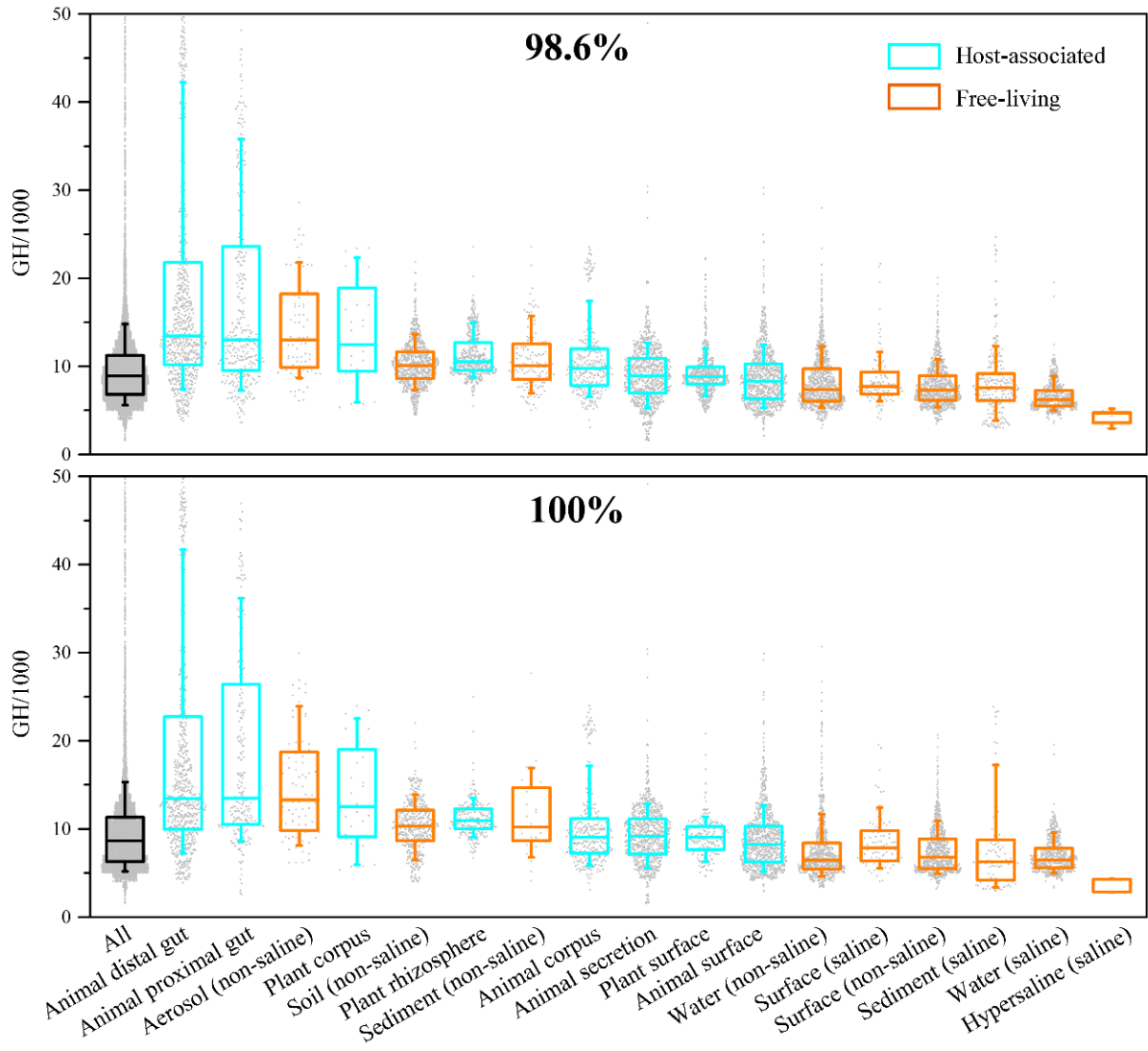


## **Supporting Information**

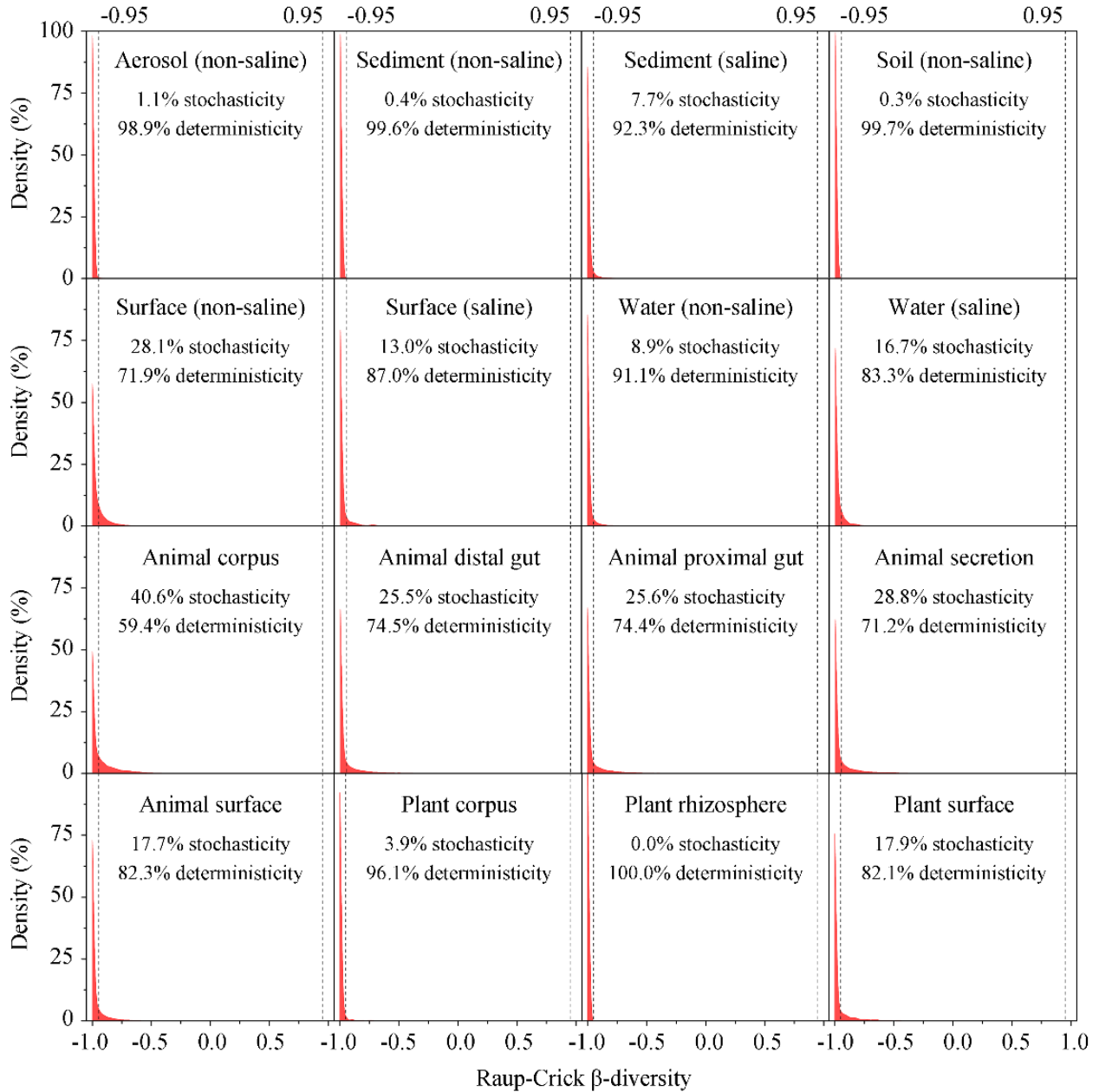
# **Deterministic assembly processes shaping habitat-specific glycoside hydrolase composition**

List:

