

PAPER**CRIMINALISTICS**

Ellen M. Jesmok,^{1,†} M.S.; James M. Hopkins,^{1,‡} M.S.; and David R. Foran,² Ph.D.

Next-Generation Sequencing of the Bacterial 16S rRNA Gene for Forensic Soil Comparison: A Feasibility Study*

ABSTRACT: Soil has the potential to be valuable forensic evidence linking a person or item to a crime scene; however, there is no established soil individualization technique. In this study, the utility of soil bacterial profiling via next-generation sequencing of the 16S rRNA gene was examined for associating soils with their place of origin. Soil samples were collected from ten diverse and nine similar habitats over time, and within three habitats at various horizontal and vertical distances. Bacterial profiles were analyzed using four methods: abundance charts and nonmetric multidimensional scaling provided simplification and visualization of the massive datasets, potentially aiding in expert testimony, while analysis of similarities and *k*-nearest neighbor offered objective statistical comparisons. The vast majority of soil bacterial profiles (95.4%) were classified to their location of origin, highlighting the potential of bacterial profiling via next-generation sequencing for the forensic analysis of soil samples.

KEYWORDS: forensic science, soil evidence, bacterial profiling, soil profiling, bacterial abundance charts, nonmetric multidimensional scaling, *k*-nearest neighbor, analysis of similarities

Soil can become evidence in a criminal investigation when a crime occurs in an outdoor location or is collected from clothing, shoes, tires, or other items associated with a crime. In these instances, soil can reveal information about the location where a crime was committed or link a suspect or victim to the scene (1). Soil was involved in criminal investigations as far back as the 1800s, when a visual comparison of sand was used to trace a barrel that had once been filled with silver to a specific train station on the Prussian railroad, resulting in a conviction (2). Traditional forensic techniques for soil analysis involve the examination of class characteristics such as grain size, pH, and moisture content (3,4). Conclusions drawn based on these characteristics may vary from one expert to the next, and unless the

soil contains very rare attributes, it often does not carry much probative value. Given such shortcomings, the need for an overarching, objective technique for soil identification and comparison is evident. Bacterial profiling, already widely used by microbiologists to characterize varied soil samples, holds the potential for the individualization of soil for forensic purposes.

Bacterial Profiling of Soil

One gram of soil contains 4×10^7 – 2×10^9 bacteria, which vary widely in species diversity and abundance (5). For over a century, attempts have been made to detect and quantify bacteria from an array of sample types, dating as far back as Koch's (6) isolation of bacteria from blood via culture on Petri plates. Culturing can be used to quantify bacteria, but only viable cells under suitable growth conditions develop colonies, resulting in a steep underestimation of diversity, limiting the applicability of bacterial culture methodologies for forensic purposes (7). Techniques developed for bacterial identification beyond culturing also generally fail to capture the diversity and abundance of bacteria present in soils, as they cannot differentiate many bacteria either because they do not detect certain species or they lack resolving power so that different species categorize as the same.

Woese and Fox (8) pioneered the use of 16S ribosomal RNA (rRNA) gene analysis for constructing bacterial phylogenies. This marker is conserved across bacteria and archaea, but contains nine variable regions that can be used for identification, often to the species level. Two decades later, Liu et al. (9) described terminal restriction fragment length polymorphism (T-RFLP) analysis, which generates 16S bacterial profiles via locus amplification, restriction enzyme digestion, and electrophoretic separation of the products. The resultant bands or peaks are compared to estimate bacterial similarities among samples. Since

¹Forensic Science Program, School of Criminal Justice, Michigan State University, 655 Auditorium Road, Room 560A, East Lansing, MI 48824.

²Forensic Science Program, School of Criminal Justice and Department of Integrative Biology, Michigan State University, 655 Auditorium Road, Room 560A, East Lansing, MI 48824.

[†]Present address: Minnesota Bureau of Criminal Apprehension Laboratory, 1430 Maryland Ave E St. Paul, MN 55106.

[‡]Present address: St. Jude Children's Research Hospital, 262 Danny Thomas Place, Mail Stop 1160, Memphis, TN 38105.

*Presented at the 67th Annual Meeting of the American Academy of Forensic Sciences, February 16–20, 2015, in Orlando, FL.

A portion of this research was supported by grant number 2013-R2-CX-K010, awarded by the National Institute of Justice, Office of Justice Programs, U.S. Department of Justice.

Points of view are those of the authors and do not necessarily represent the official position or policies of the U.S. Department of Justice.

This article was published online on [18 feb 2016]. An error was subsequently identified. This notice is included in the online and print versions to indicate that both have been corrected [29 feb 2016].

Received 6 Mar. 2015; and in revised form 19 July 2015; accepted 26 July 2015.

then, T-RFLP has been widely used in both the microbiological (e.g., 10) and forensic (e.g., 11–13) fields. Unfortunately, as noted above for other techniques, the massive number of bacterial species in soil and the limited resolving power of T-RFLP diminishes its forensic utility.

Given this, it is clear that successful forensic analysis of soils based on their bacterial makeup requires the production of more complete, higher resolution data. Next-generation sequencing of the 16S rRNA gene, first described by Jonasson et al. (14), may fill this role, as it is an extremely robust methodology for bacterial identification. Next-generation sequencing is currently being used extensively by microbiologists for bacterial community analysis (e.g., 15,16), and large bacterial sequence reference databases have been created (17,18), allowing classification of bacteria at taxonomic levels from phylum to species. Several next-generation sequencing platforms are available, including ion semiconductor sequencing (Life Technologies, Carlsbad, CA), pyrosequencing (e.g., Roche, South San Francisco, CA), and Illumina sequencing by synthesis (San Diego, CA). Loman et al. (19) compared the performance of these three platforms in sequencing an *E.coli* isolate and found the latter produced the lowest error rate and highest throughput. The Illumina methodology has been successful for generating bacterial profiles collected from a variety of environments such as the human gut and lake sediment (20), but its value for forensic characterization of soil has not been documented.

Analysis of Soil Bacterial Profiles Generated via Next-Generation Sequencing

Next-generation sequencing can produce hundreds of thousands of sequence reads per soil sample, which results in tremendous amounts of potential identification data, but at the same time creates unique challenges for forensic applications. Microbiologists have used myriad methods for analyzing these massive datasets, including bacterial abundance charts (e.g., 21,22), pairwise comparisons (e.g., 23,24), hierarchical cluster analysis (e.g., 25,26), nonmetric multidimensional scaling [NMDS (e.g., 10,13)], analysis of similarities [ANOSIM (27)], and supervised classifications (e.g., 28). All of these techniques have been useful for understanding bacterial community structure and function; however, the criminal justice system has somewhat different data analysis demands, as its main goal is to associate or differentiate soil samples. This can potentially be performed visually using next-generation sequencing data, although the 2009 National Academy of Sciences (NAS) report on forensic science (29) called for “the development and establishment of quantifiable measures of the reliability and accuracy of forensic analyses,” meaning that simply stating the bacterial profile of a questioned and known soil appear similar or dissimilar is not an adequate forensic conclusion. On the other hand, forensic scientists may be required to present their data to a judge or jury in an easily understood fashion, which is often best accomplished through demonstrative charts or diagrams, which again are visual. All factors taken together, it is not clear if the analysis techniques currently used by microbiologists will meet the demands of forensic science, but they represent viable starting points, as some produce a visual output (albeit one that is subjectively interpreted), while others result in specific, objective measures of similarity.

Hopkins (30) examined the utility of several soil bacterial profile analysis techniques for forensic application, eliminating some while finding others quite useful. Using 16S rRNA sequence data produced through pyrosequencing, it was deter-

mined that hierarchical cluster analysis and pairwise comparisons were problematic in that slight procedural modifications often resulted in different outcomes for the same set of soil samples; such results would be decidedly detrimental in a forensic context, where definitive answers are sought. In contrast, the author found that results based on three methods—bacterial abundance charts, NMDS, and supervised classification—remained consistent throughout, fulfilling a primary requirement of forensics.

