

1 **Machine learning pattern recognition and differential network**  
2 **analysis of gastric microbiome under proton pump inhibitor**  
3 **treatment or *Helicobacter pylori* infection**

4  
5 Claudio Durán<sup>1,\*</sup>, Sara Ciucci<sup>1,\*</sup>, Alessandra Palladini<sup>2,3,\*</sup>, Umer Z. Ijaz<sup>4</sup>, Antonio G. Zippo<sup>5</sup>,  
6 Francesco Paroni Sterbini<sup>6</sup>, Luca Masucci<sup>6</sup>, Giovanni Cammarota<sup>7</sup>, Gianluca Ianiro<sup>7</sup>, Pirjo  
7 Spuul<sup>8</sup>, Michael Schroeder<sup>9</sup>, Stephan W. Grill<sup>9,10</sup>, Bryony N. Parsons<sup>11</sup>, D. Mark Pritchard<sup>11,12</sup>,  
8 Brunella Posteraro<sup>6</sup>, Maurizio Sanguinetti<sup>6</sup>, Giovanni Gasbarrini<sup>7</sup>, Antonio Gasbarrini<sup>7</sup>, and  
9 Carlo Vittorio Cannistraci<sup>1,13,§</sup>

10

11 <sup>1</sup>Biomedical Cybernetics Group, Biotechnology Center (BIOTEC), Center for Molecular and  
12 Cellular Bioengineering (CMCB), Center for Systems Biology Dresden (CSBD), Department  
13 of Physics, Technische Universität Dresden, Dresden, Germany;

14 <sup>2</sup>Paul Langerhans Institute Dresden, Helmholtz Zentrum Munchen, Carl Gustav Carus,  
15 Technische Universität Dresden, Dresden, Germany;

16 <sup>3</sup>German Center for Diabetes Research (DZD e.V.), Neuherberg, Germany;

17 <sup>4</sup>Department of Infrastructure and Environment University of Glasgow, School of  
18 Engineering, Glasgow, UK;

19 <sup>5</sup>Institute of Neuroscience, Consiglio Nazionale delle Ricerche, Milan, Italy;

20 <sup>6</sup>Institute of Microbiology, Università Cattolica del Sacro Cuore, Rome, Italy;

21 <sup>7</sup>Internal Medicine and Gastroenterology Unit, Università Cattolica del Sacro Cuore, Rome,  
22 Italy;

23 <sup>8</sup>Department of Chemistry and Biotechnology, Division of Gene Technology, Tallinn  
24 University of Technology, Tallinn 12618, Estonia;

25 <sup>9</sup>Biotechnology Center (BIOTEC), Center for Molecular and Cellular Bioengineering (CMCB),  
26 Technische Universität Dresden, Dresden, Germany;

27 <sup>10</sup>Max Planck Institute of Molecular Cell Biology and Genetics, Pfotenhauer Str. 108, 01307  
28 Dresden, Germany;

29 <sup>11</sup>Department of Cellular and Molecular Physiology, Institute of Translational Medicine,  
30 University of Liverpool, Liverpool, UK;

31 <sup>12</sup>Department of Gastroenterology, Royal Liverpool and Broadgreen University Hospitals  
32 NHS Trust, Liverpool, UK;

33 <sup>13</sup>Complex Network Intelligence Lab, Tsinghua Laboratory of Brain and Intelligence, Tsinghua  
34 University, Beijing, China.

35

36 \*These authors contributed equally to this work.

37 §Correspondence should be addressed to: [kalokagathos.agon@gmail.com](mailto:kalokagathos.agon@gmail.com)

38

## 39 **Abstract**

40 The stomach is inhabited by diverse microbial communities, co-existing in a dynamic balance.  
41 Long-term use of drugs such as Proton Pump Inhibitors (PPIs), or bacterial infection such as  
42 *Helicobacter pylori*, cause significant microbial alterations. Yet, studies revealing how the  
43 commensal bacteria re-organize, due to these perturbations of the gastric environment, are in  
44 early phase and rely **principally** on linear techniques for multivariate analysis.

45 Here we disclose the importance of complementing linear dimensionality reduction techniques  
46 such as Principal Component Analysis and Multidimensional Scaling with nonlinear  
47 approaches **to unveil hidden patterns that with linear approaches remain unseen**. Then, we  
48 show the importance to complete multivariate pattern analysis with differential network  
49 analysis, to reveal mechanisms of re-organizations which emerge from microbial variations  
50 induced by a medical treatment (PPIs) or an infectious state (*H. pylori*). **Finally, we reveal**  
51 **metabolomic network alterations associated to the perturbed microbial communities**.

## 52 **Keywords**

53 Proton Pump Inhibitors – Dyspepsia – *Helicobacter pylori* – Gastric microbiota – Linear and  
54 nonlinear unsupervised methods – Minimum Curvilinear Embedding – Nonlinearity – PC-corr  
55 network – 16S rRNA

56

## 57 **Introduction**

58 The gastric environment with its microbiota is the active gate that regulates access to the whole  
59 gastrointestinal tract, and therefore it has a remarkable impact on the correct functionality of  
60 the entire human organism. Recent studies have revealed that many orally administered drugs  
61 can perturb the elegant balance of the gastric **microbiota**<sup>1,2</sup>. However, not all of them cause  
62 permanent adverse effects and particular attention should be addressed to drugs that are  
63 frequently prescribed and administered for long periods. They can cause permanent unbalance  
64 of the gastric microbiota that might generate adverse side effects for the patient's health. Since  
65 the introduction of proton pump inhibitors (PPIs) into clinical practice more than 25 years ago,  
66 PPIs have become the mainstay in the treatment of gastric-acid-related diseases<sup>3</sup>. PPIs are  
67 potent agents that block acid secretion by gastric parietal cells by binding covalently to and  
68 inhibiting the hydrogen/potassium (H<sup>+</sup>/K<sup>+</sup>)-ATPases (or proton pumps), and additionally they  
69 can bind non-gastric H<sup>+</sup>/K<sup>+</sup>-ATPases, both on human cells and on bacteria and fungi, such as  
70 *Helicobacter pylori* (*H. pylori*)<sup>4-6</sup>.

71 PPIs are drugs of first choice for peptic ulcers (PU) and their complications (e.g. bleeding),  
72 gastroesophageal reflux disease (GERD), nonsteroidal anti-inflammatory drug (NSAID)-  
73 induced gastrointestinal (GI) lesions, Zollinger-Ellison syndrome and dyspepsia<sup>3,7,8</sup>. In  
74 particular, dyspepsia is a common clinical problem characterized by symptoms (e.g. epigastric  
75 pain, burning, postprandial fullness, or early satiation) originating from the gastroduodenal  
76 region<sup>9</sup>. The potent gastric-acid suppression drugs PPIs can treat the most frequent causes of  
77 dyspepsia including GERD, medication-induced gastritis, and peptic ulcers, thus minimizing  
78 the need for costly and invasive testing, and moreover are currently recommended to eradicate  
79 *H. pylori* infection, in combination to antibiotics<sup>7,9,10</sup>. Nevertheless, some patients are resistant  
80 or partial responders to empiric PPI therapy, and continue to have dyspepsia<sup>7</sup>.

81 Additionally, there is growing evidence that these medications are associated with increased  
82 rates of pharyngitis and upper and lower respiratory tract infections<sup>11</sup>. Their long-term

83 overutilization has been associated with potential adverse effects. For instance: the  
84 development of corpus predominant atrophic gastritis in *H. pylori* positive patients (that is a  
85 precursor of gastric cancer), enteric infections (especially *Clostridium difficile*-associated  
86 diarrhoea), increased risk of fundic gland polyps, hypomagnesaemia and hypocalcaemia,  
87 osteoporosis and bone fractures, vitamin and mineral deficiency, pneumonia, acute interstitial  
88 nephritis, and increased risk of drug–drug interactions, among others <sup>7,12–15</sup>.

89 Consumption of such acid-suppressive medications has also been associated with changes in  
90 microbial composition and function of gut microbiota. More recent studies relying on amplicon-  
91 based metagenomic approaches, have shown that PPIs exert an effect on gastric, oropharyngeal,  
92 and lung microflora in children with a chronic cough <sup>11</sup>, and have a significant impact on the  
93 gut microbiome in healthy subjects, with an increase of oral and pharyngeal bacteria and  
94 potential pathogenic bacteria <sup>16,17</sup>. Furthermore, another study by Tsuda *et al.* <sup>18</sup> revealed that  
95 PPIs influence the bacterial composition of saliva, gastric fluid and stool in a cohort of adult  
96 dyspeptic patients. However, this latter study highlights how the influence of PPI administration  
97 on the fecal and gastric luminal microbiota is still controversial and further investigation is  
98 required to understand the interaction between PPIs and non-*H. pylori* bacteria. Hence, this  
99 represents the first reason that motivates the present study.

100 In fact, by irreversibly blocking H<sup>+</sup>/K<sup>+</sup>-ATPases, PPIs inhibit gastric acid secretion by gastric  
101 parietal cells, which results in a higher intragastric pH, meaning the microenvironment of this  
102 niche changes, hence allowing more bacteria to survive the gastric acid barrier <sup>4,5,16</sup>. The use of  
103 PPIs and higher gastric pH were indeed correlated with the overgrowth of non-*H. pylori*  
104 bacterial **microflora** in the stomach of patients with gastric-reflux and PPIs were shown to  
105 aggravate gastritis because of co-infection with *H. pylori* and non-*H. pylori* bacterial species  
106 <sup>4,14,19,20</sup>. However, PPIs may also affect the gastrointestinal microbiome through pH-  
107 independent mechanisms, by directly targeting the proton pumps of naturally occurring bacteria  
108 by binding P-type ATPases (e.g. *H. pylori*) <sup>4,6</sup>.

109 Attempts to detect patterns of PPI related gastrointestinal changes have been made in different  
110 studies <sup>21,22</sup> through linear multidimensional analysis techniques, such as Principal Component  
111 Analysis (PCA) and Multidimensional Scaling (MDS), also called Principal Coordinates  
112 Analysis (PCoA). Nevertheless, they failed to detect the effect of PPIs on gastric *fluid* samples  
113 <sup>21</sup>, nor any significant PPI-related modification in esophageal <sup>21</sup> and gastric <sup>22</sup> *tissue* samples.  
114 This represents the second reason that motivates our investigation. Are these controversial  
115 results due to complex patterns that cannot be detected using linear analysis?  
116 In this study, we show an unprecedented result: unlike linear approaches, Minimum Curvilinear  
117 Embedding (MCE) <sup>23</sup>, which is a technique for *nonlinear* dimension reduction, discriminated  
118 both the esophageal and the gastric tissue microbial profiles of patients taking PPI medications  
119 from untreated ones when re-analyzing the data published in the abovementioned studies. This  
120 finding demonstrates the importance of routinely integrating the use of nonlinear  
121 multidimensional techniques into clinical metagenomic studies, since addressing nonlinearity  
122 could significantly modify the results and conclusions. Indeed, the absence of separation by  
123 means of linear transformations does not imply absence of separation in general, and nonlinear  
124 techniques could prove it, especially in complex datasets such as the ones generated in  
125 metagenomics 16S rRNA. As a matter of fact, the high throughput profiling of bacteria is  
126 frequently used in clinical studies, thus posing a challenge to efficient information retrieval:  
127 understanding how microbial community structure affects health and disease can indeed  
128 contribute to better diagnosis, prevention, and treatment of human pathologies <sup>24</sup>.  
129 The common practice in unsupervised dimension reduction data analysis is to consider only the  
130 first two (or three, less used) dimensions of mapping, and the goal is to visually explore the  
131 distribution of the samples and the incidence of significant patterns <sup>25</sup>. **This type of analysis is**  
132 **advantageous to validate hypothesis or to generate new ones. In addition,** this procedure is  
133 particularly useful in case of studies with small size datasets <sup>23</sup>, **or for imbalance class**

134 **samples**, to obtain unbiased (the labels are not used) confirmation of the separation between  
135 groups of samples for which diversity is theorized or expected.

136 **In addition, we will provide an analysis with two nonlinear algorithms for dimensionality**  
137 **reduction often used in literature, namely Isomap<sup>26</sup> and t-SNE<sup>27,28</sup>. These methods,**  
138 **although unsupervised, need hyperparameters optimization. Indeed, Isomap needs as**  
139 **input a parameter related to ‘k’ number of neighbours to construct a network, whereas**  
140 **t-SNE needs the perplexity and number of dimensions (or components). Different values**  
141 **of these parameters may lead to different results, which represent a challenge in an**  
142 **unsupervised scenario where automatic and label-free selection of the best solution is**  
143 **wished. This is the reason why this study will focus mainly on parameter-free**  
144 **dimensionality reduction techniques, whereas Isomap and t-SNE results will be shortly**  
145 **considered for a specific dataset in the result section.**

146 Here, we will specifically analyse the many aforementioned 16S rRNA amplicons datasets to  
147 address the following pattern recognition questions: (1) Is PPI treatment affecting change on  
148 the microbiota of esophageal and gastric tissues in dyspeptic patients, regardless of the initial  
149 pathological infection due to *H. pylori*? (2) Is this PPI-induced change so dominant as to result  
150 in a discernible pattern in the first two dimensions of mapping by unsupervised dimension  
151 reduction? (3) Are linear techniques sufficient to bring out patterns in complex microbial data?  
152 Furthermore, using differential network analysis we will address from the systems point of view  
153 these other questions: (4) How is PPI affecting the microbiota in the gastric environment in  
154 dyspeptic patients? (5) What is the effect of *H. pylori* infection on gastric mucosal microflora?  
155 Both factors (PPI treatment and *H. pylori* infection) can influence the composition of the gastric  
156 microbiota, and this further analysis will help to understand the general (overall) behaviour of  
157 the microbial ecosystem under these conditions. Ultimately, this means that we will try to  
158 clarify and visualize via network representation how the bacterial cooperative organization is

159 systemically altered either by the use of this acid suppressant drug in the gastric environment  
160 under dyspepsia, or by *H. pylori* infection in the gastric mucosa.

161

## 162 **Methods**

### 163 ***Dataset description***

#### 164 *Amir3 (esophageal mucosa)*

165 The 16S rRNA gene sequences were generated by Amir and colleagues <sup>21</sup> and are publicly  
166 available via the MG RAST database (<http://metagenomics.anl.gov/linkin.cgi?project=5767>).

167 The dataset was obtained from 16 esophageal mucosal biopsies of eight individuals before and  
168 after eight weeks of PPI treatment. Two patients with heartburn presented normal  
169 oesophagogastroduodenoscopy (H) indicating that they present healthy oesophageal tissues but  
170 are exposed to gastric refluxate, four patients had oesophagitis (ES) and two had Barrett's  
171 oesophagus (BE). **Metagenomics data** were obtained by pyrosequencing 16S rRNA **gene**  
172 amplicons on the GS FLX system (Roche). Data were processed by replicating the  
173 bioinformatics workflow followed by Amir and colleagues <sup>21</sup> in order to obtain the matrix of  
174 the bacterial absolute abundance: sequence reads were analysed with the pipeline Quantitative  
175 Insights into Microbial Ecology (QIIME) v. 1.6.0 <sup>29</sup> using default parameters (sequences were  
176 removed if shorter than 200 nt, if they contained ambiguous bases or uncorrectable barcodes,  
177 or if the primer was missing). Operational Taxonomic Units (OTUs), that are clusters of  
178 sequences showing a pairwise similarity no lesser than 97%, were identified using the UCLUST  
179 algorithm (<http://www.drive5.com/usearch/>). The most abundant sequence in each cluster was  
180 chosen as the representative of its OTU, and this representative set of sequences was then used  
181 for taxonomy assignment by means of the Bayesian Ribosomal Database Project classifier <sup>30</sup>  
182 and aligned with PyNAST103. Chimeras, that are PCR artefacts, were identified using  
183 ChimeraSlayer <sup>31</sup> and removed. The Greengenes database, which was used for the annotation

184 of the reads, additionally identifies groups of bacteria that are supported by whole genome  
185 phylogeny, but are not yet officially recognized by the Bergeys taxonomy, which is the  
186 reference taxonomy and is based on physiochemical and morphological traits. This results in a  
187 special annotation for some taxa, like *Prevotella*, that thus appears both with the general  
188 annotation, that is *Prevotella*, and with the special annotation, that is between square brackets,  
189 [*Prevotella*].

190

#### 191 *Amir4 (gastric fluid)*

192 The dataset was generated by Amir and colleagues <sup>21</sup>, and is public and available in the MG  
193 RAST database (<http://metagenomics.anl.gov/linkin.cgi?project=5732>). It comprises eight  
194 patients, whose gastric fluid was sampled at two different time points, that is before PPI  
195 treatment and after eight weeks of PPI treatment, for a total of 16 samples. The patients are the  
196 same described in Amir3. **Metagenomics data** were obtained by pyrosequencing fragments of  
197 the 16S rRNA gene **amplicons** on the GS FLX system (Roche). Then the data were processed  
198 by replicating the same bioinformatics workflow followed by Amir and colleagues <sup>21</sup> that was  
199 described in the previous data description (Amir3), in order to obtain the matrix of the bacterial  
200 absolute abundance. As for Amir3, the Greengenes database was used for the annotation of the  
201 reads.

202

#### 203 *Paroni Sterbini (gastric mucosa)*

204 The dataset was generated by Paroni Sterbini and colleagues <sup>22</sup>, and is public and available in  
205 the NCBI Sequence Read Archive (SRA) (<http://www.ncbi.nlm.nih.gov/sra>, accession number  
206 SRP060417), where all details pertaining the sequencing experimental design are also reported.  
207 It contains 24 biopsy specimens of the gastric antrum from 24 individuals who were referred to  
208 the Department of Gastroenterology of Gemelli Hospital (Rome) with dyspepsia symptoms (i.e.  
209 heartburn, nausea, epigastric pain and discomfort, bloating, and regurgitation). Twelve of these



210 individuals (PPI1 to PPI12) had been taking PPIs for at least 12 months, while the others (S1  
211 to S12) were not being treated (naïve) or had stopped treatment at least 12 months before sample  
212 collection. In addition, 9 (5 treated and 4 untreated) were positive for *H. pylori* infection, where  
213 *H. pylori* positivity or negativity was determined by histology and rapid urease tests.  
214 **Metagenomics data** were obtained by pyrosequencing fragments of the 16S rRNA gene  
215 **amplicons** on the GS Junior platform (454 Life Sciences, Roche Diagnostics). Then the  
216 sequence data were processed by replicating the bioinformatics workflow followed by Paroni  
217 Sterbini *et al.*<sup>22</sup>, in order to obtain the matrix of the bacterial absolute abundance.

218

219 *Parsons (gastric mucosa)*

220 The dataset was generated by Parsons and colleagues<sup>32</sup>, and is public and available in the EBI  
221 short-read archive (the European Nucleotide Archive, ENA) (<https://www.ebi.ac.uk/ena>,  
222 accession number PRJEB21104). In the original study, the authors focused on the analysis of  
223 gastric biopsy samples of 95 individuals (in groups representing normal stomach, PPI treated,  
224 *H. pylori*-induced gastritis, *H. pylori*-induced atrophic gastritis and autoimmune atrophic  
225 gastritis), selected from a larger prospectively recruited cohort patients who underwent  
226 diagnostic upper gastrointestinal endoscopy at Royal Liverpool University Hospital<sup>32</sup>. RNA  
227 extracted from gastric corpus biopsies was analysed using 16S rRNA sequencing (MiSeq).  
228 Then the sequence analysis was performed, as described by the authors in the supplementary  
229 methods of the original article<sup>32</sup>. Here we focused on the analysis of gastric biopsy specimens  
230 (in total 42 samples) from normal stomach group (20 patients) and belonging to the *H. pylori*  
231 gastritis group (22 patients). As described in<sup>32</sup>, patients in the normal stomach group showed  
232 normal endoscopy, no evidence of *H. pylori* infection by histology, rapid urease test or serology,  
233 were not treated by PPI and were normogastrinaemic. Patients in the *H. pylori* gastritis group  
234 were instead positive to *H. pylori* infection by urease test, histology and serology, were not  
235 taking PPI medication and were normogastrinaemic.

236

237 ***Data exploration and visualization: the reason for unsupervised dimension***  
238 ***reduction***

239 The main reason to perform an unsupervised dimension reduction is to explore and visualize  
240 the most relevant sample patterns that should emerge in the first two dimensions of embedding  
241 (which represent the information of higher variability in the data) from the hidden  
242 multidimensional space of a dataset. The fact that the sample labels (if known) are not used for  
243 the data projection makes the analysis unsupervised. The advantage of performing an  
244 unsupervised analysis is both for data quality checking and to gather the main trends hidden in  
245 the data, independently from any hypothesis or knowledge available on the samples. This is  
246 particularly useful to discover the presence of interesting sub-groups inside the studied cohort  
247 or to detect the influence of confounding factors.

248 A final interesting advantage offered by unsupervised analysis is in small size datasets, where  
249 the number of samples  $n$  is significantly lower than the number of features  $p$ , a condition that  
250 unfortunately occurs in several metagenomic studies. When  $n \ll p$  the application of supervised  
251 approaches can become problematic, because the supervised procedure of parameter learning  
252 can suffer from overfitting<sup>23,33,34</sup>.

253 The mainstream multivariate methods to unsupervisedly explore data patterns in metagenomic  
254 studies are based on linear dimension reduction, in particular PCA<sup>35,36</sup> and MDS<sup>37,38</sup>, also  
255 known as PCoA, methods that have been used to explore and visualize data structure in many  
256 metagenomic studies, from sponge<sup>39,40</sup> to gastric tissue microbiota<sup>22</sup>. These tools perform a  
257 dimension reduction of the data either by *multidimensional variance analysis* (for instance  
258 PCA) or *dissimilarity embedding* (for instance MDS/PCoA). PCA collects uncorrelated  
259 variance in the multidimensional space, creating new synthetic orthogonal variables, which are  
260 linear combinations of the original ones, then plots the samples in a reduced space using the  
261 new variables that embody the largest orthogonal variances. MDS computes dissimilarities

262 between every pair of samples, plotting the Euclidean part of these dissimilarities as distances  
263 between every pair of points (MDS) in a reduced space, in this way the linear part of the sample  
264 relations can be represented.

265

## 266 *PCA, MDS (or PCoA) and LDA*

267 Below, we report some of the PCA major advantages and drawbacks, that were pinpointed in a  
268 recent study on multidimensional population genomics <sup>41</sup>, and of other conventional  
269 dimensional reduction techniques employed for the analysis of metagenomic data.

270 PCA is time-efficient, parameter-free and straightforward to interpret, yet it strives to resolve  
271 structure in datasets with few samples and highly numerous features, which enclose nonlinear  
272 patterns. Therefore, PCA can occasionally fail to reveal differences among samples, even when  
273 differences are known a-priori, which means it can also miss represent hidden nonlinear  
274 relations among the samples in the feature space. For instance, see the illustration of the PCA  
275 two-dimension reduction mapping of the Tripartite-Swiss-Roll dataset in Suppl. Fig. S1B. PCA  
276 clearly fails to unfold and reveal the structure of the three separated groups of samples.

