

1 **Supplementary Information**

2 **Title: Structural and functional differentiation of bacterial communities in post-coal**
3 **mining reclamation soils of South Africa: bioindicators of soil ecosystem restoration**

4 **Authors:** Obinna T. Ezeokoli^{1,2,3}, Cornelius C. Bezuidenhout¹, Mark S. Maboeta¹, Damase
5 Khasa^{3,4}, Rasheed A. Adeleke^{1,2*}

6

7 **Affiliations**

8 ¹Unit for Environmental Sciences and Management, North-West University, Potchefstroom
9 2520, South Africa.

10 ²Agricultural Research Council-Institute for Soil, Climate and Water, Arcadia, 0001, Pretoria,
11 South Africa

12 ³Institute for Integrative and Systems Biology, Universite Laval, Quebec City, QC G1V 0A6,
13 Canada

14 ⁴Centre for Forest Research, Faculty of Forestry, Geography and Geomatics, Universite Laval,
15 Quebec City, QC G1V 0A6, Canada

16

17 *Correspondence: Rasheed.Adeleke@nwu.ac.za.

18 **Supplementary Text (Methods)**

19

20 **Text S1: Brief method description of physicochemical analyses**

21 Briefly, pH was determined from a 1:2.5 soil-water suspension using a pre-calibrated pH meter
22 (pH 700, Eutech Instruments Pte Ltd, Singapore). Particle size distribution was determined by
23 the Bouyoucos method¹. Cations and exchangeable cations were determined from ammonium
24 acetate (1 M, pH 7) soil extracts using Inductively coupled plasma - optical emission
25 spectrometry. Anions were determined from water extracts using Ion-exchange
26 chromatography. BD determination was done after overnight drying (at 105°C) of soils
27 collected on-site with a BD sampler.

28

29 **Text S2: Bioinformatics Analyses-quality trimming**

30 Quality trimming was performed using Trimmomatic software as follows: poor quality trailing
31 and leading nucleotide positions were first trimmed from both forward and reverse reads.
32 Thereafter, reads with an average quality score (Phred, Q) less than 20 and read length less
33 than 250 bp were eliminated. For assembly, quality-trimmed forward and reverse reads were
34 assembled and filtered for ambiguous bases (“N”) and spurious length (assembled read length
35 $420 \text{ bp} \geq L \leq 466 \text{ bp}$) by using the Simple Bayesian algorithm and a threshold of 0.7 in
36 PANDASeq software (v. 2.10)².

37

38 **Text S3: Indicator species analysis**

39 Indicator species analysis assigns an indicator value (between 0 and 1) to each species in the
40 cluster or group based on the product of the relative abundance and relative frequency of that
41 species within the cluster (reclamation or reference). A high indicator value (close to 1)

42 suggests that a given species is highly abundant within a group compared to the other group
43 (or groups) (referred to as “specificity”) and is present in most members of that group
44 (referred to as “fidelity”)³. Tests for statistical significance of the indicator value was further
45 determined through permutation (probability) tests⁴. In this study, the indicator species
46 analysis was performed using the “indval ()” function in the labdsv package of R software ⁴
47 In this study, KO terms with FDR-adjusted $P < 0.1$ between reclamation and reference soils
48 of at least one site and with an indicator value > 0.6 was adjudged discriminant between
49 reference and reclamation soil.

50 **Supplementary Tables**51 **Table S1. Selected soil physicochemical properties**

Properties ¹	Site X		Site Y		Site Z	
	Recl.	Ref.	Recl	Ref.	Recl	Ref.
Cl ⁻ (mg kg ⁻¹)	1.05±0.28 ^a	1.42±0.94 ^a	0.94±0.32 ^a	1.06±0.07 ^a	1.04±0.49 ^a	0.88±0.72 ^a
Na (mg kg ⁻¹)	3.22±0.40 ^a	4.05±1.42 ^a	15.01±10.42 ^a	6.67±1.84 ^a	2.07±0.88 ^a	7.04±8.77 ^a
K (mg kg ⁻¹)	22.88±1.58 ^a	40.56±23.35 ^a	79.27±13.28 ^a	75.27±14.43 ^a	77.30±17.09 ^a	69.50±19.47 ^a
Ca (mg kg ⁻¹)	42.80±10.82 ^b	176.48±126.09 ^a	322.85±74.22 ^a	1371.90±7.50 ^b	237.08±45.47 ^a	226.84±120.26 ^a
Mg (mg kg ⁻¹)	10.76±2.22 ^b	37.01±20.27 ^a	147.03±35.52 ^b	358.32±138.24 ^a	57.89±32.05 ^a	43.39±22.58 ^a
Textural class	SaClLm	SaLm	SaClLm	SaLm	SaClLm	SaLm

52 ¹See also Table 1 for other physicochemical properties.