Bacterial abundance charts are generated from taxonomic data, categorizing and quantifying the bacteria that make up a profile. The charts can be depicted at any taxonomic level, however, if too many groups exist, such as when considering genera or species, the charts may be largely uninterpretable; therefore, most microbiologists build abundance charts at the phylum or class level. Such charts present an easy to understand display of bacterial profile members, and researchers have used them to assess, for instance, the influence of environmental stressors such as repeated wetting and drying (31) or diesel fuel contamination (26) on the bacterial makeup of soil. Abundance charts have also been proposed for the forensic assessment of changing bacterial levels on and within decomposing bodies (32). In court, abundance charts could potentially provide the expert witness with a useful visualization tool for a jury; however, they do not allow for a numerical or statistical measure of relative similarities among bacterial profiles; therefore, further analysis methods are necessary.

NMDS provides a measure of relative association among samples by orienting the datasets in multidimensional space based on their similarity or dissimilarity. For soil analysis, such values would be calculated between each pair of bacterial profiles in a given sample set, either based on DNA sequences or at a given taxonomic level. The location of a soil bacterial profile within a plot, represented by a single point, reveals its relative similarity to the other profiles being ordinated. Groups, or clusters, of similar profiles are identified, and further statistical analysis, such as ANOSIM (27), can be employed to compare them. Increasing the number of dimensions in a NMDS plot can tease out subtle differences among soil profiles; however, visualization becomes more difficult at these higher levels, while two-dimensional plots are easier to interpret. A goodness of ordination measure, termed stress, accompanies a NMDS plot (33), wherein the lowest stress value obtained indicates the best fit of the data. Multidimensional scaling has been employed forensically when comparing soil bacterial profiles generated via T-RFLP analysis (10,13), showing that profiles from the same habitat cluster together in multidimensional space; however, intermingling of profiles among habitat types occurred, precluding clear differentiation of them based on T-RFLP data. A drawback of NMDS is that it has subjective components (e.g., determining the number of dimensions to use or defining a cluster), but in spite of this, the data depiction that NMDS plots provide may have substantial value for jury comprehension of soil bacterial profiling.

Supervised classification techniques have the potential to produce a more objective assignment of bacterial profiles to a location of origin (34). These techniques build models from groups of known samples collectively called training sets. Unknowns are then introduced and assigned to the closest group or, depending on the technique, to no group at all. Yang et al. (28) used supervised classification techniques to assign soil microbial communities to their location of origin with approximately 90% accuracy, based on length differences in 16S rRNA variable regions 1, 2, 3, and 9. This methodology does not hold nearly the resolving power of next-generation sequencing; thus, the

accuracy might be increased; however, it does highlight the utility of supervised classifiers for bacterial profile analysis.

The research presented here was designed to examine the feasibility of using next-generation sequencing of the 16S rRNA gene to individualize soil samples for forensic purposes. Sensabaugh (35) noted that for a soil microbial profiling method to achieve forensic utility, two locations must be differentiable, the technique must be robust and repeatable, and objective statistical measures must be implemented to assess similarities and differences among samples. Minding these recommendations, in this study variation among bacterial profiles from differing habitat types, similar habitat types, and time and space within a habitat was assessed based on next-generation sequencing of the 16S rRNA gene. Bacterial abundance charts and NMDS allowed visual comparison of bacterial profiles, ANOSIM was used to statistically compare clusters of profiles from the same location, and bacterial profiles were assigned to a location of origin based on the supervised classification technique *k*-nearest neighbor (*k*-NN). The forensic utility of these analysis techniques was examined, focusing on their ability to accurately trace soil to its point of origin, their objectivity, and their capacity for facilitating jury understanding.

Materials and Methods

Soil samples were collected in 2013 and 2014 using a garden trowel that had been rinsed with deionized water and wiped with a paper towel. Approximately 100 g of surface soil from a 1-ft² area was homogenized in an 18 oz Whirl-Pak[®] bag (Nasco Fort, Atkinson, WI) and stored at -20°C until DNA extraction. Habitat types and GPS coordinates obtained from Google Maps (Mountain View, CA) are shown in Table 1.

The influence of time on bacterial profiles both within and among habitat types was examined by collecting soil samples from 10 diverse habitats every 3 months for 1 year. Soils were also collected from nine locations of the same habitat type (woodlots 1–9, Table 1) once every 2 weeks over an 8-week period in the summer of 2014, totaling five samples per location.

TABLE 1—GPS coordinates of sampling sites.

Site Name	GPS Coordinates
Marsh*	42°42'32.0"N 84°30'53.4"W
Field*	42°42'38.9"N 84°31'15.4"W
Coniferous Forest*	42°41'11.9"N 84°38'05.1"W
Beach*	42°45'13.9"N 84°24'16.5"W
Corn Agricultural Field*	42°42'33.5"N 84°28'17.5"W
Fallow Agricultural Field*	42°45'06.4"N 84°39'42.8"W
Road Side*	42°48'03.4"N 84°11'10.1"W
Dirt Road*	42°48'17.2"N 84°09'33.5"W
Yard* (Depth and Spatial)	42°42'39.0"N 84°30'53.5"W
Deciduous Woodlot* (Depth and Spatial)	42°42'33.7"N 84°31'01.3"W
Chemically Treated Yard (Spatial)	42°43'26.6"N 84°28'02.5"W
Chemically Treated Yard (Depth)	42°43'44.0"N 84°28'23.4"W
Woodlot 1	42°42'50.8"N 84°28'38.5"W
Woodlot 2	42°41'25.6"N 84°27'41.2"W
Woodlot 3	42°41'03.3"N 84°31'26.1"W
Woodlot 4	42°40'57.2"N 84°28'05.6"W
Woodlot 5	42°42'00.8"N 84°31'35.0"W
Woodlot 6	42°42'33.7"N 84°31'00.6"W
Woodlot 7	42°43'38.9"N 84°30'08.8"W
Woodlot 8	42°44'38.9"N 84°28'57.9"W
Woodlot 9	42°44'28.2"N 84°27'09.8"W

*Denotes the 10 diverse habitats

The influence of space on bacterial profiles was examined by collecting soil samples from three different habitats: a deciduous woodlot, a chemically treated (herbicides and fertilizer) yard, and an untreated yard, in October 2013 and March 2014. Surface soil was collected at a center sampling site and 5, 10, 50, and 100 ft (1.5, 3.0, 15.2, and 30.5 m) distant in each of the cardinal directions. Depth (vertical) samples were collected from 1, 2, 5, 10, 20, and 60 in (2.5, 5.1, 12.7, 25.4, 50.8, and 152.4 cm) using a soil corer and mud auger (AMS, Inc. American Falls, ID) and from the surface with a trowel in a deciduous woodlot, chemically treated yard, and untreated yard in October 2013 (excluding the treated yard vertical space samples for technical reasons) and April 2014. The deepest collection obtainable in the treated yard was 25 in (63.5 cm) due to obstructions at greater depths.

DNA Extraction, Amplification, and Quantification

DNAs were extracted using a PowerSoil[®] DNA Isolation Kit (MoBio, Carlsbad, CA) following the manufacturer's instructions with one additional ethanol wash. Bacterial 16S rRNA gene variable regions 3 and 4 were amplified with a forward primer [357F (36)] and one of 96 bar-coded reverse primers (806R) from the Caporaso et al. (37) primer set, producing a bacterial 16S rRNA gene product of approximately 450 bp. Fifteen microliter reactions contained 1 µL of template DNA, 1 µL of 10 µM forward and reverse primers, 1 U AmpliTaq Gold[®] DNA Polymerase (Applied Biosystems, Foster City, CA), 1.5 µL of the accompanying 10X PCR buffer II, 1.5 µL of 25 mM MgCl₂, 1.5 µL of 2.0 mM nucleotide triphosphates, and 1.5 µL of 4.0 µg/µL bovine serum albumin. Reactions included a 10-min 94°C initial denaturation/enzyme activation step, followed by 35 cycles of 30-s 94°C denaturation, 45-s 60°C primer annealing, and 1-min 72°C extension, followed by a final 10-min 72°C extension.

