COMPARISON OF AMPLICON SEQUENCING TO CONVENTIONAL FECAL SOURCE TRACKING TECHNIQUES IN THE NAVESINK RIVER

by

SARAH PHELAN

A thesis submitted to the

School of Graduate Studies

Rutgers, The State University of New Jersey

In partial fulfillment of the requirements

for the degree of

Master of Science

Graduate Program in Civil and Environmental Engineering

Written under the direction of

Dr. Nicole Fahrenfeld

and approved by

New Brunswick, New Jersey

October, 2018

ABSTRACT OF THE THESIS

Comparison of amplicon sequencing to conventional fecal source tracking techniques in the Navesink River By SARAH PHELAN

Thesis Director:

Dr. Nicole Fahrenfeld

Fecal indicator bacteria are commonly used to evaluate water quality and make decisions on designating and restricting use. A drawback of this method is it does not differentiate fecal sources for management purposes. Microbial source tracking is being explored as an alternative and more informative method of assessing and identifying contaminant sources. High throughput sequencing of the 16s rRNA region allows for rapid, large scale microbial community analyses and is an attractive source tracking method that has not been thoroughly explored. To investigate the effectiveness of this method, samples were collected from representative fecal sources from horse, dogs, geese, and wastewater then mixed in known quantities to create a fecal library. A series of field samples were collected during wet and dry weather from the Navesink River, Red Bank, New Jersey. The library and field samples were evaluated for fecal coliform, qPCR for fecal marker genes, and those results were compared with the microbial community fingerprints generated by amplicon sequencing. The results from all three techniques were cross compared to understand the consistency of results between methods, and will inform the application of amplicon sequencing for source tracking in surface waters impacted by

ii

fecal/bacterial contamination. Results showed significant differences in library samples with biomarkers that could be applied with success to select sites with elevated coliform results; indicate further biomarker investigation is needed in sequencing results.

Acknowledgements

Thanks to William R. Morales Medina, Disha Soni, Kris, Parker, Sheri Elsaker, and Alex Massie for sampling and analysis assistance. Thanks to NJ WRRI for the grant-in-aid funding to perform this research.

Table of Contents

Abstractii
1. Chapter (1) 1
1-1- Fecal Contamination is a widespread water quality issue
1-2 State of science for microbial source tracking methods4
1-3 The Navesink River: a microbially impaired waterway4
1-4 Research questions
2. Chapter (2)
2-1- Introduction
2-2- Materials and Methods 11
2-2-1- Fecal library samples
2-2-2- Wet and dry weather samples 12
2-2-3- Coliform analysis 14
2-2-4 Biomarker analysis14
2-2-5 Statistical analysis 16
2-3- Results
2-3-1 Performance of source tracking on fecal library 17
2-3-2 Fecal contamination and source tracking in the Navesink
2-3-3 The Relationship between fecal coliform, qPCR, and amplicon sequencing28
2-4- Discussion
2-4-1 Relationship between fecal markers, pathogen sequences, and coliform counts

2-4-2 Source tracking and implications for the Navesink	
2-5- Conclusions	
3. Chapter (3)	
3-1- Conclusion	
References	39
Supplemental Figures	43

Table of Figures:

Figure 1- Map of the Navesink showing shellfish restrictions	5
Figure 2- Map of sampling locations 1	3
Figure 3-qPCR results for library and field samples	9
Figure 4- Cluster analysis and relative abundance of community in library samples 2	2
Figure 5- Coliform results for wet and dry weather samples	23
Figure 6- Cluster analysis and relative abundance of community in field samples2	27
Table of Tables:	
Table 1-Summry of biomarker, qPCR, and coliform results for library and field	

1. Chapter (1): Broad Introduction

1-1 Fecal contamination is a widespread water quality issue

In the United States, water quality impairment due to the presence of pathogens is an ongoing and well documented human health risk. In 2017 the USEPA reported that out of 1.1 million miles of rivers and streams surveyed, approximately 614,000 miles were classified as impaired, and an additional 5,500 miles were acceptable but deteriorating (USEPA, 2017). A leading cause cited was pathogens, as indicated by the presence of E. *coli*, total coliform, fecal coliform, and *Enterococcus*. Twenty three percent of rivers reported had *Enterococci* levels that exceeded human health limits. One of the runoff sources most frequently cited in reported water quality impairment was agricultural activities. Combined sewer overflows (CSOs) are another event that results in the degradation of water quality through the release of fecal bacteria, with human health implications (Donovan et. al, 2008). Wildlife contamination from avian feces is a commonly documented case in waterways (Green et al, 2011, Eramo et al, 2017) however these sources are less likely to cause disease than human waste (Schoen, M.E. et al., 2010). A variety of human and wildlife sources have the potential to impact waterways and can increase non-specific fecal indicators.

The NJDEP 2014 Integrated Water Quality Assessment Report cited pathogens as one of the leading causes of impairment, consistent with national reporting (NJDEP 2014). In the biannual water quality assessment and in site specific assessments, the NJDEP measures total coliform, fecal coliform, and *Enterococcus* as indicators of pathogen pollution. These parameters are used in the assessment and limitation of water quality for recreational use (NJDEP, 2006). Any water quality impairments identified through water quality assessments require the calculation of a Total Maximum Daily Load (TMDL),

2

which is the maximum amount of the pollutant the water body can receive to still meet water quality standards. These loadings are site specific and are calculated based on point and nonpoint sources of pollution, background pollutant levels, and current/future use of the water body (NJDEP, 2006). TMDL limits are beneficial because they target nonpoint and stormwater point sources and the loading capacity is allocated to these sources in an effort to provide balanced and measurable reduction goals for pollutant source sites (NJDEP, 2006). A drawback of these limits is that in their development, they do not discern between sources of the indicator bacteria which can be non-specific, and rely on default fecal coliform loading rates that may not reflect actual loadings from sources (NJDEP, 2006). TMDL limits also focus on fecal indicators that do not necessarily correlate with pathogen concentrations, and may target sources that do not adversely impact human health (Gilmore et. al, 2014).

The widespread occurrence of water quality issues relating to the presence of fecal impacts highlights the need for effective and accurate characterization. Fecal indicator bacteria fall short because they do not provide information on the source of contamination. *Enterococci* for example does not always correlate to human fecal pollution- *Enterococci* are also found in environmental and animal fecal sources (Boehm and Sassoubre, 2014). Being able to differentiate between fecal sources is important because 1) it can help locate where the pollution is originating allowing for remediation, and 2) multiple epidemiological studies have shown that waterborne diseases are more likely to occur in people when sewage based fecal contamination is present (Schoen, M.E. et al., 2010).

1-2 State of science for microbial source tracking methods

Several chemical and microbial methods have been and are currently being explored to identify the source of fecal contamination (Meaysa et. al, 2004). Chemical source tracking methods such as Multiple Antibiotic Resistance differentiate sources using human and livestock associated antibiotics. Another method that has been used is an analysis for the presence of caffeine as an indicator of human pollution. This method is specific to human activity but can be costly and has shown to be non-sensitive in some studies (Hagedorn, 2001). Optical brighteners are another indicator of human wastewater contamination and are prolific in household greywaters, found in bleaches and detergents. This method is fast and inexpensive but because it is a chemical based method, it does not provide any information on the age of the contamination. Microbial source tracking methods may be cultivation based or target biomolecular compounds. Cultivation based methods include, for example, multiple Antibiotic Resistance (MAR), an established cultivation-based source tracking method that provides identification of multiple sources through profiling the antibiotic resistance profiles of isolates (Kasper, 1990). Benefits of this method are that it can discriminate from multiple sources and is fast. Drawbacks are that this method is specific to geographic regions where certain antibiotics are used, and it has been shown to be difficult to differentiate with this method in mixed samples (Hagedorn, 2001)

1-3 Navesink River: a microbially impaired waterway

The Navesink River is an impaired estuary located in Monmouth county, New Jersey, with restrictions on recreational swimming and shellfish harvesting. The river contains

4

2,290 acres of shellfish growing waters, for hard and soft clams. Ninety-Five square miles drain into the watershed and are primarily comprised of urban/suburban (39%), wetlands (20%) and agricultural (15%) land use areas. Permitted wastewater treatment discharges to the river are located



FIGURE 1: Map showing composition of the Navesink Watershed with shellfish classification restrictions (NJDEP, 2012)

either in upstream parts of the river or in the Atlantic Ocean (NJDEP, 2006). A 2006 study by the NJDEP classified the upstream portion of the Navesink (Estuary B) as impaired and classified 152 acres of the upper estuary as "prohibited" for shellfishing (Monmouth County, 2010). TMDLs were established for the Navesink River to regulate permitted stormwater discharges for pollutant sources to the river. For the shellfish impaired portion of the Navesink (Navesink Estuary B), the annual TMDL was set at 1.26×10^{15} cfu/yr. To achieve this limit, a 93% reduction in stormwater outputs from agricultural, marina, and urban sources was required (NJDEP, 2006).

In 2008, the NJDEP bureau of Marine Water Monitoring (BMWM) produced a report on storm studies, identifying that rain events correlated with high coliform results, indicating stormwater runoff was a source of coliform pollution, highest in the upstream portion of the river (NJDEP, 2008). The investigation by the DEP in 2008 included microbial source tracking using Multiple Antibiotic Resistance Sampling (MAR), optical brightening, and F+ RNA coliphage. The study identified multiple sources of elevated

coliform including both human and wildlife at various sites along the river (NJDEP, 2008). The Monmouth county health department coordinated with municipalities along the Navesink to identify and address many sources. In 2015 after additional data review, an additional 565.7 acres of river were downgraded from "Special Restricted" to "Prohibited" for shellfish harvesting. The Navesink River, in addition to the downstream Shrewsbury River, supports almost all soft clam fishery in New Jersey (NJDEP, 2006) - this downgrade has potential economic impacts within the state of New Jersey. In addition to these restrictions, the Navesink is listed in a 2012 Water body report as impaired for Fish Consumption, Primary Contact Recreation, and Aquatic Life (USEPA, 2012).

