# Evolutionary ecology of natural comammox *Nitrospira* populations

Alejandro Palomo[1], Arnaud Dechesne[1], Otto X. Cordero[2] and Barth F. Smets[1]

[1]Department of Environmental Engineering, Technical University of Denmark, Kgs. Lyngby, Denmark

[2]Ralph M. Parsons Laboratory for Environmental Science and Engineering, Department of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

## ABSTRACT

Microbial life on Earth commonly occurs in diverse and complex communities where species interact, and their genomic repertoires evolve over time. Our understanding of species interaction and evolution has increased during last decades, but most studies of evolutionary dynamics are based on single species in isolation or experimental systems composed of few interacting species. Here, we use the microbial ecosystem found in groundwater-fed sand filters as a model to avoid this limitation. In these systems, diverse microbial communities experience relatively stable conditions, and the coupling between chemical and biological processes is generally well defined. Metagenomic analysis of 12 sand filters revealed systematic co-occurrence of at least five comammox *Nitrospira* species, favoured by low ammonium concentrations. *Nitrospira* species showed intra-population sequence diversity, although possible clonal expansion was detected in few abundant local comammox populations. *Nitrospira* populations were separated by gene flow boundaries, suggesting natural and cohesive populations. They showed low homologous recombination and strong purifying selection, the latest process being especially strong in genes essential in energy metabolism. Positive selection was detected on genes related to resistance to foreign DNA and phages. Additionally, we analysed evolutionary processes in populations from different habitats. Interestingly, our results suggest that in comammox *Nitrospira* these processes are not an intrinsic feature but greatly vary depending on the habitat they inhabit. Compared to other habitats, groundwater fed sand filters impose strong purifying selection and low recombination. Together, this study improves understanding of interactions and evolution of species in the wild, and sheds light on the environmental dependency of evolutionary processes.

Microorganisms dominate the tree of life based on species number and diversity, and they play an essential role in Earth's global biogeochemical cycles. Microbial species interact with each other and with the environment (ecological processes), and also undergo changes in their genomic repertoire over time (evolutionary processes). Yet, the interaction between ecological and evolutionary processes is largely unknown, especially for complex open communities. For many years, most studies of microbial communities in open, complex environments have focused on ecological aspects as it was believed that evolutionary changes happen at much larger timescale[1]. However, in recent years, with the growth of population-genomics analysis, researchers have started to investigate both ecological and evolutionary processes in microbial communities. Yet, most studies of evolutionary dynamics are based on single species in isolation[2] or experimental systems composed of only a few interacting species[3]. Although these analyses have helped to better understand some aspects of evolutionary processes patterns, they have limitations because they lack many characteristics of actual natural populations (spatial structure, existence of microdiversity, predation, immigration, etc.). On the other hand, observing populations in the wild also has limitations because the conditions vary with little control (hence uncontrolled variation in population size, selection regime) and the typically unknown ecophysiology of retrieved genomes makes it difficult to interpret the observed patterns. Therefore, studying well-defined model microbial ecosystems can help to understand ecological and evolutionary processes in microbial communities[4].

Rapid sand filters (RSF), widely used to produce drinking water from surface- or groundwater, can be a useful model system as they are characterized by stable conditions and active growth primarily driven by the oxidation of ammonia, methane, and other inorganic compounds present at low concentration in the influent water, large populations ($\sim 10^9$ cells/g), significant mixing, continuous but limited immigration from prokaryotes in the influent water, no dispersal between separate sand filters (resulting in allopatric populations), and relatively well defined coupling between chemical and biological processes[5–7]. In addition, microbial communities inhabiting these systems have been described and show the dominance of complete ammonia oxidizers (comammox)[8,9], which are expected to have a relatively simple basic ecology (due to their chemolithoautotrophic metabolism)[10], yet are poorly studied in terms of what drives their diversity, distribution and evolution. Furthermore, as comammox bacteria occur in RSF as coexisting populations[9,11], RSF offer an appropriate opportunity for resolving fine-scale genomic heterogeneity within closely related strains, and investigate if they show similar patterns in evolutionary processes (selection, recombination, etc.).

With the RSF microbial ecosystem model, different eco-evolutionary questions can be addressed. Of particular interest is to what extent the evolutionary processes that drive the diversification of *Nitrospira* species in RSF are dependent on their environment, as opposed to intrinsic properties of the species. Environmental dependency of microbial evolution has been investigated from different perspectives. Several studies have focused on genome signatures variations (GC, tetranucleotide signatures, codon usage, purine-pyrimidine ratio) associated with different environments (reviewed in Dutta and Paul (2012)[12]). Others, have studied bacterial adaptation to shifting environments[13], or have targeted a specific evolutionary process across several lifestyles (homologous recombination[14], selection[15], etc.). Most of these studies, however, considered different species living in different environments, or closely related species with a different lifestyle
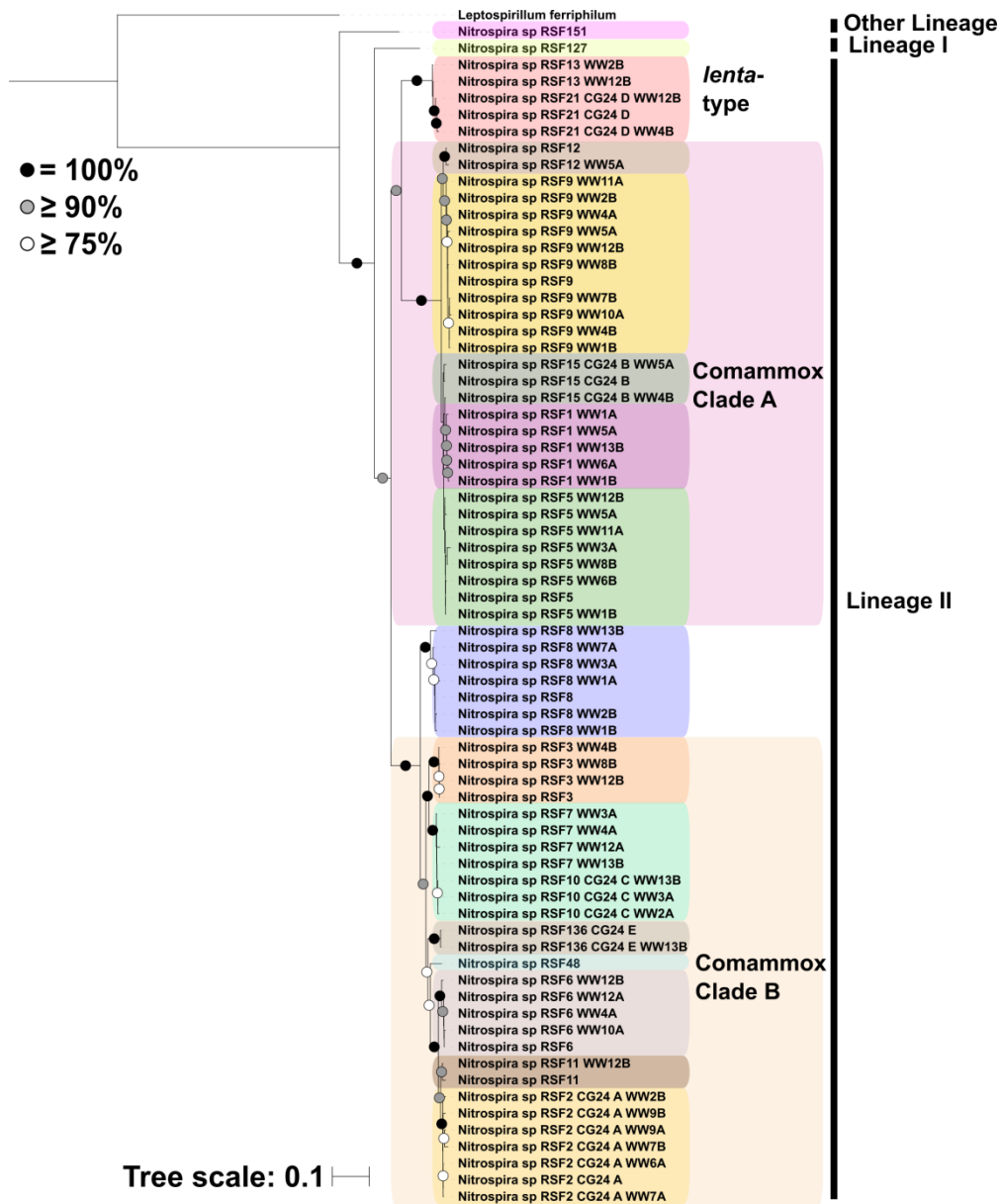
84    (i.e. free-living organisms vs pathogens). Yet, little is known about ongoing evolutionary processes
85    of species belonging to the same lineage inhabiting different open environments. In this study, taking
86    advantage of the multiple *Nitrospira* species present in several groundwater-fed RSF, we thoroughly
87    investigated evolutionary processes on this local environment, and compared these observations with
88    those in *Nitrospira* species inhabiting other open environments.