277 MDS, on the other hand, preserves the sample distances in a 2D-space based on the calculation  
278 of a distance matrix (Suppl. Fig. S1C,D). In ecology, distance (or dissimilarity) matrices are a  
279 major way to transpose the ecological information of samples in terms of their species  
280 composition and abundance <sup>42,43</sup>. In this article we will consider classical MDS (which uses  
281 Euclidean distance and is in practice equivalent to PCA <sup>44,45</sup>), and non-metric MDS (NMDS)  
282 obtained according to Sammon's Mapping <sup>46</sup>. In the latter, the elements of the multivariate  
283 space are mapped onto a lower dimensional space while retaining the original inter-point  
284 dissimilarities, by means of a nonlinear, but monotonic transformation (Sammon Mapping).  
285 Since it respects the ranking of dissimilarities, it tends to linearize the relationships between the  
286 samples. In addition, MDS will be performed also according to Bray-Curtis (MDSbc)  
287 dissimilarity and weighted UniFrac (MDSwUF) distance because they are considered the

288 reference in metagenomics studies. Bray-Curtis dissimilarity quantifies how dissimilar two sites  
289 (samples) are based on counts (bacterial abundances), where 0 means two samples are identical  
290 and 1 means that the two samples do not share any taxa<sup>47,48</sup>. Dissimilarly, the UniFrac distance,  
291 either unweighted (qualitative) or weighted (quantitative), is the most popular phylogenetic  
292 distance measure for the microbial community diversity between different samples (also known  
293 as  $\beta$ -diversity<sup>49</sup>) and, differently from the previous discussed methods, uses the phylogenetic  
294 information (which is an external knowledge not contained in the dataset) on the taxa to  
295 compare samples. In particular, its weighted-version weights the branches of a phylogenetic  
296 tree based of the taxa abundance information<sup>50-53</sup>. Hence the weighted UniFrac distance  
297 directly accounts for differences in the abundance of different kinds of bacteria, and can be  
298 crucial to describe community changes<sup>51</sup> in the studied samples.

299 We need to specify that both MDSwUF and NMDS are in practice nonlinear methods and  
300 weighted UniFrac is not a classical unsupervised technique like the others. In fact, MDSwUF  
301 adopts a distance that combines the information given by the bacterial abundance of the dataset  
302 with the supervised prior (external) knowledge regarding the known hierarchical phylogenetic  
303 relationship among the bacteria. However, like PCA, MDS can fail to detect patterns if data are  
304 not properly linearized<sup>54</sup>. For instance, see Suppl. Fig. S1C-D where MDSbc and NMDS  
305 respectively fail to resolve the Tripartite-Swiss-Roll dataset. When we consider clinical **16S**  
306 **rRNA amplicons** data, this failure potentially reduces the chances of correctly pinpointing  
307 samples which may represent clinical subspecies, and thus remain undetected and undiagnosed.  
308 In brief, these methods are not efficient to perform *hierarchical embedding* directly from the  
309 abundance value, since hierarchies preserve tree-like structures, and tree-like structures follow  
310 a hyperbolic, thus nonlinear, geometry<sup>55-57</sup>. Only MDSwUF is able to account for nonlinear  
311 hierarchical organization, yet this is not directly inferred from the abundance values, but rather  
312 forced as a constraint of prior supervised knowledge on the phylogeny of bacteria. For this  
313 reason we cannot offer a test on the Tripartite-Swiss-Roll dataset.

314 In our analysis of the Paroni Sterbini dataset, we also showed the results of a supervised  
315 technique, Linear Discriminant Analysis (LDA), which uses the labels to perform dimension  
316 reduction. LDA aims to separate the samples into groups based on hyperplanes and describe  
317 the differences between groups by a linear classification criterion that identifies decision  
318 boundaries between groups <sup>37</sup>. This technique is not congruous (and sometimes statistically  
319 invalid) for small sample size datasets. The reason is that given the reduced sample size we  
320 cannot divide the dataset in a training and test set, which is a fundamental requirement of  
321 supervised methods such as LDA.

322

### 323 ***Minimum Curvilinear Embedding (MCE)***

324 In 2010, Cannistraci *et al.* <sup>23</sup> introduced the centred version of Minimum Curvilinear  
325 Embedding (MCE), which provided notable results in: i) visualisation and discrimination of  
326 pain patients in peripheral neuropathy, and the germ-layer characterisation of human organ  
327 tissues <sup>23</sup>; ii) discrimination of microbiota in sea sponges <sup>39</sup>; iii) embedding of networks in the  
328 hyperbolic space <sup>56</sup>; iv) stage identification of embryonic stem cell differentiation based on  
329 genome-wide expression data <sup>58</sup>. In this fourth example, MCE performance ranked first on 12  
330 different tested approaches (evaluated on 10 diverse datasets). More recently in 2013 <sup>33</sup>, the  
331 non-centred version of the algorithm, named ncMCE, has been used: i) to visualise clusters of  
332 ultra-conserved regions of DNA across eukaryotic species <sup>59</sup> ; ii) as a network embedding  
333 technique for predicting links in protein interaction networks <sup>33</sup>, outperforming several other  
334 link prediction techniques; iii) to unsupervisedly reveal hidden patterns related with gender  
335 difference and metabolic-disease risk-factors in lipidomic profiles extracted from human  
336 plasma samples <sup>60</sup>; iv) to unsupervisedly infer and visualize phylogenetic (hierarchical)  
337 relations directly from individual SNP profiles in human population genetics <sup>41</sup>. Finally, also  
338 applications in non-biological problems such as the unsupervised discrimination of bad from

339 good radar signals <sup>33</sup>, represented a proof of concept of the universality of MCE for addressing  
340 nonlinear investigation of data and signals in general. Also in the case of the metagenomics  
341 studies targeting sea sponges, <sup>39,40</sup>, both MCE and its non-centred variant <sup>23,33</sup> once again proved  
342 successful in detecting structure where PCA and MDS could not, or hardly find any. This is  
343 mainly because MCE/ncMCE are unsupervised and parameter-free topological machine  
344 learning for *nonlinear* dimensionality reduction and multivariate analysis, that are able to  
345 perform a *hierarchical embedding*.

346 This study stems from the intuition that MCE/ncMCE analysis could successfully reveal  
347 undetected patterns also in esophageal and gastric metagenomics data, where only unsupervised  
348 linear methods or classical nonlinear methods such as NMDS and MDSwUF had been used and  
349 had failed to achieve any clear-cut result <sup>21,22</sup>.

350 Minimum Curvilinearity (MC) <sup>23</sup>, the principle behind MCE and ncMCE, was invented with  
351 the aim to reveal nonlinear data structures also, and especially, in the case of datasets with few  
352 samples and many features. MC principle suggests that curvilinear (nonlinear) distances  
353 between samples may be estimated as pairwise distances over their Minimum Spanning Tree  
354 (MST), constructed according to a selected distance (Euclidean, correlation-based, etc.) in a  
355 multidimensional feature space (here the metagenomic data space). In this study, we considered  
356 Pearson-correlation based distance (refer to <sup>23</sup> for details on the way to compute the distance  
357 for the MST). The collection of all nonlinear pairwise distances forms a distance matrix called  
358 the MC-distance matrix or MC-kernel, which can be used as an input in algorithms for  
359 dimensionality reduction, clustering, classification and generally in any type of machine  
360 learning. In MCE and ncMCE, the MC-kernel (which is non-centred for ncMCE) is followed  
361 by dimensionality reduction using singular value decomposition (SVD), and then by the  
362 projection of the samples onto a two-dimensional space for visualisation and analysis. Thus,  
363 MCE/ncMCE is a form of nonlinear and parameter-free kernel PCA <sup>33</sup>. In the rest of the article  
364 we will simply use the name MCE to indicate both MCE and ncMCE, since the centring

365 transformation is related to the specific data pre-processing and will be specified for each  
366 dataset as a technical detail in the respective results' tables.

367

### 368 ***MCE to unsupervisedly infer and visualize phylogenetic (hierarchical) relations***

369 A previous study by Alanis-Lobato *et al.* <sup>41</sup> showed that MCE is automatically able to  
370 unsupervisedly infer and visualize phylogenetic (hierarchical) relations directly from individual  
371 SNP profiles in human population genetics. Precisely, ncMCE detected separation between  
372 ethnic groups and provided an ordering over the discriminating dimension that was related to  
373 the phylogenetic organization of these populations.

374 This ability of MCE to infer and visualize phylogenetic (hierarchical) relationships was  
375 confirmed in our study on the Paroni Sterbini *et al.* dataset <sup>22</sup> (see Results section- '*Gastric*  
376 *tissue dataset unsupervised analysis*'). As previously mentioned (see the previous section  
377 '*PCA, MDS (or PCoA) and LDA*'), MDSwUF uses a weighted Unifrac distance that combines  
378 the prior knowledge of the bacterial phylogenetic tree with the information given by the  
379 bacterial abundance. Here we show that MCE perform better than MDSwUF on the Paroni  
380 Sterbini *et al.* dataset, due to its ability to infer the (hierarchical) phylogenetic relationship  
381 among the bacteria directly from the bacterial abundance of the dataset, by performing a  
382 hierarchical embedding. Hence, MCE can be used to compare the composition of microbial  
383 communities in the studied samples, where the phylogenetic information is instead directly  
384 inferred from bacterial abundance, differently from MDSwUF.

385

### 386 ***Procedure to evaluate the performance of the dimension reduction algorithms***

387 The performance of the mentioned dimension reduction algorithms is evaluated as the ability  
388 to separate the samples in the first two dimensions of embedding since, as discussed above,  
389 **they are related with the treatment/infection response**. In order to quantitatively evaluate  
390 the performance, we use a recently proposed index **termed Projection Separability Index**

391 **(PSI) used** for sample separation <sup>61</sup>. This index can be defined for any separation-measure and  
392 in this study we considered well-known measures: Area Under the ROC-Curve (**PSI-ROC**) and  
393 Area Under the Precision-Recall curve (**PSI-PR**), that are regularly used to quantitatively  
394 measure the performance of a binary predictor.

395 More precisely, in the 2D space a line is drawn between the centroids of the two groups that are  
396 compared, subsequently all the points are projected on this line and then AUC and AUPR are  
397 computed for the projected points. This new index can actually be applied not only in a 2D  
398 space, but in any N dimensional space. For the calculation of the centroids we consider the 2D-  
399 median of each cluster/class's group. In case more than two groups are present in a dataset, all  
400 the AUC and AUPR **values** between the possible pair-groups are computed. **Then, the**  
401 **following formula is applied:  $E/(1 + \delta)$ , where E is the mean of the pairwise PSI values and**  
402  **$\delta$  their standard deviation. Thus, the standard deviation works as a penalization in case**  
403 **of outliers PSI values, ensuring that the overall PSI is high only when all pairwise PSI**  
404 **values are close to the mean. The computed values are finally** chosen as an overall estimator  
405 of separation between the groups in the 2D reduced space. This case applies only to the Paroni  
406 Sterbini dataset, which is composed of three or, possibly, four groups of samples. All the other  
407 datasets are instead composed of two groups.

408 It is important to note that the **PSIs were** also applied to the data in the original high-  
409 dimensional (HD) space, as a reference to see how good the unsupervised dimension reduction  
410 approaches are in preserving the original group separability of the HD space.

411 All the algorithms were tested considering (when allowed by the dimension reduction method)  
412 data centring or non-centring. In addition, multiple normalization options were investigated and  
413 the datasets were considered under a certain type of normalization: division by the column -  
414 which reports the OTU - sum (indicated by DCS); division by the row - which reports the  
415 sample - sum (indicated by DRS); function  $\log_{10}(1+x)$  applied to the dataset (indicated by  
416 LOG).



417 **In order to verify that the performances obtained by our evaluations using PSI on the DR**  
418 **techniques are not obtained by chance, we calculate a measure termed trustworthiness,**  
419 **which exploits a resampling technique based on label-reshuffling to build a null model**  
420 **(Suppl. Fig. S2). The labels are reshuffled uniformly at random on the embedded points**  
421 **whose location is maintained unaltered in the reduced space. For each random reshuffling**  
422 **(the total number of reshuffling is a resampling parameter decided by the user, we**  
423 **adopted 1000 realizations), a PSI measure value is computed. The collection of all these**  
424 **values is used to draw the null model distribution. This distribution is employed to**  
425 **compute the probability to get at random a separation equal or higher than the one**  
426 **detected by using the original labels.**

427

### 428 *From Markov Clustering (MCL) to Minimum Curvilinear Markov Clustering* 429 *(MC-MCL)*

430 MCL is an unsupervised algorithm for the clustering of weighted graphs based on simulations  
431 of (stochastic) flow in graphs <sup>62</sup> (<http://micans.org/mcl/>). By varying a single parameter called  
432 inflation (with values between 1.1 and 10), clustering patterns on different scales of granularity  
433 can be detected. For clustering samples of a multidimensional dataset, the workflow starts with  
434 the computation of correlations (generally Pearson correlations) between the samples, and  
435 creates an edge between each pair of samples, where the edge-weight assumes the value of the  
436 respective pairwise positive sample correlation, or values zeros in case of negative correlations.  
437 This generates a weighted correlation graph (network), which is used as a map to simulate  
438 stochastic flows and detect the structural organization of clusters in the graph.

439 With the purpose of creating and testing a nonlinear variant of the MCL algorithm, we adopt  
440 an innovative algorithm which was recently proposed and called MC-MCL <sup>63</sup>. The idea is the  
441 following. The MC-kernel – discussed above in the MCE section - is a nonlinear distance matrix  
442 (or kernel) that expresses the pairwise relations between samples as a value of distance: small

443 samples distance indicates sample similarity, while large samples distance indicates sample  
444 dissimilarity. Here we reverse (using the following function:  $f(x) = 1 - x$ ) and after this we  
445 put to zero the negative values - **strategy already applied in the original MCL algorithm** -  
446 of the *MC-distance* kernel to get a *MC-similarity* kernel, where small values (close to zero)  
447 indicate low sample similarity and large values (close to one) indicate high sample similarity.  
448 A technical detail: for the computation of the MC-distance kernel, it is necessary to firstly  
449 square root the original distances (correlation-based) between the samples. As already  
450 investigated in <sup>23</sup>, this attenuates the estimation of large distances and amplifies the estimation  
451 of short distances; consequently it helps to regularize the nonlinear distances inferred over the  
452 MST in order to subsequently use them for message passing <sup>23</sup> (such as affinity propagation) or  
453 flow simulation (such as MCL) clustering algorithms.

454 Then, the standard stochastic flow simulations of MCL algorithm runs on the graph weighted  
455 with the values of the MC-similarity kernel (which collects pairwise *nonlinear* associations  
456 between samples) instead of the Pearson-correlation kernel (which collects pairwise *linear*  
457 associations between samples). In practice, this is a new algorithm for clustering that is a  
458 nonlinear version (based on the MC-kernel) of the classical MCL. The goal of the MC-MCL  
459 analysis is to verify whether the use of the MC-kernel improves performance, by solving  
460 nonlinearity, not only in dimension reduction (such as in MCE) but also in clustering (such as  
461 in MC-MCL).

462

### 463 ***Procedure to evaluate the performance of clustering algorithms***

464 The clustering algorithms MCL and MC-MCL were applied to the datasets, either raw, or after  
465 the same normalization procedures used before dimensionality reduction (DCS: division by  
466 column (OTU) sum; DRS: division by row (sample) sum; LOG: function  $\log_{10}(1+x)$  applied  
467 to the dataset) and their performance was evaluated by means of accuracy. The accuracy is

468 computed as the ratio of the number of samples assigned to the correct clusters over the total  
469 number of samples. For both MCL and MC-MCL, we tested Pearson and Spearman correlations  
470 to build the similarity measure to feed into the clustering methods. The Spearman correlation  
471 can also detect a subclass of nonlinear associations (which have monotonic shape function) or  
472 correct for outliers. Differently from what suggested for large gene datasets with thousands of  
473 samples in <sup>62</sup> (<http://micans.org/mcl/>), in this study we had to consider the whole set of original  
474 positive correlations without applying any threshold (cut-off) to the values. This was  
475 compulsory, since we considered datasets with few samples. In our case, to keep the graph  
476 connected, with one unique connected component, we could not introduce any kind of threshold  
477 that would otherwise alter the real graph connectivity (dividing the graph in disconnected  
478 components) and hence the clustering result. Since the MCL algorithm needs a single input  
479 parameter (inflation) to control the granularity of the output clustering, we ran it for different  
480 inflation values until we achieved the desired number of clusters. Finally, in the Paroni Sterbini  
481 *et al.* dataset <sup>22</sup> it was not clear in advance whether the correct number of clusters present in the  
482 multidimensional space was three or four. Hence, we tested the clustering algorithms  
483 considering as output both three and four clusters' configurations, and we identified as the best  
484 solution the one that offered the highest accuracy.

485

### 486 ***PC-corr network***

487 Furthermore, we investigated the effect of PPI on the microbiota of gastric fluid and gastric  
488 mucosa in dyspeptic patients, and the changes induced by *H. pylori* infection on the gastric  
489 mucosal microbiota, by means of the PC-corr approach <sup>64</sup>. PC-corr represents a simple  
490 algorithm that associates to any PCA segregation a discriminative network of features'  
491 interactions <sup>64</sup>. It is a method for linear multivariate-discriminative correlation network reverse  
492 engineering, that, thanks to its multivariate nature, can help to stress and squeeze out the  
493 underlying combinatorial and multifactorial mechanisms that generate the differences between

494 the studied conditions <sup>64</sup>. Said what PC-corr is able to do, now we offer an intuitive  
495 explanation of how it works. PCA is one of the most employed approaches to  
496 unsupervisedly map linear dissimilarities (hidden in the high-dimensional space) into a  
497 visible space of data representation. When we notice that two or more groups of samples  
498 are separated along one of the axes of this representation space, generally the first  
499 question is to discover what are the features that are contributing more to this separation.  
500 This is easily achievable by analysing the PCA loading values that are associated to the  
501 axis along which emerges the sample separation under investigation. But the loading  
502 values do not provide any information on how those features mutually interact. On the  
503 other hand, a correlation network between the features provide information on their  
504 associative relation but not on their contribution to the discrimination. PC-corr is an  
505 algorithm that is able to integrate together the discriminative information of the loadings  
506 with the combinatorial information of the correlations. Indeed, PC-corr offers as output  
507 a discriminative correlation network of features that can help to elucidate the possible  
508 associative mechanisms that are at support of the sample separation along a specific axis  
509 of PCA representation. Hence, for the studied datasets, it can be employed to point out the  
510 possible presence of bacterial alterations and their interplay, induced by a medical treatment  
511 (PPIs in dyspepsia) or infectious state (*H. pylori*).

### 512 *Bacteria-metabolite multilayer network construction and metabolite pathway* 513 *analysis*

514 We used a recently realized bacteria-metabolite bipartite network which is an open access  
515 resource <sup>65</sup> to infer the metabolic activity of the bacteria presented in the intersections of  
516 figures 6 and 7. The study <sup>65</sup> provided a large set of 9136 bacteria to metabolite  
517 interactions validated on experimental studies from mouse and human gastroenteric  
518 microbiota. It was available as a network, named NJC19, where node represented either  
519 bacteria or metabolites connected by several types of edges (e.g. production, consumption,

520 degradation). In this dataset we restricted the analysis to interactions found on human  
521 bacteria. Since the dataset identified bacteria according to the taxonomic levels of species  
522 while our data referred to the genus level, we made a new form of the NJC19 network  
523 with edges starting from the bacteria genus to metabolites. When we did not find any  
524 interactions for specific bacteria, we discarded them from further analyses. Therefore,  
525 from the list of intersected bacteria from figure 6 (*Porphyromonas*, *Capnocytophaga*,  
526 *Streptophyta*, *Granulicatella*, *Clostridiales*, *Oribacterium*, *Veillonella*, *Bulleidia*,  
527 *Fusobacterium*, *Leptotrichia*, *Campylobacter*, *Prevotella*) we dropped *Streptophyta*,  
528 *Granulicatella*, *Oribacterium*, *Bulleidia* and *Prevotella*. Similarly, from the intersected  
529 bacteria of figure 7 (*Enhydrobacter*, *Methylobacterium*, *Catonella*, *Pseudomonas*,  
530 *Acinetobacter*, *Sphingomonas*, *Propionibacterium*, *Bulleidia*) we dropped *Bulleidia*,  
531 *Catonella*, *Sphingomonas* and *Enhydrobacter*. For the graph representation, we used the  
532 color code already applied in the previous figures for the bacteria according to the  
533 taxonomic order. While for metabolites we classified them in 7 classes, assigning to each  
534 a different node shape and colour: vitamins, glycolysis, lipids, amino acids, carbohydrates,  
535 amines and miscellaneous. Furthermore, an enrichment analysis of metabolites (linked to  
536 the discriminative bacteria networks detected by PC-corr) has been conducted to unveil  
537 the metabolic pathways that might be associated to these bacteria perturbations. To this  
538 purpose, we used the framework provided by metaboloAnalyst suite <sup>66</sup>. Specifically, we  
539 performed the “Enrichment Analysis” against the KEGG database and we selected the  
540 significant pathways according to the Benjianini corrected p-values smaller than the  
541 significance level of 0.05. For the case of the *H. Pylori*-affected network, just few nodes  
542 were available and only one significant pathway was obtained from it with few metabolite  
543 hits. Therefore, the network was expanded by adding first neighbours metabolites –  
544 obtained from KEGG – from the current ones.

545 Finally, the metabolite layer network nodes were grouped according to the three most  
546 significant pathways in both the PPI- and *H. Pylori*-affected bacteria-metabolite  
547 networks. This was ensured according to the following procedure: a ranking was  
548 generated for the list of significant pathways in each of the two networks. Then, the nodes  
549 of each network were grouped according to the three pathways with the highest average  
550 ranking in the two networks, which in our study are: aminoacyl-tRNA biosynthesis;  
551 galactose metabolism; Alanine, aspartate and glutamate metabolism. A fourth group  
552 encapsulating the metabolites involved in the other significant pathways was also  
553 provided. Regarding the links considered in each bacteria and metabolite layer, the  
554 bacteria-bacteria associations were maintained from figures 6 and 7, whilst edges between  
555 metabolites were obtained by the metabolite interaction involved in the significant  
556 enriched KEGG pathways.

557 The processing pipeline has been developed in the R environment <sup>67</sup> and by using the  
558 following packages: *igraph* <sup>68</sup>, *taxize* <sup>69</sup>, *graphite* <sup>70</sup>, *RCy3* <sup>71</sup>.

559

### 560 *Computing platforms adopted to implement the algorithms*

561 Dimensionality reduction was performed in MATLAB on the abundance matrix of genus-level  
562 taxonomic assignments, with samples in rows and taxonomic assignments (OTUs) in columns.  
563 For MDSwUF, the computation of the weighted UniFrac distance was performed in R. We used  
564 the following MATLAB functions to calculate PCA, MDS and NMDS (Sammon Mapping)  
565 respectively: *svd*, *cmdscale* and *mdscale*. **For the calculation of the Theta YC distance, the**  
566 ***mothur* <sup>72</sup> approach was implemented in MATLAB.** For the calculation of Bray-Curtis  
567 dissimilarity, we used the function MATLAB *f\_braycurtis* in the Fathom Toolbox <sup>73</sup>  
568 (<http://www.marine.usf.edu/user/djones/matlab/matlab.html>). Instead, for the calculation of the  
569 weighted Unifrac distance for all sample pairs, we used the R function *UniFrac* in the phyloseq  
570 package (<https://bioconductor.org/packages/release/bioc/html/phyloseq.html>), after creating a

571 phyloseq-class object (with R function *phyloseq* in the same package) that contains both the  
572 abundance table (OTU table) and the phylogenetic tree. The MATLAB code for MCE/ncMCE  
573 is available online at: [https://sites.google.com/site/carlovittoriocannistraci/5-datasets-and-](https://sites.google.com/site/carlovittoriocannistraci/5-datasets-and-matlab-code/minimum-curvilinearity-ii-april-2012)  
574 [matlab-code/minimum-curvilinearity-ii-april-2012](https://sites.google.com/site/carlovittoriocannistraci/5-datasets-and-matlab-code/minimum-curvilinearity-ii-april-2012). For MCL clustering, we installed the  
575 MCL-edge software (<http://micans.org/mcl/>) in a Windows environment, following the  
576 procedure suggested by the authors in the software website. To apply this algorithm, we created  
577 a MATLAB function that generates automatically the input for MCL (equivalent to the  
578 `mcl` function in the software) and then uses a system call to run MCL in a UNIX-like  
579 environment (Cygwin, <https://www.cygwin.com/>). PC-corr method was performed in  
580 MATLAB on the abundance matrix of the genus-level taxonomic assignments, with samples in  
581 rows and taxonomic assignments in columns. The PC-corr algorithm is available as MATLAB  
582 function (as well as R function) at: [https://github.com/biomedical-cybernetics/PC-corr\\_net](https://github.com/biomedical-cybernetics/PC-corr_net).  
583 Then the obtained PC-corr and **bacteria-metabolite networks** were displayed by Cytoscape  
584 (<http://www.cytoscape.org/>).