1-4 Research questions

Given that fecal contamination is wide spread and that current methods have limitations with respect to identifying fecal sources, there is motivation to use the Navesink as a case study for evaluating accepted and emerging microbial source tracking methods. Identifying contamination sources is critical because it allows stakeholders (including regulators, interest groups, and point source facilities) to better understand and quantify existing contaminant sources. This will allow for the effective application of best management practices and infrastructure upgrades where they are most needed. Additionally, identifying fecal contamination from sewage (versus wildlife or domestic animal sources) is important because multiple epidemiological studies have shown that waterborne diseases are more likely to occur in people when sewage based fecal contamination is present (Schoen et al, 2010). This drives the need for source tracking methods such as amplicon sequencing to better identify, risk-rank, and manage sources. The research presented seeks to answer four critical questions:

(1) What are the current sources of fecal contamination in the Navesink during wet and dry weather?

(2) Are fecal sources from humans, wildlife, and domestic animals unique, in terms of microbial community structure, and can these differences be identified through community profiling?

(3) How can these differences be applied in field sampling to successfully identify fecal contaminant sources?

(4) Are sequencing results consistent with results obtained using conventional fecal methods (Fecal Indicator Bacteria, and qPCR) and previous environmental studies at the site?

Using the Navesink as a case study, this research drives to answer these questions and provides information to support future environmental investigation and source tracking studies.

2. Chapter (2): Introduction to the research, Experiments Description, Results and Discussion

2-1 Introduction

Fecal contamination degrades surface water quality but solving the problem is complicated by the fact that fecal pollution has multiple potential sources that are not differentiated by cultivation-based regulatory methods. Being able to differentiate between fecal contaminant sources is critical to identify and manage non-point pollution. This drives the need for better methods for source identification of fecal indicator bacteria.

A variety of methods have been and are currently being explored to identify the source of fecal contamination (Meaysa et. al, 2004). Advantages of biomolecular techniques include that they are less labor intensive and pose a lower biosafety risk than cultivation of fecal microbes. Among the available biomolecular techniques, polymerase chain reaction (PCR) and quantitative polymerase chain reactions (qPCR) for fecal marker genes are well demonstrated source tracking methods (Boehm et al. 2013). Assays exist for a variety of source-specific target genes (e.g., human, dog, horse, gull, pig, and others) and comprehensive evaluations have been performed demonstrating which assays have proven both specific and sensitive (Boehm et al. 2013). qPCR has been found to be a viable method for identifying sources of fecal contamination in water (Silkie and Nelson, 2009). However, while qPCR quantifies target fecal marker genes, these results have not been found to correlate with fecal indicator organisms (Boehm et al. 2009). Another approach is microbial community fingerprinting for which several techniques have been applied including denaturing gradient gel electrophoresis (DGGE) (Sigler and Pasutti, 2006) and terminal restriction fragment length polymorphism (t-RFLP) (Fogarty

and Voytek, 2005) but these have fallen out of favor given the lowered cost of DNA sequencing.

High throughput sequencing is a potentially attractive alternative technique for identifying the source of fecal indicator organisms. Amplicon sequencing of the 16S rRNA gene is no longer prohibitively expensive and is a feasible method for fast, large scale microbial community evaluations (Stenuit et. al, 2008, McLellan et. al, 2014). A major potential benefit of amplicon sequencing is that it provides a large amount of information the microbial community structure in a single assay with less labor than cultivation based methods. Programs such as SourceTracker (Knights et. al, 2011) and Biomarker analysis with the LEfSe galaxy database (Segata et al, 2011) use statistical, Bayesian algorithms and Linear Discriminant Analysis to process sequencing data and identify indicator organisms within sample sets. Amplicon sequencing, in conjunction with these robust statistical tools, has the potential to be used to identify a broad range of information on contaminant sources and expand the understanding of microbial communities beyond indicator microbes such as E. coli and Enterococci. Field studies have used amplicon sequencing to identify bovine and human indictor bacteria in Michigan (Wu et al. 2018) and evaluate waterfowl and human fecal waste, and compared bacterial populations to multiple site locations along Australian waterways (Henry et al., 2016). Multiple other studies conducted to date have investigated bacterial community profiling using amplicon sequencing to identify indicators of fecal contamination, with varying success (Kirs et al, 2017, Cao et al, 2013, Bradshaw et al, 2016, McCarthy et al, 2017). Researchers for these studies have not performed comparison of multiple animal

and human sources (human, horse, goose, and dog) with mixed library samples for comparison to surface water samples and validated these results using qPCR.

The aim of this work is to determine (1) if amplicon sequencing distinguishes clearly between samples of wastewater and fecal material from different animal sources in a fecal library, and (2) to compare the results of amplicon sequencing to established qPCR methods for fecal marker genes in field samples. Samples were collected from representative fecal sources from horse, dogs, geese, and wastewater then mixed in known quantities to create a fecal library. A series of field samples were collected during wet and dry weather from the Navesink river, a New Jersey waterway with known fecal contamination (NJDEP, 2015, COA, 2016). The library and field samples were evaluated for fecal coliform, qPCR for fecal marker genes, and those results were compared with the microbial community fingerprints generated by amplicon sequencing. The results from all three techniques were cross compared to understand the consistency of results between methods. The results of this will inform the application of amplicon sequencing to research on contaminated surface waters across the U.S., to evaluate and determine the source of any fecal/bacterial contamination.

2-2 Materials & Methods

2-2-1 Fecal Library Samples

To create a representative 'library' of samples from mixed sources, fecal samples were collected from multiple locations in NJ (Table S1 and Figure 1) to represent dogs, horses, and geese. All samples were stored in sterile bottles or ziplock bags on ice until processing in the laboratory up to 6 hrs later.

Wastewater influent (2L) was collected from a wastewater treatment plant located in northern NJ, as representative of human feces and sewage. Surface water samples (11L total volume) were collected from two sites on the Navesink River, which based on previous data (NJDEP Sampling data) had little to no fecal coliform present. The sites selected were NJDEP Site 32 and NJDEP Site 58 (Fig. 1). Approximately 5.5 L of sample were collected from each site in sterile bottles. The trip blank consisted of De-Ionized water was prepared and taken to and from the sampling sites during the sample events. All samples were stored on ice until processing, up to 6 hrs.

The fecal library consisted of 0.9L of surface water with varying volumes of WW (0.1L) and/or masses of homogenized feces (1g) (Table S2). The surface water used for the library creation was a 1:1 (v:v) ratio of samples from Sites 34 and 58. The library mixtures were homogenized with a 10 speed blender (Black and Decker BL2010BG). Each library mixture was prepared in duplicate and preserved for coliform or biomolecular (DNA) analysis, as described below. Between distinct mixtures, the blender was triple cleaned with a mixture of 10% bleach and alconox soap, and spatulas were flame sterilized.

2-2-2 Wet and Dry Weather Samples

Wet weather samples were collected during a storm event from surface water in various locations along the Navesink River. Sites sampled during the wet weather event were: NJDEP Site 10, Site 14, Site 34, Site 52, and Site 56 (Figure 1). From previous sampling events, samples were expected to contain a mixture of human and wildlife signatures from urban, impervious, and rural sources. At the time of sample collection (8/29/17,

1PM - 4 PM), total precipitation was approximately 0.1" of rainfall. During the rain event duplicate samples (1L) were collected from each location. Historical rainfall data was collected from Weather Underground (The Weather Company, 2017).

Dry weather samples were collected on 10/15/17 8AM-12PM from the same surface water locations along the Navesink River as the wet weather samples. Additionally, one downstream sample was collected from NJDEP Site 58. Duplicate (1L) samples were collected from each location. All samples were stored on ice until processing in the laboratory. A trip blank consisting of autoclaved deionized water was prepared and taken to the field during the trip and latter analyzed for QA/QC.



FIGURE 2: Map of sampling locations for wet weather sampling and dry weather sampling along various parts of the Navesink River.

Red Bank tidal data for each sampling event was collected from tides.mobilegeographics.com (figure S1) (Mobile Geographics.) Spike sampling, wet and dry weather sampling events were conduced roughly midway between high and low tide.

2-2-3 Coliform Analysis

Coliform analysis was performed for the library and field samples. Both the library and field samples from the wet weather event were analyzed at a certified laboratory (NJ Analytical Lab, Location) for total coliform analysis using analytical method SM9222B as described in standard methods (APHA, 2015). Samples collected during the dry weather sampling event were analyzed using analytical method EPA SM1604 (EPA, 2002). Briefly, 8-112 mL of sample from each site was filtered through 0.45 µm membrane filter (Cellulose Ester Membrane, Millipore Sigma). The filters were then placed on 5 mL plates of MI agar and incubated at 35°C for 24 hours. The colonies that grew within that time were then visually checked for the presence of blue color from the breakdown of IBDG by the E. coli enzyme 4-glucuronidase, and under UV light breakdown of MUGal by the TC enzyme 4-galactosidase. Observed colonies were counted and recorded to determine total coliform per 100 mL.

Each sample was filtered twice, using lower and higher volumes of samples to obtain a countable number of colonies (Table S3). Sample volumes were selected based on the coliform results from the wet weather sampling event.

2-2-4 Biomolecular Analyses

Field and library samples were filter concentrated (0.22 μm nitrocellulose filters (Millipore Corporation, Billerica, MA, United States) and stored at -20 C until DNA

extraction. DNA was extracted from filter concentrated samples (wet weather, dry weather, and library samples), using a commercial kit (FastDNA Spin Kit for Soil, MP Biomedicals, Hurcules, CA) following the manufacturer's directions. DNA extracts were stored at -20 C until analysis. To determine the presence of human and horse, feces, qPCR was performed on all samples using HF183 and BacHum as a human fecal indicator (Seurinck et al., 2005, Kildare et al, 2007) and HOF597 as a horse fecal indicator (Dick et al., 2004). qPCR was also performed for the 16S rRNA gene as an estimate of total bacterial population (Suzuki et al. 2000). A standard SYBR Green (5 μ L SsoFast EvaGreen, Bio-Rad, Hercules, CA, United States) chemistry with 0.4 μ M forward and reverse primers, and 1 μ L diluted (1:100) DNA extract in a 10 μ L reaction was used for all genes. A summary of the primers used is in table S4. QA/QC performed during qPCR included analyzing a no-template control on each plate, a seven-point calibration curve, and melt curve and/or gel electrophoresis to verify the specificity of qPCR products.

