

# A Systematic Analysis of Biosynthetic Gene Clusters in the Human Microbiome Reveals a Common Family of Antibiotics

Mohamed S. Donia,<sup>1</sup> Peter Cimermancic,<sup>1</sup> Christopher J. Schulze,<sup>2</sup> Laura C. Wieland Brown,<sup>3</sup> John Martin,<sup>4</sup> Makedonka Mitreva,<sup>4</sup> Jon Clardy,<sup>5</sup> Roger G. Linington,<sup>2</sup> and Michael A. Fischbach<sup>1,\*</sup>

<sup>1</sup>Department of Bioengineering and Therapeutic Sciences and the California Institute for Quantitative Biosciences, University of California, San Francisco, San Francisco, CA 94158, USA

<sup>2</sup>Department of Chemistry and Biochemistry, University of California, Santa Cruz, Santa Cruz, CA 95064, USA

<sup>3</sup>Department of Chemistry, Indiana University, Bloomington, IN 47405, USA

<sup>4</sup>The Genome Institute, Department of Medicine and Department of Genetics, Washington University School of Medicine, St. Louis, MO 63108, USA

<sup>5</sup>Department of Biological Chemistry and Molecular Pharmacology, Harvard Medical School, Boston, MA 02115, USA

\*Correspondence: [fischbach@fischbachgroup.org](mailto:fischbach@fischbachgroup.org)

<http://dx.doi.org/10.1016/j.cell.2014.08.032>

## SUMMARY

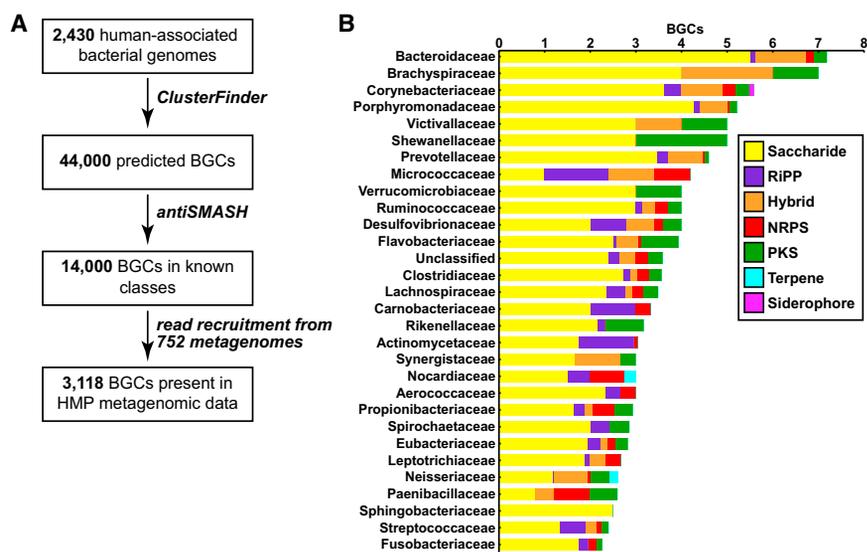
In complex biological systems, small molecules often mediate microbe-microbe and microbe-host interactions. Using a systematic approach, we identified 3,118 small-molecule biosynthetic gene clusters (BGCs) in genomes of human-associated bacteria and studied their representation in 752 metagenomic samples from the NIH Human Microbiome Project. Remarkably, we discovered that BGCs for a class of antibiotics in clinical trials, thiopeptides, are widely distributed in genomes and metagenomes of the human microbiota. We purified and solved the structure of a thiopeptide antibiotic, lactocillin, from a prominent member of the vaginal microbiota. We demonstrate that lactocillin has potent antibacterial activity against a range of Gram-positive vaginal pathogens, and we show that lactocillin and other thiopeptide BGCs are expressed *in vivo* by analyzing human metatranscriptomic sequencing data. Our findings illustrate the widespread distribution of small-molecule-encoding BGCs in the human microbiome, and they demonstrate the bacterial production of drug-like molecules in humans.

## INTRODUCTION

The human microbiome is composed of hundreds of bacterial species and thousands of strains, and its composition differs from person to person and between different body sites of the same individual (Human Microbiome Project Consortium, 2012b). During the last decade, tremendous efforts have been made to sequence isolates of the human microbiota and metagenomic samples from various body sites (Human Microbiome Project Consortium, 2012a, 2012b; Nelson et al., 2010; Qin

et al., 2010). These studies have yielded a basic understanding of the “healthy” human microbiome and have correlated deviations from the healthy state to maladies such as obesity, diabetes, bacterial vaginosis, and Crohn’s disease (Gajer et al., 2012; Gevers et al., 2012, 2014; Ravel et al., 2011; Turnbaugh et al., 2009). Several recent studies have begun to examine the human microbiome from a functional point of view, where direct molecular interactions between host and microbe are revealed (An et al., 2014; Hsiao et al., 2013; Mazmanian et al., 2005; Mazmanian et al., 2008; Nougayrède et al., 2006; Wieland Brown et al., 2013; Wyatt et al., 2010).

Diffusible and cell-associated small molecules often mediate host-microbe interactions in complex environments. Examples of small-molecule-mediated interactions have been revealed in symbioses between bacteria and insects (Oh et al., 2009), marine invertebrates (Kwan et al., 2012), nematodes (McInerney et al., 1991), and plants (Long, 2001). In addition, several studies have explored the role of small molecules in interactions between microbiota and the mammalian host. For example, *Staphylococcus aureus* pyrazinones were shown to be inducers of bacterial virulence (Wyatt et al., 2010), the *Escherichia coli* metabolite colibactin was found to contribute to colon cancer (Nougayrède et al., 2006), and polysaccharide A from *Bacteroides fragilis* has been shown to suppress the gut mucosal immune response (Mazmanian et al., 2005, 2008). Recently, we and others showed that *Bacteroides fragilis* produces the canonical CD1d ligand  $\alpha$ -galactosylceramide, revealing a specific mechanism by which the gut microbiota are capable of modulating host natural killer T cell function (An et al., 2014; Wieland Brown et al., 2013). Another recent study correlated the prevalence of hepatic cancer in mice to deoxycholic acid, a secondary bile acid produced by certain members of the gut microbiota (Yoshimoto et al., 2013). A recent metatranscriptomic study showed the expression of genes matching the clusters of orthologous groups (COG) category “secondary metabolites biosynthesis, transport, and catabolism,” which is consistent with the possibility of small-molecule production but could also indicate



**Figure 1. Overview of BGCs in the Human Microbiome**

(A) Computational and experimental workflow for the identification of BGCs from human-associated bacteria.

(B) A bar graph showing the top 30 families by average BGC abundance from the human microbiome and the average number and type of BGCs discovered in isolates of each genus using the workflow in (A). See also [Figure S1](#) and [Data S1](#) for the full data set of predicted BGCs.

the expression of catabolic and/or transport genes unrelated to biosynthesis (Leimena et al., 2013). These examples raise the question of whether there exists a much larger set of bacterially produced molecules that mediate microbiota-host interactions. Due to the complexity of the human microbiome and its vast coding potential, a more systematic approach is needed to explore small-molecule-mediated interactions between humans and their microbiota.

In this study, we explored the biosynthetic capacity of the human microbiome by performing the first systematic identification and analysis of its biosynthetic gene clusters. Unlike previous approaches that have focused on one compound or bacterial strain at a time, our approach allows the global analysis of biosynthetic gene clusters (BGCs) that encode small molecules in thousands of isolates of the human microbiota. By measuring the representation of these BGCs in human metagenomic samples, we can assess the small-molecule coding capacity of a community, generating powerful hypotheses about which molecules might mediate microbe-host and microbe-microbe interactions in a particular community and how their prevalence differs among individuals. To illustrate the utility of this approach, we used a combination of chemistry, genetics, metagenomics, and metatranscriptomics to study a family of gene clusters that is widely distributed in the human microbiome, including the characterization of its small-molecule product and the analysis of its prevalence among body sites and individuals.

## RESULTS

### A Systematic Approach Identifies 3,118 Biosynthetic Gene Clusters that Are Present in Human Metagenomic Samples

The identification of biosynthetic gene clusters in bacterial genome sequences has become a powerful tool for natural product discovery. We began by using an algorithm we recently developed, ClusterFinder, to systematically analyze 2,430 reference genomes of the human microbiota from a range of body

sites (Cimermancic et al., 2014). ClusterFinder detected >14,000 biosynthetic gene clusters (average of 6 gene clusters per genome) for a broad range of small-molecule classes, including saccharides, nonribosomal peptides (NRPs), polyketides (PKs), ribosomally encoded and posttranslationally modified peptides

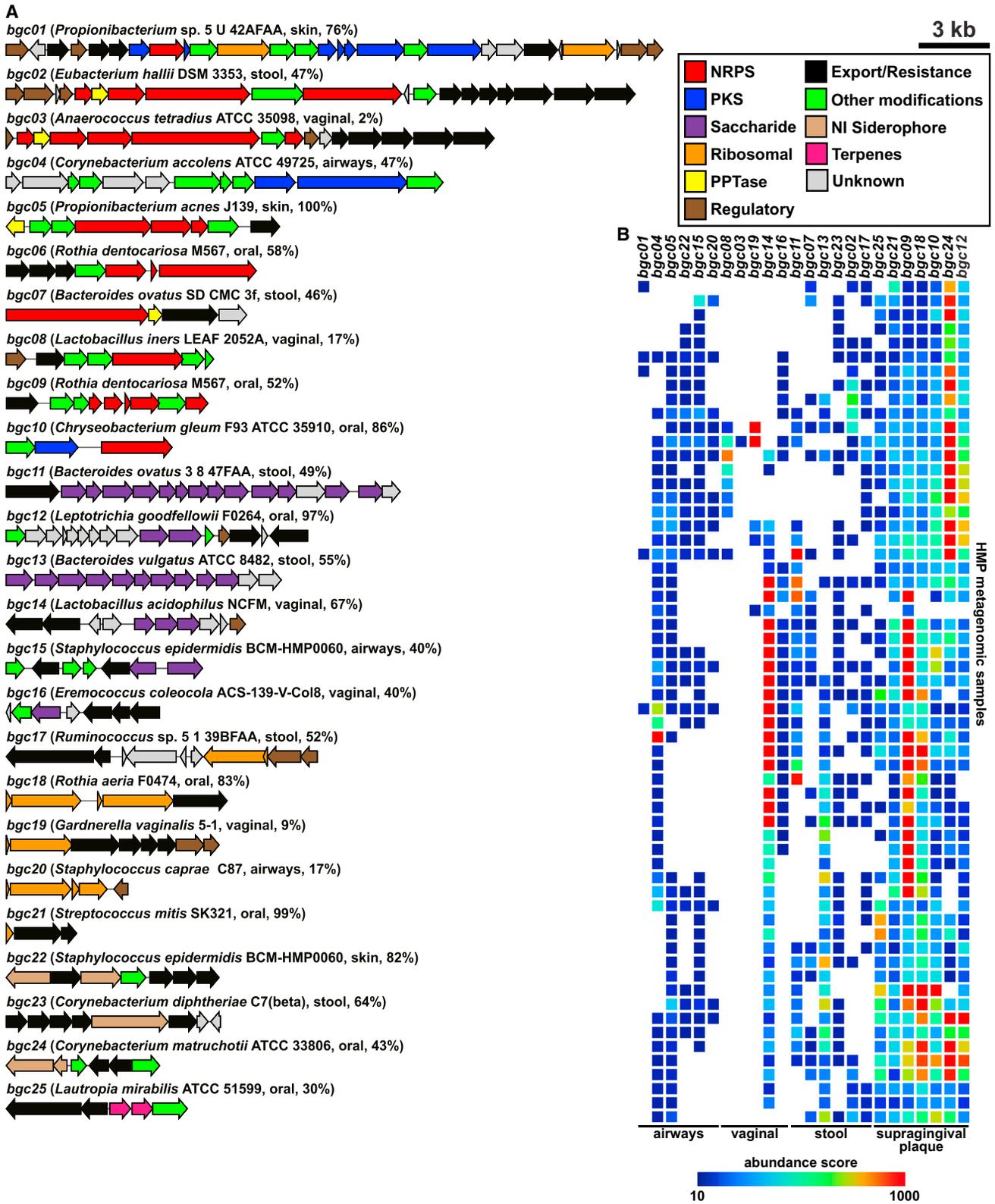
(RiPPs), NRPs-independent siderophores, and hybrids thereof (Figure 1 and Figure S1 available online).

Reasoning that the small-molecule products of BGCs that are widely distributed in the human population are most likely to mediate conserved microbe-host and microbe-microbe interactions, our next goal was to identify the subset of BGCs that is commonly found in healthy individuals. To achieve this goal, we examined the representation of these gene clusters in shotgun metagenomic sequencing data generated by the Human Microbiome Project (HMP, 752 samples collected from five main body sites of healthy subjects) (Human Microbiome Project Consortium, 2012b). Briefly, we used mblastx (Davis et al., 2013) to detect metagenomic reads that match the >14,000 predicted BGCs (see [Experimental Procedures](#) for a more detailed description of how gene clusters were quantified in metagenomic samples). Using this method, we calculated that a subset of 3,118 biosynthetic gene clusters (22%) are present in the microbiomes of healthy individuals; all subsequent analyses focused on this subset (Figures 1 and S1 and Data S1).

### An Overview of BGC Types and Distribution in the Human Microbiome

The identification of >3,000 biosynthetic gene clusters in the human microbiome was remarkable, given that almost nothing is known about their small-molecule products or biological activities. We next examined the distribution and abundance of each of these BGCs in the 752 HMP samples.

