

## ARTICLE OPEN



# Monitoring the microbiome for food safety and quality using deep shotgun sequencing

Kristen L. Beck<sup>1,2,9</sup>✉, Niina Haiminen<sup>1,3,9</sup>, David Chambliss<sup>1,2</sup>, Stefan Edlund<sup>1,2</sup>, Mark Kunitomi<sup>1,2</sup>, B. Carol Huang<sup>1,4</sup>, Nguyet Kong<sup>1,4</sup>, Balasubramanian Ganesan<sup>1,5,6</sup>, Robert Baker<sup>1,5</sup>, Peter Markwell<sup>1,5</sup>, Ban Kawas<sup>1,2</sup>, Matthew Davis<sup>1,2</sup>, Robert J. Prill<sup>1,2</sup>, Harsha Krishnareddy<sup>1,2</sup>, Ed Seabolt<sup>1,2</sup>, Carl H. Marlowe<sup>1,7</sup>, Sophie Pierre<sup>1,8</sup>, André Quintanar<sup>1,8</sup>, Laxmi Parida<sup>1,3</sup>, Geraud Dubois<sup>1,2</sup>, James Kaufman<sup>1,2</sup> and Bart C. Weimer<sup>1,4</sup>✉

In this work, we hypothesized that shifts in the food microbiome can be used as an indicator of unexpected contaminants or environmental changes. To test this hypothesis, we sequenced the total RNA of 31 high protein powder (HPP) samples of poultry meal pet food ingredients. We developed a microbiome analysis pipeline employing a key eukaryotic matrix filtering step that improved microbe detection specificity to >99.96% during in silico validation. The pipeline identified 119 microbial genera per HPP sample on average with 65 genera present in all samples. The most abundant of these were *Bacteroides*, *Clostridium*, *Lactococcus*, *Aeromonas*, and *Citrobacter*. We also observed shifts in the microbial community corresponding to ingredient composition differences. When comparing culture-based results for *Salmonella* with total RNA sequencing, we found that *Salmonella* growth did not correlate with multiple sequence analyses. We conclude that microbiome sequencing is useful to characterize complex food microbial communities, while additional work is required for predicting specific species' viability from total RNA sequencing.

npj Science of Food (2021)5:3; <https://doi.org/10.1038/s41538-020-00083-y>

## INTRODUCTION

Sequencing the microbiome of food may reveal characteristics of the associated microbial content that culturing or targeted whole-genome sequencing (WGS) alone cannot. However, to meet the various needs of food safety and quality, next-generation sequencing (NGS), and analysis techniques require additional development<sup>1</sup> with specific consideration for accuracy, speed, and applicability across the supply chain<sup>2</sup>. Microbial communities and their characteristics have been studied in relation to flavor and quality in fermented foods<sup>3–5</sup>, agricultural processes in grape<sup>6</sup> and apple fruit<sup>7</sup>, and manufacturing processes and production batches in Cheddar cheese<sup>8</sup>. However, the advantage of using the microbiome specifically for food safety and quality has yet to be demonstrated.

Currently, food safety regulatory agencies including the Food and Drug Administration (FDA), Centers for Disease Control and Prevention (CDC), United States Department of Agriculture (USDA), and European Food Safety Authority (EFSA) are converging on the use of WGS for pathogen detection and outbreak investigation. Large scale WGS of food-associated bacteria was first initiated via the 100 K Pathogen Genome Project<sup>9</sup> with the goal of expanding the diversity of bacterial reference genomes—a crucial need for foodborne illness outbreak investigation, traceability, and microbiome studies<sup>10,11</sup>. However, since WGS relies on culturing a microbial isolate prior to sequencing, there are inherent biases and limitations in its ability to describe the microorganisms and their interactions in a food sample. Such information would be very valuable for food safety and quality applications.

High-throughput sequencing of the total DNA and total RNA are promising approaches to characterize microbial niches in their native state without introducing bias due to culturing<sup>12–14</sup>. In

addition, total RNA sequencing has the potential to provide evidence of live and biologically active components of the sample<sup>14,15</sup>. It also provides accurate microbial naming, relative microbial abundance, and better reproducibility than total DNA or amplicon sequencing<sup>14</sup>. Total RNA sequencing minimizes PCR amplification bias that occurs in single gene amplicon sequencing and overcomes the decreased detection sensitivity from using DNA sequencing in metagenomics<sup>14</sup>. Total RNA metatranscriptome sequencing, however, is yet to be examined in raw food ingredients as a method to provide a robust characterization of the microbial communities and the interacting population dynamics.

From a single sequenced food microbiome, numerous dimensions of the sample can be characterized that may yield important indicators of safety and quality. Using total DNA or RNA, evidence for the eukaryotic food matrix can be examined. In Haiminen et al.<sup>16</sup>, we quantitatively demonstrated the utility of metagenome sequencing to authenticate the composition of complex food matrices. In addition, from total DNA or RNA, one can observe signatures from commensal microbes, pathogenic microbes, and genetic information for functional potential (from DNA) or biologically active function (from RNA)<sup>14,15</sup>. Detecting active transcription from live microbes in food is very important to avoid spurious microbial observations that may instead be false positives due to quiescent DNA in the sample. The use of RNA in food analytics also offers the opportunity to examine the expression of metabolic processes that are related to antibiotic resistance<sup>17,18</sup>, virulence factors, or replication genes, among others. In addition, it has the potential to define viable microbes that are capable of replication in the food and even microorganisms that stop replicating but continue to produce metabolic activity that changes food quality and safety<sup>19–24</sup>.

<sup>1</sup>Consortium for Sequencing the Food Supply Chain, San Jose, CA, USA. <sup>2</sup>IBM Almaden Research Center, San Jose, CA, USA. <sup>3</sup>IBM T.J. Watson Research Center, Yorktown Heights, Ossining, NY, USA. <sup>4</sup>University of California Davis, School of Veterinary Medicine, 100 K Pathogen Genome Project, Davis, CA 95616, USA. <sup>5</sup>Mars Global Food Safety Center, Beijing, China. <sup>6</sup>Wisdom Health, A Division of Mars Petcare, Vancouver, WA, USA. <sup>7</sup>Bio-Rad Laboratories, Hercules, CA, USA. <sup>8</sup>Bio-Rad, Food Science Division, MArnes-La-Coquette, France. <sup>9</sup>These authors contributed equally: Kristen L. Beck, Niina Haiminen. ✉email: kbeck@us.ibm.com; bcweimer@ucdavis.edu

Microorganisms are sensitive to changes in temperature, salinity, pH, oxygen content, and many other physicochemical factors that alter their ability to grow, persist, and cause disease. They exist in dynamic communities that change in response to environmental perturbation—just as the gut microbiome shifts in response to diet<sup>25–28</sup>. Shifts in microbiome composition or activity can be leveraged in the application of microbiome characterization to monitor the food supply chain. For example, Noyes et al. followed the microbiome of cattle from the feedlot to the food packaging, concluding that the microbial community and antibiotic resistance characteristics change based on the processing stage<sup>17,18,29</sup>. We hypothesize that observable shifts in microbial communities of food can serve as an indicator of food quality and safety.

In this work, we examined 31 high protein powder samples (HPP; derived from poultry meal). HPP are commonly used raw materials in pet foods. They are subject to microbial growth prior to preparation and continued survival in powder form<sup>30</sup>. We subjected the HPP samples to deep total RNA sequencing with ~300 million reads per sample. In order to process the 31 samples collected over ~1.5 years from two suppliers at a single location, we defined and calibrated the appropriate methods—from sample preparation to bioinformatics analysis—needed to taxonomically identify the community members present and to detect key features of microbial growth. First, we removed the HPP's food matrix RNA content as eukaryotic background with an important bioinformatic filtering step designed specifically for food analysis. The remaining sequences were used for relative quantification of microbiome members and for identifying shifts based on food matrix content, production source, and *Salmonella* culturability. This work demonstrates that total RNA sequencing is a robust approach for monitoring the food microbiome for use in food safety and quality applications, while additional work is required for predicting pathogen viability.

## RESULTS

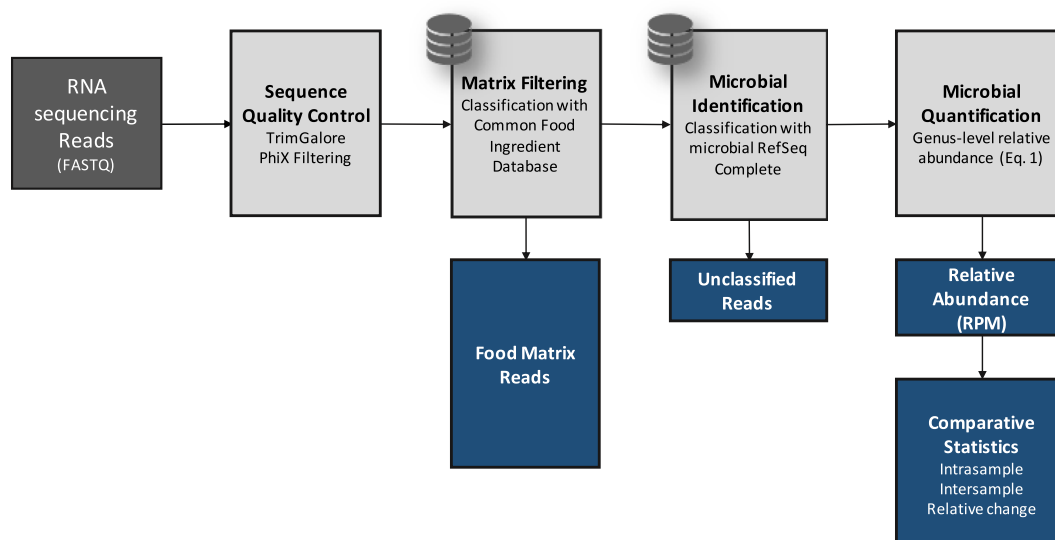
### Evaluation of microbial identification capability in total RNA and DNA sequencing

Microbial identification in microbiomes often leverages shotgun DNA sequencing; however, total RNA sequencing can provide additional information about viable bacterial activity in a community via transcriptional activity. Since using total RNA to

study food microbiomes is novel, each step of the analysis workflow (Fig. 1) was carefully designed and scrutinized for accuracy. For all analyses done in this study, we report relative abundance in reads per million (RPM) (Eq. 1) as recommended by Gloor et al.<sup>31,32</sup> and apply the conservative threshold of RPM > 0.1 to indicate presence as indicated by Langelier et al. and Illot et al.<sup>33,34</sup>. Numerically, this threshold translates to ~30 reads per genus per sample considering a sequencing depth of ~300 million reads per sample (see the section “Microbial identification”). First, we examined the effectiveness of RNA for taxonomic identification and relative quantification of microbes in the presence of food matrix reads. We observed that RNA sequencing results correlated ( $R^2 = 0.93$ ) with the genus relative quantification provided by DNA sequencing (Supplementary Fig. S1). RNA sequencing also detected more genera demonstrated by a higher  $\alpha$ -diversity than the use of DNA (Supplementary Fig. S2). In addition, from the same starting material, total RNA sequencing resulted in 2.4-fold more reads classified to microbial genera compared to total DNA sequencing (after normalizing for sequencing depth). This increase is substantial as microbial reads are such a small fraction of the total sequenced reads. Considering these results, we further examined the microbial content from total RNA extracted from 31 HPP samples (Supplementary Table 1) that resulted in an average of ~300 million paired-end 150 bp sequencing reads per sample in this study.