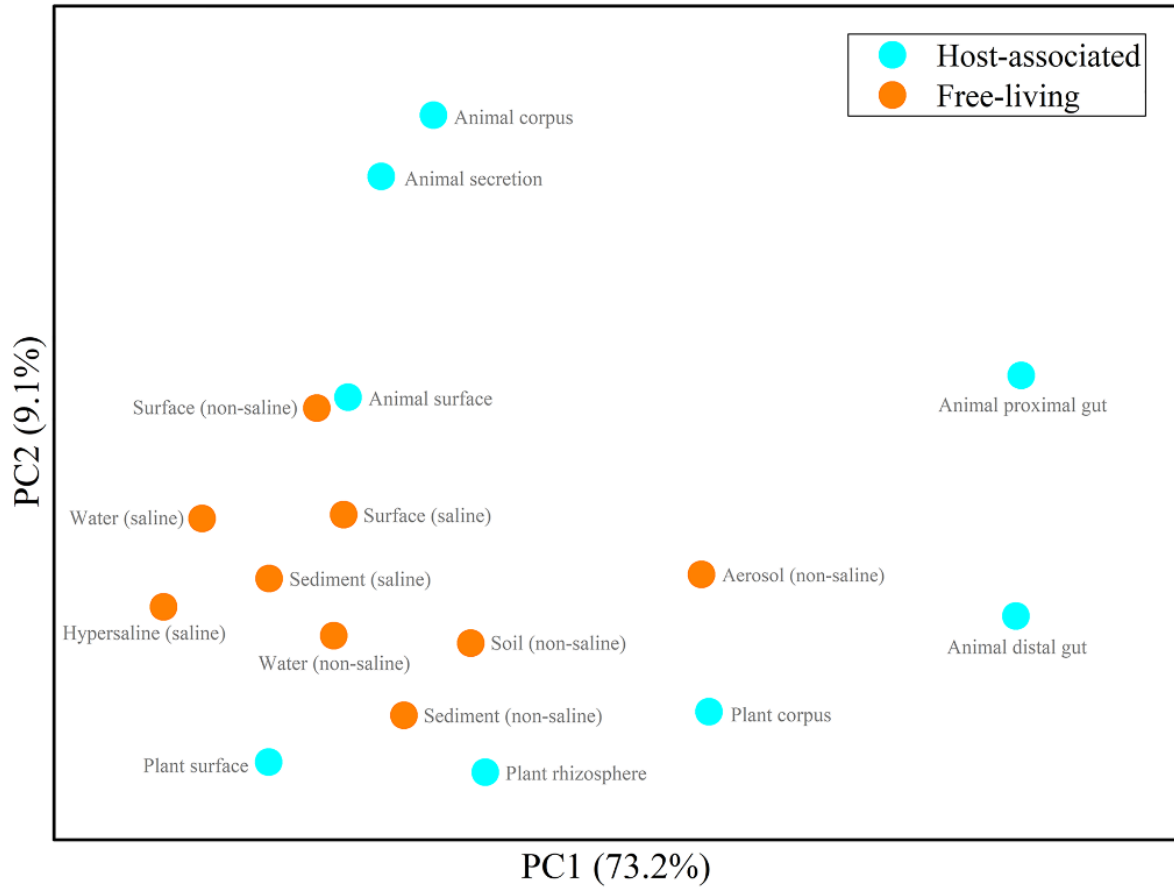
Figure S1-S13, Data S1-S4.



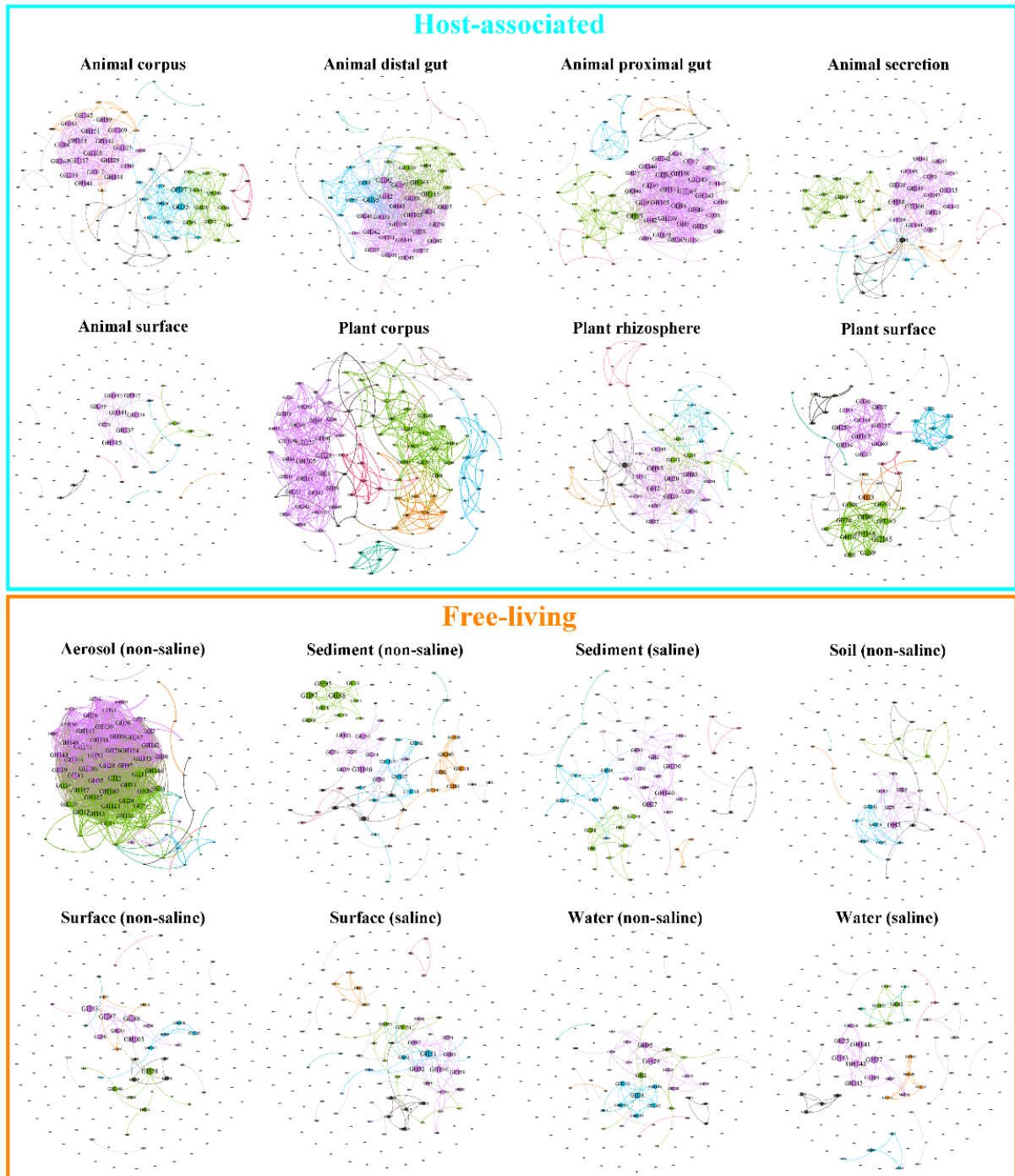
**Figure S1. Differences in GH abundance within prokaryotic communities in different environments.** The OTUs corresponding to the CAZy genomes were determined based on a 16S rRNA-V4 region identity of greater than 98.6% or 100%. The GH abundance was calculated separately for the two sample sets, with each grey point representing an EMP sample. For the boxplots, the middle line indicates the median, the box represents the 25th-75th percentile, and the whisker indicates the 10th-90th percentile of observations. EMPO classified 17 microbial environments (level 3) into free-living or host-associated (level 1), and saline or non-saline (for free-living) or animal or plant (for host-associated) (level 2). Host-associated environments are represented in cyan, while free-living environments are represented in orange.



