

A Discrete, Stochastic Model and Correction Method for Bacterial Source Tracking

MARK D. LEACH,^{†,‡}
SHIRA L. BROSCHEAT,^{*,†,‡,§} AND
DOUGLAS R. CALL^{†,§}

School of Electrical Engineering and Computer Science, Center for Integrated Biotechnology, and Department of Veterinary Microbiology and Pathology, Washington State University, Pullman, Washington 99164

Received April 20, 2007. Revised manuscript received September 18, 2007. Accepted September 19, 2007.

We have developed a model to test several underlying assumptions of bacterial source tracking (BST) when the BST method is based on detection of discrete genetic markers from source-specific bacteria. The model consists of an environment and discrete-time input signals that represent sources of contamination partitioned into marker-bearing and nonmarker-bearing units “shed” into the environment. Simulations run for different types of environmental contamination patterns indicate that if hosts shed different percentages of BST markers, the environment cannot be accurately characterized unless a correction method is used. The correction method, which requires the solution to a linear system, reduces the mean error in estimating the proportions of host contamination to below 3%. The effectiveness of the method depends on accurate knowledge of the occurrence and prevalence of markers in the various hosts; this may be a challenging task, especially if these values vary across populations in space and time. In addition, the correction method does not compensate for environments with low-density or unmixed contamination. In conclusion, our simulations highlight several fundamental challenges that may prevent absolute quantification of fecal input using discrete marker BST.

Introduction

Bacterial source tracking (BST) refers to techniques for identifying sources of fecal contamination, distinguishing between host origin, and assigning proportions to individual sources. BST techniques utilize bacterial indicators that are not themselves pathogenic, but which imply the presence of pathogenic bacteria or viruses. Recent efforts have focused on library-independent BST methods that identify genetic markers for indicator bacteria specific to particular hosts (15). In theory, the occurrence of these discrete markers should be constant through time and over a wide geographic area, but questions regarding the distribution of the markers among hosts and the validity of field applications remain unanswered.

DNA markers specific to human and nonhuman sources have been identified (3, 4, 8, 9, 11, 12, 14). A number of

techniques have been used to identify these markers including the screening of fecal samples for specific 16S rRNA sequences (1, 2, 10), using mixed-genome microarrays to identify bacterial genomic DNA markers (14), and employing suppression subtraction hybridization libraries to identify unique markers (8). Polymerase chain reaction (PCR) is then employed to detect markers in fecal and environmental samples. It is unclear, however, how many isolates need to be screened, as well as from how many hosts and from what geographic regions, to validate the markers. Our experience in screening markers suggests that they may show adequate host-specificity in PCR screening at the fecal sample testing stage, but they may be unable to distinguish between hosts in field studies because of ubiquitous distribution in water samples (D.R. Call, unpublished data). While this may be due to differences in the survival of marker strains in aqueous environments, based on laboratory screening of isolates, marker distribution is not predictable. Even if accurate quantitative methods are available, each marker may represent different masses of fecal pollution; without a priori knowledge of this representational bias, it is difficult to attribute volume of waste according to marker abundance. Further complications include our lack of knowledge of the relative health risks posed by fecal contamination from different hosts and recent evidence that BST markers correlate poorly with fecal indicators (10).

The measurement of markers in the field is closely related to the problem of marker characterization. We must establish criteria for interpreting results based on marker prevalence in samples from known sources. That is, assuming our marker characterization is correct, we need to know how to interpret measurements when there is nonspecific host occurrence or differences in the prevalence between hosts. For instance, a small nonspecific percentage of marker may lead to ambiguous results if the nonspecific host is a large contributor to the waste stream. Nevertheless, if we know the relative prevalence of the different markers among hosts, and we can make a quantitative measurement of the sample, we should be able to correct the sampled data to estimate the actual host contributions in the environment.

In this paper we develop a model for addressing some of the issues that arise with the use of discrete genetic markers from source-specific bacteria. A primary motivation is that there are no gold standards by which to assess the interpretation of BST marker data from field samples. Models provide us with a means to compare estimates with known distributions. In particular, we are interested in a model that will represent significant statistical variation with regard to the distribution of bacteria in an environment, the difference in prevalence among markers, the nonspecific occurrence of markers, and the effects of differences in survival between marker-bearing bacteria. Our objective is to derive simplified scenarios that can be used to locate primary problems in the BST methodology that need further attention. If problems arise in the simplified case, we assume they will also be present in real-world applications.