The gut and oral cavity are by far the richest in BGCs, reflecting their increased microbial abundance relative to the other body sites (see also Figure S2A, see [Experimental Procedures](#)) (Human Microbiome Project Consortium, 2012b). A typical gut harbors 599 clusters, whereas a typical oral cavity harbors 1,061 clusters; notably, the SD of each of these values is large (152 and 143, respectively), pointing to a far greater degree of gene cluster diversity in some individuals than others. Typical skin, airway, and urogenital tract samples harbor many fewer BGCs on average (101, 31, and 41, respectively, see also Table S1).



(legend on next page)

The smaller number of BGCs in the skin, airway, and urogenital tract communities could be a result of the lower microbial diversity in these communities in comparison to the oral and gut microbiota (Human Microbiome Project Consortium, 2012b). In addition, fewer samples were sequenced from skin, airways, and urogenital tract communities (27, 65, and 94 samples, respectively) than from oral and gut communities (415 and 147, respectively), possibly affecting the total number of BGCs that were found to be present in these body sites.

The gene cluster classes in the human microbiota differ in important ways from those in non-human-associated bacteria (see also Figure S2B). Even in light of the predominance of saccharides among environmental bacteria (Cimermancic et al., 2014), saccharides are significantly enriched in the microbiota (see also Figure S2B), making them by far the most abundant BGC class in the gut and oral cavity (average of 443 and 662 saccharide clusters per sample, respectively). RiPPs, which are modestly enriched, are broadly distributed in every body site. Although nonribosomal peptide synthetase (NRPS) and polyketide synthase (PKS) gene clusters are significantly depleted in the microbiota, they are still present at moderate levels (average of 57 PKS and 19 NRPS clusters in a gut sample and 129 PKS and 46 NRPS clusters in an oral cavity sample; see also Table S1), and notable exemplars are widely distributed in the healthy human population (see below). Most of the BGCs from nonhuman environmental isolates are harbored by members of the abundant genera *Streptomyces*, *Bacillus*, *Pseudomonas*, *Burkholderia*, and *Myxococcus* (Cimermancic et al., 2014); in contrast, most human-associated BGCs are harbored by members of the abundant human-associated genera *Bacteroides*, *Parabacteroides*, *Corynebacterium*, *Rothia*, and *Ruminococcus* (up to seven BGCs per genome), pointing to undermined taxa that should be rich BGC sources for experimental mining efforts (Figure 1 and see also Figure S1). Many of these genomes harbor large BGCs despite being only 2 to 3 Mb, suggesting an important ecological role for their small-molecule products. Some common genera of the human microbiota, including *Escherichia*, *Lactobacillus*, *Haemophilus*, and *Enterococcus*, harbor fewer than two BGCs per genome. Taken together, these findings indicate that the human microbiome harbors a rich and diverse array of BGCs that is mostly distinct in source and composition from that of nonhuman isolates.

The best-studied gene clusters are not widely distributed among healthy individuals. With the exception of the pyrazinone BGC from *S. aureus* (found in >9% of the airways samples), two of the most thoroughly studied BGCs from human-associated bacteria are found in <5% of the samples from their body site of origin (colibactin from *E. coli* and polysaccharide A from *Bacteroides fragilis* (Mazmanian et al., 2005; Nougayrède et al.,

2006; Wyatt et al., 2010). In contrast, nearly all of the BGCs from our data set that are widely distributed in healthy humans (present in >10% of the samples from the body site of origin) have never been studied or even described; a selected subset of these clusters is shown in Figure 2 (*bgc01–bgc25*). These results illustrate how little is known about the small-molecule products of the most common BGCs in human-associated niches.

Intriguingly, a smaller subset of 519 clusters is present in >50% of the gut samples; similar “common” BGC subsets are present in the oral cavity (582 clusters), skin (65 clusters), urogenital tract (16 clusters), and airways (11 clusters) (see also Table S1). In addition, several widely distributed families of closely related BGCs were revealed by our analysis. Among these BGC families, four examples stood out because of their predominance in a body site, resemblance to a known BGC, or cosmopolitan distribution: a family of NRPS BGCs from the gut, a pair of PKS BGCs from the oral cavity, Bacteroidetes saccharides, and RiPPs.

### A Large Family of Nonribosomal Peptide BGCs Is Widely Distributed in the Gut Community

The first BGC family we identified is a set of NRPS clusters that are found exclusively in gut isolates (hereafter *bgc26–bgc54*). Surprisingly, members of this family are harbored by species of a wide variety of Gram-positive (*Clostridium*, *Ruminococcus*, *Eubacterium*, *Lachnospiraceae*, *Blautia*) and Gram-negative genera (*Bacteroides* and *Desulfovibrio*). In addition to human gut isolates, this NRPS family is found in isolates from the chicken gut (*Bacteroides barnesiae* DSM 18169) and bovine rumen (*Methanobrevibacter ruminantium* M1) (Figure 3). Remarkably, 137/149 (92%) of the HMP stool samples contained at least one member of this family of NRPS clusters (see also Figure S3), whereas none or very few of the airway (0/94), oral (2/406), vaginal (2/65), and skin (1/27) metagenomic samples harbored a member of this family. The wide distribution and broad representation of this NRPS family suggests its importance for gut microbial inhabitants. Members of this BGC family fall into two main groups: the first consists of three NRPS modules, whereas the second is composed of only two (Figure 3). The predicted substrate of the penultimate adenylation domain is an aliphatic amino acid (Leu, Ile, Val, or Ala), whereas the other adenylation domains have predicted substrates that vary widely among the clusters, suggesting that the small-molecule products of these BGCs may be a family of closely related molecules.

### A Family of Complex Polyketides from Oral Actinobacteria Are Widely Distributed in the Oral Cavity

The most widely distributed multimodular PKS in the oral cavity is a pair of closely related ~80 kb clusters of the trans-AT type

#### Figure 2. A Selected Subset of BGCs from the Human Microbiota

(A) 25 selected BGCs from the human microbiota, spanning each of the body sites (gut, vagina, airways and skin, and oral cavity), BGC types (PKS, NRPS, RiPPs, terpenes, NI siderophores, and saccharides), and prevalent bacterial phyla (Actinobacteria, Bacteroidetes, Firmicutes, and Proteobacteria). The label of each gene cluster indicates its source organism, body site of origin, and the percentage of HMP samples harboring this cluster in its body site of origin. All but two of these BGCs are present in more than 10% of the samples from their body site of origin, indicating that they are widely distributed among healthy subjects. (B) Heat map showing the representation of BGCs from (A) in a subset of 60 selected HMP metagenomic samples from four body sites. The color of the cells in the heat map represents an abundance score ranging from 10 (blue) to 1,000 (red) (key shown to the bottom, see Extended Experimental Procedures for calculation of abundance scores and see also Figure S2 and Table S1).

(hereafter *bgc54–bgc55*). These clusters, which are found in two different species of oral Actinobacteria (*Propionibacterium propionicum* F0230a and *Actinomyces timonensis* DSM 23838), are strikingly similar in their domain architecture to a BGC from the marine isolate *Streptomyces* sp. A7248 that encodes the production of the homodimeric macrolide SIA7248 (Figure 4) (Zou et al., 2013). Despite this similarity, the enzymes encoded by the three clusters share only ~40% identity, and no mobile elements were found in either of the human-associated clusters. SIA7248 is chemically similar to the marinomycins, which are known to have potent antibacterial and antitumor activity (Kwon et al., 2006). Our computational analysis of PKS architecture and AT domain selectivity predicts that the small-molecule products of both PKS clusters will be nearly identical to the marinomycins and SIA7248, except for an alternative loading moiety activated by an acyl-CoA ligase domain that is absent in the SIA7248 cluster. Notably, *bgc54* is widely distributed in healthy subjects (34% of the HMP supragingival plaque samples), indicating that a close relative of a complex marine bacterial metabolite might be common in the oral cavity.

#### Bacteroidetes Saccharides Are Predominant and Variable in the Gut

Saccharide BGCs are the most abundant family in each of the five body sites and are particularly predominant in the gut (74% of all BGCs in a typical HMP stool sample are saccharides; see also Table S1). The phylum Bacteroidetes harbors the largest number of saccharide BGCs (Cimermancic et al., 2014), two of which have been structurally and functionally characterized (polysaccharide A and B; Baumann et al., 1992; Mazmanian et al., 2005, 2008; Tzianabos et al., 1993). Two observations were notable. First, among the most common saccharide BGCs were capsular polysaccharide loci from *B. vulgatus* and *B. ovatus* (see also Figure S4), which differ considerably from the better-known examples in *B. fragilis*, suggesting that structurally distinct capsular polysaccharides might be remarkably common in the healthy human population. Second, ten HMP stool samples that show similar taxonomic composition by MetaPhlAn profiling (Segata et al., 2012) harbor largely distinct sets of saccharide BGCs, with only a small number of BGCs common among the samples (see also Figure S4). Thus, knowledge about the taxonomic composition of gut communities from 16S rDNA analysis or metagenomic classification tools does not reveal the inherent diversity of BGCs in these communities.

#### Ribosomally Synthesized, Posttranslationally Modified Peptides Are among the Most Widely Distributed and Variable BGCs Encoded by the Human Microbiota

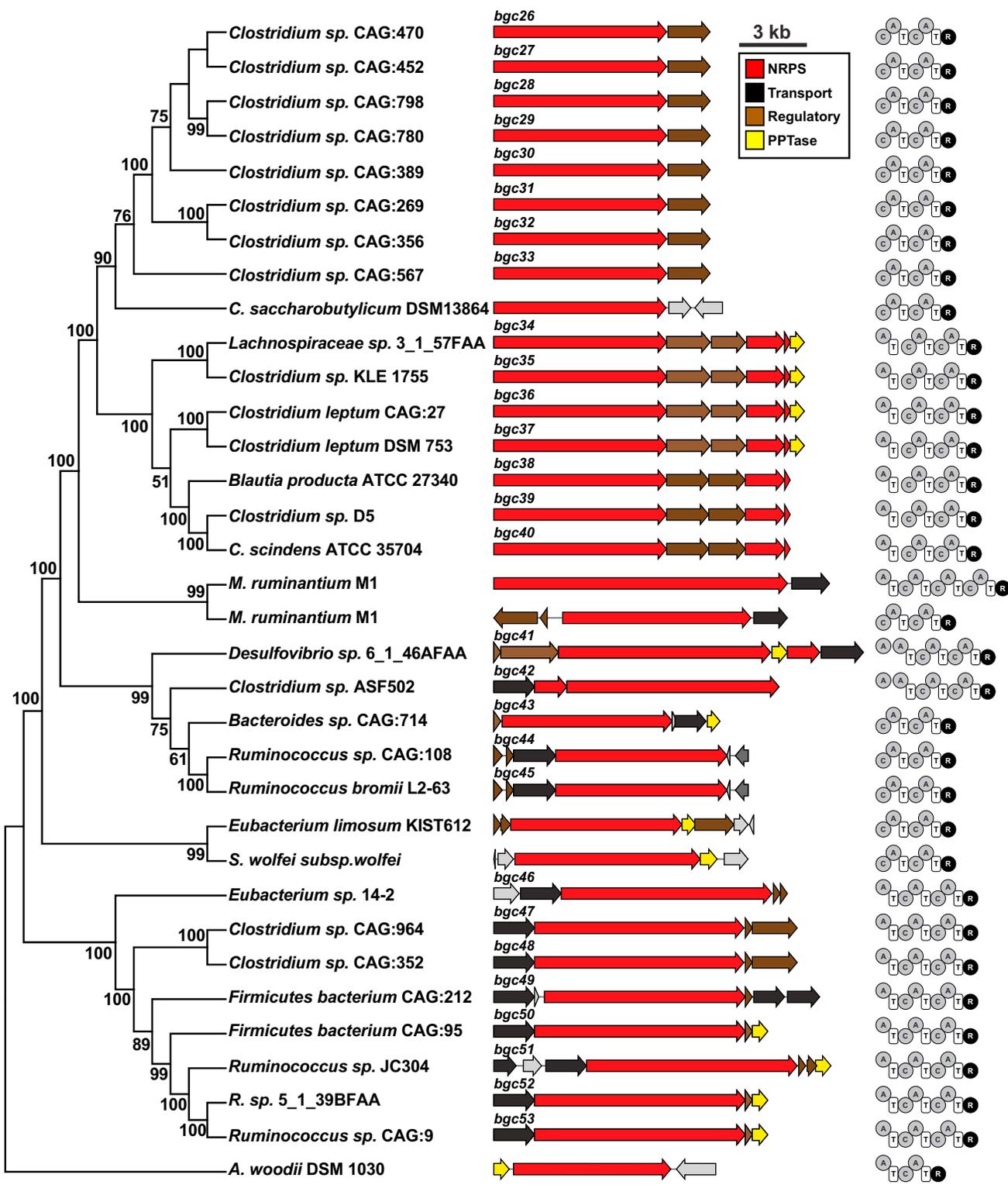
BGCs encoding modified RiPPs are widely distributed among the human microbiota. Three subclasses of RiPPs were particularly notable: lantibiotics, which were numerous, variable, and broadly distributed throughout the Firmicutes and Actinobacteria from all of the major body sites; thiazole/oxazole-modified microcins (TOMMs), which were prevalent in the oral cavity; and thiopeptides. Although a small set of lantibiotics and TOMMs have been isolated from members of the human microbiota (Chikindas et al., 1995; Gonzalez et al., 2010; Lee et al., 2008), most are produced by pathogens or rare commensals.