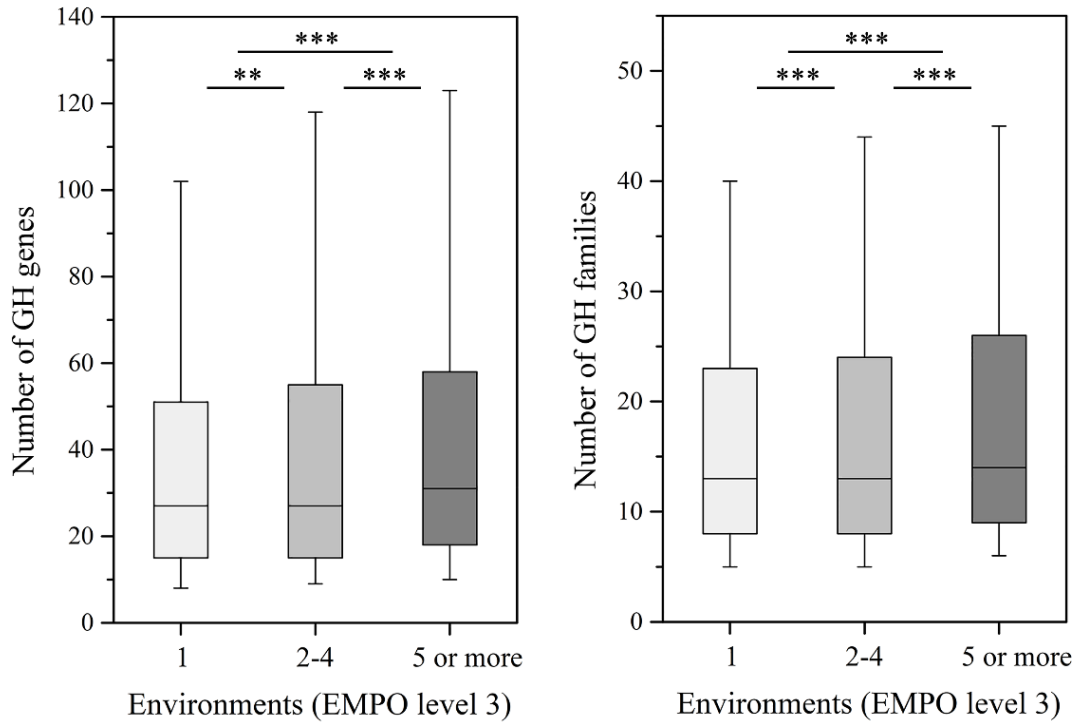
**Figure S2. Deterministic processes dominate the assembly of GH communities.** Pairwise Raup-Crick between communities indicates that the GH genes are either more similar ( $-0.95 > \text{Raup-Crick} > -1$ ) or less similar ( $0.95 < \text{Raup-Crick} < 1$ ) to each other than expected by random chance (1000 randomizations). Any other value of Raup-Crick indicates that the GH communities are stochastically assembled ( $-0.95 < \text{Raup-Crick} < 0.95$ ). The vertical dashed lines mark the positions of  $-0.95$  and  $0.95$  in panels. The contribution ratios of deterministic and stochastic processes are also marked in panels.



**Figure S3. Similarity of GH composition across various environments.** Principal coordinate analysis (PCoA) of the GH distribution characteristics in EMPO level 3 environments was performed using the PAST3 software based on the Bray-Curtis method.

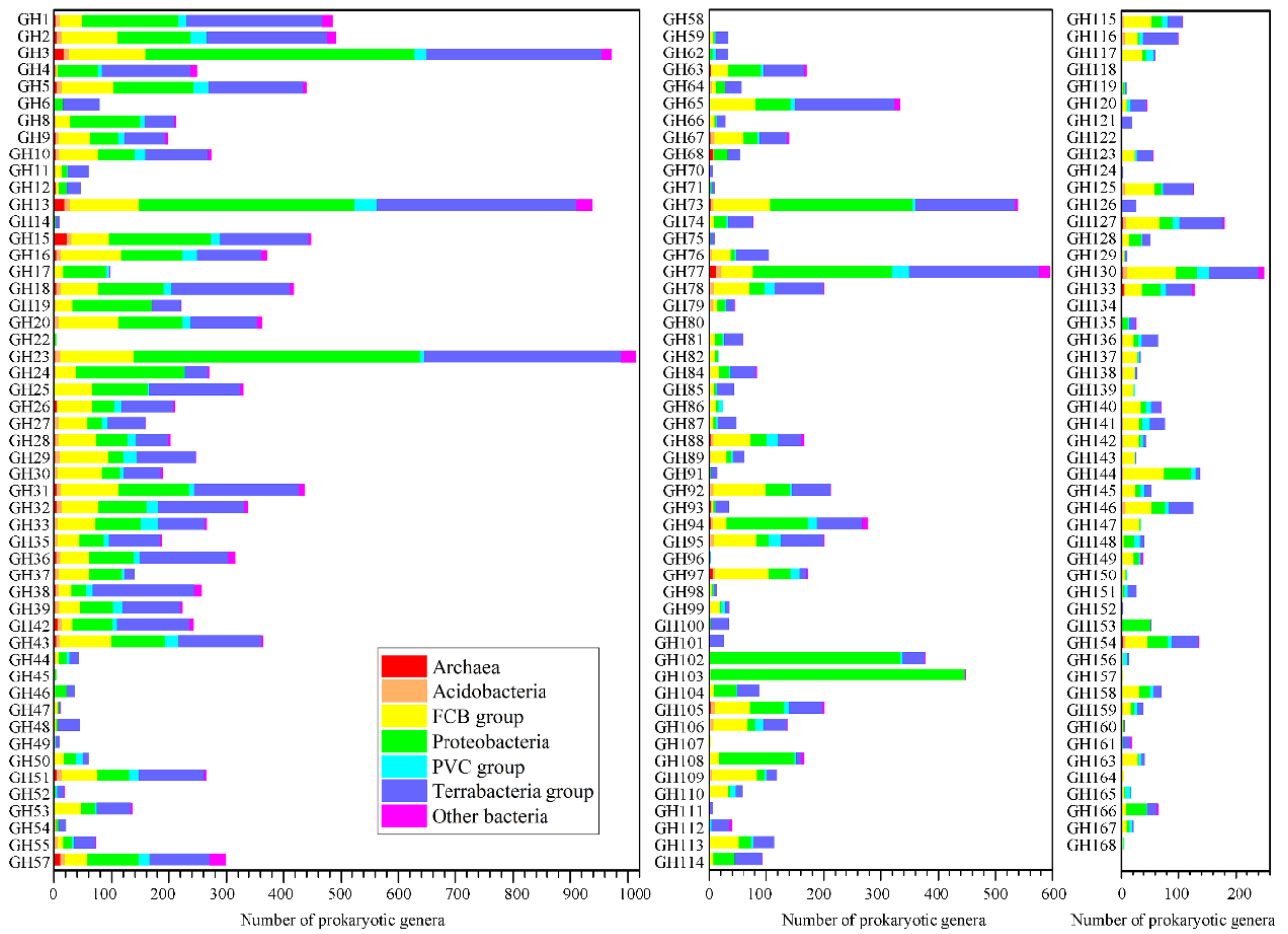


**Figure S4. A co-occurrence network diagram exhibiting the correlations between various GH families within the environment.** The EMPO classified environmental samples into 17 microbial environments. The hypersaline (saline) environmental samples were too few to be displayed. A co-occurrence network was constructed by Gephi, and the pairwise Spearman's  $\rho$  and  $p$  values were determined using the “psych” package in R. Values of Spearman's  $\rho > 0.8$  and  $p < 0.01$  were considered to indicate valid relationships (connecting lines in the network diagram).