Model Development. We modeled the BST system as a multiinput, multioutput system where discrete-time input signals represent fecal contamination entering an environment, and the output signals represent statistics derived from the measurement of markers in a sample space. Fecal contamination is partitioned into abstract units that act as cellular automata. For each time step the internal state of the environment evolves according to simple rules that govern the automata. Some of the rules define random processes

* Corresponding author phone: (509) 335-5693; Fax: (509) 335-3818; e-mail: shira@eecs.wsu.edu.

[†] School of Electrical Engineering and Computer Science.

[‡] Center for Integrated Biotechnology.

[§] Department of Veterinary Microbiology and Pathology.

and this introduces variation that simulates the uncertainty in actual environments.

Environment. The environment is an abstract representation of a waterway in which we create an environmental model that is as simple as possible but still able to simulate significant issues of variation in the BST process. In the sampling process, we test for the presence of markers in a waterway such as a stream, river, or lake. The essential character of the environment is a one-dimensional surface, such as a riverbank, where points along the surface correspond to variations in the concentration of fecal contamination. Accordingly, our model consists of a one-dimensional array of bins. The sources of contamination are distributed along the array and each source sheds contamination into proximal bins. A sample space for measurement of the environment consists of a subset of the total bins.

Hosts. The model has five different host species arbitrarily designated cow, dog, elk, goose, and human. We used these names from previous studies (8, 14); they do not reflect any biological parameters in the simulation. Each host type has the same range of possible parameter values.

In actual environments sources of fecal contamination can be classified as point and nonpoint. Point sources include agriculture and human, such as runoff from pastures or sewer overflows. Nonpoint sources can include wildlife, such as elk or waterfowl, where the entry point of contamination is not readily apparent. Each source in the model sheds contamination of a single host type and at a rate assigned to that host. Each source also has a location along the environment bin array. Nonpoint sources are modeled by distributing contamination in a uniform distribution over the entire environment, whereas point sources are modeled by distributing the contamination in a normal distribution with the mean at the source location. Each source is a constant, discrete signal $S_i(k) = K_i$ where the constant K_i is determined by the host type.

Cells. In an aqueous environment the bacterial cells are small and numerous enough to consider the level of contamination a continuous variable. In the simulation, however, a continuous variable presents a problem in terms of mixing and sampling. We need a method of sampling a “chunk” of the environment that contains random levels and types of contamination from the various hosts. By dividing the contamination into discrete units and randomly distributing them into the environment, we can simulate the mixing of bacterial contamination without specifying the environment as a whole. Thus, our approach is to specify only the behavior and traits of the discrete units (cells). We do not define the distribution of cells in the environment directly, but allow it to evolve over time through the action of a large number of these cellular automata.

Factors such as sunlight, temperature, and species affect bacteria survival and persistence in aqueous environments (5–7, 13). In the model, we assume that bacteria do not persist or replicate in the environment, but begin to die off as soon as shed from a source. We model cell decay according to the formula $N e^{-\alpha k}$ where N is the number of cells at time 0, and k is the time step. Each time step represents approximately one hour. A random lifespan is assigned to each cell from an exponential distribution and at each time step the cell ages until its age equals its lifespan; then the cell dies and is removed from the environment. From the perspective of the environment as a whole this constitutes exponential decay of cells over time. Each cell is also marked according to its host type and what marker, if any, it carries. This allows the measurement of cells in the environment, including counts of the whole environment and counts within a particular sample space.

Marker Matrix. In a best-case scenario, BST markers would be 100% specific to their assigned hosts, but in the

real world there may be nonspecific occurrence because of transient carriage or incomplete host-specificity. In our model, samples from known hosts are hypothetically tested against a group of markers, and the proportion of markers is determined for the host. For instance, analysis of cow feces may indicate this host produces indicator bacteria where 50% harbor a cow marker, 10% harbor a goose marker, and the remaining 40% harbor no recognized marker. This data can be represented in a column vector as $(0.5, 0.0, 0.0, 0.1, 0.0)^T$ where the elements represent the prevalence of a particular marker. When contamination from each host is characterized, the column vectors can be combined into a matrix representing the prevalence data. In the best case of perfect specificity, the marker matrix M can be written as the product of the identity matrix I and some constant. For instance, if each marker occurs in 70% of host contamination, then $M = 0.7 I$.