585

## 586 **Results**

587 To answer the five questions stated in the Background section, we analysed the abovementioned  
588 16S rRNA gene sequencing datasets with information on PPI consumption in dyspeptic  
589 patients, following the flowchart shown in Fig. 1. **Our study is innovative at two different**  
590 **levels. At the more general ‘methodological level’, we introduce a new computational data**  
591 **mining pipeline (Fig.1) which explains how to overcome the limits of current multivariate**  
592 **analysis of small-size microbial data. At the more specific ‘technical level’, we propose**  
593 **innovative solutions in each of the 5 steps that composes this pipeline: dimension**  
594 **reduction, clustering, PC-corr networks, multilayer bacteria-metabolite networks and**  
595 **metabolic network pathways analysis. In the dimension reduction section, we innovate by**  
596 **illustrating the benefits to apply minimum curvilinear nonlinear machine learning**

597 **methods for dimension reduction. This is a completely new technical way to perform**  
598 **nonlinear analysis in the microbial field. In the clustering section, we propose MC-MCL,**  
599 **which is the first nonlinear version of Markov clustering and represents a novel nonlinear**  
600 **clustering approach. In the PC-corr section, we show how to extract valuable and robust**  
601 **information (that would otherwise be missed using standard procedure of analysis) across**  
602 **several (4 in total) small-size microbial datasets. In the fourth and fifth step we clarify**  
603 **how to enhance the biomedical interpretation with the aim to increase the impact of the**  
604 **findings on the scientific community.**

605 It is important to underline that, in one of the three initially analysed datasets (in Paroni Sterbini  
606 *et al.*<sup>22</sup>), we have the additional information on positivity or negativity to *H. pylori* infection. A  
607 fourth dataset (Parsons *et al.*<sup>32</sup>) is used only for the validation of the PC-corr network results  
608 and it contains not only information on PPI consumption but also additional information on  
609 positivity or negativity to *H. pylori* infection.

610 Unsupervised approaches were chosen for dimension reduction, and clustering because  
611 supervised (constrained) methods have been shown to perform poorly on small datasets, as  
612 explained in the paper by Smialowski *et al.*<sup>34</sup> and the work by Zagar and colleagues<sup>58</sup>.

613 Firstly, we performed unsupervised dimension reduction, both linear and nonlinear (described  
614 in the 'Methods- PCA, MDS (or PCoA) and LDA' and 'Methods- Minimum Curvilinear  
615 Embedding (MCE)') and we focused on the first two dimensions of embedding as **they are**  
616 **significantly related with the treatment/infection response (Suppl. Table S1)**. As we will  
617 show, linear techniques will fail to bring out the patterns in the microbial datasets, related to  
618 PPI-treatment. Instead, nonlinear dimension reduction will reveal the presence of hidden  
619 patterns related to PPI treatment. In particular, in the gastric biopsies dataset (Paroni Sterbini *et*  
620 *al.*<sup>22</sup>), nonlinear dimension reduction will point out the evidence of PPI perturbation. Secondly,  
621 clustering algorithms were applied to the studied datasets to confirm that the hidden patterns  
622 detected by nonlinear dimension reduction are well posed. Finally, the PC-corr algorithm<sup>64</sup> is



623 used to find the bacteria community (features) that make the difference between the patterns or  
624 groups, allowing our understanding of the PPI-induced and *H. pylori*-induced microbial  
625 perturbations.

626

## 627 **Gastric tissue dataset unsupervised analysis**

628 According to the questions formulated in our study, we are interested in an unsupervised  
629 approach to verify whether PPI drugs cause a major change in the gastric tissue microbiota of  
630 dyspeptic patients regardless of the initial pathological infection due to *H. pylori* <sup>22</sup>.

631 In our first analysis, we focused on the Paroni Sterbini *et al.* dataset <sup>22</sup> and, to facilitate the  
632 visualization of the sample separations in the 2D reduced space, we assigned: red colour to  
633 untreated dyspeptic patients without *H. pylori* infection (**H-**); green colour to untreated  
634 dyspeptic patients with *H. pylori* infection (**H+**); and blue colour to patients treated with PPI  
635 regardless of their *H. pylori* infection (**P**). However, to help to detect also the effect of the *H.*  
636 *pylori* infection we reported the labels close to each sample, with a '**&H+**' indicating the  
637 infection (**P&H+**) or a '**&H-**' indicating the absence of infection (**P&H-**). Finally, we also  
638 tested whether this separation into three main groups (**H-, H+, P**) is more truthful, from the  
639 metagenomics data standpoint, than the one in four groups (**H-, H+, P&H-, P&H+**).

640 Figure 2 shows the results of the multivariate techniques widely employed in metagenomic  
641 studies, PCA (Fig. 2A), MDSbc (Fig. 2B), MDSwUF (Fig. 2C), and NMDS (with Sammon  
642 Mapping) (Fig. 2D) (for more detail see the corresponding method section; the plots represents  
643 the best results based on **PSI-ROC** in Suppl. Table S2), which could only differentiate the  
644 group of untreated *H. pylori* positive samples (green dots) with respect to the group of untreated  
645 *H. pylori* negative samples (red dots), and PPI treated samples (blue dots), and no further  
646 separation is significantly detectable. **Considering the PSI results, the values are high (Table**  
647 **1 and Fig. 2)** (evaluated in the 2D embedding space, for details see '*Procedure to evaluate the*  
648 *performance of the dimension reduction algorithms*'). **PCA (PSI-ROC=0.85, PSI-PR=0.91)**

649 **and NMDS (PSI-ROC=0.85, PSI-PR=0.90) exhibit the highest PSI-ROC and PSI-PR**  
650 **values, followed by MDSwUF (PSI-ROC=0.84, PSI-PR=0.88) and MDSbc (PSI-**  
651 **ROC=0.81, PSI-PR=0.86).** Indeed, in all the plots there is a visible trend of separation between  
652 PPI-treated (blue dots) and untreated (red and green dots) samples, but this is not sufficient to  
653 declare the presence of the complete separation, and a manifest ‘crowding problem’<sup>33</sup> mixes  
654 the two cohorts together (blue and red dots). According to this output, the dataset appears to be  
655 strongly influenced by the presence of *H. pylori*, which is the predominant taxon (abundance >  
656 50%, Suppl. Table S3, percent abundance sheet) in four of the untreated *H. pylori* positive  
657 patients: where *H. pylori* is predominant, sample groups are quite close to one another and far  
658 from all the other samples in all four multivariate analyses (Fig. 2). Thus, PCA and MDS mainly  
659 show us that these **16S rRNA amplicons** separate according to *H. pylori* abundance, and there  
660 is no treatment-related pattern.

661 Non-centred MCE (Figure 3A, DCS normalization) was the best performing technique, with a  
662 **PSI-ROC of 0.91 and PSI-PR of 0.96** (Table 1) (for details see Suppl. Table S2). It even  
663 outperforms the nonlinear methods NMDS (Sammon Mapping) and MDSwUF, since **MCE** is  
664 automatically able to **unsupervisedly infer from data the underlying** (hierarchical)  
665 phylogenetic relationship among the bacteria. **MCE does not receive in input any**  
666 **phylogenetic information but directly infers it** from the bacterial abundance of the dataset  
667 by performing a hierarchical embedding, as already shown in the study of Alanis-Lobato *et al.*  
668 <sup>41</sup> (see ‘*Methods- MCE to unsupervisedly infer and visualize phylogenetic (hierarchical)*  
669 *relations*’). **The gain in performance compared with the rest of the dimensionality**  
670 **reduction techniques is relevant.**

671 **Indeed, the PSI-ROC improvement from 0.85 (PCA and NMDS) to 0.91 is not trivial. We**  
672 **want to stress that in general offering an AUC-ROC result that is higher than 0.9 is**  
673 **considered relevant in all scientific literature. Furthermore, as suggested by Ammirati et**  
674 **al.<sup>74</sup>, the same level of increase becomes more significant when being close to perfect**

675 segregation. This becomes evident when “quantifying the improvement in terms of the  
676 distance from the exact predictor”. As a didactic example, let us compare the current PSI-  
677 AUC improvement of 0.06 (0.85 – 0.91) against a case with a same hypothetical  
678 improvement but closer to randomness (0.50 – 0.56). In the former the relative  
679 improvement in respect to the exact predictor is 40% (computed as  $(0.91-0.85)/(1-$   
680  $0.85)*100$ ), whereas in the latter is 12% (computed as  $(0.56-0.50)/(1-0.50)*100$ ). Similarly  
681 for PSI-PR, MCE (PSI-PR=0.96) relative improvement from PCA (PSI-PR=0.91) in  
682 respect to the perfect predictor is 56% (computed as  $(0.96-0.91)/(1-0.91)*100$ ).

683 Furthermore, the MCE performance does not depend on its centring/non-centring, in fact the  
684 centred MCE version resolves the nonlinearity in the data too. Whereas, PCA regardless of  
685 being centred or non-centred does not resolve the nonlinearity in the data.

686 While MDS and PCA are confounded by the mixture of factors characterizing the samples and  
687 do not manage to resolve the differences between treated and untreated samples, non-centred  
688 MCE is the only technique that visibly separates samples by ordering them along the second  
689 dimension into three groups, detecting a treatment-related structure in the data (Fig. 3B). This  
690 is plausible, because in any non-centred embedding the first dimension points towards the  
691 centre of the manifold<sup>33</sup>, while the second dimension in the case of non-centred MCE represents  
692 the direction of higher topological nonlinear extension of the manifold. Interestingly, untreated  
693 *H. pylori* negative samples (red dots, **H-**) gather in the upper tail of the samples’ distribution,  
694 while treated samples (blue dots, **P**), both *H. pylori* test positive (**P&H+**) and negative (**P&H-**  
695 ), are mixed and show no other internal discernible groups. Untreated *H. pylori* positive samples  
696 (green samples, **H+**) gather at the bottom of the plot (Fig. 3A). Unlike the other approaches,  
697 non-centred MCE detects a treatment-related structure in the data and separates patients into  
698 three, not four, groups: PPI-treated, untreated *H. pylori* negative and untreated *H. pylori*  
699 positive. This last group appears as a subgroup marginally discriminating from the PPI-treated  
700 group and the topology of the samples seems to suggest that PPI treatment modifies the gastric

701 microbiota of *H. pylori*-negative patients with dyspeptic symptoms and gastric mucosa  
702 inflammation, shifting their gastric ecosystem in the same direction of PPI-treated *H. pylori*-  
703 positive patients. We speculate that the fact that PPI treatment and *H. pylori* infection determine  
704 the samples to gather in a similar position (i.e. out of the PPI-untreated/HP-negative group) in  
705 the non-centred MCE reduced space, indicates that both the PPI drugs and *H. pylori* induce an  
706 ecological change in the stomach, which might be driven by similar mechanisms. As a matter  
707 of fact, *H. pylori* can colonize the acidic lumen of the stomach thanks to its ability to hydrolyse  
708 urea into carbon dioxide (CO<sub>2</sub>) and ammonia (NH<sub>3</sub>)<sup>75</sup>, thus increasing the intragastric pH. On  
709 the other hand, PPIs obtain the same result through the inhibition of acid secretion in gastric  
710 parietal cells, which blocks H<sup>+</sup>/K<sup>+</sup> -ATPases. Both processes are therefore shifting the gastric  
711 environment towards an alkaline condition. Thus, MCE provides an ordering of the groups  
712 along the second dimension that is related to pH increment (from **H-** to **P&H+**). **Furthermore,**  
713 **we contrast MCE performance on this challenging dataset versus two baseline algorithms**  
714 **for nonlinear dimension reduction: t-SNE and Isomap. As we stressed in the introduction**  
715 **these algorithms require optimal tuning of parameters (two for t-SNE and one for**  
716 **Isomap). We believe that advanced nonlinear data analysis needs adaptiveness and**  
717 **automatization, whereas methods such as t-SNE and Isomap, although in principle are**  
718 **unsupervised, in practice are applied in a supervised manner and the hypothesized class**  
719 **labels are used to learn their best parameter tuning. Unlikely, in small size datasets,**  
720 **parameter tuning is a relevant issue that may cause overfitting, especially with more than**  
721 **one parameter such as in the case of t-SNE and, to the best of our knowledge, there is not**  
722 **yet any commonly accepted solution for this. Here, with the mere intention to provide a**  
723 **proof of concept that allows to compare MCE with other nonlinear dimension reduction**  
724 **methods, we apply a supervised procedure in which the labels are used to supervisedly**  
725 **tune the internal parameters of these methods and we select the solution which offers the**  
726 **best performance. The results are shown in Suppl. Figure S3. t-SNE (PSI-ROC: 0.90, PSI-**

727 **PR: 0.94) and Isomap (PSI-ROC: 0.87, PSI-PR: 0.94) performances are lower than MCE**  
728 **performances, displaying difficulty to resolve the difference between treated and**  
729 **untreated samples, mostly for the cases of treated patients (blue points) and untreated**  
730 **patients without *H. Pylori* infection (red points). This indicates that in principle adaptive**  
731 **parameter-free algorithms such as MCE may also outperform more complex algorithms**  
732 **under difficult scenarios such as for this particular case.**

733 Similarly to the Paroni Sterbini *et al.* microbial dataset, the Tripartite-Swiss-roll dataset (that is  
734 a synthetic dataset containing nonlinear structures obtained by tri-partitioning a discrete Swiss-  
735 Roll manifold<sup>26</sup> in a three-dimensional space, for more details see the **Suppl. Methods section:**  
736 **Artificial Datasets**), presents a hierarchical-organized nonlinearity (Fig. S1A). And also in this  
737 case, similarly to the result of the Paroni Sterbini *et al.* analysis, non-centred MCE is able to  
738 perform a hierarchical embedding that orders the hidden subgroups of the dataset along the  
739 second dimension of embedding (Fig. S1E). On the contrary - as already commented in the  
740 method section - PCA, MDSbc and NMDS (Fig. S1B-D) were unable to resolve the nonlinearity  
741 of the Tripartite-Swiss-Roll: its three partitions are either superimposed (Fig. S1B, D) or twisted  
742 in a horseshoe shape (Fig. S1C). Indeed, the Tripartite-Swiss-Roll is purposely created to  
743 reproduce a manifold that is nonlinear and discontinuous (broken in three parts) such as the  
744 results of MCE analysis of Paroni Sterbini *et al.* seems to be. **Furthermore, to compare the**  
745 **different approaches in a more “realistic” scenario, a synthetic microbial-like dataset**  
746 **(which resamples the nonlinearities encountered in the Paroni Sterbini et al. data) is**  
747 **generated and analysed (for more details see the Suppl. Methods section: Artificial**  
748 **Datasets). MCE overcomes once again the other dimensionality reduction techniques and**  
749 **is very close to guarantee a separability equivalent to the one obtained in the high**  
750 **dimensional space (HD) (Suppl. Table S4). These results are similar to the ones obtained**  
751 **in the real datasets. As expected, MDS with weighted Unifrac distance is highly affected**  
752 **by the fact that phylogenetic information between synthetic features is not available and**

753 **it is directly extracted from the OTU table. Interestingly, and opposite to what already**  
754 **shown in the real dataset, MDS with Theta-YC distance obtains great performances close**  
755 **to MCE.**

756 For the Paroni Sterbini dataset, we also performed a supervised linear approach for dimension  
757 reduction, LDA (Suppl. Fig. S4), yet the cross-validation test showed that this constrained  
758 technique could re-assign samples to their groups with 54% of error (ldaCVerErr in Suppl. Table  
759 S5), confirming its statistical invalidity for the small size dataset problem.

760 Moreover, the clustering algorithms MCL and MC-MCL, that is the minimum curvilinear  
761 version of MCL were applied to the Paroni Sterbini *et al.* dataset and the best results (highest  
762 accuracies) are shown in Table 1 (bottom panel) (for more details see the methods' sections  
763 '*From Markov Clustering (MCL) to Minimum Curvilinear Markov Clustering (MC-MCL)*' and  
764 '*Procedure to evaluate the performance of clustering algorithms*'). MC-MCL performs better  
765 than the MCL (both for three and four clusters), even if their accuracies are not remarkably  
766 high, confirming that difficulties in pattern-recognition arise also from the presence of three  
767 clusters in the high-dimensional space. In addition, the hypothesis of three clusters seems more  
768 congruous than four clusters, because both MC-MCL and MCL decrease their accuracies in  
769 detecting four clusters.

770 While MC-MCL represents the minimum curvilinear version of MCL, MCE is the minimum  
771 curvilinear version of PCA, particularly valuable for small sample size datasets. The principle  
772 behind them is MC<sup>23</sup>, that suggests that curvilinear (nonlinear) distances between samples may  
773 be estimated as pairwise distances over their Minimum Spanning Tree (MST) (constructed  
774 according to a selected distance). In fact, as explained in <sup>76</sup>, to approximate nonlinear  
775 (curvilinear) distances between the points of the manifold it is not necessary to reconstruct the  
776 nearest-neighbour graph. Indeed, a greedy routing process (that exploits a norm, for instance  
777 Euclidean) between the points in the multidimensional space is enough to efficiently navigate  
778 the hidden network that approximates the manifold in the multidimensional space. And a

779 preferable greedy routing strategy, at the basis of MC-kernel, is the minimum spanning tree  
780 (MST).

781 Overall, we can conclude that both MCE in dimensionality reduction and MC-MCL in  
782 clustering perform better than the respective non-MC-based versions, and this result confirms  
783 the presence of nonlinear complexity in this dataset, generated by a three-body interaction  
784 (presence of three clusters). In addition, when considering correlation-based distances, they do  
785 not react to the presence of compositionality, since pairwise correlations are computed between  
786 samples. Compositionality instead is a problem that arises when the correlations is computed  
787 between OTUs (features) from metagenomics abundance data (which are normalized by dividing  
788 each OTU count to the total sum of counts in the sample <sup>77,78</sup>), which yields unreliable results  
789 due to dependency of microbial relative abundances.

790 Moreover, because of the discovered major nonlinear complexity in the Paroni Sterbini gastric  
791 biopsy dataset, we wanted to verify whether it was generated by multi-grouping (three-body  
792 interaction problem associated to the presence of three hidden clusters). To do so, we applied  
793 PCA to three subsampled versions of the dataset (with the best normalization originally found  
794 for the complete dataset), each corresponding to the combination of two groups (Fig. 4A-C),  
795 and PCA could find significant separation (**PSI-ROC and PSI-PR > 0.80**). To further confirm  
796 that the presence of multiple sample groups generates the data complexity, we did the same for  
797 the Tripartite Swiss-Roll (Fig. S5A-C), where we recovered the discrimination, even though  
798 two comparisons overlap to some extent (Fig. S5A and C). **Additionally, it might be argued**  
799 **that the presence of *H. pylori* only drives the difference of the microbial community,**  
800 **instead of PPI treatment. However, if this were the case, then the segregation between H+**  
801 **and H- samples would be evident as well inside the PPI treated group. However, the P-**  
802 **values are not anymore significant for this case (P-value PCA: 0.46 & P-value MCE: 1)**  
803 **and no evident segregation arises neither by eyes, as supported by the Suppl. Fig S6.**



804 In conclusion, the results confirm that linear techniques, even if supervised like LDA, are not  
805 able to resolve the differences in the data due to the presence of nonlinear complexity generated  
806 by the three-body interaction (**H-**, **H+** and **P**). Once the complexity is reduced to a two-body  
807 interaction, the problem tends to vanish and PCA can detect significant differences between the  
808 groups, as shown by the PCA pairwise comparisons.

809 Hence, the results of unsupervised analysis on Paroni Sterbini *et al.* dataset show that PPI  
810 treatment causes a major change in gastric mucosal communities of dyspeptic patients,  
811 regardless of the initial pathological infection due to *H. pylori*.

812

### 813 **Comparison of unsupervised analysis in three gastro-esophageal datasets**

814 We compared the performance of unsupervised analysis (dimensional reduction and clustering)  
815 in the Paroni Sterbini dataset <sup>22</sup> (gastric biopsies) and two additional datasets by Amir and  
816 colleagues <sup>21</sup>, that investigated the PPI influence on the esophageal microbiota (Amir3) and  
817 gastric fluid (Amir4).

818 Table 1, top panel, shows the best results in performance of unsupervised dimension reduction  
819 (PCA, MDSwUF, MDSbc, NMDS, MCE, for details see '*Methods - PCA, MDS (or PCoA) and*  
820 *LDA*' and '*Methods - Minimum Curvilinear Embedding (MCE)*') according to PSI based on  
821 AUC and AUPR, on the three different datasets (for more details on the PSI see '*Methods -*  
822 *Procedure to evaluate the performance of the dimension reduction algorithms*'). **Just the space**  
823 **of the first two dimensions of embedding were here used since they are the ones related**  
824 **with the treatment/infection –related structures (Suppl. Table S1).** The mean performance  
825 across all datasets is shown in the last column of the table for each method. The corresponding  
826 ranked performance for each method, based **PSI-ROC and PSI-PR**, is presented instead in  
827 Table 2. For the Paroni Sterbini dataset, we show the results for three different labels (untreated  
828 **H-**, untreated **H+** and PPI-treated). For the Amir datasets, the p-values were computed for two  
829 groups, identified by the presence or absence of PPI treatment. The PSI was also applied to the



830 data in the original high-dimensional (HD) space, as a reference to see how good the  
831 unsupervised dimension reduction approaches are in preserving the group separability in the  
832 HD. Moreover, the **PSI-ROC** and **PSI-PR** best results with **trustworthiness and** standard error  
833 on the **real** datasets, when applying leave-one-out-cross-validation (LOOCV), are shown in  
834 Suppl. Table S6.

835 For the Paroni Sterbini dataset, the PSI evaluation in the first two dimensions of embedding  
836 identifies MCE as the best dimension reduction technique that is able to preserve the group  
837 separability in the HD space. Surprisingly, MCE (presented in Fig. 3A, **PSI-ROC=0.91, PSI-  
838 PR=0.96**) outdoes HD in sample separation in three groups (for HD, **PSI-ROC=0.88, PSI-  
839 PR=0.94**). Similarly, in Amir4, MCE (**PSI-ROC=0.91, PSI-PR=0.920**) succeeds in preserving  
840 the separability of the original HD space (in HD, **PSI-ROC=0.98, PSI-PR=0.99**), better than  
841 the other dimension reduction methods. Finally, dimension reduction analysis on the Amir3  
842 dataset shows that esophageal biopsies were significantly different before and after PPI  
843 treatment, as shown by MDSwUF results (**PSI-ROC=1=PSI-PR**), that surpass the **PSI-ROC  
844 and PSI-PR** values in HD space (**PSI-ROC=0.95, PSI-PR=0.96**). Markedly, MDSwUF  
845 reaches a value of AUPR and AUC of 1, meaning perfect classification of the samples.