To determine the microbial community present in the fecal library and surface samples, amplicon sequencing (Illumina MiSeq, 300 bp, paired end) was performed targeting the V3-4 region of the 16S rRNA gene at a commercial lab (MR DNA, Shallowater, TX, United States). Sequences were processed using the QIIME v. 1.9.1 (Accessed 2/2018) (Kuczynski et al, 2012). Briefly, sequences were trimmed using Trimmomatic 0.36 ((Bolger et al. 2014) prior to joining the paired ends using Pandaseq (Masella et al. 2012). After demultiplexing (split_libraries.py), chimeras were removed (identify_chimeric_seq.py) with UCHIME (Edgar et al., 2011) the reads were parsed into operational taxonomic units (OTUs) within a 97% sequence identity cutoff (pick_de_novo_otus.py). Chloroplasts and mitochondria were removed by filtering in QIIME (filter_taxa_from_otu_table.py). To allow for comparison between samples with a different number of sequences, the samples were rarified to the lowest sequencing depth within the samples (32404 sequences). Rarefaction was performed in QIIME, and alpha rarefaction curves were generated (Figure S1)

2-2-5 Statistical Analyses

A Bray–Curtis similarity matrix was calculated on log-normalized subsampled (N =32,404 sequences) operational taxonomic unit data at the class level followed by cluster analysis with a SIMPROF test and non-metric multidimensional scaling (nMDS) in PRIMER 7. To determine which OTUs were preferentially associated with a given sample type from the library samples (human, horse, canine, or goose) biomarker analysis was performed on class-level relative abundance data for the library samples. The linear discriminant analysis effect size (LEfSe) tool (Segata et al., 2011) was used to identify biomarkers using relaxed parameters. A Wilcox Rank Sum test was performed for non parametric data to compare wet and dry weather sampling results for BacHum and HF183 results. BacHum and HF183 results were compared using a Kruskal-Wallis test with a post hoc-pairwise t-test with a Bonferronni correction for all library samples to compare wastewater and non-wastewater spiked samples. Coliform samples from wet and dry weather samples were compared using a Welch two sample t-test. Shannon Diversity Index results for sequencing were compared using a Wilcox Rank Sum test for library, and wet and dry field samples.

2-3 Results

2-3-1 Performance of source tracking methods on the fecal library *qPCR* for fecal marker genes

qPCR was performed human and horse fecal marker genes on the fecal library to confirm performance of these assays on samples with known fecal content. Two fecal marker genes were tested for human (HF183 and BacHum) fecal signatures. The surface water used to create the library samples was selected based on historical reports that the sites selected had no fecal coliform contamination, but 600 CFU/100 mL total coliform were observed in the samples used to generate the library. Amplification with the BacHum and HF183 assays was observed in all library samples (Fig. 3), including the two surface water samples used to create the fecal library. Both the BacHum and HF183 results were on average 2.3-2.5 log copies/mL higher samples spiked with wastewater (Table S5), however differences were not significant (BacHum p=0.3, HF183 p = 0.23). Amplification was not observed for BacHum or HF183 in the field blank, which was processed in parallel with the library samples. Comparing the two different human marker genes, BacHum results were higher than HF183 for qPCR results for 17 of 20 library samples and replicates, but differences were not significant (p = 0.32). qPCR for 16S rRNA gene copies was performed as a surrogate for total bacterial population. Amplification of the 16S rRNA gene copies was observed in all library samples without significant differences.

qPCR was also performed with one horse fecal marker gene. The qPCR primers for horse-associated fecal indicator bacteria resulted in specific and non-specific

amplification. Therefore, results were analyzed for presence/absence of the HOF597 gene following confirmation of the correct PCR product length via gel electrophoresis. The horse fecal marker was observed in five out of five samples where horse manure was added and four out of five samples where horse manure was not added including the surface water used to create the library (Fig. 3). One of the samples where horse manure was added showed amplification in only one out of two replicates. Samples where amplification was not observed contained wastewater and goose feces (SWG) and one of the wastewater, goose, horse, dog (SWGHD) library replicates.



FIGURE 3: qPCR results for a. library and b. field samples using human fecal marker genes HF183 and BacHum and 16S rRNA gene for total bacterial population. Error bars represent high and low values of replicate (N=2) samples. Presence of the horse fecal marker gene HOF597 is indicated by solid black circles, observation in one of two replicates with half filled circles, and absence by white circles. Dashed lines represent trip blank (TB-2) concentrations of HF183 and BacHum.

Amplicon sequencing for fecal finger printing

Amplicon sequencing of the V3V4 variable regions of the 16S rRNA gene was performed to study the microbial community structures of the fecal library to determine if these would provide sufficiently discriminating features to be useful for fecal source tacking. Twenty-one fecal library samples were sequenced resulted in 38,070-69,636 total sequences (Table S6). Sequences were subsampled to provide an equal number of sequences per sample for use in comparisons and for each sample 52 to 93 unique OTUs were defined at the class level (Table S6). The Shannon diversity index was calculated as an indicator of community diversity and evenness and for the library samples ranged from 7.3 to 11.4 (Table S6). To demonstrate that sequencing was performed to suitable depth, rarefaction curves at the class level are included as Fig. S2. The relative abundances of archaeal and bacterial community members at the genus level were determined for the library samples and cluster analysis was performed to help determine whether the library samples could be discriminated based upon their fecal content (Fig. 4). The microbiome of surface water samples used to create the fecal library (collected from S-34 and S-58 and homogenized) was composed primarily of Cyanobacteria (20.5%), Actinobacteria (15%), Bacteroidetes (21%), and Proteobacteria (28%). The surface water samples used to create the library were 47.1% similar to library samples spiked with fecal material, and clustered as more similar (up to 73.9%) to surface water samples collected during the field sampling study (Fig. 4). Several library samples clustered as significantly different based on fecal material added. Replicate samples from the same library type clustered with an average of 79% similarity between replicate samples, for seven out of ten library samples. Five of the library replicate samples were not significantly different from one another, based upon a SIMPROF test (all p>0.7). Library samples showing significant differences between replicates included four out of five samples with horse manure: SWH, SWDH, SDGH, SWGHD. Samples spiked with wastewater had increased levels of firmicutes, and of classes Bacteroidia and Epsilonproteobacteria compared to surface water samples.

Bacterial families Bacteroidaceae, Lachnospiraceae, Ruminococcaceae,

Porphyromonadaceae, and Prevotellaceae are associated with fecal material (Newton et. al 2015, McLellan et al. 2015) and were detected in elevated relative abundance in samples spiked with fecal material and wastewater, compared to surface water samples (Fig. S3). Within the fecal library, samples SW had the highest relative abundance of *Porphyromonadaceae* and *Provotellaceae*. SWH replicates had higher levels of

Ruminococceae, while *Lachnospiraceae* was observed at highest levels in SWD and higher than SW in many other mixed fecal samples.

Biomarker analysis was performed with the Huttenhower Lab LEfSE Linear Discriminant Analysis tool on bacterial classes, to identify potential indicator organisms of a library sample type. An initial analysis was performed using strict parameters and yielded no discriminating features. The analysis was performed on the fecal library samples using relaxed parameters (pairwise comparisons performed among subclasses with different names, one-against all multiclass analysis). From this, 13 discriminating features were determined associating microbial classes with different fecal library samples (Fig. S4). Surface water samples were characterized by *Acidiobacteriia*, Caldithrixae, Deltaproteobacteria, ABY1, and Phycisphraerae. Surface Water biomarkers were present in library samples at concentrations ranging from 0.13 to 0.65% relative abundance (Fig. S5). Library samples spiked with wastewater were characterized by Betaproteobacteria, Synergistia, and Theromicrobia, with wastewater and horse manure by Spirochaetes, T7_3, and MVP_15, and with wastewater and dog feces by Coriobacteriia. Wastewater biomarkers (11.2% of the SW sample) were present in the library samples with wastewater in relative abundances ranging from 1.6 to 8.3% and library samples without wastewater from 0.8-2.4% (Fig. S5). Four of the eight library samples containing wastewater had SW biomarkers present with a lower relative abundance than the "S" library sample, suggesting that the biomarkers identified at the relative abundance in the class level were not specific enough to distinguish samples containing wastewater from samples that did not (Fig. S6). Similar inconsistencies were identified when comparing biomarkers for surface water, wastewater and horse mixtures

(Fig. S7). Library samples containing dog feces had higher relative abundance of *Coriobacteriia* than samples that did not contain dog feces, suggesting the potential utility of this organism as an indicator of contamination with dog feces (Fig. S8).



FIGURE 4 **a**. Cluster analysis at the genus level and b. relative abundance of major (>5%) bacterial and archaeal classes in the fecal library and surface water used to dilute the fecal samples. Red bars connect samples without significant differences and black bars samples with significant differences (all p=0.001). (N=32404 sequences per sample). Note that the surface water used to create library samples clustered more closely with field samples, but is included here for comparative purposes.

2-3-2 Fecal contamination and source tracking in the Navesink

Water quality

Elevated fecal coliform counts were observed in both wet and dry weather sampling events (Fig. 5). Rainfall data is included as supplemental data (Fig. S9). The average relative percent difference for field replicates was 37.9% during wet weather sampling and 60.2% during dry weather sampling (Table S7). The highest quantifiable coliform results during both sampling events were observed at site S-10, which had 3,100 CFU/100 mL in wet weather and 2,800 CFU/100 mL in dry weather. These results were ten times higher than any other quantifiable result collected during wet weather sampling and more than three times higher than any other result during dry weather sampling. Two samples collected during the dry weather sampling event were too numerous to count: samples from site S-56 and S-14. Site S-58 was only sampled during the dry weather sample was collected the furthest downstream near where the Navesink joins Shrewsbury River (Fig. S10). All other samples had more than double the coliform observed at site S-58.



