

A Profile Hidden Markov Model to investigate the distribution and frequency of LanB-encoding lantibiotic modification genes in the human microbiome (#14436)

1

First submission

Please read the **Important notes** below, and the **Review guidance** on the next page. When ready [submit online](#). The manuscript starts on page 3.

Important notes

Editor and deadline

Ramy Aziz / 30 Nov 2016

Files

8 Figure file(s)

4 Table file(s)

Please visit the overview page to [download and review](#) the files not included in this review pdf.

Declarations

No notable declarations are present




Please in full read before you begin

How to review






When ready [submit your review online](#). The review form is divided into 5 sections. Please consider these when composing your review:

- 1. BASIC REPORTING**
- 2. EXPERIMENTAL DESIGN**
- 3. VALIDITY OF THE FINDINGS**
4. General comments
5. Confidential notes to the editor



 You can also annotate this **pdf** and upload it as part of your review

To finish, enter your editorial recommendation (accept, revise or reject) and submit.







BASIC REPORTING

-  Clear, unambiguous, professional English language used throughout.
-  Intro & background to show context. Literature well referenced & relevant.
-  Structure conforms to [PeerJ standard](#), discipline norm, or improved for clarity.
-  Figures are relevant, high quality, well labelled & described.
-  Raw data supplied (See [PeerJ policy](#)).

VALIDITY OF THE FINDINGS

-  Impact and novelty not assessed. Negative/inconclusive results accepted. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.
-  Data is robust, statistically sound, & controlled.

EXPERIMENTAL DESIGN

-  Original primary research within [Scope of the journal](#).
-  Research question well defined, relevant & meaningful. It is stated how research fills an identified knowledge gap.
-  Rigorous investigation performed to a high technical & ethical standard.
-  Methods described with sufficient detail & information to replicate.
-  Conclusion well stated, linked to original research question & limited to supporting results.
-  Speculation is welcome, but should be identified as such.

The above is the editorial criteria summary. To view in full visit <https://peerj.com/about/editorial-criteria/>

A Profile Hidden Markov Model to investigate the distribution and frequency of LanB-encoding lantibiotic modification genes in the **human microbiome**

Calum J Walsh^{1,2}, Caitriona M Guinane¹, Paul W O' Toole^{2,3}, Paul D Cotter^{Corresp. 1,3}

¹ Teagasc Food Research Centre, Moorepark, Co. Cork, Ireland

² School of Microbiology, University College Cork, Co. Cork, Ireland

³ APC Microbiome Institute, University College Cork, Co. Cork, Ireland

Corresponding Author: Paul D Cotter
Email address: paul.cotter@teagasc.ie

Background

The human microbiota plays a key role in health and disease, and bacteriocins, which are small, bacterially produced, antimicrobial peptides, are likely to have an important function in the stability and dynamics of this community. Here we examined the density and distribution of the class I lantibiotic modification protein, LanB, in human oral and stool microbiome datasets using a specially constructed profile Hidden Markov Model (HMM).

Methods

The model was validated by correctly identifying known *lanB* genes in the genomes of known bacteriocin producers more effectively than a model obtained from the Pfam database, while being sensitive enough to differentiate between different classes of lantibiotic modification proteins. This approach was compared with several existing methods to screen both genomic and metagenomic datasets obtained from the Human Microbiome Project (HMP).

Results

Of the methods evaluated, the new profile HMM identified the greatest number of putative LanB proteins in the stool and oral metagenome data while BlastP identified the fewest. In addition, the model identified more LanB proteins than the aforementioned Pfam lanthionine dehydratase model. Searching the gastrointestinal tract subset of the HMP reference genome database with the new HMM identified seven putative class I lantibiotic producers, including two members of the *Coprobacillus* genus.

Conclusions

These findings establish custom profile HMMs as a potentially powerful tool in the search for novel bioactive producers with the power to benefit human health, and reinforce the repertoire of apparent bacteriocin-encoding gene clusters that have been overlooked by culture-dependent mining efforts to date.

1 **Author Cover Page**

2 A Profile Hidden Markov Model to investigate the distribution and frequency of LanB-encoding

3 lantibiotic modification genes in the human microbiome

4

5 Calum J. Walsh^{1,3}, Caitriona M. Guinane¹, Paul W. O'Toole^{2,3}, Paul D. Cotter^{1,2*}

6

7 Teagasc Food Research Centre, Moorepark, Fermoy, Cork, Ireland¹;

8 APC Microbiome Institute, University College Cork, Cork, Ireland²;

9 School of Microbiology, University College Cork, Cork, Ireland³.

10

11 *To whom correspondence should be addressed: paul.cotter@teagasc.ie

12

13

14

15

16

17

18

19 **Abstract**

20 **Background**

21 The human microbiota plays a key role in health and disease, and bacteriocins, which are small,
22 bacterially produced, antimicrobial peptides, are likely to have an important function in the
23 stability and dynamics of this community. Here we examined the density and distribution of the
24 class I lantibiotic modification protein, LanB, in human oral and stool microbiome datasets using
25 a specially constructed profile Hidden Markov Model (HMM).

26 **Methods**

27 The model was validated by correctly identifying known *lanB* genes in the genomes of known
28 bacteriocin producers more effectively than a model obtained from the Pfam database, while
29 being sensitive enough to differentiate between different classes of lantibiotic modification
30 proteins. This approach was compared with **several existing methods** to screen both genomic and
31 metagenomic datasets obtained from the Human Microbiome Project (HMP).

32 **Results**

33 Of the methods evaluated, the new profile HMM identified the greatest number of putative LanB
34 proteins in the stool and oral metagenome data while BlastP identified the fewest. In addition,
35 the model identified more LanB proteins than the aforementioned Pfam lanthionine dehydratase
36 model. Searching the gastrointestinal tract subset of the HMP reference genome database with
37 the new HMM identified seven putative class I lantibiotic producers, including two members of
38 the *Coprobacillus* genus.

39 **Conclusions**

40 These findings establish custom profile HMMs as a potentially powerful tool in the search for
41 novel bioactive producers with the power to benefit human health, and reinforce the repertoire of
42 apparent bacteriocin-encoding gene clusters that have been **overlooked by culture-dependent**
43 mining efforts to date.

44 **Background**

45 Bacteriocins are ribosomally synthesised peptides produced by bacteria that inhibit the growth of
46 other bacteria. Some classes of bacteriocins are post-translationally modified to provide
47 structures beyond those possible by ribosomal translation alone. These modifications are
48 typically key to the peptide's functionality, stability and target recognition (Arnison et al. 2013).
49 **Lantibiotics** are one such class of small (<5 kDa) modified bacteriocins, possessing characteristic
50 thioester amino acids lanthionine or methyllanthionine (Perez et al. 2014). Lantibiotics form a
51 subgroup within the larger lantipeptide family, which also includes peptides that lack
52 antimicrobial activity. Lantipeptides can be divided into four different classes based on the
53 distinct biosynthetic enzymes responsible for their posttranslational modification (Arnison et al.
54 2013).

55 The most commonly studied lantibiotic, Nisin, is a subclass I lantibiotic, meaning that the linear
56 prepeptide is processed by a LanBC modification system (Arnison et al. 2013). Firstly, eight
57 serine and threonine residues in the core peptide are dehydrated by the dehydratase LanB to form
58 dehydroalanine and dehydrobutyrine, respectively (Xie & van der Donk 2004). Secondly, five
59 lanthionine and methyllanthionine crosslinks are formed by the nucleophilic addition of cysteinyl
60 thiols to dehydroalanine and dehydrobutyrine, respectively, by the cyclase LanC (Xie & van der
61 Donk 2004). Finally, the leader sequence, necessary for recognition by the modification enzymes

62 in the two previous steps, is removed by the protease LanP to produce the active lantibiotic (Xie
63 & van der Donk 2004). The gene-encoded nature of bacteriocins and bacteriocin-like peptides
64 makes them ideal candidates for genome mining. In the case of modified bacteriocins, the
65 structural prepeptide coding sequence often appears alongside the genes encoding proteins
66 responsible for its modification and export from the cell. However, as more bacteriocins are
67 discovered, the heterogeneous nature of these prepeptides is becoming ever more apparent. This
68 diversity, coupled with their small sequence length, makes bacteriocin prepeptides much more
69 difficult to detect using sequence-homology based searches like BLAST (Altschul et al. 1990).
70 In an effort to address these obstacles, shifting the focus to the detection of bacteriocin-
71 associated proteins opens up more avenues of discovery than simply searching for prepeptide
72 homologs. This provides opportunities to better determine the frequency with which specific
73 types of bacteriocin gene clusters can be found in different environmental niches, such as the
74 human microbiota, through the investigation of metagenomic data.

