

## 1 **Supplementary information**

# 3 **Stochastic process determines the spatial variations in microbial** 4 **community inhabiting terrestrial mud volcanoes across the Eurasian** 5 **continent**

7 Tzu-Hsuan Tu<sup>1,2,3</sup>, Li-Ling Chen<sup>2</sup>, Yi-Ping Chiu<sup>3</sup>, Li-Hung Lin<sup>3,4</sup>, Li-Wei Wu<sup>5</sup>, Francesco  
8 Italiano<sup>6</sup>, J. Bruce H. Shyu<sup>3</sup>, Seyed Naser Raisossadat<sup>7,8</sup>, and Pei-Ling Wang<sup>2,4\*</sup>

- 10 1. Department of Oceanography, National Sun Yat-sen University, Kaohsiung, Taiwan
- 11 2. Institute of Oceanography, National Taiwan University, Taipei, Taiwan
- 12 3. Department of Geosciences, National Taiwan University, Taipei, Taiwan
- 13 4. Research Center for Future Earth, National Taiwan University, Taipei, Taiwan
- 14 5. Department of Life Science, Tunghai University, Taichung, Taiwan
- 15 6. National Institute of Geophysics and Volcanology, Palermo, Italy
- 16 7. Department of Geology, University of Birjand, Birjand, Iran
- 17 8. Earth Science Research Group, University of Birjand, Birjand, Iran

## 21 **Material and methods**

### 22 **Sampling sites and procedures**

23 Muddy fluids from bubbling pools and a total of 16 sediment cores from the adjacent mud  
24 platform were retrieved from 15 MVs across the Eurasian continent during 2011 to 2017 (Fig. 1;  
25 plotted using the ggmap package (Kahle & Wickham, 2013) in R; Table S1) for geochemical and  
26 molecular analyses. In brief, bubbling fluids and cores were collected using sterilized cups and  
27 liners, respectively. The lengths of the cores ranged between 20 and 160 cm. Samples were  
28 transported to the nearby laboratory or accommodation within 5 hours after retrieval. The cores  
29 were immediately sectioned at an interval of 1.5 to 5 cm (Table S1) with the average depth of  
30 individual sectioned intervals as the representative depth. For gas geochemistry, we preserved 6  
31 mL of sediments in a 36-mL serum bottle with 10 mL of 1 M NaOH, and sealed with a butyl rubber  
32 and an aluminum ring. Following the gas sampling, 3 mL of sediments were collected in a 15-mL  
33 centrifuge tube for the determination of water content. Samples for pore water content were subject  
34 to freeze drying. The weight difference was used to calculate the water weight content or porosity  
35 assuming the density of dry sediment was  $2.5\text{g cm}^{-3}$  and the pore space was completely saturated  
36 with pore water. For aqueous geochemistry, the remaining sediments were placed in a 50-mL  
37 centrifuge tube and centrifuged at  $8,200 \times g$  for 15 minutes to collect pore water. The obtained  
38 pore water was decanted from the centrifuge tube,  $0.22\text{-}\mu\text{m}$ -filtered using syringe filters and split  
39 it into two fractions with one for ion chromatographic analyses of anion abundances and the other  
40 for dissolved inorganic carbon (DIC). For molecular analyses, sediments were placed in a 50-mL  
41 centrifuge tube. Upon arriving at the laboratory, anion and DNA samples were stored in a  $4^{\circ}\text{C}$   
42 refrigerator and a  $-80^{\circ}\text{C}$  freezer, respectively, until further analysis.

## 43 **Geochemical analyses**

44 Concentrations of gaseous hydrocarbon compounds in head space and dissolved inorganic  
45 carbon (DIC) in pore water were analyzed using a 6890N gas chromatograph (GC; Agilent  
46 Technologies, Santa Clara, CA, USA) equipped with a Porapak Q packed column (3 m) in line  
47 with a flame ionization detector and a thermal conductivity detector (6890N, Agilent Taiwan,  
48 Taipei, Taiwan). The measured partial pressure of a specific gas compound was used to calculate  
49 the equilibrium dissolved concentration with the Henry's law constant (Wiesenburg and Guinasso,  
50 1979). The total moles in headspace and dissolved phase were summed up and normalized to the  
51 volume of pore water in order to obtain the dissolved concentration.

52 The  $\delta^{13}\text{C}$  values of DIC was determined by acidifying pore water with 85% phosphoric acid  
53 for the production of  $\text{CO}_2$  from DIC. Carbon isotope compositions of both methane and DIC were  
54 measured using a MAT253 isotope ratio mass spectrometer (IRMS) connected with a GC Isolink  
55 (Thermo Fisher Scientific, Waltham, MA, USA). The isotopic compositions were reported as the  
56  $\delta$  notation ( $\delta$  value =  $(R_{\text{sample}}/R_{\text{standard}} - 1) \times 1000\text{‰}$ ), where R is the ratio of  $^{13}\text{C}$  to  $^{12}\text{C}$ , and the  
57 standard is Pee Dee Belemnite (PDB) for carbon isotopes, was used in reporting the data.

58 Two anions in pore water, chloride, and sulfate were analyzed using an ICS-3000 ion  
59 chromatograph (Thermo Fisher Scientific, Waltham, MA, USA). Concentrations of particulate  
60 total organic carbon (TOC), total inorganic carbon (TIC), total nitrogen (TN), and total sulfur (TS)  
61 were determined by an elemental analyzer (MICROcube, Elementar, Germany). The uncertainties  
62 for aqueous and gas geochemistry, elemental abundance, and  $\delta^{13}\text{C}$  value are  $\pm 2\%$ ,  $\pm 5\%$ ,  $\pm 2\%$ , and  
63  $\pm 0.3\%$ , respectively. The detectable limits for anions with the consideration of dilution were 10  
64 ppm.

65

## 66 **Microbial community compositions**

### 67 ***DNA extraction and amplification of 16S rRNA gene***

68 Crude DNA for 16S rRNA gene analyses was extracted from 2 to 5 g of fluids/sediments  
69 using the PowerSoil DNA Isolation Kit (Qiagen, Hilden, Germany). Bubbling fluids (if available)  
70 and sediments distributed across geochemical transition were selected for DNA extraction. These  
71 samples are representative of communities inhabiting the subsurface source region (represented  
72 by bubbling fluids) or subjected to the redox gradient developed after the sediment deposition  
73 (represented by cored sediments in adjacent mud platform). A total of 136 DNA extracts were  
74 obtained and stored at  $-80\text{ °C}$  for subsequent analyses. Polymerase chain reaction (PCR)  
75 amplification was applied to amplify the V4 hypervariable region of 16S rRNA genes using the  
76 primers F515 (5'-GTG CCA GCM GCC GCG GTA A-3') and R806 (5'-CCC GTC AAT TCM  
77 TTT RAG T-3') that target both bacterial and archaeal communities (Kozich et al., 2013). Sample  
78 specific barcodes and Illumina-specific adapters were appended with both forward and reverse  
79 primers. The ingredients of each PCR mixture contained 1.1–1.5 ng of purified genomic DNA, 1  
80 U of ExTaq polymerase (TaKaRa Bio, Japan), 0.2 mM of dNTPs, 0.2  $\mu\text{M}$  of each primer, and 2.5  $\mu\text{L}$   
81 of  $10 \times$  PCR buffer in a total volume of 25  $\mu\text{L}$ . The program of thermal cycling involved a  
82 denaturation step at  $94\text{ °C}$  for 3 minutes followed by 30 cycles of denaturation at  $94\text{ °C}$  for 45  
83 seconds, annealing at  $55\text{ °C}$  for 45 seconds, extension at  $72\text{ °C}$  for 90 seconds, and a final extension  
84 step at  $72\text{ °C}$  for 10 minutes. The products of three PCR runs for individual samples were pooled,  
85 analyzed by gel electrophoresis for size verification ( $\sim 350$  bp), and purified using the DNA Clean

86 and Concentrator Kit (Zymo Research, United States). Amplicons from different samples were  
87 pooled in equal quantities sufficient for sequencing on an Illumina MiSeq platform (Illumina,  
88 United States).