FIGURE 5. Land use map and bubble plots of the average a. wet and b. dry weather coliform results for field samples in the Navesink River. Location of site in NJ provided in Fig. S10.

qPCR observations for field samples

qPCR was performed on field samples to perform microbial source tracking with an established method and to allow for comparison to source assignments from the sequencing results. qPCR was performed on wet and dry weather field samples for HF183, BacHum, and HOF597 maker genes (Fig. 3b). Amplification with the HF183 and

BacHum assay was observed in all field samples (except for S-10 Wet, which did not amplify for BacHum). Amplification of both of these human maker genes was observed in the trip blank, despite the fact that fecal coliform were not observed in the trip blank. Therefore, only samples with fecal markers observed at greater concentration the blank are considered positive hits. Using this logic, HF183 was elevated in dry weather samples from sites S-10 and S-34 and wet weather samples from S-34. BacHum was elevated in dry weather samples from sites S-10, S-14, S-34, S-52 and S-56 and wet weather samples from S-52. Six of eight of the samples where elevated HF183 and BacHum were observed had elevated coliform results (>200 CFU/100 mL or TNTC). Overall, there were no significant differences observed between wet and dry weather sampling with the BacHum assay (p = 0.19), and the HF183 assay (p = 0.49). Next, the qPCR results from the fecal library were compared to the field samples. Amplification of HF183 was higher at site S-34 during wet weather than six of the library samples, including four library samples spiked with wastewater. Otherwise, the eight of the library samples spiked with wastewater had higher BacHum results than the field samples.

Amplification with the horse assay was observed in the all of the wet weather and in four of the six dry weather samples collected (Fig. 3b). Amplification of HOF597 was not observed in the trip blank. The horse marker gene was observed in during wet weather at sites S-10 and for both wet and dry weather at S-14 and S-34. The horse assay amplified for one replicate during wet weather at S-52 and S-56 and during dry weather at S-56 and S-58. Two of the samples with horse fecal marker genes did not have elevated collform results (>200 CFU/100 mL): wet weather at site S-34 and dry weather for site S-58.

Coliform results at all other locations where horse markers were identified were >200 CFU/100 mL. The horse fecal marker gene was not observed at site S-10 during dry weather, which had the highest coliform results during the dry weather sampling event. Significant differences were not observed based on weather conditions during sampling, as only two samples and one replicate from the dry weather sampling event were missing the horse marker gene. Comparing field and library results, both sample types resulted in amplification of the horse marker gene for most samples collected.

Amplicon sequencing for fecal fingerprinting

Amplicon sequencing was performed for the field samples to determine if the microbial community fingerprints could be useful in fecal source tracking. The number of sequences prior to subsampling, number of OTUs after subsampling, and rarefaction curves are included as Table S6 and Fig S2. Shannon diversity indices for the field samples ranged from 9.5-11 (Table S6). There were no significant differences between library and field samples, with respect to the Shannon diversity index (p=0.61). All of the replicate samples collected for a given field site and weather condition clustered without significant differences (average similarity for significant clusters 80.6%). No significant clusters were identified outside of replicate samples.

First, it is of interest to determine whether the microbial community structure of the field samples contained bacterial classes known to contain fecal microbes. Waterborne pathogens *E. coli, Salmonella, Vibrio,* and *Shigella* are Gammaproteobacteria, along with many other commensal organisms, and are specific pathogens of interest when discussing fecal impacts to sourcewater (Myers et al, 2014). Data were reviewed for the presence of

these organisms. Vibrio was identified in all library and field samples in relative abundances ranging from <0.1 to 13%. Vibrio was highest in field samples S-58 Dry (13% relative abundance), S-14 dry (2%), and S-56 dry (3%) (Fig. S11). In all other library/field samples, the relative abundance of Vibrio was less than 1%. Clostridium *perfrigens* is an alternate indicator bacteria of sewage contaminated waste is a part of the class Clostridia (Myers et al, 2014), which was present in 12 of the library samples and in none of the field samples at relative abundances >5%. The genus *Clostridium* was identified in all library and field samples but at low relative abundances, with the highest result identified in library sample SWD (0.25%). Library samples had higher Clostridium relative abundance than field samples. The highest field sample result was at S-10 Dry (0.05%). Enterococci and fecal Streptococci are two other indicator organisms commonly associated with fecal pollution (Myers et al, 2014). Enterococci are a subgroup of fecal Streptococci, and fecal Streptococci may be associated with non-human fecal contamination (Myers et al, 2014). Streptococci was present in relative abundances ranging from 0% to 0.02% in the field samples and 0.12% to 1.4% in library samples. The highest *Streptococci* results were observed in the SW Library samples (1.4%). The highest relative abundance of *Streptococci* in the field samples was S-10 Dry (0.02%). *Enterococci* were present library and field samples in low relative abundance [highest results were observed in sample SWGD (0.75%) and SWG (0.62%)]. All field samples contained Enterococci in abundances of less than 0.01%

Next, the fecal library and field sampling amplicon sequencing results were compared to determine whether the total bacterial community structure could be useful in source tracking. The field samples were significantly different from library samples containing

spikes of fecal samples (only 47% similarity, p=0.001). Library samples were distinct from field samples in that they contained *Firmicutes (Clostridia* and *Bacilli)* in all samples except for one at relative abundance levels >5%, none of the field samples contained these classes at relative abundance levels >5%. Ordination with nDMS was also performed on the amplicon sequencing results and demonstrated that the microbial community structures resulted in separation in 2D space based on the fecal source or mix of fecal sources (Fig. S12), and showed some separation between surface water sampling groups. Sequencing results were compared to coliform results using multidimensional scaling overlaid with coliform data (Fig. S13). Clustering results do not show significant differences between coliform results for wet and dry weather sampling events.



FIGURE 6 a. Cluster analysis at the genus level and b. Relative abundance of bacterial and archaeal classes in the wet and dry weather field samples. (N=32404 sequences per sample) Red bars connect samples without significant differences and black bars connect samples with significant differences (p=0.001).

Then, it is of interest to determine if any of common fecal marker and LeFSE biomarkers of the fecal library were elevated in the field samples. The relative abundance of biomarkers identified using linear discriminant analysis were compared for library and field samples. The relative abundance of surface water biomarkers was higher in surface water field samples than library samples, and ranged from 0.3 to 3.2% (Fig. S4). Four field samples (S-10 Wet and Dry, S-14 Wet, and S-52 Dry), contained the wastewater biomarkers at a relative abundance greater than 20%, and all other field samples were less than 10%. Coriobacteriia were observed in field samples at lower relative abundance than library samples spiked with dog feces (Fig. S8). The highest relative abundance of *Coriobacteriia* was observed in field samples from site S-10 during wet and dry weather and S-14 during wet weather at 0.01 and 0.02%. The relative abundance of the SWH biomarkers in field samples ranged from less than 0.01 to 0.7% relative abundance. Field samples from site S-10 during wet weather and sites S-10 and S-56 during dry weather had the SWH biomarkers present at levels higher than at least one library sample spiked with horse manure (Fig. S7). For all biomarkers, results were not significantly different between wet and dry weather sampling events (average p = 0.49).

2-3-3 Relationship between fecal coliform, qPCR, and amplicon sequencing

Results from coliform, qPCR, and amplicon sequencing were cross compared to determine if amplicon sequencing was an effective source tracking method within this study and to assign sources to the fecal contamination in the Navesink. Biomarkers from amplicon sequencing data were used in the comparison to total coliform and qPCR results, and a heat map was developed to identify the highest and lowest coliform, qPCR and biomarkers (Table 1). Samples S-10 Wet and S-10 Dry had the highest countable

coliform results in the field samples. The wastewater biomarkers (SW) were similarly elevated, however the human marker genes were absent (BacHum) or below the trip blank (HF183) from the wet weather sample from site S-10. Dry weather samples from S-10 were above the trip blank for BacHum and HF183, suggesting wastewater impacts were present, consistent with biomarker and elevated coliform results, and with land use maps which show urban land use in adjacent land areas. Horse Marker genes HoF597 were present in S-10 Wet but absent in S-10 dry, but horse biomarkers in both samples were higher than in other field samples. Similar comparisons in other field samples show that coliform and human fecal marker genes were in relative agreement with SW biomarker relative abundance results for dry weather samples from sites S-14, S-56, and S-58. All other samples had conflicting results due to qPCR results for only one Human biomarker being above the trip blank, and/or contradictory fecal marker gene and biomarker results., Results from Horse HoF597 marker gene and horse biomarkers are in agreement for five samples (S-10 Wet, S-52 Dry, S-56 Dry, S-52 Wet and S-34 dry), but horse biomarkers were very low (<0.01% relative abundance) or absent at five locations where HoF597 was present, suggesting that biomarkers identified at the class level for horse manure and wastewater spikes may not have been specific enough to distinguish those sources in field samples.

								/ /
		/~		5 /	/ /	/ /		ed se
		20	W St	in mi	> Xe	15 010	LIENO	2 death
		AT AT	COPIE-	diest	mana	n ^{co} nt	Jers varo	
	6	offortor	ja /		al Blound	Nest.	nation of	Her olo
	A BO	1001 / IT		O'S CIN2	Ne Ar of	Y BH	Ne Ar gione	ance
	Neva (1)	· / active	14 NOS	Nasio	ALL OF	JOISER	di og our	•
Sample ID	/ 🗸 🗸	/•			/ x ·	/ \ . \	/	
S-10 Wet	3100	0.00		37 70	X	0.07	0.01	
S-10 Met	2800	3.67	3.84	44.85	Λ	0.07	0.01	
S-14 Wet	260	3.59	2 74	21.01	X	0.01	0.02	
S-14 Drv	TNTC	2.60	3.54	6.37	X	0.00	0.00	
S-34 Wet	190	1.94	5.10	2.93	X	0.01	0.01	
S-34 Drv	203	2.89	4.05	3.48	X	0.02	0.01	
S-52 Wet	200	3.74	3.11	3.15	X	0.03	0.00	
S-52 Drv	180	3.55	3.08	21.91		0.01	0.01	
S-56 Wet	950	1.14	3.72	2.83	Х	0.01	0.01	
S-56 Dry	TNTC	3.76	3.40	7.93	Х	0.07	0.01	
S-58 Dry	87	1.91	3.05	1.37	Х	0.00	0.00	
Trip Blank	0	2.83	3.77	NA		NA	NA	
			Library Sa	mples				
S								
(Composite								
of S-34 and	000	4.07	0.00	0.07	X	0.01	0.04	
5-58)	600 TNTO	4.07	3.33	2.37	X	0.01	0.01	
SW	TNIC	5.73	4.84	11.18	X	0.20	0.10	
SWG	TNIC	5.04	5.05	4.42	V	0.04	0.04	
200 200		5.75	5.40	0.20		0.57	1.14	
SWCH		6.30	4.80	4.00	×	0.05	0.30	
SWGD		6.40	4.43 5.27	2.04		0.15	0.50	
SWDH	TNTC	6.43	5.27	1 35	X	0.01	0.00	
SWGHD	TNTC	6.42	5.23	1.64	X	0.12	0.40	
SDGH	TNTC	3.69	2 11	0.82	X	0.04	0.41	
*NA = Not Ana	lvzed	0.00	2.11	0.02		0.00	0.41	l
	.,							
	Low Relativ	e Value					High Relativ	ve Value
	Scale	- / 0.00						
	* colors sho	wn are rel	ative to hia	hest and lo	owest resul	lts within sa	mple set.	