Our method of describing the distribution of markers may not be practical. We may be able to determine the relative percentages of marker in isolates (e.g., through quantitative PCR), but not know the true proportion of bacteria that harbors no marker. In addition, we may not know the amount of marker produced per unit weight of fecal contamination. Our model does not address these ambiguities but assumes an ideal marker matrix where the proportion of “no marker” is known, and the proportions represent the total contamination entering the environment.

Sampling. Two methods are used to estimate the contribution of hosts. In the first method, we randomly select a group of bins and count all of the marker-bearing cells in those bins. This method provides an exact quantification of the markers in the sample space. We can then convert the total counts for each host type to a percentage. After applying this same procedure to the entire environment, counting both marker-bearing and nonmarker-bearing cells together, we can compare the percentages of marker-bearing cells in the sample to the total cells in the environment. For example, if each host is equally represented in the environment and each source has the same shedding and decay rates, the percentages of cells in the environment are similar for each host with some variation. In practice, however, we can only measure the marker-bearing cells in the sample space and this can create significant error in the estimates of proportional contributions of each host (Table 1). The second method of estimating host contributions is to perform a presence/absence measurement of the bins in the sample space. In this case, we register a true or false response for each host per bin and then calculate the percentage of positive bins in the sample space for each host.

We also use two types of sample spaces: contiguous and uniform. In the contiguous case, the sample space consists of a block of contiguous bins selected at a uniform random location within the environment with the restriction that the position of the block fits within the environment. This simulates a sampling scheme whereby water samples are collected near each other. In the uniform case, each water sample (bin) is selected from a uniform random location within the model environment.

Environment Characterization. For different sets of model parameters the environment evolves in a distinct manner. For instance, if the environment is small and there are a large number of cells shed per time step, the environment will “fill up” and reach an equilibrium with a large number of cells in each bin. If the cells are shed in a uniform distribution, the environment will be dense and well-mixed. Each bin is likely to contain cells of each host type and the standard deviation of cells per bin will be relatively small compared to the average. Alternatively, if the environment is large or there are few cells shed per time step, the system will evolve to a sparse equilibrium state. There may be empty

TABLE 1. Environment Statistics^a

	cells in the environment				
	cow	dog	elk	goose	human
total cell count:	3502	3812	3618	3722	3666
percentage total cells:	19.1%	20.8%	19.7%	20.3%	20.0%
	cells in sampled bins				
	cow	dog	elk	goose	human
marker-bearing cell count:	125	195	76	188	163
percentage marker-bearing cells:	16.7%	26.1%	10.2%	25.2%	21.8%
percent error:	-12.6%	25.5%	-48.2%	24.1%	9.0%

^a Statistics derive from quantitative measurement of an environment sampled after 300 time steps and show cell counts, proportional volume of cells, and percent error between environment as a whole and the sample space.

bins and bins that lack cells or markers for a particular host. When the cells are shed in a Gaussian distribution about the sources, the variation is increased, particularly for sparse environments.

To characterize the density and mixing of the environment in one parameter, the model counts the cells in all bins and calculates the mean and standard deviation for each host. The percent coefficient of variation (COV) for each host ($\sigma/\mu \times 100$) serves as a measure of the state of the total environment. For instance, there may be an average of 80 cow cells (both marker- and nonmarker-bearing) per bin with a standard deviation of 10 cells, yielding a COV of 12.5%. The COV indicates the density of pollution in the environment because when bins contain large numbers of cells, there will be less “spread” in the distribution.

Correction Method. The columns of the marker matrix determine the number of marker-bearing cells in the environment. The relationship between marker prevalence profiles, the proportion of markers in the sample, and the proportional volume of total indicator bacteria can be represented as a system of linear equations. The system of equations can be represented as a matrix equation $M \cdot \mathbf{x} = \mathbf{b}$, where the goal is to measure \mathbf{b} from an environmental sample and solve for \mathbf{x} . We would expect the diagonal elements of the marker matrix to be nonzero (i.e., cows producing cow marker, dogs producing dog marker, and so on). This is a sufficient condition for $M \neq 0$ and the existence of a unique solution for the linear system. Once \mathbf{x} is found, we can normalize its N elements so that the proportional volumes p_i can be compared to the proportional volumes of total cells in the environment.

$$p_i = \frac{x_i}{\sum_{i=1}^N x_i} \quad (1)$$

This normalized solution represents the relative volumes of cells for each host in the sample space. The model compares this solution to the cell volumes in the environment as a whole.