Thiopeptides are highly modified RiPPs that have potent antibacterial activity against Gram-positive species. A semisynthetic member of this class, LFF571 (Novartis), is currently in phase II clinical trials for treating *Clostridium difficile* infection (LaMarche et al., 2012). Although a thiopeptide BGC was previously predicted by our group and others in the genome of the skin isolate *Propionibacterium acnes* (Brzuszkiewicz et al., 2011; Wieland Brown et al., 2009), no thiopeptide BGC or small-molecule product from the human microbiome has ever been characterized experimentally. Remarkably, in the analysis reported here, we discovered thiopeptide-like BGCs in isolates from every human body site (*Lactobacillus gasseri*, urogenital; *Propionibacterium acnes*, skin; *Streptococcus downei* and *S. sobrinus*, oral; and *Enterococcus faecalis*, gut). Additionally, we discovered a thiopeptide gene cluster in the porcine gut isolate *Lactobacillus johnsonii* PF01, indicating that thiopeptide BGCs are found in the microbiota of animals other than humans (Figure 5A).

We next turned to the question of whether we could discover additional thiopeptide gene clusters in metagenomic data that were not present in any sequenced bacterial genome. To address this question, we mined HMP metagenomic assemblies for the presence of additional thiopeptide gene clusters; remarkably, we discovered eight additional thiopeptide BGCs (7 complete, 1 partial) using this strategy, increasing the number of thiopeptide BGCs identified here to 13 (*bgc56–bgc68*) (Figure 5A) (see Experimental Procedures). Of these, we found one in human gut metagenomic samples (*bgc61*) and seven in human oral metagenomic data sets (*bgc58*, *bgc59*, *bgc62*, *bgc64*, *bgc65*, *bgc67*, and *bgc68*).

An analysis of the genes flanking these clusters suggests that *bgc61* resides in a prominent member of the human gut microbiome (*Eubacterium rectale*), whereas two of the oral metagenomic thiopeptide gene clusters are found in *Actinomyces* and five others are harbored by *Streptococcus* (see also Table S2). Transposases and phage integrases are unusually prevalent among the thiopeptide BGCs (70%), and in at least one case (*bgc66*), we could show bioinformatically and experimentally that the cluster exists on a plasmid, suggesting a potential for mobility (see below, and note that the thiopeptide BGC from the porcine gut isolate *Lactobacillus johnsonii* PF01 [Lee et al., 2011] also exists on a plasmid). Four of the 13 thiopeptide BGCs are present in >20% of the HMP samples at one of the body sites, and 155/406 HMP oral samples (38%) harbor at least one thiopeptide BGC (Figure 5C and Table S2). For example, *bgc61* is present in 11% of gut samples and *bgc65* is present in 34% of oral samples, indicating that there exist widely distributed biosynthetic gene clusters in the healthy human population that are not found in the reference genome database.

In order to gain more insight into the evolution of these human-associated thiopeptide BGCs, we performed a detailed bioinformatic analysis of their biosynthetic genes and precursor peptides. Genomic and metagenomic thiopeptide BGCs fall into six subfamilies in which members of each subfamily harbor similar precursor peptides (Figures 5A and 5B). Notably, the BGC precursor peptides and enzymes that posttranslationally modify them cluster largely according to their body site of origin rather than the phylogeny of their host (Figures 5A, 5B, and S5). Moreover, the distribution of the oral thiopeptide clusters across



**Figure 3. An Abundant Family of NRPS BGCs Is Found Exclusively in Gut Isolates and Stool Metagenomes**  
 Left, a phylogenetic tree (Maximum Parsimony, MEGA5) based on the main NRPS gene of 28 clusters from human bacterial gut isolates, two from bovine rumen archaeal isolates, and four from environmental isolates (see [Experimental Procedures](#)). The numbers next to the branches represent the percentage of replicate trees in which this topology was reached in a bootstrap test of 1,000 replicates. Middle, schematic of each BGC (see also [Figure S3](#) for a full heat map showing the prevalence and abundance of each member of this NRPS family in HMP stool samples). Right, domain organization of the NRPS genes of each cluster (A, adenylation domain; C, condensation domain; T, thiolation domain; R, terminal reductase domain).



1224.15184, calculated  $[M+H]^+$  1224.15354,  $\Delta$ ppm = 1.3). This mass and empirical formula deviated from the computationally predicted mass and empirical formula for the *bgc66* product (1058.0667 and  $C_{42}H_{34}N_{12}O_8S_7$ , respectively), suggesting that the *bgc66* product harbored at least one unknown posttranslational modification.

To obtain multimilligram quantities of the *bgc66* product for structural characterization, we performed a 50 L cultivation of *L. gasseri* JV-V03. Our initial attempts to perform NMR experiments on the isolated product failed due to its low solubility in NMR solvents and apparent instability. Hypothesizing that the *bgc66* product harbored a free carboxylic acid that was partly responsible for its insolubility and instability, we treated the product extracted from a second 60 L cultivation with TMS-diazomethane to convert any free carboxylic acids to methyl esters. The resulting product had a mass consistent with the addition of a single methyl group and exhibited greatly improved solubility and stability. After extensive purification by HPLC, the purified product (0.5 mg) was characterized using a series of 1D and 2D NMR experiments, high-resolution tandem mass spectrometry (MS/MS and MS<sup>n</sup>), and isotope labeling experiments. Collectively, these data enabled us to determine the structure of the *bgc66* product, to which we assign the name lactocillin (see [Experimental Procedures](#), [Figures 6B](#) and [S6](#), and [Data S2](#)).

Lactocillin harbors a canonical 26-membered thiopeptide scaffold, with a trithiazolylpyridine core and a macrocycle with four cysteine-derived heterocycles and a single dehydrobutyrine residue. Three structural features set it apart from thiocillin. (1) The presence of an indolyl-S-cysteine residue at position 8, a feature that is structurally and regiochemically reminiscent of the 3-methylindolyl and quinaldic acid moieties of nosiheptide and thiostrepton, respectively ([Just-Baringo et al., 2014](#); [Kelly et al., 2009](#); [Liao et al., 2009](#); [Yu et al., 2009](#)). Interestingly, the enzymes predicted to install the indolyl-S-cysteine in lactocillin (*lclI*, a thiolation domain, *lclJ*, an adenylation domain, and *lclK*, an alpha-beta hydrolase) are unrelated to any genes in the nosiheptide or thiostrepton BGCs, indicating that it may have arisen by convergent evolution. Lactocillin also (2) harbors a free carboxylic acid at the C terminus, a rare feature among thiopeptides ([Just-Baringo et al., 2014](#)) and (3) lacks any oxygen-requiring posttranslational modifications, reflecting the anaerobic conditions under which *L. gasseri* thrives in the vaginal community ([Figure 6B](#)).

Next, we set out to determine whether lactocillin has a spectrum of activity similar to other thiopeptides, which are potent antibiotics against Gram-positive, but not Gram-negative, bacteria. To answer this question, we tested purified lactocillin against commonly observed pathogens and commensals from the vaginal community. Lactocillin was active against *Staphylococcus aureus*, *Enterococcus faecalis*, *Gardnerella vaginalis*, and *Corynebacterium aurimucosum*, but not *Escherichia coli*, an activity spectrum similar to that of other thiopeptides and suggestive of a potential role in protecting the vaginal microbiota against pathogen invasion. Interestingly, lactocillin was inactive against the vaginal commensals *Lactobacillus jensenii* JV-V16, *Lactobacillus crispatus* JV-V01, and *L. gasseri* SV-16A-US (a different strain from the lactocillin producer) at 425 nM, raising the possibility that the vaginal microbiota have evolved resistance against a compound they commonly encounter ([Table 1](#) and see [Experimental Procedures](#)).

### Metatranscriptomic Data Analysis Reveals that the *lcl* Cluster Is Expressed in Humans

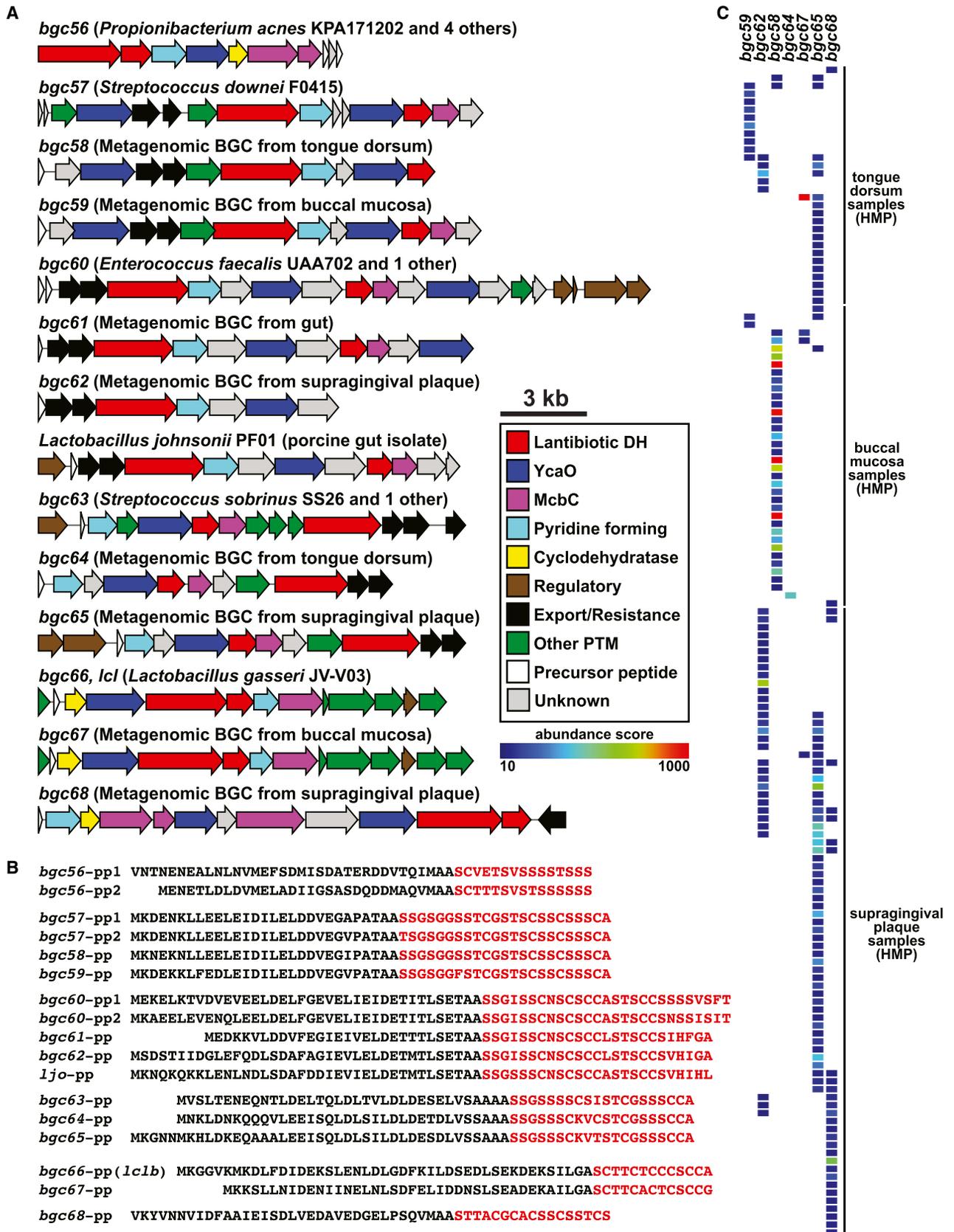
We next examined the distribution pattern of the *lcl* gene cluster in the human microbiota. In contrast to other human-associated thiopeptide BGCs ([Figures 5A](#) and [5C](#)), the *lcl* gene cluster appears to have a limited distribution in genomic and metagenomic data sets: *L. gasseri* JV-V03 is the only one of the ten sequenced *L. gasseri* genomes that harbors the cluster (see [Extended Experimental Procedures](#)), and none of the six vaginal metagenomic samples in which *L. gasseri* is the dominant strain contain reads that match the *lcl* cluster. However, building on the finding that *lcl* resides on a plasmid ([Figure 6C](#)) and that some BGCs are present in more than one body site (see also [Figure S2](#)), we broadened our search to publicly available metatranscriptomic data sets from the oral cavity (see [Experimental Procedures](#)). To our surprise, we found that two thiopeptide gene clusters—*lcl* and *bgc65*—were covered by oral metatranscriptomic reads in multiple samples (in total, reads matching *lcl* and *bgc65* were found in 3/38 and 11/38 of the supragingival plaque metatranscriptomic samples, respectively) (see [Experimental Procedures](#), [Figure 6D](#) and see also [Figure S7](#)). In addition, the rest of genes on the plasmid in which the *lcl* cluster resides were also present in oral metatranscriptomic samples, further emphasizing its mobility (see also [Figure S7](#)). These results show that thiopeptide gene clusters are actively transcribed in human samples, suggesting a potential role in mediating microbe-microbe interactions in the communities in which they are expressed.