846 Overall, when averaging across all datasets, the two metrics based on **PSI-ROC** and **PSI-PR**  
847 pointed out that MDSwUF (**PSI-ROC=0.90, PSI-PR=0.93**) gave the best results of separability  
848 compared to HD (**PSI-ROC=0.94, PSI-PR=0.96**), followed by MCE with closer results (**PSI-  
849 ROC=0.90, PSI-PR=0.92**). Then PCA is the third best result (**PSI-ROC=0.87, PSI-PR=0.90**),  
850 followed by MDStcy, NMDS and MDSbc. However, to conclude what is the best method, we  
851 considered an evaluation based on ranking (Table 2). It is important to note that MCE was the  
852 dimension reduction approach that ranked first in performance across all the datasets, followed  
853 by MDSwUF (Table 2). Hence, the results of sample separability suggest the presence of hidden  
854 patterns that emerge by applying nonlinear dimension reduction techniques like MCE and  
855 MDSwUF.

856 Then clustering algorithms, MCL and its Minimum Curvilinear version (for more information  
857 see *'Methods - From Markov Clustering (MCL) to Minimum Curvilinear Markov Clustering*  
858 *(MC-MCL)'*), were used to confirm the well-possedeness of the hidden patterns that were  
859 recognized by nonlinear dimension reduction. The best results as highest accuracies in each  
860 dataset and the mean performance across all the datasets are exhibited in Table 1, bottom panel.  
861 As already discussed in the previous section, the minimum curvilinear version of MCL (MC-  
862 MCL, acc=0.71) outperforms the MCL clustering algorithm (acc=0.67) in the Paroni Sterbini  
863 dataset, confirming the presence of underlying non-linear complexity in the data. However, the  
864 accuracy doesn't reach high values, because of the difficulty in pattern recognition generated  
865 by the three-body problem in the HD space. Curiously, the accuracies for four clusters (**H-**, **H+**,  
866 **P&H-**, **P&H+**) drop to 0.58 for MC-MCL and to 0.63 for MCL, supporting the hypothesis that  
867 three clusters are more congruous than four clusters. Notably in Amir3, MC-MCL attains high  
868 clustering accuracy (acc=0.81), compared to MCL (acc=0.69). This is the dataset for which,  
869 surprisingly, Amir and collaborators did not find significant changes in the esophageal tissue  
870 microbiota following PPI-treatment, using classical MDS unsupervised multivariate method  
871 with unweighted UniFrac distance <sup>21</sup>. Instead, in the gastric fluid dataset (Amir 4), MC-MCL  
872 and MCL got the same accuracy of 0.75, where a significant separation of samples according  
873 to PPI consumption was already proved in the original article <sup>21</sup>.  
874 However, we have to clarify that normalizations besides scaling (DRS and DCS) and log-  
875 transformation ( $\log(1+x)$ ) could potentially lead to different performance results of  
876 unsupervised analysis. Normalization is crucial to address uneven sampling depth and sparsity  
877 (high proportion of zeros) in microbiome data, like rarefying an OTU table, that is randomly  
878 sampling without replacement from each sample such that all samples have the same number  
879 of total counts (sequencing depth) <sup>79-82</sup> ([http://qiime.org/scripts/single\\_rarefaction.html](http://qiime.org/scripts/single_rarefaction.html)). This  
880 normalization is recommended to moderate the sensitivity of UniFrac distances to sequencing

881 (sampling) depth <sup>52,83</sup>, especially differences in the presence of rare OTUs <sup>50</sup>, nonetheless it is  
882 also considered statistically improper due to the omission of data <sup>83</sup>.

883 Another normalization was introduced in 2010 by Anders and colleagues for general sequence  
884 count data (function *varianceStabilizingTransformation* implemented in the Bioconductor  
885 DESeq2 package), that uses a Variance-Stabilization Transformation (VST) by modelling  
886 microbiome count data with Negative Binomial (NB) distribution <sup>80,83</sup>.

887 We also provide the results with these two different normalizations, and we further confirm that  
888 the data are segregated in the HD space when pre-processed according to them, as shown in the  
889 PSI-ROC and PSI-PR tables in Additional file (for negative binomial, Suppl. Tables S7-9; for  
890 rarefaction, Supplementary Table S12-14). Interestingly, across all the datasets MCE decreases  
891 its performance with these pre-processing techniques, remarkably with rarefied datasets, while  
892 the other linear techniques improve in performance (Suppl. Table S7 for negative binomial;  
893 Suppl. Table S12 for rarefaction), suggesting that these adjustments linearize the datasets.  
894 Indeed, since MCE is a hierarchical technique, it needs the presence of nonlinearity to perform  
895 well. In a similar way, with these two normalizations the accuracy of MC-MCL drops down  
896 (less remarkably in the rarefaction datasets), while the performance of MCL does not increment  
897 (Suppl. Table S10 for negative binomial; Supplementary Table S15 for rarefaction). It is true  
898 that some pre-processing steps such as negative binomial tend to linearize the data but, in this  
899 manner, they can also remove important nonlinear discriminative information, as we show with  
900 the results of unsupervised analysis. Therefore, some pre-processing approaches can also cancel  
901 important nonlinear discriminant information present in the analysed data.

902

903 **Network analysis clarifies the effect of PPI-treatment on the gastric**  
904 **microbiota**

905 Five major phyla have been detected in the normal gastric microbiota: *Firmicutes*,  
906 *Bacteroidetes* and *Actinobacteria* dominate the gastric fluid samples, while *Fusobacteria* and  
907 *Proteobacteria* are the most abundant phyla in gastric mucosal samples <sup>1</sup>.

908 However, the composition and abundance of gastric microbiota may be affected by many  
909 factors, such as dietary habits, *H. pylori* infection, diseases and drugs, including PPIs <sup>1</sup>.

910 Yet, although recent studies have highlighted the potential of these antacid drugs to affect the  
911 gastric microbiota, more knowledge needs to be gained about the association between PPI usage  
912 and the non-*H. pylori* bacteria in the stomach.

913 Since we wanted to investigate the effect of PPI intake on gastric microbiota in dyspepsia, we  
914 analysed: Amir4 for gastric fluid microbiota <sup>21</sup> and Paroni Sterbini et al. dataset <sup>22</sup> for gastric  
915 mucosal microflora, in the latter case restricting to PPI-treated *H. pylori*-negative (**P&H-**) and  
916 untreated *H. pylori* negative patients (**H-**). In both studies, the samples from dyspeptic patients  
917 were analysed using the same next-generation sequencing technologies for direct sequencing  
918 of 16S rRNA gene amplicons, 454 Pyrosequencing.

919 For this purpose, we employed PC-corr algorithm, that was discussed in the Methods section  
920 named: '*PC-corr network*'. In brief, PC-corr discloses the discriminative network of features  
921 that are associated to a sample separation along a principal component direction. Hence, we  
922 expect that the PC-corr network of bacteria will offer a view on how the community of  
923 bacteria respond to PPI-treatment perturbation in the gastric niche (environment), in  
924 dyspeptic patients.

925 **Up to this point, in order to assess the emergence of nonlinear patterns in data, the**  
926 **application and performance of linear and non-linear dimensionality reduction**  
927 **algorithms has been compared. Special focus was on Paroni Sterbini dataset, where the**  
928 **presence or absence of *H. pylori* infection in addition to the medical treatment (or not)**

929 **with PPI medicaments created a complex nonlinear scenario difficult to disentangle using**  
930 **linear transformations and even some nonlinear ones. Then, with the didactic help of the**  
931 **Tripartite-Swiss-roll dataset, we clarified that the origin of the Paroni Sterbini**  
932 **nonlinearity stays in the three-body problem. Indeed, considering pairwise comparisons**  
933 **of only two groups, the nonlinearity vanished. Based on these considerations, now we**  
934 **conduct only the two-group comparison of PPI treated/nontreated patients in which**  
935 **presence of *H. Pylori* was negative, since these data are available both in Paroni Sterbini**  
936 **and Amir. Such simplification of the investigation enables the application of the above**  
937 **mentioned PC-corr algorithm, since, for the binary class problem both Paroni Sterbini**  
938 **and Amir4 datasets present a significant segregation measured by MW-pvalue when**  
939 **embedded by the linear algorithm PCA.**

940 In Amir4 (gastric fluid), PCA revealed that gastric fluid samples were separated into two groups  
941 according to PPI treatment along PC2 and their difference is significant ( $p$ -value  $< 0.01$ )  
942 (Suppl. Figure S7). Hence, we built the PC-corr network<sup>64</sup> using the loadings of PC2 at cut-off  
943 0.5 (Suppl. Figure S8).

944 Similarly for the Paroni Sterbini dataset (gastric mucosa), PCA (Suppl. Figure S9) could  
945 (significantly or close to significance) separate PPI-treated *H. pylori*-negative patients from  
946 untreated *H. pylori*-negative patients along PC2 and PC15 ( $p$ -value along PC2 = 0.014,  $p$ -value  
947 along PC15=0.054). Therefore we built the PC-corr network for both PC2 and PC15  
948 discriminating dimension using 0.5 cut-off (Suppl. Figure S10, panel A and B).

949 Subsequently, to investigate how PPI is affecting the microbiota in the gastric environment, we  
950 considered the conserved **PC-corr network as an indication of bacteria behavior robustness.**

951 **It** is obtained as the union of the two PC-corr networks (obtained for PC2 and PC15) derived  
952 from the Paroni Sterbini gastric mucosa dataset intersected with the PC-corr network derived  
953 from the Amir4 gastric fluid dataset. The resulting conserved network displays the bacteria with  
954 same trend in the two datasets, i.e. either increased or decreased **abundance for patients** with

955 PPI-treatment, respectively in red and black colour, as emphasized by the violet circle at the  
956 centre of Figure 5. Figure 6 is the same as Figure 5 but here the nodes are coloured according  
957 to phylum-level taxonomy. The conserved network which arises at the overlap between the two  
958 PC-corr networks (union of Paroni Sterbini networks intersected with the Amir4 network) is  
959 statistically significant ( $p$ -value=1.00e-04), as a result of the statistical test based on trying to  
960 obtain the same conserved network by random resampling the bacteria in the two networks  
961 (Suppl. Figure S11), implying the difficulty of generating this intersection simply at random  
962 (since this intersection lies to the right of the critical value at the 0.05 level in the distribution  
963 of overlap). This is an important result because it confirms the robustness of the detected  
964 conserved network as a microbiota signature perturbed by PPI treatment. The top and bottom  
965 panels in Figure 5 and 6 show instead the remaining part of Amir4's network (top panel) and  
966 of Paroni Sterbini's network (bottom panel) that are not in the intersection, and therefore might  
967 be more specific for the gastric fluid and mucosa respectively. The PPI-perturbed conserved  
968 network is characterized by a main interconnected module with nine bacteria of four different  
969 phyla (*Bacteroidetes*, *Fusobacteria*, *Proteobacteria*, *Firmicutes*) that are positively associated  
970 (red edges) and by two single bacteria order without interactions (*Streptophyta*, *Clostridiales*),  
971 all being increased following PPI treatment, except *Streptophyta* that is instead decreased with  
972 PPI-treatment (Fig. 5 and 6). Note that a mix between genera, phyla and order of bacteria can  
973 be found in the networks. The reason behind it is the availability of detail information regarding  
974 different bacteria. Some of the spotted bacteria (*Veillonella*, *Clostridiales*, *Campylobacter*)  
975 were already observed in previous studies. The genus *Veillonella* was found increased in  
976 relation to PPI use <sup>16</sup> in the gut microbiome and has been associated with increased  
977 susceptibility to *Clostridium difficile* infection <sup>84</sup>. These Gram-negative anaerobic cocci with  
978 lactate fermenting abilities are abundant in the human microbiome and are normally found in  
979 the intestines and oral mucosa of humans <sup>85</sup>. Interestingly, they favour nitrite accumulation in  
980 the stomach during nitrate reduction, promoting a carcinogenic effect <sup>1</sup>. In addition, the order

981 *Clostridiales*, that is associated to *Clostridium difficile* infection, was also seen significantly  
982 changed in the gastrointestinal tract, however Freedberg *et al.*<sup>4</sup> found it significantly decreased  
983 during PPI use, in contrast to our results. PPIs use also increases the risk of other enteric  
984 infections, apart from *C. difficile* infection, such as campylobacteriosis, as reported in<sup>86,87</sup>.  
985 Moreover, half of the bacteria present in the network normally colonize the human oral  
986 cavity. Indeed, it is the main purpose of PPI treatment to increase the stomach pH, and the  
987 higher pH of treated patients is known to favour the growth of bacteria that usually reside in  
988 the mouth and esophagus and are not adapted to survive the normal gastric acidity<sup>6,20</sup>.  
989 Among genera usually reported as part of the normal **microbiota** of the gastrointestinal tract,  
990 only *Veillonella* is found regularly at other sites, like the mouth<sup>88</sup>. *Leptotrichia* species mostly  
991 colonize the oral cavity and they were isolated from various human infections, suggesting that  
992 they are emerging human pathogens<sup>89,90</sup>. *Oribacterium* also inhabits the mouth, besides the  
993 upper respiratory tract<sup>91</sup>. *Prevotella* is a genus of Gram-negative bacteria that tend to colonize  
994 the human gut, mouth and vagina, and may cause infections, mostly observed in the oral cavity  
995 (odontogenic infections)<sup>90</sup>. *Porphyromonas* has been found by<sup>92</sup> as part of the salivary  
996 microbiome. Both *Prevotella* and *Porphyromonas* contribute to the formation of abscesses and  
997 soft tissue infections in various part of the body and they can cause infections, including  
998 periodontal and endodontal diseases<sup>93</sup>. *Capnocytophaga* are inhabitants of the oral cavity too,  
999 and these opportunistic pathogens can cause infections (both in immunocompromised and  
1000 immunocompetent hosts), the severity of which depend on the immune status of the host<sup>94,95</sup>.  
1001 As well, *Granulicatella* are Gram-positive cocci normally found in the oral **microbiota** and are  
1002 uncommon causes of infections, nevertheless they can cause infections, including bloodstream  
1003 infection and infective endocarditis<sup>96</sup>. Besides, the genus *Fusobacterium* inhabits the mucosal  
1004 membranes of humans and all its species are parasites of humans<sup>97</sup>, and some species are found  
1005 in the oral cavity. The remaining bacteria (*Campylobacter*, *Bulleidia*) do not belong to the oral  
1006 microbiota<sup>93</sup>. The genus *Campylobacter* was increased in relation to PPI use and the increased



1007 abundance of these Gram-negative bacteria has the potential to cause diseases and infections in  
1008 humans (most commonly diarrhoea). Due to the induced increase of pH, PPI is hypothesised to  
1009 facilitate gastrointestinal infections and a study by Brophy *et al.*<sup>98</sup> reported an increased risk of  
1010 *Campylobacter* infection following PPI therapy. Moreover Campylobacteriosis, mostly caused  
1011 by eating undercooked foods derived from poultry or other warm-blooded animals or contact  
1012 with contaminated water or ice<sup>99</sup>, has been shown by the Dutch National Institute for Public  
1013 Health and the Environment to noticeably increase in incidence when PPI use grows<sup>86</sup>.

1014 Altogether, PC-corr approach was applied on gastric fluid and gastric mucosal datasets (in the  
1015 latter case, excluding the samples positive to *H. pylori* infection) to investigate how PPI is  
1016 affecting the gastric microbiota (both gastric fluid and gastric mucosal microbiota), because of  
1017 PC-corr's ability to pinpoint the combination of bacteria that play a major role in the  
1018 discrimination of the samples, in this case according to PPI intake. The PC-corr conserved  
1019 network identified eleven genera and order of bacteria, which belong to the phyla  
1020 (*Bacteroidetes*, *Fusobacteria*, *Proteobacteria*, *Firmicutes*) commonly found in the stomach  
1021 which, with exception of *Streptophyta*, demonstrated increased abundance following PPI  
1022 treatment. Mostly all the found bacteria were not reported in previous studies, except  
1023 *Veillonella*, *Clostridiales* and *Campylobacter*, but they were found as inhabitants of the oral  
1024 cavity and/or possible cause of infections and diseases in humans. Hence, and in concordance  
1025 to previous studies<sup>6,20</sup>, these results point out that PPI treatment, by increasing the intragastric  
1026 pH, favours the growth of bacteria that usually reside in the mouth and survive through the  
1027 harsh acidic conditions of the stomach. Furthermore, the results suggest that PPI-associated  
1028 increases of some bacterial populations may lead to infections and diseases or increase  
1029 susceptibility for other bacterial infections (like *Veillonella*) or promote a carcinogenic effect  
1030 (like *Veillonella*). Previous studies have highlighted that PPI intake is associated with decreased  
1031 bacterial richness<sup>16,18,100,101</sup>, increased risk of enteric and other infections (e.g. caused by  
1032 *Salmonella*, *Clostridium difficile*, *Shigella*, *Listeria*)<sup>17,102</sup>, increase in the abundance of oral and



1033 upper GI tract commensals and potential pathogenic bacteria (e.g. *Enterococcus*,  
1034 *Streptococcus*, *Staphylococcus*, and *Escherichia coli* )<sup>16,17</sup> in the gut microbiota. Nevertheless,  
1035 our analysis by means of PC-corr does not spot single bacteria perturbed in the gastric  
1036 environment by PPI treatment, but a community of bacteria is altered in abundance by PPIs and  
1037 their inter-specific bacterial interactions in the gastric niche.

1038 Therefore our study will ground the basis for further investigations that could better clarify the  
1039 effect of PPI-treatment on the human gastric microbiota and additionally verify the identified  
1040 altered bacteria, as PPIs may have possible side-effects, including increased risks of different  
1041 infections and diseases.

1042

## 1043 **Network analysis clarifies the effect of *H. pylori* infection on gastric mucosal** 1044 **microbiota**

1045 The stomach was long thought sparsely colonized by bacteria due to the gastric microbicidal  
1046 acidic barrier (pH<4.0)<sup>103</sup>. This view dramatically changed with the discovery of the Gram-  
1047 negative bacterium *H. pylori* in the 1980's by Warren and Marshall<sup>104</sup>, that is a carcinogenic  
1048 bacterial pathogen infecting the stomach of more than one-half of the world's  
1049 human population. This human pathogen is able to survive in the highly acidic environment  
1050 within the stomach by producing cytoplasmic urease that, by catalysing the hydrolysis of urea  
1051 into CO<sub>2</sub> and NH<sub>4</sub>, produces a neutralizing ammonia cloud around it<sup>19,105,106</sup>. However, most  
1052 *H. pylori* avoid the acidic environment of the gastric lumen by swimming towards the mucosal  
1053 cell surface (using their polar flagella and chemotaxis mechanisms) and may adhere and invade  
1054 the gastric mucosal epithelial cells<sup>107,108</sup>. Hence, it doesn't represent a dominant species in  
1055 gastric fluid microbiota<sup>109</sup>, but was found to generally to reside in the gastric mucosae<sup>5,107,110</sup>.  
1056 Persistent (chronic) infection with this Gram-negative bacterium induces changes in gastric  
1057 physiology and immunology, e.g. reduced gastric acidity and parietal cell mass, perturbed  
1058 nutrient availability, local innate immune responses<sup>111,112</sup>, that most probably induces shift in

1059 gastric microbiota composition <sup>111</sup>. Although *H. pylori* colonization usually persists in the  
1060 human stomach for many decades without adverse effects, the infection of this bacteria is  
1061 associated with increased risk for several diseases, including peptic ulcers, chronic gastritis,  
1062 mucosa-associated lymphoid tissue lymphoma, gastric adenocarcinoma <sup>113,114</sup>, and dyspepsia  
1063 <sup>115,116</sup>. The potential alterations induced by the *H. pylori* can in turn lead to dysbiosis and may  
1064 cause aberrant proinflammatory immune responses <sup>117</sup>, susceptibility to bacterial pathogens and  
1065 increased risk of gastric disease, including cancer <sup>1,118</sup>. However, the effect of *H. pylori*  
1066 infection on overall composition of gastric microbiota at genus level and the bacterial interplay  
1067 in presence of this widespread human infection remain unclear.

1068 **Similar to the PPI treatment network analysis in the previous section, in order to**  
1069 **investigate the influence of *H. pylori* infection on the gastric mucosal microbiota by means**  
1070 **of PC-corr**, we analysed: 1) Paroni Sterbini *et al.* <sup>22</sup> considering only PPI-untreated dyspeptic  
1071 patients, either infected (**H+**) or not by *H. pylori* (**H-**); 2) Parsons *et al.* <sup>32</sup> restricting to PPI-  
1072 untreated patients from: i) normal stomach group with no evidence of *H. pylori* infection; ii) *H.*  
1073 *pylori* gastritis group with evidence of *H. pylori* infection. Even though the same technology is  
1074 important for a comparative study, unfortunately in the literature there was no such data  
1075 available like Paroni Sterbini's one, that is 16S rRNA gene pyrosequencing data (derived from  
1076 gastric mucosal microflora in dyspeptic untreated patients either positive or negative for *H.*  
1077 *pylori*). Despite this, the two studied datasets, obtained with two different next-generation  
1078 sequencing technologies for direct sequencing of 16S rRNA gene amplicons (454  
1079 Pyrosequencing for Paroni Sterbini *et al.* and Illumina MiSeq for Parsons *et al.*) <sup>119</sup>, both contain  
1080 community profiling of gastric mucosa-associated microbiota in PPI-untreated *H. pylori*-  
1081 negative and -positive subjects. However, for the sake of clarity, we have to specify a  
1082 difference: while in Paroni Sterbini's dataset the gastric mucosal biopsy specimens were  
1083 collected from patients with dyspepsia, this is not the case for Parsons's data.

1084 To enhance the understanding of the *H. pylori*-triggered microbial perturbation in this  
1085 ecological niche, we employed again PC-corr algorithm, that is able to associate to any PCA  
1086 analysis of an omic dataset, where a sample separation emerges, a network of discriminative  
1087 features (for details see '*Methods-PC-corr network*'). The analysis of the 16S rRNA sequencing  
1088 data was restricted only the overlapping OTUs, excluding *Helicobacter* because our goal is to  
1089 investigate its impact on the rest of the microbial network.

1090 In Paroni Sterbini's dataset, since PCA could significantly separate gastric mucosal biopsy  
1091 samples of PPI-untreated patients according to *H. pylori*-positivity (p-value=0.01) along PC2  
1092 (Suppl. Fig. S12), the PC-corr network was constructed from PC2 loadings at 0.5 cut-off (Suppl.  
1093 Fig. S13). Similarly, for Parsons' dataset, since PCA (Supplementary Figure S14) could  
1094 significantly separate patients from the normal stomach group with no evidence of *H. pylori*  
1095 infection and PPI-untreated (Control) from *H. pylori* gastritis group positive to *H. pylori*  
1096 infection and not using PPIs (HPGas) along PC1 (p-value along PC1 <0.01.), the PC-corr  
1097 network was constructed from this discriminating dimension at 0.5 cut-off (Suppl. Fig. S15).