Table 1: Summary of results of biomarker, qPCR, and coliform data for library and field samples

* colors shown are relative to highest and lowest results within sample set. Library and field samples compared separately

2-4 Discussion

A combination of qPCR for fecal maker genes and amplicon sequencing analyzed using a variety of techniques (i.e., cluster analysis, nDMS, biomarker analysis) was applied to a fecal library and field samples for microbial source tracking. Cluster analysis of the total microbial community structure indicated that library samples were significantly different based on fecal source and

therefore has the potential provide information on the source of samples (Fig. 2, Fig. S6.). Significant clusters were formed by several library replicates, which is consistent with other sequencing studies that identified differences between bacterial compositions in different fecal sources (Fisher et. al, 2015, Cao et. al, 2013).

Amplicon sequencing results for field samples were compared to library samples using Brays-Curtis similarity testing, along with nDMS, and showed that field samples were significantly different from library samples. It is possible these differences were due to inherent variability in fecal samples (Fisher 2016). To reduce this source of error, several fecal sources for dog, horse, and geese were used to create the library but a broader sampling may be necessary. Other researchers have also shown temporal and hydrodynamic variation in microbial communities in wastewater, stormwater, and surface water (McCarthy 2017, Henry 2016, Fisher 2016). The lack of clustering observed here is potentially because the ratio of surface water to fecal material was higher in the library than the field samples. A criticism of amplicon sequencing is that even with deep sequencing, the presence of very diverse bacterial communities may make it difficult to identify more minor community groups that are present and potentially better for distinguishing between indicators (Cao, 2013). This underscores the importance of robust statistical approaches to identify key biomarkers within mixed samples.

Biomarkers were evaluated at the class level for mixed library samples using linear discriminant analysis (LDA) then results were plotted for all library samples to indicate which biomarkers performed consistently across the different mixtures. Biomarkers for surface water spiked with wastewater and canine feces were higher in library samples containing canine feces, suggesting potential utility of the biomarker in assessing that fecal source. Biomarkers for SWG, SWH, and SW were not consistently higher in other library samples containing those respective mixtures. Two of the biomarkers identified at the class level were Proteobacteria (*Deltaproteobacteria* for surface water, and *Betaproteobacteria* for surface water and wastewater mixtures). Other

literature has shown that Proteobacteria is prolific in environmental samples (Cloutier, 2017, Kirs, 2017) suggesting that these groups are unlikely unique enough to be effective indicators of specific fecal pollution. Additionally, biomarker analysis was performed at the class level in this study, however other researchers evaluated data at the family (McCarthy, 2017) and genus level (McCarthy, 2017, Fisher, 2015). A drawback of assessing at lower taxonomic levels is that relatively short sequences are generated from amplicon sequencing, raising questions about the ability to accurately assign OTUs (Fisher, 2015,Cao, 2013). In spite of these limitations, analysis has successfully been performed in other studies (McCarthy, 2017, Fisher, 2015) using these taxonomic groups and such analyses may be more appropriate when identifying and assigning biomarkers. Here, cluster analysis performed better at the genus than the class level at discriminating fecal library samples by source.

Biomarkers identified in the library samples were reviewed in the field sample and results were consistent with coliform, qPCR, and expected impacts based on adjacent land use at select sites. (Results from cluster analysis and nDMS were not compared to qPCR because they did not provide a measure of the association of a field sample with the library. Likewise, quantification of fecal markers and quantification of common waterborne pathogens were not compared to the qPCR because they are not source specific.) During dry weather at site S-10 (coliform >200 CFU/100 mL) the human marker signatures HF183 and BacHum were both present above the trip blank and elevated SW marker classes were observed relative to other field samples. During dry weather at site S-56 and wet weather at site S-10 horse marker signatures were observed by both qPCR and elevated SWH marker classes relative to other field samples, in addition to both sites having elevated coliform results >200 CFU/100 mL. These observations were also consistent with land use maps showing agricultural land use adjacent to S-56 and urban land use adjacent to S-10. Sites with lower biomarker results did not necessarily have consistent qPCR results and/or high coliform counts. This may be due to the fact that the samples with lower coliform counts

had lower relative abundances of biomarkers making the source more difficult to distinguish using this method. Biomarker results did not vary significantly between wet and dry weather sampling events, consistent with coliform and qPCR results. Similar to the library, these results are promising but suggest that a deeper taxonomic evaluation of biomarkers may be useful.

2-4-1 Relationship between fecal marker / pathogen sequences and coliform counts

The relative abundance of previously studied fecal bacterial families *Bacteroidaceae*, *Lachnospiraceae*, *Ruminococcaceae*, *Porphyromonadaceae*, *and Prevotellaceae* (Newton et. al 2015, McLellan et al. 2015) was compared in library and field samples. The results indicated that surface water library samples spiked with fecal samples had higher levels of these families than field samples. Field samples from site S-10 which had the highest coliform results also had the highest levels of these bacterial families (Fig. S2), with R² values suggesting a moderate to moderately strong correlation between these indicator groups and fecal contamination (R² between 0.4 and 0.8). This is consistent with recent literature demonstrating these as indicator groups (Newton et. al 2015, McLellan et al. 2015).

Field and library samples were also compared for other bacterial groups of interest were reviewed were fecal indictor bacteria and potential pathogens *Clostridium, Enterococcus, Streptococcus,* and *Vibrio* (USGS, 2014). Indicator organisms *Clostridium, Enterococcus,* and *Streptococcus* were higher in fecal library samples than in field samples. Interestingly *Vibrio* was highest in field sample during dry weather from site S-58, which had the lowest total coliform result of any field sampling event. This may be due to the longevity of DNA compared to viable fecal indicator organisms in water (Cacciabuie et al., 2016). In field samples, the highest *Streptococcus* and *Clostridium* results were observed during dry weather at site S-10 and S-34 which both had fecal coliform present at >200 CFU/100 mL (Fig. S4). These results may support the documented use of these bacterial groups as fecal indicator organisms (USGS, 2014).

2-4-2 Source tracking and implications for the Navesink

Results from wet and dry weather field sampling events from this study were compared to other studies in the Navesink river that took place between 2006 and 2016, to evaluate trends and consistencies within the various datasets. All three studies indicated different prevalence of human fecal sources at select locations in the river. Samples collected from site S-10 during wet weather sampling did not have elevated BacHum or HF193 marker genes, and qPCR results were consistent with a source tracking study conducted by the NJDEP in 2016 which reported results as having a wildlife signature absent of MAR phenotypes which would indicate human impacts, and a 2006 study from the NJDEP also showing a lack of human indicators using F+RNA Colifage (NJDEP, 2016, NJDEP, 2008). Dry weather samples collected at S-10 had multiple lines of evidence including biomarkers and marker genes indicating possible human fecal impact, consistent with a canine study conducted by Clean Ocean Action (COA) in 2016 which identified wastewater impacts at S-10 (COA, 2016). These results are consistent with the predominantly urban land use in at adjacent areas. Wet weather samples collected at S-14 had HF183 results which were elevated above the trip blank and, consistent with the NJDEP report which assigned this site human. HoF597 results at this site indicated horse marker genes were present but MAR profiling indicated wildlife (NJDEP, 2016). Site S-34 was previously characterized by the NJDEP and COA as having potential human, wildlife, and domestic wildlife signatures from wastewater impacts based on MAR and canine tracking (NJDEP, 2016, COA, 2016). Marker genes HF183 from wet and dry weather sampling support these results, along with HoF597 results which were positive. A canine study of site S-52, and MAR studies by the NJDEP of site S-56 indicate human impacts, which is in agreement with BacHum results from site S-52 for wet and dry weather sample, and dry weather samples from site S-56.

It is not surprising that some inconsistencies with previous reports were observed. Generally, multiple lines of evidence for source tracking may be recommended given that different methods

have different targets and accuracies (Stoeckel, 2011, Cao, 2013). Some differences were expected, because Antibiotic Resistance testing (MAR) used by the NJDEP is based on cultivatable cell phenotypes, whereas amplicon sequencing does not distinguish between DNA from live and dead bacteria. Results from this study were also compared with Clean Ocean Action reports which used trained dogs for source tracking. This source tracking method has been reported as non-specific and non sensitive in some studies (Boehm et. al, 2013). The studies from the NJDEP identified wildlife sources in antibiotic resistance sampling, however this paper did limited wildlife biomarker assessments and only attempted to identify and address goose biomarkers, which were not consistently associated with samples spiked with goose fecal material. Finally, temporal variation between the studies may also explain differences in identified sources due to differences in hydrodynamic conditions and the installation and maintenance (or lack thereof) of best management practices in the catchment.