89

## 90 ***Sequence processing***

91 Sequences of 16S rRNA gene amplicons were analyzed using the Mothur and QIIME2 (Schloss  
92 et al., 2009; Bolyen et al., 2018). Specifically, sequences for individual samples were binned in  
93 accordance with the barcode sequences. To minimize the effects of random-sequencing errors,  
94 reads that had two or more mismatches to the barcode sequences were discarded. The split raw  
95 FASTQ data were processed with the DADA2 (Callahan et al., 2017) implemented in the QIIME2  
96 (version 2018.8; <http://qiime2.org/>) (Bolyen et al., 2018; Caporaso et al., 2010) to calculate the  
97 amplicon sequence variants (ASVs) in each sample. After removing the sequencing adapters, the  
98 first 31 nucleotides of primer sequences were trimmed off. Due to the decrease of quality at the  
99 end of each read, forward and reverse sequences were truncated to a length of 220 and 200 base  
100 pairs, respectively, to obtain individual sequences with a quality score greater than 20.  
101 Subsequently, denoised reads were assembled to full sequences, aligned, and taxonomically  
102 assigned against the Silva v.132 reference set using Mothur. Sequences identified as chloroplasts  
103 and mitochondria were removed. The obtained sequences were deposited in GenBank with  
104 accession number PRJNA560274.

105

## 106 **Statistics**

### 107 ***Microbial community analyses***

108 Samples were rarefied to 9,413 sequences per sample through 100 sequence random re-  
109 sampling (without replacement) of the original amplicon sequence variants (ASV) table to account  
110 for the difference in sequencing depth for the calculation of alpha diversity indices such as  
111 observed ASV richness, Chao1 and Shannon indices (Hill, 1973; Chao et al., 1984). Based on the  
112 rarefied dataset, alpha diversity indices, such as observed ASV richness, Chao1 and Shannon  
113 indices were computed. For the beta diversity calculation, the entire ASV table was used and  
114 normalized using the function cumNorm from the R package metagenomeSeq (Paulson et al.,  
115 2013). A cumulative-sum scaling method was used to calculate the scaling factors equal to the sum  
116 of counts up to a particular quantile in order to normalize the read counts with uneven sequencing  
117 depth (Paulson et al., 2013). The dissimilarity matrix between samples was computed using the  
118 Bray-Curtis method (Bray et al., 1957) and visualized through the ordination of non-metric  
119 multidimensional scaling (NMDS) and constrained correspondence analysis (CCA). The  
120 significance of environmental variables relative to the CCA ordinations was computed using  
121 “envfit” and 999 permutations. All statistical analyses were performed in R using the packages  
122 *vegan*, *ggplot2*, and *phyloseq*.

123

### 124 ***Estimation of habitat similarities***

125 Approaches described in Ranjard et al. (2013) and Powell et al. (2015) were adopted with  
126 some modification. An estimation of habitat similarities was calculated from the Euclidean  
127 distances between paired 126 samples with the available concentrations of chloride, sulfate,

128 methane, TN, TS, TIC, and TOC. To reduce the effects of using large concentration scales, such  
129 as for chloride, environmental factors were normalized to their minimum and maximum values to  
130 scale the data to a fixed range between 0 to 1. The transformed dataset was used to evaluate habitat  
131 similarity using the following formula (Ranjard et al., 2013; Powell et al., 2015):

132

$$133 \quad E_d = \left(1 - \frac{Euc_d}{Euc_{max}}\right) \quad \text{Eq.1}$$

134

135 where  $Euc_d$  is the Euclidean distance, and  $Euc_{max}$  is the maximum distance between sites in the  
136 matrix. To test whether community similarity was significantly correlated with a variety of spatial  
137 components, non-parametric Mantel tests based on the Pearson correlation coefficient were  
138 applied with significance assessed based on 1000 Monte Carlo permutations. All statistical  
139 analyses were performed in R using the package *vegan*.

140

### 141 ***Distance decay relationships (DDR)***

142 To assess the DDR, pairwise community similarities between samples were calculated using  
143 the Sørensen-Dice index (Dice, 1945). The pairwise similarity was transformed in a logarithmic  
144 space to enhance the linear fitting using the following H formula:

145

$$146 \quad \log_{10}(S_{com}) = \log_{10}(a) + \beta \log_{10}(D) \quad \text{Eq.2}$$

147

148 where  $S_{com}$  is the pairwise similarity in community composition,  $D$  is the geographic and/or vertical  
149 distance between two samples, and  $\beta$  is the slope. Null values in the similarity/distance matrices  
150 were assumed to be 0.001 prior to the log-transformation. The distance between samples was  
151 aggregated from two categories for samples in separate cores or within the individual cores. For  
152 samples in separate cores, the distance represents the geographic distance between MVs and was  
153 calculated using the function *geodist* in the R package 'gmt'. For samples within the individual  
154 cores, the distance represents the depth difference between samples. Samples collected from the  
155 bubbling pools were regarded as the surface material (0 cm) of each sediment core. The DDR  
156 relationships were assessed for data encompassing all samples or either categories. The  
157 significance of  $\beta$  was tested by 1000 Monte Carlo permutations of the residuals under the full  
158 regression [25], and  $\beta$  was found to be significant for each sample surveyed ( $P < 0.001$ ).

159

## 160 **Results**

### 161 ***Physical and geochemical characteristics***

162 The pairwise distance between samples ranged from 2.5 to 160 cm within cores and 0.005 to  
163 9,924 km between cores (Fig. 1). Geochemical profiles of pore water showed various  
164 characteristics related to abiotic and microbial processes. Chloride concentrations varied highly  
165 among sites (ranging between 82 mM at SI02 in Myanmar and 4890 mM at GG01 in Iran) and  
166 generally decreased with increasing depth in individual cores (Figure S1). Exceptions occurred at

167 PA02, SH01, SI02, and LGH03, with substantial fluctuations in the middle or bottom part of the  
168 cores. Sulfate concentrations ranged from below the detectable level at SM22, AK03, GJ01, TA,  
169 PA01, PA02, and LGH03 to 288 mM at GG01, with most data clustering between 0.5 and 2 mM.  
170 Variations in sulfate concentration for cores with detectable sulfate were further categorized into  
171 three patterns, including depth-dependent decrease (DSZ01 and SYNH02C4) and increase (GG01,  
172 COM01, and SH01), and substantial fluctuation along the depth (AR01 and SI02).

173 Methane concentrations ranged between 0.006 mM (PA02) and 3.98 mM (SYMH02C4), with  
174 most data clustering between 0.2 and 1 mM (Figure S2). Methane concentrations either increased  
175 (DSZ01, SM22, GJ01, TA, SYNH02C4, and LGH03) or decreased (GG01, AR01, COM01, SH01,  
176 and SI02) with increasing depth. The  $\delta^{13}\text{C}$  values of methane clustered between -58‰ and -35‰  
177 and exhibited a trend opposite to that of methane concentration. The molar ratios of methane over  
178 ethane and propane ( $\text{C1 (methane)} / (\text{C2 (ethane)} + \text{C3 (propane)})$ ) were variable and ranged from  
179 22 (SI02) to approximately 1200 (AR01 and COM01; Figure S3).

