multiple classes. Deinococci showed elevated molecular conservation across
513 analyzed genes, with particularly low diversity in *napA* ($\pi = 0.116$), while
514 Negativicutes displayed consistently elevated conservation across multiple nitrogen
515 cycling functions. Classes exhibiting more relaxed molecular evolution included
516 Spirochaetia and Bacilli_A, with mean π values approaching those observed in

517 Clostridia. This phylogenetic variation demonstrates that evolutionary rate variation
518 operates at multiple hierarchical levels, with class-level phylogenetic background
519 contributing substantially to molecular evolutionary rates independently of gene
520 function. Hierarchical clustering revealed organizational structures reflecting both
521 functional modularity and phylogenetic relationships, with gene clustering separating
522 nitrogen cycling functions into distinct modules while class clustering grouped
523 phylogenetically related lineages regardless of their functional capabilities.

524 Environmental sensitivity analysis revealed substantial gene-specific variation in
525 environmental responsiveness across all 13 nitrogen cycling genes, with
526 environmental sensitivity quantified as the coefficient of variation in π values across
527 six environmental categories (Figure 5D-E). Beyond genes discussed extensively in
528 the main text, comprehensive analysis confirmed that environmental sensitivity does
529 not correspond to pathway membership or functional categories. Within
530 denitrification, genes exhibited divergent environmental sensitivities despite
531 belonging to the same functional pathway, and within dissimilatory nitrate reduction
532 pathways, *nrfA* displayed exceptional environmental stability (CV = 0.023) while
533 *napA* showed high environmental responsiveness (CV = 0.11), demonstrating that
534 genes within the same metabolic context can respond independently to ecological
535 gradients. Detailed examination of environment-specific π patterns (Figure 5D)
536 revealed that extreme environments consistently induced elevated diversity across
537 multiple genes, while terrestrial environments showed reduced diversity particularly
538 for nitrite reductases and periplasmic reductases, confirming that environmental
539 contexts impose gene-specific rather than pathway-general selective pressures.

540 Genus-level diversity analysis across 891 genera representing four ecological
541 strategy types revealed systematic relationships between functional organization and
542 molecular evolution (Figure S14). Intermediate genera exhibited the lowest nucleotide
543 diversity (mean $\pi = 0.19 \pm 0.085$), followed by N-Loss (mean $\pi = 0.22 \pm 0.076$),
544 Multifunctional (mean $\pi = 0.23 \pm 0.061$), and Retention (mean $\pi = 0.26 \pm 0.088$).
545 Statistical analysis confirmed significant differences among ecological types

546 (Kruskal-Wallis $\chi^2 = 54.89$, $df = 3$, $p < 0.001$). These patterns suggest systematic
547 differences between functional strategy and molecular evolutionary constraints,
548 though substantial variation within strategy types indicates complex interplay between
549 functional repertoire, taxonomic position, and environmental adaptation.
550

551 **Supplementary Figure**

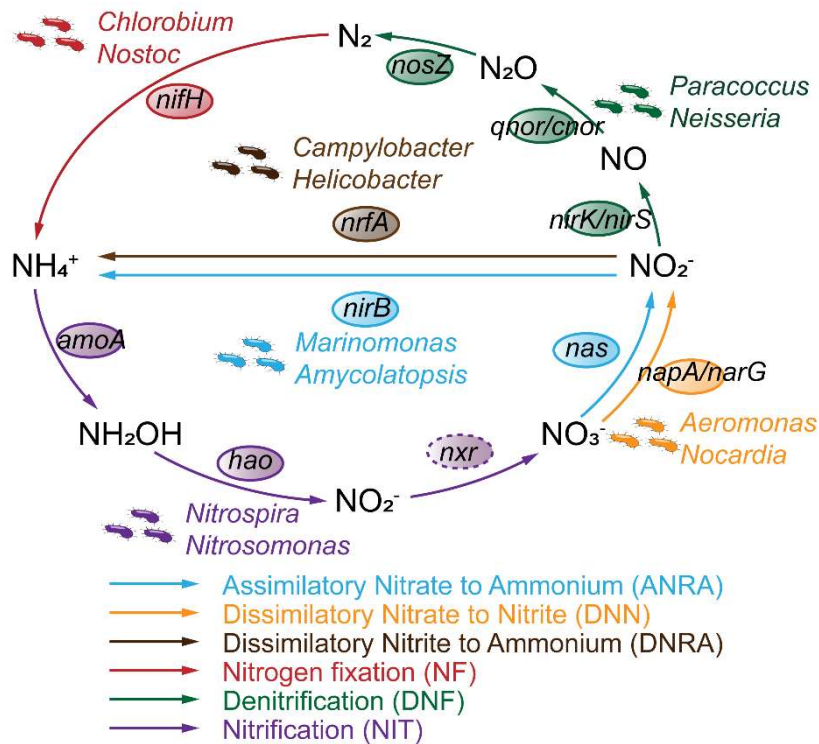


Figure S1. Conceptual framework of bacterial nitrogen cycling pathways and gene detection strategy. Six major nitrogen transformation processes: Assimilatory Nitrate to Ammonium reduction (ANRA, blue), Dissimilatory Nitrate to Nitrite reduction (DNN, orange), Dissimilatory Nitrite to Ammonium reduction (DNRA, brown), Nitrogen fixation (NF, red), Denitrification (DNF, green), and Nitrification (NIT, purple). Representative key genes are shown in gray ovals with solid lines; the *nxr* gene (dashed oval) was not included in the analysis due to annotation database limitations. The complete list of core and auxiliary genes for each pathway is provided in Supplementary Table S1. Functional capability was assigned using a dual-detection strategy requiring all core genes plus at least one auxiliary gene per pathway. Representative bacterial genera known to carry these pathways are labeled for each transformation process.