53

54 **Table S2: Community-level physiological profiling (31 Carbon substrate utilization pattern)**

Site	Soil group (Sample size)	Shannon-Weiner (<i>H'</i>)	Evenness (<i>J'</i>)
Site X	ReclX (N=3)	1.81±1.08	0.62±0.33
	RefX (N=5)	1.86±0.51	0.79±0.08
Site Y	ReclY (N=5)	1.68±0.58	0.73±0.21
	RefY (N=3)	1.91±0.20	0.69±0.07
Site Z	ReclZ (N=5)	1.37±0.89	0.61±0.30
	RefZ (N=5)	1.36±0.72	0.60±0.28

55 Values (mean ± SD). Differences are not significant ($P < 0.05$) based on a mixed effect model

56 (Random effect, variance = 0.001145, Standard deviation = 0.03384).

57

58 **Table S3:** Permutational tests for microbial community structure between reclamation and
 59 reference soils per site based on Bray-Curtis distances

Factors	Unweighted Bray (composition)			Weighted Bray (Structure)		
	PERMANOVA	PERMDISP	PERMDISP	PERMANOVA	PERMDISP	PERMDISP
	R ² (%)	<i>P</i>	<i>P</i>	R ² (%)	<i>P</i>	<i>P</i>
Pair-wise site comparison						
Site X	29.20	0.018	0.001	26.95	0.090	0.001
Site Y	53.09	0.018	0.120	51.50	0.026	0.119
Site Z	13.58	0.280	0.035	12.66	0.406	0.058
Sample-wide analyses						
Site	27.29	0.001	0.22	28.16	0.001	0.263
Soil History	6.70	0.007	0.04	8.13	0.004	0.032
Site x Soil History	13.65	0.001	ND	14.80	0.001	ND

60 ND, Not determined. PERMANOVA tests were performed by using the “adonis ()” of the
 61 vegan package of R software and are based on 999 iterations.

62 **Table S4:** Statistical test for discriminative genus-level features between reclamation and reference
 63 soil. See Figure 4.

	P values	FDR- adjusted P-values	Class	LDA score
<i>Massilia</i>	0.003	0.337	Reclamation	4.65
<i>Sporosarcina</i>	0.004	0.337	Reclamation	3.49
<i>Oryzihumus</i>	0.004	0.337	Reclamation	4.36
<i>Terrabacter</i>	0.005	0.337	Reclamation	4.09
<i>Mucilaginibacter</i>	0.006	0.337	Reclamation	3.72
<i>Oceanobacillus</i>	0.007	0.337	Reclamation	2.66
<i>Janibacter</i>	0.007	0.337	Reclamation	3.97
<i>Sphingomonas</i>	0.008	0.337	Reclamation	5.2
<i>Deinococcus</i>	0.009	0.337	Reclamation	2.74
<i>Rhodanobacter</i>	0.012	0.337	Reclamation	3.61
<i>Dokdonella</i>	0.013	0.337	Reclamation	3.43
<i>Segetibacter</i>	0.014	0.337	Reclamation	3.9
<i>Phycoccus</i>	0.014	0.337	Reclamation	3.73
<i>Dyella</i>	0.016	0.337	Reclamation	3.73
<i>Fulvimonas</i>	0.017	0.337	Reclamation	3.06
<i>Streptomyces</i>	0.019	0.337	Reclamation	4.41
<i>Clostridium sensu stricto 1</i>	0.021	0.337	Reclamation	3.08
<i>Opitutus</i>	0.026	0.380	Reclamation	3.42
<i>Arthrobacter</i>	0.026	0.380	Reclamation	3.82
<i>Flavisolibacter</i>	0.026	0.380	Reclamation	4
<i>Methylobacterium</i>	0.034	0.383	Reclamation	3.63
<i>Candidatus Koribacter</i>	0.034	0.383	Reclamation	3.8
<i>Jatrophihabitans</i>	0.034	0.383	Reclamation	4.04
<i>Clostridium sensu stricto 12</i>	0.040	0.383	Reclamation	2.75
<i>Burkholderia-Caballeronia-Paraburkholderia</i>	0.005	0.337	Reclamation	4.43
<i>Rubroacter</i>	0.003	0.337	Reference	-3.6
<i>Vicinamibacter</i>	0.007	0.337	Reference	-3.49
<i>Chitinophaga</i>	0.014	0.337	Reference	-2.97
<i>Lechevalieria</i>	0.018	0.337	Reference	-3.59
<i>Chryseolinea</i>	0.020	0.337	Reference	-2.08
<i>Sphingomicrobium</i>	0.020	0.337	Reference	-2.09
<i>Rhodopirellula</i>	0.020	0.337	Reference	-2.53
<i>Herpetosiphon</i>	0.020	0.337	Reference	-2.56
<i>Hirschia</i>	0.020	0.337	Reference	-2.8
<i>Flavitalea</i>	0.029	0.383	Reference	-3.05
FFCH5858	0.030	0.383	Reference	-2.64
<i>Virgisporangium</i>	0.030	0.383	Reference	-2.67
SWB02	0.032	0.383	Reference	-2.83
<i>Candidatus Protochlamydia</i>	0.037	0.383	Reference	-2.54