75 It has been estimated that the human microbiota comprises approximately 100 trillion bacterial
76 cells, outnumbering our own cells by a factor of 10 or more (Bäckhed et al. 2005). A recent
77 publication, however, has argued that the ratio is actually more likely to be one-to-one, with the
78 numbers being similar enough that each defecation event may alter the ratio to favour human
79 cells over bacteria (Sender et al. 2016). Of greater consequence than bacterial numbers, however,
80 is the collection of genes encoded in this metagenome, thought to be approximately 150 times
81 greater than the human gene complement, with a functional potential far broader than that of its
82 host (Qin et al. 2010). Regardless of absolute numbers, this dynamic community is thought to
83 contain 100-1000 phylotypes (Faith et al. 2013; Qin et al. 2010) and play an integral role in
84 human health and disease (Clemente et al. 2012; Flint et al. 2012). The human microbiota

85 exhibits robust temporal stability (Belstrøm et al. 2016; Jeffery et al. 2016) perhaps due, in part,
86 to the protection against invading bacteria conferred by bacteriocins and other antimicrobials
87 produced *in situ*. As such, investigation of the density and diversity of bacteriocins produced in
88 the microbiome of healthy individuals may shed light on beneficial and harmful members of this
89 community, and key organisms for maintaining typical i.e. health-associated microbiota
90 composition.

91 Mining the human microbiota, especially for antimicrobial compounds, has become a popular
92 area of research in recent years (Donia et al. 2014; Walsh et al. 2015). Due to the availability of
93 metagenomic data generated by large public funding initiatives such as the Human Microbiome
94 Project in the U.S. (The Human Microbiome Project Consortium 2012) and the European
95 MetaHIT consortium (Dusko Ehrlich 2010), *in silico* mining of data has emerged as a new tool
96 that has the potential to identify antimicrobial-producing probiotics that can modulate the gut
97 microbiota (Erejuwa et al. 2014; Walsh et al. 2014), or address the increasingly serious threat to
98 public health caused by antimicrobial resistance. There are many available tools for mining the
99 microbiome for antimicrobials, including BAGEL3 (van Heel et al. 2013), antiSMASH (Weber
100 et al. 2015), and traditional sequence-based approaches like BLAST (Altschul et al. 1990). A
101 feature commonly integrated into these tools are Hidden Markov Models (HMM) (Morton et al.
102 2015; van Heel et al. 2013; Weber et al. 2015) i.e. statistical methods often used to model
103 biological data such as speech recognition, disease interaction and changes in gene expression in
104 cancer (Gales & Young 2007; Seifert et al. 2014; Sherlock et al. 2013). Profile HMMs, a specific
105 subset of HMMs, represent the patterns, motifs and other properties of a multiple sequence
106 alignment by applying a statistical model to estimate the true frequency of a nucleotide or amino
107 acid at a given position in the alignment from its observed frequency (Yoon 2009). Profile

108 HMMs differ from general HMMs as they move strictly from left to right and do not contain any
109 cycles, a feature that makes them suitable for mimicking the actions of the ribosome during
110 translation. The profile uses three types of hidden states - match states, insert states, and delete
111 states, to describe position-specific residue frequencies, insertions, and deletions, respectively
112 (Yoon 2009). Profile HMMs are potentially more sensitive than sequence homology approaches
113 for identifying more distantly related proteins as they focus on function-dependent conserved
114 motifs that are theoretically slower-evolving, as opposed to focussing on overall sequence
115 similarity. Notably, Skewes-Cox *et al.* successfully designed an approach employing profile
116 HMMs to detect viral protein sequences in metagenomic sequence data (Skewes-Cox *et al.*
117 2014).

118 In this study we designed, validated and implemented a Profile HMM to search for putative
119 subclass I lantibiotic gene clusters in the HMP metagenomes and compared its performance to
120 some of the tools mentioned above.

121 **Methods**

122 **Data Collection**

123 HMASM (HMP Illumina WGS Assemblies) and HMRGD (HMP Reference Genomes Data)
124 were downloaded from the Data Analysis and Coordination Centre for the HMP . 835 bacterial
125 RefSeq protein sequences annotated as “lantibiotic dehydratase” were downloaded from NCBI
126 Protein website (13 Apr 2015) in FASTA format.

127 **Building and Validating the new Profile Hidden Markov Model**

128 A multiple sequence alignment was generated in the aligned-FASTA format using MUSCLE
129 (v3.8.31) (Edgar 2004), and a profile HMM was built from the MSA aligned-FASTA file using
130 the HMMER tool hmmbuild (v3.1b1 May 2013) . For comparison of the new model’s

131 performance, HMMER3's hmmsearch tool was used to search the pfam lantibiotic dehydratase
132 model PF04738 against the same stool and oral HMASM assemblies. Positive and negative
133 controls (listed in Table 1) were used to evaluate the model's ability to 1) accurately identify
134 LanB protein sequences, and 2) distinguish LanB protein sequences from other, related,
135 lantibiotic modification proteins (i.e. LanM and LanL).

136 **Target Sequence Translation**

137 The HMMER3 hmmsearch tool only accepts protein sequences as targets for comparison to
138 protein profile HMMs so a python script was created to translate the nucleotide sequences into
139 protein sequences. The DNA nucleotide sequences were translated in six frames using the
140 standard genetic code.

141 **Metagenomic Screen**

142 The HMMER3 tool hmmsearch was used to search both the new LanB profile HMM and the
143 Pfam PF04738 profile HMM (Punta et al. 2012) against the **stool and oral subsets** of the Human
144 Microbiome Project's whole metagenomic shotgun sequencing assemblies (HMASM). 139 stool
145 communities and 382 communities from eight different body sites within the oral cavity were
146 screened from the HMP database. These are listed in Table 2. As an additional comparison of
147 performance, a traditional BlastP screen was performed on the same metagenomic samples using
148 the nisin-associated lanthionine dehydratase, NisB, as the driver sequence (GenBank accession
149 number CAA79468.1).

150 **Manual Examination of Randomly Selected Gene Neighbourhoods**

151 **A subset of sixty hits were selected** and the surrounding region examined to identify other
152 proteins involved in lantibiotic biosynthesis. Open Reading Frames were identified using

153 Glimmer v3.02 (Delcher et al. 1999), which were then visualised using Artemis (Carver et al.
154 2012) and blasted against the nr database using BlastP.

155 **Genomic Screen**

156 HMMER3's hmmsearch tool was used to search the new profile HMM against the draft genomes
157 comprising the gastrointestinal tract subset of the Human Microbiome Project's reference
158 genome database.

159 **Results**

160 **Validation of the Profile Hidden Markov Model**

161 The ability of the newly developed profile HMM and the pfam lantibiotic dehydratase model
162 PF04738 to detect LanB-encoding genes were compared using the positive and negative controls
163 listed in Table 1. The positive controls selected were all previously characterised bacteriocin
164 producers for which the sequence of the relevant biosynthetic gene cluster was available. A
165 graphical representation of these clusters is presented in Figure 1. *Lactococcus lactis* subsp.
166 *lactis* KF147 was chosen as a negative control because it is of the same subspecies as three of the
167 positive controls (*Lactococcus lactis* subsp. *lactis* S0, *Lactococcus lactis* subsp. *lactis* CV56 and
168 *Lactococcus lactis* subsp. *lactis* IO-1) but does not produce a bacteriocin. *Streptococcus mutans*
169 GS-5, *Streptomyces cinnamoneus cinnamoneus* DSM 4005 and the *Lactococcus lactis* subsp.
170 *lactis* IL1835 plasmid pES2 were chosen as negative controls to evaluate the ability of the model
171 to differentiate between LanB (subclass I) proteins and the LanM proteins-from these strains,
172 which perform a similar, but distinct, function in the posttranslational modification of **class II**
173 lantibiotics. *Streptomyces venezuelae* ATCC 10712 was chosen as the final negative control as it
174 has been reported to produce a LanL-type lantipeptide (Goto et al. 2010). Examination of the
175 ATCC 10712 genome using BAGEL3 identified several other orphan lantibiotic modification

176 genes, including those encoding putative **LanL, LanM, LanD and LanB proteins**. The genome
177 also appeared to encode a **class III** lantipeptide cluster comprised of genes potentially encoding a
178 structural protein, two ABC-type transporters and a LanKC modification protein (these genes
179 and clusters are depicted in Figure 2). Notably, there have been no reports of **class I** antibiotic
180 production by this strain.