1098 The obtained microbial differential networks (Figure 7, coloured according to phylum level)  
1099 pinpointed, from the system point of view, the bacteria affected by *H. pylori* infection in the  
1100 gastric mucosa, that are precisely bacteria whose abundance is decreased in *H. pylori*-positive  
1101 patients. A presumable explanation of this trend is already pointed out in literature, where the  
1102 presence of *H. pylori* leads to a reduced gastric microbial diversity<sup>120-122</sup>. Nevertheless, in some  
1103 cases the diversity increases again, because of diverse factors that allow survival and  
1104 colonization of bacteria in the stomach<sup>1,123</sup>. Then, the preserved network of gastric mucosa  
1105 microbiota was constructed by intersecting the two PC-corr networks obtained from Paroni  
1106 Sterbini's and Parsons's dataset. Figure 8, middle panel, shows the conserved network (violet  
1107 circle), which presents the common bacteria coloured according to phylum level and their  
1108 associations. The spotted bacteria display decreased abundance with *H. pylori* infection (i.e.  
1109 increased in *H. pylori*-negative subjects) in both the two 16S rRNA gene sequencing data. By

1110 performing a statistical test based on random resampling of the bacteria in the two networks,  
1111 we verified that the shown bacterial conserved network is statistically significant and difficult  
1112 to be generated at random (p-value=1.00e-04), because getting this intersection at random is  
1113 very rare (Supplementary Figure S16). The top and bottom panels in Figure 8 show instead the  
1114 remaining part of Paroni Sterbini's network (top panel) and of Parsons's network (bottom  
1115 panel) that are not in the intersection. At the genus level, a study by Klymiuk *et al.*<sup>124</sup> identified  
1116 *Actinomyces*, *Granulicatella*, *Veillonella*, *Fusobacterium*, *Neisseria*, *Helicobacter*,  
1117 *Streptococcus*, and *Prevotella* as significantly different between the *H. pylori*-positive and *H.*  
1118 *pylori*-negative gastric samples. These bacteria do not emerge in the conserved network, while  
1119 they all (except *Neisseria*) appear altered (decreased) during *H. pylori* infection in the study by  
1120 Parsons and colleagues (present in the bottom panel of Figure 7).

1121 Our analysis pinpoints a conserved network from two independent 16S rRNA gene sequencing  
1122 data, that reveals microbial communities altered by *H. pylori* infection and their interactions in  
1123 the gastric mucosa. It revealed a main core of six associated bacteria (with positive association,  
1124 red edges) and two single nodes without any interaction with the main module, from three  
1125 different phyla (*Proteobacteria*, *Firmicutes*, *Actinobacteria*) all resulting decreased in *H.*  
1126 *pylori*-infected subjects (that is increased in non-infected subjects). The decreased abundance  
1127 of the phyla *Firmicutes* and *Actinobacteria* in *H. pylori*-positive patients with respect to *H.*  
1128 *pylori*-negative subjects was already shown in a previous study by Maldonado-Contreras *et al.*  
1129 <sup>125</sup>. In addition, other studies have demonstrated an increased colonization of *Proteobacteria* in  
1130 *H. pylori*-positive patients<sup>125,126</sup>, while the obtained conserved PC-corr network shows that the  
1131 bacteria from this phylum are instead decreased in those individuals. Among the spotted  
1132 bacteria, *Methylobacterium* is a genus of facultative methylotrophic bacteria that are commonly  
1133 found in diverse natural environments (such as leaf surfaces, soil, dust, and fresh water) and in  
1134 hospital environment due to contaminated tap water. *Methylobacterium* species can cause  
1135 health care-associated infections (mainly catheter infection), especially in

1136 immunocompromised patients <sup>127</sup>. In addition, *Sphingomonas* plays a role in human health, as  
1137 some of the sphingomonads (in particular *Sphingomonas paucimobilis*) are the cause of a range  
1138 of mostly nosocomial, non-life-threatening infections. *Sphingomonas* species are widely spread  
1139 in nature, having been isolated from many sources, from water habitats to clinical settings <sup>128</sup>,  
1140 *Pseudomonas*, due to its great metabolic versatility, can also colonize different types of niches  
1141 <sup>129</sup>, including soil and water, in addition to plant and animal associations, and includes  
1142 pathogenic species in humans <sup>130</sup>. *Acinetobacter* species are instead common, free-living  
1143 saprophytes found in soil, water, sewage and foods and are ubiquitous organisms in hospitals.  
1144 They have been increasingly identified as a key source of infection in debilitated patients in  
1145 hospitals, due to their rapid development of resistance to antimicrobials <sup>131</sup>. In particular, one  
1146 species, *Acinetobacter lwoffii*, can trigger gastritis, apart from *H. pylori* <sup>132</sup>. *Propionibacterium*,  
1147 so named for their unique ability to synthesize propionic acid by using unusual transcarboxylase  
1148 enzymes <sup>133</sup>, are primarily facultative pathogens and commensals of humans, living on the skin,  
1149 while other members are widely employed for synthesizing vitamin B<sub>12</sub>, tetrapyrrole  
1150 compounds, and propionic acid, as well as used as probiotics <sup>134</sup>. *Catonella* is another node in  
1151 the network and this bacterial genus is obligative anaerobic, non-spore-forming and non-motile,  
1152 with one known species (*Catonella morbi*) from the human gingival crevice <sup>135,136</sup>, that has been  
1153 associated with periodontitis <sup>135</sup> and endocarditis <sup>137</sup>. Besides, the bacterial genus  
1154 *Enhydrobacter* so far contains a single species, *Enhydrobacter aerosaccus*, a Gram negative  
1155 non-motile bacterium that is both oxidase and catalase positive and shows gas vacuoles <sup>138,139</sup>.  
1156 *Bulleidia*, a Gram-positive, non-spore-forming, anaerobic and non-motile genus, has one  
1157 known species too (*Bulleidia extracta*)<sup>140</sup>.

1158 In conclusion, by means of the PC-corr approach, we determined the combination of bacteria  
1159 responsible for the difference between *H. pylori*-positive and *H. pylori*-negative gastric mucosa  
1160 of untreated patients and their microbe-microbe interactions. All the bacteria, both in the  
1161 conserved network and not, were decreased in *H. pylori*-infected individuals (i.e. increased in

1162 *H. pylori*-negative group). *H. pylori*, like acid suppressing medications (for the treatment of  
1163 dyspepsia), alters the population structure of the gastric and intestinal microbiota <sup>141</sup> and  
1164 regularly, this bacterium constitutes most of the gastric microbiota <sup>123</sup>, literally depleting  
1165 bacterial biodiversity. Moreover, most of the identified bacteria represent bacteria of potential  
1166 health concern, as agents of diseases and infections.

## 1167 **Bacteria-metabolite multilayer network analysis associates possible** 1168 **metabolic pathways perturbations**

1169 **The relation between bacteria and metabolites is fundamental both to deepen the**  
1170 **understanding of mechanisms associated to diseases dysfunction and drugs action, and to**  
1171 **foster their biomedical interpretation** <sup>142–144</sup>. **For this reason, we made a quantum leap in**  
1172 **our investigation from bacteria to metabolites and we built two bacteria-metabolite**  
1173 **multilayer networks: one (Fig. 8) was derived from the PPI-affected bacteria network in**  
1174 **Fig. 6, the other (Fig. 9) was derived from the *H. pylori*-affected bacteria network in Fig.7.**  
1175 **The methodological procedure to build those multilayer networks is provided in the**  
1176 **Methods (see section: Bacteria-metabolite multilayer network construction).**  
1177 **Remarkably, by applying metabolic pathway enrichment analysis, we found that the**  
1178 **metabolite layer of the PPI-affected (Fig.8B) and *H. pylori*-affected (Fig.9B) networks**  
1179 **contain metabolites significantly involved ( $p < 0.05$  after Benjamini correction) in**  
1180 **important pathways (see full list in Suppl. Tables S18 and S19) associated with obesity** <sup>145</sup>,  
1181 **symptomatic atherosclerosis** <sup>146</sup>, **functional dyspepsia** <sup>147</sup>, **gestational diabetes mellitus** <sup>148</sup>  
1182 **Wilson's disease** <sup>149</sup>, **among others. To simplify the visualization and interpretation (for**  
1183 **the methodology of selection see Method section: Bacteria-metabolite multilayer network**  
1184 **construction) we displayed the three most significant and relevant pathways in both PPI-**  
1185 **affected (Fig.8B) and *H. pylori*-affected (Fig.9B) networks. Interestingly, the bacteria**  
1186 ***Porphyromonas* and *Fusobacterium* are highly contributing for possible perturbations on**  
1187 **the Aminoacyl-tRNA biosynthesis pathway for alterations produced by PPI (Figure 8B),**

1188 whilst *Methylobacterium* does it on the *H. Pylori* infection side (Figure 9B). Besides, N-  
1189 Acetylneuraminic acid (Suppl. Fig. S17) is a sialic acid that has been associated also with  
1190 pathogenic enteric bacteria<sup>150,151</sup> and tumors<sup>152</sup>. Overproduction of nitrites and nitrates  
1191 by the observed anaerobic bacteria have been already observed in diverse parts of the  
1192 gastrointestinal tract in patients suffering from migraine<sup>153</sup>, intestinal dysbiosis and  
1193 colorectal cancer<sup>154,155</sup>, an effect already associated with the use of PPIs such as  
1194 omeoprazole<sup>156</sup>.

## 1195 Discussion

1196 This study indicates the necessity of including nonlinear multidimensional techniques into  
1197 clinical studies based on 16S metagenomic sequencing data, since drawing a study's  
1198 conclusions by solely relying on linear techniques, such as PCA and MDS, can lead to data  
1199 misinterpretation and impair the translational path from research to diagnostic. In the era of  
1200 post-genomics and systems approaches, nonlinear dimension reduction and clustering by MCE  
1201 and MC-MCL can offer new insights into complex clinical 16S metagenomics data, like the  
1202 ones studied in this article or the presence of clinical sub-types, and serve as a valuable tool in  
1203 the run towards precision medicine. Moreover, this study shows how it is possible to  
1204 complement multivariate analysis by means of network analysis employing PC-corr algorithm,  
1205 that accounts for the bacteria responsible for the sample discrimination and their co-occurrence  
1206 relationships. Precisely, from the system point of view the obtained microbial differential  
1207 networks pinpointed marked bacteria-bacteria interactions and modules affected by PPI  
1208 treatment in the gastric environment in dyspepsia and by *H. pylori* infection in the gastric  
1209 mucosa. **Moreover, we elucidated via bacteria-metabolite multilayer networks, possible  
1210 metabolic alterations produced by the perturbed bacteria communities and the respective  
1211 metabolic pathways involved in those changes. The fact that we find significant metabolic  
1212 pathways associated to the discriminative bacteria networks, which are detected by PC-  
1213 corr, is a nontrivial finding that suggests the reliability and impact of the integrated**



1214 machine learning/network biology methodology we propose. However, some limitations  
1215 frequently present in integrative systems biology also affect our study. For instance, when  
1216 we adopt protein interaction networks in drug repositioning<sup>157</sup> or in disease analysis<sup>158</sup>,  
1217 we are aware that further information such as the contextualization of the network to the  
1218 peculiar organ, tissue, cell or cell-compartment would allow more accurate results. The  
1219 same is valid for our study, where we have to adopt a generic bacteria-metabolite gut  
1220 network, because it is the most updated resource currently available in the field. This  
1221 means that when – hopefully in future - more specialized bacteria-metabolite networks  
1222 will be available for the gastric mucosa/fluid and even in specific areas of the stomach,  
1223 then our analysis - such as many other omic analysis in integrative network biology - will  
1224 benefit of this quantum leap in the data quality and contextualization. Hence, we suggest  
1225 that our findings can be an important starting point to design new therapies that consider not  
1226 only *H. pylori* infection but also the directly associated microbial alterations as well as the  
1227 indirect alterations due to the drugs used for *H. pylori* eradication such as PPI.

1228

## 1229 **List of abbreviations**

1230 **AUC: Area Under the ROC-Curve**

1231 **AUPR: Area Under the Precision Recall Curve**

1232 LDA: Linear Discriminant Analysis

1233 MC: Minimum Curvilinearity

1234 MCE: Minimum Curvilinear Embedding

1235 MCL: Markov Clustering

1236 MC-MCL: Minimum Curvilinear Markov Clustering

1237 MDS: Multidimensional Scaling

1238 MDSbc: Multidimensional Scaling with Bray-Curtis dissimilarity

1239 MDSwUF: Multidimensional Scaling with weighted UniFrac distance



- 1240 MST: minimum spanning tree
- 1241 ncMCE: non-centred Minimum Curvilinear Embedding
- 1242 NMDS: non-metric (Sammon criterion) Multidimensional Scaling
- 1243 PC: Principal Component
- 1244 PCA: Principal Component Analysis
- 1245 PCoA: Principal Coordinate Analysis
- 1246 PPI: Proton Pump Inhibitor
- 1247 PSI: Projection-based separability index
- 1248 **PSI-ROC: Projection-based separability index applied with AUC**
- 1249 **PSI-PR: Projection-based separability index applied with AUPR**
- 1250 SVD: Singular Value Decomposition

## 1251 **Declarations**

### 1252 **Ethics approval and consent to participate**

1253 Not applicable, because the used datasets have been generated by previous biomedical  
1254 studies, for which ethics approvals and consents were formerly collected.

1255

### 1256 **Consent for publication**

1257 Not applicable

1258

### 1259 **Availability of data and materials**

1260 Not applicable.

1261

### 1262 **Code Availability**

1263 Codes for the PSI measure and MC-MCL clustering algorithm can be found in

1264 <https://github.com/biomedical-cybernetics>

1265

1266 **Competing interests**

1267 The authors declare that they have no competing interests.

1268

1269 **Funding**

1270 This work was supported by the Dresden International Graduate School for Biomedicine and  
1271 Bioengineering (DIGS-BB), granted by the Deutsche Forschungsgemeinschaft (DFG) in the  
1272 context of the Excellence Initiative. PS is supported by Estonian Research Council Starting  
1273 Grant PUT1130. CD is funded by the Research Grants – Doctoral Programs in Germany  
1274 (DAAD), Promotion program Nr: 57299294.

1275

1276 **Authors' contributions**

1277 CVC developed Minimum Curvilinearity (MCE), Minimum Curvilinear Markov Clustering  
1278 (MC-MCL) and the Projection-based Separability Index (PSI). CVC conceived all the study  
1279 and the data analysis workflow with feedbacks from MiSc and SWG. SC, CD and AP  
1280 performed the computational analysis of the data and realized the figures under CVC guidance  
1281 with help of AZ for the bacteria-metabolite analysis. SC, CD, AP together with CVC wrote  
1282 the manuscript with valuable suggestions of PS and AZ. FPS, LM, GC, GI, BP, MaSa, GG and  
1283 AG provided data and knowledge about the Paroni Sterbini *et al.* data cohort. BNP, UZI and  
1284 MP provided data and knowledge about the Parsons *et al.* data cohort. All authors discussed the  
1285 results and revised the manuscript.

1286

1287 **Acknowledgements**

1288 Not applicable

1289

1290 **References**

1291 1. Nardone, G. & Compare, D. The human gastric microbiota: Is it time to rethink the

- 1292 pathogenesis of stomach diseases? *United Eur. Gastroenterol. J.* **3**, 255–260 (2015).
- 1293 2. Quigley, E. M. M. Gut microbiome as a clinical tool in gastrointestinal disease  
1294 management: are we there yet? *Nat. Rev. Gastroenterol. Hepatol.* **14**, 315–320 (2017).
- 1295 3. Strand, D. S., Kim, D. & Peura, D. A. 25 years of proton pump inhibitors: A  
1296 comprehensive review. *Gut and Liver* vol. 11 27–37 (2017).
- 1297 4. Freedberg, D. E., Lebwohl, B. & Abrams, J. A. The impact of proton pump inhibitors  
1298 on the human gastrointestinal microbiome. *Clinics in Laboratory Medicine* vol. 34  
1299 771–785 (2014).
- 1300 5. Wu, W. M., Yang, Y. S. & Peng, L. H. Microbiota in the stomach: new insights. *J. Dig.*  
1301 *Dis.* **15**, 54–61 (2014).
- 1302 6. Vesper, B. *et al.* The Effect of Proton Pump Inhibitors on the Human Microbiota. *Curr.*  
1303 *Drug Metab.* **10**, 84–89 (2009).
- 1304 7. Scarpignato, C. *et al.* Effective and safe proton pump inhibitor therapy in acid-related  
1305 diseases ? A position paper addressing benefits and potential harms of acid  
1306 suppression. *BMC Med.* **14**, 179 (2016).
- 1307 8. Yadlapati, R. & Kahrilas, P. J. When is proton pump inhibitor use appropriate? *BMC*  
1308 *Med.* **15**, 36 (2017).
- 1309 9. Harmon, R. C. & Peura, D. A. Evaluation and management of dyspepsia. *Therap. Adv.*  
1310 *Gastroenterol.* **3**, 87–98 (2010).
- 1311 10. Malfertheiner, P. *et al.* Management of *Helicobacter pylori* infection—the Maastricht  
1312 IV/ Florence Consensus Report. *Gut* **61**, 646–664 (2012).
- 1313 11. Rosen, R. *et al.* 16S community profiling identifies proton pump inhibitor related  
1314 differences in gastric, lung, and oropharyngeal microflora. *J. Pediatr.* **166**, 917–923  
1315 (2015).
- 1316 12. Lanas, A. We are using too many PPIs, and we need to stop: A European perspective.  
1317 *American Journal of Gastroenterology* vol. 111 1085–1086 (2016).

- 1318 13. Vakil, N. Prescribing proton pump inhibitors: Is it time to pause and rethink? *Drugs* **72**,  
1319 437–445 (2012).
- 1320 14. Tran-Duy, A., Spaetgens, B., Hoes, A. W., de Wit, N. J. & Stehouwer, C. D. A. Use of  
1321 Proton Pump Inhibitors and Risks of Fundic Gland Polyps and Gastric Cancer:  
1322 Systematic Review and Meta-analysis. *Clin. Gastroenterol. Hepatol.* **14**, 1706-1719.e5  
1323 (2016).
- 1324 15. Malfertheiner, P., Kandulski, A. & Venerito, M. Proton-pump inhibitors:  
1325 Understanding the complications and risks. *Nat. Rev. Gastroenterol. Hepatol.* **14**, 697–  
1326 710 (2017).
- 1327 16. Imhann, F. *et al.* Proton pump inhibitors affect the gut microbiome. *Gut* **65**, 740–748  
1328 (2016).
- 1329 17. Jackson, M. A. *et al.* Proton pump inhibitors alter the composition of the gut  
1330 microbiota. *Gut* **65**, 749–756 (2016).
- 1331 18. Tsuda, A. *et al.* Influence of proton-pump inhibitors on the luminal microbiota in the  
1332 gastrointestinal tract. *Clin. Transl. Gastroenterol.* **6**, e89 (2015).
- 1333 19. Williams, C. & McColl, K. E. L. Review article: proton pump inhibitors and bacterial  
1334 overgrowth. *Aliment. Pharmacol. Ther.* **23**, 3–10 (2006).
- 1335 20. Sanduleanu, S., Jonkers, D., De Bruine, A., Hameeteman, W. & Stockbrügger, R. W.  
1336 Non-Helicobacter pylori bacterial flora during acid-suppressive therapy: Differential  
1337 findings in gastric juice and gastric mucosa. *Aliment. Pharmacol. Ther.* **15**, 379–388  
1338 (2001).
- 1339 21. Amir, I., Konikoff, F. M., Oppenheim, M., Gophna, U. & Half, E. E. Gastric  
1340 microbiota is altered in oesophagitis and Barrett’s oesophagus and further modified by  
1341 proton pump inhibitors. *Environ. Microbiol.* **16**, 2905–2914 (2014).
- 1342 22. Paroni Sterbini, F. *et al.* Effects of Proton Pump Inhibitors on the Gastric Mucosa-  
1343 Associated Microbiota in Dyspeptic Patients. *Appl. Environ. Microbiol.* **82**, 6633–6644

- 1344 (2016).
- 1345 23. Cannistraci, C. V., Ravasi, T., Montevocchi, F. M., Ideker, T. & Alessio, M. Nonlinear  
1346 dimension reduction and clustering by Minimum Curvilinearity unfold neuropathic  
1347 pain and tissue embryological classes. in *Bioinformatics* vol. 27 i531–i539 (2011).
- 1348 24. Kinross, J. M., Darzi, A. W. & Nicholson, J. K. Gut microbiome-host interactions in  
1349 health and disease. *Genome Med.* **3**, 14 (2011).
- 1350 25. Legendre, P. & Legendre, L. F. J. *Numerical ecology*. vol. 24 (Elsevier, 2012).
- 1351 26. Tenenbaum, J. B., de Silva, V. & Langford, J. C. A global geometric framework for  
1352 nonlinear dimensionality reduction. *Science* **290**, 2319–23 (2000).
- 1353 27. Bunte, K., Haase, S., Biehl, M. & Villmann, T. Stochastic neighbor embedding (SNE)  
1354 for dimension reduction and visualization using arbitrary divergences. *Neurocomputing*  
1355 **90**, 23–45 (2012).
- 1356 28. Maaten, L. van der & Hinton, G. Visualizing Data using t-SNE. *J. Mach. Learn. Res.* **9**,  
1357 2579–2605 (2008).
- 1358 29. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community  
1359 sequencing data. *Nat. Methods* **7**, 335–6 (2010).
- 1360 30. Wang, Q., Garrity, G. M., Tiedje, J. M. & Cole, J. R. Naïve Bayesian classifier for  
1361 rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl. Environ.*  
1362 *Microbiol.* **73**, 5261–5267 (2007).
- 1363 31. Caporaso, J. G. *et al.* PyNAST: A flexible tool for aligning sequences to a template  
1364 alignment. *Bioinformatics* **26**, 266–267 (2010).
- 1365 32. Parsons, B. N. *et al.* Comparison of the human gastric microbiota in hypochlorhydric  
1366 states arising as a result of. *PLOS Pathog.* **13**, 1–19 (2017).
- 1367 33. Cannistraci, C. V., Alanis-Lobato, G. & Ravasi, T. Minimum curvilinearity to enhance  
1368 topological prediction of protein interactions by network embedding. *Bioinformatics*  
1369 **29**, 199–209 (2013).

- 1370 34. Smialowski, P., Frishman, D. & Kramer, S. Pitfalls of supervised feature selection.  
1371 *Bioinformatics* **26**, 440–443 (2009).
- 1372 35. Ringnér. What is principal component analysis? *Nat. Biotechnol.* **26**, 303–304 (2008).
- 1373 36. Jolliffe, I. T. Principal Component Analysis. *Springer Ser. Stat.* **98**, 487 (2002).
- 1374 37. Dinsdale, E. A. *et al.* Multivariate analysis of functional metagenomes. *Front. Genet.* **4**,  
1375 41 (2013).
- 1376 38. Ramette, A. Multivariate analyses in microbial ecology. *FEMS Microbiol. Ecol.* **62**,  
1377 142–160 (2007).
- 1378 39. Moitinho-Silva, L. *et al.* Specificity and transcriptional activity of microbiota  
1379 associated with low and high microbial abundance sponges from the Red Sea. *Mol.*  
1380 *Ecol.* **23**, 1348–1363 (2014).
- 1381 40. Bayer, K. *et al.* GeoChip-based insights into the microbial functional gene repertoire of  
1382 marine sponges (high microbial abundance, low microbial abundance) and seawater.  
1383 *FEMS Microbiol. Ecol.* **90**, 832–843 (2014).
- 1384 41. Alanis-Lobato, G., Cannistraci, C. V., Eriksson, A., Manica, A. & Ravasi, T.  
1385 Highlighting nonlinear patterns in population genetics datasets. *Sci. Rep.* **5**, 8140  
1386 (2015).
- 1387 42. Legendre, P. & De Cáceres, M. Beta diversity as the variance of community data:  
1388 Dissimilarity coefficients and partitioning. *Ecol. Lett.* **16**, 951–963 (2013).
- 1389 43. Paliy, O. & Shankar, V. Application of multivariate statistical techniques in microbial  
1390 ecology. *Mol. Ecol.* **25**, 1032–1057 (2016).
- 1391 44. Zand, M. S., Wang, J. & Hilchey, S. Graphical Representation of Proximity Measures  
1392 for Multidimensional Data: Classical and Metric Multidimensional Scaling. *Math. J.*  
1393 **17**, (2015).
- 1394 45. Cox, M. A. A. & Cox, T. F. Multidimensional Scaling. *Handb. Data Vis.* (2008)  
1395 doi:10.1007/978-3-540-33037-0\_14.