89

90

## **Results and Discussion**

In this study, we examined ecological and evolutionary patterns within comammox dominated-bacterial communities inhabiting groundwater-fed rapid sand filters. To that end, we retrieved *Nitrospira* metagenome-assembled genomes (MAGs) from 12 similarly operated waterworks in Denmark using a combination of automatic and manual binning (Supplementary Table 1). These MAGs spanned 16 putative species (further on simply referred to as 'species') using a threshold average nucleotide identity (ANI) of $\geq 95\%$[16–18]. The phylogenomic analysis placed one *Nitrospira* species into lineage I, 14 into lineage II, and one into other lineages (Fig. 1). Of the 16 *Nitrospira* species, 12 were classified as comammox *Nitrospira* (5 clade A and 7 clade B) (Fig. 1).

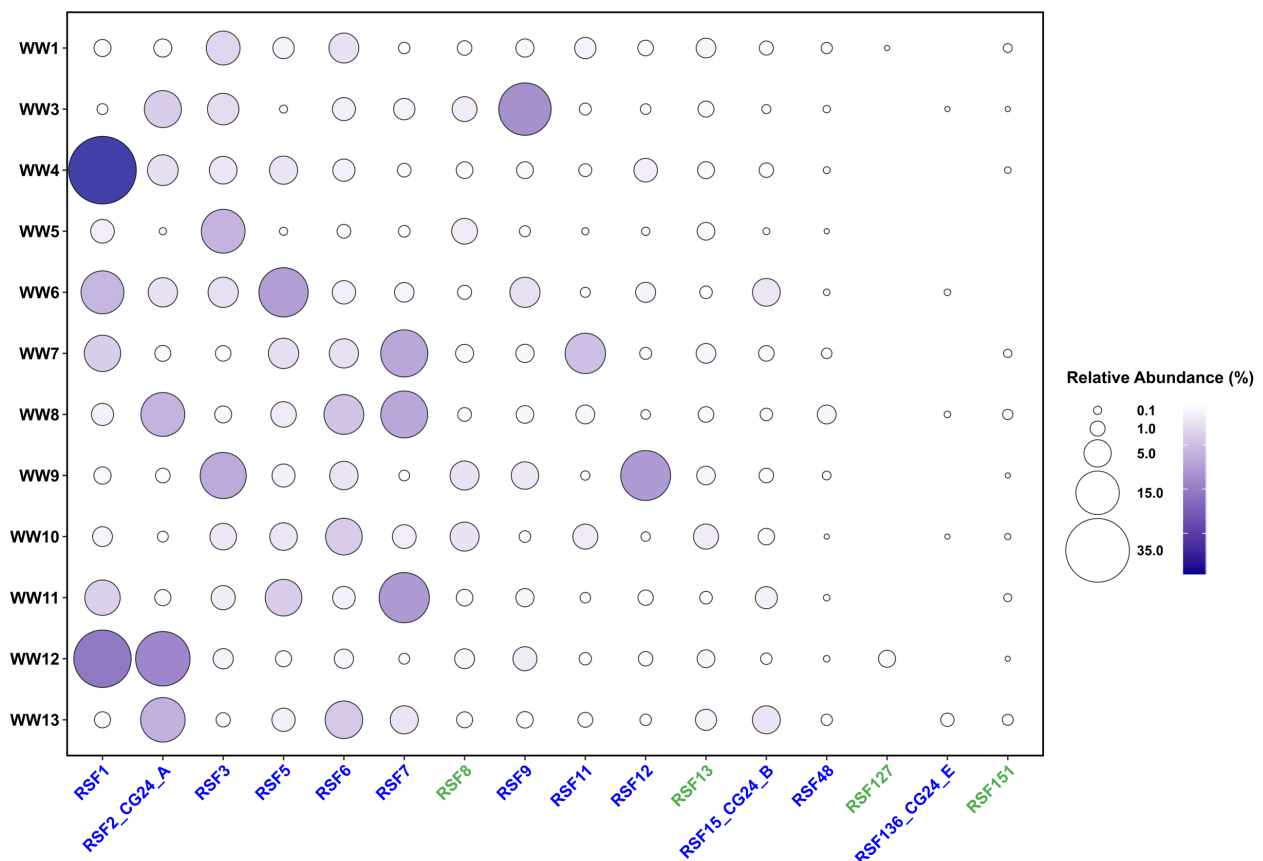**Fig. 1. Phylogenomic affiliation of *Nitrospira* MAGs retrieved from 12 waterworks.**
A phylogenetic tree was built based on the concatenation of 120 proteins. MAGs affiliated to same *Nitrospira* species (MAGS with ANI ≥ 95% are considered members of the same species) are shown with same colours. *Leptospirillum* was used to root the tree. The strength of support for internal nodes as assessed by bootstrap replicates is indicated as coloured circles (top left legend).

*Nitrospira* species comprised a large proportion of the microbial communities of the waterworks (27 - 70%), and comammox represented a large fraction of *Nitrospira* spp. (76 - 98%) (Fig. 2 and Supplementary Table 2). Multiple *Nitrospira* species (at least 5, 10 on average) co-occurred in all the waterworks (Fig. 2). In some cases, a single species constituted the majority of the community (WW4 and WW5), while, in other, two (WW9 and WW12) or more species dominated (Fig. 2). Distinct species were dominant in the different waterworks (Fig. 2). Generally, there was no significant correlation among the abundances of the comammox species, with the exception of few strong positive correlations (RSF5 and RSF15_CG24_B ($\rho$=0.84), RSF6 and RSF11 ($\rho$=0.84)), and a negative correlation between two species (RSF3 and RSF7 ($\rho$=-0.47)) (Supplementary Fig. 1). Thus, co-exclusion between comammox *Nitrospira* species seems to be rare in the studied waterworks. One reason to explain this phenomenon could be the existence of large amount of unique genes in each comammox *Nitrospira* species, which has been previously reported[11]. As in other coexisting microorganisms, these unique genes (although the function of most of them is still unknown) might promote differences in chemotactic strategies, attachment to particles strategies, secondary metabolism, defence against predation, etc. (Reviewed in Stilianos *et al.*, (2018)[19]). Furthermore, the sand grains of the investigated waterworks are very porous with bacterial cells embedded in different parts of the grain[20], which could enable fine-scale spatial separation among the comammox species. Nevertheless, more research is needed to properly characterise the reason enabling the observed co-occurrence of comammox *Nitrospira* species.

Although correlations between *Nitrospira* species were rare, positive correlations between the abundance of *Nitrospira* spp. and that of other microbes inhabiting the waterworks were more frequent (Supplementary Fig. 1). In particular, the abundance of most of the ammonia oxidizing bacterial (AOB) species positively correlated with canonical *Nitrospira* species, and interestingly, with one comammox *Nitrospira* species (RSF3) (Supplementary Fig. 1). We also observed that distinct *Nitrospira* species positively correlated with different phages retrieved from the waterworks, suggesting a specificity between phages and hosts (Supplementary Fig. 2).