181 **The newly developed LanB profile HMM correctly identified the LanB protein in all nine**
182 **positive controls, while the PF04738 profile HMM correctly identified the LanB protein in eight**
183 **of the nine positive controls, failing to detect the Bsa-associated LanB protein in *Staphylococcus***
184 ***aureus* subsp. *aureus* USA300_FPR3757. Both the LanB and PF04738 profile HMMs returned**
185 **no false positives when searched against the five negative controls used in this study, and, thus,**
186 **the orphan hypothetical LanB protein reported by BAGEL3 to be encoded in ATCC 10712**
187 **genome was correctly regarded as a negative.**

188 **Metagenomic Screen**

189 A search with the newly developed profile HMM against the HMASM database identified 399
190 hits with an E-value of less than 1×10^{-5} from the stool metagenomes and 1169 hits with an E-
191 value of less than 1×10^{-5} from the oral metagenomes. In contrast, the PF04738 model identified
192 288 hits with an E-value of less than 1×10^{-5} from the stool metagenomes and 686 with an E-value
193 of less than 1×10^{-5} from the oral metagenomes. Our model reported at least one putative
194 lantibiotic gene cluster in 81% of oral metagenomes and 86% of stool metagenomes, compared
195 to 73% and 76%, respectively, identified by the Pfam model. The distribution of hits per sample
196 is presented in Figure 3. BlastP identified 231 hits with an E-value of less than 1×10^{-5} from the
197 stool metagenomes and 374 hits with an E-value of less than 1×10^{-5} from the oral metagenomes.
198 **The results of these three approaches** were compared to ascertain what proportion of significant

199 hits was common to more than one search method. The results of this comparison are
200 summarised in Figure 4 and show that the newly developed profile HMM identified the greatest
201 number of lantibiotic modification genes in datasets from both body sites, while the BlastP
202 approach identified the fewest.

203 The overall results of these combined screening approaches, illustrated in Figure 5 and
204 summarised in Supplemental Table 1, show a higher number and density of hits in the oral
205 metagenomes than in the stool metagenomes and they also reveal a large variation in density of
206 hits between the different sites within the oral metagenomes.

207 Manual Examination of Selected Gene Neighbourhoods

208 Sixty hits were randomly selected from those identified by the new profile HMM and manually
209 examined to determine if a bacteriocin gene cluster could be identified. 42% (25/60) of these
210 were not further analyzed because the often relatively short regions assembled from the shotgun
211 data prevented the identification of a full lantibiotic gene cluster. However, of the 35 remaining
212 clusters, 28 (80%) appeared to encode multiple genes involved in the biosynthesis of bacteriocins
213 and thiopeptides. These genes encode proteins involved in posttranslational modification,
214 bacteriocin transport, leader cleavage and regulation (Supplemental Figure 1).

215 Genomic Screen

216 The draft genomes of the gastrointestinal tract subset of the HMRGD were also used as a
217 database and searched using the new profile HMM. This resulted in the identification of seven
218 hits with an E-value of less than 1×10^{-5} , including two strains of *Coprobacillus*, a potentially
219 probiotic genus (Stein et al. 2013; Yan et al. 2012) (Table 3). From these seven genomes, only
220 three lantibiotic gene clusters were identified by BAGEL3, these are illustrated in Figure 6.

221 Although this low frequency of lanthionine dehydratase proteins in the dataset contrasts with the

222 findings of the metagenome screen reported above, it is in agreement with previous reports of
223 relatively low class I lantibiotic density within the human microbiota (Walsh et al. 2015; Zheng
224 et al. 2014). A possible explanation for this is that the class I lantibiotic clusters identified in the
225 metagenomics data by the new profile HMM are **present in the genomes of rarer members** of the
226 gut microbiota, which are not represented in the HMP reference genome database.

227 **Discussion**

228 Bacteriocin production enhances the competitiveness of bacteria living in complex communities
229 and has the potential to be harnessed for the benefit of human health. The goal of this study was
230 to develop a profile HMM and to assess its ability, in comparison with several other approaches,
231 to detect putative subclass I lantibiotic gene clusters in human metagenomic datasets. Through
232 this process, it was also possible to evaluate the potential density and distribution of these
233 bacteriocin gene clusters in the human microbiota.

234 To validate the model, nine positive controls and five negative controls were selected to evaluate
235 its sensitivity and specificity. These controls were selected based on reported bacteriocin
236 production; the positive controls were all known producers of **class I lantibiotics** while the
237 negative controls produced either different classes of lantibiotics or none at all. Following
238 validation, genomic and metagenomic data corresponding to two niches within the human
239 microbiome were chosen as the focus of this study. The first of these niches was human stool and
240 was selected as the corresponding samples were most likely to yield bacteriocin producers with
241 the potential to modulate undesirable microbiota profiles associated with obesity, colorectal
242 cancer, type 2 diabetes or inflammatory bowel diseases due to their ability to survive and
243 colonise this environment. Secondly, human oral communities were examined as a previous
244 study by Zheng *et al.* showed that they contained, by far, the greatest percentage of bacteriocin

245 structural genes across a number of human metagenome samples (Zheng et al. 2014). Zheng *et*
246 *al.* reported that 80% of class I bacteriocins (lantibiotics) and 89% of all bacteriocins identified
247 using their method originated in the oral metagenomes, while the stool metagenomes contained
248 just 15% and 7%, respectively. The same study reported that 88% of samples from the oral
249 cavity and 73% of samples from the gut contained at least one bacteriocin (regardless of class),
250 while the new profile HMM reported these statistics as 81% and 83%, respectively for sub- I
251 lantibiotics alone. The *in silico* screen carried out with the profile HMM is consistent with the
252 observation by Zheng *et al* (Zheng et al. 2014) by yielding a higher number and density of hits
253 from the oral, compared to the stool, metagenomic data. Furthermore, the large variation in
254 density of hits between sites within the oral environment suggests that lantibiotic production
255 confers a greater advantage in subgingival plaque, supragingival plaque, and tongue dorsum
256 communities compared to communities from the throat and buccal mucosa. This may be due to
257 the direct benefits of antimicrobial activity but could also involve the intra- and interspecies
258 signalling roles attributed to lantibiotic peptides (Upton et al. 2001).

259 One of the most interesting observations from the study was the large variation in the numbers of
260 *lanB* genes reported by the three different approaches. The BlastP approach identified, by far, the
261 lowest number of significant hits overall and the lowest in every body site examined, except for
262 the saliva microbiome. Our model identified more than double the number of hits provided by
263 the BlastP-based approach. This is to be expected as profile HMMs are known to typically
264 outperform pairwise sequence comparison methods (such as BLAST) in the detection of distant
265 homologs (Park et al. 1998). Our model also identified a greater number of LanB proteins than
266 the Pfam PF04738 model when used to search the same data using the same parameters. While
267 the PF04738 model relates to the N-terminus of the lanthionine dehydratase protein, responsible

268 for the serine-threonine glutamylation step of lantibiotic modification (Ortega et al. 2015), the
269 newly developed profile HMM takes the full length of the LanB protein into consideration,
270 thereby providing greater predictive power.

271 Zheng *et al.*, using the same metagenomic data that was the focus of this study, identified 17
272 class I lantibiotics from stool samples and 76 from oral samples, a much lower frequency of
273 detection than in this study, probably due to the different methodologies used. That study
274 focused on searching for proteins similar to those in BAGEL3's manually curated database, an
275 approach which likely lost sensitivity because bacteriocin precursor peptides can differ
276 considerably at primary sequence level. Furthermore, the screen employed a BLAST-based
277 approach which, as demonstrated here, exhibited the lowest number of significant hits reported.