Materials and Methods

We implemented the BST model as a simulation on a Linux machine. The input to the simulator includes two files. One file contains true random integers (www.random.org) that act as seeds for the random number generator. The simulator loads the seed file, uses some of the seeds during the simulation, and stores the remainder back into the source file. The simulator uses the seeds periodically during the simulation to avoid repetition of random number sequences and to ensure unique outputs from each simulation. The

other input file is a configuration file that holds the simulation parameters.

The output of a simulation includes a statistical summary and an output file that shows the contents of each bin, one per line, where the letters c, d, e, g, and h represent cells from the five hosts. This provides a visual representation of the environment. The simulation also includes output functions to measure the COV as a function of time. We used a curve-fitting tool (Logger Pro 3, Vernier Software & Technology, Beaverton, OR) to plot the COV data and determine parameters that describe the curve.

Results

We divided simulations into categories based on environment type and sampling method. We first established a “best-case scenario” and then modified the simulation parameters to examine model sensitivity given eight additional scenarios. Our goal was to determine which factors contribute the most uncertainty to interpretation of BST data and, as a result, to enable us to better define where future research should be focused. Each simulation included three sampling methods: presence/absence, quantitative, and corrected.

Environment Equilibrium. To derive a reasonable value for the sampling time, we measured the percent COV of selected environments as they evolved over time. We determined that the typical environment reaches an equilibrium and converges as k^B where B is derived from an empirical plot of the COV vs time (Supporting Information Figure S1). For example, for the “best-case scenario” the COV converges to a small value. The equation $Ak^B + C$ provides the closest fit to this curve (Figure S1) where $A = 221.5$, $B = -0.5232$, $C = 4.398$ and the root-mean-square error (rmsE) is 1.075. The rmsE for Ak^B was 1.473, indicating a better fit with a nonzero steady-state. The fitted equation also indicates that doubling the simulation time from 300 to 600 decreases the COV from 15.6 to 12.2%, and to reduce it to 5.96% requires 13000 time steps. We did not analyze the COV for every simulation, but assumed 300 time steps would yield a low percent change (0.12% in this example) between time steps.

Best-Case Scenario (No. 1). In the best-case environment each host has the same number of sources, and each source sheds the same number of cells per time step. The cells are also distributed uniformly in the environment. These conditions produce well-mixed environments containing approximately equal numbers of cells from each host. In addition, for a best-case scenario, the markers are specific to a single host and occur in the same percentages of isolates from each host ($M = 0.7 I$). Simulation with these parameters (Supporting Information Table S1) produced dense environments with an average of 60 cells per bin. We calculated the

TABLE 2. Average Percent Error and Associated Standard Deviation for Nine Simulation Scenarios^a

sampling scheme ^b	environmental scenario ^c								
	1	2	3	4	5	6	7	8	9
presence/absence average:	0.08	1.71	103.91	0.13	68.22	21.43	5.93	0.25	9.05
standard deviation:	0.54	25.96	17.12	0.65	6.13	95.66	20.99	0.82	30.85
quantitative average:	0.50	17.25	1.31	19.81	1.05	19.12	6.28	23.62	27.01
standard deviation:	6.02	27.36	9.09	6.86	6.50	125.02	31.14	6.12	27.92
corrected ^d average:				0.66				0.51	
standard deviation:				6.92				8.35	
average COV:	28.35	141.79	43.66	28.43	35.62	147.87	147.74	28.43	28.45

^a Each scenario was executed 30 times, and the values shown here represent the average error between estimates for samples and environments. ^b Three estimation procedures are reported. The presence/absence procedure measures the number of bins that contain at least one marker-bearing cell; the quantitative procedure measures the exact number of marker-bearing cells in each bin; the corrected measurement is derived from applying the correction method to the quantitative data. The COV is the percent coefficient of variation ($\sigma/\mu \times 100$) in the number of cells per bin. ^c The nine environment scenarios include (1) best-case; (2) low-density environment; (3) unequal cell volumes; (4) unequal marker prevalence; (5) unequal decay rates; (6) unmixed environment, contiguous sample space; (7) unmixed environment, uniform sample space; (8) nonspecific markers; (9) low marker prevalence. ^d The corrected data are given only for those scenarios that include variation in marker prevalence.

mean and standard deviation of 30 runs and report the average error between the actual volumes of cells in the environment and the volumes measured from the sample space (Tables 2 and S1). Negative values indicate an underestimation and positive values an overestimation of the actual proportional volumes. Both sampling schemes produced <1% error in proportional estimates with standard deviation ranging from 0.5 (presence/absence) to 6.02 (quantitative). The percent COV is low for a best-case scenario (<30.0), which is indicative of a dense, well-mixed environment.