180

### 181 *Community structures and compositions*

182 Analyses of 16S rRNA genes yielded a total of 4,562,760 sequences. The number of observed  
183 ASVs for individual samples ranged between 58 and 1,462 with an average value of  $449 \pm 250$   
184 when singletons (presence of one sequence for an ASV at only one depth) were included. The  
185 trends of diversity indices exhibited a similar pattern (Figure S4). The lowest values of alpha  
186 diversity indices occurred at SI02 and SH01 in Myanmar, whereas the highest values were found  
187 for AR01 in Italy. The number of observed ASVs for individual MVs ranged between 204 (SI02)  
188 and 4,203 (AR01). Accumulation curves at coarse taxonomic resolution (i.e., phylum to family)  
189 revealed the sufficiency of our sequencing effort. In contrast, at the level of ASV, the accumulation  
190 curve showed a continuously increasing trend, indicating that the diversity of the entire MV  
191 community was not fully captured (Figure S5).

## Supplementary Tables

Table S1. Core length and coordinate of sampled terrestrial mud volcanoes across the Eurasian continent.

Sample name	Country	Core Length (cm)	Section interval (cm)	Longitude	Latitude
AR01C1	Italy	67	3	13.60000	37.37667
COM01C1	Italy	45	2.5	13.65194	37.44306
PA01C1	Italy	55	3	14.91972	37.54472
PA02C1	Italy	41	2.75	14.89028	37.57278
AK03C1	Georgia	49	2.5	45.91322	41.41953
GJ01C1	Georgia	44	2.5	45.79261	41.74531
QK01C1	Georgia	22	1.5	45.80564	41.28905
GG01C1	Iran	46	2.5	54.39608	37.11856
TA01C1	Iran	47	3	59.93306	25.46697
TA03C1	Iran	47	3	59.93306	25.46697
SM22C1	China	33	2	84.38722	44.18269
DSZ01C1	China	21	2	84.84636	44.30517
SH01C1	Myanmar	38	2	93.57119	19.36975
SI02C1	Myanmar	48	2	93.59169	19.39778
LGH03C4	Taiwan	160	5	121.20940	22.98306
SYNH02C4	Taiwan	52	2.5	120.409479	22.80313
SYNH02C11	Taiwan	20	2	120.409479	22.80313

Table S2. Multiple linear regression model for Shannon index versus significant geochemical parameters, after collinear variables were removed and Akaike information criterion (AIC) was applied. Variables were added to the model to generate the highest to lowest best fit from simple linear regression.

	Estimate	Std. Error	t-value	<i>P</i> -value	Signif. code
(Intercept)	4.02048	0.32306	12.445	< 0.0001	***
Methane	149.79125	76.02220	1.970	0.05106	.
TN	8.66657	3.13095	2.768	0.00652	**
TIC	0.43655	0.08123	5.374	< 0.0001	***

Multiple R<sup>2</sup>: 0.2067, adjusted R<sup>2</sup>: 0.1872

F-statistic: 10.6 on 3 and 122 degrees of freedom, p-value: 3.047E-06

Significance codes: '\*\*\*': 0-0.001; '\*\*': 0.001-0.01; '.': 0.05-0.1

Table S3. Simple linear regression of Shannon index against individual geochemical parameters.

	Slope	Std. Error	t-value	<i>P</i> -vare	R <sup>2</sup>
Sulfate	0.0004	0.0016	0.277	0.782	0.001
Chloride	0.0001	0.0001	1.808	0.0731	0.018
Methane	94.391	83.070	1.172	0.243	0.002
TN	2.733	3.267	0.837	0.404	0.002
TS	0.0985	0.1908	0.517	0.606	0.006
TIC	0.3577	0.0793	4.511	< 0.0001	0.134
TOC	0.3584	0.1862	1.925	0.0566	0.021

Table S4. Mantel test using Spearman's correlation (permutations = 999) of Bray-Curtis dissimilarities between all communities, and each geochemical parameter or geographic distance (km).

	<i>Rho</i> ( $\rho$ )	<i>p</i>	Signif. code
km	0.322	< 0.001	***
env	0.178	< 0.001	***
Chloride	0.454	< 0.001	***
Sulfate	0.258	< 0.001	***
Methane	0.068	0.026	**
TIC	0.255	< 0.001	***
TOC	-0.081	0.986	-
TN	-0.001	0.481	-
TS	0.143	< 0.001	***

Significance codes: '\*\*\*': 0-0.001; '-': 0.1-1

Table S5. Permutational multivariate analysis of variance (using continuous variables only) of beta diversity using Bray-Curtis dissimilarity.

	Df	SumsOfSqs	MeanSqs	F.Model	R <sup>2</sup>	P-value	Signif. codes
Chloride	1	1.504	1.50431	6.1690	0.02586	< 0.001	***
Sulfate	1	1.719	1.71883	7.0488	0.02954	< 0.001	***
Methane	1	0.509	0.50856	2.0856	0.00874	< 0.001	***
TIC	1	2.117	2.11708	8.6819	0.03639	< 0.001	***
TOC	1	1.377	1.37699	5.6469	0.02367	< 0.001	***
TN	1	1.828	1.82757	7.4947	0.03141	< 0.001	***
TS	1	1.177	1.17735	4	0.02024	< 0.001	***

Significance codes: '\*\*\*': 0-0.001

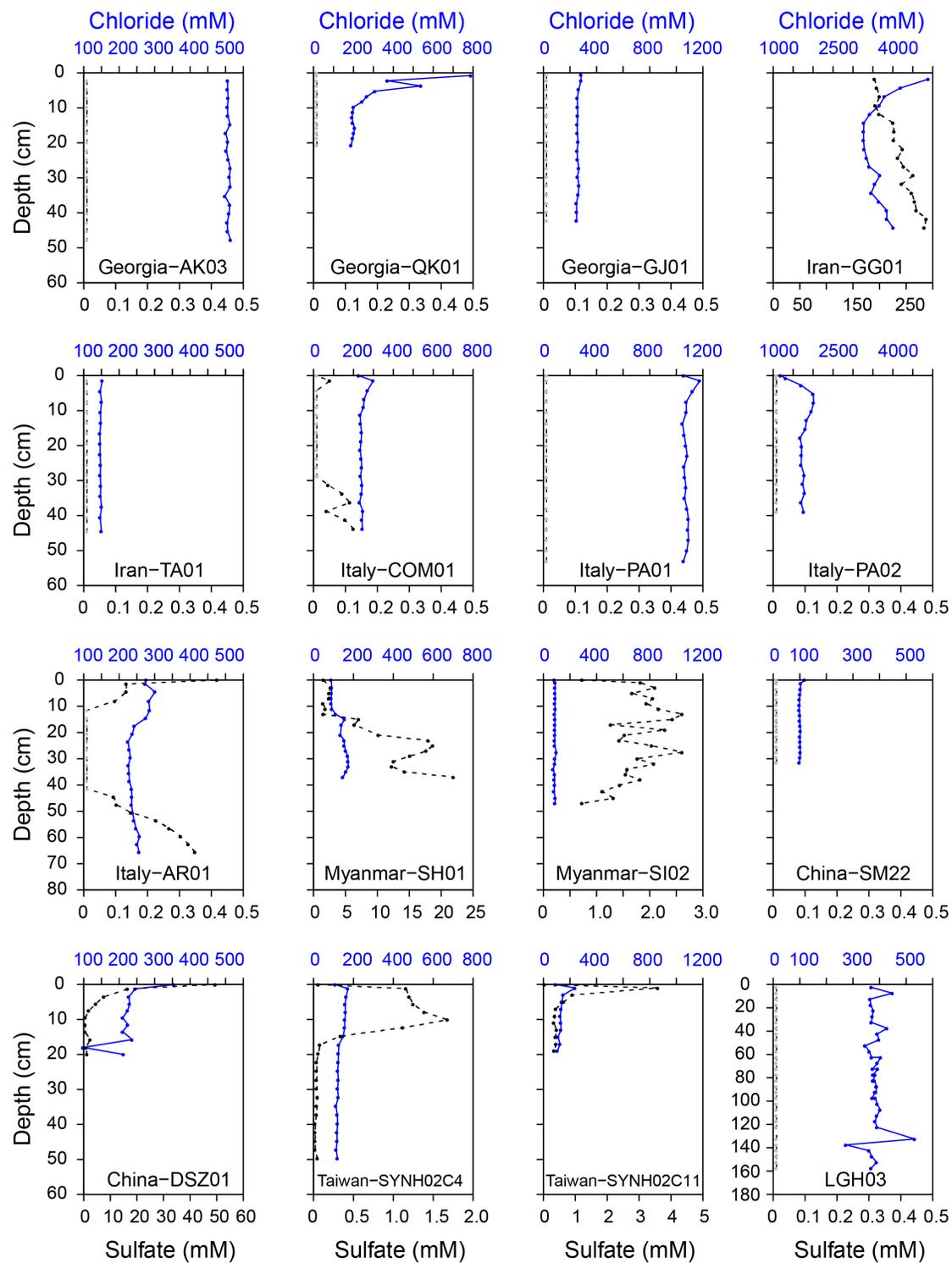
Table S6. Coefficient of variation and the Pearson's correlation coefficient (r) for chloride and sulfate concentrations.

