# Interactive comment on "Coalescence of bacterial groups originating from urban runoffs and artificial infiltration systems among aquifer microbiomes" by Yannick Colin et al.

**Yannick Colin et al.**

benoit.cournoyer@vetagro-sup.fr

Received and published: 7 May 2020

Replies (R2) to reviewer # 2 (major (Maj#) comments) (anonymous, 16 Mar 2020) (line numbers are those of the initial submission)

Maj1. The presentation of the sequencing process employed is inadequate. The current text highlights that the sequences were run on a Illumina MiSeq, without providing additional details.

Maj1a : First, the study does not mention how the nucleic acids are extracted from the samples, checked for quality, stored, and shipped to the facility. These points must be

clarified.

R2-maj1a : The following sentences were added to clarify these issues.

From L138: "About 600 mg of sediments or soils, or up to 5 L of aquifer or runoff water samples filtered using 0.22 $\mu$m polycarbonate filters, were used per DNA extraction. Total DNAs were extracted from soils/sediments or filters using the FastDNA SPIN$^{\circledR}$ Kit for Soil (MP Biomedicals, Carlsbad, France). For clay bead biofilms, microbial cells were detached by shaking at 2500 rpm for 2 min in 10 mL of 0.8 % NaCl. These suspensions were then filtered and their DNA content was extracted as indicated above. Blank samples were performed during these extractions for both the soils/sediments or filtered cells. DNAs were quantified using a nanodrop UV-Vis Spectrophotometer. Blank DNA extracts showed values below the detection limit. DNA extracts were visualized after electrophoresis at 6V/cm using a TBE buffer (89 mM Tri-borate, 89 mM boric acid, 2 mM EDTA, (pH 8.0)) through a 0.8% (w/v) agarose gel, and DNA staining with 0.4 mg.mL-1 ethidium bromide. A Gel Doc XR+ System (Bio-Rad, France) was used to observe the stained DNA, and confirm their relative quantities (between 20-120 ng/$\mu$l; median value around 40 ng/$\mu$L) and qualities. DNAs were kept at -80$^{\circ}$C, and shipped on ice within 24h to the DNA sequencing services when appropriate.

Quantitative PCR assays were performed on the DNA extracts to estimate their relative content in 16S rRNA gene copies. These assays were performed on a Bio-Rad CFX96 realtime PCR instrument with Bio-Rad CFX Manager software, version 3.0 (Marnes-la-Coquette, France). The 16S rRNA gene primers 338F and 518R described by Park and Crowley (2006) were used, together with the Brilliant II SYBR green low ROX qPCR master mix for SYBR Green qPCR. Melting T$^{\circ}$ was 60$^{\circ}$C. Linearized plasmid DNAs containing a 16S rRNA gene were used as standards, and obtained from Marti et al. (2017). Presence of inhibitors in the DNA extracts was checked by spiking known amount of plasmid harboring int2 (107 copies of plasmid per $\mu$L) in the PCR mix. Number of cycles needed to get a PCR signal was compared with wells where only plasmid DNA harboring int2 was added to the qPCR mix. When a high number of cycles was

needed to observe a signal, a 5- or 10-fold dilution of the DNA extract was done, and another round of tests was performed to confirm the absence of PCR inhibitions. Each assay was triplicated on distinct DNA extracts, and technical triplicates were performed. The 16S rRNA gene qPCR datasets are presented in Figure S1. These assays confirmed the high number of bacterial cells per compartment (Figure S1 and Table S2): (1) soils from the infiltration basin (IB) had a median content of 1.32 x 1011 16S rRNA gene copies per g dry weight; (2) sediments from the detention basin (DB) of 1.83 x 1011 16S rRNA gene copies per g dry weight, (3) the runoff waters (WS) had a median content of 4.75 x 108 16S rRNA gene copies per mL, (4) the aquifer waters (AQ_wat) of 3.10 x 106 16S rRNA gene copies per mL, and (5) the aquifer clay bead biofilms showed 1.35 x 107 16S rRNA gene copies per cm2."