64

65 **Table S5:** Tax4Fun statistics for the functional prediction of soil bacterial communities

Soil group	Average FTUs	[†] 1-Average FTUs
ReclX	0.89±0.02 ^a	0.11±0.02 ^a
RefX	0.86±0.03 ^a	0.14±0.03 ^a
ReclY	0.81±0.04 ^a	0.19±0.04 ^a
RefY	0.82±0.01 ^a	0.18±0.01 ^a
ReclZ	0.84±0.01 ^a	0.16±0.01 ^a
RefZ	0.85±0.04 ^a	0.15±0.04 ^a

66 FTU, Fraction of OTUs which were not mapped against KEGG organisms.

67 [†]Fraction of OTUs which mapped onto the KEGG organisms is obtained by subtracting FTUs from

68 1.

69 **Table S6:** Pearson correlation coefficient (*r*) for the relationship between soil physicochemical
70 properties and soil physiological data

Physico-chemical ppts.	Beta-glucosidase (P-nitrophenol μg/g/h)	Alk- phosphatase (P- nitrophenol μg/g/h)	Acid- phosphatase (P- nitrophenol μg/g/h)	Urease (NH ₄ -N μg/g/2h)	Shannon- Weiner (<i>H'</i>)	Evenness (<i>J'</i>)
	<i>r</i>	<i>r</i>	<i>r</i>	<i>r</i>	<i>r</i>	<i>r</i>
pH (H ₂ O)	0.40	0.653**	0.103	0.64**	0.223	0.003
Moisture (%)	0.439*	0.374	0.275	0.443*	0.164	0.046
Organic matter (%)	0.064	0.188	-0.024	0.253	0.091	-0.097
Bulk Density (g cm ⁻³)	-0.335	-0.133	-0.375	-0.265	-0.100	-0.239
EC (mg kg ⁻¹)	-0.299	-0.442*	-0.337	-0.256	-0.222	-0.0299
Cl ⁻ (mg kg ⁻¹)	-0.089	-0.319	-0.265	-0.178	-0.397	-0.090
NO ₃ ⁻ -N (mg kg ⁻¹)	0.015	0.178	-0.084	-0.299	0.248	0.020
NO ₂ ⁻ -N (mg kg ⁻¹)	0.334	0.418	-0.080	0.346	0.483*	0.195
PO ₄ ³⁻ -P (mg kg ⁻¹)	-0.254	-0.201	-0.059	-0.433	-0.240	-0.167
Na (mg kg ⁻¹)	-0.327	0.080	-0.532**	-0.052	0.194	0.020
K (mg kg ⁻¹)	0.460*	0.207	0.417*	0.602**	-0.100	-0.188
Ca (mg kg ⁻¹)	0.561**	0.510**	0.065	0.678**	0.144	-0.126
Mg (mg kg ⁻¹)	0.351	0.478*	0.063	0.508*	0.093	-0.091
CEC (cmol (+) kg ⁻¹)	0.162	0.273	-0.355	0.100	0.286	0.055
Sand (%)	-0.08	-0.179	0.15	-0.13	-0.25	-0.239
Silt (%)	0.34	0.1.69	0.169	0.32	0.02	0.042
Clay (%)	-0.428*	-0.369	-0.27	-0.258	-0.452*	-0.332

71 *Correlation is significant at 0.05 probability level.