2-5 Conclusions

Amplicon sequencing was evaluated as a source tracking method in comparison to qPCR using a library of fecal spiked samples from multiple sources (horse, dog, goose, and human wastewater) and field samples collected from a river impacted with fecal contamination in New Jersey. Library sequencing results indicated significant differences between many of the fecal mixtures, suggesting the potential utility of using total microbial community structure for source tracking. However, a comparison of the total microbial community in library to surface water samples showed significant differences between fecal spiked samples and field samples collected during different flow conditions, likely due the large amounts of fecal material added to the library samples. Analysis of biomarkers selected using linear discriminant analysis (LDA) demonstrated the most promise with the canine and wastewater fecal biomarkers, with less consistent

results for other biomarkers. Biomarkers identified in the library samples were reviewed in the field sample and results were consistent with coliform, qPCR, and expected impacts based on adjacent land use at select sites. Improved results may be possible with deeper taxonomic evaluation of biomarkers and/or application of other statistical methods for sequence analysis (i.e., Bayesian analysis of total community like that applied in the Source Tracker algorithm)(Knights et. al, 2011). **3.** Chapter (3): Broader Implications

3-1 Broader Implications

The results of this study can inform future application of amplicon sequencing in fecal microbial source tracking. Results indicate that discrimination of fecal mixtures from some sources (e.g., human and horse) was possible while other sources (e.g. goose and canine,) were less likely to result in significant differences in total microbial community structures. Biomarkers from this study at select sites were consistent with source tracking results from previous studies and adjacent land use, supporting the use of biomarkers in the evaluation of fecal sources. The recommendation for next steps is to explore biomarkers at a lower taxanomic level and consider other source tracking methods (such as SourceTracker, which uses a Bayesian approach (Knights et. al, 2011)) or even use multiple statistical methods to evaluate and compare biomarkers. Additionally, future fecal libraries could consider a higher ratio of surface water to fecal material to reduce differences between surface water and fecal library samples. Sampling from additional surface water locations across multiple tidal cycles could make future studies more robust and aid in the identification of biomarkers that are common during varied hydrodynamic conditions. Given the significant differences in fecal library samples with select sources, sampling a broader range of sources for a given fecal type and particularly the addition of wildlife sources could provide new insight. Results of this study continue to indicate that fecal indicator organisms are elevated in the Navesink river and sources include human wastewater and horse manure. Best management practices can be applied to help reduce these sources of fecal contamination towards resolving the impairment of this ecologically, economically, and recreationally important waterway.

References

Donovan, E., Unice, K., Roberts, J. D., Harris, M., and Finley, B. 2008. Risk of gastrointestinal disease associated with exposure to pathogens in the water of the lower Passaic River. Appl. Environ. Microbiol. 74, 994–1003.

Green HC Dick LK Gilpin B Samadpour M & Field KG (2011a) Genetic markers for rapid PCR-based identification of gull, Canada goose, duck, and chicken fecal contamination in water. Appl Environ Microbiol 78: 503–510.

Eramo, A., Delos Reyes, H., Fahrenfeld, N. L., 2017. Partitioning of Antibiotic Resistance Genes and Fecal Indicators Varies Intra and Inter-Storm during Combined Sewer Overflows, Frontiers in Microbiology, 8, 2024

2014 New Jersey Integrated Water Quality Assessment Report, 2017, NJDEP Division of Water Monitoring and Standards

Five Total Maximum Daily Loads for Total Coliform to Address Shellfish-Impaired Waters in Watershed Management Area 12 Atlantic Coastal Water Region NJDEP, 2006

M.S. Gilmore, D.B. Clewell, Y. Ike, N. Shankar (Eds.), 2014. Enterococci as indicators of environmental fecal contamination Enterococci: from commensals to leading causes of drug resistant infection, Massachusetts Eye and Ear Infirmary, Bostonpp. 1-17

Kaspar, C.W., Burgess, J.L., Knight, I.T., & Colwell, R.R. 1990. Antibiotic resistance indexing of Escherichia coli to identify sources of fecal contamination in water. Canadian Journal of Microbiology, 36(12), 891–894.

Hagedorn, C., 2001 Bacterial Source Tracking. Crop and Soil Environmental Sciences, Va Tech. http://www.bsi.vt.edu/biol_4684/BST/BST. html

Monmouth County Health Department 2010. Restoring the Shellfish Impaired Waters of the Navesink Estuary Waters of the Navesink Estuary.ppt

NJDEP Bureau of Marine Water Monitoring, 2017.Upper Navesink River Stormwater Study: Microbial Source Tracking

United States Environmental Protection Agency 2017. National Water Quality Inventory Report to Congress (305(b) report)

Schoen, M.E. and Ashbolt, N.J. (2010) Assessing pathogen risk to swimmers at non sewage

impacted recreational beaches. Environ. Sci. Technol. 44, 2286-2291

Boehm A.B., Ashbolt, N.J., Colford, J.M., Dunbar, L.E., Fleming, L. E., Gold, M.A., Hansel, J.A., Hunter, P.R., Ichida, A.M., McGee, C.D., Soller, J.A., Weisberg, S.B. 2009. A sea change ahead for recreational water quality criteria Water Health, 7, 9-20

Sigler V., PasuttiL. 2006. Evaluation of denaturing gradient gel electrophoresis to differentiate *Escherichia coli* populations in secondary environments Environ. Microbiol., 8 (10) 1703-1711.

Fogarty L.R., Voytek M.A. 2005. Comparison of *Bacteroides-Prevotella* 16S rRNA genetic markers for fecal samples from different animal species Appl. Environ. Microbiol., 71 (10) 5999-6007.

Knights, D.; Kuczynski, J.; Charlson, E. S.; Zaneveld, J.; Mozer, M. C.; Collman, R. G.; Bushman, F. D.; Knight, R.; Kelley, S. T. 2011. Bayesian community-wide cultureindependent microbial source tracking Nature Methods 8 (9), 761-763

Segata, N., Izard, J., Waldron, L., Gevers, D., Miropolsky, L., Garrett, W.S., Huttenhower, C. (2011) Metagenomic biomarker discovery and explanation. Genome Biol . 12: R60.

The Weather Company (2017). LLC Weather Underground Historical Weather. Available at: <u>https://www.wunderground.com/history/</u> [accessed June 5, 2018].

Mobile Geographics. (2018) Mobile Geographics, LLC, Available at <u>http://tides.mobilegeographics.com/locations/5256.html</u> [accessed June 5, 2018]

American Public Health Association, Washington, DC (2005) APHAStandard Methods for the Examination of Water and Wastewater (21st ed.)

EPA 2002. Method 1604: total coliforms and *Escherichia coli* in water by membrane filtration using a simultaneous detection technique (mi medium)

Seurinck, S., Defoirdt, T., Verstraete, W., Siciliano, S. D., 2005. Detection and quantification of the human-specific HF183 Bacteroides 16S rRNA genetic marker with real-time PCR for assessment of human faecal pollution in freshwater Env. Microbiology, 7 (2) 249-259

Kildare, B. J., Leutenegger, C. M., McSwain, B. S., Bambic, D. G., Rajal, V. B., and Wuertz, S. (2007). 16S rRNA-based assays for quantitative detection of universal, human-, cow-, and dog-specific fecal Bacteroidales: a Bayesian approach. Water Res. 41, 3701–3715.

Dick, L.K., Bernhard, A. E, Brodeur, T., J, Santo Domingo, J. W., Simpson, J. M.,
Walters, S.P., Field, K.G., 2005. Host distributions of uncultivated
fecal *Bacteroidales* reveal genetic markers for fecal source identification Appl, and Env.
Microbiology, 71 (6) 3184-3191

Suzuki, M. T., Taylor, L. T., and DeLong, E. F. (2000). Quantitative analysis of smallsubunit rRNA genes in mixed microbial populations via 5'-nuclease assays. *Appl. Environ. Microbiol.* 66, 4605–4614

Kuczynski, J.; Walters, W.A.; Stombaugh, J.; Knight, R.; González, A.; Caporaso, J.G. 2012. Using QIIME to analyze 16s rRNA gene sequences from microbial communities, Current Protocols in Microbiology, (SUPPL.27)

Bolger A.M., Lohse M., Usadel B., 2014. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinf 30:2114–2120

Masella A.P., Bartram A.K., Truszkowski J.M., Brown D.G., Neufeld J.D., 2012. PANDAseq: paired-end assembler for Illumina sequences. BMC Bioinf 13:31

Edgar, R.C., Haas, B.J., Clemente, J. C., Quince, C., Knight, R., 2011. UCHIME improves sensitivity and speed of chimera detection Bioinf, 27 (16) 2194-2200

Segata, N., Izard, J., Waldron, L., Gevers, D., Miropolsky, L., Garrett, W. S., et al. (2011). Metagenomic biomarker discovery and explanation. *Genome Biol*.12:R60.

Newton RJ, McLellan SL, Dila DK, Vineis JH, Morrison HG, Eren AM, Sogin ML. 2015. Sewage reflects the microbiomes of human populations. mBio 6(2):e02574 –14.

Fisher, J., Eren, A., Grenn, H, Shanks, O, Morrison, H., Vineis, J, Sogin, M., McLellan, S.L., 2015., Comparison of Sewage and Animal Fecal Microbiomes by Using Oligotyping Reveals Potential Human Fecal Indicators in Multiple Taxonomic Groups Applied and Env. Microbio., 81 (20) 7023

Myers, D.N., Stoeckel, D.M., Bushon, R.N., Francy, D.S., and Brady, A.M.G., 2014, Fecal indicator bacteria (ver. 2.1): U.S. Geological Survey Techniques of Water-Resources Investigations, book 9, chap. A7, sec. 7.1, May 2014, accessed June 6, 2018, from <u>http://pubs.water.usgs.gov/twri9A7/</u>

Fisher, J.C., Eren, M.A., Green, H.C., Shanks, O.C., Morrison, H.G., Vineis, J.H., Sogin, M.L., McLellan, S.L. 2015. Comparison of sewage and animal fecal microbiomes using oligotyping reveals potential human fecal indicators in multiple taxonomic groups. Appl. Environ. Microbiol., 81, 7023–7033.