### Evaluation and application of in silico filtering of eukaryotic food matrix reads

Sequenced reads from the eukaryotic host or food matrix may lead to false positives for microbial identification in microbiome studies<sup>35</sup>. This may occur partly due to reads originating from low complexity regions of eukaryotic genomes, e.g., telomeric and centromeric repeats, being misclassified as spurious microbial hits<sup>36</sup>. In total DNA or RNA sequencing of clinical or animal or even plant microbiomes, eukaryotic content may often comprise >90% of the total sequencing reads. This presents an important bioinformatic challenge that we addressed by filtering matrix content using a custom-built reference database of 31 common food ingredient and contaminant genomes (Supplementary Table 2) using the *k*-mer classification tool Kraken<sup>37</sup>. This step allows for rapidly classifying all sequenced reads (~300 million reads for each of



**Fig. 1 Bioinformatic pipeline schematic for processing microbiome samples in the presence of matrix content.** Description of the bioinformatic steps (light gray) applied to high protein powder metatranscriptome samples (dark gray). Black arrows indicate data flow and blue boxes describe outputs from the pipeline.

31 samples) as matrix or non-matrix. The matrix filtering process yielded an estimate of the total percent matrix content for a sample. See our work in Haiminen et al.<sup>16</sup> on quantifying the eukaryotic food matrix components with further precision.

To validate the matrix filtering step, we constructed in silico mock food microbiomes with a high proportion of complex food matrix content and low microbial content (Supplementary Table 3). We then computed the true positive, false positive, and false-negative rates of observed microbial genera and sequenced reads (Table 1). False-positive viral, archaeal, and eukaryotic microbial genera (as well as bacteria) were observed without matrix filtering, although bacteria were the only microbes included in the simulated mixtures. Introducing a matrix filtering step to the pipeline improved read classification specificity to >99.96% (from 78 to 93% without filtering) in both simulated food mixtures while maintaining zero false negatives. With this level of demonstrated accuracy, we used bioinformatic matrix filtering prior to further microbiome analysis.

### HPP microbiome ecology

After filtering eukaryotic matrix sequences, we applied the remaining steps in the bioinformatic workflow (Fig. 1) to examine the shift in the HPP microbiome membership and to quantify the relative abundance of microbes at the genus level. Genus is the first informative taxonomic rank for food pathogen identification that can be considered accurate given the current incompleteness of reference databases<sup>11,38–41</sup> and was therefore used in subsequent analyses. Overall, between 98 and 195 microbial genera (avg. 119) were identified (RPM > 0.1) per HPP sample (Supplementary Table 4). When analyzing  $\alpha$ -diversity i.e., the number of microbes detected per sample, inter-sample comparisons may become skewed unless a common number of reads is considered since deeper sequenced samples may contain more observed genera merely due to a greater sampling depth<sup>42,43</sup>. Thus, we utilized bioinformatic rarefaction i.e., subsampling analysis to showcase how microbial diversity was altered by sequencing depth. Examination of  $\alpha$ -diversity across a range of in silico subsampled sequencing depths showed that the community diversity varied across samples (Fig. 2a). One sample (MFMB-04) had 1.7 times more genera (195) than the average across

other samples (avg. 116, range 98–143) and exhibited higher  $\alpha$ -diversity than any other sample at each in silico sampled sequencing depth (Fig. 2a). Rarefaction analysis further demonstrated that when considering fewer than ~67 million sequenced reads, the observable microbial population was not saturated (median elbow calculated as indicated in Satopää, et al.<sup>44</sup>). This observation suggests that deeper sequencing or more selective sequencing of the HPP microbiomes will reveal more microbial diversity.

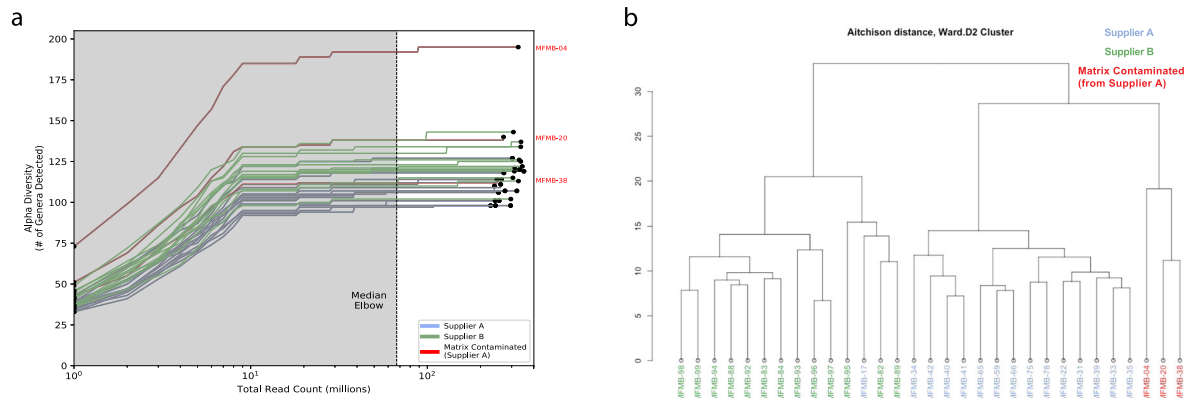
Notably, between 2 and 4% (~5,000,000–14,000,000) of reads per sample remained unclassified as either eukaryotic matrix or microbe (Supplementary Fig. S3). However, the unclassified reads exhibited a GC (guanine plus cytosine) distribution similar to reads classified as microbial (Supplementary Fig. S4), indicating these reads may represent microbial content that is absent or sufficiently divergent from existing references.

We calculated  $\beta$ -diversity to study inter-sample microbiome differences and to identify any potential outliers among the sample collection. The Aitchison distances<sup>45</sup> of microbial relative abundances were calculated between samples (as recommended for compositional microbiome data<sup>31,32</sup>), and the samples were hierarchically clustered based on the resulting distances (Fig. 2b). The two primary clades were mostly defined by the supplier (except for MFMB-17). Samples were collected over several months with Supplier A contributing three batches over time and Supplier B contributing one shipment batch (Supplementary Table 1); despite time point differences, the microbiome composition still clusters into separate clades by the supplier. In Haiminen et al.<sup>16</sup>, we reported that three of the HPP samples contained unexpected eukaryotic species. We hypothesized that the presence of these contaminating matrix components (beef identifiable as *Bos taurus* and pork identifiable as *Sus scrofa*) would alter the microbiome as compared to chicken (identifiable as *Gallus gallus*) alone. Clustering HPP samples using their microbiome membership led to a distinctly different group of the matrix-contaminated samples, supporting this hypothesis (Fig. 2b). These observations indicate that samples can be discriminated based on their microbiome content for originating source and supplier, which is necessary for source tracking potential hazards in food.

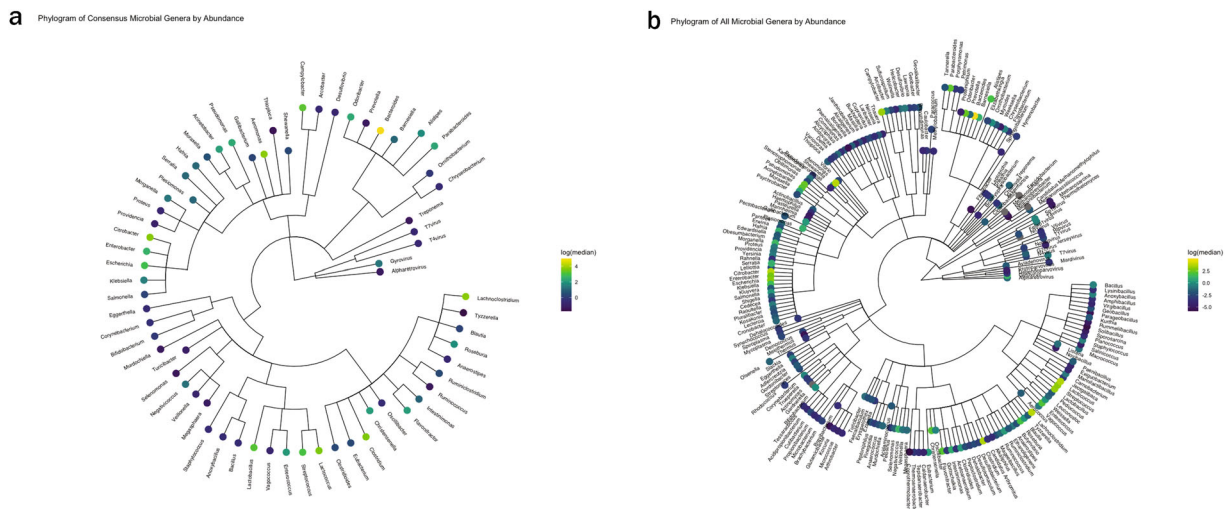
**Table 1.** Accuracy of microbial identification using two in silico constructed simulated food mixtures.