Four microliters of the PCR product were electrophoresed on a 1% agarose gel alongside 2 µL of a 2-log ladder (New England Biolabs, Ipswich, MA), followed by ethidium bromide staining and UV visualization. A Quant-iT[™] PicoGreen[®] dsDNA Assay Kit (Life Technologies) was used for quantification and PCR products were pooled so that each soil DNA sample totaled 6 ng/µL. Pooled products were purified using Agencourt AMPure XP[™] beads (Beckman Coulter, Brea, CA) at 60% of the total volume.

Next-generation Sequencing of Soil Bacterial DNA

The quality of pooled, purified PCR product was assessed on a 2100 BioAnalyzer (Agilent Technologies, Santa Clara, CA). Samples were sequenced if the DNA integrity number was ≥7. Pooled sequences were loaded on an Illumina MiSeq v2 flow cell along with a 10% PhiX Control (Illumina) and sequenced using a 2x250 bp v2 Reagent Kit (Illumina). Base calling was completed with Real Time Analysis software v1.18.54 (Illumina), and the output was demultiplexed and converted to FastQ files with Bcl2fastq Conversion Software v1.8.4 (Illumina).

Sequence Analysis

The open source software mothur (38) was used to make a single sequence file from the forward and reverse reads of all sequence libraries with a MiSeq quality score >25. Sequences were filtered to remove all ambiguous bases, trimmed to a maxi-

num length of 475 bp, and aligned to the SILVA v119 reference bacterial sequences (18). Sequences were subsampled to 3000 and any bacterial libraries containing fewer than 3000 sequences were excluded from further analysis. Sequences were binned into operational taxonomic units (OTUs) at a 97% similarity cutoff. Dissimilarity values for NMDS and *k*-NN were calculated from OTUs using both Bray–Curtis dissimilarity (39) and the Sørensen–Dice coefficient (40,41). Square symmetric matrices were entered into the XLSTAT Pro (Addinsoft, New York, NY) expansion for Microsoft Excel (Microsoft, Redmond, WA). Two-dimensional NMDS was run with 500 iterations, each stopping at a convergence of 0.00001, using the scaling by majorizing a complicated function algorithm. Random starting configurations were used, and Kruskal’s stress (42) was measured for each plot. Habitat clusters were statistically compared using ANOSIM ($\alpha = 0.05$) in PAST v.2013 (43). *k*-NN was run in Pirouette 4.0 (Infometrix Inc, Bothell, WA) using the square symmetric matrices as the input. Training and test sets used for *k*-NN are outlined in Table 2. OTUs were also grouped at the taxonomic class level based on SILVA reference sequences (18), and abundance charts of each soil bacterial profile were produced in Excel.

Results

Sequencing Efficiency, Taxonomic Class Diversity, and Dissimilarity Matrices

An average of 134,000 sequence reads of the target DNA was produced per soil sample processed. The dirt road sample collected in February 2014 was the only one that did not produce the minimum of 3000 sequences necessary to continue alignment and thus was excluded from further analysis. The remaining dirt road samples contained the fewest bacterial taxonomic classes with an average of 26.8, while the yard samples contained the most with an average of 56.2. The roadside samples also contained relatively few, having an average of 40.6 bacterial classes. The remaining locations averaged 47.0–52.8 classes. Unclassified sequences made up 2–4% of bacterial abundance charts.

Overall, Sørensen–Dice dissimilarity matrices performed better than Bray–Curtis when used as input for NMDS, resulting in tighter clusters of bacterial profiles from the same location and

TABLE 2—Training and test sets for *k*-NN.

	Training Set	Test Set
10 Diverse Habitats*	<i>N</i> = 4 per habitat	<i>N</i> = 1
9 Similar Woodlots*	<i>N</i> = 4 per woodlot	<i>N</i> = 1
Deciduous Woodlot, Yard, and Treated Yard Distance	Center, 5 ft N, 5 ft S, 5 ft W, 5 ft E Center point plus 5 ft E, 10 ft N, 50 ft W, and 100 ft S [†]	All other distance samples
Deciduous Woodlot, Yard, and Treated Yard Depth	Center 100 ft N, 100 ft S, 100 ft W, 100 ft E Surface, 2 in, 10 in, and 60 in	All other distance samples

*Analyzed via the Jackknife resampling method (55) in which each of the five samples was systematically left out and tested against the other four samples

[†]Three additional spiral training sets were developed at the other starting points (i.e., 5 ft W, N, and S), always in a counterclockwise circle.

clearer separation of profiles from different locations. Sørensen–Dice matrices also produced higher overall classification accuracy when used for *k*-NN analysis (95.4%) compared to Bray–Curtis (91.0%); thus, only results using Sørensen–Dice dissimilarities are presented. A summary of the *k*-NN results for all experiments conducted is displayed in Table 3, with assignment accuracies ranging from 87.5 to 100%.

Diverse Habitat Differentiation

Diverse habitat soils (Fig. 1) shared the same major bacterial classes up to approximately 75% based on their abundance charts, with the exception of the dirt road, which had substantially higher levels of *Gammaproteobacteria*, *Flavobacteria*, *Clostridia*, and *Bacilli*, and lower levels of *Acidobacteria* and *Betaproteobacteria* relative to the other habitats. Bacterial profiles generated from soil collected within a habitat clustered together in NMDS plots (Fig. 2), but some intermingling occurred among the ten habitats. When only three habitats were oriented at a time, the clusters separated in all cases (e.g., Fig. 3) and stress was reduced (e.g., from 0.165 in Fig. 2 to 0.133 in Fig. 3). Complete ANOSIM pairwise results can be found in Table S1. All habitat clusters were significantly different except the yard cluster compared to the deciduous woodlot and field, and the fallow agricultural field cluster compared to the marsh. *k*-NN exhibited an 88% assignment accuracy (Table 3) when all habitats were analyzed together. Misassignments occurred between the marsh and fallow agricultural field bacterial profiles and between the deciduous woodlot and yard profiles, although all were correctly assigned when analyzed as pairs in a *k*-NN model.

Similar Habitat Differentiation

The nine woodlot soils shared bacterial classes up to approximately 80% abundance (Fig. 4). Bacterial profiles generated

TABLE 3—*k*-NN summary findings using Sørensen–Dice matrices.

Study	Training Set	Accuracy, %	Misclassified Samples
Diverse Habitats	All Habitat Samples*	88	Marsh with Fallow Ag [†] Field and Deciduous Woodlot with Yard
	Marsh and Fallow Ag Field*	100	–
	Deciduous Woodlot and Yard*	100	–
Similar Habitats	All Location Samples*	87.5	All samples from Woodlot 5 and one sample from Woodlot 2
	Woodlots 1, 3, 4, 6, 7, 8, and 9	100	–
Horizontal Space	Center, 5 ft N, 5 ft S, 5 ft W, 5 ft E	94.4	Yard 100 ft N, Woodlot 100 ft S
	Center point plus 5 ft E, 10 ft N, 50 ft W, and 100 ft S [‡]	97.2	Woodlot 100 ft N
Vertical Space	Center 100 ft N, 100 ft S, 100 ft W, 100 ft E	94.4	Yard 5 ft E, Yard 10 ft N
	Surface, 2 in, 10 in, and 60 in	100	–

*Analyzed via Jackknife Method (55).

[†]Ag = Agricultural

[‡]Other spiral training sets produced similar results.