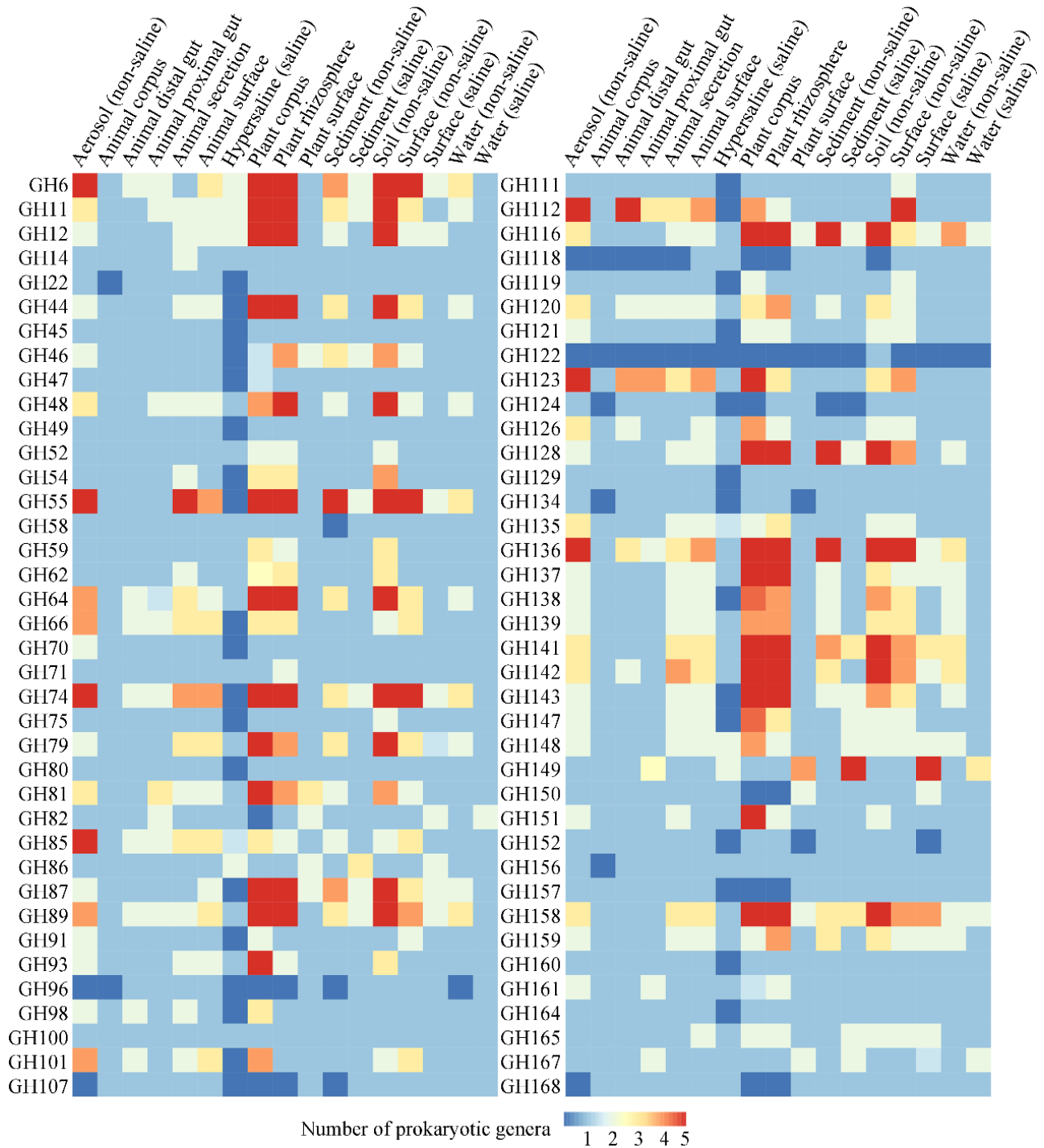


**Figure S5. Potential impact of prokaryotic-encoded GHs on their environmental adaptability.**

Comparisons between bins were analyzed using the Wilcoxon rank test. \*\*  $p < 0.01$ , \*\*\*  $p < 0.001$ .

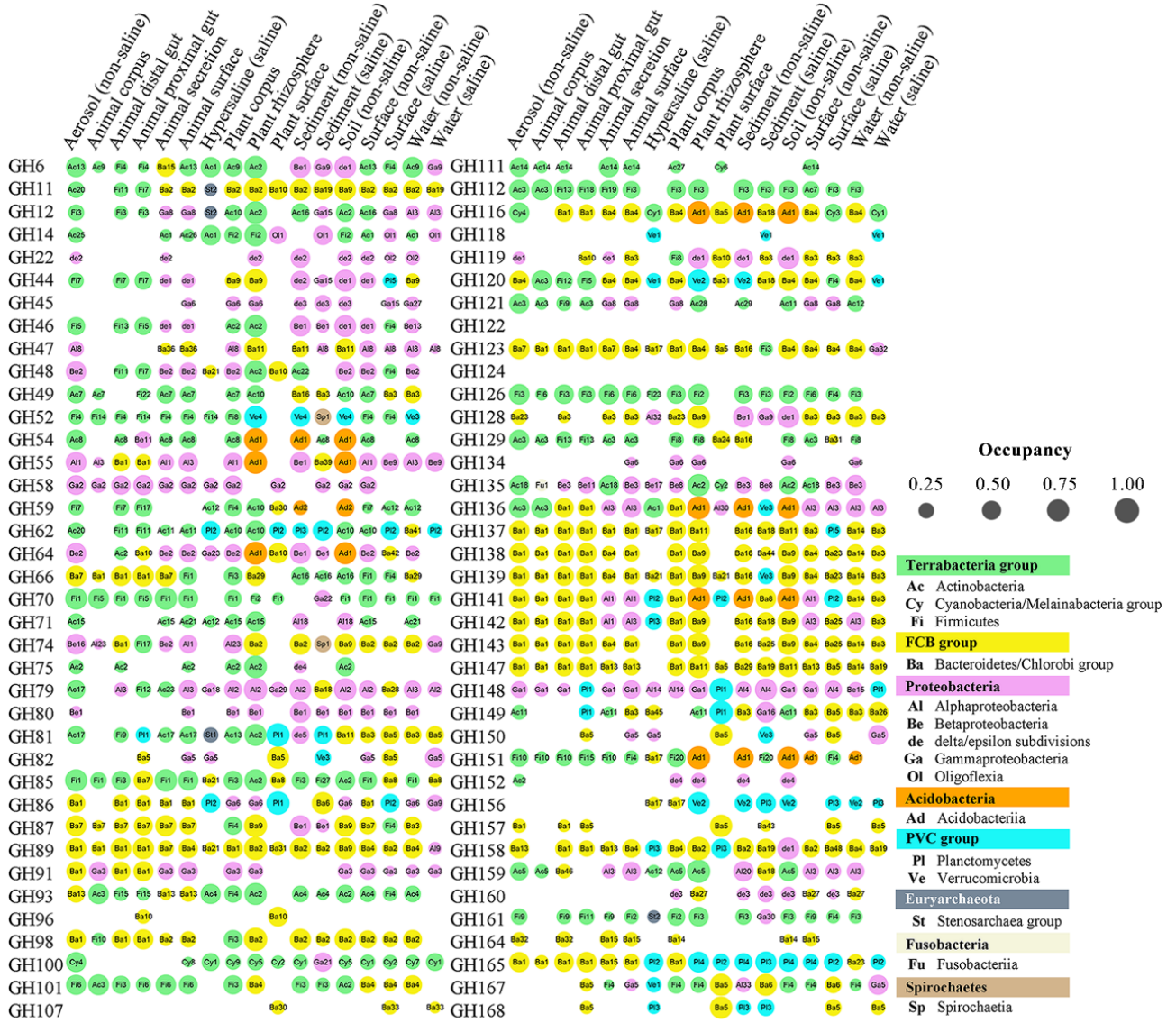


**Figure S6. Taxonomic source statistics of GH families in all EMP samples.** The number of prokaryotic genera encoding GHs in EMP samples was counted using the NCBI Taxonomy database. The different colors in the figure indicate the prokaryotic genera from distinct taxonomic groups.