	Mixture 1, no MF # GENERA	Mixture 1, MF # GENERA	Mixture 1, no MF # READS	Mixture 1, MF # READS	Mixture 2, no MF # GENERA	Mixture 2, MF # GENERA	Mixture 2, no MF # READS	Mixture 2, MF # READS
Bacteria in mixture (expected Content)	14	14	15,000	15,000	14	14	15,000	15,000
Observed bacteria	34	18	13,700	13,517	33	15	13,999	13,551
Observed viruses	9	0	563	0	4	0	328	0
Observed archaea	1	0	1	0	1	0	3	0
Observed Eukaryota	4	0	104	0	4	0	799	0
Total observed (TO)	48	18	14,368	13,517	42	15	15,129	13,551
True positives (TP)	14	14	13,571	13,511	14	14	13,623	13,548
TP as % of TO	29%	78%	94.45%	99.96%	33%	93%	90.05%	99.98%
False positives (FP)	34	4	797	6	28	1	1506	3
FP as % of TO	71%	22%	5.55%	0.04%	67%	7%	9.95%	0.02%
FP removed with MF	–	30	–	791	–	27	–	1503
% FP removed with MF	–	88.2%	–	99.2%	–	96.4%	–	99.8%

The Simulated Food Mixtures (Mixture 1 and Mixture 2, see Supplementary Table 3) contain food matrix and microbial sequences. Microbial identification results are shown without matrix filtering (no MF) and with matrix filtering (MF). The number of observed genera (# GENERA) and observed genus-assigned reads (# READS) are shown for each category and summarized as the total observed (TO) counts. True positive (TP) and false-positive (FP) counts and fractions of TO are shown. The last two rows show the counts and percentages of false positives removed with matrix filtering.



**Fig. 2 Ecological metrics of microbiome community.** **a** Alpha diversity (number of genera) for all ( $n=31$ ) high protein powder metatranscriptomes is compared to the total number of sequenced reads for a range of in silico subsampled sequencing depths. The dashed vertical line indicates the medial elbow (at  $\sim 67$  million reads). **b** Hierarchical clustering of Aitchison distance values of poultry meal samples based on microbial composition. Samples were received from Supplier A (blue and red) and Supplier B (green). Matrix-contaminated samples are additionally marked in red.



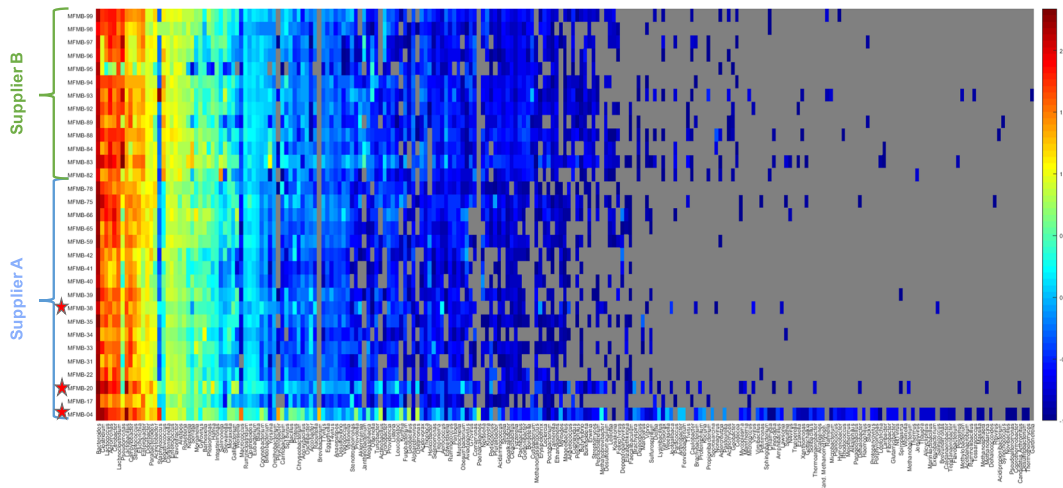
**Fig. 3 Microbial genera detected in high protein powder samples.** **a** Phylogram of the 65 microbial genera present in all samples with RPM  $> 0.1$ . **b** Phylogram of microbes observed in any sample. Log of the median RPM value across samples is indicated. Gray indicating a median RPM value of 0.

### Comparative analysis of HPP microbiome membership and composition

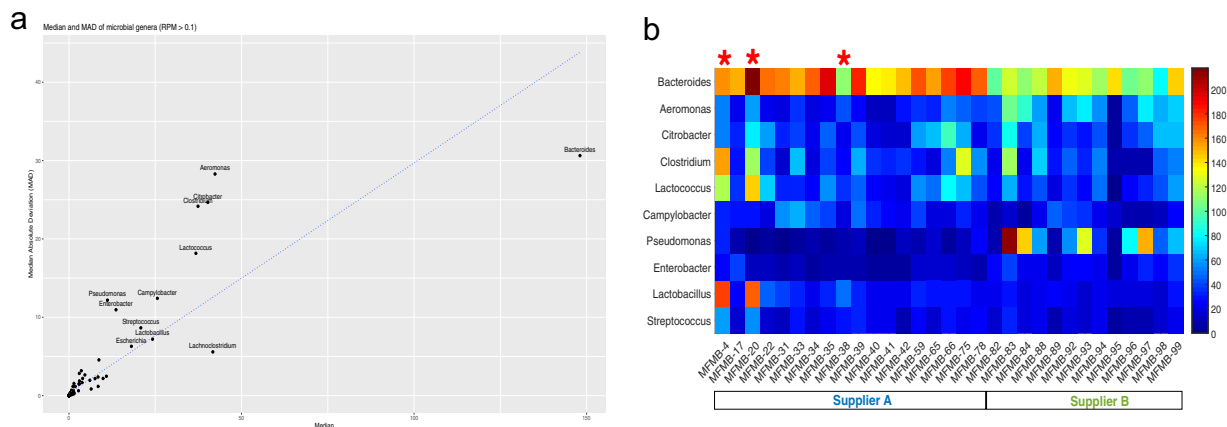
We identified 65 genera present in all HPP samples (Fig. 3a), whose combined abundance accounted for between 88 and 99% of the total abundances of detected genera per sample. *Bacteroides*, *Clostridium*, *Lactococcus*, *Aeromonas*, and *Citrobacter* were the five most abundant of these microbial genera. The identified microbial genera also included viruses, the most abundant of which was *Gyrovirus* ( $<10$  RPM per sample). *Gyrovirus* represents a genus of non-enveloped DNA viruses responsible for chicken anemia which is ubiquitous in poultry. While there were only 65 microbial genera identified in all 31 HPP samples, the  $\alpha$ -diversity per sample was on average twofold greater as previously indicated.

Beyond the collection of 65 microbes observed in all samples, there were an additional 164 microbes present in various HPP samples. Together, we identified a total of 229 genera among the 31 HPP samples tested (Figs 3b and 4, Supplementary Table 4). In order to identify genera that were most variable between samples, we computed the median absolute deviation (MAD)<sup>46</sup> using the

normalized relative abundance of each microbe (Fig. 5a). The abundance of *Bacteroides* was the most variable among samples (median = 148.1 RPM, MAD = 30.6) and showed increased abundance in samples from Supplier A (excluding samples with known host contamination) compared to Supplier B (Benjamini–Hochberg adjusted  $P < 0.00005$ ). In total, there were 55 genera with significant differences in abundance between Supplier A and Supplier B (adjusted  $P < 0.01$ ). Of the ten most variable genera based on MAD, *Aeromonas*, *Enterobacter*, *Pseudomonas*, and *Lactobacillus* also had significant differences between Supplier A and B (adjusted  $P < 0.01$  with their relative abundances shown in Fig. 5b). In addition, *Clostridium* (median = 37.4 RPM, MAD = 24.2), *Lactococcus* (median = 36.8 RPM, MAD = 18.2), and *Lactobacillus* (median = 24.2, MAD = 7.2) were also highly variable and threefold to fourfold more abundant in samples MFMB-04 and MFMB-20 compared to other samples (Fig. 5b). *Pseudomonas* (median = 11.1 RPM, MAD = 12.2) was markedly more abundant in MFMB-83 than any other sample (Fig. 5b). These genera highlight variability between microbiomes based on supplier origin or food source and may provide insights into other dissimilarities in these samples.



**Fig. 4 High protein powder (HPP) microbial composition and relative abundance per sample.** Heatmap ( $\log_{10}$ -scale) of HPP microbial composition and relative abundance (RPM) where absence ( $\text{RPM} < 0.1$ ) is indicated in gray. Genera are ordered by summed abundance across samples. Samples were received from Supplier A (blue) and Supplier B (green). Red stars indicate matrix-contaminated samples (from Supplier A).



**Fig. 5 Variability of microbial genera relative abundance.** **a** All identified microbial general are plotted with median value and median absolute deviation (MAD) of RPM abundance. Genera with  $\text{MAD} > 5$  are labeled with the genus name and a linear fit is indicated by a blue dotted line. **b** Heatmap ( $\log_{10}$ -scale) of ten microbial genera with the largest median absolute deviation (MAD) across samples. Genera are ordered by decreasing MAD from top to bottom. Samples were received from Supplier A (blue) and Supplier B (green). Red stars indicate matrix-contaminated samples (from Supplier A).

### Microbiome shifts in response to changes in food matrix composition

We tested the hypothesis that the microbiome composition will shift in response to changes in the food matrix and can be a unique signal to indicate contamination or adulteration. In 28 of the 31 HPP samples, >99% of the matrix reads were determined in our related work<sup>16</sup> to originate from poultry (*Gallus gallus*), which was the only ingredient expected based on ingredient specifications. However, three samples had higher pork and beef content compared to all other HPP samples: MFMB-04 (7.74% pork, 8.99% beef), MFMB-20 (0.53% pork, 1.00% beef), and MFMB-38 (0.92% pork, 0.29% beef) compared to the highest pork (0.01%) and beef (0.00%) content among the other 28 HPP samples (Supplementary Data by Haiminen et al.<sup>16</sup>). The microbiomes of these matrix-contaminated samples, each coming from Supplier A, also clustered into a separate sub-cluster (Fig. 2b). This demonstrated that a shift in the food matrix composition was associated with an observable shift in the food microbiome.