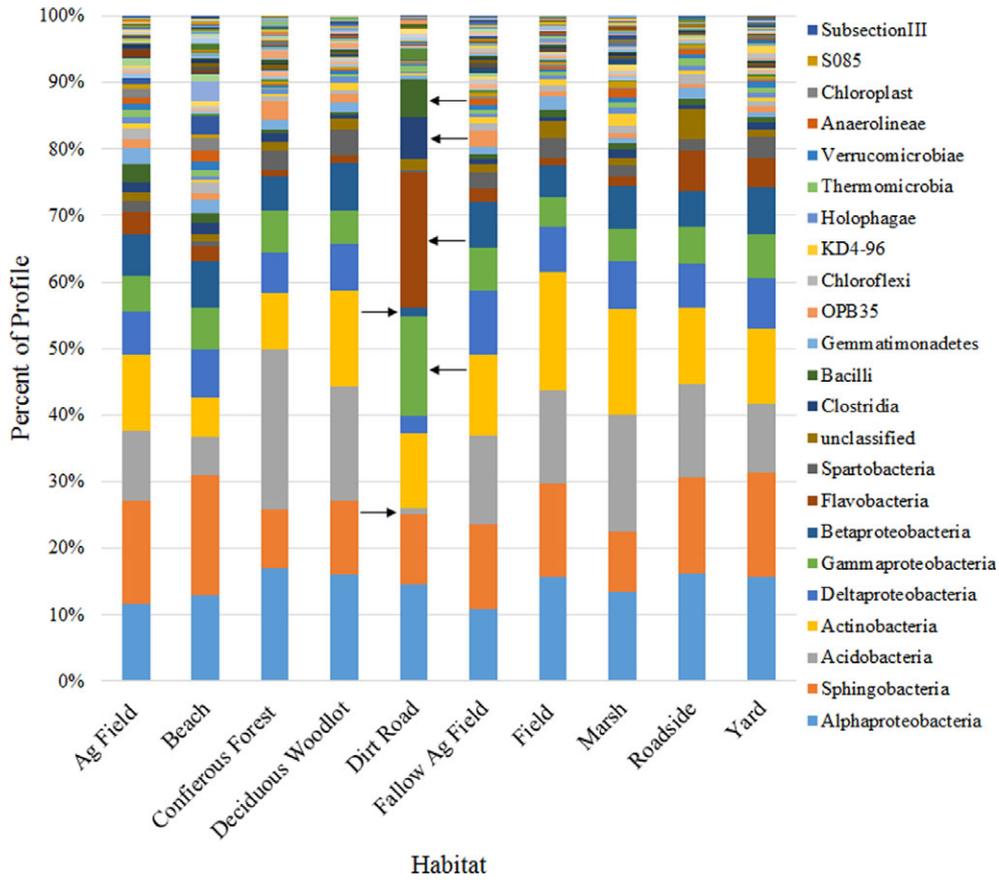


FIG. 1—Average (n = 5) bacterial class abundance of ten diverse habitats. The dirt road clearly differed from the other habitats, containing higher levels of Gammaproteobacteria, Flavobacteria, Clostridia, and Bacilli (denoted by arrows in ascending order on the right) along with lower levels of Acidobacteria and Betaproteobacteria (denoted by arrows in ascending order on the left). The 23 most abundant bacterial classes are listed on the right. Ag = Agricultural. (An interactive, color version of this chart is available as Fig. S1.)

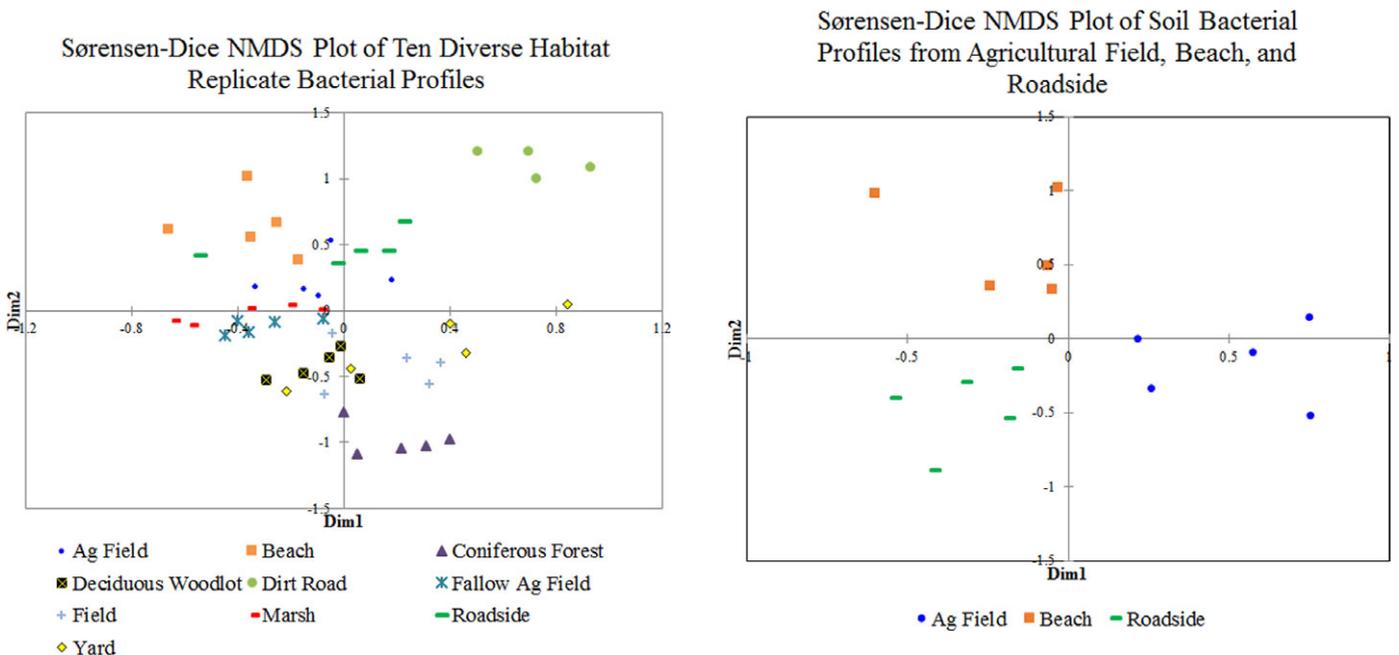


FIG. 2—NMDS plot of the ten diverse habitat bacterial profiles generated via a Sørensen–Dice dissimilarity matrix. Profiles from the same habitat formed clusters, but intermingling occurred among some of the habitats. Kruskal’s stress was 0.165. Ag = Agricultural.

FIG. 3—NMDS plot of the agricultural field, beach, and roadside bacterial profiles generated via a Sørensen–Dice dissimilarity matrix. Profiles from these locations were intermingled when all habitats were ordinated together, but were resolved when analyzed as pairs or triads in NMDS plots. Kruskal’s stress was 0.133. Ag = Agricultural.

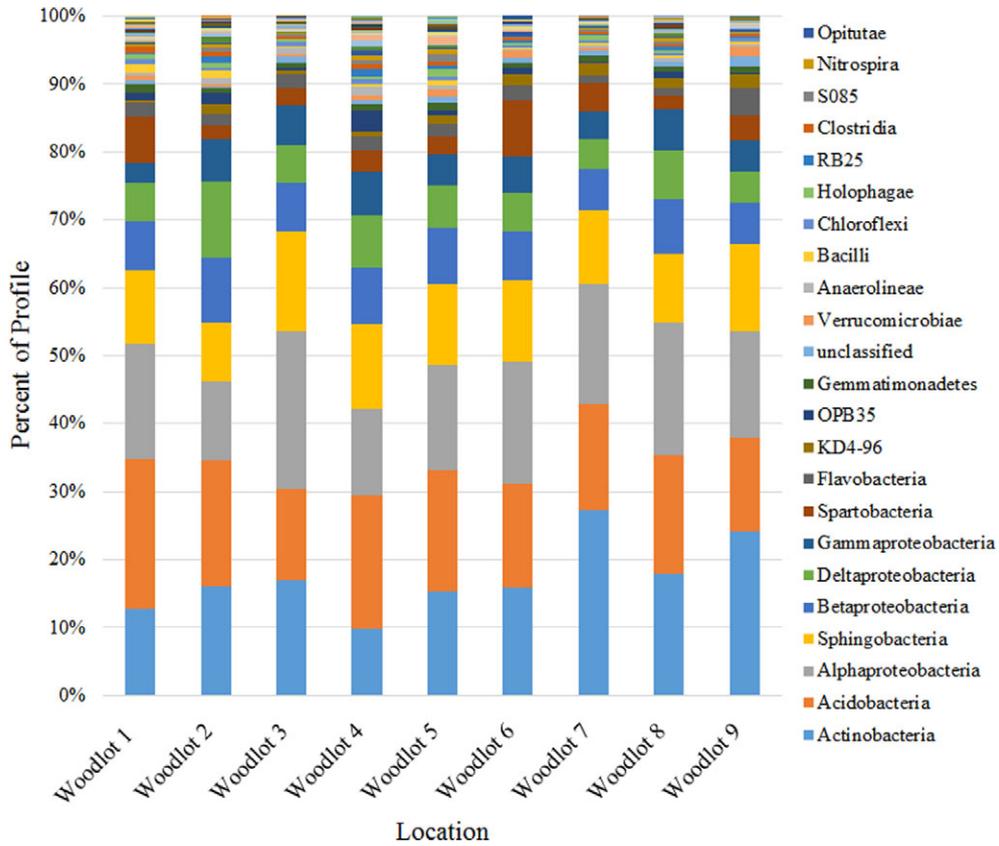


FIG. 4—Average (n = 5) bacterial class abundance of woodlot locations. The soils shared major bacterial classes up to approximately 80%. The 23 most abundant bacterial classes are listed on the right. (An interactive, color version of this chart is available as Fig. S2.)