Cloutier, Danielle D. 2017 Microbial Communities and the Diverse Ecology of Fecal Indicators at Lake Michigan Beaches (Doctoral Dissertation) Retrieved from ProQuest Dissertations Publishing (Access No. 10255660)

Kirs, M., Kisand, V., Wong, M., Caffaro-Filho, R. A., Moravcik, P., Harwood, V. J., Yoneyama, B., Fujioka, R. S. 2017 Multiple lines of evidence to identify sewage as the cause of water quality impairment in an urbanized tropical watershed, Water Res., 116, 23-33 McCarthy, D.T., Jovanovic D., Lintern A., Teakl, I., Barnes, M., Coleman, R., Rooney, G., Prosser, T., Coutts, S., Hipsey, M.R., Deletic, A., Bruce, L.C., Henry R. 2017. Source tracking using microbial community fingerprints: Method comparison with hydrodynamic modeling, Water Res., 109, 253-265

Fisher, J.C., Newton, R.J., Dila, D.K., McLellan, S.L. 2015. Urban microbial ecology of a freshwater estuary of Lake Michigan Elem Sci Anth; 3:64

Cao, Y., Van De Werfhorst, L.C., Dubinsky, E.A., Badgley, B.D., Sadowsky, M.J., Andersen, G.L., Griffith, J.F., Holden, P.A., 2013 Evaluation of molecular community analysis methods for discerning fecal sources and human waste. Water Res. 47 (18)

Fisher, J. C., Newton, R. J., Dila, D. K., and McLellan, S. L. (2015). Urban microbial ecology of a freshwater estuary of Lake Michigan. Elementa 3:000064

Gutierrez-Cacciabue, D.; Cid, A. G.; Rajal, V. B. 2016. How long can culturable bacteria and total DNA persist in environmental waters? The role of sunlight and solid particles. Sci Total Environ 539, 494-502.

Sanitary Survey of Shellfish Growing area NE 2 Navesink River, 2015, NJDEP Bureau of Marine Water Monitoring

Upper Navesink River Stormwater Study: Microbial Source Tracking, 2008, NJDEP Bureau of Marine Water Monitoring,

Canine Investigations for Source Tracking of Human Sewage Contamination in the Navesink River, Monmouth County NJ, 2016, Environmental Canine Services,

Stoeckel D., Stelzer A., Stogner R., Mau D. 2011. Semi-quantitative evaluation of fecal contamination potential by human and ruminant sources using multiple lines of evidence. Water Research 45:3225–44.

Supplemental Figures

Sample	Sample Location	Description
type		
Horse	Horse Farm in Middlesex	Fresh Manure from 2 horses from two
	County	stables
Horse	Rutgers Horse Farm	Fresh Manure Samples from multiple horses
		in farm
Goose	Park in Red Bank, New Jersey	Multiple fresh samples collected adjacent to
		park pond
Goose	Park in Somerset, New Jersey	Fresh samples collected from park
Dog	4 Domestic dogs from 3	Fresh samples provided by owners
	owners – mixed breeds	

Table S1: Fecal spike sample location and description for samples used in the creation of the fecal library.

Table S2: Fecal Library Composition by volume (surface water and wastewater) and weight (fecal material spikes.

Sample		Surface	XX /	Contraction		Dee
name	Contents	water	wastewater	Goose	Horse	Dog
Sa, Sb	Surface Water	900 mL	0 mL	0 g	0 g	0 g
SWa,	Surface Water,					
SWb	Wastewater	900 mL	100 mL	0 g	0 g	0 g
SWGa,	Surface Water,					
SWGb	Wastewater, Goose	900 mL	100 mL	1 g	0 g	0 g
SWHa,	Surface Water,					
SWHb	Wastewater, Horse	900 mL	100 mL	0 g	1 g	0 g
SWDa,	Surface Water,					
SWDb	Wastewater, Dog	900 mL	100 mL	0 g	0 g	1 g
	Surface Water,					
SWGHa,	Wastewater, Goose,					
SWGHb	Horse	900 mL	100 mL	1 g	1 g	0 g
SWGDa,	Surface Water,					
SWGDb	Wastewater, Goose, Dog	900 mL	100 mL	1 g	0 g	1 g
SWDHa,	Surface Water,					
SWDHb	Wastewater, Dog, Horse	900 mL	100 mL	0 g	1 g	1 g
SWGHD						
a,	Surface Water,					
SWGHD	Wastewater, Goose,					
b	Horse, Dog	900 mL	100 mL	1 g	1 g	1 g
SDGHa,	Surface Water, Dog,					
SDGHb	Goose, Horse	900 mL	0 mL	1 g	1 g	1 g
В	De-Ionized Water	0 mL	0 mL	0 g	0 g	0 g



Figure S1: Tidal Cycles for sample wet (a) and dry (b)weather events

Sample ID	Filtered Volume (mL)
Control	69
S-10-2a	6
S-10-2a	16
S-10-2b	8
S-10-2b	19
S-14-2a	53
S-14-2a	112
S-14-2b	45
S-14-2b	107
S-34-2a	36
S-34-2a	97
S-34-2b	48
S-34-2b	111
S-52-2a	50
S-52-2a	97
S-52-2b	41
S-52-2b	101
S-56-2a	9
S-56-2a	24
S-56-2b	20
S-56-2b	28
S-58-2a	55
S-58-2a	97
S-58-2b	47
S-58-2b	99
TB-2	99

Table S3: Filter volumes for analysis of total coliform by method EPA SM1604

*all samples normalized to counts per 100 mL

Table S4: Primer information

Target Fecal Source	Primer/Probe name	Primer/Probe Sequence 5' to 3'	Та	Reference Paper		
	Bac708R	CAATCGGAGTTCTTCGTG	53°	Dick et		
Horse	HoF597F	oF597F CCAGCCGTAAAATAGTCGG				
Human	HF183F	ATCATGAGTTCACATGTCCG	53°	Seurinck et al		
	Bac242R	TACCCCGCCTACTATCTAATG	С	2005		
	BacHum160f	TGAGTTCACATGTCCGCATGA	60°	Kildare et		
Human	BacHum241r	CGTTACCCCGCCTACTATCTAAT		al, 2007		



Figure S2: a-c Rarefaction curves for library samples at the class level



Figure S2: (continued) d-f Rarefaction curves for field samples

	Average	Standard Deviation				
HF183 wastewater	5.05	0.35				
HF183 no wastewater*	2.72	0.86				
BacHum Wastewater	5.98	0.27				
Bachum no wastewater*	3.45	0.33				
*Library samples that were not spiked with human wastewater were "S"						
and "SDGH"						

Table S5: Average and standard deviation of qPCR results (log copies/mL) for library samples spiked and not spiked with wastewater.

Table S6:	Library a	and field	sample	sequencing	information
-----------	-----------	-----------	--------	------------	-------------

	Sam ple ID	Sample Name	Number of Sequenc es	Number of OTUs at the class level**	Shannon Index*		Sample ID	Sample Name	Number of Sequences	Number of OTUs at the class level**	Shannon Index*
	1a	S	60512	85	9.8		S.10	S.10	38997	131	10.6
	1b	Sdup	48839	93	10.2		S.10d	S.10d	46356	133	10.6
	2a	SW	38070	83	11.3		S.10.2	S.10.2	62272	118	10.5
	2b	SWdup	46261	85	11.4		S.10d.2	S.10d.2	54542	112	10.5
	3a	SWG	51676	60	10.7		S.14	S.14	60325	89	10.3
	3b	SWGdup	50644	74	10.9		S.14d	S.14d	61086	95	10.4
	4a	SWH	63107	75	11.0		S.14.2	S.14.2	49723	81	10.1
	4b	SWHdup	53462	77	11.3		S.14d.2	S.14d.2	52282	59	9.7
SS	5a	SWD	60816	66	10.5		S.34	S.34	53117	78	10.6
mpla	5b	SWDdup	53998	70	11.0	es	S.34d	S.34d	53775	73	10.4
y Sa	6a	SWGH	49075	65	11.4	ampl	S.34.2	S.34.2	50952	116	10.6
brar	6b	SWGHdup	68312	58	10.7	eld S	S.34d.2	S.34d.2	47716	117	10.5
Ľ	7a	SWGD	54991	63	11.0	Εï	S.52	S.52	46295	116	10.7
	7b	SWGDdup	51986	71	10.8		S.52d	S.52d	51613	125	11.0
	8a	SWDH	64088	56	10.0		S.52.2	S.52.2	54426	100	10.6
	8b	SWDHdup	48282	57	9.9		S.52d.2	S.52d.2	40910	99	10.5
	9a	SWGDH	69636	54	9.8		S.56	S.56	39462	82	10.5
	9b	SWGDHdu p	54634	55	10.6		S.56d	S.56d	38070	67	10.5
	10a	SDGH	57054	53	10.3		S.56.2	S.56.2	54252	105	10.8
	10b	SDGHdup	62959	52	10.2		S.56d.2	S.56d.2	58018	118	11.0
	11	Blank	54511	57	7.3		S.58.2	S.58.2	48014	63	9.7
							S.58d.2	S.58d.2	32404	57	9.5



Figure S3: Heatmap showing relative number of sequences for bacterial families previously reported to be associated with fecal material: *Bacteroidaceae, Lachnospiraceae, Ruminococcaceae, Porphyromonadaceae, and Prevotellaceae*



Figure S4: Linear discriminant analysis results from galaxy LEfSE analysis tool for class level data. a. Biomakers for library samples. b. Cladogram visually illustrates demonstrates relationship between biomarker species identified in figure a.



Figure S5: Relative abundance of surface water indicator bacteria in library and field samples, *Acidobacteriia*, *Caldithrixae*, ABY1, *Phycisphaerai*, *Deltaproteobacteria*. Error bars represent high and low values of replicate samples (N = 2)



Figure S6: Relative abundance of wastewater spiked surface water indicator bacteria in library and field samples *Thermomicrobia*, C6, *Betaproteobacteria*, *Synergistia*. Error bars represent high and low values of replicate samples (N = 2)



Figure S7: Relative abundance of horse and wastewater spiked surface water indicator bacteria in library and field samples MVP-15, *Spirochaetes*, and TM7-3 Error bars represent high and low values of replicate samples (N = 2)



Figure S8: Relative abundance of dog and wastewater spiked surface water indicator bacteria in library and field samples (*Coriobacteriia*). Error bars represent high and low values of replicate samples (N = 2)



Figure S9: Precipitation (a., c.), discharge, and gage height (b., d.) data for Red Bank New Jersey during wet weather (a., b.) and dry weather (c., d) sampling events. Precipitation from Weather Underground, discharge and gage height from USGS *Station 01407500 Swimming River near Red Bank NJ*

Wet W	Wet Weather Samples Coliform Results							
Sample ID	CFU/100 mL	RPD (%)						
S-10 Wet	3100	71						
S-14 Wet	260	15						
S-34 Wet	190	11						
S-52 Wet	200	40						
S-56 Wet	950	53						
Dry W	eather Samples Coliform	Results						
Sample ID	CFU/100 mL	RPD (%)						
S-10 Dry	2800	N/A*						
S-14 Dry	TNTC**	N/A						
S-34 Dry	203	N/A*						
S-52 Dry	180	105						
S-56 Dry	TNTC**	N/A						
S-58 Dry	87	15						
TB-2	0	N/A						
Libi	Library Samples Coliform Results							
S***	600	N/A						

Table S7: Wet and dry weather coliform results and relative percent differences (RPD) for replicate samples

* RPD could not be calculated because duplicate sample was TNTC. Reported results are single sample

** TNTC: Too numerous to count

***S: Surface water sample, mixture of sample from S-34 and S-58



Figure S10: Location of the Navesink and Shrewsbury Rivers, with Watershed. Source



Figure S11 Heatmap showing relative number of sequences for bacterial families containing microbes used as fecal indicators: *Streptococcus, Clostridium, Enterococcus, Vibrio*



Figure S12: Nonmetric multidimensional scaling of microbial communities in field and library samples. Results of SIMPROF test showing no significant differences is overlaid.