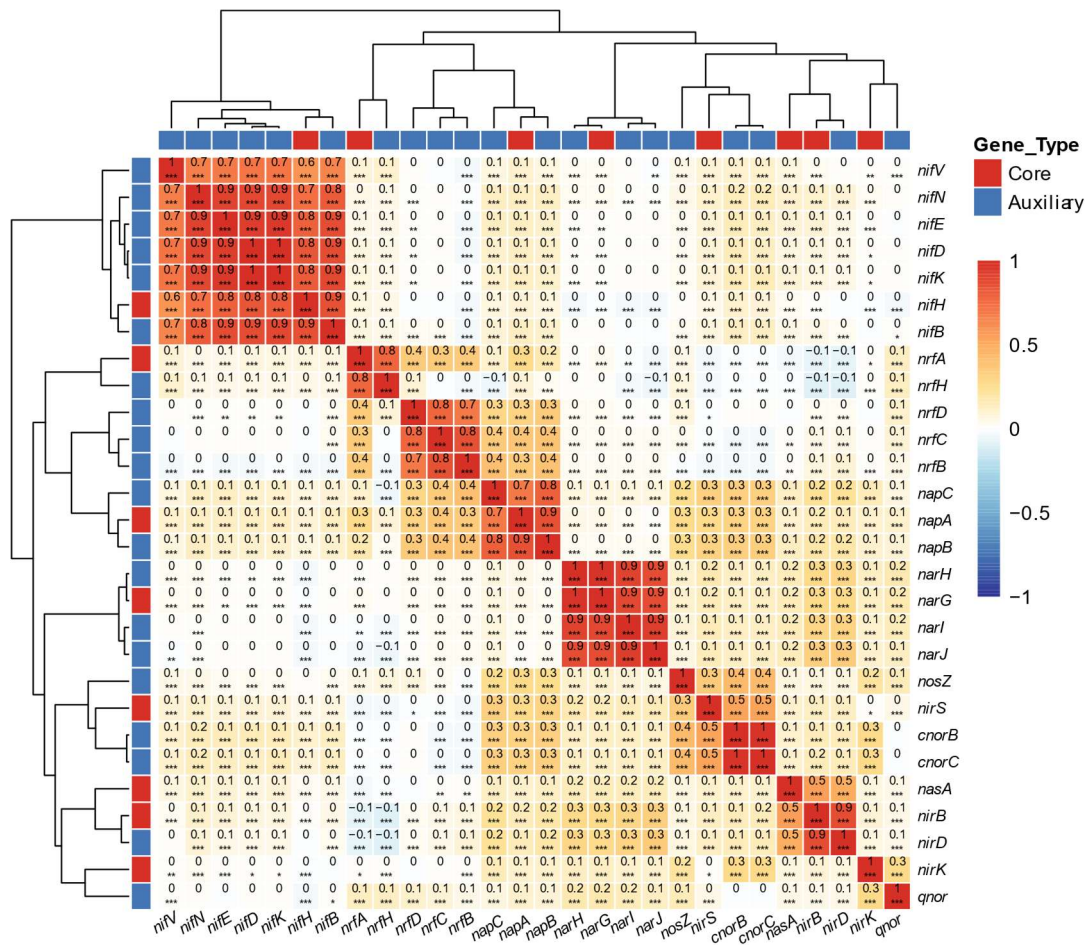
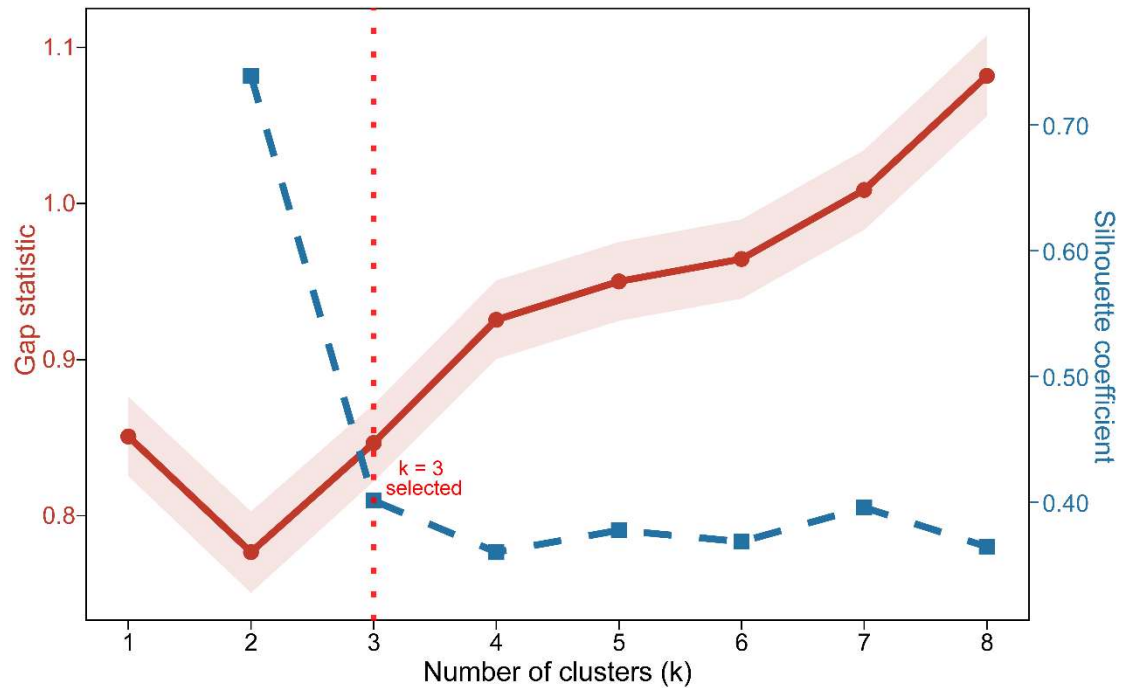


Figure S2. Core-auxiliary gene correlation analysis for nitrogen cycling pathways. Spearman correlation matrix of nitrogen cycling genes with frequency $\geq 5\%$ across 73,472 GTDB representative genomes. Gene categories: Core (red), Auxiliary (blue), Regulatory (light blue), Other (gray). Correlation coefficients shown with statistical significance: *** $P < 0.001$, ** $P < 0.01$, * $P < 0.05$ (FDR-corrected). Hierarchical clustering reveals pathway-specific correlation patterns between core and auxiliary genes.



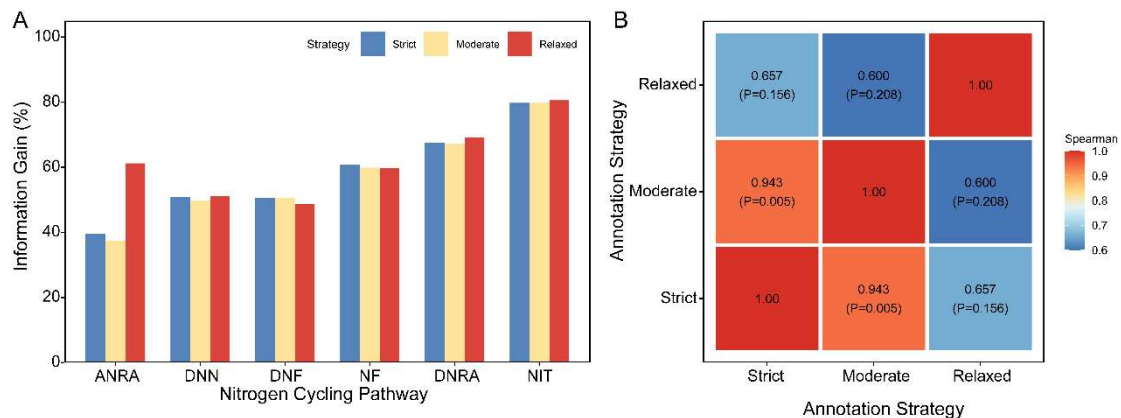
571

572 Figure S3. Cluster number validation for Ward D2 hierarchical clustering of 77 bacterial classes.

573 Gap statistic (red solid line, left y-axis) and mean silhouette coefficient (blue dashed line, right y-

574 axis) across $k = 1-8$. Vertical dotted line indicates selected $k = 3$. Point size represents genome

575 count per class.



