

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

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

18 <sup>5</sup>Internal Medicine and Gastroenterology Unit, Università Cattolica del Sacro Cuore, Rome,  
19 Italy;

20 <sup>6</sup>Department of Chemistry and Biotechnology, Division of Gene Technology, Tallinn  
21 University of Technology, Tallinn 12618, Estonia.

22 <sup>7</sup>Biotechnology Center (BIOTEC), Center for Molecular and Cellular Bioengineering (CMCB),  
23 Technische Universität Dresden, Dresden, Germany;

24 <sup>8</sup>Max Planck Institute of Molecular Cell Biology and Genetics, Pfotenhauer Str. 108, 01307  
25 Dresden, Germany

26 <sup>9</sup>Department of Cellular and Molecular Physiology, Institute of Translational Medicine,  
27 University of Liverpool, Liverpool, UK

28 <sup>10</sup>Department of Infrastructure and Environment University of Glasgow, School of  
29 Engineering, Glasgow, UK

30 <sup>11</sup>Department of Gastroenterology, Royal Liverpool and Broadgreen University Hospitals  
31 NHS Trust, Liverpool, UK

32 <sup>12</sup>Brain Bio-Inspired Computing (BBC) Lab, IRCCS Centro Neurolesi “Bonino Pulejo”,  
33 Messina, Italy

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

36

37 \*These authors contributed equally to this work.

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

39

## 40 **Abstract**

41 Although long thought to be a sterile and inhospitable environment, the stomach is inhabited  
42 by diverse microbial communities, co-existing in a dynamic balance. Long-term use of orally  
43 administered drugs such as Proton Pump Inhibitors (PPIs), or bacterial infection such as  
44 *Helicobacter pylori*, cause significant microbial alterations. Yet, studies revealing how the  
45 commensal bacteria re-organize, due to these perturbations of the gastric environment, are in  
46 the early phase. They mainly focus on the most prevalent taxa and rely on linear techniques for  
47 multivariate analysis.

48 Here we disclose the importance of complementing linear dimensionality reduction techniques  
49 such as Principal Component Analysis and Multidimensional Scaling with nonlinear  
50 approaches derived from the physics of complex systems. Then, we show the importance to  
51 complete multivariate pattern analysis with differential network analysis, to reveal mechanisms  
52 of re-organizations which emerge from combinatorial microbial variations induced by a  
53 medical treatment (PPIs) or an infectious state (*H. pylori*).

## 54 **Keywords**

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

58

## 59 Introduction

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

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

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

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

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

102 In fact, by irreversibly blocking H<sup>+</sup>/K<sup>+</sup>-ATPases, PPIs inhibit gastric acid secretion by gastric  
103 parietal cells, which results in a higher intragastric pH, meaning the microenvironment of this  
104 niche changes, hence allowing more bacteria to survive the gastric acid barrier <sup>4,5,16</sup>. The use of  
105 PPIs and higher gastric pH were indeed correlated with the overgrowth of non-*H. pylori*  
106 bacterial flora in the stomach of patients with gastric-reflux and PPIs were shown to aggravate  
107 gastritis because of co-infection with *H. pylori* and non-*H. pylori* bacterial species <sup>4,14,19,20</sup>.

108 However, PPIs may also affect the gastrointestinal microbiome through pH-independent  
109 mechanisms, by directly targeting the proton pumps of naturally occurring bacteria by binding  
110 P-type ATPases (e.g. *H. pylori*) <sup>4,6</sup>.

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

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

152

## 153 **Methods**

### 154 ***Dataset description***

#### 155 *Amir3 (esophageal mucosa)*

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

158 The dataset was obtained from 16 esophageal mucosal biopsies of eight individuals before and  
159 after eight weeks of PPI treatment. Two patients with heartburn presented normal  
160 oesophagogastroduodenoscopy (H) indicating that they present healthy oesophageal tissues but  
161 are exposed to gastric refluxate, four patients had oesophagitis (ES) and two had Barrett's