Low Density Scenario (No. 2). To create a sparse environment, we used the parameters from the best-case scenario but decreased the signal strengths to 1 and increased the environment size to 2000 bins. This is equivalent to less pollution dispersed over a broader area. As expected, a sparse environment generated significantly greater variance and resulted in underestimates for contributions from each host. From a sampling perspective the variance was higher for both methods, but average error was considerably greater (17%) for quantitative sampling compared to presence/absence (1.7%) (Table 2 and Supporting Information Table S2).

Unequal Cell Volume Scenario (No. 3). In the previous example, the measure for presence/absence of specific markers performed well because each marker is equally distributed among host bacteria. Because this measurement registers presence for any marker-bearing cells within a bin, it tends to estimate equal representation of hosts in uniform environments. When proportional cell volume is changed (signal strength), however, estimates of presence/absence for different host-specific markers are overestimated for hosts with a lower volume and underestimated for those cases with a higher volume of cells with the net result being considerable error (100%) (Tables 2 and Supporting Information S3). Unequal cell volume is equivalent to having unequal contributions of fecal contamination between hosts. The quantitative measure performs very well in this scenario (mean error 1.31%) although the standard deviation is elevated compared to the best-case scenario (Table 2).

Unequal Marker Prevalence (No. 4). For this scenario we included best-case parameters with 100% host-specificity and an equal volume of cells shed by each host, but the prevalence of markers in the indicator bacteria population was variable (Table 3). This is equivalent to the case when a marker for one host is found frequently in the GI flora, whereas a marker for another host is found rarely. This scenario produced considerable error in the quantitative measure (19.8%) while the presence/absence estimates remained quite robust (0.13%) (Tables 2 and Supporting

TABLE 3. Marker Prevalence Matrices^a

Non-uniform prevalence	Non-specific occurrence
$\begin{pmatrix} 0.5 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.8 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.4 & 0.0 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.6 & 0.0 \\ 0.0 & 0.0 & 0.0 & 0.0 & 0.7 \end{pmatrix}$	$\begin{pmatrix} 0.5 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 0.5 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.1 & 0.5 & 0.1 & 0.1 \\ 0.0 & 0.0 & 0.0 & 0.5 & 0.0 \\ 0.1 & 0.0 & 0.1 & 0.0 & 0.5 \end{pmatrix}$

^a Matrix columns indicate the relative percentages of markers found in five hosts (cow, dog, elk, goose, human), and rows correspond to five distinct markers. These matrices correspond to the marker prevalence in simulation scenarios 4 and 8, respectively.

Information S4). Nevertheless, when one or more of the prevalence parameters falls below 10%, the presence/absence estimates begin to show high error rates due to low marker-bearing cell density (data not shown). When quantitative estimates are corrected for the underlying prevalence matrix, the estimated proportions are corrected accordingly (Table 2). That is, when marker prevalence is known, quantitative estimates can be corrected.

Unequal Decay Rates (No. 5). The cells decay exponentially according to the parameter α in e^{α} where $\alpha \leq 0$. Altering the values of α (via the exponential distribution parameters) had an effect similar to altering the signal strength. The error increased dramatically for presence/absence sampling, causing an overestimation of lower host volumes and an underestimation of higher host volumes, whereas quantitative estimates were very robust (Tables 2 and Supporting Information S5). This scenario might arise if two different species of bacteria with differential survivorship were chosen as fecal markers for two different host populations.

Unmixed Environment (No. 6). To simulate an unmixed environment, we ran the simulation with the cell distribution following a Gaussian pattern with a standard deviation of 20 bins and a sample space of 30 contiguous bins. Sampling in an unmixed environment such as this results in large errors regardless of sampling method because the sample space is likely to occur in an area dominated by a single host (Tables 2 and Supporting Information S6). The correction factor has no effect in this case because the variation is not dependent on the marker matrix. This scenario, while producing

unsatisfactory estimates, probably represents a closer approximation to reality than a uniformly mixed model.