576

577 Figure S4. Sensitivity analysis of genus-level Information Gain across annotation strategies. (A)

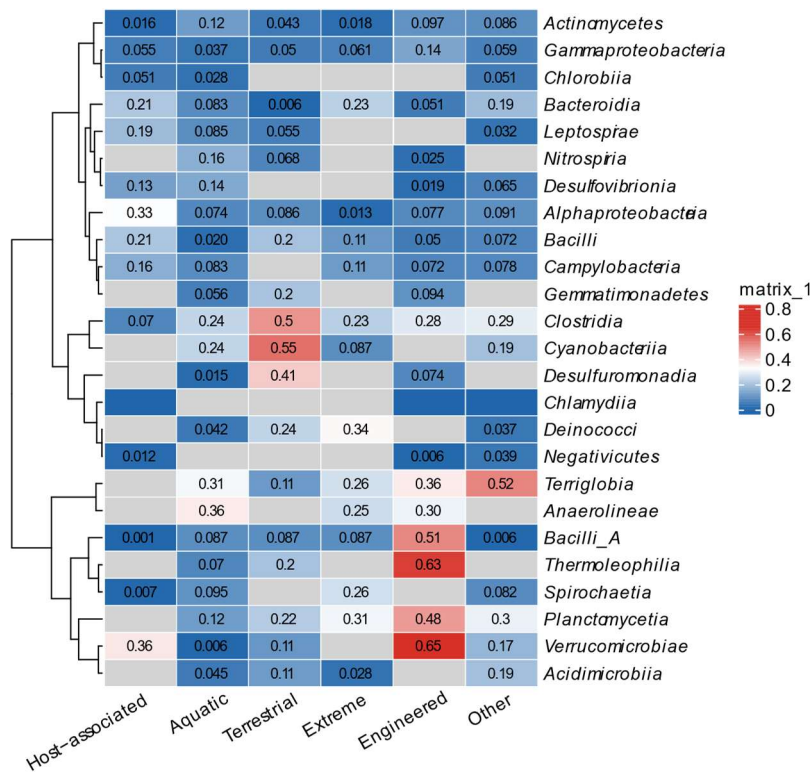
578 Comparison of Information Gain values for six nitrogen cycling pathways under Strict (core plus

579 auxiliary genes), Moderate (core genes only), and Relaxed (any pathway-related gene) annotation

580 strategies. (B) Pairwise Spearman rank correlations between annotation strategies across six

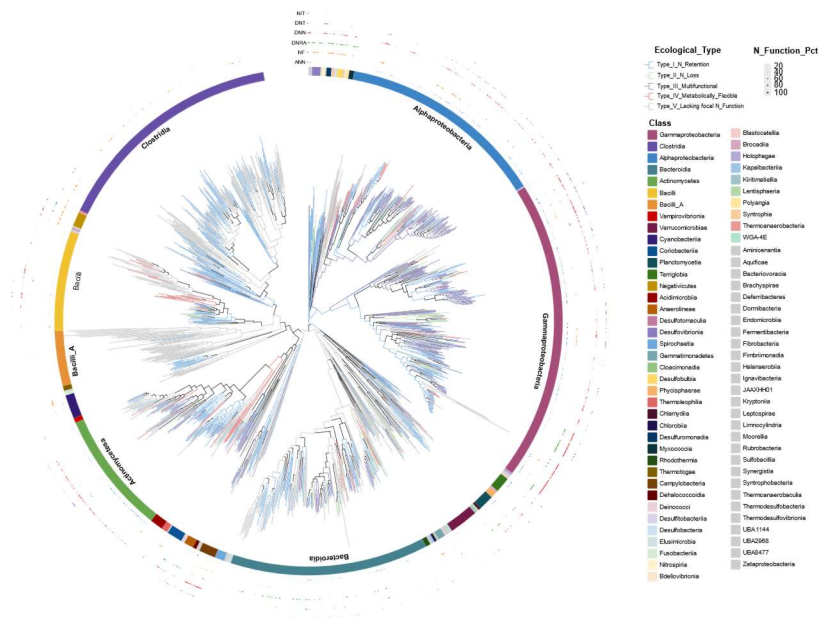
581 pathways.

582



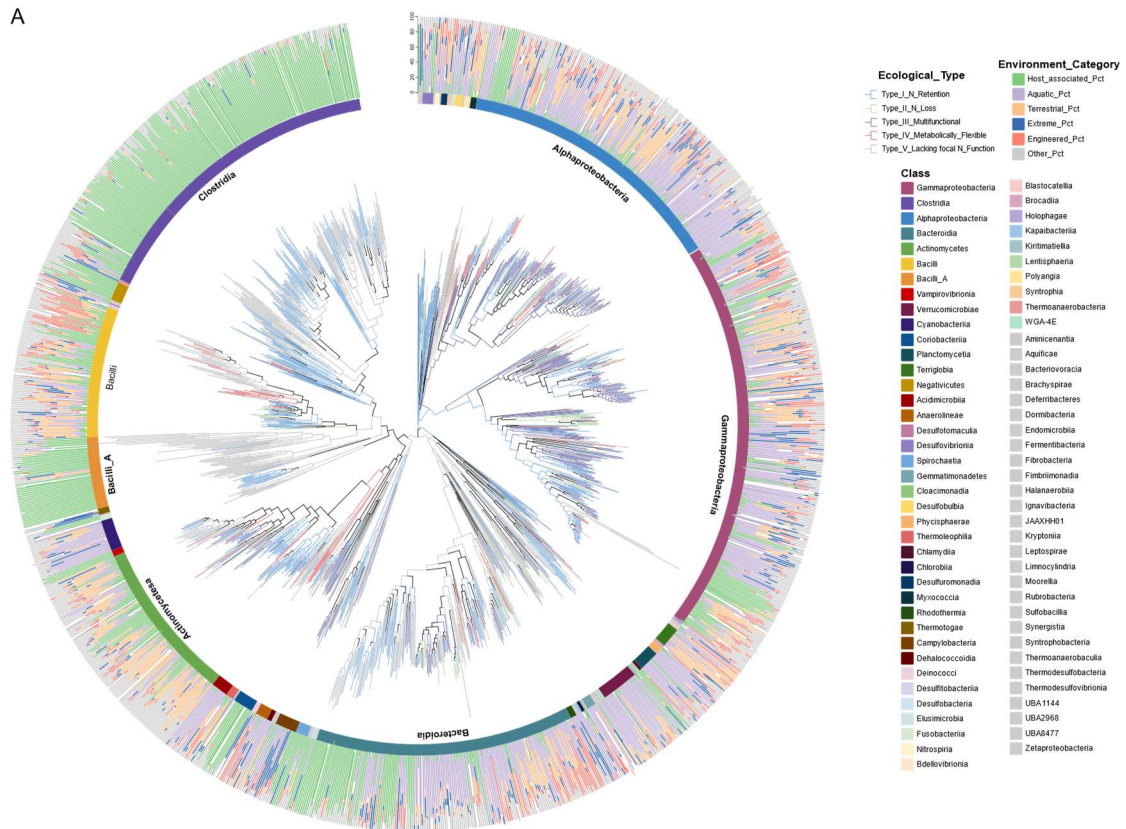
583

584 Figure S5. Environmental adaptation patterns across bacterial classes. Heatmap displaying
 585 adaptation strength calculated as Euclidean distance between class genome source strategy and
 586 average strategy. Higher values indicate greater genome source plasticity



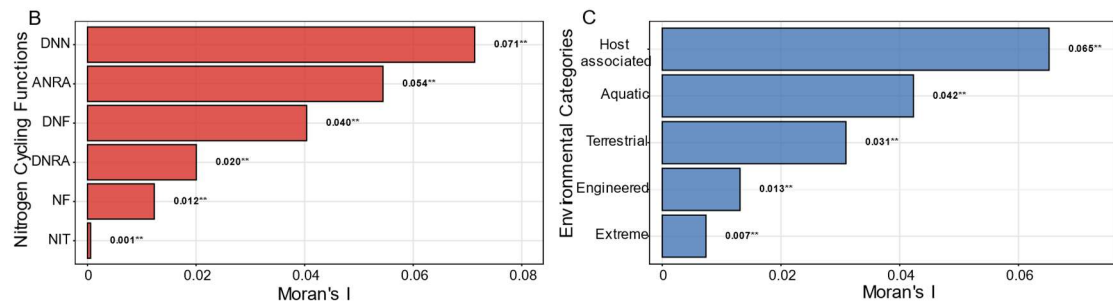
587

588 Figure S6. Phylogenetic distribution of ecological strategies across bacterial genera. Maximum-
 589 likelihood phylogenetic tree of 1,281 qualified genera with branches colored by ecological
 590 strategy types. Outer rings show prevalence patterns of six nitrogen cycling functions within each
 591 genus. Tree construction based on concatenated alignment of 16 ribosomal proteins.