Maj1b : Second, the study must clarify within section 2.2 several key points with respect to the sequencing protocol: (1) a citation for the primers used to target the 16S rRNA gene, (2) the protocol followed by the laboratory must be unambiguously indicated or referenced (TruSeq, Nextera, etc.), (3) the target length of the sequences, and (4) whether the sequence reads were paired-end or single.

R2-maj1b : After the text added for comment R2-maj1a, the following sentences were added to clarify the Maj1b issues:

Sequencing of V5-V6 16S rRNA gene (rrs) PCR products were performed by MrDNA DNA sequencing services (Shallowater, Texas, USA) on an Illumina Miseq. The PCR products were generated using DNA primers 799F (barcode + ACCMGGATTAGATAC-CCKG) and 1193R (CRTCCMCACCTTCCTC) reported by Beckers et al. (2016). PCR amplifications were performed using the HotStarTaq Plus Master Mix Kit (Qiagen, USA) using the following temperature cycles: 94 °C for 3 min, followed by 28 cycles of 94 °C for 30 s, 53 °C for 40 s, and 72 °C for 1 min, with a final elongation step at 72 °C for 5 min. PCR products and blank control samples were verified using a 2% agarose gel and following the electrophoretic procedure described above. PCR products obtained from field samples showed sizes around 430 bp but blanks did not show detectable

C3

and quantifiable PCR products. Dual-index adapters were ligated to the PCR fragments using the TruSeq® DNA Library Prep Kit which also involved quality controls of the ligation step (Illumina, Paris, France). Illumina Miseq DNA sequencings of the PCR products were paired-end, and set up to obtain around 40K reads per sample. The tpm DNA libraries were also sequenced by the Illumina MiSeq V3 technology but by the Biofidal DNA sequencing services (Vaulx-en-Velin, France). PCR products were generated using the following mix of degenerated PCR primers: ILMN-PTCF2 (5'- P5 adapter tag + universal primer + GTGCCGYTRTGYGGCAAGA-'3), ILMN-PTCF2m (5'- P5 adapter tag + universal primer + GTGCCCYTRTGYGGCAAGT-'3), ILMN-PTCR2 (5'- P7 adapter tag + universal primer + ATCAKYGCGGCGCGGTCRTA-'3), and ILMN-PTCR2m (5'- P7 adapter tag + universal primer + ATGAGBGCTGCCCTGTCRTA-'3) targeting conserved regions defined by FavreâĂŘBonté et al. (2005). The universal primer was 5'-AGATGTGTATAAGAGACAG-'3. The P5 adapter tag was : 5'-TCGTCGGCAGCGTC-'3. The P7 adapter tag was : 5'- GTCTCGTGGGCTCGG-'3. PCR reactions were performed using the 5X Hot BIOAmp® master mix (Biofidal, France) containing 12,5 mM MgCl2, and 10% DMSO and 50 ng sample DNA final concentrations. PCR cycles were as follow: (1) a hot start at 94°C for 5 min, (2) 35 cycles consisting of 94°C for 30 s, 58°C for 30 s and 72°C for 30 s, and (3) a ïňĄnal extension of 5 min at 72°C. The mix had two carefully optimized enzymes, the HOT FIREPol® DNA polymerase and a proofreading polymerase. This enzyme blend has both 5'→ 3' exonuclease and 3'→ 5' proofreading activities. This mix exhibits an increased fidelity (up to five fold) compared to a regular Taq polymerase. PCR products and blank control samples were verified using a 2% agarose gel and following the electrophoretic procedure described above. PCR products obtained from field samples showed sizes around 320 bp but blanks did not show detectable and quantifiable PCR products. Index and Illumina P5 or P7 DNA sequences were added by Biofidal through a PCR procedure using the same Hot BIOAmp® master mix and the above temperatures, but limited to 15 PCR cycles. Indexed P5/P7 tagged PCR products were purified