The chemical characteristics of the water explained 31% of the variance in *Nitrospira* composition (permutation test: p < 0.001; Supplementary Fig. 3), suggesting that water chemistry is a strong filter for the assembly of these nitrifying communities. Among the measured water constituents, the influent ammonium concentration was the variable that best explained the *Nitrospira* distribution (explained 18%; permutation test: p = 0.002). Higher comammox species richness was detected in waterworks treating lower ammonium concentration (Supplementary Fig. 4, $R^2$ = 0.54, p < 0.01). In contrast, most of canonical *Nitrospira* and canonical ammonia oxidizers were related with waterworks containing higher ammonium concentration in their influent (Supplementary Fig. 3). These observations are in line whit a previous study which predicted that higher ammonium concentration favours emergence of division of labour (canonical ammonia and nitrite oxidiser)[21].

147   Nevertheless, we observed that a few comammox species seems to also cope with slightly higher
148   ammonium concentrations (Supplementary Fig. 3). Moreover, as previously reported[22], no clear
149   distinction in comammox species distribution was observed based on their clade affiliation (which
150   depends on the phylogeny of ammonia monooxygenase subunit A) (Supplementary Fig. 3). Different
151   from the water chemistry composition, the distribution patterns of *Nitrospira* species across the
152   waterworks were not related to their geographic distance (Mantel test: r statistics = 0.08 and
153   significance > 0.05 vs. r statistics = 0.36 and significance < 0.001 for water chemistry).
154
155



156   **Fig. 2. Abundance of *Nitrospira* species across 12 waterworks**.
157   Relative abundance of 16 *Nitrospira* species in 12 waterworks. Comammox and canonical *Nitrospira* species
158   are denoted in blue and green, respectively.
159
160
161
162
163
164
165
166
167

## Microdiversity within *Nitrospira* species

Strain-level analysis across the waterworks revealed that the *Nitrospira* populations contained intra-population sequence diversity. We exploited the shotgun metagenomic data from the different waterworks to perform strain-level analyses using single nucleotide polymorphisms (SNPs). The number of SNPs/Mbp in the populations across the waterworks ranged from 14,437 to 45,664 (Supplementary Table 3). Looking into the populations at local scale (species within waterworks), the number of SNPs/Mbp ranged from 249 to 37,663 (Supplementary Table 4). We observed a wide range of microdiversity (measured as nucleotide diversity ($\pi$)) among populations (Fig. 3A): canonical *Nitrospira* RSF8 was the most diverse species, with three times more nucleotide diversity than the less diverse *Nitrospira* species of our study (RSF1 and RSF12) (Fig. 3A). Both a homogeneous degree of microdiversity across the waterworks (e.g., RSF5 and RSF8), as well as a high microdiversity variation depending on the waterworks (e.g., RSF1, RSF9 and RSF11) was detected among the *Nitrospira* populations (Supplementary Fig. 5). Based on the observations done at species level, we hypothesised that microdiversity would be higher at low ammonium concentrations, where we observed higher comammox species diversity. However, this was not the case, as for each species, the correlations of microdiversity with ammonium concentration or comammox species richness were not significant ($p > 0.05$).

In contrast to the observation across all the waterworks (i.e.: all species showed significant microdiversity across waterworks), we did detect a few highly abundant comammox populations with almost no microdiversity at local scale (e.g., comammox *Nitrospira* RSF1 in WW4, and comammox *Nitrospira* RSF12 in WW9) (Supplementary Fig. 6), which suggests local clonal expansions of these comammox populations in specific waterworks. In the same line, the analysis of major allele frequencies of common SNPs (for each species, SNPs present in all the strains present in each waterworks) revealed that only a fraction of the subspecies diversity is found locally (for a specific species, we detected different subspecies among waterworks, some of them being genetically homogenous. E.g.: comammox RSF3 (WW10B vs WW5) (Supplementary Fig. 7)). These results contrast with what we observed at species level, where all the diversity was represented in each waterworks (as they all contain most of the *Nitrospira* spp.; Fig. 2).

Similar to what we observed at species level, there was no significant correlation between similarity in subspecies composition and the geographic distance of the waterworks, with exception of *Nitrospira* sp. RSF2 ($p < 0.01$) and *Nitrospira* sp. RSF8 ($p < 0.05$) (Supplementary Fig. 8). However, we observed a geographic organisation of the genetic structure at most loci across the studied genomes, indicating that the *Nitrospira* populations were more similar within than between waterworks. We calculated pairwise fixation indexes $F_{ST}$ (which measures differences in allele frequencies between populations of the same species found in two distinct waterworks) for each gene between allele frequencies from the twelve waterworks. The mean gene $F_{ST}$ values were $\geq 15\%$ for all *Nitrospira* populations (Supplementary Fig. 9), the most extreme case being the RSF5 population ($F_{ST} > 40\%$) (Supplementary Fig. 9). In few populations (RSF2, RSF7 and RSF8), a higher dispersal ($F_{ST} < 20\%$) of most alleles between waterworks was observed (Supplementary Fig. 9). These observations differ from observations from soil bacterial populations across a meadow, where most

209  of the populations had mean gene $F_{ST}$ values < 5%[23]. These contrasting results support the notion that
210  populations in the waterworks are much more allopatric than the ones from the mentioned meadow.
211       We also investigated local regions of the *Nitrospira* genomes with significantly higher
212  $F_{ST}$ values, as this is characteristic of local population-specific (here in each waterworks) selective
213  pressures acting on specific loci[23]. Few loci were found with unusually high $F_{ST}$ in some of the
214  *Nitrospira* populations (Supplementary Fig. 10 and Supplementary Table 5), one of them being
215  related with genes involved in nitrogen assimilation (*Nitrospira* sp. RSF2) (Supplementary Table 5).
216  However, only two of these loci with unusually high site-specific differentiation of alleles (high $F_{ST}$)
217  also had fewer recombinant events, and lower nucleotide diversity (Supplementary Table 5), which
218  can be considered as a signal of recent selective sweep[23]. These results suggest that contrary to what
219  it has been observed in several natural populations[23–26], gene-specific sweeps seems to play a minor
220  role in the evolution of *Nitrospira* spp. inhabiting the waterworks. A possible explanation could be
221  the low recombination rate that characterised the waterworks *Nitrospira* populations (Fig. 3B,
222  discussed below), as opposed to genome-wide sweeps that are associated to low recombination
223  rates[27], gene-specific sweeps are expected to occur with high recombination rates[27].
224       Overall, across the 12 waterworks, all species present significant genomic microdiversity but
225  this diversity was not always represented locally, with a few occurrences of patterns consistent with
226  clonal expansion. The reason for the difference of within-species diversity across waterworks is
227  unknown but the allopatric nature of the communities likely contributes to their persistence.

228
229

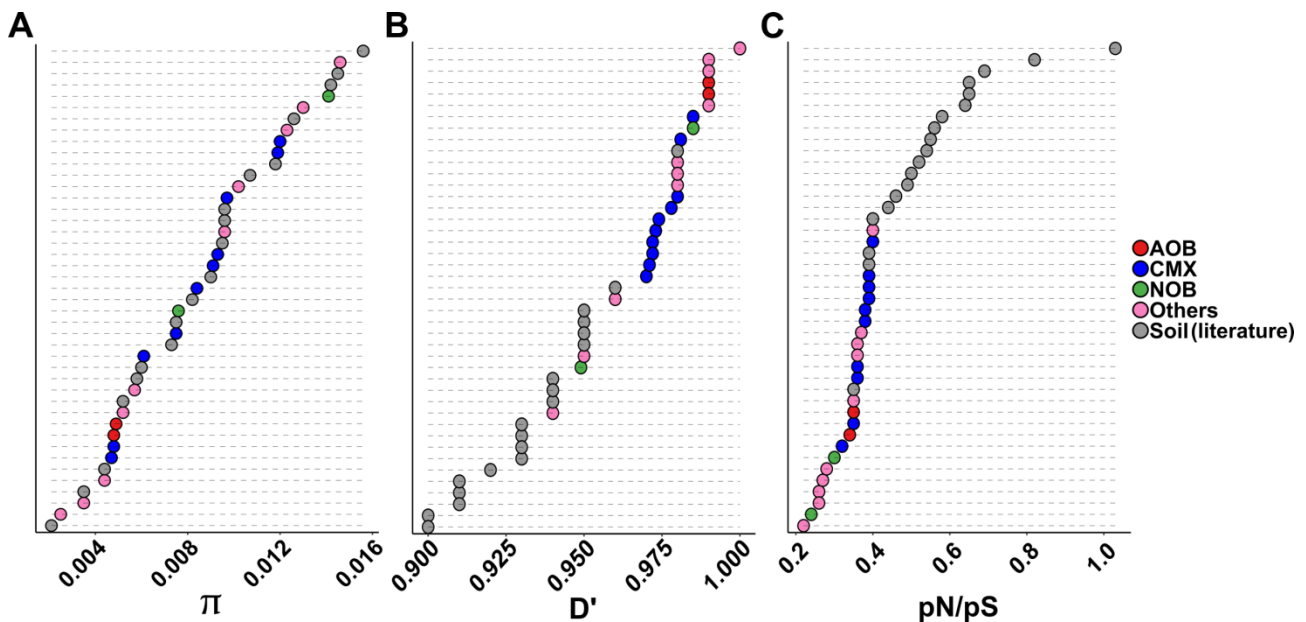## Evolutionary processes at whole-genome level