278 To investigate the areas surrounding the LanB-encoding genes identified by our model we
279 randomly selected thirty positive hits from the oral and stool metagenome screens for manual
280 examination. This approach revealed that several of the hits were on scaffolds that were either
281 too small to contain a full gene or did not contain the gene's start codon. This was most likely as
282 a consequence of the fragmented nature of the metagenomic data, as opposed the identification
283 of true false positives by the model and would probably occur regardless of the method
284 employed. 42% (25/60) of hits selected for manual examination were discarded based on these
285 criteria. It also revealed that a considerable number of hits exhibited low (~30%) similarity to
286 putative thioesterases in the nr protein sequence database, highlighting that lanthionine
287 dehydratases are relatively-closely related to proteins involved in the posttranslational
288 modification of thiopeptides, most likely those responsible for dehydration of serine and
289 threonine residues (Garg et al. 2013). The similarity between these dehydratase proteins suggests
290 a possible common ancestor protein (Kelly et al. 2009). Another possible explanation relates to

291 the fact that all of the proteins annotated as thiopeptide modification proteins are putative
292 annotations and none, to our knowledge, have been confirmed as such *in vitro*. It is possible,
293 therefore, that these may simply be lanthionine dehydratases which have been incorrectly
294 annotated due to automatic software and incomplete/under-curated databases. The majority of
295 clusters identified contained genes encoding both LanB and LanC modification proteins as well
296 as a leader cleavage and **activation peptidase**, and **ABC transporter proteins** for export of the
297 mature peptide, suggesting that these have the potential to encode a functional lantibiotic.

298 To evaluate the model's performance in a genomic context we applied it to the gastrointestinal
299 tract subset of the HMP's reference genome database and compared the results to our previously
300 published study which used the online bacteriocin genome mining tool BAGEL3 (van Heel et al.
301 2013) to screen this same database (Walsh et al. 2015). The results of the two screens were
302 startlingly different and served to highlight the variation in results that can arise from applying
303 different methods to the same data.

304 **Conclusions**

305 Across the oral and stool communities examined, this study identified **2007** unique putative
306 subclass I lantibiotic biosynthetic gene clusters, further emphasising the tremendous potential
307 that the human microbiota has as a source of therapeutic compounds. The next challenge lies in
308 correctly identifying those elements with the ability to desirably modulate the microbiota and
309 utilizing them in the treatment of microbiota-associated disease.

310 **Acknowledgements**

311 The authors would like to thank Manimozhiyan Arumugam for helpful discussion.

312 **List of Abbreviations**

Abbreviation	Description
HMASM	Human Microbiome Project's Illumina Whole Genome Shotgun Assemblies
HMM	Hidden Markov Model
HMP	Human Microbiome Project
HMRGD	Human Microbiome Project's Reference Genome Data

313

314 **References**

- 315 Hmmer. Available at <http://hmmer.org>.
- 316 Human Microbiome Project; Data Analysis and Coordination Center. Available at <http://hmpdacc.org/>.
- 317 13 Apr 2015. National Centre for Biotechnology Information; Standard Protein BLAST. Available at
- 318 <http://www.ncbi.nlm.nih.gov/protein> (accessed 13 Apr 2015).
- 319 Altschul SF, Gish W, Miller W, Myers EW, and Lipman DJ. 1990. Basic local alignment search tool. *J Mol*
- 320 *Biol* 215:403-410. 10.1016/s0022-2836(05)80360-2
- 321 Arnison PG, Bibb MJ, Bierbaum G, Bowers AA, Bugni TS, Bulaj G, Camarero JA, Campopiano DJ, Challis
- 322 GL, Clardy J, Cotter PD, Craik DJ, Dawson M, Dittmann E, Donadio S, Dorrestein PC, Entian K-D,
- 323 Fischbach MA, Garavelli JS, Goransson U, Gruber CW, Haft DH, Hemscheidt TK, Hertweck C, Hill
- 324 C, Horswill AR, Jaspars M, Kelly WL, Klinman JP, Kuipers OP, Link AJ, Liu W, Marahiel MA,
- 325 Mitchell DA, Moll GN, Moore BS, Muller R, Nair SK, Nes IF, Norris GE, Olivera BM, Onaka H,
- 326 Patchett ML, Piel J, Reaney MJT, Rebuffat S, Ross RP, Sahl H-G, Schmidt EW, Selsted ME,
- 327 Severinov K, Shen B, Sivonen K, Smith L, Stein T, Sussmuth RD, Tagg JR, Tang G-L, Truman AW,
- 328 Vederas JC, Walsh CT, Walton JD, Wenzel SC, Willey JM, and van der Donk WA. 2013.
- 329 Ribosomally synthesized and post-translationally modified peptide natural products: overview
- 330 and recommendations for a universal nomenclature. *Natural Product Reports* 30:108-160.
- 331 10.1039/C2NP20085F
- 332 Bäckhed F, Ley RE, Sonnenburg JL, Peterson DA, and Gordon JI. 2005. Host-Bacterial Mutualism in the
- 333 Human Intestine. *Science* 307:1915-1920. 10.1126/science.1104816
- 334 Belstrøm D, Holmstrup P, Bardow A, Kokaras A, Fiehn N-E, and Paster BJ. 2016. Temporal Stability of the
- 335 Salivary Microbiota in Oral Health. *PLoS ONE* 11:e0147472. 10.1371/journal.pone.0147472
- 336 Carver T, Harris SR, Berriman M, Parkhill J, and McQuillan JA. 2012. Artemis: an integrated platform for
- 337 visualization and analysis of high-throughput sequence-based experimental data. *Bioinformatics*
- 338 28:464-469. 10.1093/bioinformatics/btr703
- 339 Clemente Jose C, Ursell Luke K, Parfrey Laura W, and Knight R. 2012. The Impact of the Gut Microbiota
- 340 on Human Health: An Integrative View. *Cell* 148:1258-1270.
- 341 <http://dx.doi.org/10.1016/j.cell.2012.01.035>
- 342 Delcher AL, Harmon D, Kasif S, White O, and Salzberg SL. 1999. Improved microbial gene identification
- 343 with GLIMMER. *Nucleic Acids Res* 27:4636-4641.
- 344 Donia Mohamed S, Cimerancic P, Schulze Christopher J, Wieland Brown Laura C, Martin J, Mitreva M,
- 345 Clardy J, Linington Roger G, and Fischbach Michael A. 2014. A Systematic Analysis of