C4

using the SPRIselect procedure (Beckman Coulter, Roissy, France). PCR products and blank control samples were verified using the QIAxcel DNA kit (Qiagen, France), and band sizes around 400 bp were observed but not in the blank samples. Quantification of PCR products by the picogreen approach using the Quantifluor dsDNA kit (Promega, France) and a Qubit$^®$ 2.0 Fluorometer (Thermo Fisher Scientific, France) was performed, and showed low values among the blanks which were at the limit of detection (around 0,07 ng/$\mu$l). Still, tpm harboring bacteria being in low number among a bacterial community (about 2-3%), these controls were run during the Miseq DNA sequencing of the PCR products. Illumina Miseq DNA sequencings of the tpm PCR products were paired-end, and set up to obtain around 40K reads per sample. Blank samples generated low numbers of tpm reads (blank 1 = 24 reads; blank 2 = 3 reads, blank 4 = 1028 reads, and blank 5 = 1 read), and these have been listed in Table S3. These reads mainly belonged to unknown species (86%). However, reads from P. fluorescens (from OTUs not found in the field samples), P. xanthomarina (17 reads over all blanks) and P. fragi (n=3 reads over all blanks) were recovered but did not have any impact on the coalescence analysis.

Maj1c : Third, the presented study does not mention either positive mock community or negative comparison controls (and how those samples are incorporated into the analyses to remove contaminating sequences). The authors must present these controls.

R2-maj1c : As indicated above in replies "R2-maj1a" and R2-maj1b, several blanks and lab controls were performed all over the investigations. Blanks were run during the DNA extractions, and did not yield detectable contaminant DNAs. Furthermore, the 16S rRNA gene qPCR datasets (Table S2) confirmed that high bacterial numbers were found among each compartment investigated in this study as indicated in reply "R2-maj1a". In fact, blanks were performed during the 799F - 1193R PCR amplifications of the V5-V6 16S rRNA gene regions, and DNA yields were found below the detection limit (<0,05 ng/$\mu$l). Any contaminant DNA would thus be highly diluted and not expected to have major incidence on this 16S rRNA gene-based meta-barcoding community

coalescence analysis. However, it is to be noted that the bacterial tpm community being expected to be in lower number per sample, blank samples for the tpm meta-barcoding sequencing scheme were sequenced. As indicated in "R2-maj1b", low number of tpm reads were obtained and their matching OTUs were listed in Table S3. These reads did not match tpm OTUs transferred from the above ground environments down into the aquifer.

To further clarify these issues, the following sentences were added:

From L294: It is to be noted that blank samples sequenced during the tpm meta-barcoding assay revealed 23 Pseudomonas OTUs coming from the DNA extraction kit or generated during the PCR product Illumina sequencing process (Table S3). Only OTU00573 was found in high number (867 reads) but this contaminant did not have an impact on the coalescence analysis because of its absence in the below ground datasets. Other contaminant OTUs did not represent more than 10 times the ones observed in the field samples for identical OTUs, a criterium used to distinguish significant contaminants (Lukasik et al., 2017; doi.org/10.1111/mec.14140). In fact, only seven OTUs found among the blanks matched OTUs recovered from the environmental samples, and only two of these could be related to well-defined species i. e. P. xanthomarina (17 reads among all blanks) and P. fragi (three reads among all blanks). These reads matched a single OTU over eleven allocated to P. xanthomarina in the environmental samples, and one OTU over 52 for P. fragi.

Maj2. The results of the sequencing campaign additionally requires a more comprehensive presentation. L193-194 presents the total sequencing reads, but must present the average and range of reads per sample. A supplemental table must be provided with the raw and processed sequencing counts for each sample.

R2-maj2 : These features are now indicated in Table S2, and cited in the text. From Line 193, the following sentence was added: "The analysis of the 16S rRNA V5-V6 gene libraries yielded 2,124,272 high-quality sequences distributed across 103 samples, as

described in Table S2.

Maj3. Additionally, to explore quantitatively the mixing ratios and why certain communities are providing more biomass, the actual concentration of the community within these compartments should be mentioned or addressed as to why these measurements were neglected.

R2-maj3 : The 16S rRNA gene qPCR datasets are now shown in Figure S1 and Table S2. They confirmed a lower number of bacterial cells among the aquifer than the runoff waters.

From L343, the following sentence was added: "...These results were confirmed by qPCR estimations of 16S rRNA gene copies per compartment. These values were much lower in the aquifer waters than the runoffs."