**Figure S7. Diverse taxonomic sources of GH families in prokaryotic communities.** Statistical analyses were carried out at EMPO level 3. The colors in the figure indicate the median number of source prokaryotic genera of GH families in all EMP samples from the same environment, with increasing values from blue to red. Only 76 GH families with an occupancy of less than 0.75 in all EMP samples are displayed, and the other 76 GH families are shown in Figure 6.





**Figure S8. Most common prokaryotic genera of GH families in EMPO environments.** The most common prokaryotic genus was defined as the genus that appeared in the most samples and had an occupancy of at least 0.05 in all samples. The EMP samples were classified into 17 microbial environments by EMPO. Different colors represent different taxonomic groups. The circle size indicates the occupancy of the genus in all samples within a single environment. Only 76 GH families with an occupancy of less than 0.75 in all EMP samples are displayed, and the other 76 GH families are shown in Figure 7. Specific information on the prokaryotic genera can be found in Figure S9.

**Proteobacteria****Alphaproteobacteria**

A11	<i>Sphingomonas</i>
A12	<i>Bradyrhizobium</i>
A13	<i>Sphingobium</i>
A14	<i>Smorhizobium</i>
A15	<i>Celeribacter</i>
A16	<i>Sulfobacter</i>
A17	<i>Halocynthiibacter</i>
A18	<i>Brevundimonas</i>
A19	<i>Caulobacter</i>
A110	<i>Ruegeria</i>
A111	<i>Nordella</i>
A112	<i>Pseudolabrys</i>
A113	<i>Agrobacterium</i>
A114	<i>Fansifer</i>
A115	<i>Paracoccus</i>
A116	<i>Erythrobacter</i>
A117	<i>Pseudohalocynthiibacter</i>
A118	<i>Ancylobacter</i>
A119	<i>Barionella</i>
A120	<i>Devosia</i>
A121	<i>Methylotinus</i>
A122	<i>Rhodobacteraceae</i>
A123	<i>Sphingosinithalassobacter</i>
A124	<i>Acetobacter</i>
A125	<i>Boseongicola</i>
A126	<i>Commensalibacter</i>
A127	<i>Fluviobacterium</i>
A128	<i>Grandibacter</i>
A129	<i>Methylobacterium</i>
A130	<i>Novosphingobium</i>
A131	<i>Pelagibaca</i>
A132	<i>Rhodovulum</i>
A133	<i>Sphingopyxis</i>
A134	<i>Beijerinckia</i>
A135	<i>Labrenzia</i>
A136	<i>Mesorhizobium</i>
A137	<i>Rhodobacter</i>
A138	<i>Starkeya</i>

**Betaproteobacteria**

Be1	<i>Mitsuaria</i>
Be2	<i>Janthinobacterium</i>
Be3	<i>Acidovorax</i>
Be4	<i>Methylthium</i>
Be5	<i>Lautropia</i>
Be6	<i>Rhodoferrax</i>
Be7	<i>Snodgrassella</i>
Be8	<i>Delftia</i>
Be9	<i>Methylotrophus</i>
Be10	<i>Neisseria</i>
Be11	<i>Paraburkholderia</i>
Be12	<i>Bordetella</i>
Be13	<i>Candidatus Nitrotoga</i>
Be14	<i>Cupriavidus</i>
Be15	<i>Herbaspirillum</i>
Be16	<i>Massilia</i>
Be17	<i>Ralstonia</i>
Be18	<i>Alicyclicphilius</i>
Be19	<i>Betaproteobacteria</i>
Be20	<i>Dechloromonas</i>
Be21	<i>Rhizobacter</i>

**delta/epsilon subdivisions**

de1	<i>Corallibacterium</i>
de2	<i>Labilithrix</i>
de3	<i>Pseudonithyromyxa</i>
de4	<i>Alnicystis</i>
de5	<i>Chondromyces</i>
de6	<i>Desulfobacula</i>
de7	<i>Geobacter</i>
de8	<i>Pelobacter</i>

**Gammaproteobacteria**

Ga1	<i>Pseudomonas</i>
Ga2	<i>Escherichia</i>
Ga3	<i>Raoultella</i>
Ga4	<i>Acinetobacter</i>
Ga5	<i>Pseudocitronomonas</i>
Ga6	<i>Cellvibrio</i>
Ga7	<i>Gilliamella</i>
Ga8	<i>Xanthomonas</i>
Ga9	<i>Microbulbifer</i>
Ga10	<i>Aleromonas</i>
Ga11	<i>Kineobacterium</i>
Ga12	<i>Methyloticaldum</i>
Ga13	<i>Haemophilus</i>
Ga14	<i>Stenotrophomonas</i>
Ga15	<i>Teredinibacter</i>
Ga16	<i>Vibrio</i>
Ga17	<i>Colwellia</i>
Ga18	<i>Dyella</i>
Ga19	<i>Sahmmonas</i>
Ga20	<i>Spiribacter</i>
Ga21	<i>Thioalkalivibrio</i>
Ga22	<i>Azoibacter</i>
Ga23	<i>Candidatus Methylospiria</i>
Ga24	<i>Granulosicoccus</i>
Ga25	<i>Halomonas</i>
Ga26	<i>Immundisolibacter</i>
Ga27	<i>Legionella</i>
Ga28	<i>Marichromatium</i>
Ga29	<i>Marinomonas</i>
Ga30	<i>Methylomonas</i>
Ga31	<i>Salinivibrio</i>
Ga32	<i>Serratia</i>
Ga33	<i>Marinobacter</i>

**Oligoflexia**

O11	<i>Halobacteriovorax</i>
O12	<i>Proteobacteria</i>

**Terrabacteria group****Actinobacteria**

Ac1	<i>Herbiconiux</i>
Ac2	<i>Sireptomyces</i>
Ac3	<i>Bifidobacterium</i>
Ac4	<i>Leifsonia</i>
Ac5	<i>Pseudarthrobacter</i>
Ac6	<i>Corynebacterium</i>
Ac7	<i>Actinomyces</i>
Ac8	<i>Curtobacterium</i>
Ac9	<i>Mycolicibacterium</i>
Ac10	<i>Micromonospora</i>
Ac11	<i>Cellulomonas</i>
Ac12	<i>Microbacterium</i>
Ac13	<i>Nocardoides</i>
Ac14	<i>Schaalia</i>
Ac15	<i>Kineococcus</i>
Ac16	<i>Tetrasphaera</i>
Ac17	<i>Brachybacterium</i>
Ac18	<i>Rothia</i>
Ac19	<i>Candidatus Planktiophila</i>
Ac20	<i>Nocardopsis</i>
Ac21	<i>Rhodococcus</i>
Ac22	<i>Acidothermus</i>
Ac23	<i>Brevibacterium</i>
Ac24	<i>Candidatus Aquiluna</i>
Ac25	<i>Geodermatophilus</i>
Ac26	<i>Marmoricola</i>
Ac27	<i>Pauljensia</i>
Ac28	<i>Plantacinospira</i>
Ac29	<i>Verrucosipora</i>
Ac30	<i>Arthrobacter</i>
Ac31	<i>Humatobacter</i>
Ac32	<i>Mycobacterium</i>
Ac33	<i>Nonomuraea</i>

**Chloroflexi**

Ch1	<i>Caldlinea</i>
Ch2	<i>Chloroflexus</i>
Ch3	<i>Roseiflexus</i>

**Cyanobacteria/Melainobacteria group**

Cy1	<i>Synechococcus</i>
Cy2	<i>Chondrocystis</i>
Cy3	<i>Cyanobium</i>
Cy4	<i>Microcoleus</i>
Cy5	<i>Nostoc</i>
Cy6	<i>Acaryochloris</i>
Cy7	<i>Anabaena</i>
Cy8	<i>Chroococcidiopsis</i>
Cy9	<i>Leptolyngbya</i>
Cy10	<i>Calothrix</i>
Cy11	<i>Oscillatoria</i>
Cy12	<i>Thermoleptolyngbya</i>