592

593



594

595 Figure S7. Multi-scale phylogenetic constraints in nitrogen cycling functions and genome source

596 preferences. (A) Phylogenetic distribution of ecological strategies and genome source

597 associations. Maximum-likelihood phylogenetic tree of 1,281 genera with tree branches colored

598 by ecological strategy types. Inner ring displays strategy classifications; outer ring shows genome

599 source category distributions. (B) Phylogenetic signal in genome source distributions. Moran's I

600 values for six genome source categories across 1,281 genera showing weak phylogenetic signals.

601 (C) Phylogenetic signal in nitrogen cycling gene frequencies. Moran's I values for six nitrogen

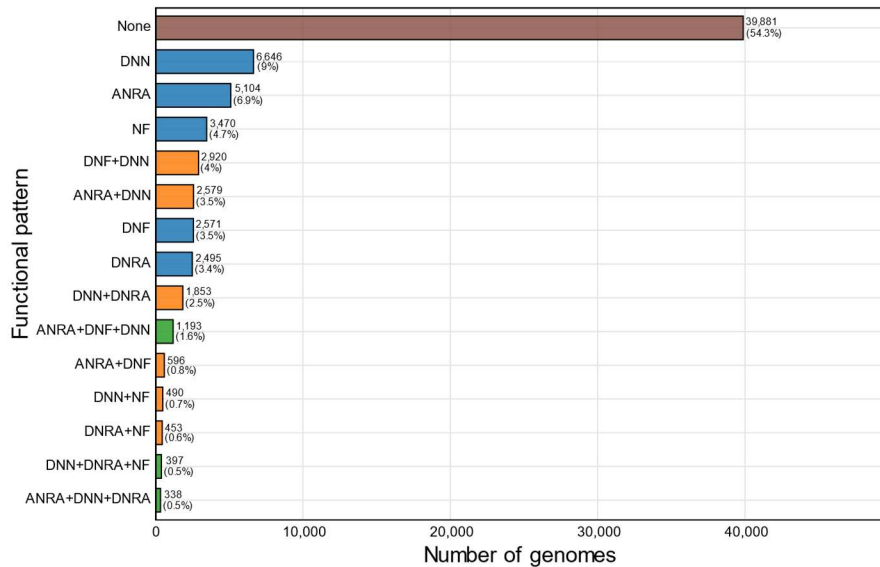
602 cycling functions across 1,281 genera ordered by signal strength. Statistical significance: **P <

603 0.01.

604

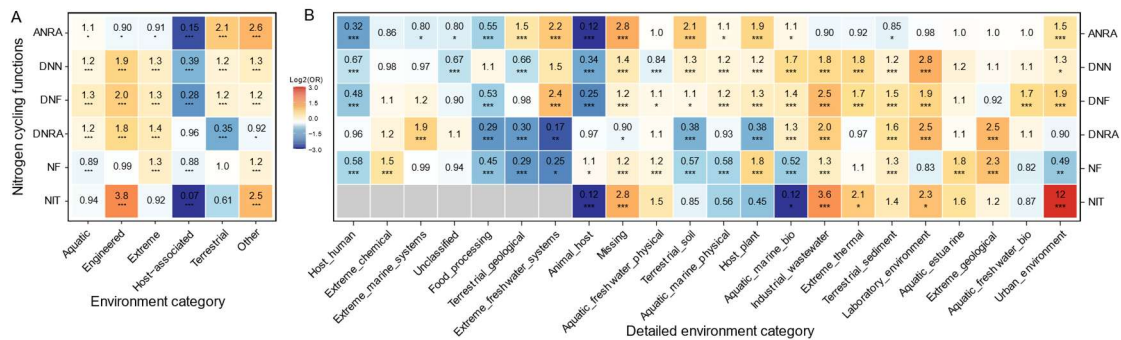
605

606 **Supplementary Note Figure**



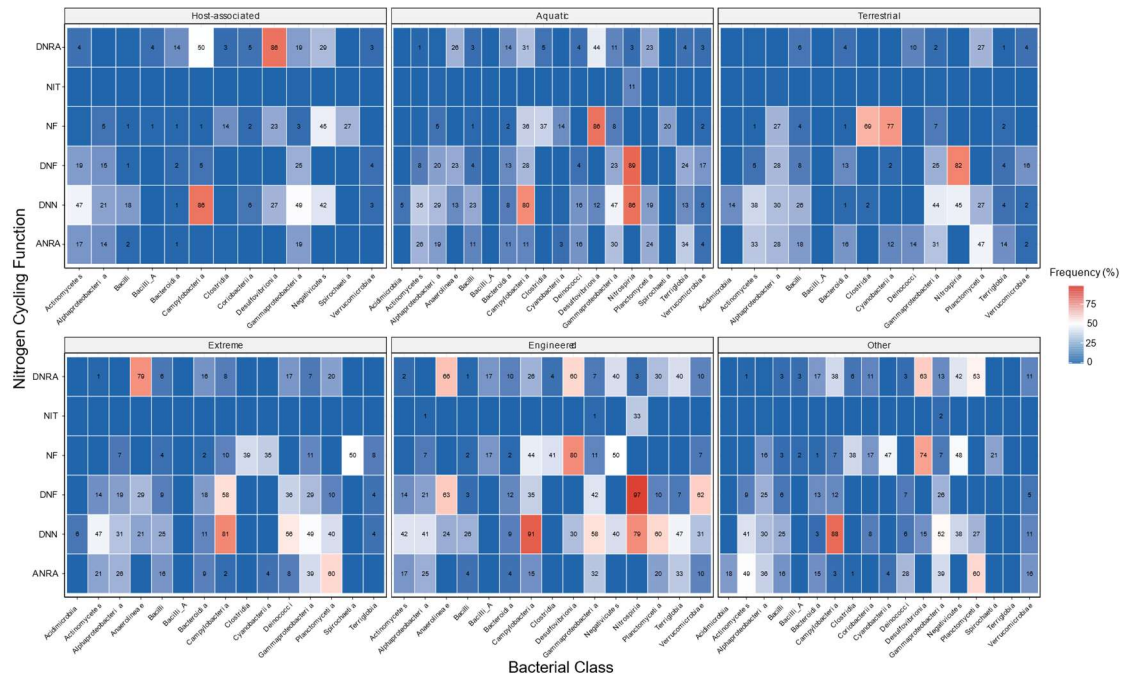
607

608 Figure S8. Distribution of nitrogen cycling functional patterns across bacterial genomes. The 15
 609 most common nitrogen cycling function combinations among 73,472 GTDB genomes. "None"
 610 indicates no detectable function markers.



611

612 Figure S9. Environmental associations of nitrogen cycling functions. (A) Function enrichment
 613 across major genome source categories. (B) Enrichment patterns across detailed source subtypes.
 614 Red: enrichment (OR>1.5), blue: depletion (OR<0.67), grey: function absent (no genomes with
 615 the function detected in that environment). Statistical significance: ***P<0.001, **P<0.01,
 616 *P<0.05.