Sørensen-Dice NMDS Plot of Soil Bacterial Profiles from Nine Woodlots

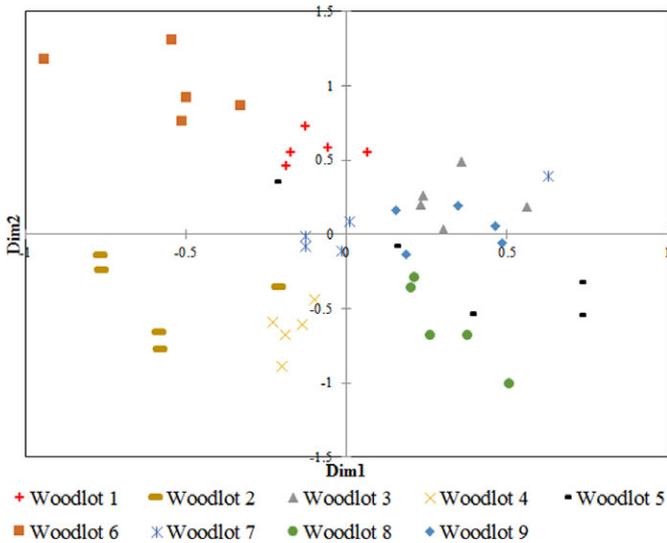


FIG. 5—NMDS plot of the nine woodlot location bacterial profiles generated from a Sørensen–Dice dissimilarity matrix. Profiles from the same location formed clusters, but intermingling occurred among some of the location clusters. Woodlot five profiles clustered relatively poorly, intermingling with several other woodlots. Kruskal’s stress was 0.165.

Sørensen-Dice NMDS Plot of Soil Bacterial Profiles from Woodlots 3, 7, and 9

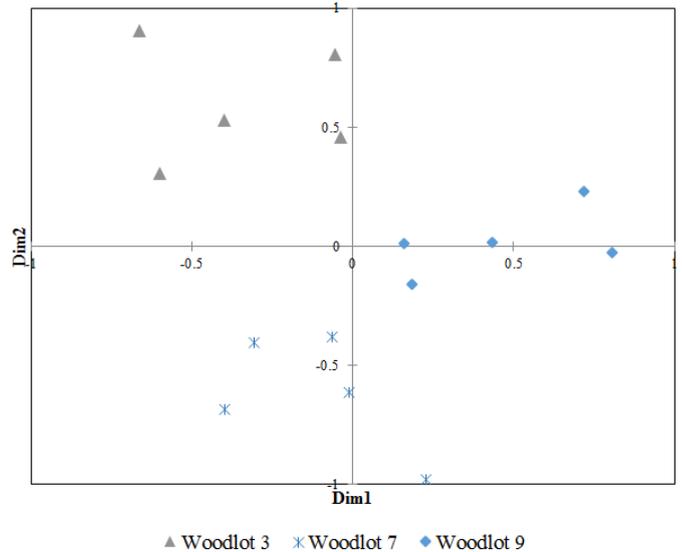


FIG. 6—NMDS plot of woodlots 3, 7, and 9 bacterial profiles generated from a Sørensen–Dice dissimilarity matrix. These clusters were intermingled when all woodlots were ordinated together, but were resolved when they were analyzed alone in a NMDS plot. Kruskal’s stress was 0.107.

from soil collected in the same woodlot clustered together in NMDS plots, but intermingling occurred among several of the clusters (Fig. 5). The most substantial overlap involved woodlot

5, whose profiles were interspersed with several other clusters. By ordinating the woodlot profiles in pairs or triads, separation of woodlots occurred (e.g., Fig. 6) and stress was reduced (e.g.,

from 0.165 in Fig. 5 to 0.107 in Fig. 6). All woodlot clusters were significantly different except woodlot 5, which did not differ significantly from woodlots 3, 7, 8, and 9. *k*-NN was accurate in its assignment of the woodlot bacterial profiles 87.5% of the time (Table 3). All profiles from woodlot 5 and one profile from woodlot 2 were incorrectly assigned, and when removed from the model, 100% assignment accuracy was achieved.

The Influence of Horizontal and Vertical Space on Bacterial Profiles

Surface soils varying distances apart shared the same major bacterial classes up to approximately 80% within and between habitats (data not shown). Bacterial profiles generated from soil collected from the same habitat across horizontal space loosely clustered together in NMDS plots (data not shown), with profiles 50 and 100 ft from the center sampling site being the farthest from the center of the clusters. The treated yard cluster was completely separated from the woodlot and yard, but the latter two intermingled slightly, although all three clusters differed significantly. *k*-NN accurately assigned bacterial profiles 94.4–97.2% of the time (Table 3) depending on the profiles used for the training set (see Materials and Methods). The most accurate classification occurred when using the center plus one profile each from 5, 10, 50, and 100 ft distances in a counterclockwise spiral pattern. Misclassifications were always at least 90 ft distant from the nearest training sample regardless of the starting position of the 5-ft training sample.

Abundance charts generated from the depth samples revealed taxonomic class differences as depth increased (Fig. 7), although

the number of classes remained similar. The most substantial class abundance differences, which existed in all habitats, were higher amounts of *Clostridia*, *Nitrospira*, and *SHA-26* as depth increased. Bacterial profiles generated from the treated yard clustered separately in NMDS plots, while the deciduous woodlot and yard profiles intermingled (data not shown). A trend existed when habitats were ordinated individually, where bacterial profiles were progressively farther away from the surface sample in multivariate space as depth increased (data not shown). April and October clusters within each habitat were not significantly different. Clusters among the three habitats were significantly different with the exception of the yard in both months compared to the deciduous woodlot cluster in April. *k*-NN accurately assigned 100% of the bacterial profiles when the surface, 2, 10, and 60 inch profiles made up the training set.

Discussion

The three criteria described by Sensabaugh (35) that must be met for the application of any microbial-based technique in the forensic analysis of soil were all examined in this study. The first was the differentiation of two or more locations, in which a soil bacterial profile from one area is unique enough that it can be distinguished from another. Second, the technique must have high enough discriminatory power such that similar ecological habitats can be differentiated, but not so high that bacterial heterogeneity results in soil samples from the same location being deemed unique. And third, analysis techniques used to assess similarities or differences between soil bacterial profiles should possess a high level of objectivity. All three criteria were