Maj4. The bioinformatic processing pipeline requires additional information. First, the approach presented divides the 16S rRNA amplicons into 97% OTUs. However, current best practices recommends utilizing the amplicon sequencing variants (ASV) approach (Knight et al., 2018).

maj4a : The authors should either update their approach to the ASV methodology or provide a concise defense as to why they selected the OTU approach.

maj4b : Second, a rarefaction analysis is presented to subsample the dataset at 20,624 sequences. This approach has been recently called into question for more directed comparisons (McMurdie and Holmes 2014). The authors should present a concise defense as to why rarefaction was employed. To bolster this defense, Figure S1 should display the rarefaction curve for the raw data, not the previously subsampled 20,624 dataset (this comment connects with Maj2 in the need to present additional information).

R2-maj4a and 4b : Figure S1 was replaced by Figure S2 which is now showing both the OTU rarefaction curves before and after having performed a sub-sampling

at 20,624 reads per sample. OTUs were defined at a 97% identity cut-off to collapse reads into groups that reduce the incidence of sequencing errors on the dataset as suggested by several authors including Eren et al. (2013; PLOS ONE 8, doi : 10.1371/journal.pone.0066643), and Johnson et al. (2019; Nat. Commun. 10:5029, doi: 10.1038/s41467-019-13036-1).

It is to be noted that the original paper by Knights et al. (2011) describing the development of the SourceTracker made use of OTU contingency tables built with a 97% identity cut-off. This was also the case of the paper describing a "reliability" test for the source tracker inferences (Henry et al., 2016; https://doi.org/10.1016/j.watres.2016.02.029). Looking at recently published papers on the SourceTracker, one can find that most research groups have maintained a use of OTU-based contingency tables e. g. O'Dea et al. (2019, https://doi.org/10.1016/j.watres.2019.114967); Han et al. (2020, https://doi.org/10.1016/j.watres.2020.115469), Chen et al. (2019, https://doi.org/10.1038/s41598-019-42548-5), Bi et al. (2019, doi:10.1111/1462-2920.14614), and so on. Still, we confirm that a few papers have used the ASV approach to build their contingency tables for the SourceTracker and for other purposes e. g. Karstens et al. 2019, https:// doi.org/10.1128/mSystems.00290-19, and Caruso et al., 2019; https:// doi.org/10.1128/mSystems.00163-18. We recognize that the ASV approach is reliable to identify conserved ASV among datasets showing variable number of reads. However, the ASV approach also has its weaknesses. For our actual application of the SourceTracker, and according to other papers, the OTU-based contingency table was thus kept for our downstream analyses. Nevertheless, we've now cited articles on ASV in order to make sure that future readers of this paper will be aware of this approach, and might consider using it for the SourceTracker analyses.

The sub-sampling performed at 20,624 reads allowed to reduce the incidence of the variable number of reads obtained per sample. An uneven sequencing depth (ranging from 6,062 to 181,207 reads per sample) was recorded, and found to be related to

technical DNA sequencing problems. In fact, the qPCR datasets on 16S rRNA gene copies supported this conclusion. No correlation was observed between the 16S rRNA gene copy numbers (biomass) and the number of reads obtained per sample (see Table S2). In this context, we've decided to sub-sample our dataset to compensate for these discrepancies. In our opinion, sub-sampling datasets remain a good standardization technique to mitigate sample library size artifacts, especially for very unequal library sizes between groups. In accordance with this, our sub-sampled dataset (20,624 reads per sample) led to a very good separation of samples according to their origin (i.e. WS, DB, IB, AQ_wat and AQ_bio) (see Fig. 3).