Sample name	Chloride	Sulfate	r	Signif. codes
AR01C1	8.97%	117.09%	0.15	-
COM01C1	8.37%	122.09%	0.05	-
PA01C1	3.61%	NA	NA	NA
PA02C1	13.24%	NA	NA	NA
AK03C1	7.55%	NA	NA	NA
GJ01C1	5.61%	NA	NA	NA
QK01C1	63.13%	NA	NA	NA
GG01C1	14.84%	17.49%	-0.35	-
TA01C1	2.15%	NA	NA	NA
SM22C1	7.08%	NA	NA	NA
DSZ01C1	25.927%	186.04%	0.46	-
SH01C1	15.58%	102.65%	0.91	***
SI02C1	3.45%	29.34%	0.24	.
LGH03C4	6.56%	NA	NA	NA
SYNH02C4	15.02%	125.62%	0.88	***
SYNH02C11	28.11%	138.05%	0.96	***

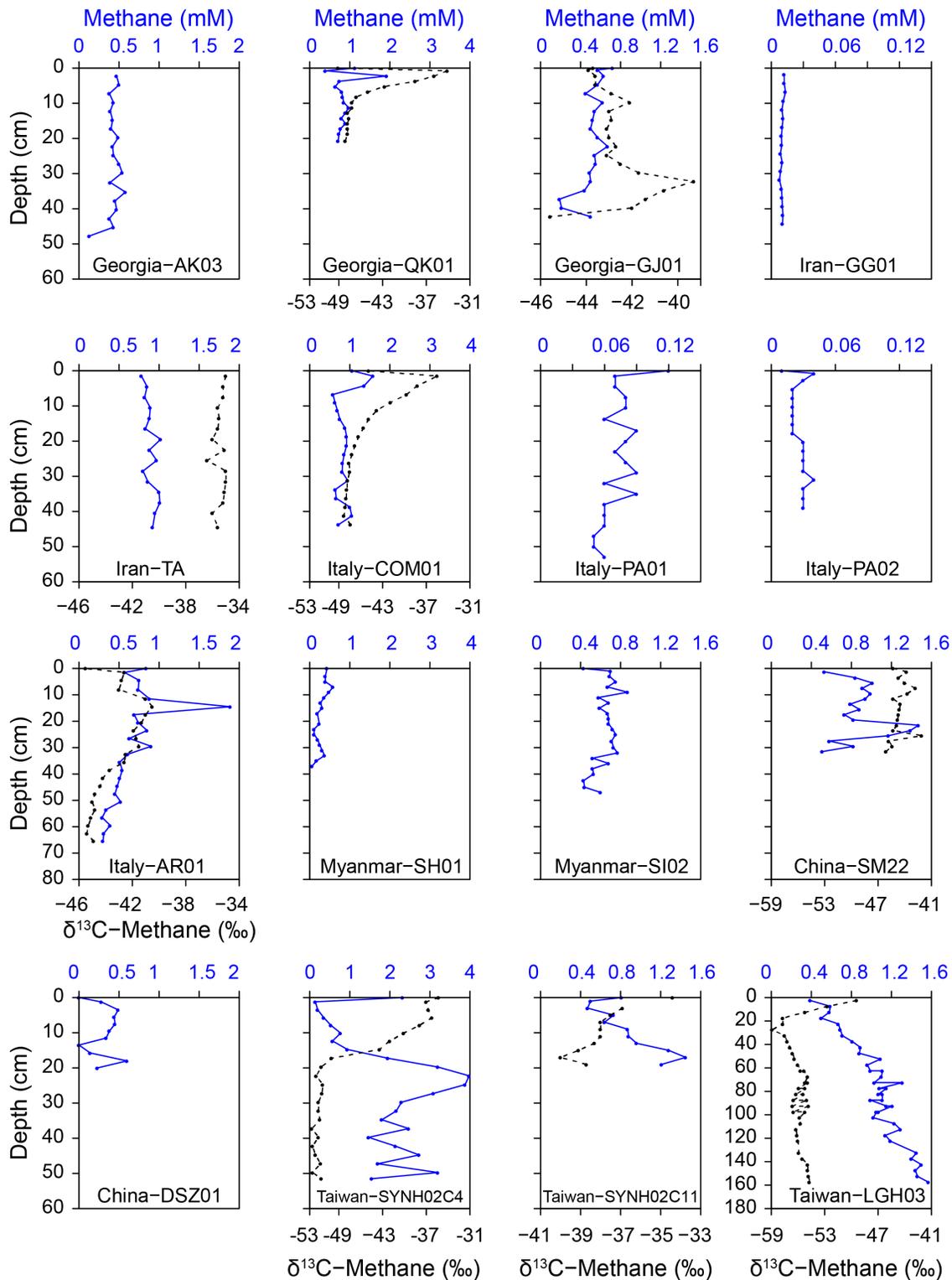
Significance codes: '\*\*\*': 0-0.001; '.': 0.05-0.1; '-': 0.1-1