72 ** Correlation is significant at the 0.01 probability level.

73 **Table S7:** Significance of terms (physicochemical properties) in the CCA model of Fig. 9.

Constraints	F	Pr. (>F)
Sand	2.2023	0.136
Silt	5.5790	0.013*
Clay	1.6856	0.199
BD	4.3089	0.022*
pH	25.7254	0.001***
Moist	1.0265	0.360
OM	1.8720	0.164
EC	5.4034	0.015*
Cl	1.5691	0.221
NO ₃	2.0760	0.174
NO ₂	1.4687	0.260
PO ₄	1.3124	0.272
Na	3.4884	0.041*
K	0.9829	0.393
Ca	4.5304	0.026*
Mg	1.0584	0.376
CEC	1.5399	0.222

74 * Significant at the 0.05 probability level

75 **Significant at the 0.01 probability level

76 ***Significant at the 0.001 probability level

77 **Supplementary Figure**

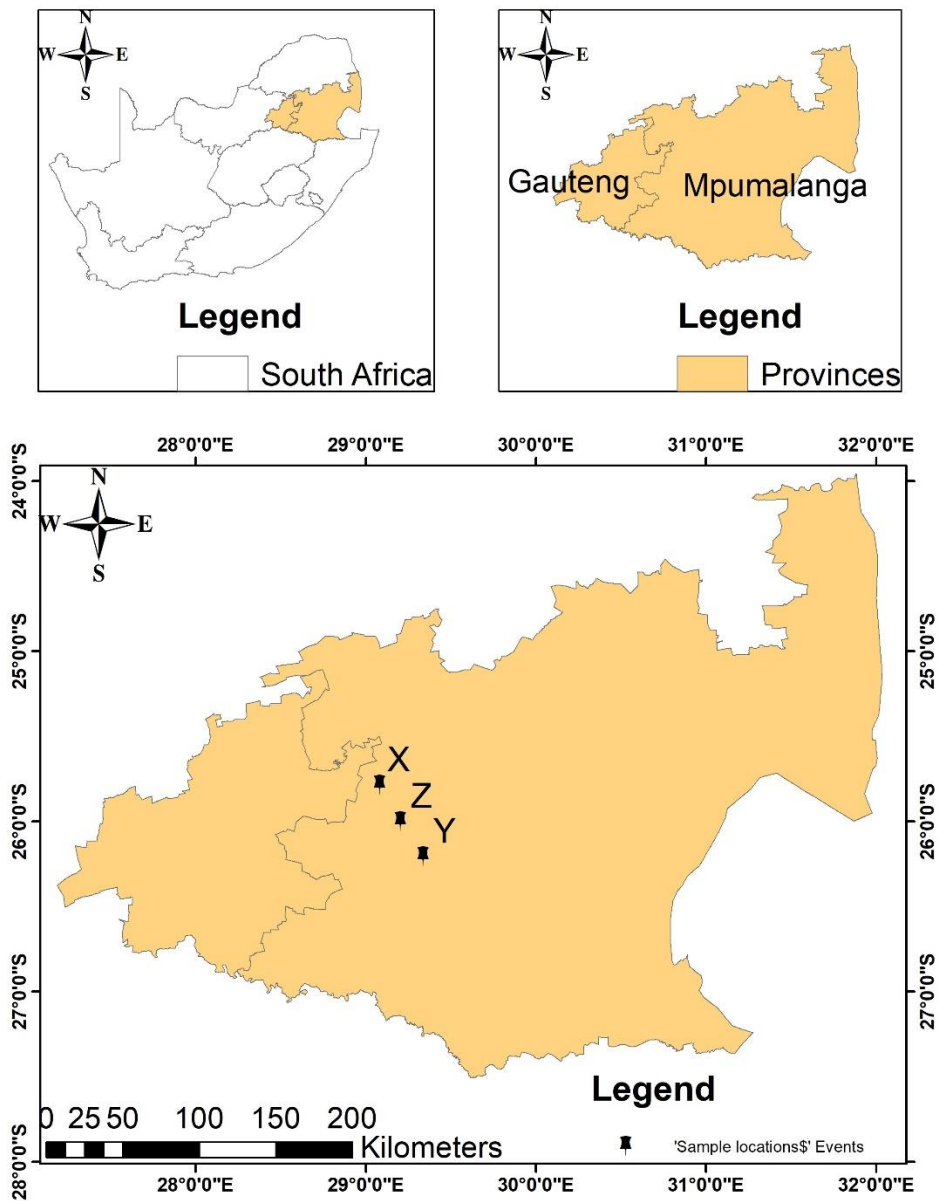


Figure. S1

78

79 **Figure S1:** Map of the sampling area. Site Y is approximately 67 km to the south of site X; site
 80 Z is approximately 32 km to the south of site X, while site Z is approximately 41 km to the
 81 north of site Y. Map was generated using the ArcMap software (v. 10.5; Esri, Redlands, CA,
 82 USA) using a shapefile obtained from the North-West University University’s
 83 (www.nwu.ac.za) local Geo-Database archives.

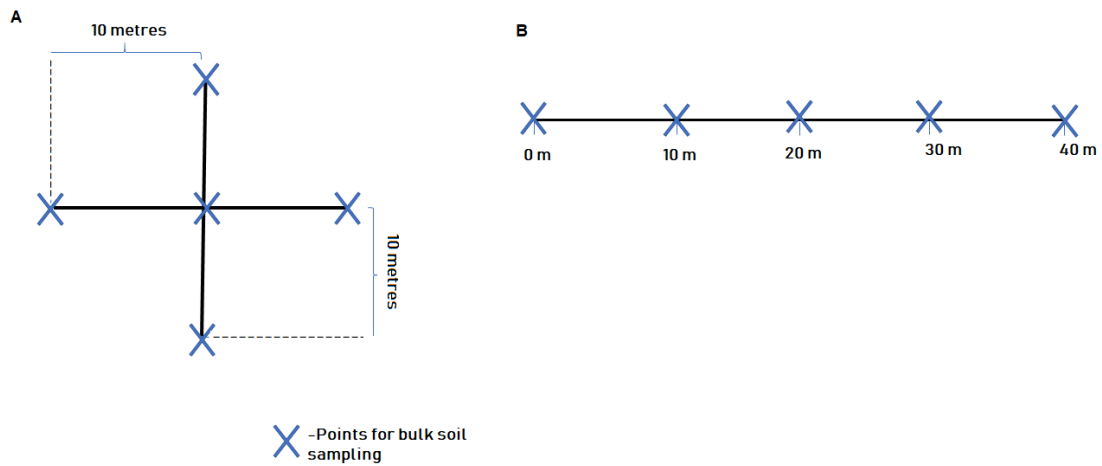
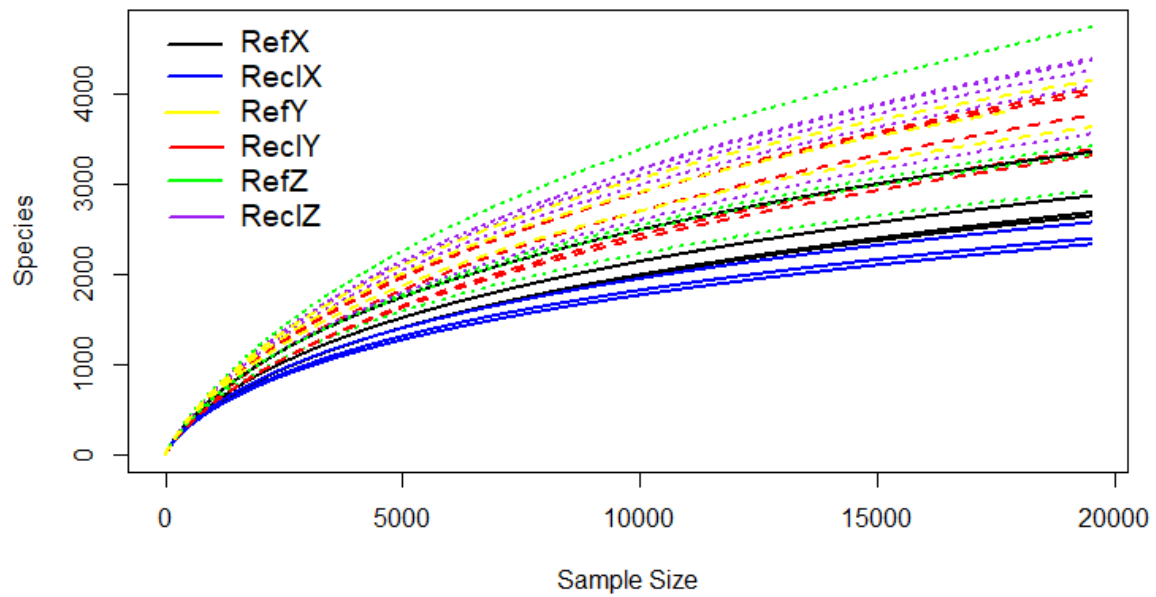