From 155, the following sentences were added to clarify these issues: Variability in the number of cleaned reads per sample was observed but not correlated with variations in the number of 16S rRNA gene sequences (Table S2). These variations were thus considered to be due to the DNA sequencing process. Therefore, a sub-sampled dataset (20,624 reads per sample; with exclusion of samples with total reads below this threshold) was used to mitigate the artifact of sample library sizes. Operational Taxonomic Units (OTUs) were defined using a 97% identity cut-off as recommended by several authors in order to collapse sequences into groups that reduce the incidence of sequence errors on the datasets (e. g., Eren et al. 2013; and Johnson et al. 2019). It is to be noted that amplicon sequence variants (ASV) could also be used to build contingency tables (e. g., Callahan et al. 2016; Karstens et al. 2019). However, exact sequence variants can generate uncertainties when using 16S rRNA gene sequences because of variations among species and strains due to the presence of multiple copies per genome (Johnson et al. 2019). Figure S2 shows the OTU rarefaction curves for the full and the sub-sampled datasets. This sub-sampled dataset was used for all downstream analyses except those of the SourceTracker Bayesian approach.

Maj5. In the SourceTracker default code, the rarefied sample is then rarefied further to 1000. This procedure should be repeated to draw those 1000 reads from the full dataset, not the previously rarefied data.

R2-maj5 : We agree with this comment. Analyses were thus re-run using the cleaned but not re-sampled 16S rRNA gene reads, and the matching OTU contingency table (the one used to build Figure S2a). We then used the default SourceTracker code, including a sub-sampling of 1,000 reads as recommended by Henry et al. (2016). This analysis was run 3 times, and the coefficient of variation (i.e. Relative Standard Deviation) was used as a gauge to evaluate confidence on the computed values as suggested by Henry et al. (2016) and McCarthy et al. (2017). Table 1 was modified according to these computings.

Maj6. L319-337 presents a great overview of the study that is more appropriate for the abstract rather than the discussion. This section should be removed in its entirety.

R2-maj6 : This paragraph was deleted but a few sentences kept to facilitate the understanding of the discussion

Maj7. Throughout the text, the presence of a specific 16S rRNA transcript often is utilized to state the presence of a specific function within the community, notably within the abstract (e.g., L25, L27). Whereas the 16S taxonomical assignment is a good indicator that a specific function is likely encoded on the metagenome of the community, the linkage is not directly shown through the 16S survey and must be caveated by "likely", "putative", or "predicted to be". This is recognized more consistently within the discussion of the results, but must be maintained throughout the text to recognize that the assignment provided by FAPROTAX is a hypothesis.

R2-maj7 : Ok, this was clarified over the text.

Maj8. The authors commendably provided the raw data as publicly available datasets through EBI. Additionally, the authors should provide all code utilized to process these data as a part of the supplemental materials to allow future readers to reconstruct the presented results.

R2-maj8 : From L149, the following sentences were added so that future readers

can reproduce the results generated in this work : All paired-end MiSeq reads were processed using Mothur 1.40.4 by following a standard operation protocol (SOP) for MiSeq-based microbial community analysis (Schloss et al., 2009; Kozich et al.(2013), so-called MiSeq SOP available at http://www.mothur.org/wiki/MiSeq_SOP. Due to the large number of sequences to be processed, the cluster.split command was used to assign sequences to OTUs.

Maj9. The authors are encouraged to focus on improving the English language and grammar associated with the presented article. A non-exhaustive list of suggested grammar improvements is provided in the final section of this review, but additional editing services are recommended to enhance the clarity and accuracy of the text.

R2-maj9 : we did a complete grammar review and rewrote some sentences to clarify certain formulations.

Minor Min1. The bulk physical and chemical properties of the sampling sites should be presented or directly cited such as pH, temperature, electroconductivity etc.

reply : fixed; the most significant chemical datasets are now indicated in the paper from L365; and a selection of papers was cited so that readers can complete their knowledge of the investigated sites through analysis of these papers which present pH, electrical conductivity, soil properties, and many other datasets. See replies to reviewer 1 for this issue.

Additionally, please replace "for which physico-chemical and biological monitorings have been implemented" with "that records both physico-chemical and biological properties."

reply : fixed accordingly

Min2. L34 – Please clarify what is meant by "DNA imprints allocated"

reply : was changed for "Some tpm sequence types of . . ."