231

232       The *Nitrospira* populations were characterised by a low degree of homologous recombination.
233  We investigated homologous recombination in the *Nitrospira* populations based on linkage
234  disequilibrium (D' is only < 1 if all possible combinations of a pair of biallelic sites are observed[28].
235  Lower D' values indicates higher levels of homologous recombination; Fig. 3B) and linkage decay
236  (Supplementary Fig. 11). Similar results were observed for other abundant populations inhabiting the
237  waterworks (Fig. 3B). In general, recombination was lower in the waterworks populations than in
238  populations inhabiting a grassland meadow[23], where a similar study was conducted (Fig. 3B). To
239  further examine the relative effect of homologous recombination on the genetic diversification of
240  populations, we measured the rates at which nucleotides become substituted as a result of
241  recombination versus mutation using the *r/m* ratio. Most of the *Nitrospira* populations had a relatively
242  low *r/m*  (*r/m* < 2) compared to recombinogenic species reported in literature (*r/m* > 4)[29]
243  (Supplementary Fig. 12), although in one case (RSF15_CG24_B) the rate was similar to the value
244  reported for a *S. flavogriseus* population (*r/m* = 28) considered to be approaching panmixia[30]. Overall,
245  these results suggest a low effect of recombination in the *Nitrospira* population inhabiting the studied
246  waterworks. Increasing recombination rate has been associated with fluctuating environments as a
247  source of variation which can accelerate adaptation favouring survival in this type of
248  environments[31,32]. On the other hand, constant environments - as the waterworks studied here - tend
249  to reduce recombination rates of inhabiting microbes[31].

250        A phylogenetic analysis to assess the impact of recombination across the different *Nitrospira*
251   species showed that these ones were separated by gene flow boundaries, consistent with the notion
252   that these species represent cohesive populations (i.e. gene flow within species is higher than between
253   species). We performed a quartet-based phylogenetic analysis for each *Nitrospira* species by building
254   tree quartets of whole-genome orthologous genes ("gene trees") and core genes ("species trees") with
255   four strains of the same species or with two pairs of strains from two different species. In most of the
256   pairwise species comparisons, the percentage of gene trees supporting the species phylogeny was >
257   98% (Supplementary Fig. 13). In few cases, this percentage was lower (down to 67% in the
258   comparison between RSF1 and RSF15_CG24_B), but these numbers were always above the within-
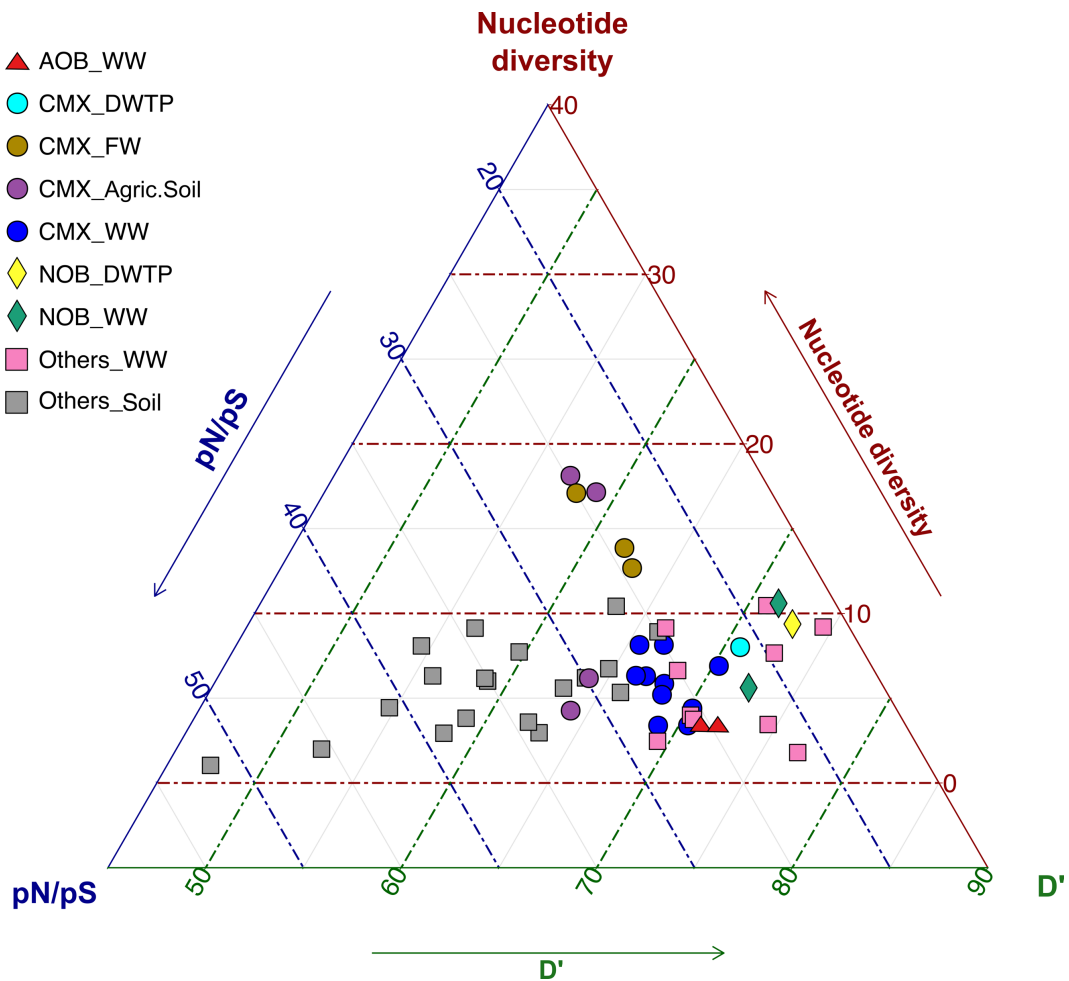259   species analysis (Supplementary Fig. 13).

260        The *Nitrospira* populations were characterised by strong purifying selection. We used the
261   relation between non-synonymous and synonymous polymorphisms (pN/pS) to investigate this
262   evolutionary process. We detected pN/pS < 1, indicating purifying selection, for all *Nitrospira* species
263   (Fig. 3C). Similar results were observed for other abundant populations inhabiting the waterworks
264   (Fig. 3C). Purifying selection has frequently been observed in wild populations, and it was the case
265   for populations inhabiting a grassland meadow[23] (pN/pS = 0.56 ± 0.17, n = 19), but this process seems
266   to be especially strong in the waterworks populations (pN/pS = 0.34 ± 0.05, n = 24) (Two-Sample t-
267   test, p < 0.0001) (Fig. 3C), which suggests prior evolution to optimal adaptation to this stable
268   environment, with purging of non-synonymous mutations.



269
270 **Fig. 3. Evolutionary metrics of *Nitrospira* populations across 12 waterworks.**
271 A Right) Genetic diversity of most abundant bacterial populations across 12 waterworks. A left) It also
272 includes most abundant bacterial populations across grassland meadow[23]. Microdiversity is measured as
273 nucleotide diversity (π). B Right) Homologous recombination (D') of most abundant bacterial populations
274 across 12 waterworks. B Left) It includes most abundant bacterial populations across grassland meadow. C
275 left) Selection (pN/pS ratio) of most abundant bacterial populations across 12 waterworks. C right) It includes
276 most abundant bacterial populations across grassland meadow. Colour legends are displayed on the right of
277 each figure.