We further computed pairwise Spearman's correlation between all samples, using the RPM vectors for the 229 detected genera as input (Supplementary Fig. S5). Here, we exclude MFMB-04, MFMB-20, and MFMB-38 from the group "Supplier A samples" and consider them as a separate matrix-contaminated group. The mean correlation between Supplier A samples was 0.946, while the mean correlation between Supplier B samples was 0.816. The mean correlation between Supplier A and Supplier B samples was 0.805, lower than either within-group correlation. Contrasted with this, the mean correlation between MFMB-04 and Supplier A samples was 0.656, analogously for MFMB-20 the mean correlation was 0.866, and for MFMB-38 it was 0.885. The increasing correlation values correspond with decreasing percentages of cattle and pork reads in the matrix-contaminated samples (16.7% in MFMB-04, 1.5% in MFMB-20, and 1.2% in MFMB-38), indicating a trend toward the microbial baseline with decreasing matrix contamination.

MFMB-04 and MFMB-20 had the highest percentage of microbial reads compared to other samples (Supplementary Fig.

S3). They also exhibited an increase in *Lactococcus*, *Lactobacillus*, and *Streptococcus* relative abundances compared to other samples (Fig. 5b), also reflected at respective higher taxonomic levels above genus (Supplementary Fig. S6).

There were 53 genera identified uniquely in MFMB-04 and/or MFMB-20 i.e., RPM values above the aforementioned threshold in these samples but not present in any other sample. (MFMB-38 had a very low microbial load and contributed no uniquely identified genera above the abundance threshold.) MFMB-04 contained 44 unique genera (Fig. 4) with the most abundant being *Macrocooccus* (35.8 RPM), *Psychrobacter* (23.8 RPM), and *Brevibacterium* (18.1 RPM). In addition, *Paenalcaligenes* was present only in MFMB-04 and MFMB-20 with an RPM of 6.4 and 0.3, respectively, compared to a median RPM of 0.004 among other samples. Notable differences in the matrix-contaminated samples' unique microbial community membership compared to other samples may provide microbial indicators associated with unanticipated pork or beef presence.

### Genus-level identification of foodborne microbes

We evaluated the ability of total RNA sequencing to identify genera of commonly known foodborne pathogens within the microbiome. We focused on fourteen pathogen-containing genera including *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*, *Cronobacter*, *Escherichia*, *Helicobacter*, *Listeria*, *Salmonella*, *Shigella*, *Staphylococcus*, *Vibrio*, and *Yersinia* that were found to be present in the HPP samples with varying relative abundances. Of these genera, *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*, *Escherichia*, *Salmonella*, and *Staphylococcus* were detected in every HPP with median abundance values between 0.58–48.31 RPM (Fig. 6a). This indicated that a baseline fraction of reads can be attributed to foodborne microbes when using NGS. Of those genera appearing in all samples, there was observed sample-to-sample variation in their abundance with some genera exhibiting longer tails of high abundance, e.g., *Staphylococcus* and *Salmonella*, whereas others exhibit very low abundance barely above the threshold of detection, e.g., *Bacillus* and *Yersinia* (Fig. 6a). None of the pathogen-containing genera were consistent with higher relative abundances due to differences in food matrix composition. *Bacillus* and *Corynebacterium* exhibited slightly higher relative abundances in sample MFMB-04 which contained 7.7% pork and 9.0% beef (Fig. 6b). Yet while

MFMB-04 contained higher cumulative levels of these foodborne microbes, the next highest sample was MFMB-93 which was not associated with altered matrix composition, and both MFMB-04 and MFMB-93 contained higher levels of *Staphylococcus* (Fig. 6b). Thus, matrix composition alone did not explain variations of these pathogen-containing genera.

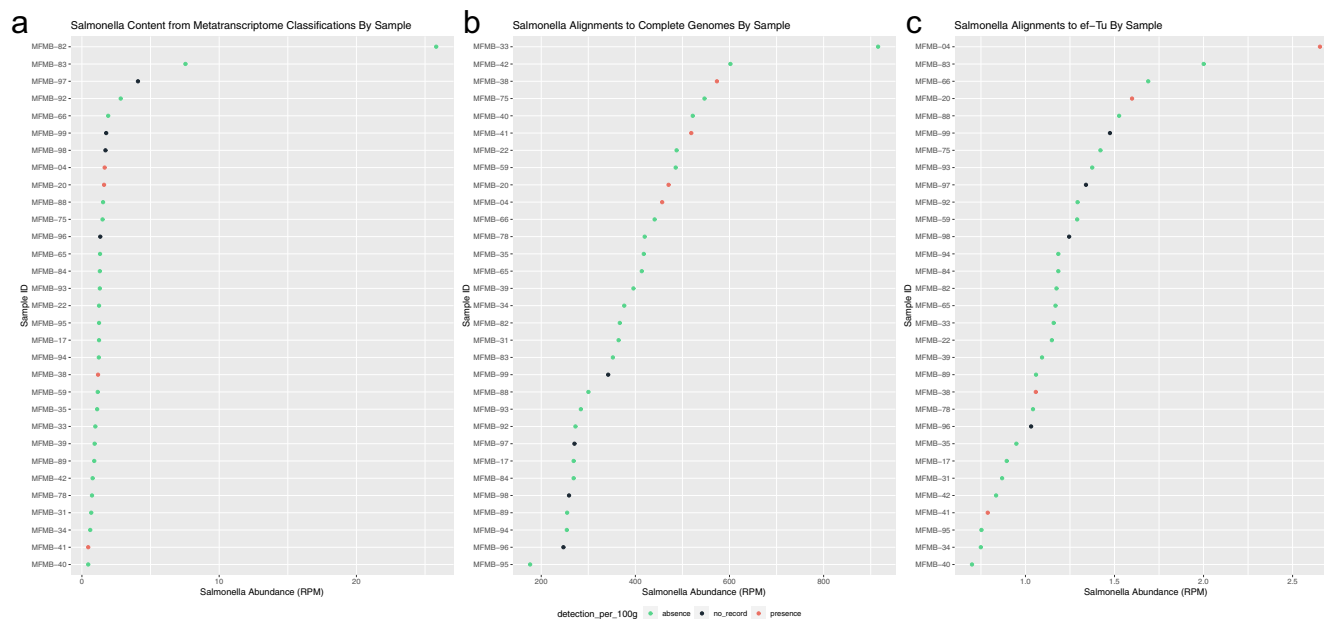
Interestingly, low to moderate levels of *Salmonella* were detected within all 31 HPP microbiomes (Fig. 6a). The presence of *Salmonella* in HPP is expected but the viability of *Salmonella* is an important indicator of safety and quality. Thus, we further sought to delineate *Salmonella* growth capability within these microbiomes by comparing culturability with multiple established bioinformatic NGS methods for *Salmonella* relative abundances in the samples.

### Assessment of *Salmonella* culturability and total RNA sequencing

Total RNA sequencing of food microbiomes has the potential to provide additional sensitivity beyond standard culture-based food safety testing to confirm or reject the presence of potentially pathogenic microbes. In all of the examined HPP samples, some portion of the sequenced reads were classified as belonging to pathogen-containing genera (Fig. 6); however, the presence of RNA transcripts does not necessarily indicate the current growth of the organism itself. We further inspected one pathogen of interest, *Salmonella*, to determine the congruence between sequencing-based and culturability results. Of the 31 samples examined with total RNA sequencing, *Salmonella* culture testing was applied to 27 samples, of which four were culture-positive. Surprisingly, *Salmonella* culture-positive samples were not among those with the highest relative abundance of *Salmonella* from sequencing (Fig. 7a). When ranking the samples by decreasing *Salmonella* abundance, the culture-positive samples were not enriched for higher ranks ( $P = 0.86$  from Wilcoxon rank-sum test indicating that the distributions are not significantly different, Table 2). To confirm that the microbiome analysis pipeline did not miss *Salmonella* reads present, we completed two orthogonal analyses on the same dataset used in the microbial identification step. The reference genomes relevant to these additional analyses were publicly available and closed high-quality genomes available from the sources indicated below.



**Fig. 6** Relative abundance for fourteen pathogen-containing genera. **a** Relative abundance distribution of genera with high relevance to food safety and quality from high protein powder (HPP) total RNA sequenced microbiomes. The width of the violin plot indicates the density of samples with relative abundance at that value. Observation threshold of RPM = 0.1 is indicated with the horizontal black line. **b** The relative abundances of those same genera are shown across samples of HPP total RNA sequenced samples.



**Fig. 7** *Salmonella* culture-positive status vs. high-throughput sequencing read abundance. Read abundance (RPM) shown from **a** *k*-mer classification to NCBI Microbial RefSeq Complete, **b** alignments to 1447 *Salmonella* genomes, and **c** alignments to 4846 *ef-Tu* gene sequences. *Salmonella* presence (red) indicates culture-positive result, absence (green) indicates culture-negative result, and no record (black) indicates samples for which no culture test was completed.