**Deinococcus-Thermus**

Dt1	<i>Meiothermus</i>
-----	--------------------

**Firmicutes**

Fi1	<i>Streptococcus</i>
Fi2	<i>Bacillus</i>
Fi3	<i>Clostridium</i>
Fi4	<i>Paenibacillus</i>
Fi5	<i>Lactobacillus</i>
Fi6	<i>Enterococcus</i>
Fi7	<i>Ruminococcus</i>
Fi8	<i>Cohnella</i>
Fi9	<i>Lachnospiraceae</i>
Fi10	[ <i>Ruminococcus</i> ]
Fi11	<i>Herbinix</i>
Fi12	<i>Roseburia</i>
Fi13	[ <i>Clostridium</i> ]
Fi14	<i>Anoxybacillus</i>
Fi15	<i>Blautia</i>
Fi16	<i>Lactococcus</i>
Fi17	[ <i>Eubacterium</i> ]
Fi18	<i>Intestinimonas</i>
Fi19	<i>Lachnospira</i>
Fi20	<i>Paenoclostridium</i>
Fi21	<i>Staphylococcus</i>
Fi22	<i>Desulfobacterium</i>
Fi23	<i>Exiguobacterium</i>
Fi24	<i>Pseudobutyrvibrio</i>
Fi25	<i>Romboutsia</i>
Fi26	<i>Selenomonas</i>
Fi27	<i>Turicibacter</i>

**Thermotogae**

Thermotogae	
Th1	<i>Fervidobacterium</i>
Th2	<i>Pseudothermotoga</i>

**Fusobacteria**

Fusobacteriia	
Fu1	<i>Leptotrichia</i>
Fu2	<i>Fusobacterium</i>

**Spirochaetes**

Spirochaetia	
Sp1	<i>Spirochaeta</i>

**Aquificae**

Aquificae	
Aq1	<i>Thermocrinis</i>

**FCB group****Bacteroidetes/Chlorobi group**

Ba1	<i>Bacteroides</i>
Ba2	<i>Paraflovitalea</i>
Ba3	<i>Flavobacterium</i>
Ba4	<i>Pedobacter</i>
Ba5	<i>Polaribacter</i>
Ba6	<i>Flavobacteriaceae</i>
Ba7	<i>Prevotella</i>
Ba8	<i>Maribacter</i>
Ba9	<i>Niastella</i>
Ba10	<i>Aquimarina</i>
Ba11	<i>Chitinophaga</i>
Ba12	<i>Psychroserpens</i>
Ba13	<i>Hymenobacter</i>
Ba14	<i>Mucilaginibacter</i>
Ba15	<i>Capnocytophaga</i>
Ba16	<i>Niabella</i>
Ba17	<i>Petrimonas</i>
Ba18	<i>Prolixibacteraceae</i>
Ba19	<i>Urechidicola</i>
Ba20	<i>Tenacibaculum</i>
Ba21	<i>Flavivirga</i>
Ba22	<i>Dokdonia</i>
Ba23	<i>Flavisolibacter</i>
Ba24	<i>Cellulophaga</i>
Ba25	<i>Cyclobacterium</i>
Ba26	<i>Formosa</i>
Ba27	<i>Cytophaga</i>
Ba28	<i>Fuzhyella</i>
Ba29	<i>Pseudobacter</i>
Ba30	<i>Wonyingzhungia</i>
Ba31	<i>Zobellia</i>
Ba32	<i>Alitapes</i>
Ba33	<i>Bernardella</i>
Ba34	<i>Oceanohabitans</i>
Ba35	<i>Pontibacter</i>
Ba36	<i>Sphingobacterium</i>
Ba37	<i>Chryseobacterium</i>
Ba38	<i>Cloacibacterium</i>
Ba39	<i>Draconibacterium</i>
Ba40	<i>Imiticia</i>
Ba41	<i>Filimonas</i>
Ba42	<i>Kordia</i>
Ba43	<i>Labilibaculum</i>
Ba44	<i>Mariniflexile</i>
Ba45	<i>Murticola</i>
Ba46	<i>Odaribacter</i>
Ba47	<i>Tannerella</i>
Ba48	<i>Winogradskyella</i>
Ba49	<i>Algoriphagus</i>
Ba50	<i>Arcticibacterium</i>
Ba51	<i>Gramella</i>
Ba52	<i>Ignavibacteriae</i>
Ba53	<i>Pseudarcicella</i>
Ba54	<i>Ramella</i>

**Gemmatimonadetes**

Gem1	<i>Gemmatimonas</i>
------	---------------------

**PVC group****Planctomycetes**

P11	<i>Mariniblastus</i>
P12	<i>Rosemaritima</i>
P13	<i>Planctomycetes</i>
P14	<i>Gemmata</i>
P15	<i>Rhytopirellula</i>
P16	<i>Isosphaera</i>
P17	<i>Limnoglobus</i>

**Verrucomicrobia**

Ve1	<i>Corallimargarita</i>
Ve2	<i>Opitutus</i>
Ve3	<i>Verrucomicrobia</i>
Ve4	<i>Lacmispheara</i>
Ve5	<i>Akkermansia</i>

**Acidobacteria****Acidobacteriia**

Ad1	<i>Candidatus Solibacter</i>
Ad2	<i>Acidiscaria</i>
Ad3	<i>Edaphobacter</i>
Ad4	<i>Granulicella</i>

**Blastocatellia**

B11	<i>Chloracidobacterium</i>
-----	----------------------------

**Vicinamibacteria**

V11	<i>Lentitalea</i>
-----	-------------------

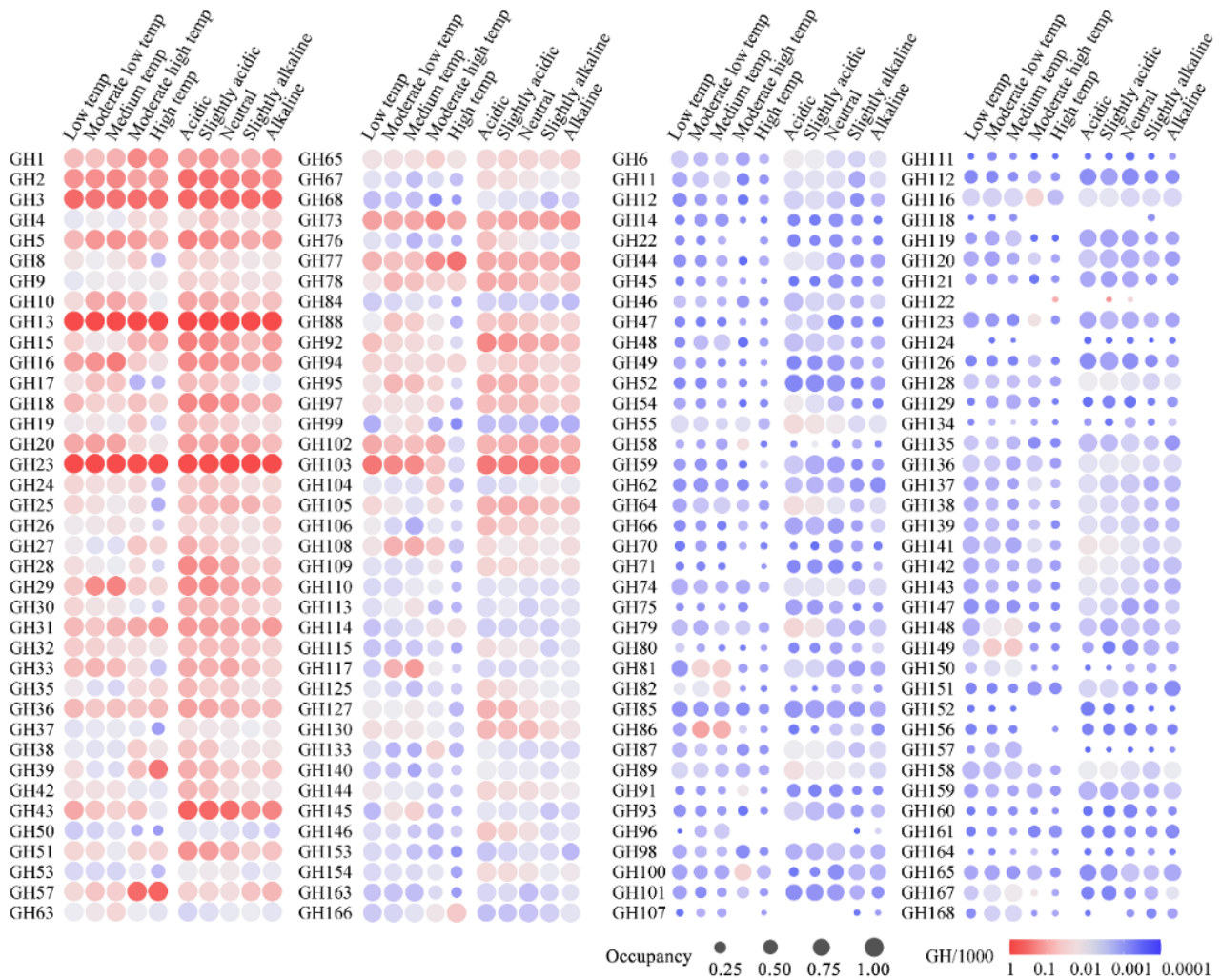
**Euryarchaeota****Stenosarchaea group**