162 oesophagus (BE). Metagenomes were obtained by pyrosequencing 16S rRNA amplicons on the  
163 GS FLX system (Roche). Data were processed by replicating the bioinformatics workflow  
164 followed by Amir and colleagues <sup>21</sup> in order to obtain the matrix of the bacterial absolute  
165 abundance: sequence reads were analysed with the pipeline Quantitative Insights into Microbial  
166 Ecology (QIIME) v. 1.6.0 <sup>26</sup> using default parameters (sequences were removed if shorter than  
167 200 nt, if they contained ambiguous bases or uncorrectable barcodes, or if the primer was  
168 missing). Operational Taxonomic Units (OTUs), that are clusters of sequences showing a  
169 pairwise similarity no lesser than 97%, were identified using the UCLUST algorithm  
170 (<http://www.drive5.com/usearch/>). The most abundant sequence in each cluster was chosen as  
171 the representative of its OTU, and this representative set of sequences was then used for  
172 taxonomy assignment by means of the Bayesian Ribosomal Database Project classifier <sup>27</sup> and  
173 aligned with PyNAST103. Chimeras, that are PCR artefacts, were identified using  
174 ChimeraSlayer <sup>28</sup> and removed. The Greengenes database, which was used for the annotation  
175 of the reads, additionally identifies groups of bacteria that are supported by whole genome  
176 phylogeny, but are not yet officially recognized by the Bergeys taxonomy, which is the  
177 reference taxonomy and is based on physiochemical and morphological traits. This results in a  
178 special annotation for some taxa, like *Prevotella*, that thus appears both with the general  
179 annotation, that is *Prevotella*, and with the special annotation, that is between square brackets,  
180 [*Prevotella*].

181

#### 182 *Amir4 (gastric fluid)*

183 The dataset was generated by Amir and colleagues <sup>21</sup>, and is public and available in the MG  
184 RAST database (<http://metagenomics.anl.gov/linkin.cgi?project=5732>). It comprises eight  
185 patients, whose gastric fluid was sampled at two different time points, that is before PPI  
186 treatment and after eight weeks of PPI treatment, for a total of 16 samples. The patients are the  
187 same described in Amir3. Metagenomes were obtained by pyrosequencing fragments of the

188 16S rRNA gene on the GS FLX system (Roche). Then the data were processed by replicating  
189 the same bioinformatics workflow followed by Amir and colleagues <sup>21</sup> that was described in  
190 the previous data description (Amir3), in order to obtain the matrix of the bacterial absolute  
191 abundance. As for Amir3, the Greengenes database was used for the annotation of the reads.

192

193 *Paroni Sterbini (gastric mucosa)*

194 The dataset was generated by Paroni Sterbini and colleagues <sup>22</sup>, and is public and available in  
195 the NCBI Sequence Read Archive (SRA) (<http://www.ncbi.nlm.nih.gov/sra>, accession number  
196 SRP060417), where all details pertaining the sequencing experimental design are also reported.

197 It contains 24 biopsy specimens of the gastric antrum from 24 individuals who were referred to  
198 the Department of Gastroenterology of Gemelli Hospital (Rome) with dyspepsia symptoms (i.e.  
199 heartburn, nausea, epigastric pain and discomfort, bloating, and regurgitation). Twelve of these  
200 individuals (PPI1 to PPI12) had been taking PPIs for at least 12 months, while the others (S1  
201 to S12) were not being treated (naïve) or had stopped treatment at least 12 months before sample  
202 collection. In addition, 9 (5 treated and 4 untreated) were positive for *H. pylori* infection, where  
203 *H. pylori* positivity or negativity was determined by histology and rapid urease tests.

204 Metagenomes were obtained by pyrosequencing fragments of the 16S rRNA gene on the GS  
205 Junior platform (454 Life Sciences, Roche Diagnostics). Then the sequence data were  
206 processed by replicating the bioinformatics workflow followed by Paroni Sterbini *et al.* <sup>22</sup>, in  
207 order to obtain the matrix of the bacterial absolute abundance.

208

209 *Parsons (gastric mucosa)*

210 The dataset was generated by Parsons and colleagues <sup>29</sup>, and is public and available in the EBI  
211 short-read archive (the European Nucleotide Archive, ENA) (<https://www.ebi.ac.uk/ena>,  
212 accession number PRJEB21104). In the original study, the authors focused on the analysis of  
213 gastric biopsy samples of 95 individuals (in groups representing normal stomach, PPI treated,



214 *H. pylori*-induced gastritis, *H. pylori*-induced atrophic gastritis and autoimmune atrophic  
215 gastritis), selected from a larger prospectively recruited cohort patients who underwent  
216 diagnostic upper gastrointestinal endoscopy at Royal Liverpool University Hospital<sup>29</sup>. RNA  
217 extracted from gastric corpus biopsies was analysed using 16S rRNA sequencing (MiSeq).  
218 Then the sequence analysis was performed, as described by the authors in the supplementary  
219 methods of the original article <sup>29</sup>. Here we focused on the analysis of gastric biopsy specimens  
220 (in total 42 samples) from normal stomach group (20 patients) and belonging to the *H. pylori*  
221 gastritis group (22 patients). As described in <sup>29</sup>, patients in the normal stomach group showed  
222 normal endoscopy, no evidence of *H. pylori* infection by histology, rapid urease test or serology,  
223 were not treated by PPI and were normogastrinaemic. Patients in the *H. pylori* gastritis group  
224 were instead positive to *H. pylori* infection by urease test, histology and serology, were not  
225 taking PPI medication and were normogastrinaemic.