**Table 2.** *Salmonella* analyses.

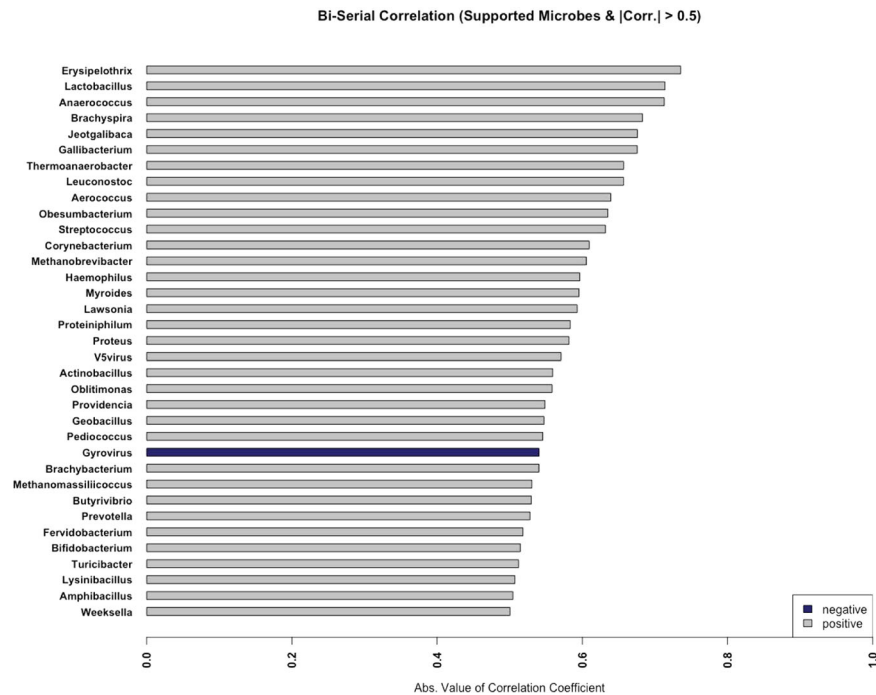
Salmonella-positive sample	<i>k</i> -mer classification	Whole-genome alignment	<i>ef-Tu</i> alignment
MFMB-04	8th	10th	1st
MFMB-20	9th	9th	4th
MFMB-38	20th	3rd	21st
MFMB-41	30th	6th	28th
Rank-sum test	$P = 0.86$	$P = 0.06$	$P = 0.56$
$P$ value			

The ranks for *Salmonella*-positive samples and the associated  $P$  values from Wilcoxon rank-sum test are shown for high-throughput sequencing read abundance (in RPM) for multiple analyses: *k*-mer classification to NCBI Microbial RefSeq Complete (left), alignments to 1447 *Salmonella* genomes (middle), and alignments to 4846 *ef-Tu* gene sequences (right). The corresponding *Salmonella* relative abundances are shown in Fig. 7a–c.

First, for targeted analysis, we aligned the sequenced reads using a different tool, Bowtie2<sup>47</sup>, to an augmented *Salmonella*-only reference database. This reference was comprised of the 264 *Salmonella* genomes extracted from NCBI RefSeq Complete (used in our previous microbial identification step) as well as an additional 1183 public *Salmonella* genomes which represent global diversity within the genus<sup>48</sup>. The number of reads that aligned to the *Salmonella*-only reference was on average 370-fold higher than identified as *Salmonella* by Kraken using the multi-microbe NCBI RefSeq Complete. In this additional analysis, the culture-positive samples had overall higher ranks compared to culture-negative samples ( $P = 0.06$ , Table 2), indicating that additional *Salmonella* genomic data in the reference significantly improved discriminatory identification power. *Salmonella* culture-positive samples were still not the most abundant (Fig. 7b), but with an enriched database, sequencing positioned all four culturable samples within the top ten rankings.

The second additional analysis examined the alignment of the reads to a specific gene required<sup>49</sup> for replication and protein production in actively dividing *Salmonella*—elongation factor Tu (*ef-Tu*). This was done by aligning the reads to 4846 gene sequences for *ef-Tu* extracted for a larger corpus of *Salmonella* genomes from the Functional Genomics Platform (formerly OMXWare<sup>50</sup>). The relative abundances of this transcript in culture-positive samples were still comparable to culture-negative samples (Fig. 7c). Culture-positive samples did not exhibit higher ranks compared to culture-negative samples ( $P = 0.56$ , Table 2), indicating that *ef-Tu* relative abundance alone was not sufficient to improve the lack of concordance in culturability vs sequencing. These two orthogonal analyses demonstrated that results from carefully developed culture-based testing and those from current high-throughput sequencing technologies, whether assessed at overall reads aligned or specific gene abundances, were not conclusively in agreement when detecting active *Salmonella* in food samples (Fig. 7 and Table 2). However, the use of a reference database enriched in whole-genome sequences of the specific organism of interest was found appropriate for food safety applications.

Since microbes compete for available resources within an environmental niche and therefore impact one another<sup>51</sup>, we investigated *Salmonella* culture results in conjunction with co-occurrence patterns of other microbes in the total RNA sequencing data (Fig. 8). Point-biserial correlation coefficients ( $r_{pb}$ ) were calculated between *Salmonella* culturability results (presence or absence which were available for 27 of the 31 samples) and microbiome relative abundance. We observed 31 genera that positively correlated and with *Salmonella* presence ( $r_{pb} > 0.5$ ). *Erysipelothrix*, *Lactobacillus*, *Anaerococcus*, *Brachyspira*, and *Jeotgallibaca* exhibited the largest positive correlations. *Gyrovirus* was negatively correlated with *Salmonella* growth ( $r_{pb} = -0.54$ ). In three of the four *Salmonella*-positive samples (MFMB-04, MFMB-20, and MFMB-38), food matrix contamination was also observed (Supplementary Data in Haiminen et al.<sup>16</sup>). The concurrency of *Salmonella* growth and matrix contamination was affirmed by the microbial co-occurrence (specifically *Erysipelothrix*, *Brachyspira*, and *Gyrovirus*). This highlights the complex dynamic and



**Fig. 8** *Salmonella* status correlations with genus relative abundances. Only those genera with the absolute value of the correlation coefficient  $>0.5$  are shown. Positive and negative correlations are indicated in gray and blue, respectively.

community co-dependency of food microbiomes, yet shows that multiple dimensions of the data (microbiome composition, culture-based methods, and microbial load) will signal anomalies from typical samples when there is an issue in the supply chain.

## DISCUSSION

Accurate and appropriate tests for detecting potential hazards in the food supply chain are key to ensuring consumer safety and food quality. Monitoring and regular testing of raw ingredients can reveal fluctuations within the supply chain that may be an indicator of an ingredient's quality or of a potential hazard. Such quality is assessed by standardized tests for the chemical and microbial composition to meet legal requirements and specifications from government agencies throughout the world. For raw materials or finished products to meet these bounds of safety and quality, their composition must usually have a low microbiological load (except in fermented foods) and be chemically identical in macro-components such as carbohydrate, protein, and fat. Methods in this space must avoid false-negative results that could endanger consumers, while also minimizing false positives which could lead to unnecessary recalls and food loss.

Existing microbial detection technologies used in food safety today such as pulse-field gel electrophoresis (PFGE) and WGS require microbial isolation. This provides biased outcomes as it removes microbes from their native environment where other biotic members also subsist and select microbes by culturability alone. Amplicon sequencing, while a low-cost alternative to metagenome or metatranscriptome sequencing for bacteria, also imparts PCR amplification bias and reduces detection sensitivity due to reliance on a single gene (16S ribosomal RNA)<sup>14,52,53</sup>. We, therefore, investigated the utility of total RNA sequencing of food microbiomes and demonstrated that from this single test, we are able to yield several pertinent results about food safety and quality.

For this evaluation, we developed a pipeline to characterize the microbiome of typical food ingredient samples and to detect potentially hazardous outliers. Special considerations for food

samples were made as computational pipelines for human or other microbiome analyses are not sufficient for applications in food safety without modification. In food, the eukaryotic matrix needs to be confirmed, maybe mixed, and, as we and others have shown, affects the identification accuracy of microbes that are present<sup>35,36</sup>. By filtering food matrix sequence data properly, we avoid incorrect microbial identification and characterization of the microbiome<sup>36</sup> while also increasing the computational efficiency for downstream processing. The addition of this filtering step in the pipeline removed  $\sim 90\%$  of false-positive genera and provided results at 99.96% specificity when evaluating simulated mixtures of food matrix and microbes (Table 1).

Through the analysis of 31 HPP total RNA sequencing samples, we demonstrated the pipeline's ability to characterize food microbiomes and indicate outliers. In this sample collection, we identified a core catalog of 65 microbial genera found in all samples where *Bacteroides*, *Clostridium*, and *Lactococcus* were the most abundant (Supplementary Table 4). We also demonstrated that in these food microbiomes the overall diversity was twofold greater than the core microbe set. Fluctuations in the microbiome can indicate important differences between samples as observed here, as well as in the literature for grape berry<sup>6</sup> and apple fruit microbiomes (pertaining to organic versus conventional farming)<sup>7</sup> or indicate inherent variability between production batches or suppliers as observed here and during cheddar cheese manufacturing<sup>8</sup>. Specifically, we observed a shift in the microbial composition (Fig. 2b) and the microbial load (Supplementary Fig. S3) in HPP samples (derived from poultry meal) where unexpected pork and beef were observed. Matrix-contaminated samples were marked by increased relative abundances of specific microbes including *Lactococcus*, *Lactobacillus*, and *Streptococcus* (Fig. 5b). This work shows that the microbiome shifts with observed food matrix contamination from sources with similar macronutrient content and thus, the microbiome alone is a likely signal of compositional change in food.

Beyond shifts in the microbiome, we focused on a set of well-defined foodborne-pathogen-containing genera and explored their relative abundances observed from total RNA sequencing.



Of these genera, *Aeromonas*, *Bacillus*, *Campylobacter*, *Clostridium*, *Corynebacterium*, *Escherichia*, *Salmonella*, and *Staphylococcus* were detected in every HPP sample. This highlights that when using NGS there may be an observable baseline of sequences assigned to potentially pathogenic microbes. For this ingredient type, this result lends a range of normalcy of relative abundance generated by NGS. Further work is needed to establish a definitive and quantitative range of typical variation in samples of a particular food source and the degree of an anomaly for a new sample or genus abundance. However, preliminary studies of this nature can inform the development of guidelines when working with increasingly sensitive shotgun metagenomic or metatranscriptomic analysis.