- 346 Biosynthetic Gene Clusters in the Human Microbiome Reveals a Common Family of Antibiotics.
347 *Cell* 158:1402-1414. 10.1016/j.cell.2014.08.032
- 348 Dusko Ehrlich S. 2010. [Metagenomics of the intestinal microbiota: potential applications]. *Gastroenterol*
349 *Clin Biol* 34 Suppl 1:S23-28. 10.1016/s0399-8320(10)70017-8
- 350 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
351 *Acids Res* 32:1792-1797. 10.1093/nar/gkh340
- 352 Erejuwa OO, Sulaiman SA, and Ab Wahab MS. 2014. Modulation of gut microbiota in the management
353 of metabolic disorders: the prospects and challenges. *Int J Mol Sci* 15:4158-4188.
354 10.3390/ijms15034158
- 355 Faith JJ, Guruge JL, Charbonneau M, Subramanian S, Seedorf H, Goodman AL, Clemente JC, Knight R,
356 Heath AC, Leibel RL, Rosenbaum M, and Gordon JI. 2013. The Long-Term Stability of the Human
357 Gut Microbiota. *Science* 341. 10.1126/science.1237439
- 358 Flint HJ, Scott KP, Louis P, and Duncan SH. 2012. The role of the gut microbiota in nutrition and health.
359 *Nat Rev Gastroenterol Hepatol* 9:577-589.
- 360 Gales M, and Young S. 2007. The application of hidden Markov models in speech recognition. *Found*
361 *Trends Signal Process* 1:195-304. 10.1561/20000000004
- 362 Garg N, Salazar-Ocampo LMA, and van der Donk WA. 2013. In vitro activity of the nisin dehydratase
363 NisB. *Proceedings of the National Academy of Sciences of the United States of America*
364 110:7258-7263. 10.1073/pnas.1222488110
- 365 Goto Y, Li B, Claesen J, Shi Y, Bibb MJ, and van der Donk WA. 2010. Discovery of Unique Lanthionine
366 Synthetases Reveals New Mechanistic and Evolutionary Insights. *PLoS Biol* 8:e1000339.
367 10.1371/journal.pbio.1000339
- 368 Jeffery IB, Lynch DB, and O'Toole PW. 2016. Composition and temporal stability of the gut microbiota in
369 older persons. *ISME J* 10:170-182. 10.1038/ismej.2015.88
- 370 Kelly WL, Pan L, and Li C. 2009. Thiostrepton biosynthesis: prototype for a new family of bacteriocins. *J*
371 *Am Chem Soc* 131:4327-4334. 10.1021/ja807890a
- 372 Morton JT, Freed SD, Lee SW, and Friedberg I. 2015. A large scale prediction of bacteriocin gene blocks
373 suggests a wide functional spectrum for bacteriocins. *BMC Bioinformatics* 16:1-9.
374 10.1186/s12859-015-0792-9
- 375 Ortega MA, Hao Y, Zhang Q, Walker MC, van der Donk WA, and Nair SK. 2015. Structure and mechanism
376 of the tRNA-dependent lantibiotic dehydratase NisB. *Nature* 517:509-512. 10.1038/nature13888
- 377 Park J, Karplus K, Barrett C, Hughey R, Haussler D, Hubbard T, and Chothia C. 1998. Sequence
378 comparisons using multiple sequences detect three times as many remote homologues as
379 pairwise methods. *J Mol Biol* 284:1201-1210. 10.1006/jmbi.1998.2221
- 380 Perez RH, Zendo T, and Sonomoto K. 2014. Novel bacteriocins from lactic acid bacteria (LAB): various
381 structures and applications. *Microbial Cell Factories* 13:S3-S3. 10.1186/1475-2859-13-S1-S3
- 382 Punta M, Coghill PC, Eberhardt RY, Mistry J, Tate J, Boursnell C, Pang N, Forslund K, Ceric G, Clements J,
383 Heger A, Holm L, Sonnhammer ELL, Eddy SR, Bateman A, and Finn RD. 2012. The Pfam protein
384 families database. *Nucleic Acids Res* 40:D290-D301. 10.1093/nar/gkr1065
- 385 Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N, Levenez F, Yamada T,
386 Mende DR, Li J, Xu J, Li S, Li D, Cao J, Wang B, Liang H, Zheng H, Xie Y, Tap J, Lepage P, Bertalan
387 M, Batto J-M, Hansen T, Le Paslier D, Linneberg A, Nielsen HB, Pelletier E, Renault P, Sicheritz-
388 Ponten T, Turner K, Zhu H, Yu C, Li S, Jian M, Zhou Y, Li Y, Zhang X, Li S, Qin N, Yang H, Wang J,
389 Brunak S, Dore J, Guarner F, Kristiansen K, Pedersen O, Parkhill J, Weissenbach J, Bork P, Ehrlich
390 SD, and Wang J. 2010. A human gut microbial gene catalogue established by metagenomic
391 sequencing. *Nature* 464:59-65.
392 http://www.nature.com/nature/journal/v464/n7285/supinfo/nature08821_S1.html

- 393 Seifert M, Abou-El-Ardat K, Friedrich B, Klink B, and Deutsch A. 2014. Autoregressive Higher-Order
394 Hidden Markov Models: Exploiting Local Chromosomal Dependencies in the Analysis of Tumor
395 Expression Profiles. *PLoS ONE* 9:e100295. 10.1371/journal.pone.0100295
- 396 Sender R, Fuchs S, and Milo R. 2016. Revised estimates for the number of human and bacteria cells in
397 the body. *bioRxiv*. 10.1101/036103
- 398 Sherlock C, Xifara T, Telfer S, and Begon M. 2013. A coupled hidden Markov model for disease
399 interactions. *Journal of the Royal Statistical Society Series C, Applied Statistics* 62:609-627.
400 10.1111/rssc.12015
- 401 Skewes-Cox P, Sharpton TJ, Pollard KS, and DeRisi JL. 2014. Profile Hidden Markov Models for the
402 Detection of Viruses within Metagenomic Sequence Data. *PLoS ONE* 9:e105067.
403 10.1371/journal.pone.0105067
- 404 Stein RR, Bucci V, Toussaint NC, Buffie CG, Räscht G, Pamer EG, Sander C, and Xavier JB. 2013. Ecological
405 Modeling from Time-Series Inference: Insight into Dynamics and Stability of Intestinal
406 Microbiota. *PLoS Comput Biol* 9:e1003388. 10.1371/journal.pcbi.1003388
- 407 The Human Microbiome Project Consortium. 2012. Structure, function and diversity of the healthy
408 human microbiome. *Nature* 486:207-214. 10.1038/nature11234
- 409 Upton M, Tagg JR, Wescombe P, and Jenkinson HF. 2001. Intra- and Interspecies Signaling between
410 *Streptococcus salivarius* and *Streptococcus pyogenes* Mediated by SalA and SalA1 Lantibiotic
411 Peptides. *J Bacteriol* 183:3931-3938. 10.1128/JB.183.13.3931-3938.2001
- 412 van Heel AJ, de Jong A, Montalban-Lopez M, Kok J, and Kuipers OP. 2013. BAGEL3: Automated
413 identification of genes encoding bacteriocins and (non-)bactericidal posttranslationally modified
414 peptides. *Nucleic Acids Res* 41:W448-453. 10.1093/nar/gkt391
- 415 Walsh CJ, Guinane CM, Hill C, Ross RP, O'Toole PW, and Cotter PD. 2015. In silico identification of
416 bacteriocin gene clusters in the gastrointestinal tract, based on the Human Microbiome Project's
417 reference genome database. *BMC Microbiol* 15:183. 10.1186/s12866-015-0515-4
- 418 Walsh CJ, Guinane CM, O'Toole PW, and Cotter PD. 2014. Beneficial modulation of the gut microbiota.
419 *FEBS Lett* 588:4120-4130. 10.1016/j.febslet.2014.03.035
- 420 Weber T, Blin K, Duddela S, Krug D, Kim HU, Bruccoleri R, Lee SY, Fischbach MA, Müller R, Wohlleben W,
421 Breitling R, Takano E, and Medema MH. 2015. antiSMASH 3.0—a comprehensive resource for
422 the genome mining of biosynthetic gene clusters. *Nucleic Acids Res* 43:W237-W243.
423 10.1093/nar/gkv437
- 424 Xie L, and van der Donk WA. 2004. Post-translational modifications during lantibiotic biosynthesis. *Curr*
425 *Opin Chem Biol* 8:498-507. 10.1016/j.cbpa.2004.08.005
- 426 Yan X, Gurtler JB, Fratamico PM, Hu J, and Juneja VK. 2012. Phylogenetic identification of bacterial MazF
427 toxin protein motifs among probiotic strains and foodborne pathogens and potential
428 implications of engineered probiotic intervention in food. *Cell & Bioscience* 2:1-13.
429 10.1186/2045-3701-2-39
- 430 Yoon B-J. 2009. Hidden Markov Models and their Applications in Biological Sequence Analysis. *Current*
431 *Genomics* 10:402-415. 10.2174/138920209789177575
- 432 Zheng J, Gänzle MG, Lin XB, Ruan L, and Sun M. 2014. Diversity and dynamics of bacteriocins from
433 human microbiome. *Environ Microbiol*:n/a-n/a. 10.1111/1462-2920.12662

434

Figure 1

BAGEL3 output of putative bacteriocin gene clusters identified in positive controls used in validation of our new profile HMM.