We can compensate for the large mean error rate in this scenario by changing from a contiguous to uniform random sampling strategy (scenario no. 7). By sampling bins from the entire environment, we are compensating for the effect of sampling in a localized area of the environment (contiguous bins), which is dominated by cells from a nearby contamination source. In this case, both presence/absence and quantitative errors were reasonable (<6.3%), but standard deviations remained high (Tables 2 and Supporting Information S7).

Nonspecific Markers (No. 8). Markers are likely to have some level of nonspecific occurrence, either due to transient carriage or incomplete specificity (Table 3). In this scenario, we introduce nonspecific markers at the 10% level; that is, while a marker may be present in 50% of host-specific isolates, it can also be found in up to 10% of isolates from nonspecific hosts. For example, the marker matrix can be set so that elk and human markers are found in other hosts (Supporting Information Table S12), and this leads to significant error in the quantitative measure (Tables 2 and Supporting Information S8). The presence/absence estimation remains close to best-case, but this is not surprising because at the density employed in this particular simulation, there is likely to be at least one cell from each host in each bin (Table 2). Importantly, this scenario demonstrates that if the underlying distribution matrix is known, the correction algorithm generally recovers accurate quantitative estimates.

Low Marker Prevalence (No. 9). Our simulations also show low marker shedding rates can yield significant problems. Soule et al. (14) reported that as few as 2% of *Enterococcus* tested from a given host were marker-bearing isolates. To simulate this shedding level, we used a marker matrix $M = 0.02 I$. The results show significant error in the quantitative measure and unsatisfactory performance for presence/absence sampling (Tables 2 and Supporting Information S9). In effect, this error is due to sparse marker occurrence within an otherwise dense environment.

Discussion

Our simulations highlight two significant sources of error when applying genetic BST markers: first, errors that arise from variation in the distribution and density of fecal contamination in the environment; second, errors that arise from variation in the levels of marker occurrence and specificity among hosts. In the first case we derived significant error from environments that are sparse (low-density) in terms of total amount of contamination or in marker-bearing contamination and from environments that are not well-mixed. Given an accurate method for quantitative measurement and a sufficiently large sample space, environments that vary in terms of total amount of contamination per host did not show significant error. We also found a more pronounced error in environments that were both sparse and unmixed. In practice, we might hypothesize that real environments are likely to have one or both of these features. Furthermore, our coefficient of variation measurement (which indicates the degree of mixing and density of contamination) would likely be a time-varying parameter; if samples are obtained over a period of days, the environment may be well-mixed and dense on one occasion and not on another. Such variation is also likely at different sampling locations. Further experiments are needed to determine how actual environments correspond to models that include sparseness and mixing as variables. Without experimental data regarding the validity of these parameters, it is difficult to estimate the significance of the second type of error. At the very least our simulations suggest that more studies are

needed before we can have confidence in measurements that rely on quantification.

In terms of variation in the levels of marker occurrence, we found that significant error arises when markers do not represent the actual amount of fecal contamination due to each host. Another obvious source of error arises when there is nonspecific marker occurrence. The former variation is perhaps the most significant error suggested by our simulations in that, in real-world applications, we would be highly unlikely to find markers that are shed at exactly the same percentage of total contamination for all hosts. For this reason, our results indicate that without accurate knowledge of the levels of marker occurrence, confidence for any quantitative measurements will be low. To apply discrete marker BST one would need to first determine the amount of marker produced by each host and then determine quantitative measurements of environmental samples. For instance, for a given amount of contamination from cow hosts one would need to determine what percentage contains no marker, cow marker, dog marker, etc. Soule et al. (14) suggested quantitative PCR as a possible method but highlighted some complications related to its use. For instance, because of the presence of PCR inhibiting factors in environmental samples, a standard curve needs to be determined for each sample. Determining the percentage of "no marker" in a sample may be difficult for quantitative PCR (this assumes availability of a genus, species, or lineage marker with a high degree of specificity). However, our model and correction method will work even when relative percentages rather than absolute quantification of markers are estimated. For either a relative or absolute quantification approach, it is clear that an accurate measure of marker occurrence is critical for generating accurate estimates.