- 1396 46. Sammon, J. W. A Nonlinear Mapping for Data Structure Analysis. *IEEE Trans.*  
1397 *Comput.* **C18**, 401–409 (1969).
- 1398 47. Beals, E. W. Bray-curtis ordination: An effective strategy for analysis of multivariate  
1399 ecological data. in *Advances in Ecological Research* vol. 14 1–55 (1984).
- 1400 48. Bray, J. R. & Curtis, J. T. An Ordination of the Upland Forest Communities of  
1401 Southern Wisconsin. *Ecol. Monogr.* **27**, 325–349 (1957).
- 1402 49. Whittaker, R. H. Vegetation of the Siskiyou Mountains, Oregon and California. *Ecol.*  
1403 *Monogr.* **30**, 279–338 (1960).
- 1404 50. Lozupone, C., Lladser, M. E., Knights, D., Stombaugh, J. & Knight, R. UniFrac: An  
1405 effective distance metric for microbial community comparison. *ISME J.* **5**, 169–172  
1406 (2011).
- 1407 51. Lozupone, C. A., Hamady, M., Kelley, S. T. & Knight, R. Quantitative and qualitative  
1408 beta diversity measures lead to different insights into factors that structure microbial  
1409 communities. *Appl. Environ. Microbiol.* **73**, 1576–85 (2007).
- 1410 52. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing  
1411 microbial communities. *Appl. Environ. Microbiol.* **71**, 8228–35 (2005).
- 1412 53. Chen, J. *et al.* Associating microbiome composition with environmental covariates  
1413 using generalized UniFrac distances. *Bioinformatics* **28**, 2106–13 (2012).
- 1414 54. Podani, J. & Miklós, I. Resemblance Coefficients and the Horseshoe Effect in Principal  
1415 Coordinates Analysis. *Ecology* **83**, 3331–3343 (2002).
- 1416 55. Papadopoulos, F., Psomas, C. & Krioukov, D. Network mapping by replaying  
1417 hyperbolic growth. *IEEE/ACM Trans. Netw.* **23**, 198–211 (2015).
- 1418 56. Muscoloni, A., Thomas, J. M., Ciucci, S., Bianconi, G. & Cannistraci, C. V. Machine  
1419 learning meets complex networks via coalescent embedding in the hyperbolic space.  
1420 *Nat. Commun.* **8**, 1615 (2017).
- 1421 57. Muscoloni, A. & Cannistraci, C. V. Minimum curvilinear automata with similarity

- 1422 attachment for network embedding and link prediction in the hyperbolic space. (2018).
- 1423 58. Zagar, L. *et al.* Stage prediction of embryonic stem cell differentiation from genome-  
1424 wide expression data. *27*, 2546–2553 (2011).
- 1425 59. Ryu, T., Seridi, L. & Ravasi, T. The evolution of ultraconserved elements with  
1426 different phylogenetic origins. *BMC Evol. Biol.* **12**, 236 (2012).
- 1427 60. Sales, S. *et al.* Gender, Contraceptives and Individual Metabolic Predisposition Shape a  
1428 Healthy Plasma Lipidome. *Sci. Rep.* **6**, 27710 (2016).
- 1429 61. Acevedo, A., Ciucci, S., Kuo, M. J., Durán, C. & Cannistraci, C. V. Measuring group-  
1430 separability in geometrical space for evaluation of pattern recognition and embedding  
1431 algorithms. *ArXiv:1912.12418* 1–20 (2019).
- 1432 62. van Dongen, S. Graph clustering by flow simulation. *Graph Stimul. by flow Clust.*  
1433 (2000) doi:10.1016/j.cosrev.2007.05.001.
- 1434 63. Duran, C., Acevedo, A., Ciucci, S., Muscoloni, A. & Cannistraci, C. Nonlinear Markov  
1435 Clustering by Minimum Curvilinear Sparse Similarity. *ArXiv:1912.12211* 1–17 (2019).
- 1436 64. Ciucci, S. *et al.* Enlightening discriminative network functional modules behind  
1437 Principal Component Analysis separation in differential-omic science studies. 1–24  
1438 (2017) doi:10.1038/srep43946.
- 1439 65. Lim, R. *et al.* Large-scale metabolic interaction network of the mouse and human gut  
1440 microbiota. *Sci. Data* (2020) doi:10.1038/s41597-020-0516-5.
- 1441 66. Pang, Z., Chong, J., Li, S. & Xia, J. Metaboanalyst 3.0: Toward an optimized  
1442 workflow for global metabolomics. *Metabolites* (2020) doi:10.3390/metabo10050186.
- 1443 67. R Core Team. R: A language and environment for statistical computing. *R Foundation*  
1444 *for Statistical Computing* (2019).
- 1445 68. Csardi, G. & Nepusz, T. The igraph software package for complex network research.  
1446 *InterJournal Complex Syst.* (2006).
- 1447 69. Chamberlain, S. A. & Szöcs, E. Taxize: Taxonomic search and retrieval in R.



- 1448 *F1000Research* (2013) doi:10.12688/f1000research.2-191.v2.
- 1449 70. Sales, G., Calura, E., Cavalieri, D. & Romualdi, C. Graphite - a Bioconductor package  
1450 to convert pathway topology to gene network. *BMC Bioinformatics* (2012)  
1451 doi:10.1186/1471-2105-13-20.
- 1452 71. Gustavsen, J. A., Pai, S., Isserlin, R., Demchak, B. & Pico, A. R. RCy3: Network  
1453 biology using Cytoscape from within R. *F1000Research* (2019)  
1454 doi:10.12688/f1000research.20887.3.
- 1455 72. Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-independent,  
1456 community-supported software for describing and comparing microbial communities.  
1457 *Appl. Environ. Microbiol.* (2009) doi:10.1128/AEM.01541-09.
- 1458 73. Jones, D. L. The Fathom Toolbox for Matlab: multivariate ecological and  
1459 oceanographic data analysis. *Coll. Mar. Sci. Univ. South Florida, St. Petersburg, FL,*  
1460 *USA* (2014).
- 1461 74. Ammirati, E. *et al.* Patterns in ST-Elevation Acute Myocardial Infarction. *Circ. Res.*  
1462 **111**, 1336–1348 (2012).
- 1463 75. Montecucco, C. & Rappuoli, R. Living dangerously: how *Helicobacter pylori* survives  
1464 in the human stomach. *Nat. Rev. Mol. Cell Biol.* **2**, 457–466 (2001).
- 1465 76. Boguñá, M., Krioukov, D. & Claffy, K. C. Navigability of complex networks. *Nat.*  
1466 *Phys.* **5**, 74–80 (2008).
- 1467 77. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data.  
1468 *PLoS Comput. Biol.* **8**, (2012).
- 1469 78. Kurtz, Z. D. *et al.* Sparse and Compositionally Robust Inference of Microbial  
1470 Ecological Networks. *PLoS Comput. Biol.* **11**, e1004226 (2015).
- 1471 79. Wong, R. G., Wu, J. R. & Gloor, G. B. Expanding the UniFrac toolbox. *PLoS One* **11**,  
1472 e0161196 (2016).
- 1473 80. Weiss, S. *et al.* Normalization and microbial differential abundance strategies depend

- 1474 upon data characteristics. *Microbiome* **5**, 27 (2017).
- 1475 81. Navas-Molina, J. A. *et al.* Advancing our understanding of the human microbiome  
1476 using QIIME. in *Methods in Enzymology* vol. 531 371–444 (2013).
- 1477 82. Hughes, J. B. & Hellmann, J. J. The application of rarefaction techniques to molecular  
1478 inventories of microbial diversity. in *Methods in Enzymology* vol. 397 292–308 (2005).
- 1479 83. McMurdie, P. J., Holmes, S., Hoffmann, C., Bittinger, K. & Chen, Y. Waste Not, Want  
1480 Not: Why Rarefying Microbiome Data Is Inadmissible. *PLoS Comput. Biol.* **10**,  
1481 e1003531 (2014).
- 1482 84. Antharam, V. C. *et al.* Intestinal dysbiosis and depletion of butyrogenic bacteria in  
1483 *Clostridium difficile* infection and nosocomial diarrhea. *J. Clin. Microbiol.* **51**, 2884–  
1484 2892 (2013).
- 1485 85. Vesth, T. *et al.* Veillonella, Firmicutes: Microbes disguised as Gram negatives. *Stand.*  
1486 *Genomic Sci.* **9**, (2013).
- 1487 86. Bouwknegt, M., van Pelt, W., Kubbinga, M., Weda, M. & Havelaar, A. Potential  
1488 association between the recent increase in campylobacteriosis incidence in the  
1489 Netherlands and proton-pump inhibitor use – an ecological study. *Eurosurveillance* **19**,  
1490 20873 (2014).
- 1491 87. Leonard, J., Marshall, J. K. & Moayyedi, P. Systematic review of the risk of enteric  
1492 infection in patients taking acid suppression. *Am. J. Gastroenterol.* **102**, 2047–2056  
1493 (2007).
- 1494 88. Allaker, R. P. Non-sporing anaerobes: Wound infection; periodontal disease; abscess;  
1495 normal flora. *Med. Microbiol. Eighteenth Ed.* 359–364 (2012) doi:10.1016/B978-0-  
1496 7020-4089-4.00051-2.
- 1497 89. Eribe, E. R. K. & Olsen, I. Leptotrichia species in human infections II. *J. Oral*  
1498 *Microbiol.* **9**, 1368848 (2017).
- 1499 90. Liu, D. *Molecular detection of human bacterial pathogens.* (CRC press, 2011).

- 1500 91. Carlier, J.-P. *Oribacterium*. in *Bergey's Manual of Systematics of Archaea and*  
1501 *Bacteria* 1–5 (John Wiley & Sons, Ltd, 2015).  
1502 doi:10.1002/9781118960608.gbm00649.
- 1503 92. Wang, K. *et al.* Preliminary analysis of salivary microbiome and their potential roles in  
1504 oral lichen planus. *Sci. Rep.* **6**, 22943 (2016).
- 1505 93. Torok, E., Moran, E. & Cooke, F. *Oxford Handbook of Infectious Diseases and*  
1506 *Microbiology*. (Oxford University Press, 2009).  
1507 doi:10.1093/med/9780198569251.001.0001.
- 1508 94. Jolivet-Gougeon, A., Sixou, J.-L., Tamanai-Shacoori, Z. & Bonnaure-Mallet, M.  
1509 Antimicrobial treatment of Capnocytophaga infections. *Int. J. Antimicrob. Agents* **29**,  
1510 367–373 (2007).
- 1511 95. Piau, C., Arvieux, C., Bonnaure-Mallet, M. & Jolivet-Gougeon, A. Capnocytophaga  
1512 spp. involvement in bone infections: a review. *Int. J. Antimicrob. Agents* **41**, 509–515  
1513 (2013).
- 1514 96. Cargill, J. S., Scott, K. S., Gascoyne-Binzi, D. & Sandoe, J. A. T. Granulicatella  
1515 infection: Diagnosis and management. *J. Med. Microbiol.* **61**, 755–761 (2012).
- 1516 97. Hofstad, T. The Genus *Fusobacterium*. in *The Prokaryotes* 1016–1027 (Springer New  
1517 York, 2006). doi:10.1007/0-387-30747-8.
- 1518 98. Brophy, S. *et al.* Incidence of Campylobacter and Salmonella Infections Following  
1519 First Prescription for PPI: A Cohort Study Using Routine Data. *Am. J. Gastroenterol.*  
1520 **108**, 1094–1100 (2013).
- 1521 99. Allos, B. M. Campylobacter infections. in *Bacterial Infections of Humans:*  
1522 *Epidemiology and Control* 189–211 (Springer US, 2009). doi:10.1007/978-0-387-  
1523 09843-2\_9.
- 1524 100. Lee, C. & Hong, S. N. Does long-term proton pump inhibitor therapy affect the health  
1525 of gut microbiota? *Gut and Liver* vol. 10 865–866 (2016).

- 1526 101. Seto, C. T., Jeraldo, P., Orenstein, R., Chia, N. & DiBaise, J. K. Prolonged use of a  
1527 proton pump inhibitor reduces microbial diversity: Implications for *Clostridium*  
1528 *difficile* susceptibility. *Microbiome* **2**, (2014).
- 1529 102. Bavishi, C. & DuPont, H. L. Systematic review: The use of proton pump inhibitors and  
1530 increased susceptibility to enteric infection. *Alimentary Pharmacology and*  
1531 *Therapeutics* vol. 34 1269–1281 (2011).
- 1532 103. Olbe, L. *Proton pump inhibitors*. (Birkhäuser, 2012).
- 1533 104. Warren, J. R. & Marshall, B. Unidentified curved bacilli on gastric epithelium in active  
1534 chronic gastritis. *Lancet* **321**, 1273–1275 (1983).
- 1535 105. Ha, N. *et al.* Supramolecular assembly and acid resistance of *Helicobacter pylori*  
1536 urease. *Nat. Struct. Biol.* **8**, 505–509 (2001).
- 1537 106. Berger, A. Scientists discover how helicobacter survives gastric acid. *Br. Med. J.* **29**,  
1538 268 (2000).
- 1539 107. Amieva, M. R. & El-Omar, E. M. Host-Bacterial Interactions in *Helicobacter pylori*  
1540 Infection. *Gastroenterology* **134**, 306–323 (2008).
- 1541 108. Scott Merrell, D. *et al.* Adhesion and Invasion of Gastric Mucosa Epithelial Cells by  
1542 *Helicobacter pylori*. *Front. Cell. Infect. Microbiol* **6**, 1593389–159 (2016).
- 1543 109. von Rosenvinge, E. C. *et al.* Immune status, antibiotic medication and pH are  
1544 associated with changes in the stomach fluid microbiota. *ISME J.* **7**, 1354–1366 (2013).
- 1545 110. Eun, C. S. o. *et al.* Differences in gastric mucosal microbiota profiling in patients with  
1546 chronic gastritis, intestinal metaplasia, and gastric cancer using pyrosequencing  
1547 methods. *Helicobacter* **19**, 407–416 (2014).
- 1548 111. Cao, L. & Yu, J. Effect of *Helicobacter pylori* Infection on the Composition of Gastric  
1549 Microbiota in the Development of Gastric Cancer. *Gastrointest. tumors* **2**, 14–25  
1550 (2015).
- 1551 112. Brawner, K. M., Morrow, C. D. & Smith, P. D. Gastric microbiome and gastric cancer.

- 1552 *Cancer J.* **20**, 211–6 (2014).
- 1553 113. Cover, T. L. & Blaser, M. J. Helicobacter pylori in health and disease.  
1554 *Gastroenterology* **136**, 1863–73 (2009).
- 1555 114. Sanders, M. K. & Peura, D. A. Helicobacter pylori-Associated Diseases. *Curr.*  
1556 *Gastroenterol. Rep.* **4**, 448–54 (2002).
- 1557 115. Talley, N. J. Helicobacter pylori and dyspepsia. *Yale J. Biol. Med.* **72**, 145–51 (1999).
- 1558 116. Shadwell, J. Helicobacter pylori–associated dyspepsia. *2016*.
- 1559 117. Noto, J. M. & Peek, R. M. The gastric microbiome, its interaction with Helicobacter  
1560 pylori, and its potential role in the progression to stomach cancer. *PLoS Pathogens* vol.  
1561 13 (2017).
- 1562 118. Schwabe, R. F. & Jobin, C. The microbiome and cancer. *Nature Reviews Cancer* vol.  
1563 13 800–812 (2013).
- 1564 119. Fraher, M. H., O’Toole, P. W. & Quigley, E. M. M. Techniques used to characterize  
1565 the gut microbiota: a guide for the clinician. *Nat. Rev. Gastroenterol. Hepatol.* **9**, 312–  
1566 322 (2012).
- 1567 120. Andersson, A. F. *et al.* Comparative Analysis of Human Gut Microbiota by Barcoded  
1568 Pyrosequencing. *PLoS One* **3**, e2836 (2008).
- 1569 121. Bik, E. M. Molecular analysis of the bacterial microbiota in the human stomach. *Proc.*  
1570 *Natl. Acad. Sci. USA* **103**, 732–737 (2006).
- 1571 122. Llorca, L. *et al.* Characterization of the gastric microbiota in a pediatric population  
1572 according to Helicobacter pylori status. in *Pediatric Infectious Disease Journal* vol. 36  
1573 173–178 (2017).
- 1574 123. Jo, H. J. The effect of H. pylori infection on the gastric microbiota. in *Helicobacter*  
1575 *pylori* (ed. Kim, N.) 529–533 (Springer Singapore, 2016). doi:10.1007/978-981-287-  
1576 706-2\_54.
- 1577 124. Klymiuk, I. *et al.* The Human Gastric Microbiome Is Predicated upon Infection with

- 1578 Helicobacter pylori. *Front. Microbiol.* **8**, 2508 (2017).
- 1579 125. Maldonado-Contreras, A. *et al.* Structure of the human gastric bacterial community in  
1580 relation to Helicobacter pylori status. *ISME J.* **5**, 574–579 (2011).
- 1581 126. Aviles-Jimenez, F., Vazquez-Jimenez, F., Medrano-Guzman, R., Mantilla, A. &  
1582 Torres, J. Stomach microbiota composition varies between patients with non-atrophic  
1583 gastritis and patients with intestinal type of gastric cancer. *Sci. Rep.* **4**, 4202 (2015).
- 1584 127. Kovaleva, J., Degener, J. E. & van der Mei, H. C. Methylobacterium and its role in  
1585 health care-associated infection. *J. Clin. Microbiol.* **52**, 1317–21 (2014).
- 1586 128. White, D. C., Sutton, S. D. & Ringelberg, D. B. The genus Sphingomonas: physiology  
1587 and ecology. *Curr. Opin. Biotechnol.* **7**, 301–306 (1996).
- 1588 129. Madigan, M., Martinko, J., Stahl, D. and Clark, D. Brock Biology of Microorganisms.  
1589 321 (2012).
- 1590 130. Özen, A. I. & Ussery, D. W. Defining the Pseudomonas genus: where do we draw the  
1591 line with Azotobacter? *Microb. Ecol.* **63**, 239–48 (2012).
- 1592 131. Towner, K. The genus Acinetobacter. in *The Prokaryotes* 545–577 (Springer New  
1593 York, 2006). doi:10.1007/978-3-642-30194-0.
- 1594 132. Rathinavelu, S., Zavros, Y. & Merchant, J. L. Acinetobacter lwoffii infection and  
1595 gastritis. *Microbes Infect.* **5**, 651–657 (2003).
- 1596 133. Cheung, Y. F., Walsh, C. & Fung, C. H. Stereochemistry of Propionyl-Coenzyme A  
1597 and Pyruvate Carboxylations Catalyzed by Transcarboxylase. *Biochemistry* **14**, 2981–  
1598 2986 (1975).
- 1599 134. Piwowarek, K., Lipińska, E., Hać-Szymańczuk, E., Kieliszek, M. & Ścibisz, I.  
1600 Propionibacterium spp.—source of propionic acid, vitamin B12, and other metabolites  
1601 important for the industry. *Applied Microbiology and Biotechnology* vol. 102 515–538  
1602 (2018).
- 1603 135. Moore, L. V. H. & Moore, W. E. C. Oribaculum catoniae gen. nov., sp. nov.; Catonella

- 1604 morbi gen. nov., sp. nov.; *Hallella seregens* gen. nov., sp. nov.; *Johnsonella ignava* gen.  
1605 nov., sp. nov.; and *Dialister pneumosintes* gen. nov., comb. nov., nom. rev., Anaerobic  
1606 Gram-Negative Bacilli from. *Int. J. Syst. Bacteriol.* **44**, 187–192 (1994).
- 1607 136. Willems, A. & Collins, M. D. *Catonella*. in *Bergey's Manual of Systematics of*  
1608 *Archaea and Bacteria* 1–7 (John Wiley & Sons, Ltd, 2015).  
1609 doi:10.1002/9781118960608.gbm00641.
- 1610 137. Menon, T. & Kumar, V. N. *Catonella morbi* as a cause of native valve endocarditis in  
1611 Chennai, India. *Infection* **40**, 581–582 (2012).
- 1612 138. Balows, A., Truper, H., Dworkin, M., Harder, W. & Schleifer, K. *The Prokaryotes. A*  
1613 *Handbook on the Biology of Bacteria: Proteobacteria: Gamma subclass. The*  
1614 *prokaryotes* (Springer, 1991). doi:10.1007/0-387-30745-1.
- 1615 139. Staley, J. T., Irgens, R. L. & Brenner, D. J. *Enhydrobacter aerosaccus* gen. nov., sp.  
1616 nov., a Gas-Vacuolated, Facultatively Anaerobic, Heterotrophic Rod. *Int. J. Syst.*  
1617 *Bacteriol.* **37**, 289–291 (1987).
- 1618 140. Wade, W. G. & Downes, J. *Bulleidia*. *Bergey's Manual of Systematics of Archaea and*  
1619 *Bacteria* (2015) doi:doi:10.1002/9781118960608.gbm00760.
- 1620 141. Kienesberger, S. *et al.* Gastric *Helicobacter pylori* Infection Affects Local and Distant  
1621 Microbial Populations and Host Responses. *Cell Rep.* **14**, 1395–1407 (2016).
- 1622 142. Amato, S. M. *et al.* The role of metabolism in bacterial persistence. *Frontiers in*  
1623 *Microbiology* (2014) doi:10.3389/fmicb.2014.00070.
- 1624 143. Li, Z. *et al.* Effects of metabolites derived from gut microbiota and hosts on pathogens.  
1625 *Frontiers in Cellular and Infection Microbiology* (2018)  
1626 doi:10.3389/fcimb.2018.00314.
- 1627 144. Vojinovic, D. *et al.* Relationship between gut microbiota and circulating metabolites in  
1628 population-based cohorts. *Nat. Commun.* (2019) doi:10.1038/s41467-019-13721-1.
- 1629 145. Del Chierico, F. *et al.* Gut microbiota markers in obese adolescent and adult patients:

- 1630 Age-dependent differential patterns. *Front. Microbiol.* (2018)  
1631 doi:10.3389/fmicb.2018.01210.
- 1632 146. Karlsson, F. H. *et al.* Symptomatic atherosclerosis is associated with an altered gut  
1633 metagenome. *Nat. Commun.* (2012) doi:10.1038/ncomms2266.
- 1634 147. Luo, L. *et al.* Association between metabolic profile and microbiomic changes in rats  
1635 with functional dyspepsia. *RSC Adv.* (2018) doi:10.1039/c8ra01432a.
- 1636 148. Ma, S. *et al.* Alterations in Gut Microbiota of Gestational Diabetes Patients During the  
1637 First Trimester of Pregnancy. *Front. Cell. Infect. Microbiol.* (2020)  
1638 doi:10.3389/fcimb.2020.00058.
- 1639 149. Cai, X. *et al.* Altered Diversity and Composition of Gut Microbiota in Wilson's  
1640 disease. *Sci. Rep.* 1–10 (2020) doi:10.21203/rs.2.24572/v1.
- 1641 150. Severi, E., Hood, D. W. & Thomas, G. H. Sialic acid utilization by bacterial pathogens.  
1642 *Microbiology* (2007) doi:10.1099/mic.0.2007/009480-0.
- 1643 151. Vimr, E. R., Kalivoda, K. A., Deszo, E. L. & Steenbergen, S. M. Diversity of  
1644 Microbial Sialic Acid Metabolism. *Microbiol. Mol. Biol. Rev.* (2004)  
1645 doi:10.1128/mmbr.68.1.132-153.2004.
- 1646 152. Zhou, X., Yang, G. & Guan, F. Biological Functions and Analytical Strategies of Sialic  
1647 Acids in Tumor. *Cells* (2020) doi:10.3390/cells9020273.
- 1648 153. Gonzalez, A. *et al.* Migraines Are Correlated with Higher Levels of Nitrate-, Nitrite-,  
1649 and Nitric Oxide-Reducing Oral Microbes in the American Gut Project Cohort.  
1650 *mSystems* (2016) doi:10.1128/msystems.00105-16.
- 1651 154. Kobayashi, J. Effect of diet and gut environment on the gastrointestinal formation of  
1652 N-nitroso compounds: A review. *Nitric Oxide - Biology and Chemistry* (2018)  
1653 doi:10.1016/j.niox.2017.06.001.
- 1654 155. Hughes, R. & Rowland, I. R. Metabolic activities of the gut microflora in relation to  
1655 cancer. *Microb. Ecol. Health Dis.* (2000) doi:10.1080/089106000750060431.