Figure S2

85

86 **Figure S2.** Schematic diagram for soil sampling designs. (A) Cross design. (B) Transect
 87 method. The sampling design was aimed at obtaining representative samples and differed based
 88 on the topography and dimension of the sites. Transect method was used in site X, while cross
 89 designs were applied to site Y and Z. Samples were collected from each sampled area using 3-
 90 5 transects or crosses which served as replicates.



91

92 **Figure S3.** Rarefaction curve. Bacterial communities (97% 16S RNA gene similarity OTUs)
 93 were subsampled at a depth of 19500 sequences per sample. The absence of a plateau in most
 94 sites suggest that the richness of the community is underrepresented.

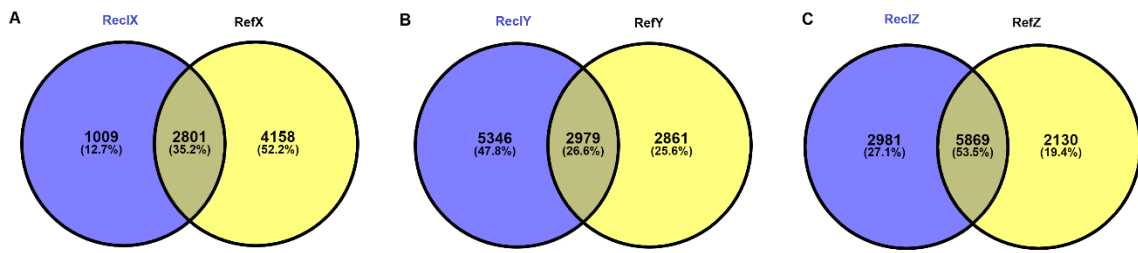


Figure S4

95

96 **Figure S4.** Unique and shared OTUs between and within sites. (A) Site X. (B) Site Y. (C) Site

97 Z. Total number of OTUs per set (or soil cluster) is the sum of unique OTUs in all replicates

98 for each soil group (e.g. site or history). The proportion (expressed in percentage) of OTUs

99 within each subset with respect to the total number of OTUs for any given sets are provided in

100 parenthesis. Venn diagram was constructed by using the online Venny 2.1 software available