NA means concentration of sulfate below the level of detection

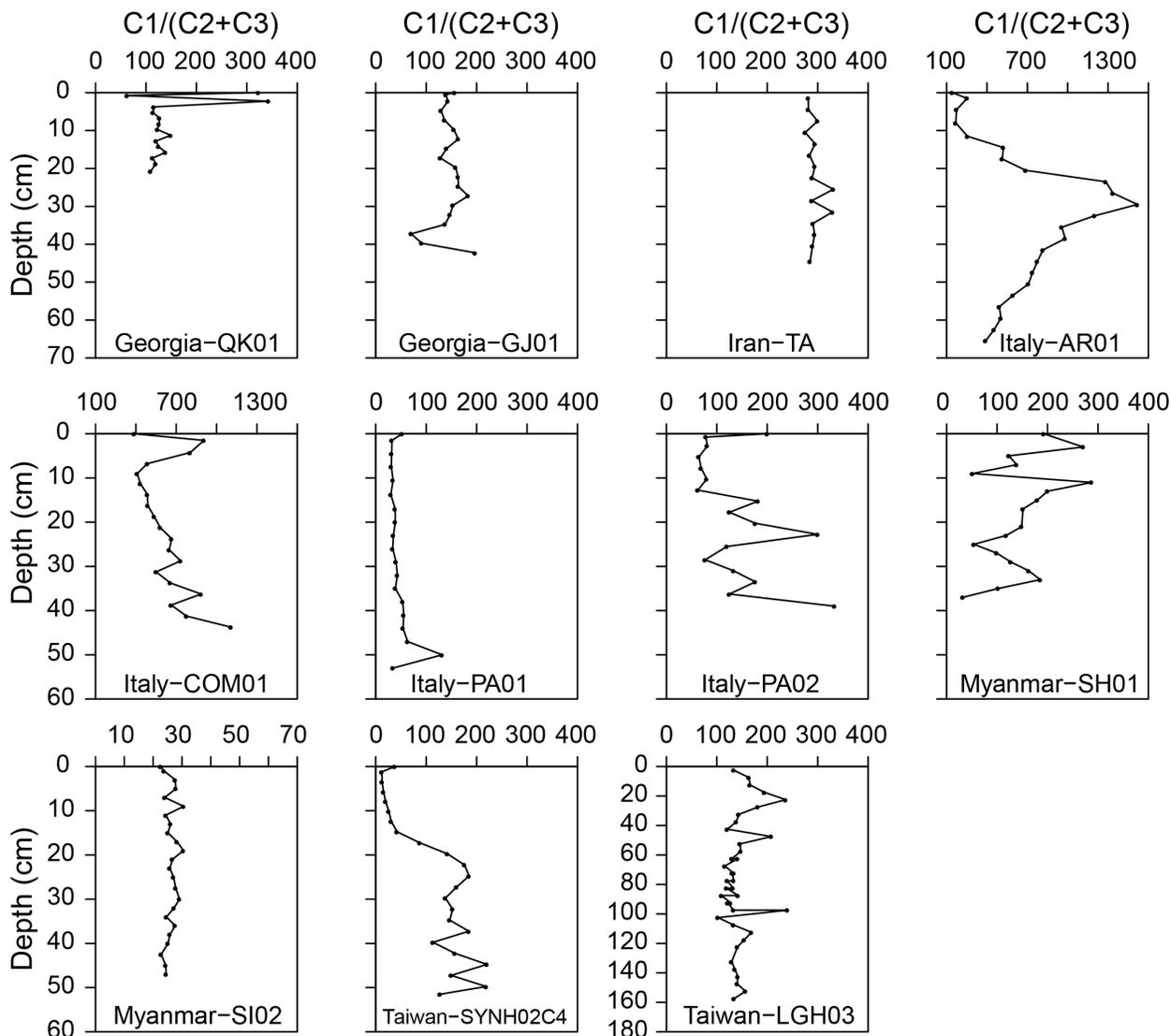
## Supplementary Figures



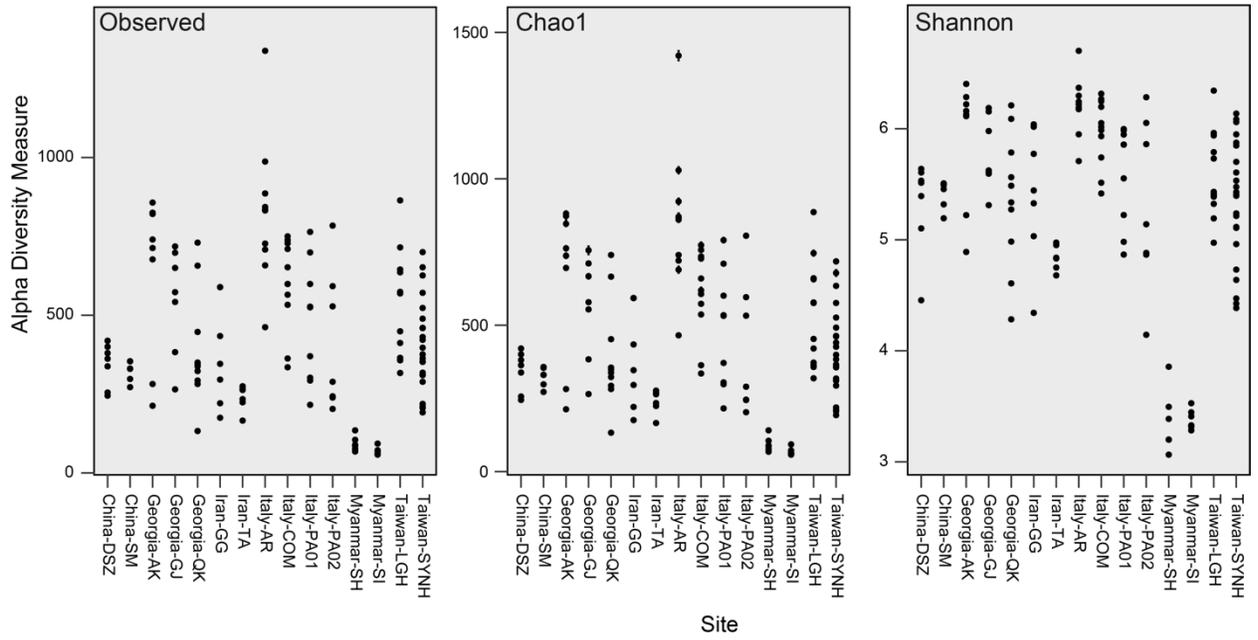
**Figure S1: Chloride (in blue) and sulfate (in black) concentration profiles. Sulfate concentrations lower than the limit of detection (0.01 mM) are shown in gray.**



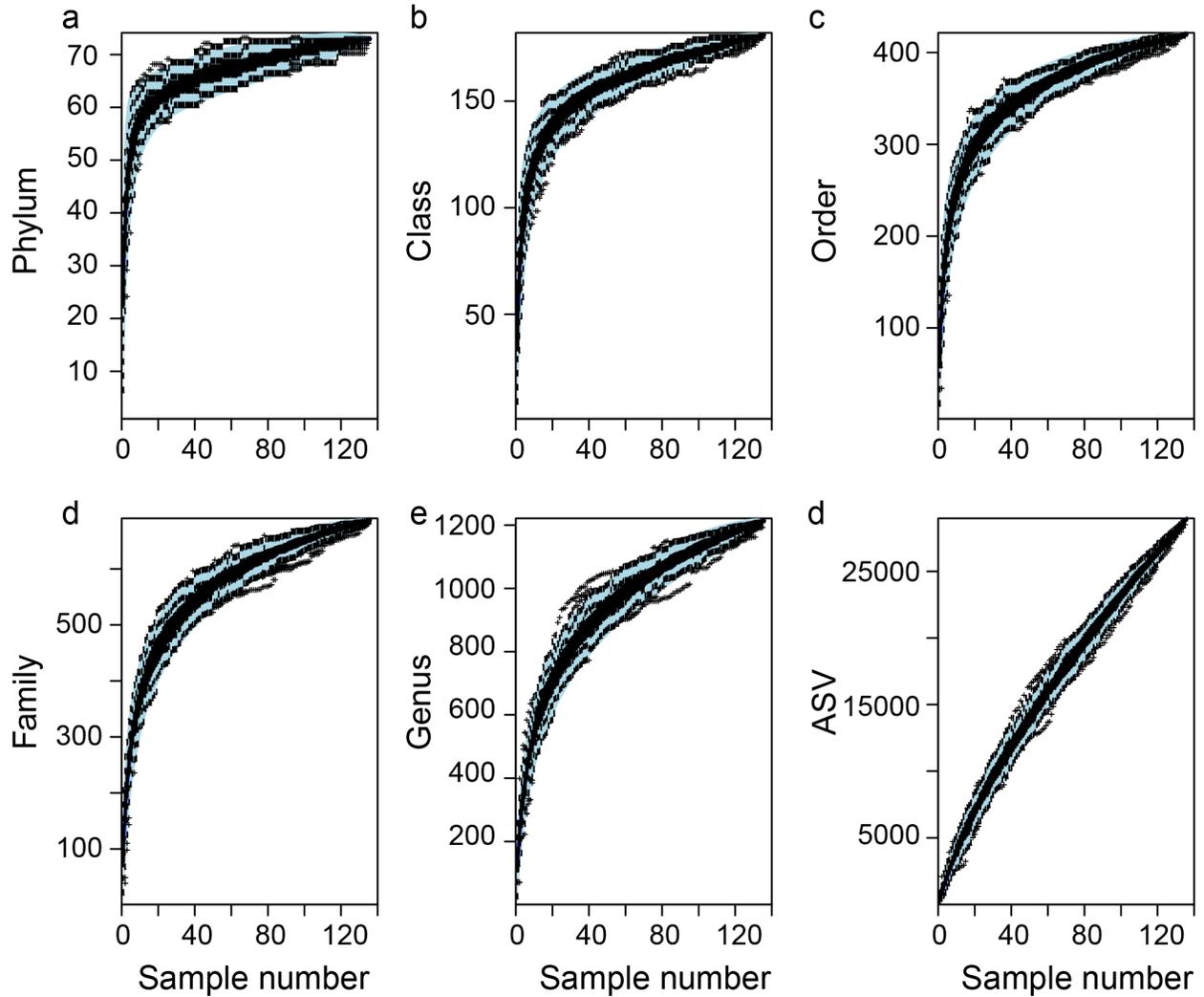
**Figure S2: Profiles of methane concentrations (in blue) and  $\delta^{13}\text{C}$  values of methane (in black). Samples collected from Georgia are not sufficient for porosity measurement. Therefore, their concentrations are normalized to the weight of wet sediments. Assuming that the weight proportion of pore water is 0.5, the concentration of 1 mmole g-1 could be converted to 200 mM.**