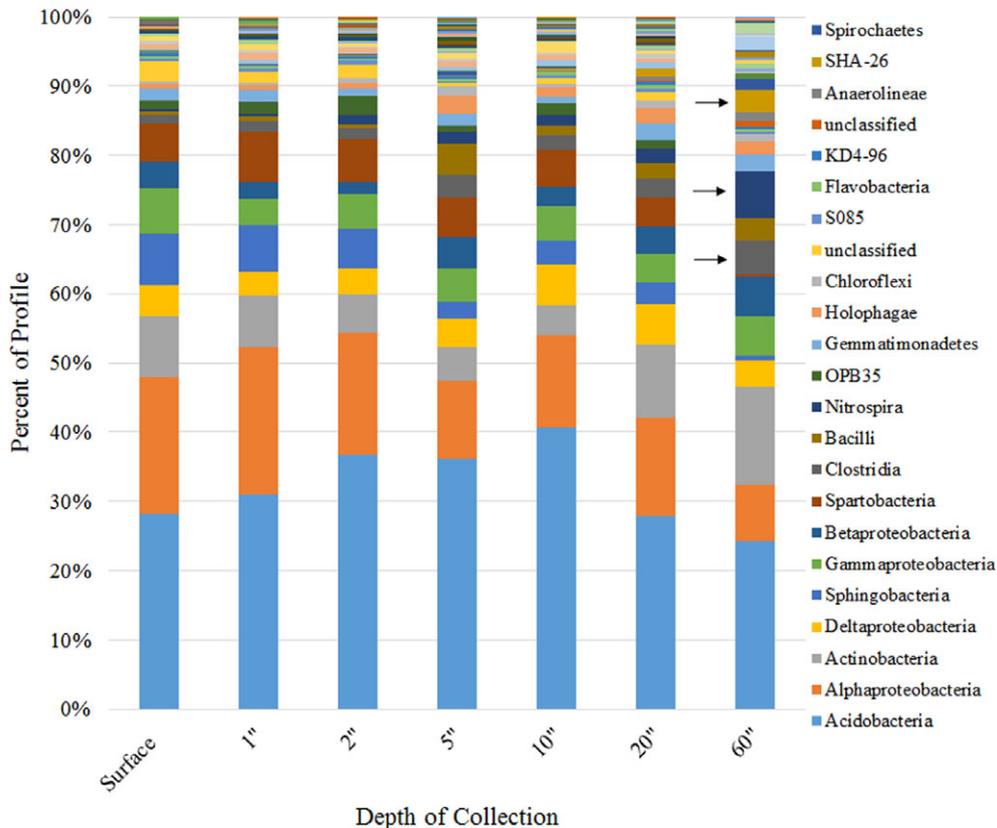


FIG. 7—Bacterial class abundance of woodlot depth samples in October 2013. As depth increased, substantial differences in *Clostridia*, *Nitrospira*, and *SHA-26* (denoted by arrows in ascending order) existed in all habitats. The 23 most abundant bacterial classes are listed on the right. (An interactive, color version of this chart is available as Fig. S3.)

met through the research presented here, in that next-generation sequencing of the bacterial 16S rRNA gene from soil samples produced highly discriminating data both within and among habitat types, which were conducive to objective analyses.

The most basic requirement of a bacterial profiling technique for forensic soil analysis is that it differentiates diverse habitats. In previous soil studies, researchers were often able to distinguish a small number of habitats in multidimensional space or through the presence or absence of specific bacteria (10,13,44–46), but overlap among habitats often occurred. In the current research, similar distinguishability and overlap in NMDS plots existed when examining many habitats simultaneously, but it was resolved when pairs or triads were compared. This increased resolution was accompanied by a decrease in stress, signifying the results better represented the bacterial profile associations that existed.

Once it was evident that diverse habitat types could be differentiated based on bacterial profiles, the next step was to determine whether the same was true for similar habitats. This presents a much larger challenge, as similar habitats are likely to share many of the same physical and chemical characteristics that can affect bacterial communities. Most of the researchers who have collected soil from similar habitats (e.g., 10,44–46) generally pooled ecologically similar habitats together for diverse habitat comparison rather than attempting to differentiate the similar ones. Although pooling soils is useful for microbiologists attempting to define basic bacterial properties of soil types, forensic scientists need to distinguish both ecologically diverse and similar locations. The latter, more challenging task was largely achieved in the current research through differentiation of several woodlot sites within close proximity. The exception was woodlot 5, which exhibited substantial intralocation variation over the eight-week period. Further investigation revealed that this location was directly adjacent to a large gravel pit that had been recently converted to a park (47), which may explain why it exhibited such a large amount of spatial heterogeneity.

With next-generation sequencing of soil bacteria's potential for differentiating diverse and similar habitats established, factors that could influence bacterial profiles within a location were considered. It is fundamentally impossible to collect known soil samples precisely when a crime occurs; consequently, temporal changes in bacterial makeup must be assessed. Past studies of temporal change have shown substantial differences in bacterial profiles collected over time based on T-RFLP or pyrosequencing of the 16S locus (12, 45 respectively), resulting in intermingling of habitats in multidimensional space and different levels of bacterial diversity over time, although there was no indication of seasonal or other temporal trends. Temporal fluctuations were evident in the current study as well, but again there were no predictable seasonal changes. More importantly, bacterial profiles produced through next-generation sequencing remained stable enough to correctly classify to their location of origin the vast majority of the time. These results indicate that bacterial profiles generated from soil collected weeks or months after a crime occurred will likely be representative of the location where the soil transfer took place, allowing for accurate association between evidentiary profiles and a location of origin.

It is also unlikely that known soil samples will be collected from the exact spot to which the evidentiary item was exposed, but instead could be collected feet, yards, or greater distances away, highlighting the importance of considering spatial variability of soil within a location. Bacterial variation over small distances has been attributed to microenvironmental factors such as foliage, pH, and nutrient supply (48,49), although in reality, an

almost unlimited number of factors could come into play. In a forensic study, Meyers and Foran (12) described variability, sometimes substantial, among soil samples collected 10 feet away from one another within habitats. Similar variation or patchiness was reflected in soil samples collected in the current research; however, this was overcome through the collection of multiple samples across a habitat to capture the variability that inherently exists. In this study, misclassified soils were always 90 or more feet away from all or four of the five training profiles, with the highest classification success utilizing a range of samples in a spiral pattern across the surface of a habitat. Representative training sets were also essential when analyzing soils collected at different depths within a habitat. Depth-based bacterial profiles were visually different in both abundance charts and NMDS plots, but accurate assignment to a habitat of origin was still achieved using *k*-NN when the training set was built from a range of soil depths. In contrast, training sets comprised of bacterial profiles from the top three or bottom three depth samples resulted in lower assignment accuracy of the remaining depth samples from that site (data not shown), again accenting the importance of collecting a range of samples to capture bacterial variation for a training set.

The second major goal of this research was to investigate varied techniques for analyzing next-generation sequencing-based bacterial profile data, which ideally possess objective qualities to meet the recommendations in the 2009 NAS report (29), while being comprehensible to a trier of fact. The massive amount of data produced via next-generation sequencing makes meeting these requirements challenging, as it is impossible to simply look at hundreds of thousands of DNA sequences and come to a definitive conclusion regarding the soil's origin. Therefore, techniques that can sort and display these datasets are vital. In reality, it is unlikely a single analysis technique will encompass both forensic needs; thus, it is quite possible more than one technique will be necessary for effective forensic purposes.

The most basic strategy for simplifying the massive datasets produced in studies such as this is to group them based on bacterial reference sequences and visualize them via abundance charts. These charts provide a graphic quantification of what bacteria are present in a profile, which should facilitate expert witness testimony and draw attention to profile variances. For instance, the dirt road ordinated separately from all other habitats in NMDS plots, but the extent and bacterial cause of that difference was not realized until abundance charts were examined. Upon further investigation, it was established that the road was treated with calcium chloride twice each summer to reduce dust levels (Shiawasse County Road Commission, personal communication, 2015), which apparently strongly impacted the bacterial community, and indeed some of the bacterial classes that existed at unusually high levels (e.g., *Gammaproteobacteria* and *Flavobacteria*) are known to thrive in such halophilic conditions (50,51). Scientists and nonscientists alike would be able to easily perceive these distinct differences based on abundance charts, allowing the expert witness to better explain their results and conclusions. However, we cannot rely solely on visual comparisons of abundance charts to associate soil bacterial profiles, as such assessments would be subjective, a forensic science weakness emphasized in the 2009 NAS report (29).

An attempt to objectively analyze and compare bacterial profiles may begin with numerically measuring how similar or dissimilar the profiles are. In this research, two dissimilarity indices, Sørensen–Dice and Bray–Curtis, were calculated for use in NMDS and *k*-NN analysis. The former index outperformed in all facets of analysis, resulting in tighter location of origin clus-

tering in NMDS plots and higher classification accuracy with k -NN. This likely resulted from how the two indices' are generated, where Sørensen–Dice is calculated based only on how many unique sequences are shared across two profiles, while Bray–Curtis is calculated based on both the number and abundance of shared sequences between two bacterial profiles. Misassignments of bacterial profiles using Bray–Curtis may have resulted from small bacterial abundance fluctuations that occur over time and space, an obvious detriment for forensic analyses. Sørensen–Dice is not sensitive to bacterial abundance fluctuations, which in this study resulted in more accurate assignment of soils to their location of origin.