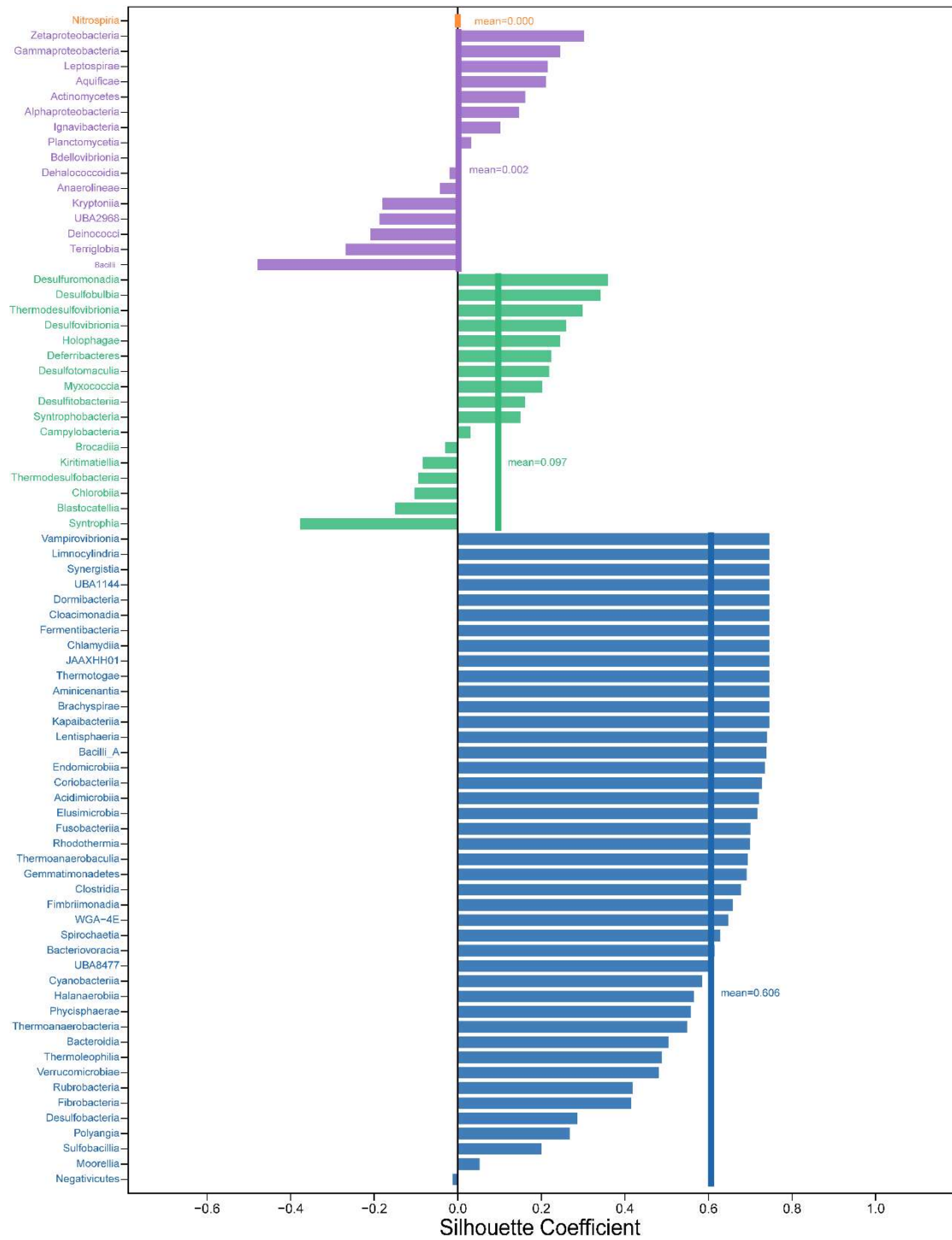
617

618 Figure S10. Class-level functional profiles across major bacterial lineages. Stacked horizontal bars

619 showing distribution of individual nitrogen cycling function frequencies within major bacterial

620 classes (≥500 genomes).

621



622

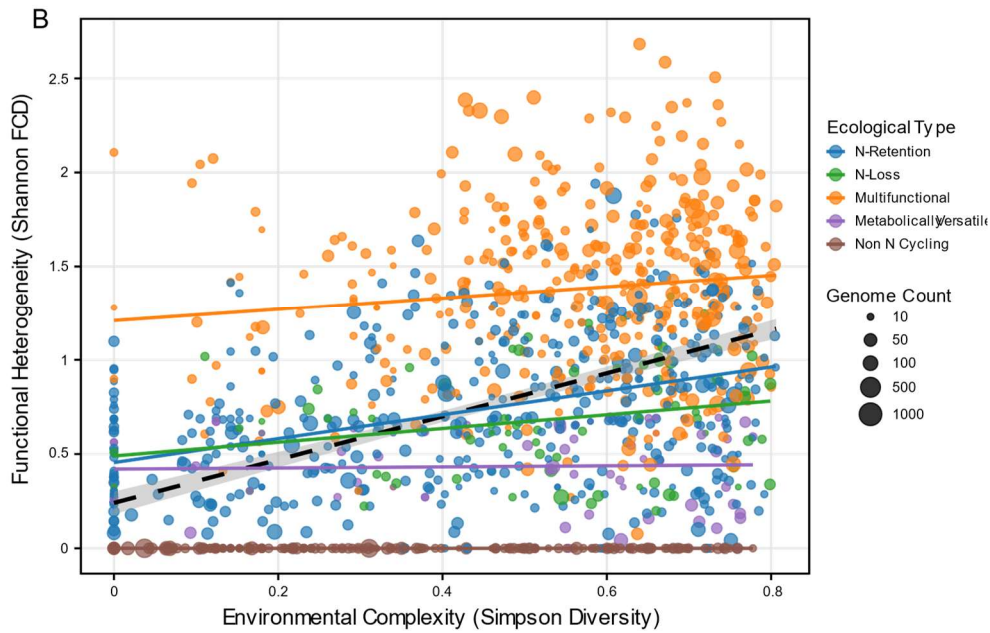
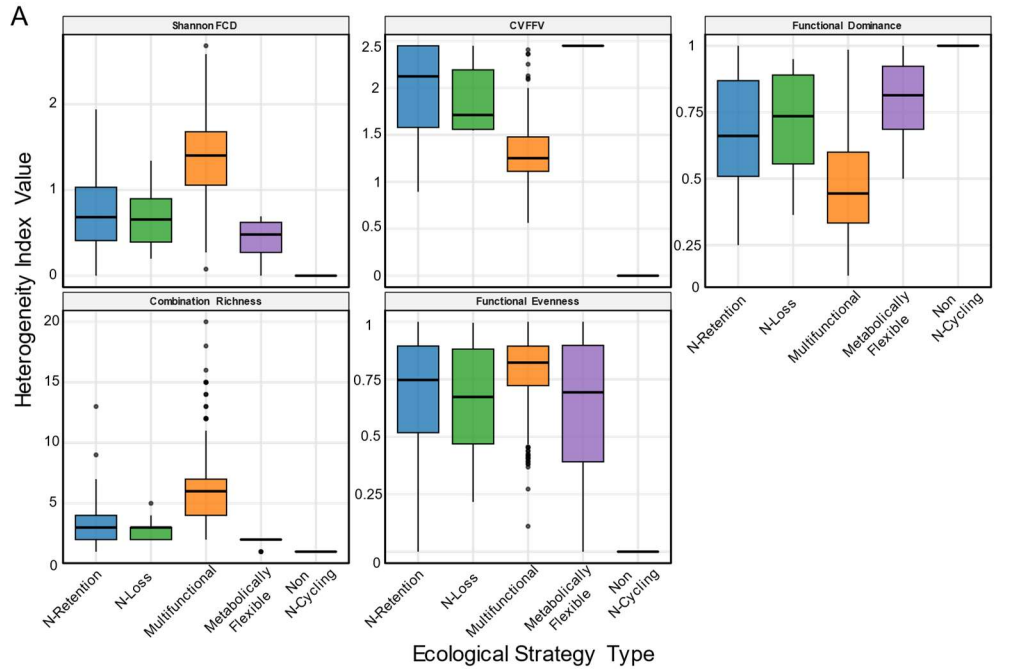
623 Figure S11. Silhouette analysis of Ward D2 hierarchical clustering (k = 3) for 77 bacterial classes.