Figure S13: Nonmetric multidimensional scaling results with bubble plots for coliform data. Samples without bubbles were too numerous to count.

	r for r Clar											
Bir	marker Biomarker	ŝ	SA	SAC	SWH	EWA	SWG	A SMC	D SWDH	- Star	AD SDGH	
S	Acidobacteriia	46	0	0	0	0	0	0	0	0	0	
S	Caldithrixae	46	0	0	0	0	0	0	0	0	0	
S	ABY1	463	401	216	154	93	31	93	46	31	0	
S	Phycisphaerae	185	93	0	77	31	15	62	15	0	0	
S	Deltaproteobacteria	5709	5894	2376	2901	1574	895	1636	1312	849	417	
SW	Thermomicrobia	0	154	46	46	31	0	15	0	0	0	
SW	C6	0	31	0	0	0	0	0	0	0	0	
SW	Betaproteobacteria	23654	110449	43667	82274	40304	18347	21957	13471	16294	8178	
SW	Synergistia	31	1204	447	494	309	93	231	77	77	0	
SWD	Coriobacteriia	93	1003	447	1420	14350	2963	6589	4752	1497	4120	
SWH	M VP-15	0	293	139	926	185	77	0	123	15	31	
SWH	Spirochaetes	62	957	46	3348	185	1173	77	849	293	386	
SWH	TM7-3	15	756	262	1420	154	293	62	278	108	139	
	All data mu	ultiplied b	y a factor	of 10 ⁶ fo	r ease of	review.	Data sh	own is r	elative a	bundan	се	
		Low rel	ative abui	ndance]	High rela	tive abu	ndance	
	AN (1855)											
	vertie ver											
	$\nabla \chi \nabla \chi$	1 th	A	<u>ر م</u>	.A	1 e	ہے۔ ہ		هر الخ	\ .\\$	م ج	مر \
o iom.	resionarts	, 10 Wet	, 10 Dry	AWet	APry	3ª Wet	· 34 Dr	S.S.W	et 52 Dry	5676	a	58 Dr
Biom	Biomart	5.10 Wet	5.10 Dr.3	5.14 Wet	SIADRY	5.34 Web	5-34 Dr	5.524	st Salary	5.56 W	a Stip Dry	5.58 Dr
Biom	Biomart	5-10 ^{Wet} 3100	5-10 Dry 2800	5-14 ^{Wet} 260	S.IADIY TNTC	5.34 Web	5.34 Drs 203	3.52 M	^{عل} ي بتركي المعلم الم معلم المعلم معلم	5:50 N° 5:50	A Stophy Stophy TNTC	\$:58 JF
Bion.	Acidobacteriia	5.19 ^{Wet} 3100 62	5-10 Dr. 9 2800 108	5-14 ^{Wet} 260 0	S-IADIY TNTC 0	5.34 Web 190	5.34 Dr. 203 15	200 0	^{عل} ي: ج:ج: کی کی کی 180	950 950	et ₅ .56 DFC TNTC 46	\$58 JP
Bionit S S S	Acidobacteriia Caldithrixae	3100 62 31 2685	2800 108 0	÷-14 ^{wet} 260 0 15	5-14 DE ³ TNTC 0 31	5-3-4 We	203 15 447	3,52 ⁴⁴ 200 0 278	et 5:52 DF3 180 0 0	950 15 31	5 ,56 DTNTC 46 62	87 0 15
Bionit S S S S	Acidobacteriia Caldithrixae ABY1 Physisphearae	3100 62 31 2685	5.10 DY 2800 108 0 2176	514W ^{et} 260 0 15 324 1975	÷14.045 ÷14.045 TNTC 0 31 463 31	5 ³⁴ ¹⁹⁰ 190 0 62 262	\$-74 m ² \$-74 m ² 203 15 447 432 231	200 278 787 201	262 577 180 0 262	950 15 31 201	⁴ 5 ⁵⁶ p ¹⁰ TNTC 46 62 1342 154	87 0 31
B inner S S S S S	Acidobacteriia Caldithrixae ABY1 Phy cisphaerae	3100 3200 62 311 2685 108 12653	2800 108 2176 62	260 0 15 324 1975	51400 31 463 31	190 0 262 93	203 5 447 432 231 20552	200 278 787 201	2 ² 52 M ² 180 0 262 77 5771	950 950 15 31 201 139	TNTC 46 62 1342 154	87 (15 31 (22895
Bioner S S S S S S S W	Acidobacteriia Caldithrixae ABY1 Phycisphaerae Deltaproteobacteria	 4.100 Wet 3100 62 311 2685 108 12653 15 	2800 108 108 2176 62 9598	260 0 15 324 1975 5817	TNTC 0 31 463 31 2639	190 0 62 262 93 4969	203 203 15 447 432 231 30552	200 278 787 201 24009	262 777 5771	950 950 15 31 201 139 2423	TNTC 1342 14875	87 (15 31 (2885
Bioner S S S S S S W S W	Acidobacteriia Caldithrixae ABY1 Phy cisphaerae Deltaproteobacteria Thermomicrobia	3100 3100 62 31 62 31 12685 12653 15 31	2800 2800 108 0 2176 62 9598 0	260 0 15 324 1975 5817 77	51111111111111	190 190 0 0 0 262 93 4969 0 0	203 203 15 447 432 231 30552 62	200 278 787 201 24009 15	2, 52, 77, 78, 78, 79, 79, 79, 79, 79, 79, 79, 79, 79, 79	950 950 15 31 201 139 2423 0	TNTC 46 62 1342 154 14875 15	87 (15 31 (2885 (
Bioner S S S S S S W S W S W	Acidobacteriia Caldithrixae ABY1 Phy cisphaerae Deltap roteobacteria Thermomicrobia C6 Bet aproteobacteria	 3100 3100 62 311 2685 108 12653 15 31 377855 	2800 108 0 2176 62 9598 0 0 0 448479	260 260 0 15 324 1975 5817 77 0 0	TNTC 0 31 463 31 2639 0 15 63680	190 190 0 262 93 4969 0 0 29348	203 203 15 447 432 231 30552 62 139 34564	200 278 787 201 24009 15 170 31277	218954	950 950 15 31 201 139 2423 0 0 0 28330	 a spin b spin c spin <lic li="" spin<=""> c spin c spin</lic>	87 (15 31 (2885 (31 (31 (31 (31 (31))(31)) (31)) (31))(31))(31))(31))(31))(31))(31))(31))(31))((31))(31))(31))(31))(31))(31)(31
Bioner S S S S S S W S W S W S W	Acidobacteriia Caldithrixae ABY1 Phy cisp haerae Deltaproteobacteria C6 Betaproteobacteria Sv nergistia	 3100 3100 62 311 2685 108 12653 15 311 377855 15 	2800 2800 108 2176 22176 62 9598 0 448479	260 260 0 15 324 1975 5817 77 77 0 209974	TNTC 0 31 463 31 2639 0 0 15 63680 15	190 190 0 262 93 4969 0 0 29348 0	203 203 15 447 432 231 30552 62 139 34564	200 278 787 201 24009 15 170 31277 31	2 3 3 3 3 3 3 3 3 3 3 3 3 3	950 950 15 31 201 139 2423 0 0 28330	tinterimentation of the second	87 (15 31 (2885) (31 13625) (13625)
Bioner S S S S S S S W S W S W S W S W D	Acidobacteriia Caldithrixae ABY1 Phy cisp haerae Deltap roteobacteria Thermomicrobia C6 Betap roteobacteria Synergistia	 3100 3100 622 311 26855 108 126533 155 311 3778555 158 108 	2800 2800 108 0 2176 62 9598 0 448479 0 247	260 260 0 15 324 1975 5817 77 0 209974 15	TNTC 0 31 463 31 2639 0 15 63680 15	190 190 0 262 262 93 4969 0 0 29348 0	203 203 15 447 432 231 30552 62 139 34564 15	200 278 787 201 24009 15 170 31277 31	2 3 3 3 3 3 3 3 3 4 3 3 4 3 4 5 5 7 7 5 7 7 5 7 7 5 7 7 5 7 7 5 7 7 1 80 2 62 7 7 5 7 7 1 80 2 62 7 7 5 7 7 1 80 2 62 7 7 5 7 7 1 8 9 3 0 0 2 62 7 7 5 7 7 1 8 9 3 0 0 2 62 7 7 5 7 7 1 9 3 0 0 2 62 7 7 5 7 7 1 9 3 0 0 0 2 6 7 7 1 5 7 7 1 9 3 0 0 2 1 8 9 5 4 1 5 7 1 9 3 0 0 2 1 8 9 5 4 1 5 7 1 1 5 7 1 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 5 7 1 7 1 7 1 7 7 1 7 7 7 1 7 7 7 7 7 7 7 7 7 7 7 7 7	950 950 15 31 201 139 2423 00 28330 00 28330	TNTC 46 62 1342 154 14875 15 79234 0 62	87 (15 31 (2885 (31 13625 ()

All data multiplied by a factor of 10^6 for ease of review. Data shown is relative

Low relative abundance

SWH Spirochaetes

SWH TM7-3

High relative abundance

Figure S14: Heat maps for biomarkers identified using Linear Discriminant Analysis