Min3. L70-75 – Please provide citations in support of these claims.

reply : fixed

Min4. L78-L79 – Replace "The tested hypotheses were that" with "Two hypotheses were tested:".

reply : fixed

Because these statements are presenting the underlying hypotheses (supported or rejected), all qualifiers for the verbs must be removed. Therefore, remove L78 "should" and L79 "could also". L79 – Replace "but" with ", and".

reply : fixed accordingly

Similarly with L88-90, please replace "was likely to be" with "will be"

reply : fixed accordingly

Min5. L291-307 – The long list of species mapped to the Pseudomonas genera is difficult to interpret in the currently presented form. Please condense this section for readability.

reply : we've tried to simplify this text but citing all these species is important for specialists; several of these species had never been described in these environmental contexts or in Europe

Min6. Throughout the text, ensure that a comma appears after Latin abbreviations such as i.e., and e.g.,

reply : fixed accordingly

Min7. Figure 1, please italicize the names of the phyla.

reply : fixed accordingly

Grammar / reply: all grammar issues raised by this reviewer were considered and fixed.

Please also note the supplement to this comment:
https://www.hydrol-earth-syst-sci-discuss.net/hess-2020-39/hess-2020-39-AC2-supplement.pdf

Interactive comment on Hydrol. Earth Syst. Sci. Discuss., https://doi.org/10.5194/hess-2020-39, 2020.

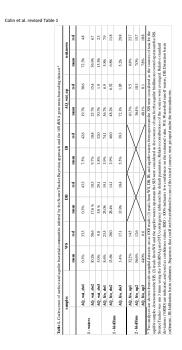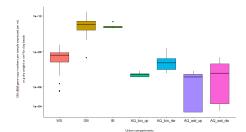Colin et al. revised Table 1

Table 1. Coalescence of surface and aquifer bacterial communities inferred by the SourceTracker Bayesian approach and the 16S rRNA gene meta-barcoding dataset*

| samples | | WS | | DB | | IB | | AQ_wat_up | | unknown | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | mean | rsd | mean | rsd | mean | rsd | mean | rsd | mean | rsd |
| 1 - waters | AQ_wat_dw1 | 0.3% | 333 | 0.3% | 433 | 7.5% | 426 | 19.7% | 306 | 72.3% | 48 |
| | AQ_wat_dw2 | 10.2% | 506 | 17.6 % | 103 | 9.7% | 188 | 25.7% | 154 | 36.9% | 67 |
| | AQ_wat_dw3 | 5.0% | 9.0 | 5.0 % | 291 | 3.8% | 320 | 70.7% | 1.9 | 15.8% | 23 |
| 1 - biofilms | AQ_bio_dw1 | 8.6% | 235 | 25.0% | 191 | 3.9% | 74.1 | 56.7% | 6.9 | 5.8% | 79 |
| | AQ_bio_dw2 | 13.6% | 280 | 28.4% | 14.1 | 2.9% | 460 | 48.2% | 6.52 | 6.8% | 118 |
| | AQ_bio_dw3 | 3.4% | 17.1 | 13.9% | 18.4 | 5.5% | 393 | 72.8% | 1.88 | 5.2% | 298 |
| 2 - biofilms | AQ_bio_up1 | 32.2% | 145 | | | | | 61.5% | 9.5 | 6.8% | 237 |
| | AQ_bio_up2 | 56.6% | 126 | | | | | 36.6% | 18.3 | 7.0% | 157 |
| | AQ_bio_up3 | 44.0% | 6.6 | | | | | 48.1% | 8.1 | 7.8% | 108 |

* Two analyses are shown from sub-sampled datasets set at 1000 reads: (1) reads from WS, DB, IB, and aquifer waters from upstream the SIS were considered as the sources of taxa for the aquifer sample downstream the SIS; (2) reads from WS and the aquifer waters upstream the SIS were considered as the sources of taxa for the aquifer biofilms recovered upstream the SIS. SourceTracker was run 3 times using the 16S rRNA pre-OTU contingency table and the default parameters. Relative contributions of the sources were averaged. Relative standard deviations (%RSD) are indicated, and used as confidence values. RSD > 100% indicates low confidence on the estimated value. WS: Watershed runoff waters; DB: Detention basin sediments; IB: Infiltration basin sediments. Sequences that could not be attributed to one of the tested sources were grouped under the term unknown.