St1	<i>Halorubrum</i>
St2	<i>Halorhabdus</i>
St3	<i>Haloarcula</i>

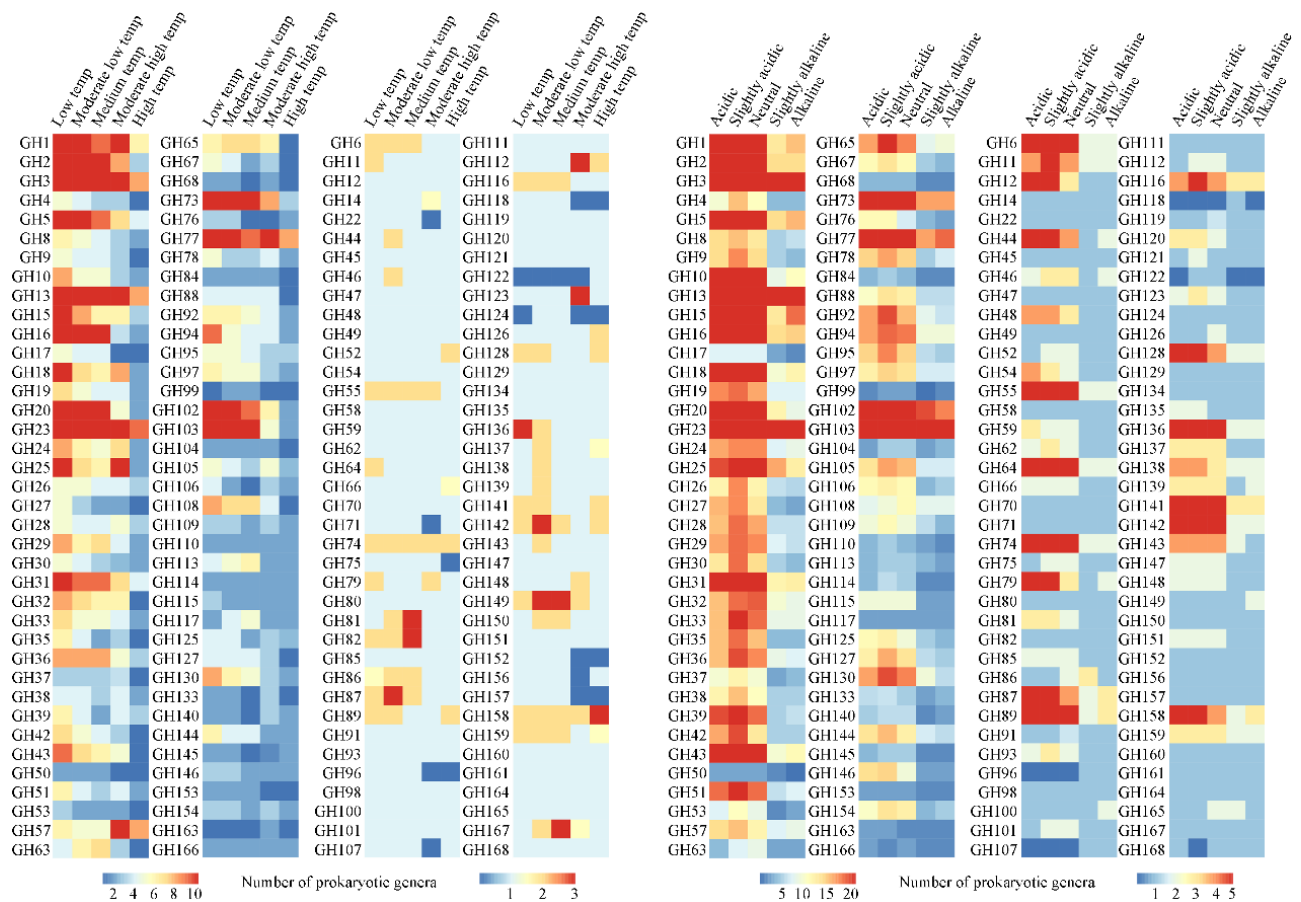
**Nitrospirae**

Nitrospirae	
Ni1	<i>Nitrospira</i>

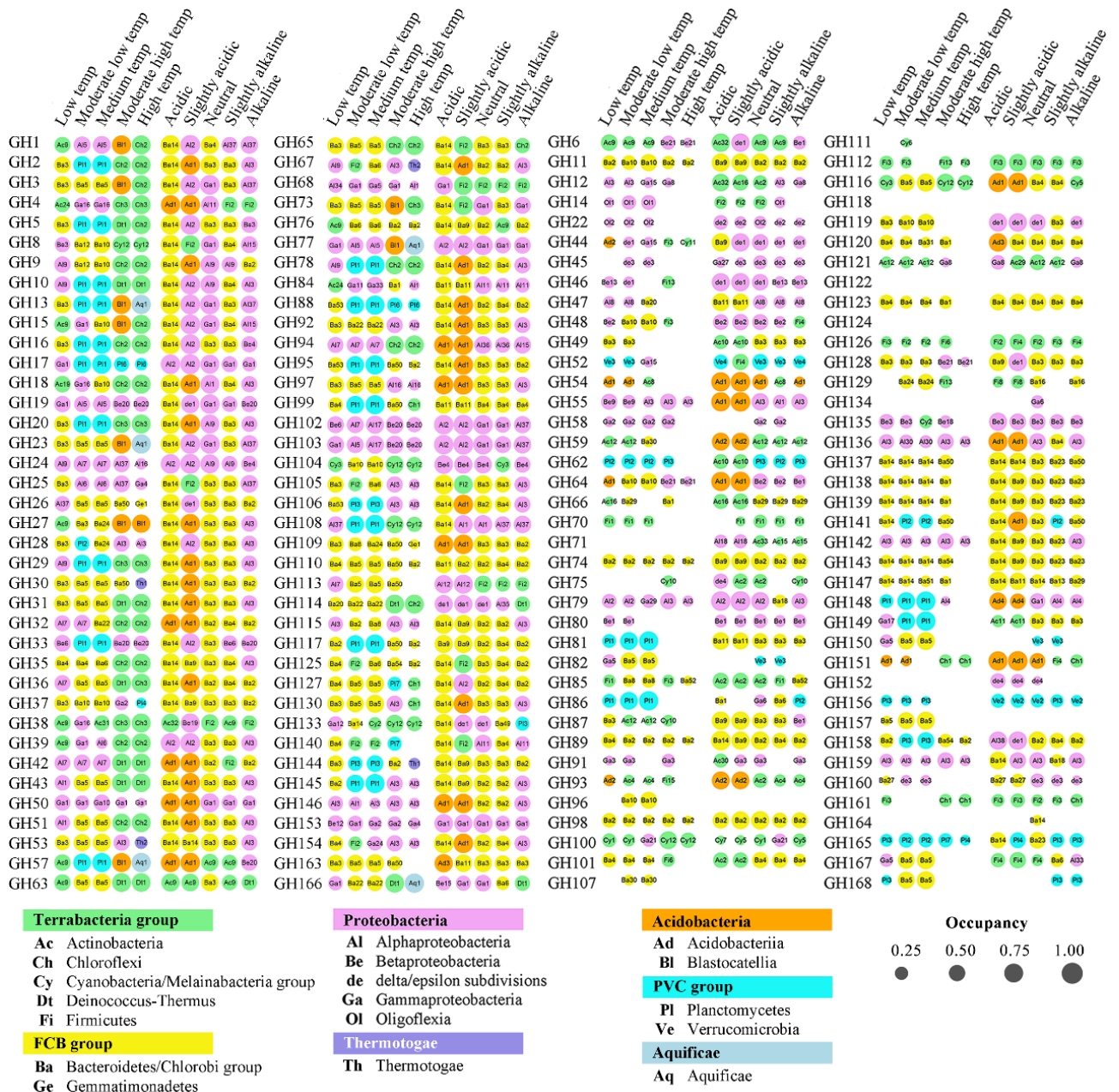
**Figure S9. Specific information on the most common prokaryotic genera.** Based on NCBI Taxonomy statistics, different colors represent prokaryotic genera from different taxonomic groups, and subgroups are shown in bold. The genus codes correspond to those shown in Figure 7, Figure S8, and Figure S12.