Furthermore, sequenced DNA or RNA alone does not imply microbial viability. Therefore, we investigated the relatedness of culture-based tests and total RNA sequencing for the pathogenic bacterium *Salmonella* in the HPP samples. As has been reported for human gut<sup>54</sup> and deep sea<sup>55</sup> microbiomes, we also did not detect a correlation between *Salmonella* read abundance and culturability (Fig. 7 and Table 2). Sequence reads matching *Salmonella* references were observed for all samples (both culture-positive and culture-negative) as determined by multiple analysis techniques: microbiome classification, alignment to *Salmonella* genomes, and targeted growth gene analysis. When ranking the HPP samples based on *Salmonella* abundance from whole-genome alignments, the culture-positive samples were enriched for higher ranks ( $P = 0.06$ ). However, the culture-positive samples were still intermixed in ranking with culture-negative samples. This indicated that there was no clear minimum threshold of sequence data as evidence for culturability and that this analysis alone is not predictive of pathogen growth. One possible reason for this is that the culture-positive variant of *Salmonella* is missing from existing reference data sets. Potentially, *Salmonella* attained a nonculturable state wherein it was detected by sequencing techniques yet remained nonculturable from the HPP sources. Successful isolation of total RNA and DNA and gene expression analysis from experimentally known nonculturable bacteria has been demonstrated by Ganesan et al. in multiple studies in other genera<sup>19,22</sup>. The physiological state should thus be taken under consideration when benchmarking sequencing technologies in comparison with culture-based methods. Thus, total RNA sequencing of food samples may identify shifts that standard food testing does not, but the incongruity between sequencing read data and culture-based results highlights the need to perform more benchmarking in food microbiome analysis for pathogen detection.

The characterization of HPP food microbiomes leveraged current accepted public reference databases, yet it is known that these databases are still inadequate<sup>1,2,11,56,57</sup>. Furthermore, when considering congruence between *Salmonella* culturability and NGS read mapping techniques, the genetic breadth and depth of multi-genome reference sequences are essential. For example, focusing on *ef-Tu*, a known marker gene for *Salmonella* growth was not sufficient to mirror the viability of in vitro culture tests. This highlights the limitations of single-gene approaches for identification. When the sequenced reads were examined in the context of an augmented reference collection of *Salmonella* genomes, we observed improved ranking and read mapping rate for culture-positive samples (yet we did not achieve complete concordance). This improvement underlined the increased analytical robustness yielded from a multi-genome reference. We also recognize that the read mapping rate may be exaggerated as reads from non-*Salmonella* genomes could map to *Salmonella* in the absence of any other reference genomes. Overall for robust analysis and applicability to food safety and quality, microbial references must be expanded to include more genetically diverse representatives of pathogenic and spoilage organisms.

Description of food microbiomes will only improve as additional public sequence data is collected and leveraged.

In our sample collection, 2–4% (effectively 5 to 14 million) of reads remain unclassified. The GC content distribution of unclassified reads matched microbial GC content distribution (Supplementary Fig. S4) suggesting that these reads may have been derived from microbes missing from the current reference database that have not yet been isolated or sequenced. By sequencing the microbiome, we sampled environmental niches in their native state in a culture-independent manner and therefore collected data from diverse and potentially never-before-seen microbes. Tracking unclassified reads will also be essential for monitoring food microbiomes. The inability to provide a name from existing references does not eliminate the possibility that the sequence is from an unwanted microbe or indicates a hazard. In addition to tracking known microbes, quantitative or qualitative shifts in the unclassified sequences might be used to detect when a sample is different from its peers.

We demonstrated the potential utility of analyzing food microbiomes for food safety using raw ingredients. This study resulted in the detection of shifts in the microbiome composition corresponding to unexpected matrix contaminants. This signifies that the microbiome is likely an important and effective hazard indicator in the food supply chain. While we have used total RNA sequencing for the detection of microbiome membership, the technology has future applicability for the detection of antimicrobial resistance, virulence, and biological function for multiple food sources, and for other sample types. Notably, while this pipeline was developed for food monitoring, with applicable modifications and identification of material-specific indicators, it can be applied to other microbiomes including human and environmental.

## METHODS

### Sample collection, preparation, and sequencing

HPP (HPP, 2.5 kg) samples were each collected from a train car in Reno, NV, USA between April 2015 and February 2016 in four batches from two suppliers. HPP sample was composed of five sub-samples from random locations within the train car prior to shipment. Each HPP was shipped to the Weimer laboratory at UC Davis (Davis, CA) with 2-day delivery. Upon arrival, each HPP was aliquoted into at least three tubes containing Trizol for long term storage and use in sequencing studies (see extraction section for further processing before sequencing). The remaining HPP was sealed in the plastic bag it arrived in. Those bags were put in closed storage tubs that were stored at room temperature (~25 °C) for the remainder of the study. Sample preparation, total RNA extraction, and integrity confirmation, cDNA construction, and library construction for the sample material used was described in our companion publication<sup>16</sup>.

Sequencing was performed by BGI@UC Davis (Sacramento, CA) using Illumina HiSeq 4000 (San Diego, CA) with 150 paired-end chemistry for each sample except the following: HiSeq 3000 with 150 paired-end chemistry was used for MFMB-04 and MFMB-17. All total RNA sequencing data are available via the 100K Pathogen Genome Project BioProject (PRJNA186441) at NCBI (Supplementary Table 1).

For evaluation of total RNA sequencing for microbial classification in paired processing steps, total RNA and total DNA were extracted from the same sample and denoted as MFMB-03 and MFMB-08, respectively. The total RNA was extracted and sequenced as described above. The total DNA was extracted and sequenced as described elsewhere<sup>10,48,58–62</sup>. The Illumina HiSeq 2000 with 100 paired-end chemistry was used for MFMB-03 and MFMB-08.

### Sequence data quality control

Illumina Universal adapters were removed and reads were trimmed using Trim Galore<sup>63</sup> with a minimum read length parameter 50 bp. The resulting reads were filtered using Kraken<sup>37</sup>, as described below in Section 4.3, with a custom database built from the PhiX genome (NCBI Reference Sequence: NC\_001422.1). Removal of PhiX content is suggested as it is a common

contaminant in Illumina sequencing data<sup>64</sup>. Trimmed non-PhiX reads were used in subsequent matrix filtering and microbial identification steps.

### Matrix filtering process and validation

Kraken<sup>37</sup> with a  $k$ -mer size of 31 bp (optimal size described in the Kraken reference publication) was used to identify and remove reads that matched a pre-determined list of 31 common food matrix and potential contaminant eukaryotic genomes (Supplementary Table 2). These food matrix organisms were chosen based on preliminary eukaryotic read alignment experiments of the HPP samples as well as high-volume food components in the supply chain. Due to the large size of eukaryotic genomes in the custom Kraken<sup>37</sup> database, a random  $k$ -mer reduction was applied to reduce the size of the database by 58% using Kraken-build with option `max-db-size`, in order to fit the database in 188 GB for in-memory processing. A conservative Kraken score threshold of 0.1 was applied to avoid filtering microbial reads. The matrix filtering database includes low complexity and repeats regions of eukaryotic genomes to capture all possible matrix reads. This filtering database with the score threshold was also used in the matrix filtering in silico testing as described below.

Matrix filtering was validated by constructing synthetic paired-end reads (150 bp) using DWGSIM<sup>65</sup> with mutations from reference sequences using the following parameters: base error rate ( $e$ ) = 0.005, the outer distance between the two ends of a read pair ( $d$ ) = 500, rate of mutations ( $r$ ) = 0.001, a fraction of indels ( $R$ ) = 0.15, probability an indel is extended ( $X$ ) = 0.3. Reference sequences are detailed in Supplementary Table 3. We constructed two in silico mixtures of sequencing reads by randomly sampling reads from eukaryotic reference genomes. Simulated Food Mixture 1 was comprised of nine species with the following number of reads per genome: 2 M cattle, 2 M salmon, 1 M goat, 1 M lamb, 1 M tilapia (transcriptome), 962 K chicken (transcriptome), 10 K duck, 1 K horse, and 1 K rat totaling 7.974 M matrix reads. Simulated Food Mixture 2 contained 5 M soybean, 4 M rice, 3 M potato, 2 M corn, 200 K rat, and 10 K drain fly reads, totaling 14.210 M matrix reads. Both simulated food mixtures included 1000 microbial sequence reads generated from 15 different microbial species for a total of 15 K sequence reads (Supplementary Table 3).

### Microbial identification

Remaining reads after quality control and matrix filtering were classified using Kraken<sup>37</sup> against a microbial database with a  $k$ -mer size of 31 bp to determine the microbial composition within each sample. NCBI RefSeq Complete<sup>66</sup> genomes were obtained for bacterial, archaeal, viral, and eukaryotic microorganisms (~7800 genomes retrieved April 2017). Low complexity regions of the genomes were masked using Dustmasker<sup>67</sup> with default parameters. A threshold of 0.05 was applied to the Kraken score in an effort to maximize the F-score of the result (as demonstrated in Kraken's operating manual<sup>37</sup>). Taxa-specific sequence reads were used to calculate a relative abundance in reads per million (RPM; Eq. 1), where  $R_T$  represents the reads classified per microbial entity (e.g., the genus *Salmonella*) and  $R_Q$  represents the number of sequenced reads remaining after quality control (trimming and PhiX removal) for an individual sample, including any reads classified as eukaryotic:

$$\text{RPM} = \frac{R_T}{R_Q} \times 1,000,000. \quad (1)$$

This value provides a relative abundance of the microbial entity of interest and was used in comparisons of taxa among samples. Genera with a conservative threshold of  $\text{RPM} > 0.1$  were defined as present, as previously applied by others in the contexts of human infectious disease and gut microbiome studies<sup>33,34</sup>. Pearson correlation of resulting microbial genus counts was computed.

### Community ecology analysis

Rarefaction analysis at multiple subsampled read depths  $R_D$  was performed by multiplying the microbial genus read counts with  $R_D/R_Q$  and rounding the results down to the nearest integer to represent observed read counts. Here,  $R_Q$  is the total number of reads in the sample after quality control (including microbial, matrix, and unclassified reads). Resulting  $\alpha$ -diversity at read depth  $R_D$  was computed as the number of genera with resulting  $\text{RPM} > 0.1$  and plotted at five million read intervals:  $R_D = 5 \text{ M}, 10 \text{ M}, 15 \text{ M}, \dots, R_Q$ . If, due to random sampling and rounding effects, the computed  $\alpha$ -diversity was lower than the diversity computed at any

previous depth, the previous higher  $\alpha$ -diversity was used for plotting. The median elbow was calculated as indicated using the R package `kneed`<sup>44</sup>.