278       Interestingly, the degree of recombination and diversity across different *Nitrospira*
279  populations varied substantially with habitat (Fig. 4). High variability of recombination in closely
280  related bacterial species has occasionally been reported[33], and lifestyle appear as one of the most
281  relevant factors to explain this variability[14,33]. Our analysis of evolutionary processes in populations
282  from different habitats (*Nitrospira* from drinking water treatment plants (DWTP), freshwaters and
283  soils) suggests that the environment also influences ongoing evolutionary processes: different
284  bacterial types inhabiting the same environment tended to share similar features (Fig. 4), while the
285  evolutionary characteristics of comammox *Nitrospira* populations differed depending on the
286  environment where they were retrieved (Fig. 4). Comammox species in the studied waterworks and
287  in other DWTP were characterised by low recombination, strong purifying selection and moderate
288  microdiversity (Fig. 4). On the other hand, comammox present in freshwater or, especially, in soils
289  had higher microdiversity and recombination rate, and weaker purifying selection (Fig. 4).
290  Intriguingly, we consistently observed that in drinking water treatment systems canonical *Nitrospira*
291  species showed features similar to those of comammox *Nitrospira* but with even stronger purifying
292  selection (Fig. 4). This feature, together with the much lower richness observed in canonical
293  *Nitrospira* compared to comammox bacteria (Fig. 2), suggests that competition can play a more
294  intense role in canonical *Nitrospira*, which might select for few species optimally adapted to this type
295  of stable environment. However, a broader analysis is required to confirm this hypothesis.
296

**Fig. 4. Impact of environment and microbial type in evolutionary metrics**
Triplot composed of the nucleotide diversity, pN/pS ratio and D' values for the bacterial populations of this study (WW) as well as most abundant bacterial populations across grassland meadow[23] (Soil), and other *Nitrospira* populations abundant in other systems (Supplementary Table 6; CMX_FW, CMX_DWTP, NOB_DWTP, and CMX_Agric.Soil). Colour legends are displayed on the left of the figure.

## Evolutionary processes at the gene level

In addition to a genome-wide analysis, we investigated the evolutionary processes at the gene level. Genes involved in nitrification (ammonia monooxygenase: *amoA* and *amoB*; hydroxylamine dehydrogenase: *haoA* and *haoB*; nitrite oxidoreductase: *nxrA* and *nxrB*) in the studied *Nitrospira* populations generally had a similar nucleotide diversity ($\pi$) (Fig. 5A) and homologous recombination rate (D') (Fig. 5B) compared to the rest of the genome, but with higher levels of purifying selection (pN/pS) (Fig. 5C). The nucleotide diversities of genes related to nitrification were very similar with the exception of *amoB*, which had a significantly lower nucleotide diversity than *nxrB* ($p < 0.05$) (Fig. 5A). A similar pattern was detected for the recombination, but in this case *amoA*, as well as *amoB*, had significantly lower recombination than *nxrB* ($p < 0.05$) (Fig. 5B). We observed a very strong purifying selection for most of the nitrifying genes, especially for *amoA*, *nxrA*, and *nxrB* ($p < 0.01$) (Fig. 5C). In the case of *nxrB*, not a single non-synonymous mutation was found in most of the *Nitrospira* species (0-1 non-synonymous sites vs 17-66 synonymous sites), even though this gene had a higher nucleotide diversity and homologous recombination (Fig. 5A and Fig. 5B). Our observations on selection are in line with previous studies, as generally, bacterial essential genes and enzymes catalysing reactions that are difficult to by-pass through alternative pathways are subject to higher purifying selection compared to nonessential ones[34–37].
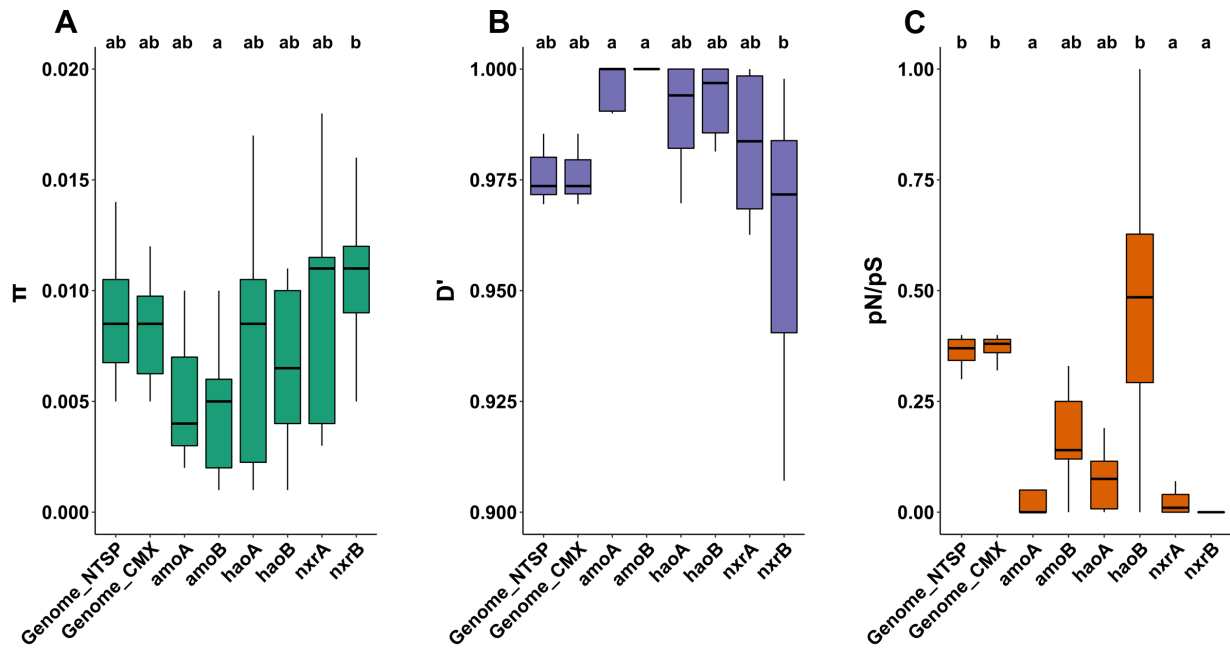
Even though the average pN/pS values were below 1 in all *Nitrospira* species (Fig. 3), indicating purifying selection, genes with pN/pS values above 1 and significantly higher than the genomic average were detected in each species (Supplementary Fig. 14). Many of those genes were related to defence mechanism against phages (e.g. genes putatively involved in phage entry into cells, ribonucleases, genes coding proteins associated to restriction-modification systems, genes related to toxin-antitoxin systems, etc (Supplementary Table 7)). Comparable findings were made in other abundant species from the waterworks (Supplementary Table 8), in the additional *Nitrospira* species retrieved from other environments (Supplementary Table 8), as well as by Petersen *et al.* (2007) and by Rabbi *et al.* (2015) in *E. coli* and *Vibrio* sp. strains, respectively[38,39]. These observations suggest that positive selection in phage-related genes is widespread across bacteria, and highlights the evolutionary arms race occurring between phages and bacteria as an important driver in bacterial ecology and evolution[27,40,41]. Additionally, we found nondefense, mobile genetic elements, such as transposons and integrases, with significantly higher pN/pS values than the genome average in the *Nitrospira* spp. (Supplementary Table 9).