Assuming a valid method of quantification can be identified, there are still hurdles to overcome with regard to determining levels of marker occurrence. One question that needs to be addressed is whether marker prevalence is constant in space and time. If the prevalence is not constant in space, for instance, a new set of marker occurrence data may be needed for each geographical area; this would make the BST method time-consuming and cost ineffective (16). Yet if these problems can be overcome, we can make up for errors in measurement by applying the correction method outlined in this study. Furthermore, this method can compensate for any level of nonspecific occurrence or difference in marker prevalence. Thus, less effort could be spent finding markers with high specificity as long as the prevalence between host populations is accurately measured. Our correction method only applies to errors arising from marker prevalence and occurrence, not from variation due to mixing and density of the environment.

The distribution of sampling points (localized vs random over a wide area) can also affect estimates. Scenarios producing high estimation variance can be mitigated to some extent by distributing sample collections over a wide spatial domain as was demonstrated by scenarios 6 and 7, given unmixed environments. Quantification via presence/absence measurement performed well in the low-density environments, but only for equal volumes of host cells which may be unreasonable in practice.

The preceding discussion applies to the quantitative method of measurement, the main focus of our study. We also examined the presence/absence method of measurement, which is a likely format for discrete genetic marker BST. This method showed significant error in two cases: first, when there were unequal amounts of contamination among hosts but at least some contamination from each host in nearly every sample; second, when there were equal volumes of contamination among hosts but significant number of samples that lacked contamination from one or more hosts

(e.g., due to low density in marker-bearing contamination). Thus, presence/absence only performed better than the quantitative measure in dense environments that have equal proportions of contamination among hosts, an unlikely scenario.

In addition to our model parameters, there are other aspects of variation we did not examine. For instance, our model does not include time-varying rates of contamination input to the environment, which we expect from real-world sources. Transient events such as a rain storm or a disturbance in a riverbed may produce environments that are dense and rich in a particular marker that then disperses over time. We also assumed simplified survival dynamics for bacteria in an aqueous environment. We assumed that all cells of a particular host followed a constant decay rate once they entered the environment, but the decay rate is likely to be time-varying, and the survival characteristics of marker-bearing bacteria may differ from those of other bacteria from the same host. Sampling methods that require bacterial culture may also produce variance due to differential recovery rates. Another source of variation we excluded is the amount of marker produced per unit weight of fecal material. In the marker prevalence matrix, we model the amount of non-marker producing cells, but this refers only to the indicator bacteria (e.g., *Enterococcus*). The marker-per-unit-weight variation is similar to the variation we simulated with unequal marker prevalence and, thus, could have a similar effect on measurement error.

Our assumption is that the unmodeled sources of variation mentioned above would lower our confidence in quantitative measurements beyond that implied by modeled variation. This combined with the lack of a method for correcting for sparse, unmixed environments, and the uncertainty in deriving accurate levels of marker occurrence per total amount of contamination leads us to conclude that absolute quantification of fecal contamination may not be possible using discrete marker BST. To confirm or modify our conclusions we suggest that future research be directed to two areas: first, developing methods to accurately quantify a group of markers in a given sample; second, clearly defining marker prevalence across hosts, space, and time to assess the feasibility of "correcting" estimates from field samples. If the variance in measurement of marker prevalence is found to be low, research could then focus on determining whether real-world environments are suitable for use with the correction method introduced in this study.

Acknowledgments

This project was funded by USDA NRI contract 2002-35102-12374 and by the Agricultural Animal Health Program at the College of Veterinary Medicine, Washington State University, Pullman, WA. The authors would also like to express their gratitude for the support of the Carl M. Hansen Foundation.

Supporting Information Available

Detailed scenario simulation results, the model parameter summary, best-case scenario parameters, scenario marker prevalence matrices, and an example plot of percent coef-

ficient of variation vs time. This information is available free of charge via the Internet at <http://pubs.acs.org>.