In compositional data analysis<sup>31</sup>, non-zero values are required when computing  $\beta$ -diversity based on Aitchison distance<sup>45</sup>. Therefore, reads counts assigned to each genus were pseudo-counted by adding one in advance of computation of RPM (Eq. 1) prior to calculating the Aitchison distance for the microbial table.  $\beta$ -diversity was calculated using the R package `robCompositions`<sup>68</sup> and hierarchical clustering was performed using base R function `hclust` using the "ward.D2" method as recommended for compositional data analysis<sup>31</sup>.

Pairwise Spearman's correlation was computed between all samples (with the Matlab function `corr`) using the RPM vectors for the 229 detected genera as input. For the purpose of comparing correlation values within and between suppliers, the samples MFMB-04, MFMB-20, and MFMB-38 have excluded from the group "Supplier A samples" and considered as a separate matrix-contaminated group. In addition, a two-sample  $t$  test was calculated per genus on the RPM abundances from samples from Supplier A (excluding MFMB-04, MFMB-20, and MFMB-38 due to known non-poultry matrix content) and Supplier B using base R with a Benjamini-Hochberg adjustment for multiple hypothesis testing.

### Unclassified read analysis

The GC percent distributions of the matrix (from matrix filtering), microbial, and remaining unclassified reads per sample were computed using `FastQC`<sup>69</sup> and collated across samples with `MultiQC`<sup>70</sup>.

### Analysis of *Salmonella* culturability

Growth of *Salmonella* was determined using a real-time quantitative PCR method for the confirmation of *Salmonella* isolates for presumptive generic identification of foodborne *Salmonella*. Testing was performed fully in concordance with the Bacteriological Analytical Manual (BAM) for *Salmonella*<sup>71,72</sup> for this approach that is also AOAC-approved. All samples with positive results for *Salmonella* were classified as containing actively growing *Salmonella*. To compare culture results with those from total RNA sequencing, *Salmonella* RPM values were parsed from the genus-level microbe table (described in the section "Microbial identification").

Two additional approaches were employed to examine *Salmonella* read mapping with a more sensitive tool and broader reference databases. Quality controlled matrix-filtered reads were aligned using `Bowtie2`<sup>47</sup> with very-sensitive-local-mode to (1) an expanded collection of whole *Salmonella* genomes and (2) a curated growth gene reference for elongation factor Tu (*ef-Tu*). For results from both complete genome and *ef-Tu* gene alignments, the relative abundance (RPM) was computed as shown in Eq. 1.

For whole-genome alignments, a reference was constructed from 1183 recently published *Salmonella* genomes<sup>48</sup> in addition to the 264 *Salmonella* genomes extracted from the aforementioned NCBI RefSeq Complete collection (see the section "Microbial identification").

To construct a curated growth gene (*ef-Tu*) reference, gene sequences annotated in *Salmonella* genomes as "elongation factor Tu", "EF-Tu" or "eftu" (case insensitive) were retrieved from the Functional Genomics Platform (formerly OMXWare)<sup>50</sup> using its Python package. This query yielded 4846 unique gene sequences from a total of 36,242 *Salmonella* genomes which were assembled or retrieved from the NCBI Sequence Read Archive or RefSeq Complete Sequences as previously indicated<sup>50</sup>. The retrieved *ef-Tu* gene sequences were subsequently used to build a custom `Bowtie2`<sup>47</sup> reference. Read alignment was completed with very-sensitive-local-mode.

The read counts for each sample were ranked and the Wilcoxon rank-sum test was computed between the rank vectors of 4 *Salmonella*-positive and 23 *Salmonella*-negative samples. The four samples with unknown *Salmonella* status were excluded from the rankings.

Point-biserial correlation coefficients ( $r_{pb}$ ) were calculated between *Salmonella* growth indicated by culture results (+1 and -1 for presence and absence, respectively) and observed relative abundance from total RNA sequencing results using the R package `ltm`<sup>73</sup>. The point-biserial correlation is a special case of the Pearson correlation that is better suited for a binary variable e.g., when *Salmonella* is reported as present or absent (a sample's *Salmonella* status).

## DATA AVAILABILITY

All high protein powder (HPP) poultry meal sequences are available through the 100K Pathogen Genome Project (PRJNA186441) in the NCBI BioProject (see Supplementary Table 1 for a complete list of accession numbers).

## CODE AVAILABILITY

The pipeline and microbial or matrix references were constructed from publicly available tools and reference sequences as described in "Methods". Automated usability of this pipeline is available through membership in the Consortium for Sequencing the Food Supply Chain.

Received: 14 January 2020; Accepted: 24 November 2020;

Published online: 08 February 2021

## REFERENCES

- Kovac, J., Bakker, H. den, Carroll, L. M. & Wiedmann, M. Precision food safety: a systems approach to food safety facilitated by genomics tools. *TrAC Trends Anal. Chem.* <https://doi.org/10.1016/j.trac.2017.06.001> (2017).
- Weimer, B. C. et al. Defining the food microbiome for authentication, safety, and process management. *IBM J. Res. Dev.* **60**, 1 (2016).
- Walsh, A. M. et al. Microbial succession and flavor production in the fermented dairy beverage kefir. *mSystems* **1**, e00052–16 (2016).
- Walsh, A. M. et al. Species classifier choice is a key consideration when analysing low-complexity food microbiome data. *Microbiome* **6**, 50 (2018).
- Duru, I. C. et al. Metagenomic and metatranscriptomic analysis of the microbial community in Swiss-type Maasdam cheese during ripening. *Int. J. Food Microbiol.* **281**, 10–22 (2018).
- Martins, G. et al. Grape berry bacterial microbiota: Impact of the ripening process and the farming system. *Int. J. Food Microbiol.* **158**, 93–100 (2012).
- Abdelfattah, A., Wisniewski, M., Droby, S. & Schena, L. Spatial and compositional variation in the fungal communities of organic and conventionally grown apple fruit at the consumer point-of-purchase. *Hortic. Res.* **3**, 16047 (2016).
- Williams, A. G., Choi, S.-C. & Banks, J. M. Variability of the species and strain phenotype composition of the non-starter lactic acid bacterial population of cheddar cheese cheese manufactured in a commercial creamery. *Food Res. Int.* **35**, 483–493 (2002).
- Weimer, B. C. 100K pathogen genome project. *Genome Announc.* **5**, e00594–17 (2017).
- Emond-Rheault, J.-G. et al. A Syst-OMICS approach to ensuring food safety and reducing the economic burden of salmonellosis. *Front. Microbiol.* **8**, 996 (2017).
- Kaufman, J. H. et al. Insular microbiogeography: Three pathogens as exemplars. *Curr. Issues Mol. Biol.* **36**, 89–108 (2020).
- Bashiardes, S., Zilberman-Schapira, G. & Elinav, E. Use of metatranscriptomics in microbiome research. *Bioinform. Biol. Insights* **10**, 19–25 (2016).
- McGrath, K. C. et al. Isolation and analysis of mRNA from environmental microbial communities. *J. Microbiol. Methods* **75**, 172–176 (2008).
- Cottier, F. et al. Advantages of meta-total RNA sequencing (MeTRS) over shotgun metagenomics and amplicon-based sequencing in the profiling of complex microbial communities. *npj Biofilms Microbiomes* **4**, 2 (2018).
- Macklaim, J. M. et al. Comparative meta-RNA-seq of the vaginal microbiota and differential expression by *Lactobacillus iners* in health and dysbiosis. *Microbiome* **1**, 12 (2013).
- Haiminen, N. et al. Food authentication from shotgun sequencing reads with an application on high protein powders. *npj Sci. Food* **3**, 1–11 (2019).
- Lakin, S. M. et al. MEGARes: an antimicrobial resistance database for high throughput sequencing. *Nucleic Acids Res.* **45**, D574–D580 (2016).
- Noyes, N. R. et al. Resistome diversity in cattle and the environment decreases during beef production. *eLife* **5**, e13195 (2016).
- Ganesan, B., Dobrowolski, P. & Weimer, B. C. Identification of the leucine-to-2-methylbutyric acid catabolic pathway of *Lactococcus lactis*. *Appl. Environ. Microbiol.* **72**, 4264–4273 (2006).
- Ganesan, B., Seefeldt, K., Koka, R. C., Dias, B. & Weimer, B. C. Monocarboxylic acid production by lactococci and lactobacilli. *Int. Dairy J.* **14**, 237–246 (2004).
- Ganesan, B., Seefeldt, K. & Weimer, B. C. Fatty acid production from amino acids and -keto acids by *Brevibacterium linens* BL2. *Appl. Environ. Microbiol.* **70**, 6385–6393 (2004).
- Ganesan, B., Stuart, M. R. & Weimer, B. C. Carbohydrate starvation causes a metabolically active but nonculturable state in *Lactococcus lactis*. *Appl. Environ. Microbiol.* **73**, 2498–2512 (2007).
- Ganesan, B. et al. Probiotic bacteria survive in Cheddar cheese and modify populations of other lactic acid bacteria. *J. Appl. Microbiol.* **116**, 1642–1656 (2014).
- Ganesan, B. & Weimer, B. C. *Cheese: Chemistry, Physics, and Microbiology* (Elsevier, 2004).
- Shefflin, A. M., Melby, C. L., Carbonero, F. & Weir, T. L. Linking dietary patterns with gut microbial composition and function. *Gut Microbes* **8**, 113–129 (2017).
- McDonald, D. et al. American gut: an open platform for citizen science microbiome research. *mSystems* **3**, e00031–18 (2018).
- Clemente, J. C., Ursell, L. K., Parfrey, L. W. & Knight, R. The impact of the gut microbiota on human health: an integrative view. *Cell* **148**, 1258–1270 (2012).
- Richards, J. L., Yap, Y. A., McLeod, K. H., Mackay, C. R. & Mariño, E. Dietary metabolites and the gut microbiota: an alternative approach to control inflammatory and autoimmune diseases. *Clin. Trans. Immunol.* **5**, e82 (2016).
- Yang, X. et al. Use of metagenomic shotgun sequencing technology to detect foodborne pathogens within the microbiome of the beef production chain. *Appl. Environ. Microbiol.* **82**, 2433–2443 (2016).
- Hofacre, C. L. et al. Characterization of antibiotic-resistant bacteria in rendered animal products. *Avian Dis.* **45**, 953–961 (2001).
- Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J. Microbiome datasets are compositional: and this is not optional. *Front. Microbiol.* **8**, 2224 (2017).
- Gloor, G. B. & Reid, G. Compositional analysis: a valid approach to analyze microbiome high-throughput sequencing data. *Can. J. Microbiol.* **62**, 692–703 (2016).
- Langelier, C. et al. Integrating host response and unbiased microbe detection for lower respiratory tract infection diagnosis in critically ill adults. *Proc. Natl Acad. Sci. USA* **115**, E12353–E12362 (2018).
- Ilott, N. E. et al. Defining the microbial transcriptional response to colitis through integrated host and microbiome profiling. *ISME J.* **10**, 2389–2404 (2016).
- Ripp, F. et al. All-Food-Seq (AFS): a quantifiable screen for species in biological samples by deep DNA sequencing. *BMC Genomics* **15**, 639 (2014).
- Lee, A. Y., Lee, C. S. & Gelder, R. N. Van. Scalable metagenomics alignment research tool (SMART): a scalable, rapid, and complete search heuristic for the classification of metagenomic sequences from complex sequence populations. *BMC Bioinforma.* **17**, 292 (2016).
- Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biol.* **15**, R46 (2014).
- Hug, L. A. et al. A new view of the tree of life. *Nat. Microbiol.* **1**, 16048 (2016).
- Wu, D. et al. A phylogeny-driven genomic encyclopaedia of bacteria and archaea. *Nature* **462**, 1056–1060 (2009).
- Kyrpides, N. C. et al. Genomic encyclopedia of bacteria and archaea: sequencing a myriad of type strains. *PLoS Biol.* **12**, e1001920 (2014).
- Kyrpides, N. C., Eloe-Fadrosh, E. A. & Ivanova, N. N. Microbiome data science: understanding our microbial planet. *Trends Microbiol.* **24**, 425–427 (2016).
- Thompson, L. R. et al. A communal catalogue reveals Earth's multiscale microbial diversity. *Nature* **551**, 457 (2017).
- Nayfach, S. & Pollard, K. S. Toward accurate and quantitative comparative metagenomics. *Cell* **166**, 1103–1116 (2016).
- Satopaa, V., Albrecht, J., Irwin, D. & Raghavan, B. Finding a "Kneedle" in a Haystack: Detecting knee points in system behavior. In *31st International Conference on Distributed Computing Systems Workshops*, Minneapolis, USA, 20–24 June 2011, pp 166–171 (2011).
- Aitchison, J., Barceló-Vidal, C., Martín-Fernández, J. A. & Pawlowsky-Glahn, V. Logratio analysis and compositional distance. *Math. Geol.* **32**, 271–275 (2000).
- Di Palma, M. A. & Gallo, M. A co-median approach to detect compositional outliers. *J. Appl. Stat.* **43**, 2348–2362 (2016).
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
- Kong, N. et al. Draft genome sequences of 1,183 *Salmonella* strains from the 100K pathogen genome project. *Genome Announc.* **5**, e00518–17 (2017).
- Tubulekas, I. & Hughes, D. A single amino acid substitution in elongation factor Tu disrupts interaction between the ternary complex and the ribosome. *J. Bacteriol.* **175**, 240–250 (1993).
- Seabolt, E. et al. IBM functional genomics platform, a cloud-based platform for studying microbial life at scale. *IEEE/ACM Trans. Comput. Biol. Bioinforma.* 1–1. <https://doi.org/10.1109/tcbb.2020.3021231> (2020).
- Zelezniak, A. et al. Metabolic dependencies drive species co-occurrence in diverse microbial communities. *Proc. Natl Acad. Sci. USA* **112**, 6449–6454 (2015).
- Jones, M. B. et al. Library preparation methodology can influence genomic and functional predictions in human microbiome research. *Proc Natl Acad Sci USA.* <https://doi.org/10.1073/pnas.1519288112> (2015).
- Pollock, J., Glendinning, L., Wisedchanwet, T. & Watson, M. The madness of microbiome: attempting to find consensus 'best practice' for 16S microbiome studies. *Appl. Environ. Microbiol.* AEM.02627-17. <https://doi.org/10.1128/AEM.02627-17> (2018).