**Fig. 5. Evolutionary metrics of nitrification genes in *Nitrospira* populations across 12 waterworks.**
Left) Boxplot of nucleotide diversity of *Nitrospira* bacterial populations for whole genome (all *Nitrospira* and comammox *Nitrospira*) and nitrification genes. Differences between the mean nucleotide diversities were assessed by a Dunn's test; same letter have means not significantly different from each other ($p < 0.05$). Middle) Boxplot of linkage disequilibrium of *Nitrospira* bacterial populations for whole genome (all *Nitrospira* and comammox *Nitrospira*) and nitrification genes. Differences between the mean linkage disequilibriums were assessed by a Dunn's test; same letter have means not significantly different from each other ($p < 0.05$). Right) Boxplot of pN/pS ratios of *Nitrospira* bacterial populations for whole genome (all *Nitrospira* and comammox *Nitrospira*) and nitrification genes. Differences between the mean pN/pS ratios were assessed by a Dunn's test; same letter have means not significantly different from each other ($p < 0.01$).

## Conclusions

A major unresolved question is how the relationship between ecology and evolution shapes complex communities in wild environments. Here we use a model microbial system to examine this question. Strain-level analyses enabled us to decipher the degree of intra-population diversity in wild dominant comammox *Nitrospira* inhabiting rapid sand filters and estimate important evolutionary processes such as recombination and selection in these populations. We showed that compared to other environments the *Nitrospira* populations in rapid sand filters are characterized by low recombination and strong purifying selection therefore, we conclude that the evolutionary processes that drive the diversification of *Nitrospira* are dependent on the local environment, as opposed to intrinsic properties of the species.

# Methods

## Sampling, sequencing, and metagenomic assembled genomes recovery

The sampling description, DNA extraction and sequencing have been previously described[22]. Briefly, filter material was collected from two locations at the top of the filters of 12 Danish waterworks. DNA was extracted from 0.5 g of sand material using the MP FastDNA Spin Kit (MP Biomedicals LLC, Solon, USA). DNA libraries were generated using the 24 extracted DNA with the Nextera XT DNA library preparation kit (Illumina Inc.) according to the manufacturer's instructions. The samples were sequenced in one lane, with 2×150 paired read sequencing on the Illumina HiSeq4000 at BGIs facility in Copenhagen. As previously described[22], high-quality reads were used for metagenomic assembled genomes (MAGs) recovery using a combination of automatic and manual binning followed by filtering and refinement steps to improve the quality of the MAGs. The resulted MAGs were dereplicated using dRep[42] with the secondary clustering threshold -sa 0.99. Dereplicated MAGs completeness and contamination was evaluated using CheckM[43]. MAGs with completeness > 50% and contamination < 10% were kept for downstream analyses.

## Species abundance estimation

A 95% average nucleotide identity (ANI) cut-off was used to define species as proposed by Klappenbach *et al.* (2007)[16]. The retrieved MAGs were dereplicated using dRep with the secondary clustering threshold set at 95% gANI. Among the genomes classified as belonging to the same species, the one with higher quality was chosen as representative genome for that species. The species abundance and coverage of each representative genome across the studied metagenomes was assessed using MIDAS [44]. Briefly, MIDAS uses reads mapped to 15 universal single-copy gene families (with ability to accurately recruit metagenomic reads to the correct species [44]) to estimate the abundance and coverage of bacterial species from a shotgun metagenome. We used the species retrieved in this study to build the database of universal-single-copy genes.

## Genome classification and annotation

MAGs were classified (Supplementary Table 10) using the classify workflow of the GTDB-Tk v.0.1.3 tool[45]. Open reading frames were predicted using Prodigal v. 2.63[46], and annotated using blastp[47] against NCBI nr[48], UniProt[49], KEGG[50], PFAM[51] and eggNOG[52]. Genes were assigned to antiphage defense systems using the strategy described in Doron *et al.* (2018)[53].

## Phylogenetic analysis

Phylogenetic analyses of *Nitrospira* genomes were conducted with the GTDB-Tk v.0.1.3 tool[45] using the *de novo* workflow with a set of 120 single copy marker proteins and the genome taxonomy database (GTDB)[54]. Concatenated alignments were used to construct a maximum likelihood tree using RAxML v. 8.2.11[55] with 400 rapid bootstraps (determined using the autoMRE option) and the LG likelihood model of amino acid substitution with gamma distributed rates and fraction of invariant sites (-m PROTGAMMAILGF; best model determined using ProtTest v. 3.4.2[56]). The tree was rooted using two *Leptospirillum* species as outgroup. The rooted tree was visualized using the online web tool from the Interactive Tree of Life (iTol)[57].

**Quartet analysis**

To retrieve waterworks-specific MAGs for each *Nitrospira* species, the reassembly module of metaWRAP[58] was used with individual reads from each sample and the representative *Nitrospira* species MAG. Resulted MAGs were kept for the quartet analysis if the completeness did not vary in more than 10% compared to the representative *Nitrospira* species MAG and the contamination remained < 5%. A species phylogenetic tree of the resulting *Nitrospira* MAGs was constructed as described above. For each quartet analysis we selected four *Nitrospira* MAGs. Four from the same *Nitrospira* species for within species analysis (the most phylogenetically distant MAGs), and two from one *Nitrospira* species and two from another one for between species analysis (the most phylogenetically distant MAGs). Orthofinder v. 2.3.3[59] was used to identify orthologous genes among each set of four genomes, retaining for subsequent analyses only single-copy orthologous genes. Orthofinder v. 2.3.3 with the options -M msa -T raxml was also used to produce phylogenetic trees for each orthologous gene. For each within or between species analysis, topological differences between each orthologous gene and the species tree were assessed by calculating the Robinson-Foulds (RF) distance[60] with the R function RF.dist of the phangorn package[61]. This analysis allowed to obtain, for each quartet, the percentage of single-copy orthologous genes phylogenetic trees which did not support the species phylogenetic tree topology.

**Read mapping, SNP calling, and population genomic analysis**

The population genomic analysis was done following the approach described in Crits-Christoph et al (2020)[23]. High-quality reads were mapped to an indexed database of the 176 species MAGs recovered from the waterworks using BWA-MEM[62]. Resulted alignments were filtered using samtools[63] view -q30 to remove reads with mapping quality less than 30, and also with the script filter_reads.py[23] (with the options: -m 96 to retain reads with a percent identity of at least 96% to the reference; and -q 2 to assure uniquely best mapping read pairs in the index). Downstream analysis was performed for 24 species genomes (all the 12 *Nitrospira* ones, and 12 other abundant species genomes). For each of these species genomes, we analysed its data in samples that passed a cutoff of at least 50% of the genome being covered with at least 5× coverage. 149 out of 576 sample genome comparisons (24 genomes × 24 samples) passed this minimum requirement. Sample read mappings were pooled by each waterworks and by all samples across the waterworks. Nucleotide diversity ($\pi$), linkage disequilibrium (D') and pN/pS ratio were calculated for each sample, each waterworks and across all the waterworks as described elsewhere[23] using the provided scripts. $F_{ST}$ was calculated following the same procedure but on sites segregating across two waterworks being compared (for all the possible waterworks comparisons). As Crits-Christoph *et al.* (2020)[23] recommended, only sites with a coverage of at least 20× in each waterworks was used to calculate $F_{ST}$. In addition, genes with coverages in a waterworks outside of the range of two standard deviations were excluded from the analysis. As previously suggested[23], a two-sample Wilcoxon test was conducted to find out if average linkage of highly differentiated loci differed from the genomic average for each species. Similarly, a two-sample *t*-tests was used to conclude if average nucleotide diversity of highly differentiated loci differed from the genomic average. Both sets of tests were corrected for multiple hypotheses using the Benjamini–Hochberg method.