- 1656 156. Verdu, E. *et al.* Effect of omeprazole on intragastric bacterial counts, nitrates, nitrites,  
1657 and N-nitroso compounds. *Gut* (1994) doi:10.1136/gut.35.4.455.
- 1658 157. Durán, C. *et al.* Pioneering topological methods for network-based drug–target  
1659 prediction by exploiting a brain-network self-organization theory. *Brief. Bioinform.* 1–  
1660 20 (2017) doi:10.1093/bib/bbx041.
- 1661 158. Muscoloni, A., Abdelhamid, I., Decano, J. L., Souza, E. & Maiorino, E. Hyperedge  
1662 entanglement in high-order multilayer networks. (2020)  
1663 doi:10.20944/preprints202012.0500.v1.

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674

1675

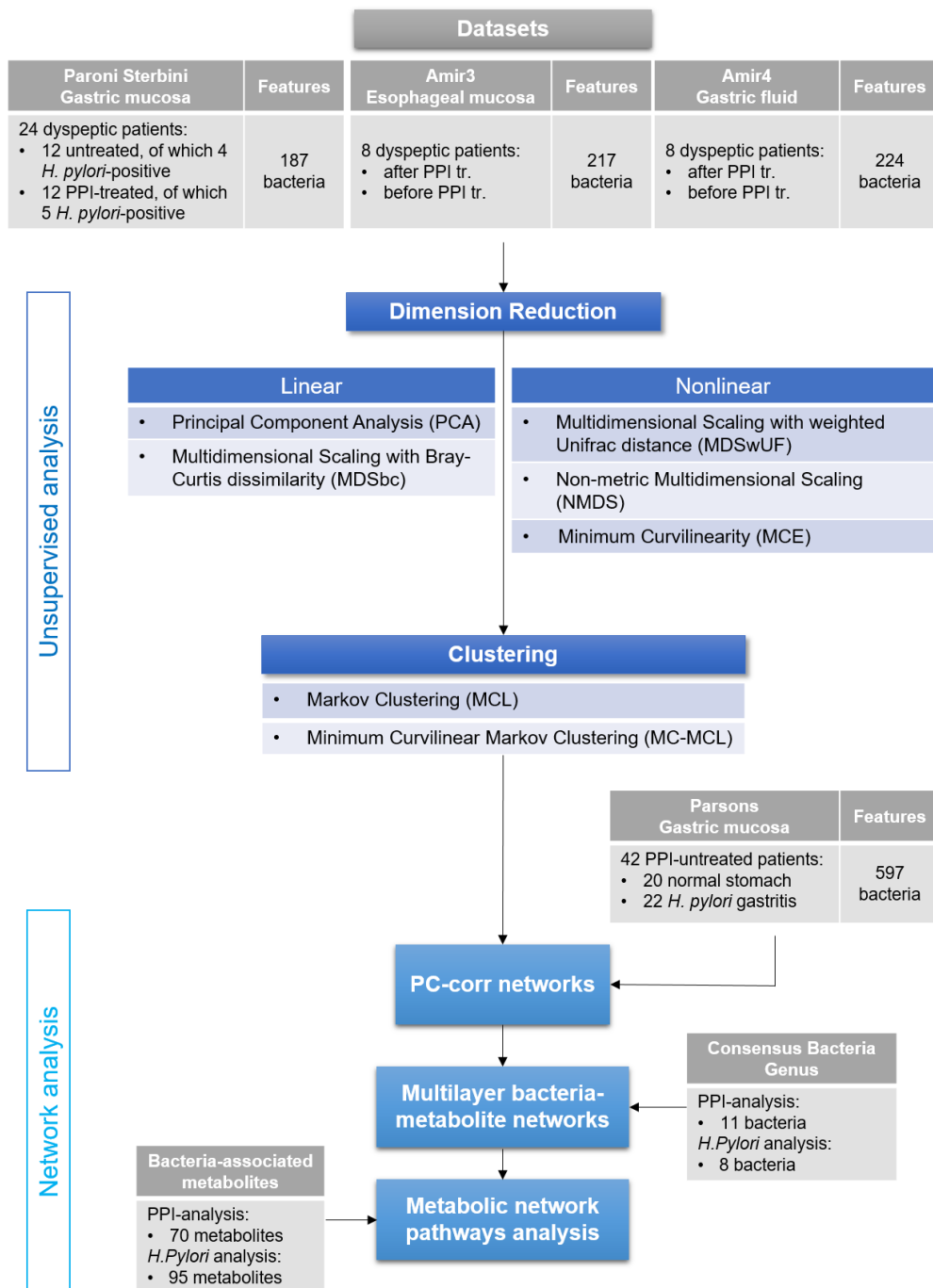
1676

1677

1678

1679 **Figures and tables**

1680



1681

1682 **Figure 1. Flowchart of the data analysis.** To answer the five questions under investigation in our study,

1683 we implemented a workflow based on machine learning tools. Following the flowchart shown in the

1684 figure, we analysed three 16S rRNA gene sequencing datasets with information on PPI use in dyspeptic

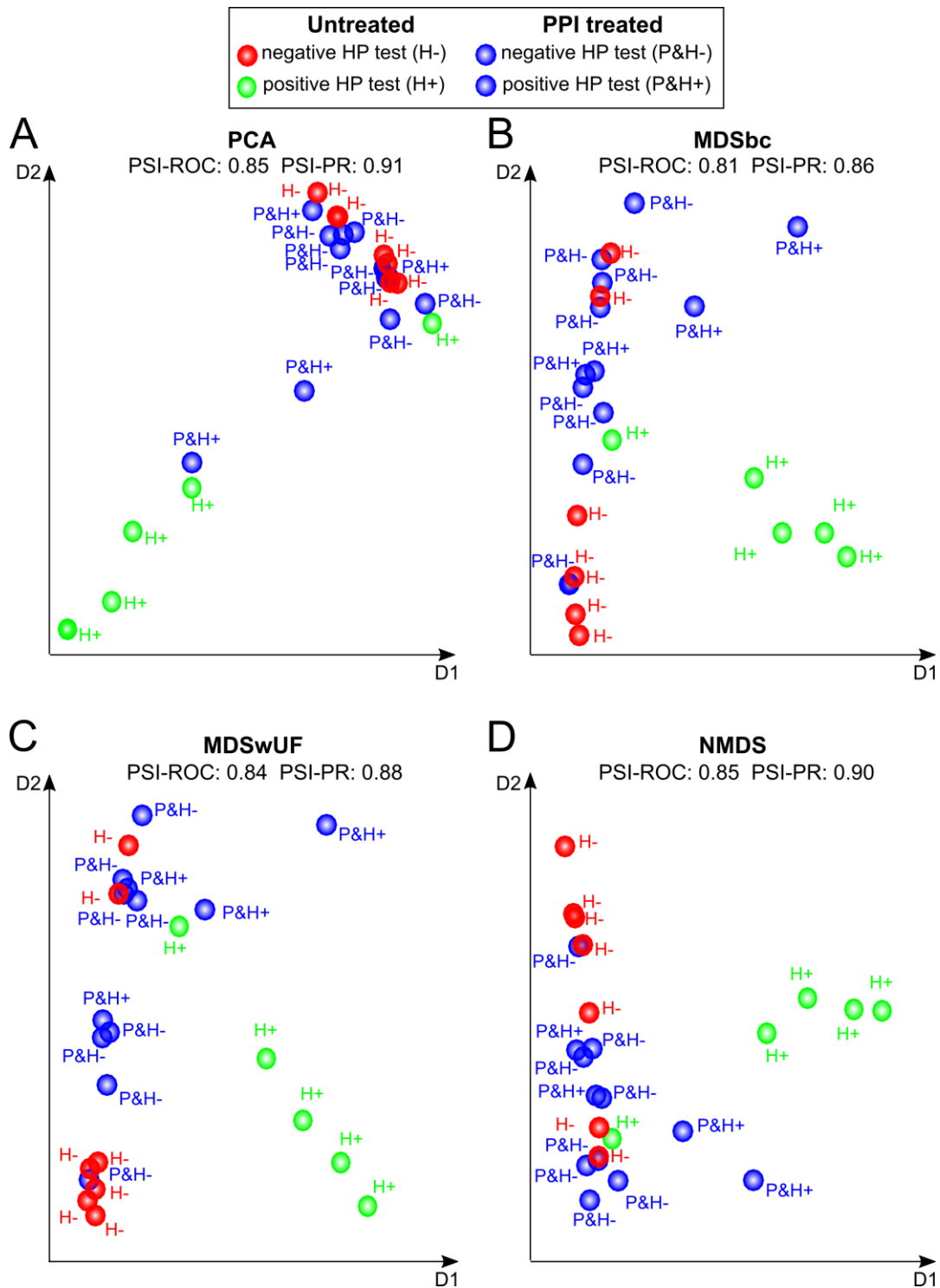
1685 patients; for one of the datasets (Paroni Sterbini *et al.*<sup>22</sup>), patients were also determined to be positive

1686 or negative to *H. pylori* infection.

1687 Firstly, we performed unsupervised dimension reduction, both linear and nonlinear, in the first two

1688 dimensions of embedding. Nonlinear dimension reduction will show the presence of hidden patterns,

1689 in the form of sample groups. Secondly, nonlinear clustering was applied to confirm the well-  
1690 possessedness of the hidden patterns found by nonlinear dimension reduction. Furthermore, our workflow  
1691 ends with the network analysis. It starts with the use of the PC-corr algorithm, that reveals which  
1692 combination of bacteria (features) are responsible for the identified differences between the groups of  
1693 samples. A fourth dataset (Parsons *et al*<sup>32</sup>.) is used only for the validation of the PC-corr network results  
1694 and it contains information of PPI treatment and *H. pylori* infection. From the consensus bacteria found  
1695 in each PC-corr network, a bacteria-metabolite multilayer analysis that lastly end with the metabolite  
1696 pathway enrichment analysis that introduces evidence to possible perturbed biological mechanisms.



1697

1698 **Figure 2. Dimension reduction techniques usually employed in metagenomic data analysis and**

1699 **applied to the Paroni Sterbini dataset.** The plots represent the best PCA and MDS results based on

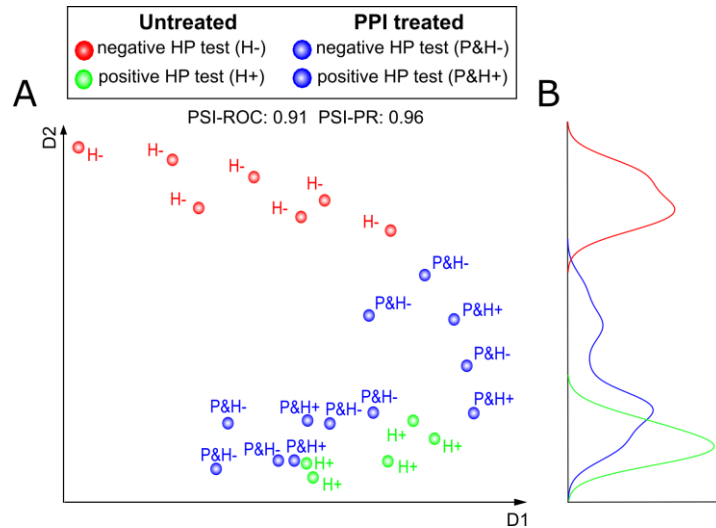
1700 (average) p-value projection-based separability index (PSI) for the three different labels (PPI-treated,

1701 untreated H+ and untreated H-), evaluated in the 2D embedding space. Moreover, also the average

1702 values of all pairwise PSI-ROC and PSI-PR are reported as overall estimators of separation between the

1703 groups in the 2D reduced space. (A) PCA; (B) MDS with Bray-Curtis dissimilarity (MDSbc); (C) MDS

1704 with weighted UniFrac distance (MDSwUF); **(D)** non-metric MDS with Sammon Mapping (NMDS).  
 1705 Blue dots represent PPI-treated samples, while red and green dots are the untreated samples which  
 1706 resulted either negative (red) or positive (green) to the *H. pylori* test (histological observation and urease  
 1707 test).



1708 **Figure 3. MCE, a topological machine learning for nonlinear and hierarchical dimension**  
 1709 **reduction.** **A)** Results on the Paroni Sterbini et al.<sup>22</sup> dataset. The shown best MCE result is based on  
 1710 PSI-PR projection-based separability index (PSI) for the three different labels (P-treated, untreated H+  
 1711 and untreated H-), evaluated in the 2D embedding space under the DCS normalization. The PSI-ROC  
 1712 and PSI-PR are reported as overall estimators of separation between the groups in the 2D reduced space.  
 1713 Blue dots represent PPI-treated samples, while red and green dots are the untreated samples which  
 1714 resulted either negative (red) or positive (green) to the *H. pylori* test (histological observation and urease  
 1715 test). **B)** The curves in three different colours (red, blue and green) highlight the different distributions  
 1716 of the three groups on the second dimension.

**Table 1. Results of unsupervised analysis on the real datasets.** Best results of unsupervised dimension reduction techniques (top panel) and of clustering (bottom panel).

**(Top panel):** Best results of unsupervised dimension reduction techniques according to the PSI indices for sample separation in the space of the first two dimensions of embedding. HD (no dimension reduction) represents the reference results to see how good the separability present in the high dimensional space is preserved by dimension reduction techniques. Results are ordered from the best (top) to the worst (bottom) method. For the Paroni Sterbini dataset, we show the results for three different labels (PPI-treated, untreated H+ and untreated H-). For the Amir datasets, the PSI measures were computed for two groups, identified by the presence or absence of PPI treatment. For each PSI value, a respective trustworthiness was calculated.

**(Bottom panel):** Best results of clustering (highest accuracies, regardless of the normalization and type of correlation) MCL and MC-MCL, in each of the three studied datasets (Paroni Sterbini, Amir3 and Amir4), and the mean performance (mean of the highest accuracies) across all the datasets.

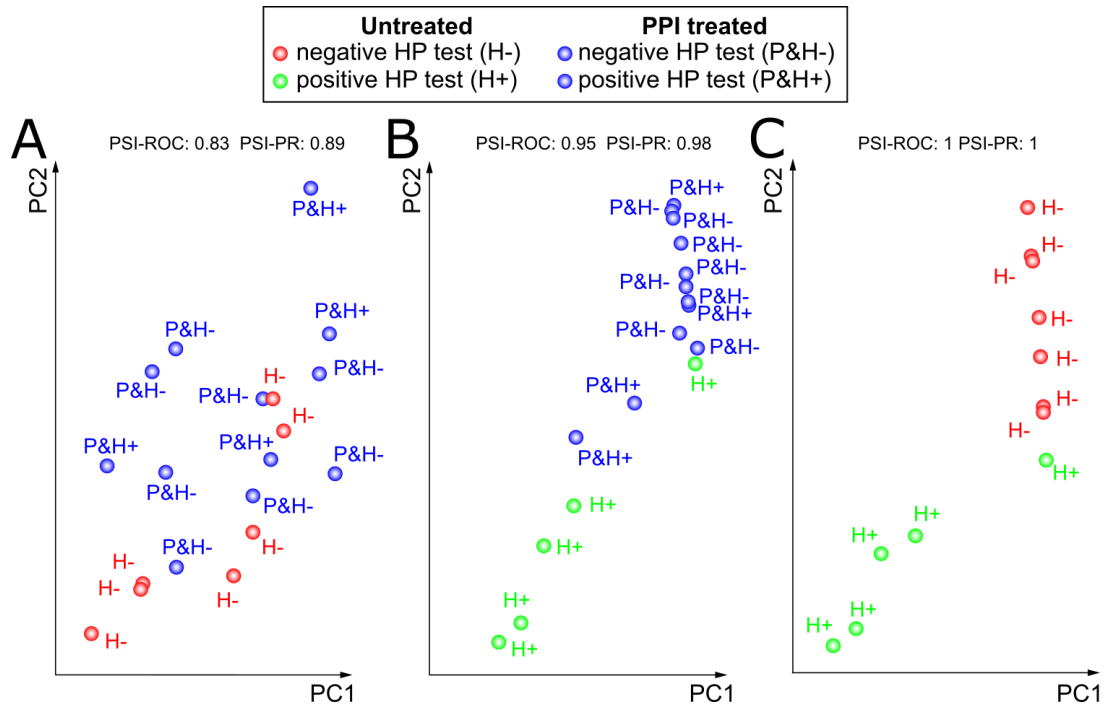
For Paroni Sterbini dataset, we show the results for three clusters (PPI-treated, untreated H+ and untreated H-) and in brackets the results for four clusters (P&H+, P&H-, untreated H+ and untreated H-). Instead, for Amir datasets, the accuracies were computed for two groups, identified according to

the presence or absence of PPI treatment.

		PSI-ROC						
		Method	Paroni Sterbini	Trust	Amir3	Trust	Amir4	Trust
Dimension Reduction	HD	0.88	0.0036	0.95	0.0009	0.98	0.0009	0.94
	MDSwUF	0.84	0.0089	1.00	0.0009	0.88	0.0329	0.90
	MCE	0.91	0.0036	0.88	0.0329	0.91	0.0009	0.90
	PCA	0.85	0.0063	0.91	0.0009	0.86	0.0169	0.87
	MDStyc	0.84	0.0076	0.88	0.0009	0.84	0.0249	0.85
	nMDS	0.85	0.0036	0.86	0.0169	0.84	0.0089	0.85
	MDSbc	0.81	0.0183	0.86	0.0089	0.84	0.0189	0.84
		PSI-PR						
		Method	Paroni Sterbini	Trust	Amir3	Trust	Amir4	Trust
Dimension Reduction	HD	0.94	0.0009	0.96	0.0009	0.99	0.0009	0.96
	MDSwUF	0.88	0.0036	1.00	0.0009	0.90	0.0089	0.93
	MCE	0.96	0.0009	0.89	0.0089	0.92	0.0039	0.92
	PCA	0.91	0.0039	0.90	0.0009	0.88	0.0089	0.90
	MDStyc	0.88	0.0116	0.90	0.0009	0.88	0.0089	0.89
	MDSbc	0.86	0.0116	0.89	0.0009	0.90	0.0009	0.88
	nMDS	0.90	0.0036	0.87	0.0089	0.87	0.0009	0.88
		Clustering						
		Accuracy	Paroni Sterbini	Amir3	Amir4	Mean performance		
Clustering	MC-MCL	0.71 (0.58)	0.81	0.75	0.76			
	MCL	0.67 (0.63)	0.69	0.75	0.70			

Note: all PSI-ROC and PSI-PR values can be found in Supplementary Table S2, while all the accuracies can be found in Supplementary Table S17.

**Abbreviations:** HD: High Dimension; MCE: Minimum Curvilinear Embedding; MDSbc: Multidimensional Scaling with Bray-Curtis dissimilarity; MDSwUF: Multidimensional Scaling with weighted UniFrac distance; NMDS: Non-metric Multidimensional Scaling; MDStyc: Multidimensional Scaling with Theta-YC distance; PCA: Principal Component Analysis; MCL: Markov Clustering; MC-MCL: Minimum Curvilinear Markov Clustering; PSI-ROC: Projection Separability Index measured by Area Under the Curve; PSI-PR: Projection Separability Index measured by Area Under the Precision Recall; Trust: Trustworthiness.



**Figure 4. Pairwise PCA of Paroni Sterbini's gastric samples.** PCA was applied to three subsampled versions of the Paroni Sterbini dataset (keeping the best normalization found for the original dataset), each corresponding to the combination of two groups: (A) PPI-treated and untreated *H. pylori* negative samples; (B) PPI-treated and untreated *H. pylori* positive samples; (C) untreated *H. pylori* negative and untreated *H. pylori* positive samples. The PSI-ROC and PSI-PR are reported as well as overall estimators of separation between the groups in the 2D reduced space.

**Table 2. Ranked performance of unsupervised dimension reduction techniques on the real datasets.** The table shows the ranked performance of unsupervised dimension reduction techniques according to the PSI indices for sample separation (PSI-ROC and PSI-PR) in the space of the first two dimensions of embedding, for the three studied datasets (Paroni Sterbini, Amir3 and Amir4). Each rank is related to the results obtained in Table 1, top panel. The results are ordered by the mean performance

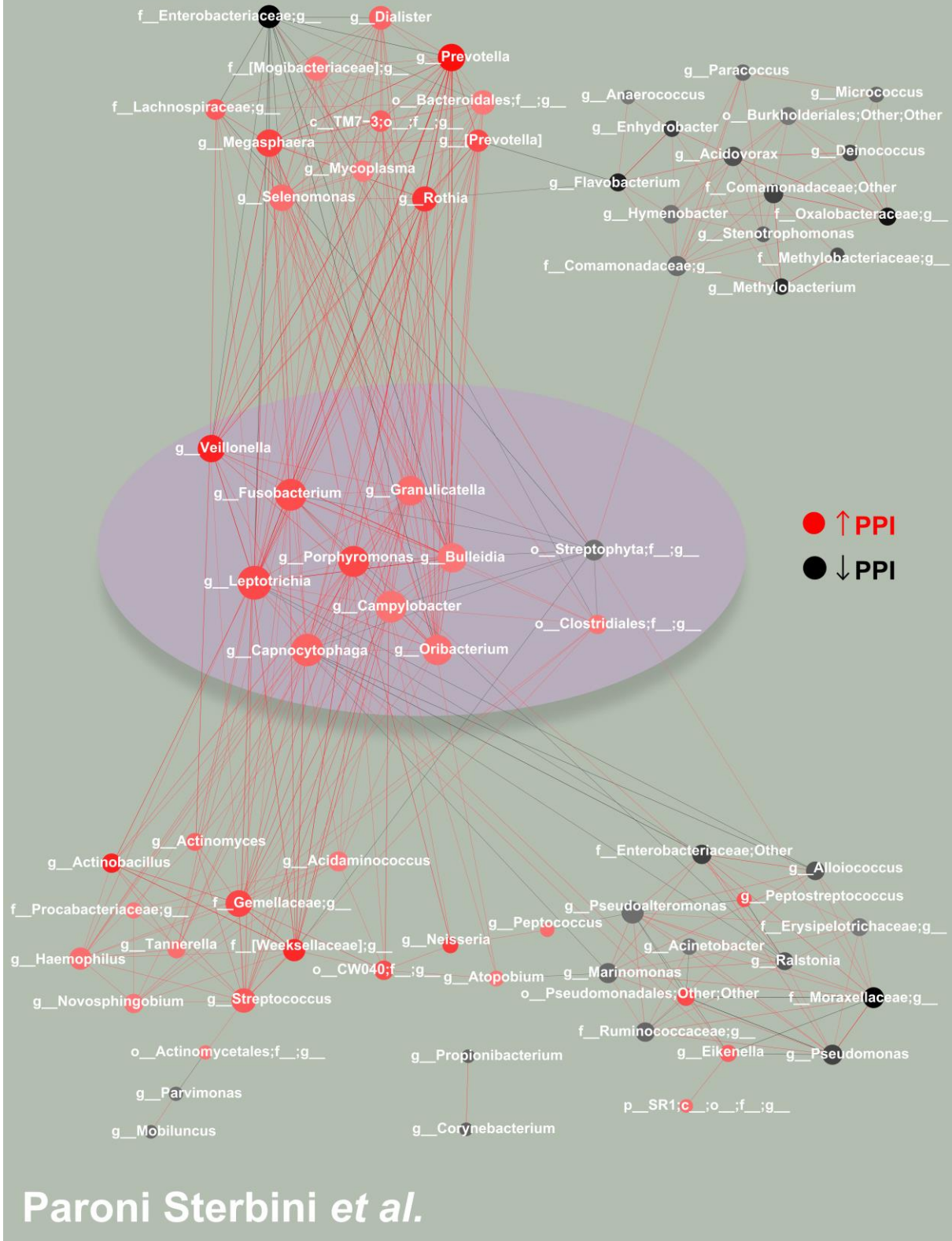


(fourth column) from the best (top) to the worst (bottom) method.