101 from <http://bioinfogp.cnb.csic.es/tools/venny/>.⁵

102

103

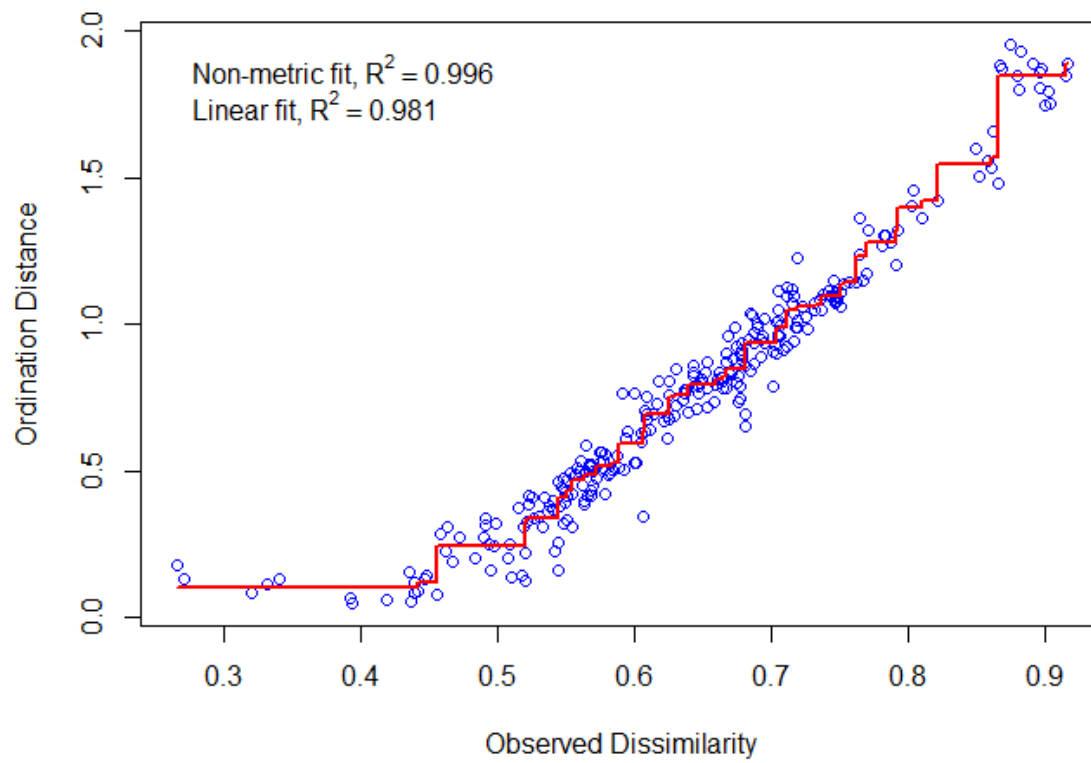


Figure S5

104

105 **Figure S5.** Stress plot for the non-metric multidimensional scaling plot of Figure 2a. Stress

106 plot was generated using the “stressplot ()” function of the Vegan package of R software^{6,7}.

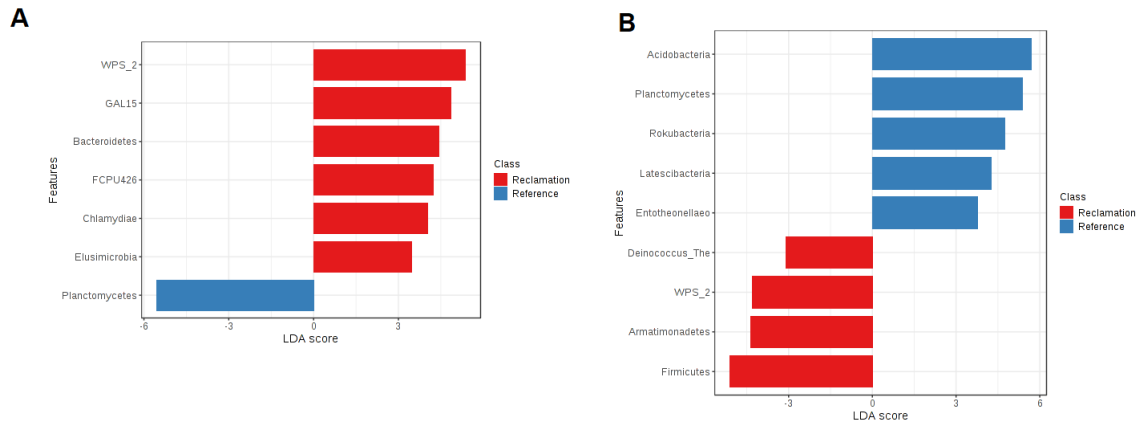


Figure S6

107

108 **Figure S6.** Differentially abundant phyla between bacterial communities of reclamation and
 109 reference soils (A) Differentially abundant phyla (LDA score > 2.0, FDR-adjusted P -value <
 110 0.1) in site X. (B) Differentially abundant phyla (LDA score > 2.0, FDR-adjusted P -value <
 111 0.3) in Site Y. Differential abundance and bar plots were determined and generated,
 112 respectively, using LefSe⁸ via the Microbiome Analyst (www.microbiomeanalyst.ca)⁹.