NMDS was highly beneficial for analyzing soil data based on its success in clustering bacterial profiles from a given location, while simultaneously distinguishing profiles that were distinct from that cluster, both of which were easily visualized. However, NMDS clusters are formed via the rank order of dissimilarities for all soil bacterial profiles being compared, meaning a single highly dissimilar profile or set of profiles can force unrelated samples together, potentially resulting in misleading indications of similarity among them. This was exemplified in the diverse habitat study, in which the dirt road profiles clustered away from the other habitats, while apparently forcing all the rest into close proximity. This is a potential danger of comparing too many habitats at a time using NMDS, and indeed when fewer habitats were ordinated, they readily formed separate clusters, and stress was lowered. The same was true in the similar habitat study, with the exception of woodlot 5, as discussed above. Thus, a caveat for evaluating these types of data via NMDS is that analyzing too many locations at once has the potential to skew results, particularly if one location plots well away from all the others.

Another benefit of NMDS is that statistical analysis such as ANOSIM can be performed based on the clusters identified in the plots, providing a more objective statistical measure of soil profile similarity/difference. However, designation of a cluster is itself subjective, a characteristic that ideally is avoided. Also, ANOSIM and related statistical measures require groups of samples, which will often not be available for evidentiary materials. Thus, it seems best to use NMDS as a visualization method, while using other strategies for purely objective analysis of the data.

Supervised classification methods have the potential to act as an objective measure for comparing soil bacterial profiles; thus, a baseline version (52), k -NN, was tested in this research. k -NN classification correctly assigned soils to their place of origin 95.4% of the time, and similar to NMDS, k -NN-based assignments were less likely to be in error when known from fewer locations were used to produce the model. A drawback of k -NN itself is that it is a hard classifier, meaning an assignment will always be made for an unknown soil profile, even if it is not similar to any samples in the training set. For forensic purposes, this is certainly problematic, as it would result in an assignment even when the actual evidentiary soil location of origin was not represented in the knowns. In this regard, other supervised classifiers that provide a stronger statistical goodness of fit measure, such as decision trees (53) or soft independent modeling by class analogy (54), are likely better suited for forensic soil analysis. However, the extremely high success rate of k -NN classification in this study indicates these more complicated supervised classifiers could be expected to perform very well for forensic soil analysis.

There are obvious pros and cons associated with each of the soil bacterial profile analysis techniques tested here, and based on this research, it may be worthwhile to utilize more than one when examining forensic soil evidence, integrating both objective and visual interpretation of the data. Clear visual representations of bacterial profiles that could aid the jury's understanding of soil evidence were generated through abundance charts and NMDS plots, which acted in a complementary manner wherein the former provided a categorization and quantification of the copious sequences and the latter produced information on soil sample associations. Together, these or similar data visualization techniques can then be used by the forensic scientist to explain results obtained using more objective techniques, such as supervised classifiers. These varied analysis techniques demonstrate how a combination of methods can provide both visual and statistical interpretation of data, offering an optimized avenue for forensic soil analysis to enter the courtroom.

Based on the research presented here, next-generation sequencing of the bacterial 16S rRNA marker shows tremendous potential for forensic soil analysis. Key to this was establishing that soil samples could be reliably and accurately differentiated or associated based on bacterial profiles. This was achieved for diverse habitats, as well as for the vast majority of very similar habitats that were all within close proximity. Further, spatial and temporal factors, which will undoubtedly come into play for forensic soil comparisons, had little negative influence on soil traceability using this methodology. Overall, the DNA-based identification strategies presented here demonstrate the utility of next-generation sequencing in producing soil bacterial profiles, helping to link a suspect, victim, or evidentiary item to a crime scene.

Acknowledgments

The authors thank current and former members of the Michigan State University Forensic Biology Laboratory, particularly those who worked on earlier forensic soil studies. We thank the Fenner Nature Center for allowing soil collection on their grounds. Also, thank you to the members of the Michigan State University Genomics Core Facility, especially Jeff Landgraf, for their help with this research.

References

1. Dawson L, Hillier S. Measurement of soil characteristics for forensic applications. *Surf Interface Anal* 2010;42(5):363–77.
2. Scientific American. Curious use of the microscope. *Science and Art* 1856;11:240.
3. Murray RC, Solebello LP. Forensic examination of soil. In: Saferstein R, editor. *Forensic science handbook*, Vol. 1, 2nd ed. Boston, MA: Pearson Education Inc, 2002;615–33.
4. Ruffell A. Forensic pedology, forensic geology, forensic geoscience, geo-forensics and soil forensics. *Forensic Sci Int* 2010;202(1–3):9–12.
5. Daniel R. The metagenomics of soil. *Nat Rev Microbiol* 2005;3:470–8.
6. Koch R. Zur Untersuchung von pathogenen Organismen. *Aus: Mitteilungen aus dem Kaiserl. Gesundheitsamte*, Bd. I 1881;1:112–63.
7. Al-Awadhi H, Dashti N, Khanafer M, Al-Mailem D, Ali N, Radwan S. Bias problems in culture-independent analysis of environmental bacterial communities: a representative study on hydrocarbonoclastic bacteria. *SpringerPlus* 2013;2:369.
8. Woese CR, Fox GE. Phylogenetic structure of the prokaryotic domain: the primary kingdoms. *Proc Natl Acad Sci USA* 1977;74(11):5088–90.
9. Liu WT, Marsh TL, Cheng H, Forney LJ. Characterization of microbial diversity by determining terminal restriction fragment length polymorphisms of genes encoding 16S rRNA. *Appl Environ Microbiol* 1997;63(11):4516–22.
10. Fierer N, Jackson JA. The diversity and biogeography of soil bacterial communities. *Proc Natl Acad Sci USA* 2006;103(3):626–31.