624 Each bar represents one class, ordered by archetype then silhouette value. Vertical lines indicate

625 archetype mean silhouette values. Nitrification Specialist silhouette set to 0 by convention for

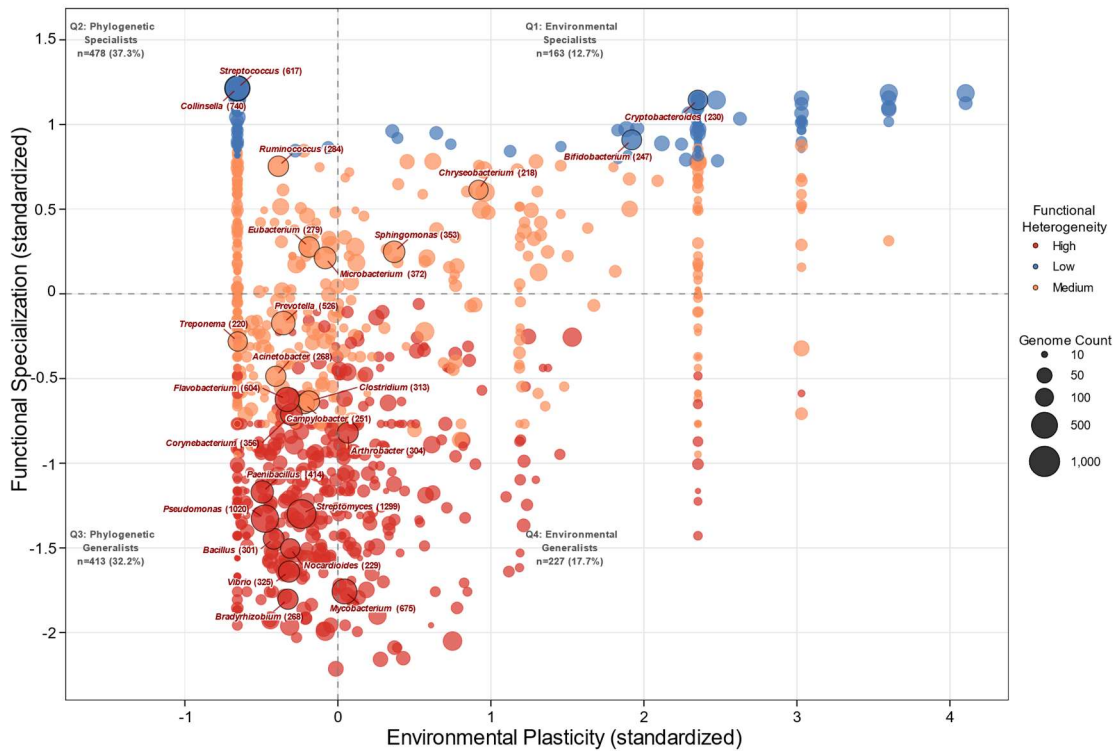
626 single-member clusters.

627



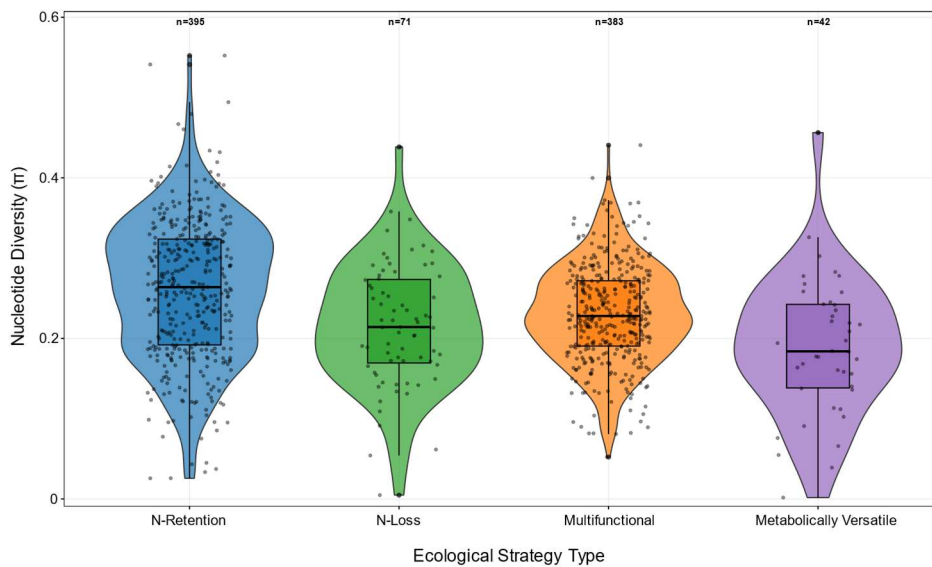
628

629 Figure S12. Functional heterogeneity patterns across ecological strategy types. (A) Functional
 630 heterogeneity indices by ecological strategy type. Five-panel boxplot displays Shannon Functional
 631 Combination Diversity, Coefficient of Variation, Functional Dominance, Combination Richness,
 632 and Functional Evenness across qualified genera. Kruskal-Wallis tests demonstrated significant
 633 variation across all indices. (B) Genome source distribution breadth versus functional
 634 heterogeneity across ecological strategy types. Point colors indicate strategy types. Linear
 635 regression: $R^2 = 0.68$, $p < 0.001$.



636

637 Figure S13. Four-quadrant integration framework for evolutionary strategies. Scatterplot
 638 positioning qualified genera along Genome Source Plasticity (x-axis) and Functional
 639 Specialization (y-axis). Point colors represent functional heterogeneity levels; quadrant analysis
 640 reveals four distinct adaptation strategies.



641

642 Figure S14. Conservation indices across ecological strategy types. Violin plots showing
 643 conservation index distributions across ecological strategy types (n = 814 genera). Box plots
 644 indicate quartiles. Kruskal-Wallis test: $\chi^2 = 54.96$, $P < 0.001$.