449 Same strain-level analysis as the one described above was conducted in *Nitrospira* MAGs previously
450 recovered[22] that passed a cutoff of at least 50% of the genome being covered with at least 5× coverage
451 in any of the metagenomes were *Nitrospira* MAGs were found to be present[22].
452
453 **Statistical Analyses.**
454 All statistical tests were performed using R v3.5.2[64]. For all statistical analyses, species abundances
455 data was treated as followed: zeros were replaced with an estimate value using the Count Zero
456 Multiplicative approach with the zCompositions R package[65], and data were further centred log-ratio
457 transformed. *Nitrospira* community dissimilarities were calculated using the Jaccard index. The
458 correlation between the *Nitrospira* community dissimilarities and geographic distances was
459 calculated using the Mantel test (significance obtained after 100,000 permutations). Same analysis
460 was used to assess the correlation between the *Nitrospira* community dissimilarities and the water
461 composition dissimilarity, as well as the correlation between major allele dissimilarities and
462 geographic distances.
463 Proportionality between abundances of the species across the 24 metagenomes were calculated using
464 the propr R package[66] (with the options metric = "rho", ivar = "clr") and visualised using the corrplot
465 R package[67]. For the network analysis, the function getNetwork from propr R package was used to
466 retain proportionalities > 0.56 (FDR < 5%). The network was visualised using the igraph R package[68].
467 Same approach was used to build the network including phages but, in this case, proportionalities >
468 0.51 (FDR < 5%) were retained.
469 Redundancy analysis (RDA) was performed in R package vegan[69]. RDA was conducted using centred
470 log-ratio transformed *Nitrospira* species abundances and chemical data of influent water. The
471 constrained ordination model and the variable significance were determined by permutation tests
472 (1000 permutations) with anova.cca in vegan. Principal components analysis (PCA) was performed
473 in R package factoextra[70] using the nucleotide diversity, pN/pS ratio and D' values for the bacterial
474 populations retrieved from the waterworks, as well as most abundant bacterial populations across
475 grassland meadow[23], and other *Nitrospira* populations abundant in other systems.
476 Differences between the mean nucleotide diversities of the nitrifying genes, whole *Nitrospira*
477 genomes, and whole comammox *Nitrospira* genomes were assessed using Kruskal–Wallis ANOVA
478 followed by Dunn's test with the Holm-Bonferroni correction. Same analysis was performed for
479 linkage disequilibrium and pN/pS ratios.
480
481 **Recovery of draft phage genomes and abundance estimation**
482 MARVEL[71] was used to recover draft phage genomes from the co-assembly generated from the
483 waterworks (describe above). As recommended by MARVEL, Metabat[72] was run using the
484 parameters -m 1500 -s 10000 to produce bins with contigs of at least 1500 bp and with a minimum
485 total size of 10 kbp. Then, we executed MARVEL to identify phage bins from the 1026 bins generated
486 with Metabat. The abundance in each sample of the 43 identified draft phage genomes was estimated
487 using the quant module from metaWRAP[73] (Supplementary Table 11) .
488
489
490

**Chemical analysis of influent water**

Ammonium was measured using a standard colorimetric salicylate and hypochlorite method[74]. Iron and manganese content was determined by ICP-MS (7700x, Agilent Technologies). NVOC analysis was performed using a wet chemical TOC-analyser TOC-V WP (Shimadzu, Kyoto, Japan).

## Data availability

All raw sequence data and *Nitrospira* genomes retrieved from the Danish rapid sand filters have been deposited at NCBI under the project PRJNA384587. The rest of the retrieved draft genomes from the Danish rapid sand filters are available on figshare (https://doi.org/10.6084/m9.figshare.12962075).

## Acknowledgements

## Authors contributions

A.P conceived the study and performed the bioinformatic analyses. A.P, O.X.C and A.D led interpretation of the results supported by B.F.S. A.P drafted the manuscript with input from A.D, O.X.C and B.F.S. All authors contributed to manuscript revision, and approved the final version of the manuscript.

## References

1. Hairston, N. G., Ellner, S. P., Geber, M. A., Yoshida, T. & Fox, J. A. Rapid evolution and the convergence of ecological and evolutionary time. *Ecol. Lett.* **8**, 1114–1127 (2005).

2. Barroso-Batista, J. *et al.* The First Steps of Adaptation of Escherichia coli to the Gut Are Dominated by Soft Sweeps. *PLoS Genet.* **10**, e1004182 (2014).

3. Lawrence, D. *et al.* Species Interactions Alter Evolutionary Responses to a Novel Environment. *PLoS Biol.* **10**, e1001330 (2012).

4. Denef, V. J., Mueller, R. S. & Banfield, J. F. AMD biofilms: using model communities to study microbial evolution and ecological complexity in nature. *ISME J.* **4**, 599–610 (2010).

5. Palomo, A. *et al.* Metagenomic analysis of rapid gravity sand filter microbial communities suggests novel physiology of Nitrospira spp. *ISME J.* **10**, 2569–2581 (2016).

6. Gülay, A. *et al.* DNA- and RNA-SIP Reveal Nitrospira spp. as Key Drivers of Nitrification in Groundwater-Fed Biofilters. *MBio* **10**, (2019).

7. Hu, W. *et al.* Metagenomics Unravels Differential Microbiome Composition and Metabolic Potential in Rapid Sand Filters Purifying Surface Water Versus Groundwater. *Environ. Sci. Technol.* **54**, 5197–5206 (2020).

8. Tatari, K. *et al.* Density and distribution of nitrifying guilds in rapid sand filters for drinking water production: Dominance of Nitrospira spp. *Water Res.* **127**, 239–248 (2017).

9. Fowler, S. J., Palomo, A., Dechesne, A., Mines, P. D. & Smets, B. F. Comammox Nitrospira are abundant ammonia oxidizers in diverse groundwater-fed rapid sand filter communities.

534      *Environ. Microbiol.* **20**, 1002–1015 (2018).

535  10. Koch, H., Kessel, M. A. H. J. van & Lücker, S. Complete nitrification: insights into the
536      ecophysiology of comammox Nitrospira. *Appl. Microbiol. Biotechnol.* 1–13 (2018).
537      doi:10.1007/s00253-018-9486-3

538  11. Palomo, A. *et al.* Comparative genomics sheds light on niche differentiation and the
539      evolutionary history of comammox Nitrospira. *ISME J.* **12**, 1779–1793 (2018).

540  12. Dutta, C. & Paul, S. Microbial Lifestyle and Genome Signatures. *Curr. Genomics* **13**, 153–
541      162 (2012).

542  13. Scheuerl, T. *et al.* Bacterial adaptation is constrained in complex communities. *Nat.*
543      *Commun.* **11**, 754 (2020).

544  14. González-Torres, P., Rodríguez-Mateos, F., Antón, J. & Gabaldón, T. Impact of Homologous
545      Recombination on the Evolution of Prokaryotic Core Genomes. *MBio* **10**, (2019).

546  15. Martinez, J. L. The role of natural environments in the evolution of resistance traits in
547      pathogenic bacteria. *Proc. R. Soc. B Biol. Sci.* **276**, 2521–2530 (2009).

548  16. Klappenbach, J. A. *et al.* DNA–DNA hybridization values and their relationship to whole-
549      genome sequence similarities. *Int. J. Syst. Evol. Microbiol.* **57**, 81–91 (2007).

550  17. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
551      throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat.*
552      *Commun.* **9**, 5114 (2018).

553  18. Olm, M. R. *et al.* Consistent Metagenome-Derived Metrics Verify and Delineate Bacterial
554      Species Boundaries. *mSystems* **5**, (2020).

555  19. Louca, S. *et al.* Function and functional redundancy in microbial systems. *Nat. Ecol. Evol.*
556      (2018). doi:10.1038/s41559-018-0519-1

557  20. Gülay, A. *et al.* Internal porosity of mineral coating supports microbial activity in rapid sand
558      filters for groundwater treatment. *Appl. Environ. Microbiol.* **80**, 7010–7020 (2014).

559  21. Costa, E., Pérez, J. & Kreft, J.-U. Why is metabolic labour divided in nitrification? *Trends*
560      *Microbiol.* **14**, 213–219 (2006).

561  22. Palomo, A., Dechesne, A. & Smets, B. F. Genomic profiling of Nitrospira species reveals
562      ecological success of comammox Nitrospira. *bioRxiv* (2019). doi:10.1101/612226

563  23. Crits-Christoph, A., Olm, M. R., Diamond, S., Bouma-Gregson, K. & Banfield, J. F. Soil
564      bacterial populations are shaped by recombination and gene-specific selection across a
565      grassland meadow. *ISME J.* 1–25 (2020). doi:10.1038/s41396-020-0655-x

566  24. Shapiro, B. J. *et al.* Population Genomics of Early Events in the Ecological Differentiation of
567      Bacteria. *Science (80-. ).* **336**, 48–51 (2012).