54. Browne, H. P. et al. Culturing of 'unculturable' human microbiota reveals novel taxa and extensive sporulation. *Nature* **533**, 543–546 (2016).
55. Eilers, H., Pernthaler, J., Glöckner, F. O. & Amann, R. Culturability and in situ abundance of pelagic bacteria from the North Sea. *Appl. Environ. Microbiol.* **66**, 3044–3051 (2000).
56. Hinchliff, C. E. et al. Synthesis of phylogeny and taxonomy into a comprehensive tree of life. *Proc. Natl Acad. Sci. USA* **112**, 12764–12769 (2015).
57. Knight, R. et al. Best practices for analysing microbiomes. *Nat. Rev. Microbiol.* **16**, 410–422 (2018).
58. Weis, A. M. et al. Genomic comparison of campylobacter spp. and their potential for zoonotic transmission between birds, primates, and livestock. *Appl. Environ. Microbiol.* **82**, 7165 LP–7167175 (2016).
59. Miller, B. et al. A novel, single-tube enzymatic fragmentation and library construction method enables fast turnaround times and improved data quality for microbial whole-genome sequencing. *Kapa Biosyst. Appl. Note* 1–8. <https://doi.org/10.13140/RG.2.1.4534.3440> (2015).
60. Lüdeke, C. H. M., Kong, N., Weimer, B. C., Fischer, M. & Jones, J. L. Complete genome sequences of a clinical isolate and an environmental isolate of *Vibrio parahaemolyticus*. *Genome Announc.* **3**, e00216–15 (2015).
61. Jeannotte, R. et al. High-throughput analysis of foodborne bacterial genomic DNA using Agilent 2200 TapeStation and genomic DNA ScreenTape system. *Agil. Appl. Note* 1–8. <https://doi.org/10.6084/m9.figshare.1372504> (2015).
62. Arabyan, N. et al. Salmonella degrades the host glycocalyx leading to altered infection and glycan remodeling. *Sci. Rep.* **6**, 1–11 (2016).
63. Krueger, F. TrimGalore: A wrapper around Cutadapt and FastQC to consistently apply adapter and quality trimming to FastQ files, with extra functionality for RRBS data. *GitHub*. Available online at: <https://github.com/FelixKrueger/TrimGalore> (2018). Accessed 28 Jun 2018.
64. Mukherjee, S., Huntemann, M., Ivanova, N., Kyrpides, N. C. & Pati, A. Large-scale contamination of microbial isolate genomes by Illumina PhiX control. *Stand. Genom. Sci.* **10**, 18 (2015).
65. Homer, N. DWGIM: Whole genome simulator for next-generation sequencing. *GitHub*. <https://github.com/nh13/DWGSIM> (2011). Accessed 14 Jun 2017.
66. O'Leary, N. A. et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745 (2016).
67. Morgulis, A., Gertz, E. M., Schäffer, A. A. & Agarwala, R. A fast and symmetric DUST implementation to mask low-complexity DNA sequences. *J. Comput. Biol.* **13**, 1028–1040 (2006).
68. Templ, M., Hron, K. & Filzmoser, P. robCompositions: an R-package for Robust Statistical Analysis of Compositional Data. In: Buccianti A. & Pawlowsky-Glahn V. *Compositional Data Analysis*, John Wiley & Sons, Ltd, pp 341–355 (2011).
69. Andrews, S. FastQC: A quality control tool for high throughput sequence data. *Babraham Bioinformatics*. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/> (2010). Accessed 01 Oct 2018.
70. Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048 (2016).
71. Andrews, W. H., Wang, H., Jacobson, A. & Hammack, T. Bacteriological analytical manual (BAM) Chapter 5: Salmonella. In *Bacteriological Analytical Manual U.S. Food and Drug Administration* (2018). Accessed 21 Jun 2019.
72. Grim, C. J. et al. High-resolution microbiome profiling for detection and tracking of *Salmonella enterica*. *Front. Microbiol.* **8**, 1587 (2017).
73. Rizopoulos, D. ltm: an R package for latent variable modeling and item response theory analyses. *J. Stat. Softw.* **17**, 1–25 (2006).

## ACKNOWLEDGEMENTS

The authors like to acknowledge the IBM Research Functional Genomics Platform team (formerly OMXWare) for their data management support and availability for the retrieval and processing of microbial genomes. This research project was financially supported by the Consortium for Sequencing the Food Supply Chain. Funding for the

total RNA sequencing of high protein powder factory ingredients was provided by Mars, Incorporated to B.C.W. with a specific interest in metagenomics of the food microbiome.

## AUTHOR CONTRIBUTIONS

K.L.B. and N.H. conceived of the experimental design, developed the approach, completed and oversaw the experiments, performed analyses, and wrote the paper and are represented as co-first authors; D.C., S.E., M.K., B.K., M.D., R.P., H.K., and E.S. developed the approach, analyzed the data, and revised the paper; B.C.H. completed nucleic acid extraction method development and sequencing library construction, and contributed to data analysis and writing; N.K. coordinated sample collection and processing, nucleic acid extraction, and contributed to writing; R.B. and P.M. conceived of the experimental design, developed the approach, and reviewed the paper; B.G. contributed to the experimental design, developed the approach, and wrote the paper; G.D., C.H.M., S.P., and A.Q. participated to the conception of the experimental design and to the review of the paper; L.P. conceived of the experiment, contributed to the data analysis, and wrote the paper; J.H.K. conceived of the experiment, developed the approach, and wrote the paper; B.C.W. conceived of the experimental design, developed the approach, oversaw the experiments, performed analyses, and wrote the paper.

## COMPETING INTERESTS

The authors were employed by private or academic organizations as described in the author affiliations at the time this work was completed. IBM Corporation, Mars Incorporated, and Bio-Rad Laboratories are members of the Consortium for Sequencing the Food Supply Chain. The authors declare no other competing interests.

## ADDITIONAL INFORMATION

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41538-020-00083-y>.

**Correspondence** and requests for materials should be addressed to K.L.B. or B.C.W.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021