11. Horswell J, Cordiner SJ, Maas EW, Martin TM, Sutherland KB, Speir TW, et al. Forensic comparison of soils by bacterial community DNA profiling. *J Forensic Sci* 2002;47(2):350–3.
12. Meyers MS, Foran DR. Spatial and temporal Influences on bacterial profiling of forensic soil samples. *J Forensic Sci* 2008;53(3):652–60.
13. Lenz EJ, Foran DR. Bacterial profiling of soil using genus-specific markers and multidimensional scaling. *J Forensic Sci* 2010;55(6):1437–42.
14. Jonasson J, Olofsson M, Monstein HJ. Classification, identification and subtyping of bacteria based on pyrosequencing and signature matching of 16S rDNA fragments. *Apmis* 2002;110(3):263–72.
15. Luo C, Rodriguez RLM, Johnston E, Wu L, Cheng L, Xue K, et al. Soil microbial community responses to a decade of warming as revealed by comparative Metagenomics. *Appl Environ Microbiol* 2014;80:1777–86.
16. Kravchenko AN, Negassa WC, Guber AK, Hildebrandt B, Marsh TL, Rivers ML. Intra-aggregate pore structure influences phylogenetic composition of bacterial community in macroaggregates. *Soil Sci Soc Am J* 2014;78(6):1924–39.
17. Cole JR, Chai B, Farris RJ, Wang Q, Kulam SA, McGarrell DM, et al. The Ribosomal Database Project (RDP-II): sequences and tools for high-throughput rRNA analysis. *Nucleic Acids Res* 2004;33(1 Suppl):D294–6.
18. Quast C, Priesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res* 2013;41(D1):D590–6.
19. Loman NJ, Misra RV, Dallman TJ, Constantinidou C, Gharbia SE, Wain J, et al. Performance comparison of benchtop high-throughput sequencing platforms. *Nat Biotechnol* 2012;30(5):434–9.
20. Caporaso GJ, Lauber CL, Walters WA, Berg-Lyons D, Huntley J, Fierer N, et al. Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J* 2012;6:1621–4.
21. Pietraszkiewicz BU. Exploring soil bacterial communities for forensic applications: a genomics approach. Washington, DC: U.S. Department of Justice 2009-IJ-CX-0021 Report, 2012.
22. Sutton NB, Maphosa F, Morillo JA, Al-Soud WA, Langenhoff AAM, Grotenhuis T, et al. Impact of long-term diesel contamination on soil microbial community structure. *Appl Environ Microbiol* 2013;79(2):619–30.
23. Korenblum E, Bastos Souza D, Penna M, Seldin L. Molecular analysis of the bacteria communities in crude oil samples from two Brazilian offshore petroleum platforms. *Int J Microbiol* 2012;2012:1–8.
24. Magnabosco C, Tekere M, Lau MCY, Linage B, Kuloyo O, Erasmus M, et al. Comparisons of the composition and biogeographic distribution of the bacterial communities occupying South African thermal springs with those inhabiting deep subsurface fracture water. *Front Microbiol* 2014;5:679.
25. Pye K, Blott SJ, Croft DJ, Carter JF. Forensic comparison of soil samples: assessment of small-scale spatial variability in elemental composition, carbon and nitrogen isotope ratios, colour, and particle size distribution. *Forensic Sci Int* 2006;163(1–2):59–80.
26. Heath LE, Saunders VA. Assessing the potential of bacterial DNA Profiling for forensic soil comparisons. *J Forensic Sci* 2006;51(5):1062–8.
27. Clarke KR. Non-parametric multivariate analyses of changes in community structure. *Aust J Ecol* 1993;18:117–43.
28. Yang C, Mills D, Mathee K, Wang Y, Jayachandran K, Sikaroodi M, et al. An ecoinformatics tool for microbial community studies: supervised classification of Amplicon Length Heterogeneity (ALH) profiles of 16S rRNA. *J Microbiol Methods* 2006;65(1):49–62.
29. National Research Council. Strengthening forensic science in the United States: a path forward. Washington, DC: National Academies Press, 2009.
30. Hopkins JM. Forensic soil bacterial profiling using 16S rRNA gene sequencing and diverse statistics [thesis]. East Lansing, MI: Michigan State U, 2014.
31. Barnard RL, Osborne CA, Firestone MK. Responses of soil bacterial and fungal communities to extreme desiccation and rewetting. *ISME J* 2013;7:2229–41.
32. Hyde ER, Haarmann DP, Lynne AM, Bucheli SR, Petrosino JF. The living dead: bacterial community structure of a cadaver at the end of the bloat stage of decomposition. *PLoS ONE* 2013;8(10):e77733.
33. Ramette A. Multivariate analyses in microbial ecology. *FEMS Microbiol Ecol* 2007;62(2):142–60.
34. Mohri M, Rostamizadeh A, Talwalkar A. Foundations of machine learning. Cambridge, MA: The MIT Press, 2012.
35. Sensabaugh GF. Microbial community profiling for the characterization of soil evidence: forensic considerations. In: Ritz K, Dawson L, Miller D, editors. Criminal and environmental soil forensics. New York, NY: Springer, 2009;49–60.
36. Haas BJ, Gevers D, Earl AM, Feldgarden M, Ward DV, Giannoukos G, et al. Chimeric 16S rRNA sequence formation and detection in Sanger and 454-pyrosequenced PCR amplicons. *Genome Res* 2011;21(3):494–504.
37. Caporaso GJ, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA, Turnbaugh PJ, et al. Global patterns of 16S rRNA diversity at a depth of millions of sequences per sample. *Proc Natl Acad Sci USA* 2010;108(1 Suppl):4516–22.
38. Schloss PD, Westcott SI, Ryabin T, Hall JR, Hartmann M, Hollister EB. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol* 2009;75(23):7537–41.
39. Bray JR, Curtis JT. An ordination of the upland forest communities of southern Wisconsin. *Ecol Monogr* 1957;27(4):325–49.
40. Dice LR. Measures of the amount of ecologic association between species. *Ecology* 1945;26(3):297–302.
41. Sørensen T. A method of establishing groups of equal amplitude in plant sociology based on similarity of species and its application to analyses of the vegetation on Danish commons. *Kongelige Danske Videnskaberne Selskab* 1948;5(4):1–34.
42. Kruskal JB. Nonmetric multidimensional scaling: a numerical method. *Psychometrika* 1964;29(2):115–29.
43. Hammer Ø, Harper DAT, Ryan PD. PAST: paleontological statistics software package for education and data analysis. *Palaeontologia Electronica* 2001;4(1):1–9.
44. Lauber CL, Hamady M, Knight R, Fierer N. Pyrosequencing-based assessment of soil pH as a predictor of soil bacterial community structure at the continental scale. *Appl Environ Microbiol* 2009;75(15):5111–20.
45. Lauber CL, Ramirez KS, Aanderud Z, Lennon J, Fierer N. Temporal variability in soil microbial communities across land-use types. *ISME J* 2013;7(8):1641–50.
46. Kim M, Boldgiv B, Singh D, Chun J, Lkhagva A, Adams JM. Structure of soil bacterial communities in relation to environmental variables in a semi-arid region of Mongolia. *J Arid Environ* 2013;89:38–44.
47. <http://pk.ingham.org/Parks/HawkIsland.aspx>
48. Ettema CH, Wardle DA. Spatial soil ecology. *Trends Ecol Evol* 2002;17(4):177–83.
49. Eichorst SA, Breznak JA, Schmidt TM. Isolation and characterization of soil bacteria that define terriglobus gen. nov., in the Phylum Acidobacteria. *Appl Environ Microbiol* 2007;73(8):2708–17.
50. Sorokin DY, Kovaleva OL, Tourova TP, Muyzer G. *Thiohalobacter thio-cyanaticus* gen. nov., sp. nov., a moderately halophilic, sulfur-oxidizing gammaproteobacterium from hypersaline lakes, that utilizes thiocyanate. *Int J Syst Evol Micr* 2010;60:444–50.
51. Ventosa A, Nieto JJ, Oren A. Biology of moderately halophilic aerobic bacteria. *Microbiol Mol Biol R* 1998;62(2):504–44.
52. Lavine BK, Davidson CE. Classification and pattern recognition. In: Gemperline P, editor. Practical guide to chemometrics, 2nd edn. Boca Raton, FL: CRC Press, 2006;339–78.
53. Quinlan JR. Induction of decision trees. *Mach Learn* 1986;1:81–106.
54. Gold S, Sjostrom M. SIMCA: a method for analyzing chemical data in terms of similarity. In: Kowalski BR, editor. Chemometrics: theory and application. Washington DC: American Chemical Society, 1977;243–82.
55. Tukey JW. Bias and confidence in not-quite large samples. *Ann Math Stat* 1958;29:614–15.

Additional information and reprint requests:

David R. Foran, Ph.D.

School of Criminal Justice and Department of Integrated Biology

560 Baker Hall

655 Auditorium Road

Michigan State University

E. Lansing, MI 48824

E-mail: foran@msu.edu

Supporting Information

Additional Supporting Information may be found in the online version of this article:

Figure S1. Interactive, color version of Fig. 1. Average ($n = 5$) bacterial class abundance of ten diverse habitats. The dirt road clearly differed from the other habitats, containing higher levels of *Gammaproteobacteria*, *Flavobacteria*, *Clostridia*, and *Bacilli* (denoted by arrows in ascending order on the right) along with lower levels of *Acidobacteria* and *Betaproteobacteria*

(denoted by arrows in ascending order on the left). The 23 most abundant bacterial classes are listed on the right. Hover over the colored segments to reveal bacterial class identity. Ag = Agricultural.

Figure S2. Interactive, color version of Fig. 4. Average ($n = 5$) bacterial class abundance of woodlot locations. The soils shared major bacterial classes up to approximately 80%. The 23 most abundant bacterial classes are listed on the right. Hover over the colored segments to reveal bacterial class identity.

Figure S3. Interactive, color version of Fig. 7. Bacterial class abundance of woodlot depth samples in October 2013. As depth increased, substantial differences in *Clostridia*, *Nitrospira*, and *SHA-26* (denoted by arrows in ascending order) existed in all habitats. The 23 most abundant bacterial classes are listed on the right. Hover over the colored segments to reveal bacterial class identity.

Table S1. ANOSIM pairwise p -values for diverse habitats (A), similar habitats (B), and depth study habitats (C).