**Fig. 1.** revised Table 1

**Figure S1. Boxplot representation of the 16S rRNA gene copy numbers measured by quantitative PCR per DNA extracts of runoff waters (WS), sediments from the detention basin (DB), soils from the infiltration basin (IB), aquifer waters (AQ_waters) or aquifer clay beads biofilms (AQ_bio).** Values were expressed per g of dry weight soil or sediment, or per mL, or per surface for the clay bead biofilms.

**Fig. 2.** new Suppl. Fig S1

C15



**Figure S2. Rarefaction curves showing the relation between the number of V5-V6 16S rRNA (*rrs*) gene reads analyzed and OTU numbers per compartment of the Mi-plaine watershed of Chassieu (France).** (a) without sub-sampling and (b) with a sub-sampling performed at 20,624 reads per sample.

**Fig. 3.** new Suppl. Fig S2

C16

**Fig. 4.** new Suppl. Table S2

C17

Table S3. Number of tpm reads among blank samples run during the tpm meta-barcoding procedure, and their taxonomic allocation and relatedness to OTUs recovered from the environmental samples. *: restricted to above ground samples; **: not considered in the coalescence analysis, see Table S8.

| blank sample OTU | total number of reads | identical OTU sequence among the environmental samples | maximum % identity with environmental tpm sequences | genus | species | blank 1 (soil) | blank 2 (soil) | blank 3 (water) | blank 4 (water) | blank 5 (water) |
|---|---|---|---|---|---|---|---|---|---|---|
| Otu01 | 867 | Otu00573* | 100 | Pseudomonas | unclassified | 0 | 0 | 0 | 867 | 0 |
| Otu02 | 118 | | 99 | Pseudomonas | fluorescens | 0 | 0 | 0 | 118 | 0 |
| Otu03 | 21 | | 99 | Pseudomonas | fluorescens | 21 | 0 | 0 | 0 | 0 |
| Otu04 | 17 | | no match | unclassified | unclassified | 1 | 0 | 15 | 0 | 0 |
| Otu05 | 17 | Otu00151* | 100 | Pseudomonas | xanthomarina | 0 | 0 | 8 | 9 | 0 |
| Otu06 | 13 | | no match | unclassified | unclassified | 1 | 0 | 12 | 0 | 0 |
| Otu07 | 10 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 10 | 0 |
| Otu08 | 7 | | no match | unclassified | unclassified | 0 | 1 | 6 | 0 | 0 |
| Otu09 | 7 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 7 | 0 |
| Otu10 | 6 | | no match | unclassified | unclassified | 1 | 0 | 5 | 0 | 0 |
| Otu11 | 5 | | no match | unclassified | unclassified | 0 | 0 | 5 | 0 | 0 |
| Otu12 | 4 | Otu01054** | 100 | unclassified | unclassified | 0 | 0 | 0 | 3 | 0 |
| Otu13 | 3 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 3 | 0 |
| Otu14 | 3 | Otu00069 | 100 | Pseudomonas | fragi | 0 | 0 | 3 | 0 | 0 |
| Otu15 | 3 | Otu00002* | 100 | Pseudomonas | unclassified | 0 | 0 | 2 | 1 | 0 |
| Otu16 | 2 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 2 | 0 |
| Otu17 | 2 | | no match | unclassified | unclassified | 0 | 0 | 0 | 1 | 0 |
| Otu18 | 2 | | 98 | unclassified | unclassified | 0 | 0 | 2 | 0 | 0 |
| Otu19 | 2 | Otu00519** | 100 | unclassified | unclassified | 0 | 0 | 0 | 1 | 1 |
| Otu20 | 2 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 2 | 0 |
| Otu21 | 2 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 2 | 0 |
| Otu22 | 2 | Otu00556** | 100 | Pseudomonas | unclassified | 0 | 2 | 0 | 0 | 0 |
| Otu23 | 2 | | 99 | Pseudomonas | unclassified | 0 | 0 | 0 | 2 | 0 |

**Fig. 5.** new Suppl. Table S3