568  25. Rosen, M. J., Davison, M., Bhaya, D. & Fisher, D. S. Fine-scale diversity and extensive
569      recombination in a quasisexual bacterial population occupying a broad niche. *Science (80-. ).*
570      **348**, 1019–1023 (2015).

571  26. Bendall, M. L. *et al.* Genome-wide selective sweeps and gene-specific sweeps in natural
572      bacterial populations. *Isme J* **10**, 1589–1601 (2016).

573  27. Shapiro, B. J., Leducq, J.-B. & Mallet, J. What Is Speciation? *PLOS Genet.* **12**, e1005860
574      (2016).

575  28. VanLiere, J. M. & Rosenberg, N. A. Mathematical properties of the measure of linkage
576      disequilibrium. *Theor. Popul. Biol.* **74**, 130–137 (2008).

577  29. Lin, M. & Kussell, E. Inferring bacterial recombination rates from large-scale sequencing
578      datasets. *Nat. Methods* **16**, 199–204 (2019).

579  30. Doroghazi, J. R. & Buckley, D. H. Widespread homologous recombination within and
580      between Streptomyces species. *ISME J.* **4**, 1136–1143 (2010).

581  31. Carja, O., Liberman, U. & Feldman, M. W. Evolution in changing environments: Modifiers

582          of mutation, recombination, and migration. *Proc. Natl. Acad. Sci.* **111**, 17935–17940 (2014).

583    32.   Hanage, W. P. Not So Simple After All: Bacteria, Their Population Genetics, and
584         Recombination. *Cold Spring Harb. Perspect. Biol.* **8**, a018069 (2016).

585    33.   Didelot, X. & Maiden, M. C. J. Impact of recombination on bacterial evolution. *Trends*
586         *Microbiol.* **18**, 315–322 (2010).

587    34.   Luo, H., Gao, F. & Lin, Y. Evolutionary conservation analysis between the essential and
588         nonessential genes in bacterial genomes. *Sci. Rep.* **5**, 13210 (2015).

589    35.   Dilucca, M., Cimini, G. & Giansanti, A. Essentiality, conservation, evolutionary pressure and
590         codon bias in bacterial genomes. *Gene* **663**, 178–188 (2018).

591    36.   Aguilar-Rodríguez, J. & Wagner, A. Metabolic Determinants of Enzyme Evolution in a
592         Genome-Scale Bacterial Metabolic Network. *Genome Biol. Evol.* **10**, 3076–3088 (2018).

593    37.   Zhong, C. *et al.* Pan-genome analyses of 24 Shewanella strains re-emphasize the
594         diversification of their functions yet evolutionary dynamics of metal-reducing pathway.
595         *Biotechnol. Biofuels* **11**, 193 (2018).

596    38.   Petersen, L., Bollback, J. P., Dimmic, M., Hubisz, M. & Nielsen, R. Genes under positive
597         selection in Escherichia coli. *Genome Res.* **17**, 1336–1343 (2007).

598    39.   Rabby, A. *et al.* Identification of the positively selected genes governing host-pathogen arm
599         race in Vibrio sp. through comparative genomics approach. *Biojournal Sci. Technol.* **2**,
600         (2015).

601    40.   Rodriguez-Valera, F. *et al.* Explaining microbial population genomics through phage
602         predation. *Nat. Rev. Microbiol.* **7**, 828–836 (2009).

603    41.   Cordero, O. X. & Polz, M. F. Explaining microbial genomic diversity in light of evolutionary
604         ecology. *Nat. Rev. Microbiol.* **12**, 263–273 (2014).

605    42.   Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. DRep: A tool for fast and accurate
606         genomic comparisons that enables improved genome recovery from metagenomes through
607         de-replication. *ISME J.* **11**, 2864–2868 (2017).

608    43.   Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
609         assessing the quality of microbial genomes recovered from isolates, single cells, and
610         metagenomes. *Genome Res.* **25**, 1043–55 (2015).

611    44.   Nayfach, S., Rodriguez-Mueller, B., Garud, N. & Pollard, K. S. An integrated metagenomics
612         pipeline for strain profiling reveals novel patterns of bacterial transmission and
613         biogeography. *Genome Res.* **26**, 1612–1625 (2016).

614    45.   Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to
615         classify genomes with the Genome Taxonomy Database. *Bioinformatics* (2019).
616         doi:10.1093/bioinformatics/btz848

617    46.   Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
618         identification. *BMC Bioinformatics* **11**, 119 (2010).

619    47.   Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
620         search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

621    48.   Sayers, E. W. *et al.* GenBank. *Nucleic Acids Res.* **47**, D94–D99 (2019).

622    49.   UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res.* **47**, D506–D515 (2019).

623    50.   Ogata, H. *et al.* KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **27**,
624         29–34 (1999).

625    51.   Punta, M. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **40**, D290–D301
626         (2012).

627    52.   Huerta-Cepas, J. *et al.* eggNOG 4.5: a hierarchical orthology framework with improved
628         functional annotations for eukaryotic, prokaryotic and viral sequences. *Nucleic Acids Res.* **44**,
629         D286–D293 (2016).

630  53.  Doron, S. *et al.* Systematic discovery of antiphage defense systems in the microbial
631       pangenome. *Science (80-. ).* **359**, eaar4120 (2018).
632  54.  Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome phylogeny
633       substantially revises the tree of life. *Nat. Biotechnol.* (2018). doi:10.1038/nbt.4229
634  55.  Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
635       phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
636  56.  Darriba, D., Taboada, G. L., Doallo, R. & Posada, D. ProtTest 3: fast selection of best-fit
637       models of protein evolution. *Bioinformatics* **27**, 1164–1165 (2011).
638  57.  Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the display and
639       annotation of phylogenetic and other trees. *Nucleic Acids Res.* **44**, W242–W245 (2016).
640  58.  Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP—a flexible pipeline for genome-
641       resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
642  59.  Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
643       genomics. *Genome Biol.* **20**, 238 (2019).
644  60.  Robinson, D. F. & Foulds, L. R. Comparison of phylogenetic trees. *Math. Biosci.* **53**, 131–
645       147 (1981).
646  61.  Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592–593 (2011).
647  62.  Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
648       *Bioinformatics* **26**, 589–95 (2010).
649  63.  Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078–
650       2079 (2009).
651  64.  Team, R. C. R: A language and environment for statistical computing. (2014).
652  65.  Palarea-Albaladejo, J. & Martín-Fernández, J. A. zCompositions — R package for
653       multivariate imputation of left-censored data under a compositional approach. *Chemom.*
654       *Intell. Lab. Syst.* **143**, 85–96 (2015).
655  66.  Quinn, T. P., Richardson, M. F., Lovell, D. & Crowley, T. M. propr: An R-package for
656       Identifying Proportionally Abundant Features Using Compositional Data Analysis. *Sci. Rep.*
657       **7**, 16252 (2017).
658  67.  Wei, T. & Simko, V. R Package 'corrplot': visualization of a correlation matrix. (2017).
659  68.  Gabor, C. & Tamas, N. The igraph software package for complex network research.
660       *InterJournal* **Complex Sy**, 1695 (2006).
661  69.  Oksanen, J. *et al.* Package 'vegan' Title Community Ecology Package. *Community Ecol.*
662       *Packag.* **2**, 1–297 (2019).
663  70.  Kassambara, A. & Mundt, F. factoextra: Extract and Visualize the Results of Multivariate
664       Data Analyses. *https://CRAN.R-project.org/package=factoextra* (2017).
665  71.  Amgarten, D., Braga, L. P. P., da Silva, A. M. & Setubal, J. C. MARVEL, a Tool for
666       Prediction of Bacteriophage Sequences in Metagenomic Bins. *Front. Genet.* **9**, (2018).
667  72.  Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for accurately
668       reconstructing single genomes from complex microbial communities. *PeerJ* **3**, e1165 (2015).
669  73.  Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP—a flexible pipeline for genome-
670       resolved metagenomic data analysis. *Microbiome* **6**, 158 (2018).
671  74.  Bower, C. E. & Holm-Hansen, T. A Salicylate–Hypochlorite Method for Determining
672       Ammonia in Seawater. *Can. J. Fish. Aquat. Sci.* **37**, 794–798 (1980).
673
674