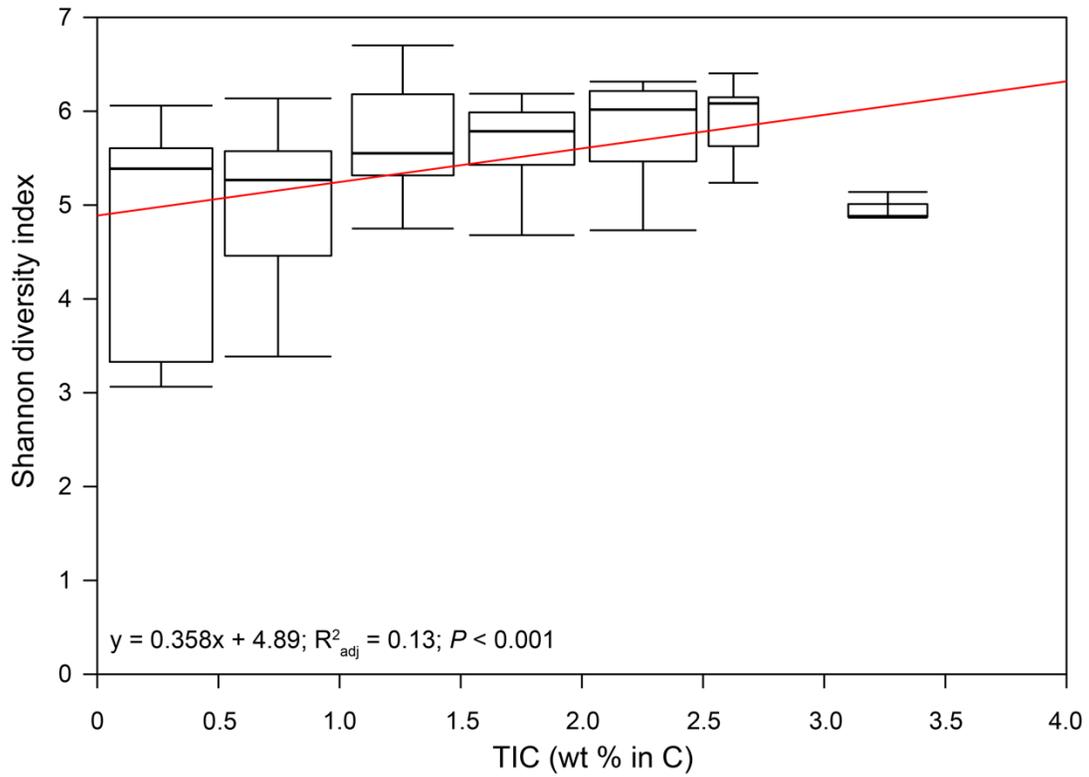
**Figure S3: Molar ratio of methane over the sum of ethane and propane ( $C1 / C2 + C3$ ).**



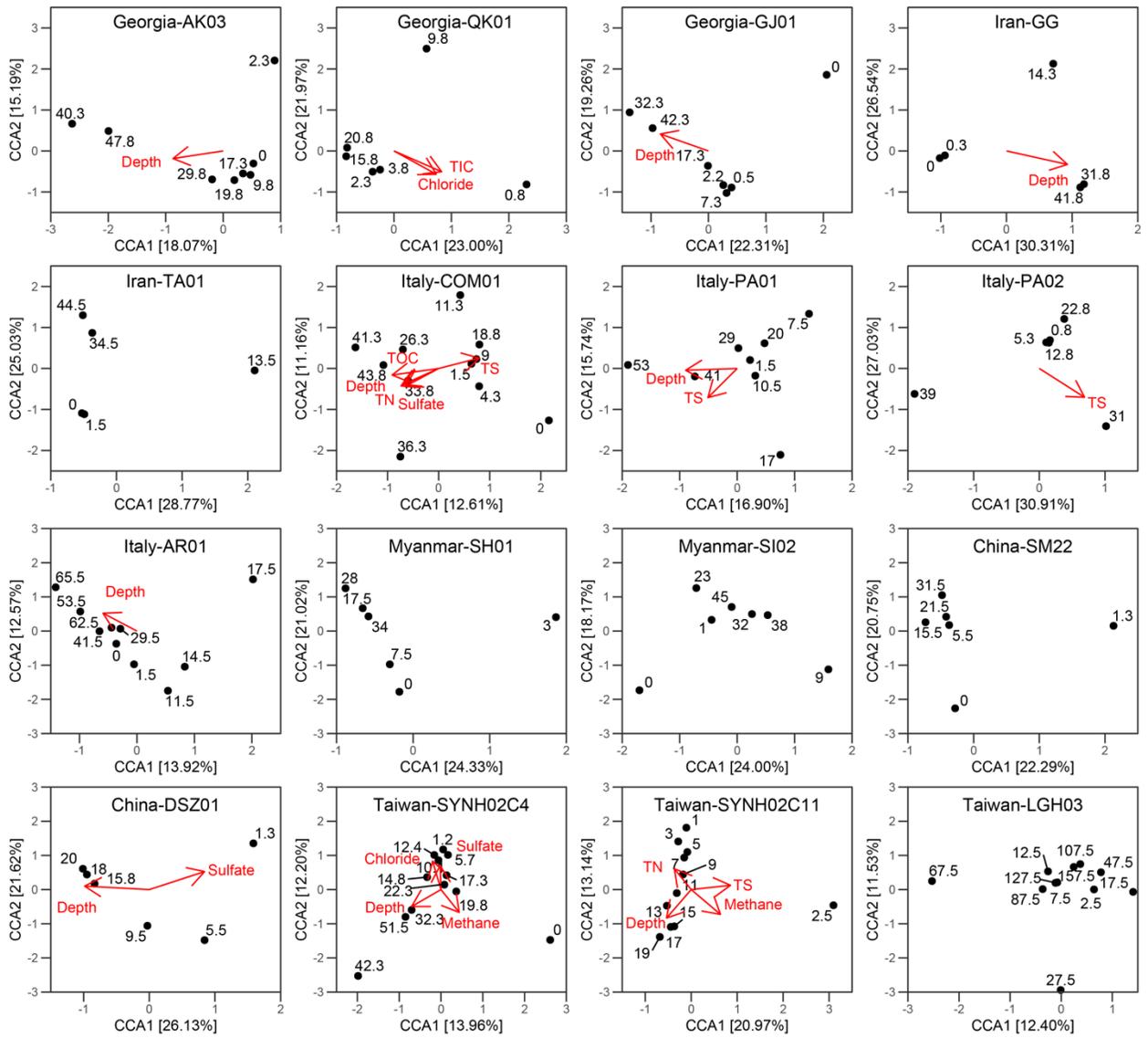
**Figure S4: Alpha diversity indices calculated based on the rarefied dataset (n=9,413). No significant relationship was found between richness and sampling depth (Spearman's  $\rho = 0.17$ ,  $P > 0.01$  for the three indices).**