### DISCUSSION

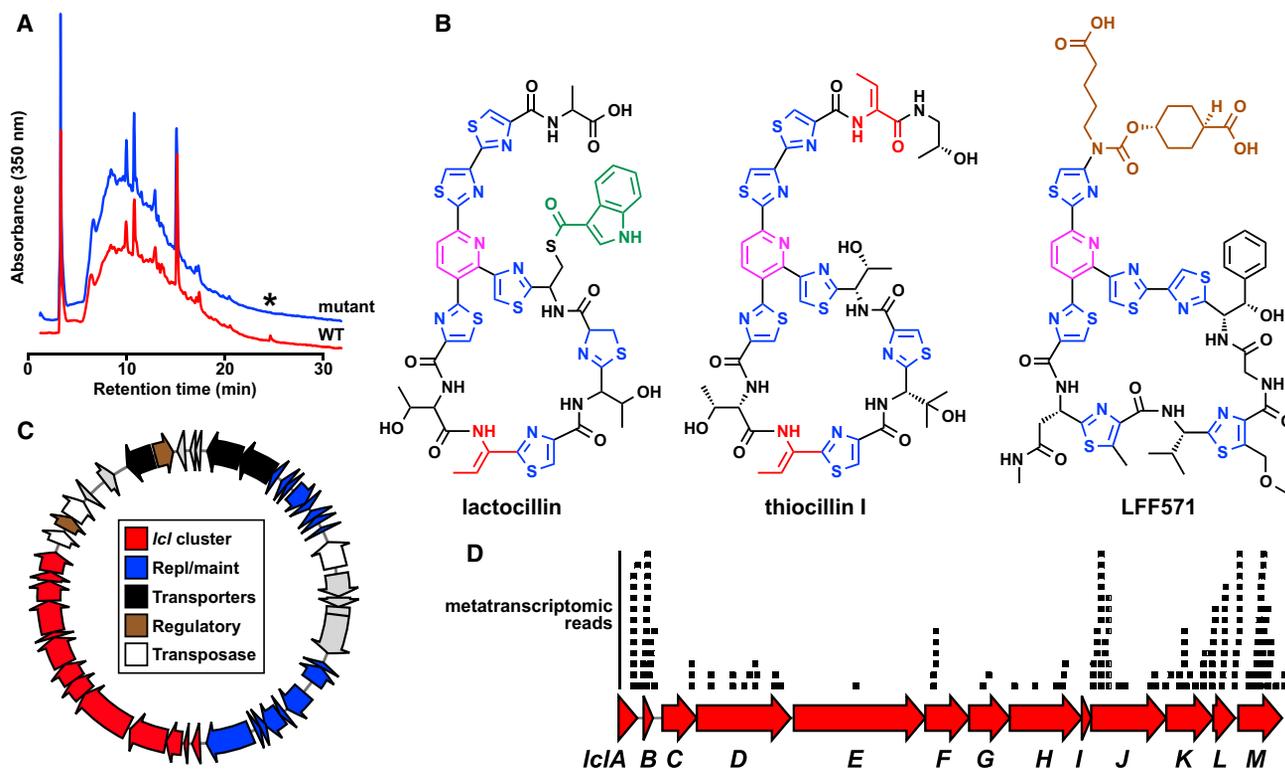
Here, we systematically identify BGCs in the genomes of human isolates, examine their representation in metagenomic and metatranscriptomic samples, identify widely distributed BGC families, and characterize the structure and activity of their small-molecule products. Our approach introduces a method for identifying and prioritizing BGCs for experimental characterization and can easily be generalized to other BGC families. The resulting database of BGCs that were detected in human samples (see also [Data S1](#)) will serve as a resource for future studies that aim to discover small-molecule-mediated interactions in the human microbiota.

Our approach has three important limitations. First, we rely on sequenced isolates of the human microbiota to predict the initial set of BGCs that we later use to recruit reads from whole-genome shotgun metagenomic samples. Although the ~2,400 bacterial genomes that we analyze here span all of the human-associated bacterial phyla and represent each major body site, they are biased toward gut and oral residents and easily cultivated species, for example.

Second, in examining the representation of BGCs in human metagenomic samples, we relied on the HMP whole-genome shotgun metagenomic data sets, which were sampled from fewer body sites and subsites than the larger 16S data set. Fewer samples were sequenced from skin than other body sites (27 versus at least 60 for all other body sites), and these samples were all obtained from one skin subsite (retroauricular crease). We anticipate that the number of skin-associated BGCs will



(legend on next page)



**Figure 6. Characterization of an Antibiotic, Lactocillin, from a Human Vaginal Isolate**

(A) HPLC analysis of organic extracts of cell pellets from wild-type (red) and *lclD* insertional mutant (blue) strains of *L. gasseri* JV-V03, monitored at 350 nm. An asterisk indicates the HPLC peak corresponding to lactocillin.

(B) Planar structure of lactocillin (see [Experimental Procedures](#) and see also [Figure S6](#) and [Data S2](#) for details about its purification and structural elucidation), the *Bacillus cereus* antibiotic thiocillin, and the clinical candidate LFF571 (Phase II, Novartis). Note the structural similarities among the three thiopeptides.

(C) Plasmid harboring the lactocillin BGC (see [Experimental Procedures](#) for details of the experimental closure of the circular plasmid). The lactocillin gene cluster (red) occupies ~30% of the plasmid; other elements on the plasmid include plasmid replication and maintenance genes (blue), transporters (black), transcription regulators (brown), and transposases and phage integrases (white).

(D) Raw oral metatranscriptomic reads were recruited to the *lcl* cluster using blastn and aligned using Geneious. Each black bar represents one read. Note that the precursor peptide *lclB* is among the most deeply covered genes in the *lcl* cluster, as anticipated for a RiPP pathway (see [Experimental Procedures](#) and see also [Figure S7](#) for metatranscriptomic analysis of the whole *lcl* plasmid and of *bgc65*).

rise as the number and diversity of whole-genome shotgun metagenomic data sets increases.

Third, we show that metatranscriptomic reads matching the *lcl* cluster exist in RNA-seq data from the human oral cavity, suggesting that the *lcl* cluster is expressed in humans. To prove conclusively that lactocillin is produced in humans, it would need to be directly detected in human vaginal or oral samples using sensitive analytical techniques. The activity of lactocillin in its native context will need to be further explored in colonization experiments in which a lactocillin producer is compared to an isogenic strain deficient in lactocillin production.

Landmark studies of gene clusters encoding catabolic pathways from the microbiota have shown the power of connecting genes to microbiome-related functions. [Sonnenburg et al. \(2010\)](#) and [Larsbrink et al. \(2014\)](#) have dissected the function of gene clusters from *Bacteroides* that catabolize fructans and xyloglucans, respectively. Both studies lend important insights into the role of specialized catabolic modules in competition for nutrient niches in the gut community, and they demonstrate that a mechanistic understanding of gene cluster function yields predictive power into how members of the gut community will respond to a change in diet. The current

**Figure 5. Thiopeptide BGCs Are Widespread in Isolates and Metagenomes of All Main Human Body Sites**

(A) Five thiopeptide BGCs from human isolates, eight thiopeptide BGCs from human metagenomes, and one thiopeptide BGC from a pig isolate are shown. The label of each BGC indicates its body site of origin.

(B) Precursor peptides corresponding to the thiopeptide BGCs shown in (A). Note that the precursor peptides fall into six subgroups of nearly identical sequences. The structural portion of the precursor peptide is shown in red (see also [Figure S5](#) for a phylogenetic analysis of thiopeptide BGCs).

(C) A heat map showing the representation and abundance of six oral thiopeptide BGCs in HMP metagenomic oral samples (see also [Table S2](#) for the quantification of all thiopeptide BGCs in all HMP samples). Note that, although each BGC shown in the heat map is well represented in the oral cavity, most samples harbor only one thiopeptide BGC.

**Table 1. Minimum Inhibitory Concentration of Lactocillin against Vaginal and Oral Pathogens and Commensals**

Strain	Phylum	Description	Lactocillin MIC (nM)
<i>Staphylococcus aureus</i>	Firmicutes	pathogen	42 <sup>a</sup>
<i>Escherichia coli</i>	Proteobacteria	commensal	NA <sup>b</sup>
<i>Enterococcus faecalis</i>	Firmicutes	pathogen	425
<i>Lactobacillus jensenii</i>	Firmicutes	vaginal commensal	NA <sup>b</sup>
<i>Lactobacillus gasseri</i> <sup>c</sup>	Firmicutes	vaginal commensal	NA <sup>b</sup>
<i>Lactobacillus crispatus</i>	Firmicutes	vaginal commensal	NA <sup>b</sup>
<i>Corynebacterium aurimucosum</i>	Actinobacteria	vaginal pathogen	42 <sup>a</sup>
<i>Finegoldia magna</i>	Firmicutes	vaginal pathogen	NA <sup>b</sup>
<i>Gardnerella vaginalis</i>	Actinobacteria	vaginal pathogen	212
<i>Streptococcus sanguinis</i>	Firmicutes	oral commensal	212
<i>Streptococcus sobrinus</i>	Firmicutes	oral commensal/pathogen	85
<i>Streptococcus mutans</i>	Firmicutes	oral commensal/pathogen	425

See [Experimental Procedures](#) for details. *Streptococcus sanguinis* and *Gardnerella vaginalis* were partially inhibited at 80 nM lactocillin but not completely until the next concentration tested (212 nM). *Streptococcus mutans* was partially inhibited at 212 nM but not completely until the next concentration tested (425 nM).

<sup>a</sup>Lowest concentration tested was 42 nM.

<sup>b</sup>Highest concentration tested was 425 nM.

<sup>c</sup>Note that the strain of *L. gasseri* tested was SV-16A-US, which is a different strain from the lactocillin producer (JV-V03).

state of knowledge of how BGCs in the human microbiota function is comparatively less well developed but holds great promise in yielding similar insights into microbe-host and microbe-microbe interactions.

The small-molecule products of BGCs are widely used in the clinic, and they constitute much of the chemical language of interspecies interactions. Our data highlight the fact that there exist hundreds of widely distributed BGCs of unknown function in the human microbiome, and they provide a template for future experimental efforts to discover biologically active small molecules from the microbiota (see also [Data S1](#) for a full data set of human-associated BGCs). These molecules represent a promising starting point for studying microbe-host interactions at the level of molecular mechanism and potentially a rich source of therapeutics.

## EXPERIMENTAL PROCEDURES

### Computational Analysis of BGCs from the Human Microbiome

Genome sequences of 2,430 bacterial strains isolated from humans were obtained from JGI-IMG, and BGCs were predicted using ClusterFinder ([Cimermancic et al., 2014](#)). The initial set of 44,000 putative BGCs was then analyzed

by antiSMASH ([Medema et al., 2011](#)), which classified a subset of 14,000 BGCs into a known category. A database containing the amino acid sequence of each gene in these 14,000 BGCs was constructed and queried against the processed (post-QC) reads of 752 HMP metagenomic samples using mblastx ([Davis et al., 2013](#)). Abundance scores of genes and gene clusters were calculated, and 3,118 BGCs were found to be present in at least one metagenomic sample (see also [Extended Experimental Procedures](#) and [Data S3](#) for details about abundance score calculations, thresholds, and additional computational analyses).

### Identification of Thiopeptide BGCs from HMP Metagenomic Data

A database containing the amino acid sequence of 24 homologs of TcIM (the enzyme responsible for pyridine ring formation in thiocillin) was constructed. This data set was then used as a query for tblastn searches against metagenomic assemblies of the 752 HMP samples. Contigs that generated hits to multiple thiopeptide genes were analyzed using antiSMASH, and results were verified manually. A thiopeptide BGC was called only when genes for all the essential posttranslational modifications were identified in addition to the precursor peptide (see also [Extended Experimental Procedures](#)).

### Generation of an Insertional Mutant in *L. gasseri* JV-V03 and Verification of Lactocillin Production

*L. gasseri* JV-V03 was cultivated in an anaerobic chamber (80% N<sub>2</sub>, 10% CO<sub>2</sub>, and 10% H<sub>2</sub>) in MRS broth (BD) at 37°C. A 1,000 bp fragment of *lciD* was PCR amplified and cloned into pKM082, a suicide vector harboring an erythromycin resistance gene. *L. gasseri* JV-V03 was transformed with this vector by electroporation, and transformants were selected on erythromycin and verified using PCR (see also [Extended Experimental Procedures](#) and [Table S3](#)). Cell pellets from the wild-type strain and insertional mutants were grown in 1 L MRS broth, extracted with methanol, and analyzed using HPLC and LC-MS monitoring at 350 and 220 nm.

### Large-Scale Production and Derivatization of Lactocillin

A 60 L culture of *L. gasseri* JV-V03 was grown for 2 days at 37°C. Cell pellets were harvested and extracted with 3 L methanol. Organic extracts were dried using rotary evaporation and fractionated on a C18 Sep Pak column using a step gradient of 20%, 40%, 60%, 80%, and 100% methanol in water. Lactocillin eluted in the 80% and 100% fractions, as determined by LC-MS. Semipurified lactocillin from the 80% fraction was methylated using TMS-diazomethane, and the methylated lactocillin product was purified by preparative HPLC (see also [Extended Experimental Procedures](#)).

### Structural Elucidation of Lactocillin Methyl Ester

The structure of purified lactocillin methyl ester was solved using a combination of MS experiments (HRMS, HRMS/MS, and HRMS<sup>n</sup>) and 1D and 2D NMR spectroscopy experiments (1D <sup>1</sup>HNMR, gCOSY, HSQC, HMBC, and ROESY). In addition, heavy isotope feeding experiments were performed to confirm the identity of the indolyl acyl group in the structure (see also [Extended Experimental Procedures](#)).