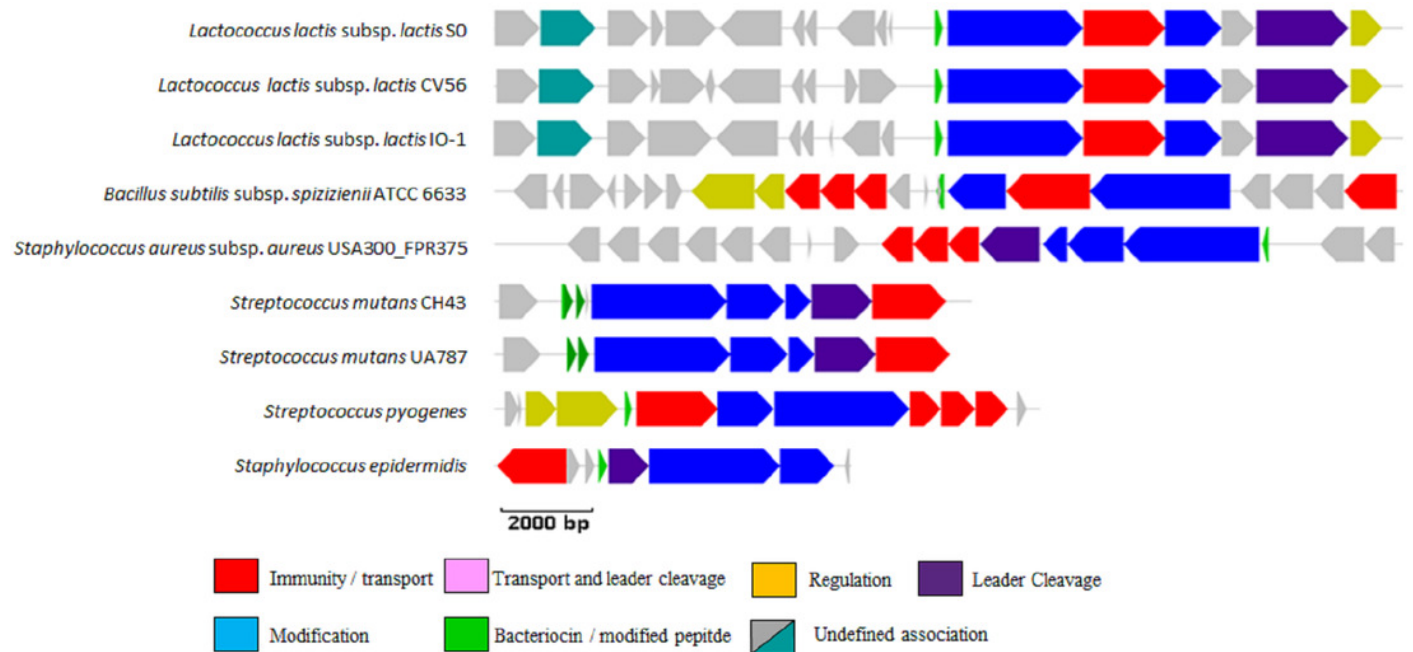


Figure 2

BAGEL3 output of putative bacteriocin gene clusters identified in negative controls used in validation of our new profile HMM.

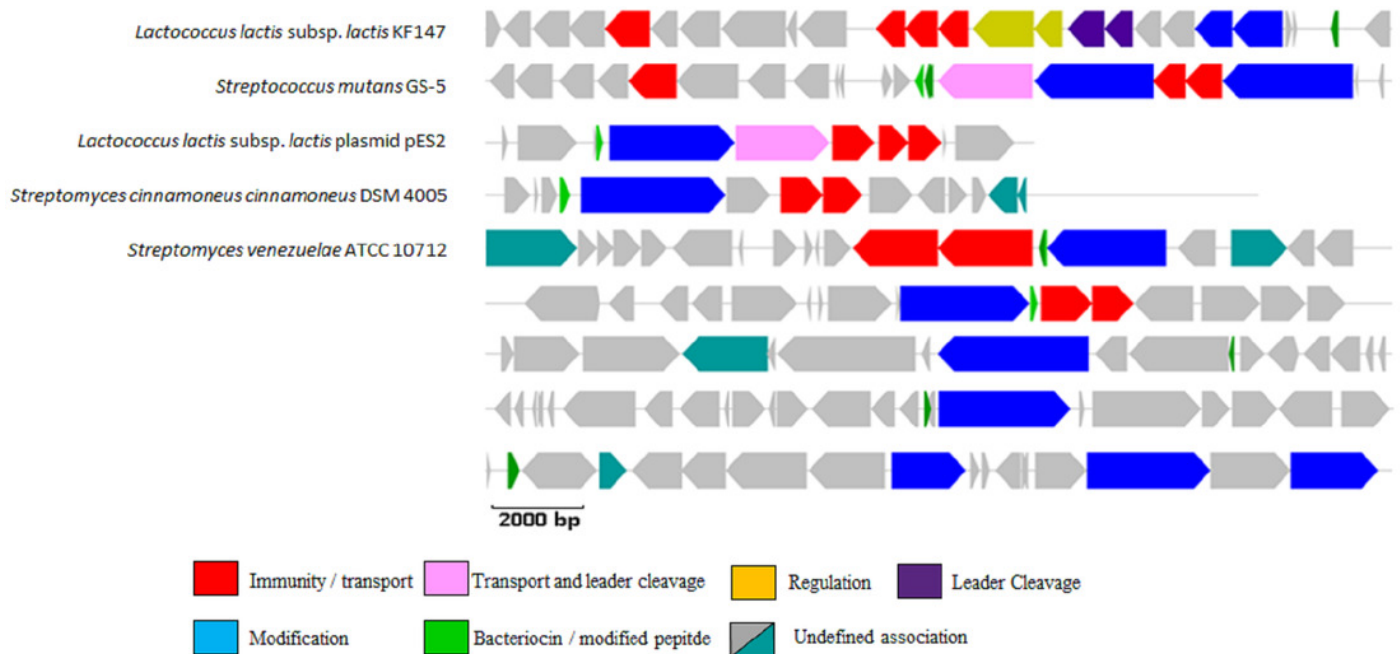


Figure 3

Distributions of lanthionine dehydratase proteins per sample identified by our new profile HMM.