PSI-ROC					PSI-PR				
Method	Paroni Sterbini	Amir3	Amir4	mean	Method	Paroni Sterbini	Amir3	Amir4	mean
HD	2	2	1	1.67	HD	2	2	1	1.67
MCE	1	4	2	2.33	MCE	1	5	2	2.67
MDSwUF	5	1	3	3.00	MDSwUF	5	1	3	3.00
PCA	3	3	4	3.33	PCA	3	3	5	3.67
nMDS	3	6	5	4.67	MDStyc	5	3	5	4.33
MDStyc	5	4	5	4.67	MDSbc	7	5	3	5.00
MDSbc	7	6	5	6.00	nMDS	4	7	7	6.00

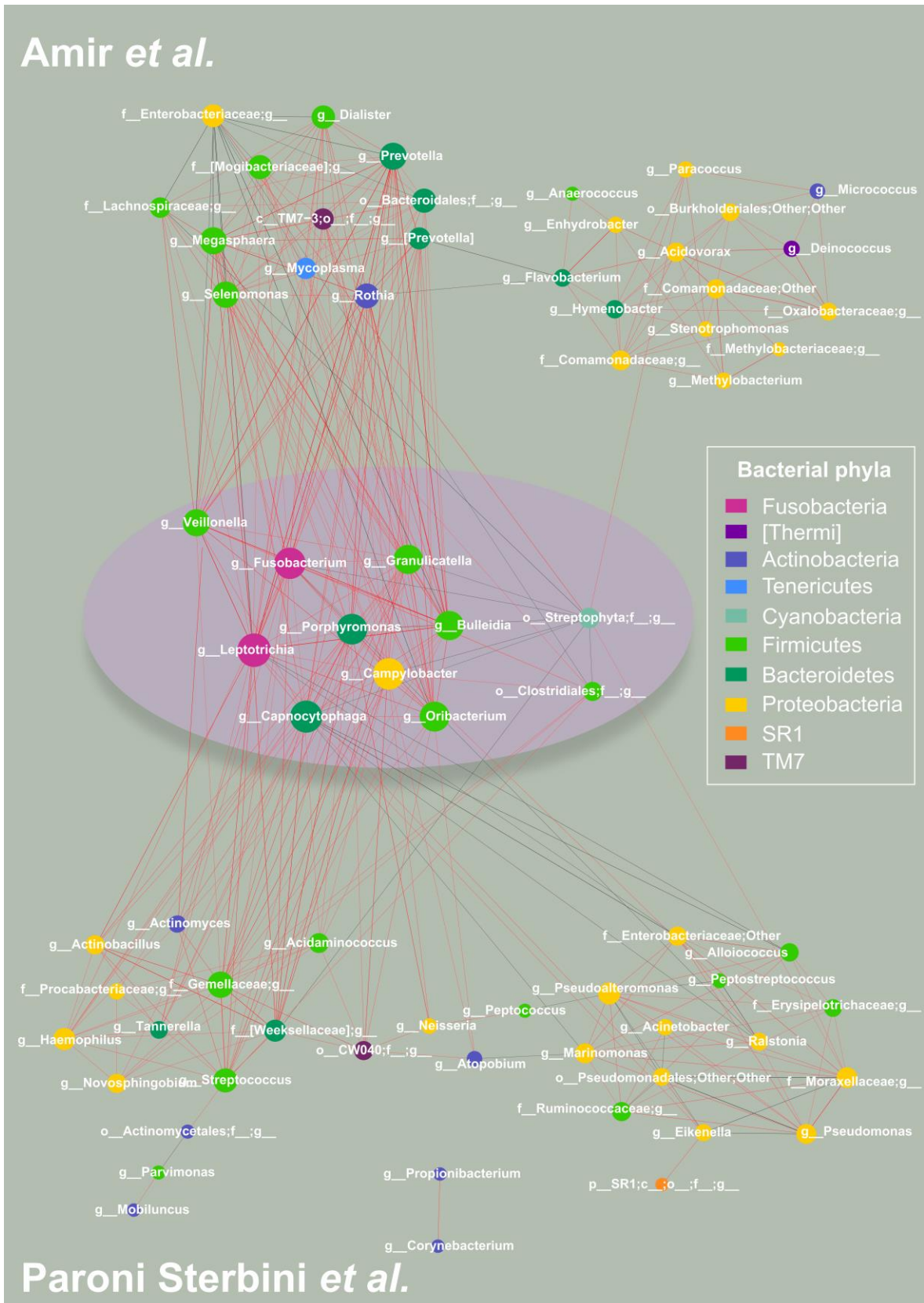
**Abbreviations:** HD: High Dimension; MCE: Minimum Curvilinear Embedding; MDSbc: Multidimensional Scaling with Bray-Curtis dissimilarity; MDSwUF: Multidimensional Scaling with weighted UniFrac distance; nMDS: Non-metric Multidimensional Scaling; MDStyc: Multidimensional Scaling with Theta-YC distance; PCA: Principal Component Analysis; PSI-ROC: Projection Separability Index measured by Area Under the Curve; PSI-PR: Projection Separability Index measured by Area Under the Precision Recall.

Amir et al.



1717 **Figure 5. PC-corr method to unveil how PPI is affecting the microbiota in gastric environment in**  
1718 **dyspeptic patients. (Middle panel)** To investigate the effect of PPIs on the gastric microbiota in  
1719 dyspeptic patients, we constructed the conserved PC-corr network at 0.5 cut-off, by merging the PC-

1720 corr networks obtained from the gastric mucosa (Paroni Sterbini *et al.* <sup>22</sup>) and the gastric fluid (Amir *et*  
1721 *al.* <sup>21</sup>). To do so, we firstly considered the union of the two PC-corr networks obtained from the gastric  
1722 tissue dataset and then we intersected it with the PC-corr network from the gastric fluid dataset. All the  
1723 bacteria spotted in the conserved PC-corr network (violet circle) were found increased with PPI use. In  
1724 both the two studied datasets, red nodes indicate bacteria whose abundance is increased with PPI-  
1725 treatment, while black nodes indicate bacteria with lower abundance following treatment with this acid  
1726 suppressing medication. The common bacteria that showed an opposite trend in the two datasets, i.e.  
1727 microbial abundance increased in one dataset and decreased in the other dataset, were removed from the  
1728 network. (**Top panel**) The top panel shows the obtained Amir4's network, not in common with the  
1729 Paroni Sterbini's network. The module on the left side (except *Enterobacteriaceae*) include bacteria  
1730 more abundant following PPI-treatment in Amir4's data, while the module on the right (and  
1731 *Enterobacteriaceae*) is composed of decreased bacteria in abundance under PPI therapy in Amir4's data.  
1732 (**Bottom panel**) The bottom panel represents the part of Paroni Sterbini's network (union of the two  
1733 PC-corr network), that is not shared with Amir4's one. As in the top and middle panels, the colour of  
1734 the nodes represents if the bacteria display higher (red nodes) or lower abundance (black nodes) in PPI-  
1735 treated samples of Paroni Sterbini's dataset.

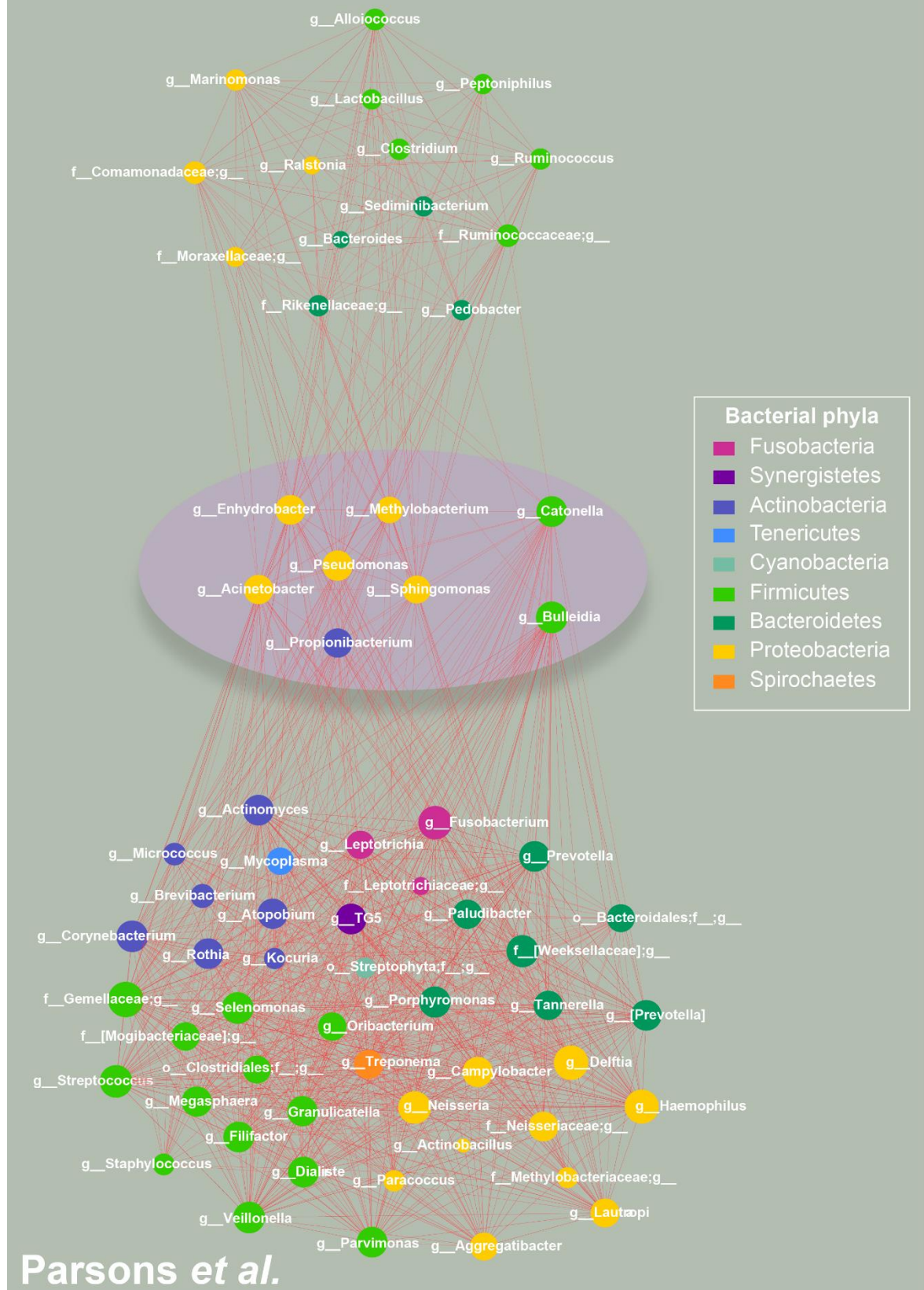


1736 **Figure 6. PC-corr networks to unveil how PPI is affecting the microbiota in gastric environment**  
 1737 **in dyspeptic patients, coloured according to phylum-level taxonomy.** To investigate the effect of  
 1738 PPIs on the gastric microbiota in dyspeptic patients, we constructed the conserved PC-corr network at

1739 0.5 cut-off, by merging the PC-corr networks obtained from the gastric mucosa (Paroni Sterbini *et al.*  
1740 <sup>22</sup>) and the gastric fluid (Amir *et al.* <sup>21</sup>). To do so, we firstly considered the union of the two PC-corr  
1741 networks obtained from the gastric tissue dataset and then we intersected it with the PC-corr network  
1742 from the gastric fluid dataset. All the bacteria spotted in the conserved PC-corr network (violet circle)  
1743 were found increased with PPI use. (**Top panel**) The top panel shows the obtained Amir4's network,  
1744 not in common with the Paroni Sterbini's network. The module on the left side (except  
1745 *Enterobacteriaceae*) include bacteria more abundant following PPI-treatment in Amir4's data, while the  
1746 module on the right (and *Enterobacteriaceae*) is composed of decreased bacteria in abundance under PPI  
1747 therapy in Amir4's data. (**Bottom panel**) The bottom panel represents the part of Paroni Sterbini's  
1748 network (union of the two PC-corr network), that is not shared with Amir4's one. As in the top and  
1749 middle panels, nodes are coloured according to bacterial phylum level.



# Paroni Sterbini *et al.*



# Parsons *et al.*

1750

1751

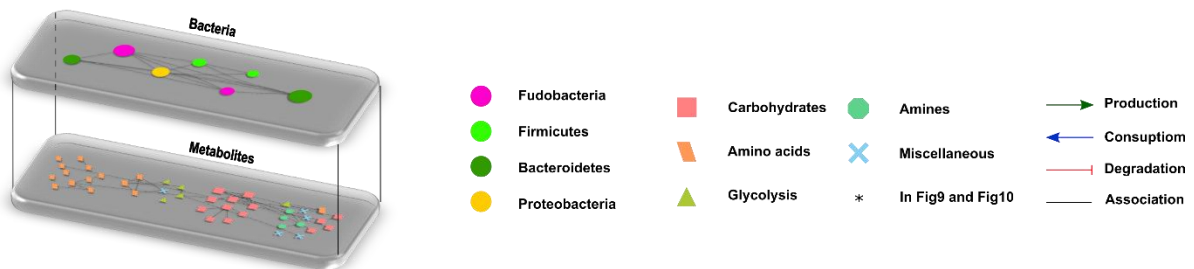
**Figure 7. PC-corr network to investigate the effect of *H. pylori* infection on the gastric mucosal microbiota, coloured according to phylum-level taxonomy. (Middle panel) To investigate the effect**

1752

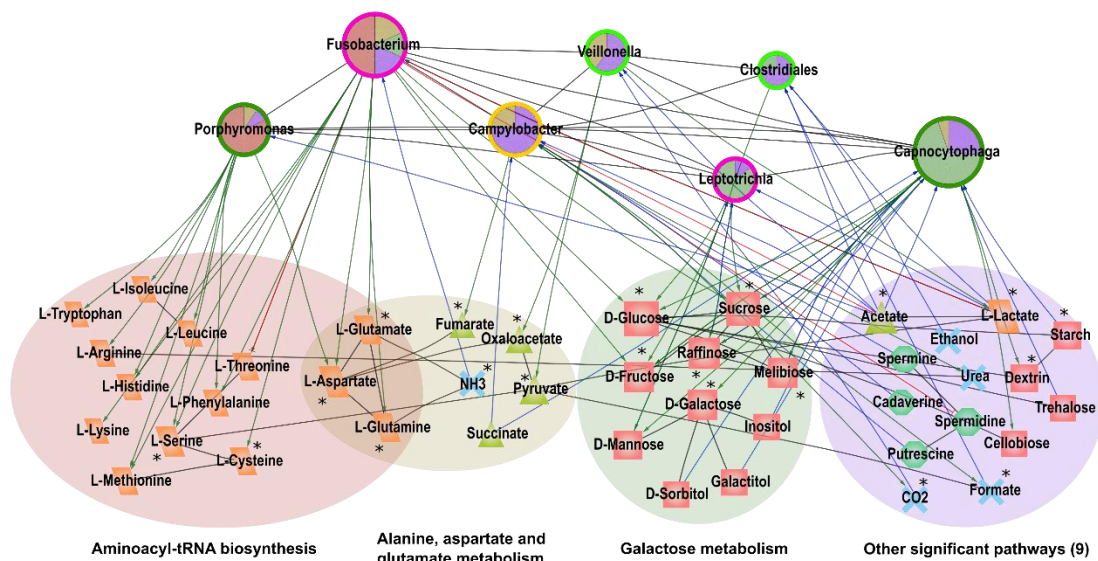
1753 of *H. pylori* infection on the gastric mucosal microbiota, we constructed the conserved PC-corr network  
 1754 at 0.5 cut-off, by intersecting the PC-corr networks obtained from Paroni Sterbini *et al.*<sup>22</sup> and Parsons  
 1755 *et al.*<sup>32</sup> dataset. All the bacteria spotted in the conserved PC-corr network (violet circle) were found  
 1756 decreased in abundance with *H. pylori* infection. The common bacteria that showed an opposite trend  
 1757 in the two datasets, i.e. microbial abundance increased in one dataset and decreased in the other dataset,  
 1758 were removed from the network. **(Top panel)** The top panel show the obtained Paroni Sterbini's  
 1759 network, not in common with the Parsons's network. It contains all bacteria whose abundance is  
 1760 decreased in *H. pylori*-positive patients in Paroni Sterbini *et al.* dataset. **(Bottom panel)** The bottom  
 1761 panel represent the part of Parsons's network that is not shared with Paroni Sterbini's one. As in the top  
 1762 and middle panels, it includes bacterial communities decreased in *H. pylori*-infected patients.

1763

A



B



1764

1765 **Figure 8. PPI-affected bacteria-metabolite network in gastric environment of dyspeptic patients.**

1766 (A) Multilayer (bacteria-metabolite) network representation: the first layer is derived from Fig.6 and  
 1767 represents the consensus network (confirmed in two datasets: gastric mucosa from Paroni Sterbini *et al.*

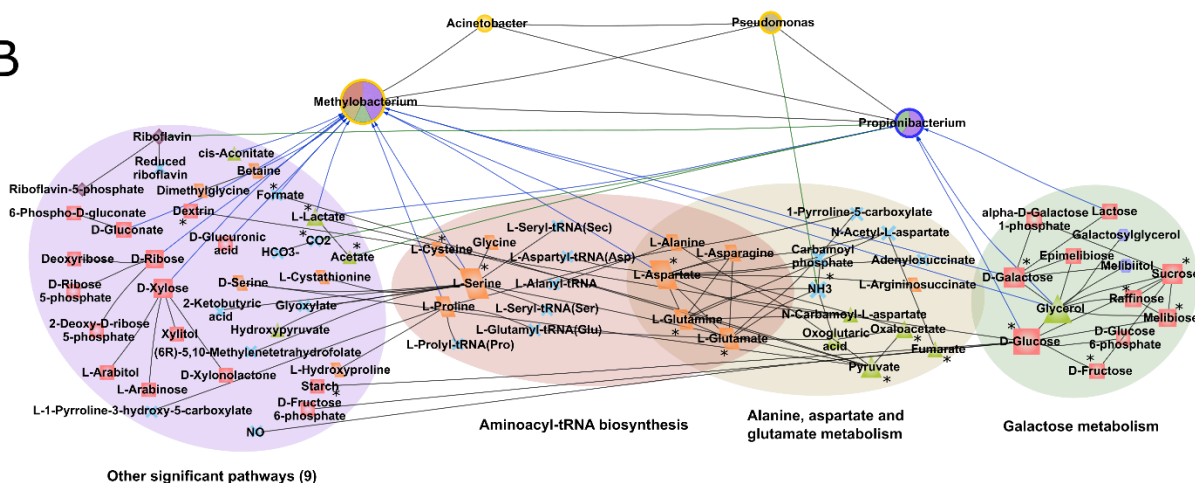
1768 <sup>22</sup> and gastric fluid from Amir et al. <sup>21</sup>) with PPI-affected bacteria nodes that present information on  
 1769 metabolite interaction in <sup>65</sup>. The second layer represents the network whose nodes are the metabolites in  
 1770 <sup>65</sup> interacting with the bacteria network in the first layer; different node shapes and colours refer to  
 1771 different metabolite classes (carbohydrates, amino acids, glycolysis, amines, miscellaneous). **(B)** In  
 1772 depth visualization of the bacteria-metabolite network interactions. The metabolites are grouped  
 1773 according to their involvement in significant pathways. For discernibility, the metabolites are arranged  
 1774 according to three significant pathways ( $p < 0.05$  after Benjamini correction as result of a metabolite  
 1775 pathway enrichment analysis) and a fourth group that encloses altogether nodes associated to other  
 1776 significant pathways (please refer to the method section: Bacteria-metabolite multilayer network  
 1777 construction and metabolite pathway analysis); note that only metabolites present in significant  
 1778 pathways are here displayed. For more information, please refer to figure S17 and table S18. The  
 1779 bacteria node stroke color is associated to the phyla information as in Figure 6, whereas the different  
 1780 colours in the inner fill are associated to the different pathways and their extent is proportional to the  
 1781 number of metabolites that the bacterium connects with in the different displayed pathways.

1782

**A**



**B**



1783



1784 **Figure 9. *H. pylori*-affected bacteria-metabolite network in gastric environment of dyspeptic**  
1785 **patients. (A)** Multilayer (bacteria-metabolite) network representation: the first layer is derived from  
1786 Fig.7 and represents the consensus network (confirmed in two different datasets of gastric mucosa:  
1787 Paroni Sterbini et al. <sup>22</sup> and Parsons et al. <sup>32</sup>) with *H. pylori*-affected bacteria nodes that present  
1788 information on metabolite interaction in <sup>65</sup>. The second layer represents the network whose nodes are  
1789 the metabolites in <sup>65</sup> interacting with the bacteria network in the first layer; different node shapes and  
1790 colours refer to different metabolite classes (carbohydrates, amino acids, glycolysis, lipids, vitamins,  
1791 miscellaneous). **(B)** In depth visualization of the bacteria-metabolite network interactions. The  
1792 metabolites are grouped according to their involvement in significant pathways. For discernibility, the  
1793 metabolites are arranged according to three significant pathways ( $p < 0.05$  after Benjamini correction as  
1794 result of a metabolite pathway enrichment analysis) and a fourth group that encloses altogether nodes  
1795 associated to other significant pathways (please refer to the method section: Bacteria-metabolite  
1796 multilayer network construction and metabolite pathway analysis); note that only metabolites present in  
1797 significant pathways are here displayed. For more information, please refer to figure S18 and table S19.  
1798 The bacteria node stroke color is associated to the phyla information as in Figure 7, whereas the different  
1799 colours in the inner fill are associated to the different pathways and their extent is proportional to the  
1800 number of metabolites that the bacterium connects with in the different displayed pathways.  
1801