### Purification and MIC Determination of Lactocillin

Lactocillin was purified by preparative HPLC from the 100% Sep Pak fraction and quantified using HPLC. The minimal inhibitory concentration of purified lactocillin against a suite of commensal and pathogenic vaginal and oral isolates was tested using a range of concentrations from 42 to 425 nM. A vehicle-alone control was used in each experiment (see also [Extended Experimental Procedures](#)).

### Metatranscriptomic Analysis of Thiopeptide BGCs

A database containing the nucleotide sequences of all human-associated thiopeptide BGCs was generated. Raw Illumina reads from 38 oral metatranscriptomic samples were compared to this database using blastn, and hits with expectation values < 1e<sup>-10</sup> were recruited to the clusters using Geneious.

## SUPPLEMENTAL INFORMATION

Supplemental Information includes Extended Experimental Procedures, seven figures, three data files, and three tables and can be found with this article online at <http://dx.doi.org/10.1016/j.cell.2014.08.032>.

## AUTHOR CONTRIBUTIONS

M.S.D., C.J.S., R.G.L., and M.A.F. designed the research and analyzed the data, with substantial input from P.C., M.M., and J.C. M.S.D., C.J.S., and L.C.W.B. performed the experimental research. M.S.D., P.C., and J.M. performed the computational research. M.S.D. and M.A.F. wrote the manuscript.

## ACKNOWLEDGMENTS

We are indebted to Marnix Medema (Max Planck Institute for Marine Microbiology), Amrita Pati (JGI), and members of the Fischbach Group for helpful advice. We thank Krishna Parsawar and Chad Nelson (University of Utah) and Jeff Johnson (UCSF) for help with mass spectrometry experiments and Mark Kelly (UCSF) for help with NMR experiments. This work was supported by a Howard Hughes Medical Institute Predoctoral Fellowship (P.C.), NIH grant TW006634 (R.G.L.), a Medical Research Program Grant from the W.M. Keck Foundation (M.A.F.), a Fellowship for Science and Engineering from the David and Lucile Packard Foundation (M.A.F.), DARPA award HR0011-12-C-0067 (M.A.F.), Program for Breakthrough Biomedical Research (M.A.F.), and NIH grants OD007290, AI101018, AI101722, and GM081879 (M.A.F.). M.A.F. is on the scientific advisory boards of NGM Biopharmaceuticals and Warp Drive Bio.

Received: April 17, 2014

Revised: May 30, 2014

Accepted: August 20, 2014

Published: September 11, 2014

## REFERENCES

- An, D., Oh, S.F., Olszak, T., Neves, J.F., Avci, F.Y., Erturk-Hasdemir, D., Lu, X., Zeissig, S., Blumberg, R.S., and Kasper, D.L. (2014). Sphingolipids from a symbiotic microbe regulate homeostasis of host intestinal natural killer T cells. *Cell* **156**, 123–133.
- Baumann, H., Tzianabos, A.O., Brisson, J.R., Kasper, D.L., and Jennings, H.J. (1992). Structural elucidation of two capsular polysaccharides from one strain of *Bacteroides fragilis* using high-resolution NMR spectroscopy. *Biochemistry* **31**, 4081–4089.
- Brzuszkiewicz, E., Weiner, J., Wollherr, A., Thürmer, A., Hüpeden, J., Lomholt, H.B., Kilian, M., Gottschalk, G., Daniel, R., Mollenkopf, H.J., et al. (2011). Comparative genomics and transcriptomics of *Propionibacterium acnes*. *PLoS ONE* **6**, e21581.
- Chikindas, M.L., Novák, J., Driessen, A.J., Konings, W.N., Schilling, K.M., and Caufield, P.W. (1995). Mutacin II, a bactericidal antibiotic from *Streptococcus mutans*. *Antimicrob. Agents Chemother.* **39**, 2656–2660.
- Cimermancic, P., Medema, M.H., Claesen, J., Kurita, K., Wieland Brown, L.C., Mavrommatis, K., Pati, A., Godfrey, P.A., Koehrsen, M., Clardy, J., et al. (2014). Insights into secondary metabolism from a global analysis of prokaryotic biosynthetic gene clusters. *Cell* **158**, 412–421.
- Davis, C., Kota, K., Balchandapani, V., Gong, W., Abubucker, S., Becker, E., Martin, J., Wylie, K.M., Khetani, R., Hudson, M.E., et al. (2013). mBLAST: Keeping up with the sequencing explosion for (meta)genome analysis. *J. Data Mining Genomics Proteomics* **4**, 3.
- Engelhardt, K., Degnes, K.F., and Zotchev, S.B. (2010). Isolation and characterization of the gene cluster for biosynthesis of the thiopeptide antibiotic TP-1161. *Appl. Environ. Microbiol.* **76**, 7093–7101.
- Gajer, P., Brotman, R.M., Bai, G., Sakamoto, J., Schutte, U.M., Zhong, X., Koenig, S.S., Fu, L., Ma, Z.S., Zhou, X., et al. (2012). Temporal dynamics of the human vaginal microbiota. *Sci. Transl. Med.* **4**, 132ra152.
- Gevers, D., Knight, R., Petrosino, J.F., Huang, K., McGuire, A.L., Birren, B.W., Nelson, K.E., White, O., Methé, B.A., and Huttenhower, C. (2012). The Human Microbiome Project: a community resource for the healthy human microbiome. *PLoS Biol.* **10**, e1001377.
- Gevers, D., Kugathasan, S., Denson, L.A., Vázquez-Baeza, Y., Van Treuren, W., Ren, B., Schwager, E., Knights, D., Song, S.J., Yassour, M., et al. (2014). The treatment-naïve microbiome in new-onset Crohn's disease. *Cell Host Microbe* **15**, 382–392.
- Gonzalez, D.J., Lee, S.W., Hensler, M.E., Markley, A.L., Dahesh, S., Mitchell, D.A., Bandeira, N., Nizet, V., Dixon, J.E., and Dorrestein, P.C. (2010). Clostridiolysin S, a post-translationally modified biotoxin from *Clostridium botulinum*. *J. Biol. Chem.* **285**, 28220–28228.
- Hsiao, E.Y., McBride, S.W., Hsien, S., Sharon, G., Hyde, E.R., McCue, T., Codelli, J.A., Chow, J., Reisman, S.E., Petrosino, J.F., et al. (2013). Microbiota modulate behavioral and physiological abnormalities associated with neurodevelopmental disorders. *Cell* **155**, 1451–1463.
- Human Microbiome Project Consortium (2012a). A framework for human microbiome research. *Nature* **486**, 215–221.
- Human Microbiome Project Consortium (2012b). Structure, function and diversity of the healthy human microbiome. *Nature* **486**, 207–214.
- Just-Baringo, X., Albericio, F., and Álvarez, M. (2014). Thiopeptide antibiotics: retrospective and recent advances. *Mar. Drugs* **12**, 317–351.
- Kelly, W.L., Pan, L., and Li, C. (2009). Thiostrepton biosynthesis: prototype for a new family of bacteriocins. *J. Am. Chem. Soc.* **131**, 4327–4334.
- Kwan, J.C., Donia, M.S., Han, A.W., Hirose, E., Haygood, M.G., and Schmidt, E.W. (2012). Genome streamlining and chemical defense in a coral reef symbiosis. *Proc. Natl. Acad. Sci. USA* **109**, 20655–20660.
- Kwon, H.C., Kauffman, C.A., Jensen, P.R., and Fenical, W. (2006). Marinomycins A–D, antitumor-antibiotics of a new structure class from a marine actinomycete of the recently discovered genus “marinispora”. *J. Am. Chem. Soc.* **128**, 1622–1632.
- LaMarche, M.J., Leeds, J.A., Amaral, A., Brewer, J.T., Bushell, S.M., Deng, G., Dewhurst, J.M., Ding, J., Dzik-Fox, J., Gamber, G., et al. (2012). Discovery of LFF571: an investigational agent for *Clostridium difficile* infection. *J. Med. Chem.* **55**, 2376–2387.
- Larsbrink, J., Rogers, T.E., Hemsworth, G.R., McKee, L.S., Tausin, A.S., Spadiut, O., Kliner, S., Pudlo, N.A., Urs, K., Koropatkin, N.M., et al. (2014). A discrete genetic locus confers xyloglucan metabolism in select human gut Bacteroidetes. *Nature* **506**, 498–502.
- Lee, S.W., Mitchell, D.A., Markley, A.L., Hensler, M.E., Gonzalez, D., Wohlrab, A., Dorrestein, P.C., Nizet, V., and Dixon, J.E. (2008). Discovery of a widely distributed toxin biosynthetic gene cluster. *Proc. Natl. Acad. Sci. USA* **105**, 5879–5884.
- Lee, J.H., Chae, J.P., Lee, J.Y., Lim, J.S., Kim, G.B., Ham, J.S., Chun, J., and Kang, D.K. (2011). Genome sequence of *Lactobacillus johnsonii* PF01, isolated from piglet feces. *J. Bacteriol.* **193**, 5030–5031.
- Leimena, M.M., Ramiro-Garcia, J., Davids, M., van den Bogert, B., Smidt, H., Smid, E.J., Boekhorst, J., Zoetendal, E.G., Schaap, P.J., and Kleerebezem, M. (2013). A comprehensive metatranscriptome analysis pipeline and its validation using human small intestine microbiota datasets. *BMC Genomics* **14**, 530.
- Liao, R., Duan, L., Lei, C., Pan, H., Ding, Y., Zhang, Q., Chen, D., Shen, B., Yu, Y., and Liu, W. (2009). Thiopeptide biosynthesis featuring ribosomally synthesized precursor peptides and conserved posttranslational modifications. *Chem. Biol.* **16**, 141–147.
- Long, S.R. (2001). Genes and signals in the rhizobium-legume symbiosis. *Plant Physiol.* **125**, 69–72.
- Malcolmson, S.J., Young, T.S., Ruby, J.G., Skewes-Cox, P., and Walsh, C.T. (2013). The posttranslational modification cascade to the thiopeptide berninamycin generates linear forms and altered macrocyclic scaffolds. *Proc. Natl. Acad. Sci. USA* **110**, 8483–8488.
- Mazmanian, S.K., Liu, C.H., Tzianabos, A.O., and Kasper, D.L. (2005). An immunomodulatory molecule of symbiotic bacteria directs maturation of the host immune system. *Cell* **122**, 107–118.

- Mazmanian, S.K., Round, J.L., and Kasper, D.L. (2008). A microbial symbiosis factor prevents intestinal inflammatory disease. *Nature* 453, 620–625.
- McInerney, B.V., Gregson, R.P., Lacey, M.J., Akhurst, R.J., Lyons, G.R., Rhodes, S.H., Smith, D.R., Engelhardt, L.M., and White, A.H. (1991). Biologically active metabolites from *Xenorhabdus* spp., Part 1. Dithiopyrrolone derivatives with antibiotic activity. *J. Nat. Prod.* 54, 774–784.
- Medema, M.H., Blin, K., Cimermancic, P., de Jager, V., Zakrzewski, P., Fischbach, M.A., Weber, T., Takano, E., and Breitling, R. (2011). antiSMASH: rapid identification, annotation and analysis of secondary metabolite biosynthesis gene clusters in bacterial and fungal genome sequences. *Nucleic Acids Res.* 39 (Web Server issue), W339–W346.
- Nelson, K.E., Weinstock, G.M., Highlander, S.K., Worley, K.C., Creasy, H.H., Wortman, J.R., Rusch, D.B., Mitreva, M., Sodergren, E., Chinwalla, A.T., et al.; Human Microbiome Jumpstart Reference Strains Consortium (2010). A catalog of reference genomes from the human microbiome. *Science* 328, 994–999.
- Nougayrède, J.P., Homburg, S., Taieb, F., Boury, M., Brzuszkiewicz, E., Gottschalk, G., Buchrieser, C., Hacker, J., Dobrindt, U., and Oswald, E. (2006). *Escherichia coli* induces DNA double-strand breaks in eukaryotic cells. *Science* 313, 848–851.
- Oh, D.C., Poulsen, M., Currie, C.R., and Clardy, J. (2009). Dentigerumycin: a bacterial mediator of an ant-fungus symbiosis. *Nat. Chem. Biol.* 5, 391–393.
- Qin, J., Li, R., Raes, J., Arumugam, M., Burgdorf, K.S., Manichanh, C., Nielsen, T., Pons, N., Levenez, F., Yamada, T., et al.; MetaHIT Consortium (2010). A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* 464, 59–65.
- Ravel, J., Gajer, P., Abdo, Z., Schneider, G.M., Koenig, S.S., McCulle, S.L., Karlebach, S., Gorle, R., Russell, J., Tacket, C.O., et al. (2011). Vaginal microbiome of reproductive-age women. *Proc. Natl. Acad. Sci. USA* 108 (Suppl 1), 4680–4687.
- Segata, N., Waldron, L., Ballarini, A., Narasimhan, V., Jousson, O., and Huttenhower, C. (2012). Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat. Methods* 9, 811–814.
- Smillie, C.S., Smith, M.B., Friedman, J., Cordero, O.X., David, L.A., and Alm, E.J. (2011). Ecology drives a global network of gene exchange connecting the human microbiome. *Nature* 480, 241–244.
- Sonnenburg, E.D., Zheng, H., Joglekar, P., Higginbottom, S.K., Firbank, S.J., Bolam, D.N., and Sonnenburg, J.L. (2010). Specificity of polysaccharide use in intestinal bacteroides species determines diet-induced microbiota alterations. *Cell* 141, 1241–1252.
- Turnbaugh, P.J., Hamady, M., Yatsunenko, T., Cantarel, B.L., Duncan, A., Ley, R.E., Sogin, M.L., Jones, W.J., Roe, B.A., Affourtit, J.P., et al. (2009). A core gut microbiome in obese and lean twins. *Nature* 457, 480–484.
- Tzianabos, A.O., Onderdonk, A.B., Rosner, B., Cisneros, R.L., and Kasper, D.L. (1993). Structural features of polysaccharides that induce intra-abdominal abscesses. *Science* 262, 416–419.
- Wieland Brown, L.C., Acker, M.G., Clardy, J., Walsh, C.T., and Fischbach, M.A. (2009). Thirteen posttranslational modifications convert a 14-residue peptide into the antibiotic thiocillin. *Proc. Natl. Acad. Sci. USA* 106, 2549–2553.
- Wieland Brown, L.C., Penaranda, C., Kashyap, P.C., Williams, B.B., Clardy, J., Kronenberg, M., Sonnenburg, J.L., Comstock, L.E., Bluestone, J.A., and Fischbach, M.A. (2013). Production of  $\alpha$ -galactosylceramide by a prominent member of the human gut microbiota. *PLoS Biol.* 11, e1001610.
- Wyatt, M.A., Wang, W., Roux, C.M., Beasley, F.C., Heinrichs, D.E., Dunman, P.M., and Magarvey, N.A. (2010). *Staphylococcus aureus* nonribosomal peptide secondary metabolites regulate virulence. *Science* 329, 294–296.
- Yoshimoto, S., Loo, T.M., Atarashi, K., Kanda, H., Sato, S., Oyadomari, S., Iwakura, Y., Oshima, K., Morita, H., Hattori, M., et al. (2013). Obesity-induced gut microbial metabolite promotes liver cancer through senescence secretome. *Nature* 499, 97–101.
- Yu, Y., Duan, L., Zhang, Q., Liao, R., Ding, Y., Pan, H., Wendt-Pienkowski, E., Tang, G., Shen, B., and Liu, W. (2009). Nosiheptide biosynthesis featuring a unique indole side ring formation on the characteristic thiopeptide framework. *ACS Chem. Biol.* 4, 855–864.
- Zou, Y., Yin, H., Kong, D., Deng, Z., and Lin, S. (2013). A trans-acting ketoreductase in biosynthesis of a symmetric polyketide dimer SIA7248. *ChemBioChem* 14, 679–683.