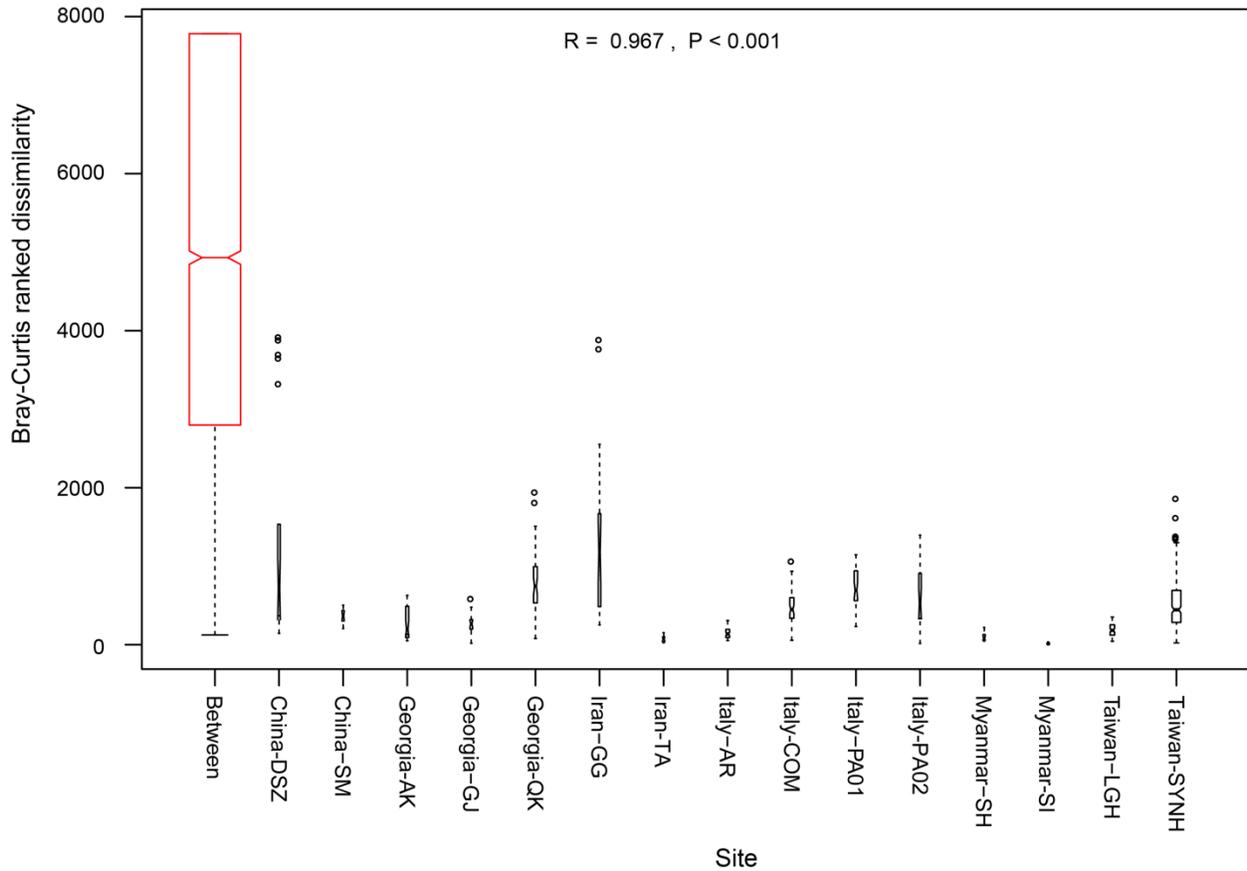
**Figure S5: Accumulation curves based on different taxonomic units: (a) phylum (b) class, (c) order, (d) family, (e) genus, and (f) ASV. Boxplots show a summary of 100 permutations calculated with random subsampling. Absolute singletons were incorporated for comparison. Blue area depicts the 95% confidence interval.**



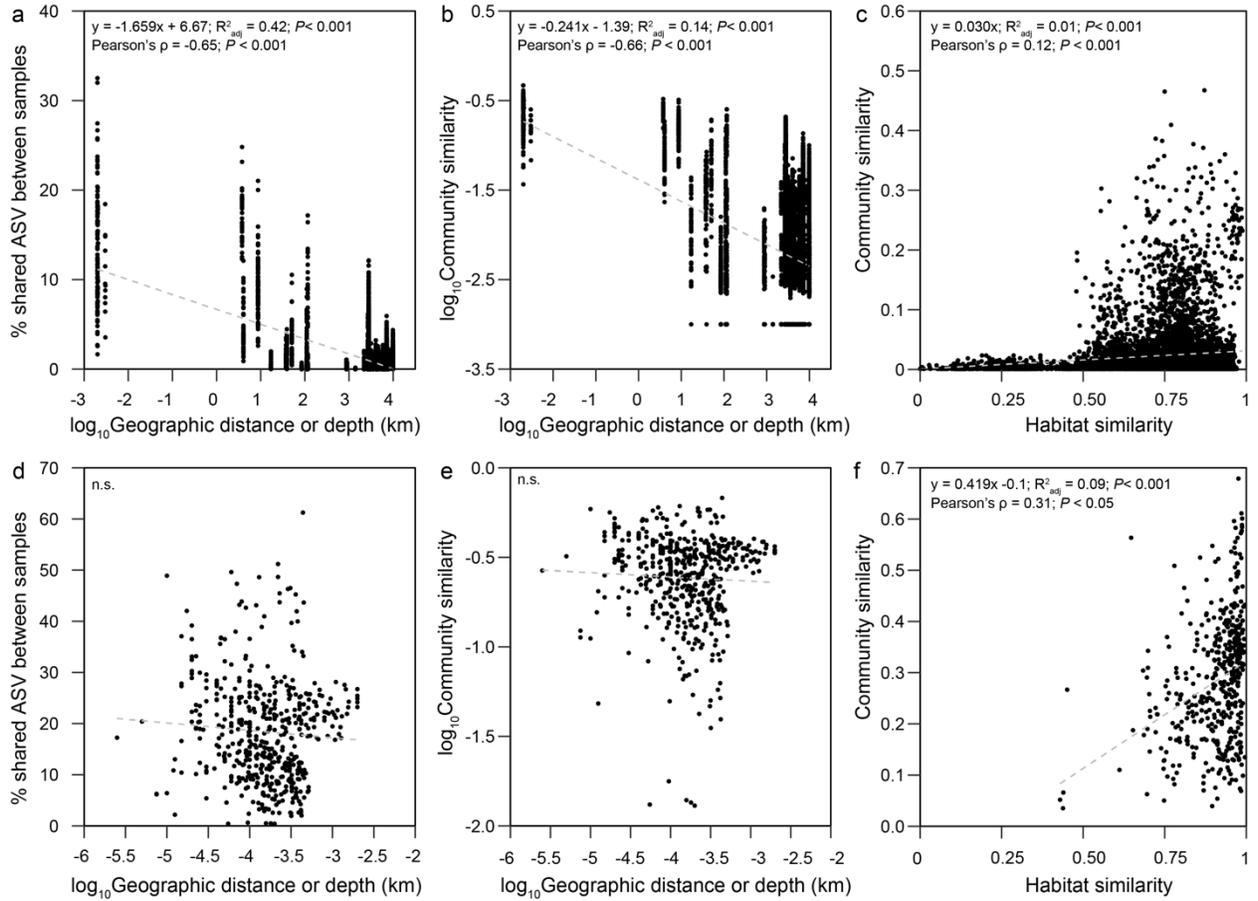
**Figure S6: Plot of Shannon diversity versus TIC. Linear regression is shown in red (n=126). Box demonstrates the interquartile range that includes the first (25%), median (50%), and third quartiles. Lower and upper whiskers are the first and third quartiles minus and plus 1.5 times interquartile range, respectively.**