645 **Supplementary Tables**646 Supplementary Table S1. Core and auxiliary genes used in nitrogen transformation pathway
647 identification

Pathway	Gene Symbol	KEGG Annotation	RAST Annotation	Core Gene
Assimilatory Nitrate Reduction (ANRA)	<i>nasA</i>	Assimilatory nitrate reductase electron transfer subunit (NasB)	Assimilatory nitrate reductase large subunit (EC 1.7.99.4)	Yes
		Nitrite reductase (NADH) large subunit (NirB)	Nitrite reductase [NAD(P)H] large subunit (EC 1.7.1.4)	Yes
		Nitrite reductase (NADH) small subunit (NirD)	Nitrite reductase [NAD(P)H] small subunit (EC 1.7.1.4)	No
Dissimilatory Nitrate to Nitrite (DNN)	<i>napA</i>	Periplasmic nitrate reductase (NapA)	Periplasmic nitrate reductase (EC 1.7.99.4)	Yes
		Cytochrome c550-type protein (NapB)	Nitrate reductase cytochrome c550-type subunit	No
	<i>napC</i>	Cytochrome c-type protein NapC	Cytochrome c-type protein NapC	No
		Nitrate reductase delta subunit (NarG)	Respiratory nitrate reductase alpha chain (EC 1.7.99.4)	Yes
	<i>narH</i>	Nitrate reductase gamma subunit (NarH)	Respiratory nitrate reductase beta chain (EC 1.7.99.4)	No
		Nitrate reductase alpha subunit (NarI)	Respiratory nitrate reductase gamma chain (EC 1.7.99.4)	No
	<i>narJ</i>	Nitrate reductase delta subunit (NarJ)	Respiratory nitrate reductase delta chain (EC 1.7.99.4)	No
		Cytochrome b-561 (NarC)	Respiratory nitrate reductase subunit, conjectural (EC 1.7.99.4)	No
Dissimilatory Nitrite to Ammonia (DNRA)	<i>nrfA</i>	Nitrite reductase cytochrome c-552 (NrfA)	Cytochrome c552 precursor (EC 1.7.2.2)	Yes
		Protein NrfH (NrfH)	Cytochrome c nitrite reductase, small subunit NrfH	No
	<i>nrfB</i>	Cytochrome c-type protein (NrfB)	Cytochrome c-type protein NrfB precursor	No

Denitrification (DNF)	<i>nrfC</i>	Nitrite reductase protein NrfC (NrfC)	NrfC protein	No
	<i>nrfD</i>	Protein nrfD (NrfD)	NrfD protein	No
	<i>nirK</i>	Copper-containing nitrite reductase (NirK)	Copper-containing nitrite reductase (EC 1.7.2.1)	Yes
	<i>nirS</i>	Cytochrome cd1 nitrite reductase (NirS)	Cytochrome cd1 nitrite reductase (EC 1.7.2.1)	Yes
	<i>qnor</i>	Nitric-oxide reductase quinol-dependent (qNor)	Nitric-oxide reductase (EC 1.7.99.7), quinol- dependent	No
	<i>cnorB</i>	Nitric-oxide reductase subunit B (cNor-B)	Nitric-oxide reductase subunit B (EC 1.7.99.7)	No
	<i>cnorC</i>	Nitric-oxide reductase subunit C (cNor-C)	Nitric-oxide reductase subunit C (EC 1.7.99.7)	No
	<i>nosZ</i>	Nitrous-oxide reductase (NosZ)	Nitrous-oxide reductase (EC 1.7.99.6)	No
Nitrogen Fixation(NF)	<i>nifH</i>	Nitrogenase reductase (NifH)	Nitrogenase (molybdenum-iron) reductase and maturation protein NifH	Yes
	<i>nifD</i>	Nitrogenase molybdenum-iron protein subunit alpha (NifD)	Nitrogenase (molybdenum-iron) alpha chain (EC 1.18.6.1)	No
	<i>nifK</i>	Nitrogenase molybdenum-iron protein subunit beta (NifK)	Nitrogenase (molybdenum-iron) beta chain (EC 1.18.6.1)	No
	<i>nifB</i>	Nitrogenase molybdenum-iron cofactor biosynthesis protein (NifN)	Nitrogenase FeMo- cofactor synthesis FeS core scaffold and assembly protein NifB	No
	<i>nifE</i>	Nitrogenase molybdenum-cofactor biosynthesis protein (NifE)	Nitrogenase FeMo- cofactor scaffold and assembly protein NifE	No
	<i>nifN</i>	Nitrogenase iron protein 2 (NifH2)	Nitrogenase FeMo- cofactor scaffold and assembly protein NifN	No
	<i>nifV</i>	Homocitrate synthase (NifV)	Homocitrate synthase (EC 2.3.3.14)	No
	Nitrification (NIT)	<i>amoA</i>	Ammonia monooxygenase (Amo)	Ammonia monooxygenase

657 Supplementary Table S4. Kruskal–Wallis test results for inter-archetype differences across
 658 functional metrics.

Metric	H	p	η^{2a}
Participation rate	60.59	< 0.001	0.789
Nitrogen retention index	58.97	< 0.001	0.767
DNRA frequency	42.05	< 0.001	0.535
DNN frequency	31.77	< 0.001	0.394
NF frequency	25.70	< 0.001	0.311
DNF frequency	24.58	< 0.001	0.296
Nitrogen loss index	24.58	< 0.001	0.296
ANRA frequency	20.67	< 0.001	0.242

659 ^a $\eta^2 = (H - k + 1)/(n - k)$; calculated across Functionally Inactive (n = 25), Functionally Moderate
 660 (n = 34), and N-Retention Dominant (n = 17). Nitrification Specialist (n = 1) excluded from
 661 Kruskal–Wallis analysis.

662
 663

665 **Supplementary References:**

- 666 Ahn, S., Cho, M., Sadowsky, M.J., Jang, J. (2025) Dissimilatory nitrate reductions in soil *Neobacillus* and
667 *Bacillus* strains under aerobic condition. *Journal of Microbiology* 63, e2411019.
- 668 Arnold, B.J., Huang, I.T., Hanage, W.P. (2022) Horizontal gene transfer and adaptive evolution in
669 bacteria. *Nature Reviews Microbiology* 20, 206-218.
- 670 Benjamini, Y., Hochberg, Y. (1995) CONTROLLING THE FALSE DISCOVERY RATE - A
671 PRACTICAL AND POWERFUL APPROACH TO MULTIPLE TESTING. *Journal of the Royal*
672 *Statistical Society Series B-Statistical Methodology* 57, 289-300.
- 673 Canty, A., Ripley, B. (2017) boot: Bootstrap R (S-Plus) Functions. R package version 1.3-20. CRAN R
674 Project.
- 675 Capella-Gutiérrez, S., Silla-Martínez, J.M., Gabaldón, T. (2009) trimAl: a tool for automated alignment
676 trimming in large-scale phylogenetic analyses. *Bioinformatics* 25, 1972-1973.
- 677 Daims, H., Lückner, S., Wagner, M. (2016) A new perspective on microbes formerly known as nitrite-
678 oxidizing bacteria. *Trends in microbiology* 24, 699-712.
- 679 Dunn, O.J. (1964) Multiple comparisons using rank sums. *Technometrics* 6, 241-252.
- 680 Durand, S., Guillier, M. (2021) Transcriptional and post-transcriptional control of the nitrate respiration
681 in bacteria. *Frontiers in molecular biosciences* 8, 667758.
- 682 Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput.
683 *Nucleic Acids Research* 32, 1792-1797.
- 684 Garnier, S., Ross, N., Rudis, R., Camargo, P.A., Sciaini, M., Scherer, C. (2021) viridis-Colorblind-
685 friendly color maps for R. R package version 0.6 2.
- 686 Gómez-Rubio, V. (2017) ggplot2-elegant graphics for data analysis. *Journal of Statistical Software* 77,
687 1-3.
- 688 Graf, D.R.H., Jones, C.M., Hallin, S. (2014) Intergenomic comparisons highlight modularity of the
689 denitrification pathway and underpin the importance of community structure for N₂O emissions. *PLoS*
690 *One* 9, e114118.
- 691 Guerrero, M.G., Lara, C. (1987) Assimilation of inorganic nitrogen.
- 692 Holm, S. (1979) A simple sequentially rejective multiple test procedure. *Scandinavian journal of statistics*,
693 65-70.
- 694 Isobe, K., Allison, S.D., Khalili, B., Martiny, A.C., Martiny, J.B.H. (2019) Phylogenetic conservation of
695 bacterial responses to soil nitrogen addition across continents. *Nature Communications* 10.
- 696 Jombart, T., Balloux, F., Dray, S. (2010) adephylo: new tools for investigating the phylogenetic signal in
697 biological traits. *Bioinformatics* 26, 1907-1909.
- 698 Katoh, K., Standley, D.M. (2013) MAFFT Multiple Sequence Alignment Software Version 7:
699 Improvements in Performance and Usability. *Mol Biol Evol* 30, 772-780.
- 700 Kolde, R., (2019) pheatmap: Pretty Heatmaps. R package version 1.0. 12.
- 701 Kosakovsky Pond, S.L., Frost, S.D.W. (2005) Not so different after all: a comparison of methods for
702 detecting amino acid sites under selection. *Mol Biol Evol* 22, 1208-1222.
- 703 Kruskal, W.H., Wallis, W.A. (1952) Use of ranks in one-criterion variance analysis. *Journal of the*
704 *American statistical Association* 47, 583-621.
- 705 Kuypers, M.M.M., Marchant, H.K., Kartal, B. (2018) The microbial nitrogen-cycling network. *Nature*
706 *Reviews Microbiology* 16, 263-276.