## EXTENDED EXPERIMENTAL PROCEDURES

### Computational Analysis of BGCs from the Human Microbiome

#### Details of ClusterFinder Analysis and antiSMASH Classification

ClusterFinder was run on a data set of 2,430 genomes obtained from JGI-IMG and selected as having the host as “human,” including finished and draft genome sequences. We defined a BGC as a set of consecutive genes with at least one Pfam domain with posterior BGC probability of  $> 0.4$  (Cimermancic et al., 2014). BGC candidates that were shorter than 2 kb and that did not contain a single class-specific domains were filtered out. The remaining putative BGCs were then classified using antiSMASH as described previously (Cimermancic et al., 2014).

#### Representation of BGCs in HMP Samples

A database containing the protein sequences of BGCs predicted from the 2,430 genomes was constructed. Processed metagenomic reads of 752 HMP samples were used as a query for *tblastx* searches against the database of BGCs, using a filter of an *e*-value of  $1e-05$  (Davis et al., 2013). Only the top gene hit per read was selected. A custom Python script was used to calculate abundance scores per gene per sample. Because each BGC is composed of multiple genes, a calculation method was implemented to ensure accuracy in predicting the frequency and abundance of each BGC in a metagenomic sample. First, a list of the Pfams in the BGCs predicted by ClusterFinder and antiSMASH was generated, then all non-biosynthetic Pfams that are commonly found in BGCs were removed manually (e.g., transporters, transcriptional regulators, transposases, etc.) (see also [Data S1](#) for a list of excluded Pfams). A BGC was considered present in a metagenomic sample when at least 50% of the biosynthetic genes in it are detected at an abundance score  $\geq 10$ . When a BGC is present, it is assigned an “overall abundance score,” which is the average of all abundance scores of its genes that were detected in the sample (and passed the individual-gene-score threshold of 10). 1000 was considered the maximum abundance score of a given gene, and any values  $> 1000$  were given a score of 1000. To assess the accuracy of this method of score calculation, 36 BGC/sample pairs that gave scores ranging from 0-1000 were analyzed manually. A database was constructed containing the nucleotide sequence of the BGC in question, then metagenomic reads from the paired sample were used as a query in a *blastn* search using the default parameters. Reads that matched a certain BGC were then recruited to it using Geneious with the following parameters (minimum overlap: 50 bp, minimum percent identity at overlap: 90%, and maximum percentage of mismatch per read: 20%). Coverage plots were generated based on reads that pass these thresholds, and evaluated manually in comparison to the scores that were given (see also [Data S3](#) for selected examples of this verification analysis).

A matrix was generated with 752 HMP samples as rows and 3,118 BGCs (using the thresholds described above) as columns (see [Data S1](#) for the parent matrix). Hierarchical clustering was performed using MEV (Saeed et al., 2003), Pearson correlation or Pearson uncentered as the distance metric, and average linkage method.

#### Details of “Human” and “Nonhuman” Comparison

3,091 BGCs that were present in at least one metagenomic sample and 9,681 BGCs from a set of 2,252 of non-human organisms were classified using antiSMASH into PKS, NRPS, RiPPs, Saccharide, NRPS-Independent siderophores, terpenes and hybrids. Counts of each class of BGCs in the two groups (Human and non-Human) were bootstrapped with 20% of the number of BGCs per group being resampled 1,000 times. Bacterial phyla of source organisms from Human and non-Human sets were compared in a similar manner (see also [Figure S2](#)).

#### Analysis of “Common” BGCs, Body Site Overlap, and Average Number of BGCs per Body Site

At each of the five main body sites, the frequency of each of the 3,118 BGCs was calculated. This frequency was measured as the percentage of HMP samples at a body site that harbored a certain BGC. For this purpose, when samples from multiple visits of the same subject existed, only samples from the first visit were analyzed to prevent skewing of the data toward a particular subject. BGCs that were found in  $> 50\%$  of the samples at a certain body site were considered “common” (see also [Table S1](#)). Using a similar analysis, BGCs were identified that were present at two or three body sites (at  $\geq 20\%$  of the samples from each of the body sites), and Venn diagrams were generated to display the extent of overlap between different body sites (see also [Figure S2](#)). The total number of BGCs per sample was calculated, and the mean and standard deviation were calculated for each body site and for each class of BGCs at a certain body site (see also [Table S1](#)).

#### Details of Phylogenetic Analysis for the Gut Family of NRPS BGCs

“Common” BGCs in the HMP stool samples contained at least two members of a family of related NRPS BGCs (*bgc45* and *bgc52*). The large multimodular NRPS genes of these clusters were used as queries for *blastp* searches against the NCBI protein database to identify related BGCs. In total, 35 related BGCs were identified, and an alignment was generated using the large NRPS gene in the clusters (ClustalX) for all except one (*Bacteroides barnesiae* DSM 18169, because the large NRPS gene was split between two open reading frames). A maximum parsimony phylogenetic tree (MEGA5) was then generated from this alignment implementing a bootstrap test of 1000 replicates.

#### Details of Computational Thiopeptide BGC Analysis

A database containing the amino acid sequence of 24 homologs of TcIM (the enzyme responsible for pyridine ring formation in thiocillin) was constructed. This data set was then used as a query for *blastp* searches against the NCBI protein database to identify thiopeptide BGCs in sequenced genomes from the human microbiome, which led to the identification of *bgc56*, *bgc57*, *bgc60*, *bgc63* and *bgc66*. Identified BGCs were then added to the same database, which was then used in *tblastn* searches against metagenomic assemblies of the 752 HMP samples. Contigs that generated hits to multiple thiopeptide genes were analyzed using

antiSMASH, and results were verified manually. A candidate thiopeptide BGC was predicted only when genes encoding enzymes for all the essential posttranslational modifications were identified (two lantibiotic dehydratase homologs; one YcaO homolog and one MbcC homolog; and one TcIM homolog for dehydrobutyryne, thiazole, and pyridine-ring formation, respectively) in addition to the precursor peptide. This method led to the identification of *bgc58*, *bgc59*, *bgc62*, *bgc64*, *bgc65*, *bgc67*, and *bgc68* (see also Table S2). Amino acid sequences of the longer lantibiotic dehydratase (TcIK homolog) and the pyridine-ring-forming enzyme (TcIM homolog) genes from all 14 thiopeptide BGCs were aligned using ClustalX. A maximum likelihood phylogenetic tree was calculated using MEGA5 and a bootstrap test of 1,000 replicates (see also Figure S5 for a phylogenetic tree of TcIK homologs, and note that similar topology was obtained using TcIM homologs).

#### Details of Metatranscriptomic Analysis of Thiopeptide BGCs

A database containing the nucleotide sequences of all human-associated thiopeptide BGCs was generated. Raw Illumina reads from 38 oral metatranscriptomic samples obtained from the HMP DACC were compared to this database using blastn, and hits with expectation values  $< 1e^{-10}$  were recruited to the clusters and displayed using Geneious. Almost every gene in the *bgc65* and *bgc66* (*lcl*) clusters was covered by metatranscriptomic reads in at least one sample (see also Figure S7). Other thiopeptide BGCs had only sporadic coverage of genes and were not considered “transcribed” in this analysis. It is important to note that the bacterial community in supragingival plaque is very diverse, and the strain harboring one of these may have been a minor member of the community; in this case, we would not expect to find every gene in the cluster covered by metatranscriptomic reads. The same analysis was also done on the whole plasmid harboring the *bgc66* (*lcl*) cluster, and almost every gene on the plasmid was covered by metatranscriptomic reads, indicating that the same plasmid most likely harbors the lactocillin BGC in the oral cavity as well as in the vaginal community (see also Figure S7). The species carrying the plasmid could not be identified, however, due to the complexity of the supragingival plaque microbial community.

### Experimental Characterization of Lactocillin

#### Generation of an Insertional Mutant in *bgc66*, Termed *lcl*, and Closure of the Circular Plasmid

*L. gasseri* JV-V03 was obtained from BEI Resources (HM-104), and was cultivated in an anaerobic chamber (80% N<sub>2</sub>, 10% CO<sub>2</sub> and 10% H<sub>2</sub>) in MRS broth (BD) at 37°C. Genomic DNA was extracted using Zymo Fungal/Bacterial DNA Mini Kit. To generate an insertional mutant, a ~1000 bp fragment of *lclD* was PCR-amplified (using primers LCW080 and LCW081, see also Table S3) and cloned into pKM082, a suicide vector harboring an erythromycin resistance gene (provided gently by David Rudner, Harvard Medical School). *L. gasseri* JV-V03 was transformed by electroporation with this vector using a previously developed protocol (Tangney et al., 1994), and transformants were selected on MLS (1000X MLS: 100 mg erythromycin (EM) and 2.5 g lincomycin (LM) in 100 ml 50% EtOH) and verified using PCR (using primers LCW090/LCW091 and LCW092/LCW093, see also Table S3). The contig harboring *bgc66* (*lcl*) (41,479 bp, GenBank: ACGO02000004.1) contained multiple plasmid elements, but was not assembled into a circular element in the final assembly of *L. gasseri* JV-V03. To verify that it is indeed a circular plasmid, primers were designed that faced outward on both ends of the contig (Plas-Gasseri-F and Plas-Gasseri-R, see also Table S3). PCR was performed using these primers and genomic DNA from *L. gasseri* JV-V03 as template. Gel electrophoresis showed one major band of 600 bp. This band was excised from the gel, gel-purified using QIAGEN Gel Extraction Kit, and sequenced using the PCR primers (Plas-Gasseri-F and Plas-Gasseri-R). Assembly of these reads with the rest of the contig led to the closure of the plasmid into one circular element (41,709 bps).

#### General Chemical Procedures

Solvents used for HPLC and LC-MS chromatography were HPLC grade and were used without further purification. NMR spectroscopy and high-resolution mass spectrometry data were obtained for all compounds. A single quadrupole LC-MS (Agilent 6130) was used for the analysis of organic extracts. NMR spectra were acquired on a 600 MHz Varian spectrometer equipped with a 5 mm HCN triple resonance cryoprobe and a 800 MHz Bruker Avance spectrometer and referenced to residual solvent proton and carbon signals. HRMS data were acquired using University of Utah Mass Spectrometry Facility’s multimode electrospray ionization (ESI) Fourier transform mass spectrometer (FTMS) or an Orbitrap Velos mass spectrometer.

#### Chemical Extraction and Wild-Type/Insertional Mutant Comparison

1 L of wild-type *L. gasseri* JV-V03 and 1 L of the *lclD* insertional mutant were grown anaerobically for two days at 37°C in MRS broth and MRS broth + 1X MLS (1000X MLS: 100 mg erythromycin (EM) and 2.5 g lincomycin (LM) in 100 ml 50% EtOH), respectively. Cell pellets were harvested and extracted twice with 200 ml methanol. After rotary evaporation, organic extracts of cell pellets were compared using HPLC and LC-MS monitoring at 220 and 350 nm. The major difference between the two extracts was a single peak that had an absorbance profile similar to thiocillin (Figures 6 and S6A) and that corresponded to a single mass ion of  $m/z$  1224 [M+H]<sup>+</sup> (see also Figure S6B).