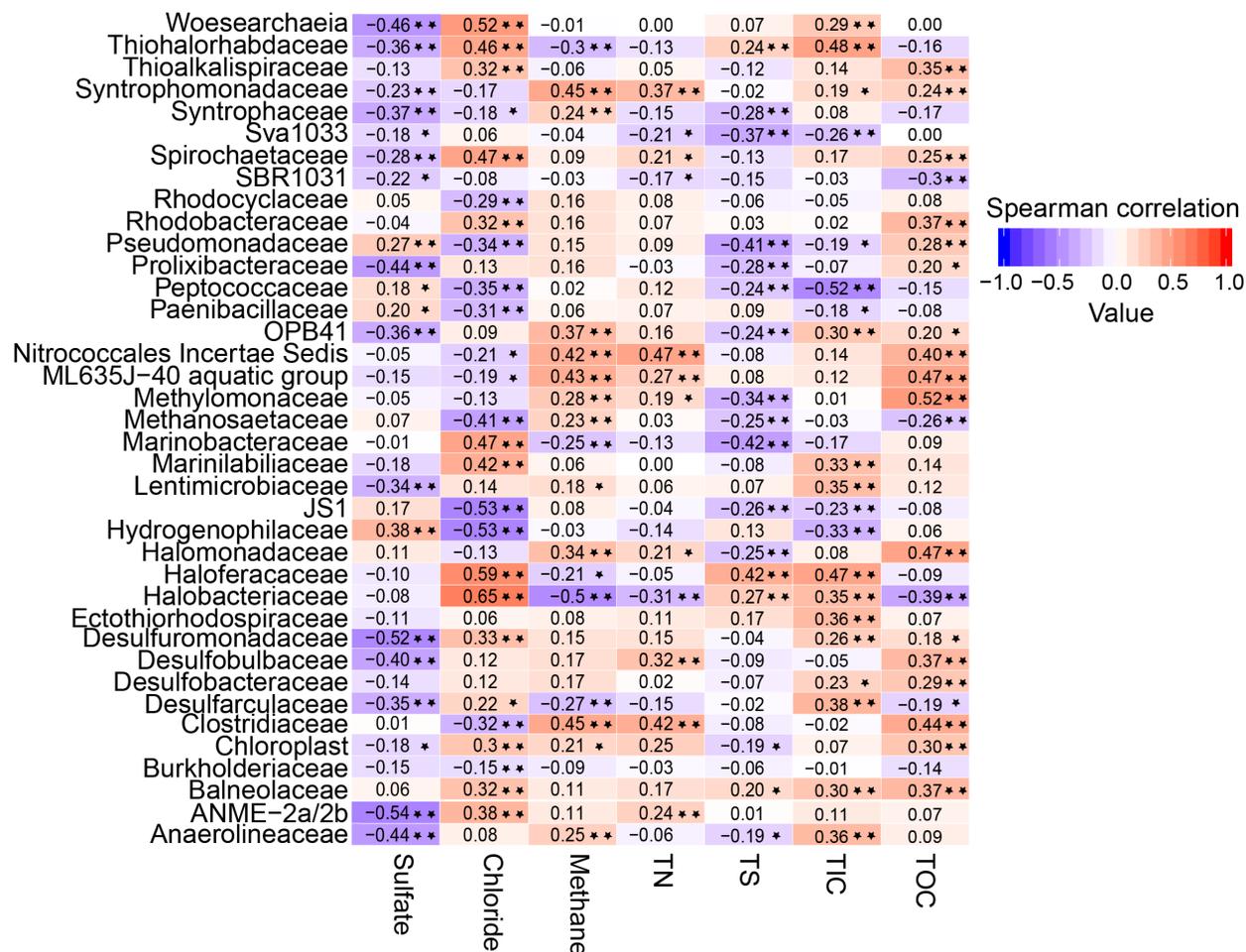
**Figure S7: Constrained correspondence analysis of community relatedness quantified by the Bray–Curtis distance with the overlay of ordination for significant environmental parameters. Numbers next to each data point indicate sampling depth (in centimeter).**



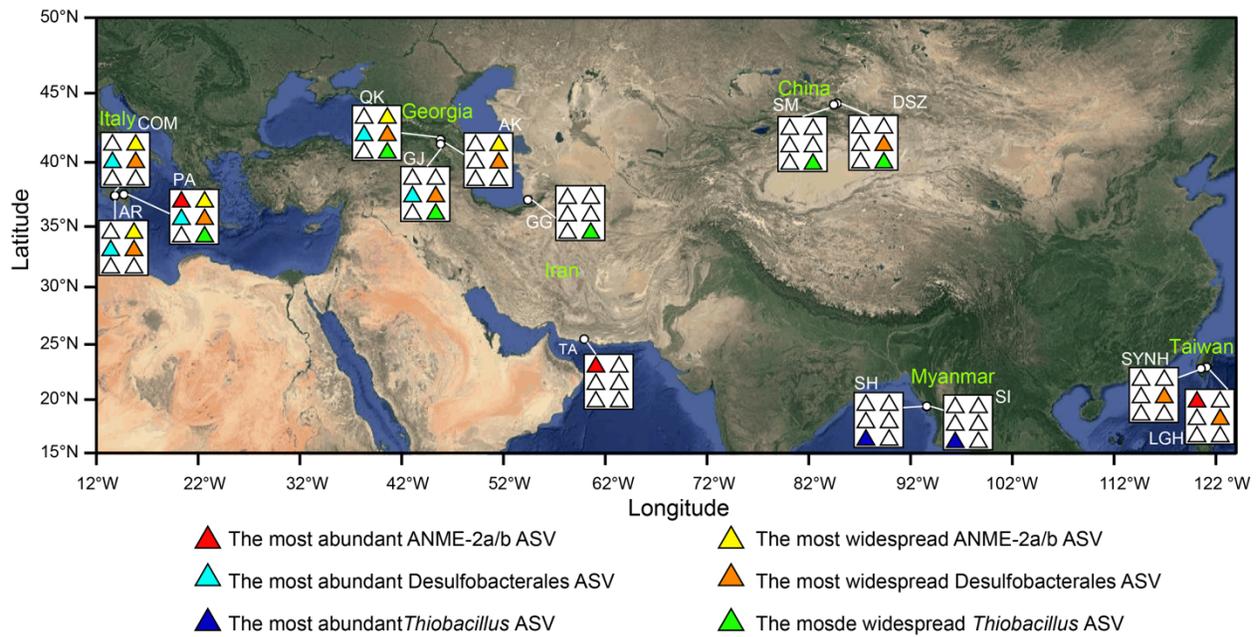
**Figure S8: Analysis of similarities (ANOSIM: |R|) for community dissimilarity between all sites (in red) and within individual sites (in black). Lower and upper whiskers are first and third quartiles minus and plus 1.5 times interquartile range, respectively.**



**Figure S9: Distance–decay and geographic patterns of microbial communities across cores (a)–(c), and within cores (d)–(f).**



**Figure S10: Correlation between concentrations of geochemical parameters and abundances of 38 major families. Color code represents the relative Spearman correlation coefficient. Major families are selected based on the top 50 most abundant families. \* and \*\* denote *P* values less than 0.01 and 0.05, respectively.**



**Figure S11: Occurrence of ASVs affiliated key taxa likely involved in methane and sulfur cycling.** Each sub-panel consists of six color codes indicating the presence or absence of six key ASVs likely involved in methane and sulfur cycling. These target ASVs include (1) the most abundant ANME-2a (in red), Desulfobacterales (in blue-green), and *Thiobacillus* (in blue) ASVs, and the most widespread ANME-2a (in yellow), Desulfobacterales (in orange), and *Thiobacillus* (in light green) ASVs. The basal map is from Google Maps © Google Maps 2021.