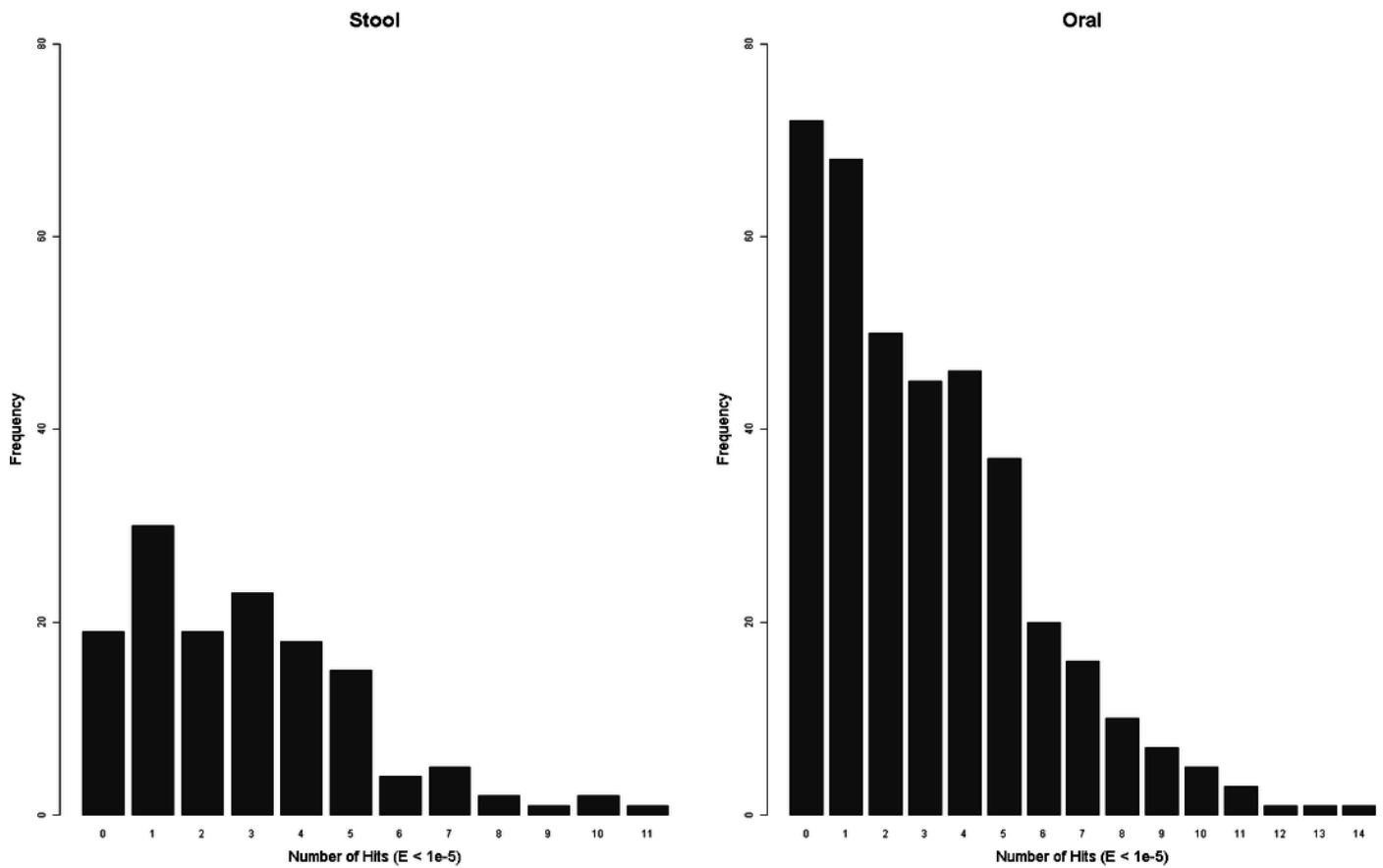


Figure 4

Numbers of lanthionine dehydratase proteins reported by single and multiple methods.

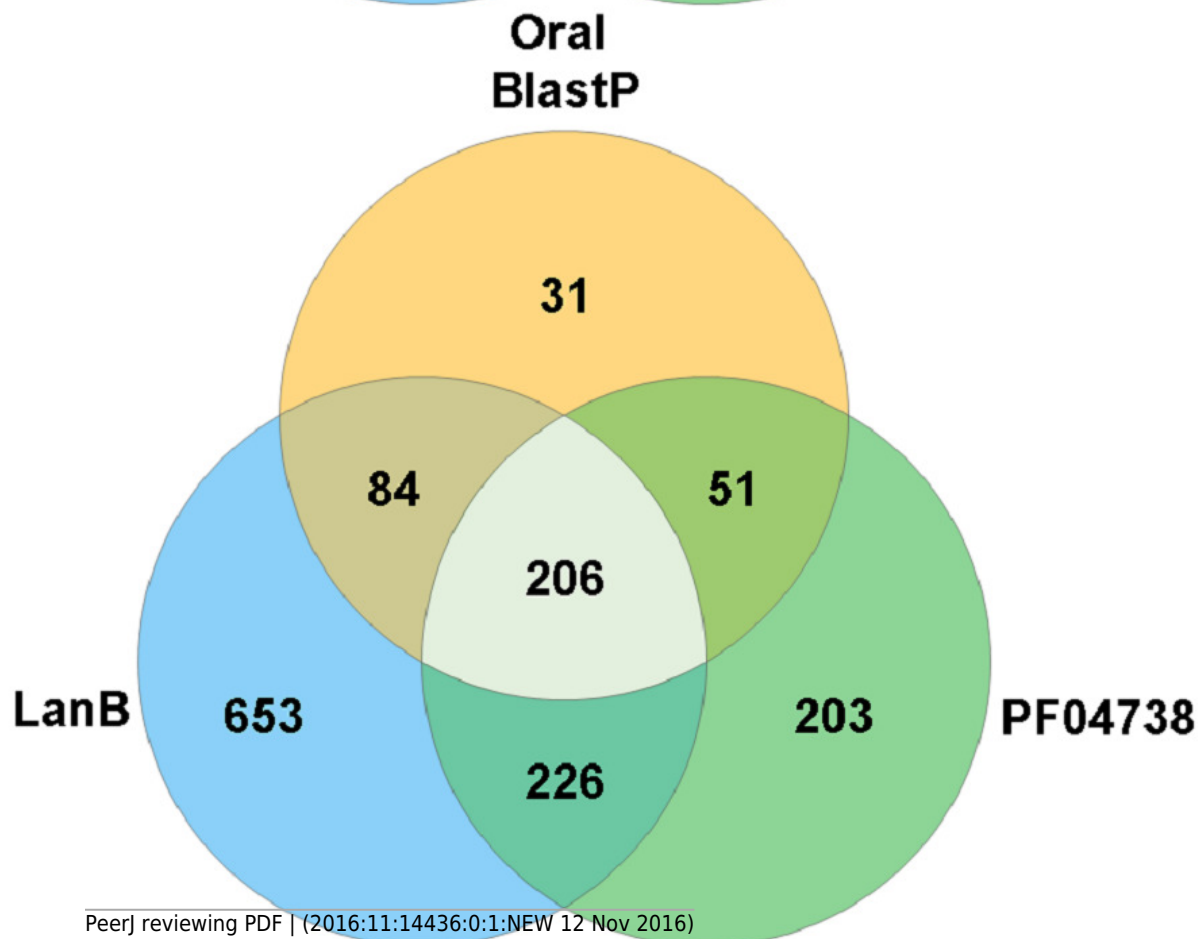


Figure 5

Comparison of lanthionine dehydratase density by body site reported by all three methods. Insert shows overall comparison between stool and oral environments.