707 Liu, S., Dai, J., Wei, H., Li, S., Wang, P., Zhu, T., Zhou, J., Qiu, D. (2021) Dissimilatory nitrate reduction
708 to ammonium (DNRA) and denitrification pathways are leveraged by cyclic AMP receptor protein (CRP)
709 paralogues based on electron donor/acceptor limitation in *Shewanella loihica* PV-4. *Applied and*
710 *Environmental Microbiology* 87, e01964-01920.

711 Lückner, S., Wagner, M., Maixner, F., Pelletier, E., Koch, H., Vacherie, B., Rattei, T., Damsté, J.S.S.,
712 Spieck, E., Le Paslier, D. (2010) A *Nitrospira* metagenome illuminates the physiology and evolution of
713 globally important nitrite-oxidizing bacteria. *Proceedings of the National Academy of Sciences* 107,
714 13479-13484.

715 Lynch, M. (2007) The evolution of genetic networks by non-adaptive processes. *Nature Reviews*
716 *Genetics* 8, 803-813.

717 Martiny, A.C., Treseder, K., Pusch, G. (2013) Phylogenetic conservatism of functional traits in
718 microorganisms. *Isme Journal* 7, 830-838.

719 Ming, Y., Al, M.A., Zhang, D., Zhu, W., Liu, H., Cai, L., Yu, X., Wu, K., Niu, M., Zeng, Q. (2024)
720 Insights into the evolutionary and ecological adaptation strategies of nirS- and nirK-type denitrifying
721 communities. *Molecular Ecology* 33, 13.

722 Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., Von Haeseler, A., Lanfear,
723 R. (2020) IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era.
724 *Mol Biol Evol* 37, 1530-1534.

725 Morris, J.J., Lenski, R.E., Zinser, E.R. (2012) The Black Queen Hypothesis: Evolution of Dependencies
726 through Adaptive Gene Loss. *Mbio* 3.

727 Mosley, O.E., Gios, E., Close, M., Weaver, L., Daughney, C., Handley, K.M. (2022) Nitrogen cycling
728 and microbial cooperation in the terrestrial subsurface. *Isme Journal* 16, 2561-2573.

729 Murali, R., Pace, L.A., Sanford, R.A., Ward, L.M., Lynes, M.M., Hatzenpichler, R., Lingappa, U.F.,
730 Fischer, W.W., Gennis, R.B., Hemp, J. (2024) Diversity and evolution of nitric oxide reduction in bacteria
731 and archaea. *Proceedings of the National Academy of Sciences* 121, e2316422121.

732 Nguyen, L.T., Schmidt, H.A., von Haeseler, A., Minh, B.Q. (2015) IQ-TREE: A Fast and Effective
733 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol* 32, 268-274.

734 Paradis, E., Schliep, K. (2019) ape 5.0: an environment for modern phylogenetics and evolutionary
735 analyses in R. *Bioinformatics* 35, 526-528.

736 Parks, D.H., Chuvpochina, M., Waite, D.W., Rinke, C., Skarshewski, A., Chaumeil, P.A., Hugenholtz, P.
737 (2018) A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of
738 life. *Nature Biotechnology* 36, 996-+.

739 Potter, L., Angove, H., Richardson, D., Cole, J. (2001) Nitrate reduction in the periplasm of gram-
740 negative bacteria.

741 Schimel, J.P., Bennett, J. (2004) Nitrogen mineralization: challenges of a changing paradigm. *Ecology*
742 85, 591-602.

743 Schurig-Briccio, L.A., Venkatakrisnan, P., Hemp, J., Briccio, C., Berenguer, J., Gennis, R.B. (2013)
744 Characterization of the nitric oxide reductase from *Thermus thermophilus*. *Proceedings of the National*
745 *Academy of Sciences* 110, 12613-12618.

746 Seitzinger, S., Harrison, J.A., Böhlke, J.K., Bouwman, A.F., Lowrance, R., Peterson, B., Tobias, C.,
747 Drecht, G.V. (2006) Denitrification across landscapes and waterscapes: a synthesis. *Ecological*
748 *applications* 16, 2064-2090.

749 Shade, A., Gilbert, J.A. (2015) Temporal patterns of rarity provide a more complete view of microbial
750 diversity. *Trends in microbiology* 23, 335-340.

751 Shannon, C.E. (1948) A mathematical theory of communication. The Bell system technical journal 27,
752 379-423.

753 Sparacino-Watkins, C., Stolz, J.F., Basu, P. (2014) Nitrate and periplasmic nitrate reductases. Chemical
754 Society Reviews 43, 676-706.

755 Templer, P.H., Mack, M.C., Iii, F.S.C., Christenson, L.M., Compton, J.E., Crook, H.D., Currie, W.S.,
756 Curtis, C.J., Dail, D.B., D'Antonio, C.M. (2012) Sinks for nitrogen inputs in terrestrial ecosystems: a
757 meta-analysis of ¹⁵N tracer field studies. Ecology 93, 1816-1829.

758 Tomczak, M., Tomczak, E. (2014) The need to report effect size estimates revisited. An overview of some
759 recommended measures of effect size.

760 Tu, Q.C., Lin, L., Cheng, L., Deng, Y., He, Z.L. (2019) NCycDB: a curated integrative database for fast
761 and accurate metagenomic profiling of nitrogen cycling genes. Bioinformatics 35, 1040-1048.

762 Wei, T., Simko, V., Levy, M., Xie, Y., Jin, Y., Zemla, J. (2017) Package 'corrplot'. Statistician 56, e24.

763 Yoon, S., Cruz-García, C., Sanford, R., Ritalahti, K.M., Löffler, F.E. (2015) Denitrification versus
764 respiratory ammonification: environmental controls of two competing dissimilatory NO₃⁻/NO₂⁻
765 reduction pathways in *Shewanella loihica* strain PV-4. The ISME journal 9, 1093-1104.

766 Yu, G., Smith, D.K., Zhu, H., Guan, Y., Lam, T.T.Y. (2017) ggtree: an R package for visualization and
767 annotation of phylogenetic trees with their covariates and other associated data. Methods in Ecology and
768 Evolution 8, 28-36.

769