Literature Cited

- (1) Bernhard, A. E.; Field, K. G. A PCR assay to discriminate human and ruminant feces on the basis of host differences in *Bacteroides-Prevotella* genes encoding 16S rRNA. *Appl. Environ. Microbiol.* **2000**, *66*, 4571–4574.
- (2) Bernhard, A. E.; Field, K. G. Identification of nonpoint sources of fecal pollution in coastal waters by using host-specific 16S ribosomal DNA genetic markers from fecal anaerobes. *Appl. Environ. Microbiol.* **2000**, *66*, 1587–1594.
- (3) Dick, L. K.; Bernhard, A. E.; Brodeur, T. J.; Santo Domingo, J. W.; Simpson, J. M.; Walters, S. P.; Field, K. G. Host distributions of uncultivated fecal *Bacteroidales* bacteria reveal genetic markers for fecal source identification. *Appl. Environ. Microbiol.* **2005**, *71*, 3184–3191.
- (4) Dick, L. K.; Field, K. G. Rapid estimation of numbers of fecal *Bacteroidetes* by use of a quantitative PCR assay for 16S rRNA genes. *Appl. Environ. Microbiol.* **2004**, *70*, 5695–5697.
- (5) Durán, A. E.; Muniesa, M.; Méndez, X.; Valero, F.; Lucena, F.; Jofre, J. Removal and inactivation of indicator bacteriophages in fresh waters. *J. Appl. Microbiol.* **2002**, *92*, 338–347.
- (6) Ferguson, C. M.; Coote, B. G.; Ashbolt, N. J.; Stevenson, I. M. Relationships between indicators, pathogens and water quality in an estuarine system. *Water Res.* **1996**, *30*, 2045–2054.
- (7) Gabutti, G.; De Donno, A.; Bagordo, F.; Montagna, M. T. Comparative survival of faecal and human contaminants and use of *Staphylococcus aureus* as an effective indicator of human pollution. *Mar. Pollut. Bull.* **2000**, *40*, 697–700.
- (8) Hamilton, M. J.; Yan, T.; Sadowsky, M. J. Development of goose- and duck-specific DNA markers to determine sources of *Escherichia coli* in waterways. *Appl. Environ. Microbiol.* **2006**, *72*, 4012–4019.
- (9) Scott, T. M.; Jenkins, T. M.; Lukasik, J.; Rose, J. B. Potential use of a host associated molecular marker in *Enterococcus faecium* as an index of human fecal pollution. *Environ. Sci. Technol.* **2005**, *39*, 283–287.
- (10) Shanks, O. C.; Nietch, C.; Simonich, M.; Younger, M.; Reynolds, D.; Field, K. G. Basin-wide analysis of the dynamics of fecal contamination and fecal source identification in Tillamook Bay, OR. *Appl. Environ. Microbiol.* **2006**, *72*, 5537–5546.
- (11) Shanks, O. C.; Santo Domingo, J. W.; Lamendella, R.; Kelty, C. A.; Graham, J. E. Competitive metagenomic DNA hybridization identifies host-specific microbial genetic markers in cow fecal samples. *Appl. Environ. Microbiol.* **2006**, *72*, 4054–4060.
- (12) Shanks, O. C.; Santo Domingo, J. W.; Lu, J.; Kelty, C. A.; Graham, J. E. Identification of bacterial DNA markers for the detection of human fecal pollution in water. *Appl. Environ. Microbiol.* **2007**, *73*, 2416–2422.
- (13) Sinton, L. W.; Hall, C. H.; Lynch, P. A.; Davies-Colley, R. J. Sunlight inactivation of fecal indicator bacteria and bacteriophages from waste stabilization pond effluent in fresh and saline waters. *Appl. Environ. Microbiol.* **2002**, *68*, 1122–1131.
- (14) Soule, M.; Kuhn, E.; Loge, F.; Gay, J.; Call, D. R. Using DNA microarrays to identify library-independent markers for bacterial source tracking. *Appl. Environ. Microbiol.* **2006**, *72*, 1843–1851.
- (15) Stewart, J. R.; Santo Domingo, J. W.; Wade, T. J. *Microbial Source Tracking*; Santo Domingo, J. W., Sadowsky, M. J., Eds.; ASM Press: Washington, DC, 2007.
- (16) Stoeckel, D. M.; Mathes, M. V.; Hyer, K. E.; Hagedorn, C.; Kator, H.; Lukasik, J.; O'Brien, T. L.; Fenger, T. W.; Samadpour, M.; Strickler, K. M.; Wiggins, B. A. Comparison of seven protocols to identify fecal contamination sources using *Escherichia coli*. *Environ. Sci. Technol.* **2004**, *38*, 6109–6117.

ES070943X