113

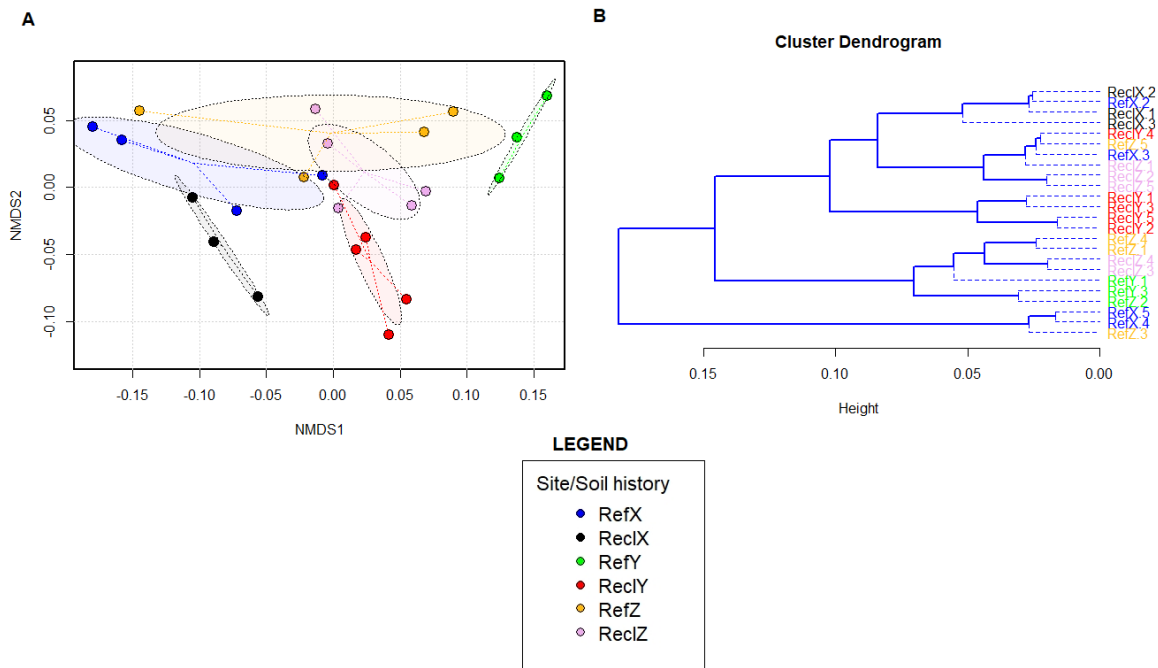


Figure S7

114

115 **Figure S7.** Bray-Curtis dissimilarity for predicted functional profile of soil bacterial
 116 communities. (A). Non-metric dimensional scaling plot. (B). UPGMA hierarchical cluster
 117 dendrogram. Dotted lines in the nMDS plot show the distance of every sample to its group
 118 centroids in multivariate space, while ellipses show 95% confidence intervals (standard error)
 119 in multivariate space around group centroids. The stress of the nMDS plot is 0.03. Differences
 120 in multivariate space are significant for site and history interactions (PERMANOVA $R^2 =$
 121 7.91%, $P = 0.045$; PERMDISP $P = 0.016$).

122 **References**

- 123 1 Bouyoucos, G. J. Hydrometer method improved for making particle size analyses of
 124 soils 1. *Agron. J.* **54**, 464-465 (1962).
- 125 2 Masella, A. P., Bartram, A. K., Truszkowski, J. M., Brown, D. G. & Neufeld, J. D.
 126 PANDAseq: paired-end assembler for illumina sequences. *BMC Bioinformatics* **13**,
 127 31, doi:10.1186/1471-2105-13-31 (2012).
- 128 3 Dufrene, M. & Legendre, P. Species assemblages and indicator species: the need for a
 129 flexible asymmetrical approach. *Ecol. Monogr.* **67**, 345-366 (1997).
- 130 4 Roberts, D. W. labdsv: Ordination and multivariate analysis for ecology. *R package*
 131 *version 1* (2007).
- 132 5 Oliveros, J. C. *Venny. An interactive tool for comparing lists with Venn's diagrams*
 133 <http://bioinfogp.cnb.csic.es/tools/venny/index.html> (2007-2015).
- 134 6 Oksanen, J. *et al. vegan: Community Ecology Package* [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
 135 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan). (2019).
- 136 7 R: A language and environment for statistical computing (R Foundation for Statistical
 137 Computing, Vienna, Austria, 2013).
- 138 8 Segata, N. *et al.* Metagenomic biomarker discovery and explanation. *Genome Biol.*
 139 **12**, R60 (2011).
- 140 9 Dhariwal, A. *et al.* MicrobiomeAnalyst: a web-based tool for comprehensive
 141 statistical, visual and meta-analysis of microbiome data. *Nucleic Acids Res.* **45**,
 142 W180-W188 (2017).