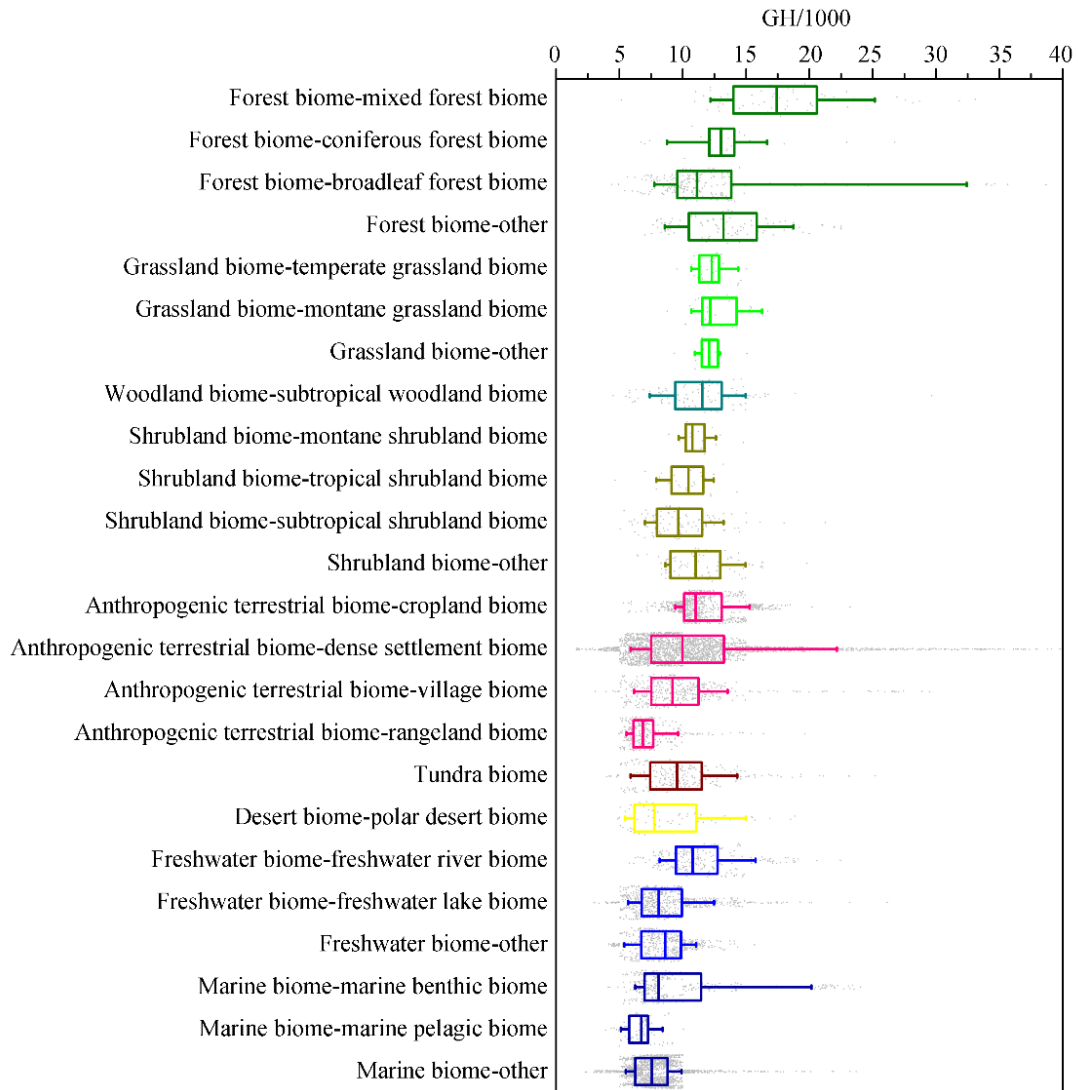
**Figure S10. Environmental temperature and pH affect the occupancy and abundance of GH families in prokaryotic communities.** Statistical analyses were carried out based on 2381 EMP samples with recorded temperature information and 1183 samples with recorded pH values. The EMP samples were categorized into 5 groups according to environmental temperature and pH: low temp ( $\leq 10^{\circ}\text{C}$ ), moderate low temp ( $> 10^{\circ}\text{C}$  and  $\leq 20^{\circ}\text{C}$ ), medium temp ( $> 20^{\circ}\text{C}$  and  $< 30^{\circ}\text{C}$ ), moderate high temp ( $\geq 30^{\circ}\text{C}$  and  $< 45^{\circ}\text{C}$ ), and high temp ( $\geq 45^{\circ}\text{C}$ ); acidic ( $\leq 5$ ), slightly acidic ( $> 5$  and  $\leq 6.5$ ), neutral ( $> 6.5$  and  $< 7.5$ ), slightly alkaline ( $\geq 7.5$  and  $< 9$ ) and alkaline ( $\geq 9$ ). The occupancy of each GH family in samples within different temperature or pH groups is represented by the circle size. Colors represent the median abundance of the GH family across all EMP samples where the family was detected within the same group, with values shown as GH/1000, increasing from blue to red.



**Figure S11. Diverse taxonomic sources of GH families in prokaryotic communities under different environmental temperature and pH conditions.** The EMP samples were categorized into 5 groups according to environmental temperature and pH. The colors in the figure indicate the median number of source prokaryotic genera of GH families in all EMP samples from the same environment, with increasing values from blue to red.



**Figure S12. Most common prokaryotic genera of GH families in various environments.** The most common prokaryotic genus was defined as the genus that appeared in the most samples and had an occupancy of at least 0.05 in all samples. The EMP samples were divided into 5 groups according to environmental temperatures and pH. Different colors represent different taxonomic groups. The circle size indicates the occupancy of the genus in all samples from a single environment. Specific information on the prokaryotic genera can be found in Figure S9.



**Figure S13. Variations in the abundance of GHs encoded by prokaryotic communities across different ENVO environments.** The GH abundance is expressed as the GH/1000 value, i.e., the number of GH genes per thousand prokaryotic genes in a sample. Each grey point represents an EMP sample. For the boxplots, the middle line indicates the median, the box represents the 25th-75th percentile, and the whisker indicates the 10th-90th percentile of observations.

**Data S1. Data on GH abundance in EMP samples.**

**Data S2. Abundance of 152 GH families in EMP samples.**

**Data S3. Correlations of different GH families in various environments.**

**Data S4. Information on taxonomic sources of GH families in prokaryotic communities.**