#### Growth Conditions/Time Point Optimization

To optimize the production of lactocillin by *L. gasseri* JV-V03, multiple growth conditions (aerobic and anaerobic), multiple media (MRS unconditioned, MRS pH 4, MRS pH 8, 0.5X MRS) and multiple time points (days 1–8) were tested. The only noticeable difference was observed in the time course analysis: lactocillin is produced at highest titer at day 2 (70 µg/l) and the yield is dramatically reduced afterward (see also Figure S6C). From this point onward, all extractions were performed after two days of growth.

#### Solid-Phase Initial Purification of Lactocillin

The crude organic extract was subjected to solid phase extraction using an Agilent BondElut C<sub>18</sub> cartridge (60 g) and eluted using a MeOH/H<sub>2</sub>O step gradient (500 ml; 10% MeOH, 20% MeOH, 40% MeOH, 60% MeOH, 80% MeOH, 100% MeOH) to afford six

fractions. The 10%, 20%, 40% MeOH fraction were discarded and the remaining three fractions were concentrated to dryness in vacuo. Aliquots of each fraction were solubilized in 50% MeOH and the presence of lactocillin was verified by LC-MS.

#### Derivatization of Lactocillin

Based on LC-MS quantification, the majority of lactocillin eluted in the 80% and 100% fractions. The 80% fraction was then partitioned into three separate vials and dried under N<sub>2</sub> and placed under vacuum overnight for subsequent chemical derivatizations. To each vial, 8 ml of anhydrous benzene, 2 ml of anhydrous methanol and 2 ml of 2.0 M (trimethylsilyl)diazomethane in hexane (CAS# 18107-18-1, Sigma) were added. The reaction was stirred for 16 hr at room temperature under an argon atmosphere and concentrated to dryness under a stream of N<sub>2</sub>. Reaction completion was verified by LC-MS (see also [Figure S6D](#)).

#### Identification and Purification of Lactocillin Methyl Ester

The crude reaction mixtures were then purified by RP-HPLC (Phenomenex Kinetex XB-C18, 100x4.6 mm, 2.6 μm, 100 Å) using a gradient of ACN:H<sub>2</sub>O with 0.02% formic acid (40%–53% over 17 min, 1 ml min<sup>-1</sup>, t<sub>R</sub> = 11.2 min) monitoring UV wavelength of 350 nm to yield lactocillin methyl ester (~500 μg).

#### NMR Characterization of Lactocillin Methyl Ester

NMR data for lactocillin methyl ester were acquired in DMSO-*d*<sub>6</sub> in a 3 mm NMR tube. One-dimensional and two-dimensional homonuclear experiments (<sup>1</sup>H, gCOSY, TOCSY) were performed on a 600 MHz Varian spectrometer equipped with a 5 mm HCN triple resonance cryoprobe. Additional two-dimensional data (HSQC, HMBC, ROESY) were collected on a 800 MHz Bruker Avance spectrometer.

#### Structural Elucidation of Lactocillin Methyl Ester

We started the structural elucidation process by analysis of the (+) HR-ESI Orbitrap data (observed [M+H]<sup>+</sup> at *m/z* 1238.16916, calculated for [M+H]<sup>+</sup> at *m/z* 1238.16916, Δ ppm = - 0.02) (see also [Figure S6E](#)), which was consistent with the empirical formula C<sub>52</sub>H<sub>47</sub>N<sub>13</sub>O<sub>10</sub>S<sub>7</sub>. Based on a computational analysis of the lactocillin BGC and the UV absorbance profile of lactocillin (see also [Figure S6A](#)), we predicted it would have a similar core structure to other 26-membered thiopeptides in this class (e.g., thiocillin) in which the C-terminal portion of the precursor peptide (H<sub>2</sub>N-SCTTCTCCSCCA-COOH) has undergone a series of dehydration, cyclization, and pyridine ring-forming posttranslational modifications.

We used one- and two-dimensional NMR experiments to elucidate the structure of this molecule, and to analyze any further modification of residues that would not be evident by analysis of the biosynthetic gene cluster. The one-dimensional <sup>1</sup>H NMR data indicated the presence of five singlet signals between δ 8.22–8.59 ppm, indicating the presence of five thiazoles in the molecule. We also observed an additional seven aromatic signals. Two of these signals were coupled aromatic doublets and were assigned as being attached to the pyridine ring based on their diagnostic chemical shifts and comparison to similar compounds (Position 14 8.39, 141.1; Position 15 8.27, 118.5, see also [Figure S6F](#)).

Analysis of two-dimensional homonuclear data including COSY and TOCSY revealed the presence of an additional six spin systems, corresponding to two threonine residues, a thiazoline, a dehydrobutyrine, a cysteine and an alanine. In addition, a key HMBC correlation allowed for the placement of the methyl ester on the C terminus of the alanine residue (see also [Figure S6F](#) and [Data S2](#)). These six spin systems, in addition to the thiazoles and the pyridine ring, accounted for all of the residues expected to be in the core sequence of the thiopeptide. However, the remaining chemical formula of C<sub>9</sub>H<sub>6</sub>NO had yet to be accounted for. Analysis of the two dimensional NMR data (COSY, TOCSY, HSQC, and HMBC) revealed that the five aromatic <sup>1</sup>H NMR signals that had yet to be assigned constituted an indole-3-carboxylic acid moiety (see also [Data S2](#)). This presence of the indole fragment was further confirmed by the incorporation of L-tryptophan-*d*<sub>5</sub> (indole-*d*<sub>5</sub>) into the structure by feeding experiments (see details in [Extended Experimental Procedures](#) and [Figure S6G](#)). We hypothesized that this fragment was attached to the core structure through a thioester linkage to the sole unmodified cysteine residue. In order to verify this hypothesis and determine the position of the modified cysteine and threonine residues, we performed extensive MS<sup>n</sup> experiments using a HR Orbitrap Velos (see also [Figure S6H](#)). Two key fragmentation events confirmed the attachment of the indole to the uncyclized cysteine through a thioester linkage. The first corresponded to fragmentation of the thioester bond, resulting in a loss of 143.0365 (consistent with the loss of the C<sub>9</sub>H<sub>6</sub>NO indole fragment, calculated 143.0371, Δ ppm = 4). In addition, we observed desulfonation of the cysteine after this loss (see also [Data S2](#)). Key fragmentations corresponding to the loss of one residue at a time in the macrocycle enabled the sequence of modified residues to be determined (in particular, the indolyl-S-cysteine at position 8 and thiazoline at position 7) and corroborated the overall residue order that was predicted by the core peptide sequence (see also [Data S2](#) for detailed HR mass spectrometry fragmentation analysis). The fragmentation pattern of lactocillin methyl ester was largely reminiscent of that of thiocillin and derivatives thereof ([Acker et al., 2009](#); [Bowers et al., 2012](#)).

#### Labeling Experiments

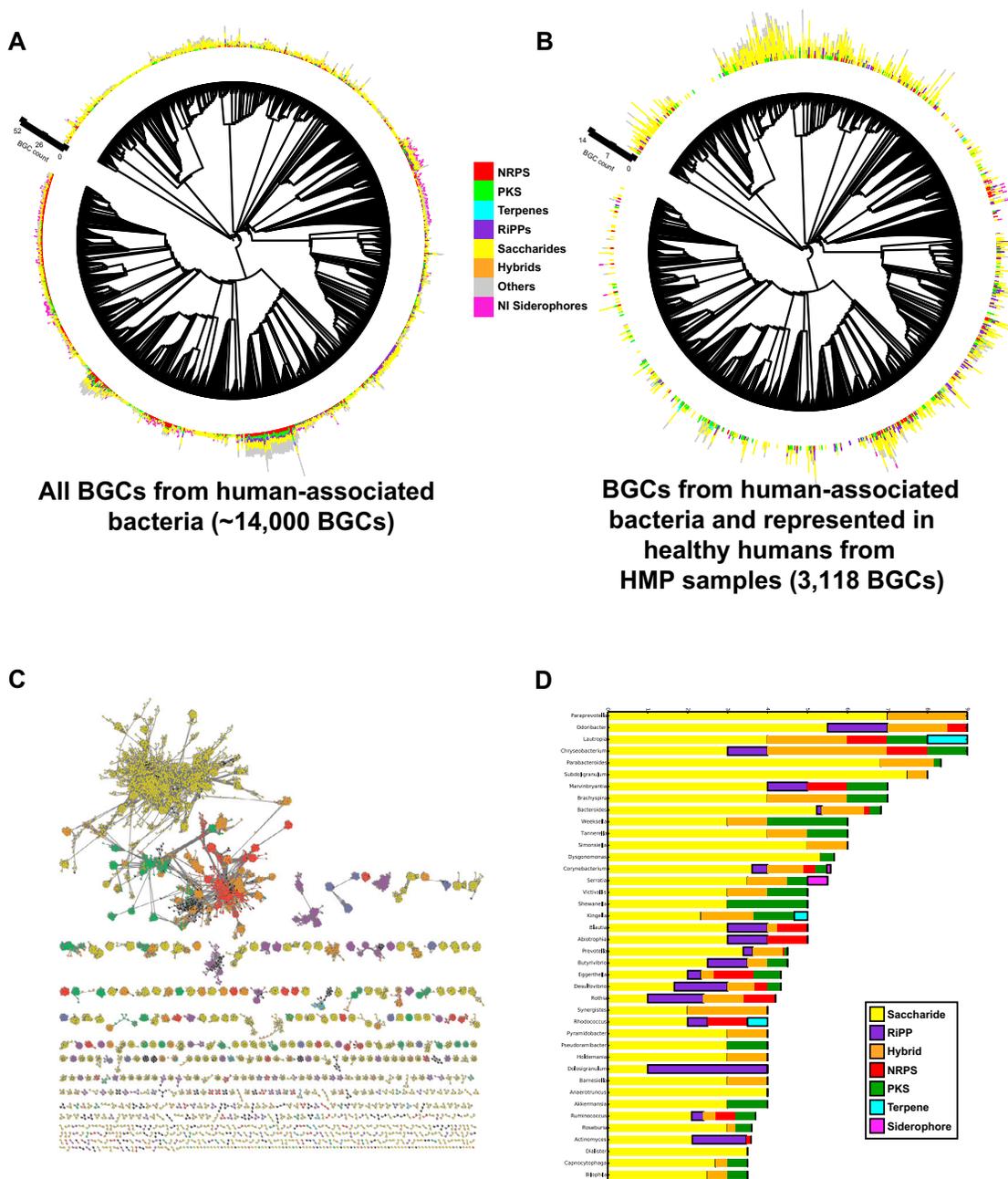
To determine the biosynthetic origin of the indole moiety in lactocillin, and further support the structure determined by NMR, heavy isotope variants of phenylalanine, tyrosine, tryptophan and alanine (as a positive control) were obtained (Cambridge Isotopes, CAS numbers 30811-19-9, 63546-27-0 and 56253-90-8). Amino acids were dissolved in MRS broth to a concentration of 20 mM, then sterile filtered using a 0.2 μm filter. 1 ml of an overnight *L. gasseri* JV-V03 culture was then added to 50 ml of MRS broth containing the heavy isotope amino acid and cultures were grown anaerobically for two days at 37°C. Cells were then harvested and extracted using 5 ml of methanol, and then analyzed using LC-MS. Only in the case of alanine-*d*<sub>3</sub> and tryptophan-*d*<sub>5</sub> was a difference in the isotope pattern of lactocillin observed, indicating that both alanine and tryptophan were incorporated into lactocillin, while phenylalanine and tyrosine were not ([Figure S6G](#)).

### **Purification of Native Lactocillin from the 100% Sep Pak Fraction and MIC Determination**

Lactocillin was purified by preparative HPLC from the 100% Sep Pak fraction (using a Phenomenex Luna 5  $\mu\text{m}$  C18 100 Å column, 250 X 10 mm) and quantified using HPLC. To determine the minimum inhibitory concentration of lactocillin against a panel of strains, four concentrations were tested for each strain (42.5, 85, 212.5 and 425 nM) in addition to a vehicle-only control (methanol). Lactocillin samples were prepared by dissolving the appropriate amount of lactocillin in methanol and adding it to 970  $\mu\text{l}$  of the appropriate growth medium (MRS for *Lactobacillus* sp., LB for *S. aureus* and *E. coli*, BHI for *Streptococcus* spp. and *C. aurimucosum*, and Casman's medium base with 5% sheep blood for *G. vaginalis*). The total volume of methanol added in all cases was 20  $\mu\text{l}$ . After adding 20  $\mu\text{l}$  of lactocillin or methanol to the medium, 10  $\mu\text{l}$  of an appropriate dilution of an overnight culture of the test organism was added. Cultures were then incubated anaerobically at 37°C for 48 hr, and growth was evaluated visually in the different concentrations of lactocillin in comparison to the vehicle-only control.