226

### 227 ***Data exploration and visualization: the reason for unsupervised dimension*** 228 ***reduction***

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

238 A final interesting advantage offered by unsupervised analysis is in small size datasets, where  
239 the number of samples  $n$  is significantly lower than the number of features  $m$ , a condition that

240 unfortunately occurs in several metagenomic studies. When  $n \ll m$  the application of  
241 supervised approaches can become problematic, because the supervised procedure of parameter  
242 learning can suffer from overfitting<sup>23,30,31</sup>.  
243 The mainstream multivariate methods to unsupervisedly explore data patterns in metagenomic  
244 studies are based on linear dimension reduction, in particular PCA<sup>32,33</sup> and MDS<sup>34,35</sup>, also  
245 known as PCoA, methods that have been used to explore and visualize data structure in many  
246 metagenomic studies, from sponge<sup>36,37</sup> to gastric tissue microbiota<sup>22</sup>. These tools perform a  
247 dimension reduction of the data either by *multidimensional variance analysis* (for instance  
248 PCA) or *dissimilarity embedding* (for instance MDS/PCoA). PCA collects uncorrelated  
249 variance in the multidimensional space, creating new synthetic orthogonal variables, which are  
250 linear combinations of the original ones, then plots the samples in a reduced space using the  
251 new variables that embody the largest orthogonal variances. MDS computes dissimilarities  
252 between every pair of samples, plotting the Euclidean part of these dissimilarities as distances  
253 between every pair of points (MDS) in a reduced space, in this way the linear part of the sample  
254 relations can be represented.

255

### 256 ***The Tripartite-Swiss-Roll dataset***

257 In order to test and visualize how the algorithms could detect nonlinearity, we performed the  
258 analyses on the Tripartite-Swiss-Roll dataset: an artificial dataset characterized by nonlinear  
259 structures and generated as discretization of the manifold associated to a Swiss-Roll function<sup>38</sup>  
260 in a three-dimensional (3D) space. Indeed, it is a synthetic dataset obtained as the partition in  
261 three sections of a discrete Swiss-Roll manifold depicted in a three-dimensional space<sup>38</sup>. It  
262 reproduces the typical nonlinearity (given by the Swiss-Roll shape) and the discontinuity (given  
263 by the tripartition of the manifold), that we do not see and that are often hidden in the  
264 multidimensional representation of our samples. See the illustration in the original 3D-space of  
265 the Tripartite-Swiss-Roll dataset in Fig. 1A. This dataset is useful to introduce readers, not

266 expert with nonlinear data analysis, to the basic concepts of nonlinear dimension reduction and  
267 therefore to facilitate their understanding of the new proposed methodologies for nonlinear  
268 dimension reduction.

269

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

271 Below, we report some of the PCA major advantages and drawbacks, that were pinpointed in a  
272 recent study on multidimensional population genomics <sup>39</sup>, and of other conventional  
273 dimensional reduction techniques employed for the analysis of metagenomic data.

274 PCA is time-efficient, parameter-free and straightforward to interpret, yet it strives to resolve  
275 structure in datasets with few samples and highly numerous features, which enclose nonlinear  
276 patterns. Therefore, PCA can occasionally fail to reveal differences among samples, even when  
277 differences are known a-priori, which means it can also miss represent hidden nonlinear  
278 relations among the samples in the feature space. For instance, see the illustration of the PCA  
279 two-dimension reduction mapping of the Tripartite-Swiss-Roll dataset in Fig. 1B. PCA clearly  
280 fails to unfold and reveal the structure of the three separated groups of samples.

281 MDS, on the other hand, preserves the sample distances in a 2D-space based on the calculation  
282 of a distance matrix (Fig. 1C,D). In ecology, distance (or dissimilarity) matrices are a major  
283 way to transpose the ecological information of samples in terms of their species composition  
284 and abundance <sup>40,41</sup>. In this article we will consider classical MDS (which uses Euclidean  
285 distance and is in practice equivalent to PCA <sup>42,43</sup>), and non-metric MDS (NMDS) obtained  
286 according to Sammon's Mapping <sup>44</sup>. In the latter, the elements of the multivariate space are  
287 mapped onto a lower dimensional space while retaining the original inter-point dissimilarities,  
288 by means of a nonlinear, but monotonic transformation (Sammon Mapping). Since it respects  
289 the ranking of dissimilarities, it tends to linearize the relationships between the samples. In  
290 addition, MDS will be performed also according to Bray-Curtis (MDSbc) dissimilarity and  
291 weighted UniFrac (MDSwUF) distance because they are considered the reference in