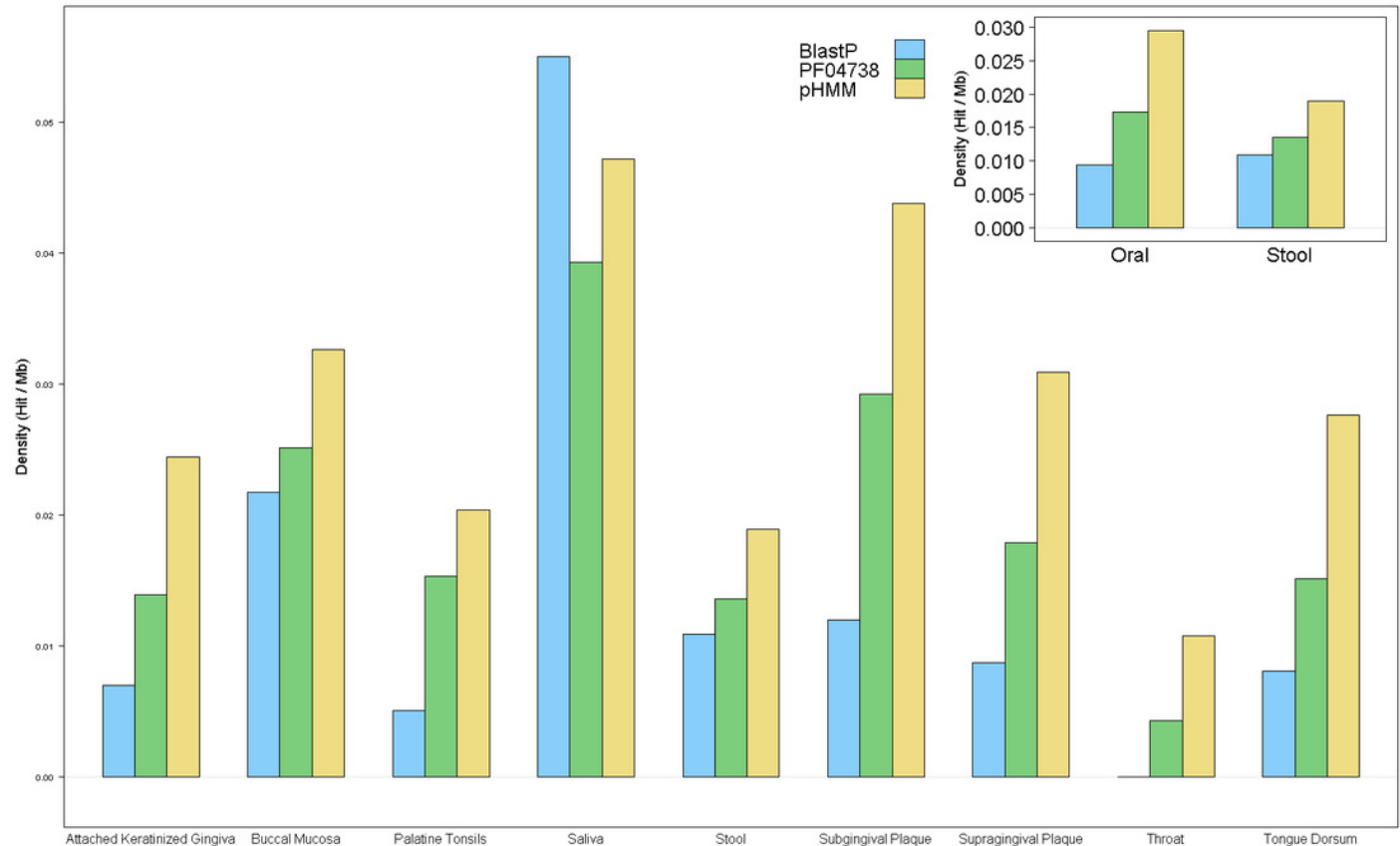


Figure 6

BAGEL3 output of three putative bacteriocin gene clusters identified from the gastrointestinal tract subset of the Human Microbiome Project's reference genome database by our new profile HMM.

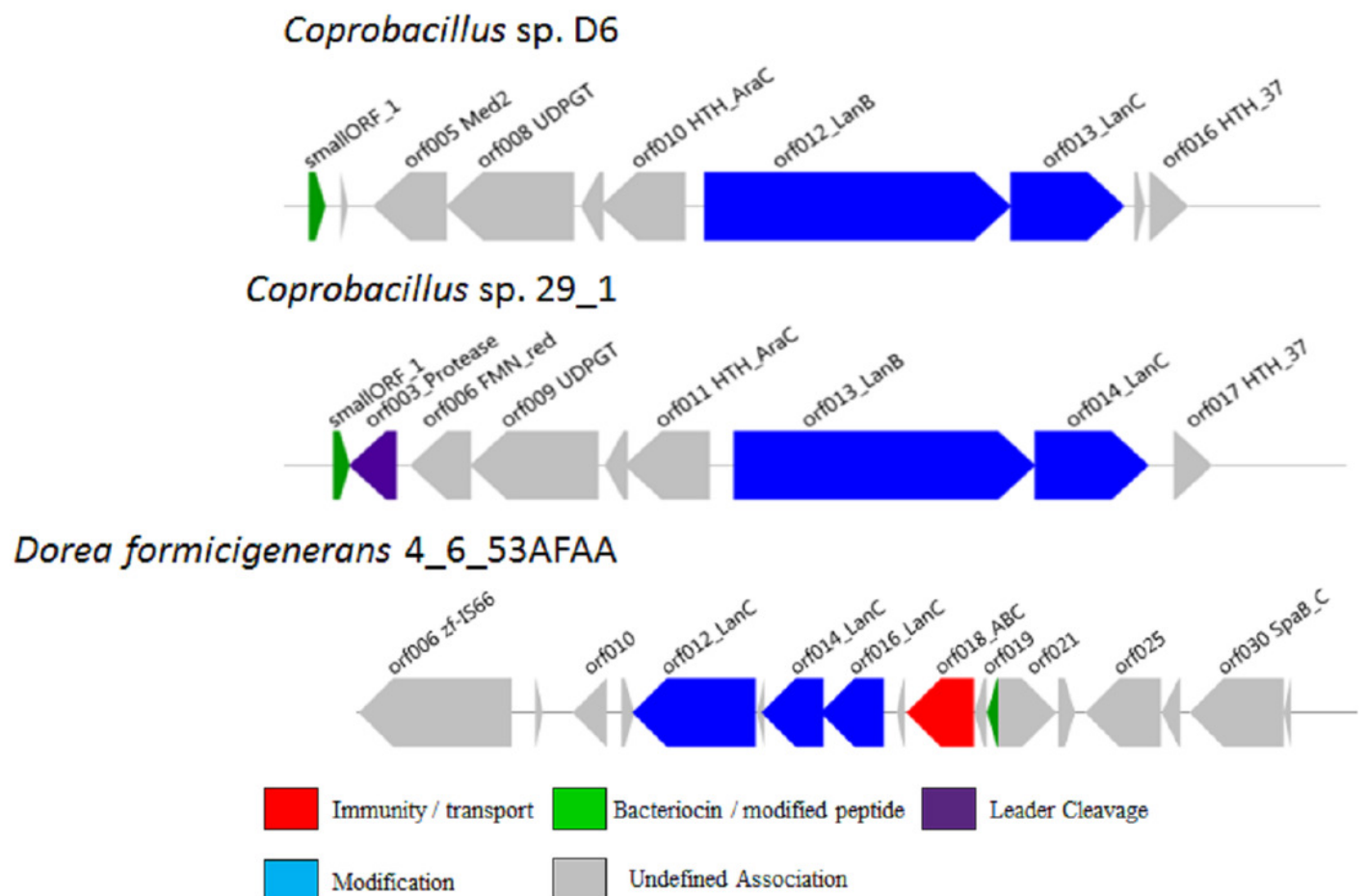


Table 1 (on next page)

Controls used in validation of the profile HMM.

^a Relevant lanthionine dehydratase protein was correctly identified by our model

^b Relevant lanthionine dehydratase protein was correctly identified by PF04738

^c No lanthionine dehydratase protein identified by either model

Strain	Bacteriocin	Class
<i>Lactococcus lactis</i> ssp. <i>lactis</i> S0 ^{a,b}	Nisin Z	LanB
<i>Lactococcus lactis</i> ssp. <i>lactis</i> CV56 ^{a,b}	Nisin A	LanB
<i>Lactococcus lactis</i> ssp. <i>lactis</i> IO-1 ^{a,b}	Nisin Z	LanB
<i>Bacillus subtilis</i> subsp. <i>spizizenii</i> ATCC 6633 ^{a,b}	Subtilin	LanB
<i>Staphylococcus aureus</i> subsp. <i>aureus</i> USA300_FPR3757 ^a	Bsa	LanB
<i>Streptococcus mutans</i> CH43 ^{a,b}	Mutacin I	LanB
<i>Streptococcus mutans</i> UA787 ^{a,b}	Mutacin III	LanB
<i>Streptococcus pyogenes</i> ^{a,b}	Streptin	LanB
<i>Staphylococcus epidermidis</i> ^{a,b}	Pep5	LanB
<i>Lactococcus lactis</i> subsp. <i>lactis</i> KF147 ^c	None	-
<i>Streptococcus mutans</i> GS-5 ^c	Mutacin GS-5	LanM
<i>Lactococcus lactis</i> subsp. <i>lactis</i> plasmid pES2 ^c	Lacticin 481	LanM
<i>Streptomyces cinnamoneus cinnamoneus</i> DSM 4005 ^c	Cinnamycin	LanM
<i>Streptomyces venezuelae</i> ATCC 10712 ^c	Venezuelin	LanL

Table 2 (on next page)

Samples per body site screened.

Site	Number of Samples
Attached Keratinized Gingiva	6
Buccal Mucosa	107
Palatine Tonsils	6
Saliva	3
Stool	139
Subgingival Plaque	7
Supragingival Plaque	118
Throat	7
Tongue Dorsum	128

1

Table 3 (on next page)

Lanthionine dehydratase proteins identified in the gastrointestinal tract subset of the Human Microbiome Project's reference genome database using our profile HMM.

Accession	Strain	E Value
JH414709	<i>Bacillus</i> sp. 7_6_55CFAA_CT2	9.0E-16
GL636578	<i>Coprobacillus</i> sp. 29_1	3.7E-67
AKCB01000002	<i>Coprobacillus</i> sp. D6	4.5E-68
JH126516	<i>Dorea formicigenerans</i> 4_6_53AFAA	2.3E-81
ACEP01000029	<i>Eubacterium hallii</i> DSM3353	9.4E-27
KI391961	<i>Fusobacterium nucleatum</i> subsp. <i>animalis</i> 3_1_33	2.2E-09
GG657999	<i>Fusobacterium</i> sp. 4_1_13	7.1E-09

1