### **SUPPLEMENTAL REFERENCES**

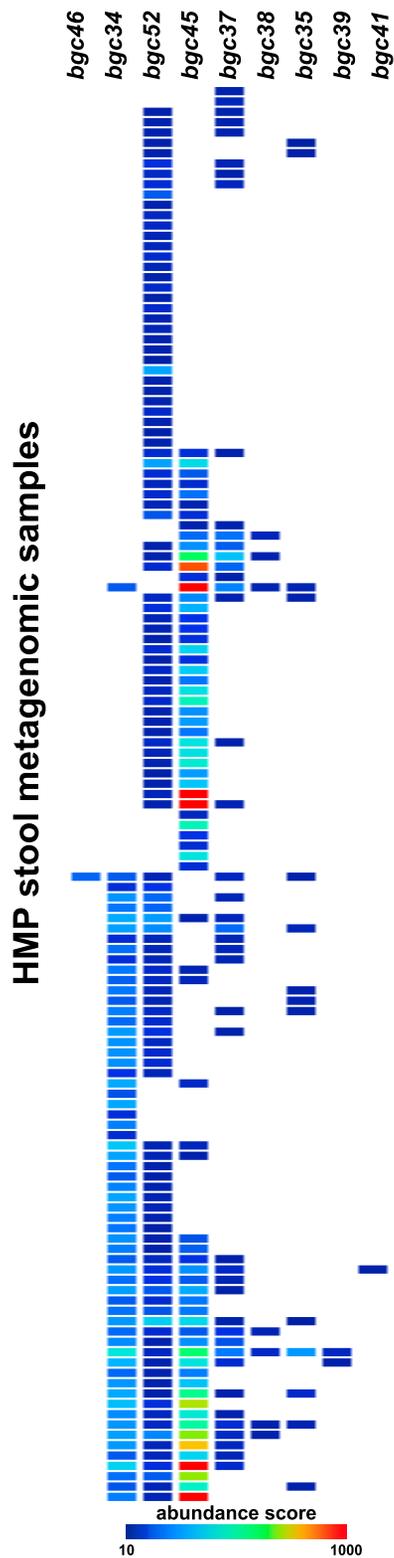
- Acker, M.G., Bowers, A.A., and Walsh, C.T. (2009). Generation of thiocillin variants by prepeptide gene replacement and in vivo processing by *Bacillus cereus*. *J. Am. Chem. Soc.* *131*, 17563–17565.
- Bowers, A.A., Acker, M.G., Young, T.S., and Walsh, C.T. (2012). Generation of thiocillin ring size variants by prepeptide gene replacement and in vivo processing by *Bacillus cereus*. *J. Am. Chem. Soc.* *134*, 10313–10316.
- Letunic, I., and Bork, P. (2007). Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* *23*, 127–128.
- Saeed, A.I., Sharov, V., White, J., Li, J., Liang, W., Bhagabati, N., Braisted, J., Klapa, M., Currier, T., Thiagarajan, M., et al. (2003). TM4: a free, open-source system for microarray data management and analysis. *Biotechniques* *34*, 374–378.
- Smoot, M.E., Ono, K., Ruscheinski, J., Wang, P.L., and Ideker, T. (2011). Cytoscape 2.8: new features for data integration and network visualization. *Bioinformatics* *27*, 431–432.
- Tangney, M., Diderichsen, B., and Priest, F.G. (1994). A method for electrotransformation of *Bacillus licheniformis* NCIB 6346 by plasmid DNA. *Biotechnol. Tech.* *8*, 463–466.



**Figure S1. Overview of BGCs Predicted in Human-Associated Isolates, Related to Figure 1**

- (A) Phylogenetic tree representing 2,430 analyzed genomes, with the corresponding number and class of predicted BGCs shown as bars (total of ~14,000 BGCs). The tree was generated from 16S rDNA sequences downloaded from JGI-IMG and the final figure was generated using iTOL (Letunic and Bork, 2007).
- (B) The same tree showing only the 3,118 BGCs detected in at least one HMP metagenomic sample.
- (C) A network view of all BGCs predicted from human-associated bacteria. Edges were calculated based on similarities of Pfam compositions between a pair of BGCs, with a cutoff of 0.4 (Cimermanic et al., 2014). The network was generated using Cytoscape (Smoot et al., 2011). The color code is indicated in the legend in (D).
- (D) Average number and class of BGCs detected in HMP samples by genus. The 40 genera with the highest average number of BGCs are shown.





**Figure S3. A Family of Gut NRPS Is Widely Distributed in HMP Stool Samples, Related to Figure 3**

A heat map showing the representation of nine members of the gut NRPS family in HMP stool samples (samples shown here represent 92% of the HMP stool samples).

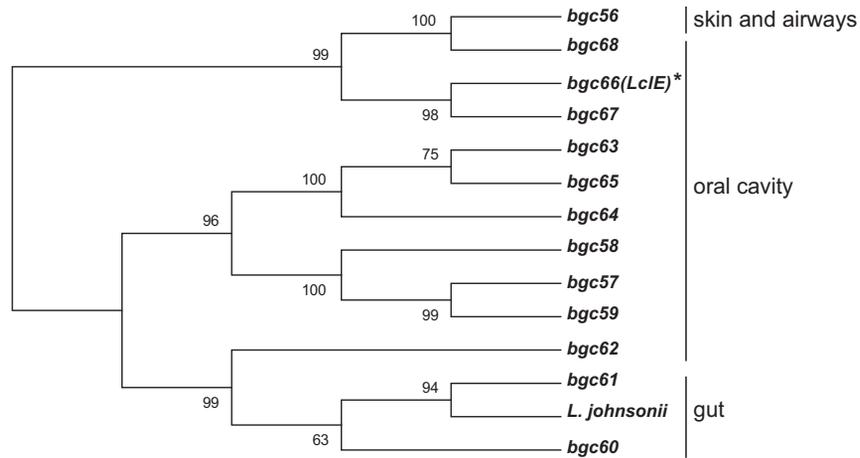


---

**Figure S4. Bacteroidetes Saccharide BGCs Are Variable in Metagenomic Stool Samples with Similar Microbial Composition, Related to Experimental Procedures**

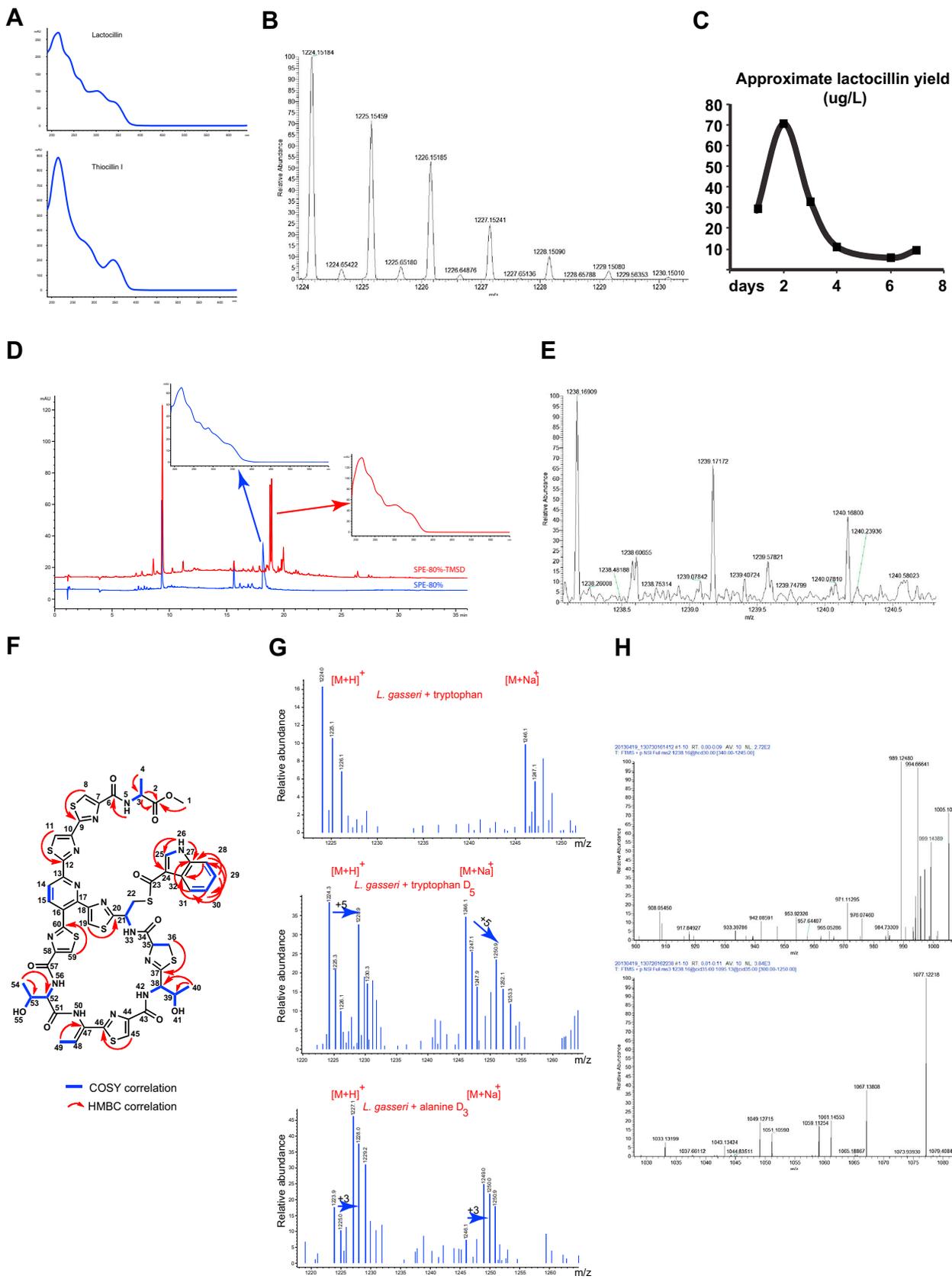
(A) Ten HMP stool samples that have similar Bacteroidetes compositions based on their MetaPhlAn (Segata et al., 2012) classifications were selected (SRS015369, SRS012902, SRS023176, SRS013158, SRS024388, SRS016989, SRS016267, SRS013215, SRS058723, and SRS049900). Only members of the phylum Bacteroidetes are shown in the heat map displayed here.

(B) A heat map showing the representation of saccharide BGCs predicted from members of the phylum Bacteroidetes and detected in the same ten samples shown in A. Although some of the BGCs displayed here show similar patterns of abundance among samples, a large subset varies substantially across the sample set.



**Figure S5. Phylogenetic Analysis of the Thiopeptide BGCs Described Here, Related to Figure 5**

A maximum likelihood phylogenetic tree of the lantibiotic dehydratase enzyme of all thiopeptide BGCs described here, constructed using MEGA5. The numbers on the branches indicate the percentage of times this topology was reached in a bootstrap test of 1,000 replicates. A similar tree was obtained using TcIM homologs (responsible for pyridine ring formation).

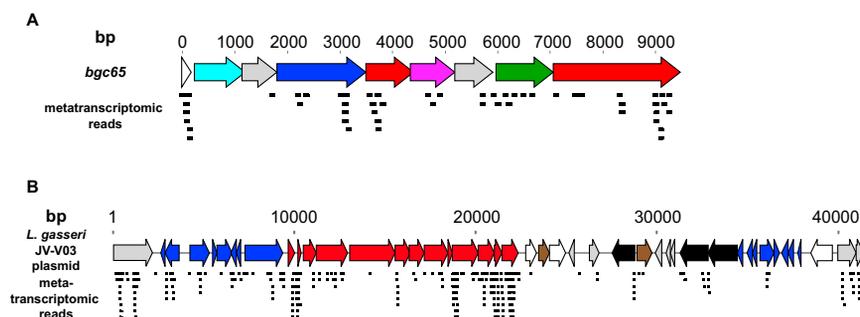


(legend on next page)

---

**Figure S6. Experimental Characterization of Lactocillin, Related to Figure 6**

- (A) UV-Vis absorbance profiles of lactocillin and thiocillin I (note the similarity between profiles).
- (B) (+) HR-FTMS of lactocillin, consistent with the empirical formula  $C_{51}H_{45}N_{13}O_{10}S_7$  (observed  $[M+H]^+$  at  $m/z$  1224.15184, calculated  $[M+H]^+$  1224.15354,  $\Delta$  ppm = 1.3)
- (C) Time course of lactocillin production by *L. gasserii* JV-V03 in MRS broth under anaerobic conditions at 37°C.
- (D) HPLC-MS chromatogram of the trimethylsilyldiazomethane derivatization reaction (red) in comparison to the fraction containing native lactocillin (blue).
- (E) (+) HR-ESI Orbitrap data of the resulting reaction showing a mass shift of 14 amu corresponding to the addition of one methyl group to the parent lactocillin and consistent with the empirical formula  $C_{52}H_{47}N_{13}O_{10}S_7$  (observed  $[M+H]^+$  at  $m/z$  1238.16916, calculated for  $[M+H]^+$  at  $m/z$  1238.16916,  $\Delta$  ppm = - 0.02)
- (F) Selected HMBC and COSY correlations observed and used to determine the structure of lactocillin methyl ester (see also [Data S2](#) for NMR spectra and table of chemical shifts)
- (G) HPLC-MS analysis of *L. gasserii* feeding experiments with 20 mM L-tryptophan (Top), 20 mM L-tryptophan  $D_5$  (middle), 20 mM L-alanine  $D_3$  (bottom). Mass spectra of the lactocillin peaks are shown in all three cases, where the +5 and +3 amu mass shifts are apparent in the case of L-tryptophan  $D_5$  and L-alanine  $D_3$ , respectively. H) Selected MS/MS and  $MS^n$  spectra of lactocillin methyl ester (See also [Data S2](#) for an illustration of the detected fragments).



**Figure S7. Metatranscriptomic Analysis of Oral Thiopeptide BGCs, Related to Table 1**

(A) Mapping of metatranscriptomic reads from a human supragingival plaque sample (HMP DACC) to *bgc65* (colors indicate the same key as in Figure 5).

(B) Mapping of metatranscriptomic reads from a human supragingival plaque sample (HMP DACC) to the *L. gasseri* JV-V03 plasmid harboring *lcl* (colors indicate the same key as in Figure 6).