292 metagenomics studies. Bray-Curtis dissimilarity quantifies how dissimilar two sites (samples)  
293 are based on counts (bacterial abundances), where 0 means two samples are identical and 1  
294 means that the two samples do not share any taxa<sup>45,46</sup>. Dissimilarly, the UniFrac distance, either  
295 unweighted (qualitative) or weighted (quantitative), is the most popular phylogenetic distance  
296 measure for the microbial community diversity between different samples (also known as  $\beta$ -  
297 diversity<sup>47</sup>) and, differently from the previous discussed methods, uses the phylogenetic  
298 information (which is an external knowledge not contained in the dataset) on the taxa to  
299 compare samples. In particular, its weighted-version weights the branches of a phylogenetic  
300 tree based of the taxa abundance information<sup>48-51</sup>. Hence the weighted UniFrac distance  
301 directly accounts for differences in the abundance of different kinds of bacteria, and can be  
302 crucial to describe community changes<sup>49</sup> in the studied samples.

303 We need to specify that both MDSwUF and NMDS are in practice nonlinear methods and  
304 weighted UniFrac is not a classical unsupervised technique like the others. In fact, MDSwUF  
305 adopts a distance that combines the information given by the bacterial abundance of the dataset  
306 with the supervised prior (external) knowledge regarding the known hierarchical phylogenetic  
307 relationship among the bacteria. However, like PCA, MDS can fail to detect patterns if data are  
308 not properly linearized<sup>52</sup>. For instance, see Fig. 1C-D where MDSbc and NMDS respectively  
309 fail to resolve the Tripartite-Swiss-Roll dataset. When we consider clinical metagenomic data,  
310 this failure potentially reduces the chances of correctly pinpointing samples which may  
311 represent clinical subspecies, and thus remain undetected and undiagnosed. In brief, these  
312 methods are not efficient to perform *hierarchical embedding* directly from the abundance value,  
313 since hierarchies preserve tree-like structures, and tree-like structures follow a hyperbolic, thus  
314 nonlinear, geometry<sup>53-55</sup>. Only MDSwUF is able to account for nonlinear hierarchical  
315 organization, yet this is not directly inferred from the abundance values, but rather forced as a  
316 constraint of prior supervised knowledge on the phylogeny of bacteria. For this reason we  
317 cannot offer a test on the Tripartite-Swiss-Roll dataset.

318 In our analysis of the Paroni Sterbini dataset, we also showed the results of a supervised  
319 technique, Linear Discriminant Analysis (LDA), which uses the labels to perform dimension  
320 reduction. LDA aims to separate the samples into groups based on hyperplanes and describe  
321 the differences between groups by a linear classification criterion that identifies decision  
322 boundaries between groups <sup>34</sup>. This technique is not congruous (and sometimes statistically  
323 invalid) for small sample size datasets. The reason is that given the reduced sample size we  
324 cannot divide the dataset in a training and test set, which is a fundamental requirement of  
325 supervised methods such as LDA.

326

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

328 In 2010, Cannistraci *et al.* <sup>23</sup> introduced the centred version of Minimum Curvilinear  
329 Embedding (MCE), which provided notable results in: i) visualisation and discrimination of  
330 pain patients in peripheral neuropathy, and the germ-layer characterisation of human organ  
331 tissues <sup>23</sup>; ii) discrimination of microbiota in sea sponges <sup>36</sup>; iii) embedding of networks in the  
332 hyperbolic space <sup>54</sup>; iv) stage identification of embryonic stem cell differentiation based on  
333 genome-wide expression data <sup>56</sup>. In this fourth example, MCE performance ranked first on 12  
334 different tested approaches (evaluated on 10 diverse datasets). More recently in 2013 <sup>30</sup>, the  
335 non-centred version of the algorithm, named ncMCE, has been used: i) to visualise clusters of  
336 ultra-conserved regions of DNA across eukaryotic species <sup>57</sup>; ii) as a network embedding  
337 technique for predicting links in protein interaction networks <sup>30</sup>, outperforming several other  
338 link prediction techniques; iii) to unsupervisedly reveal hidden patterns related with gender  
339 difference and metabolic-disease risk-factors in lipidomic profiles extracted from human  
340 plasma samples <sup>58</sup>; iv) to unsupervisedly infer and visualize phylogenetic (hierarchical)  
341 relations directly from individual SNP profiles in human population genetics <sup>39</sup>. Finally, also  
342 applications in non-biological problems such as the unsupervised discrimination of bad from

343 good radar signals<sup>30</sup>, represented a proof of concept of the universality of MCE for addressing  
344 nonlinear investigation of data and signals in general. Also in the case of the metagenomics  
345 studies targeting sea sponges,<sup>36,37</sup> both MCE and its non-centred variant<sup>23,30</sup> once again proved  
346 successful in detecting structure where PCA and MDS could not, or hardly find any. This is  
347 mainly because MCE/ncMCE are unsupervised and parameter-free topological machine  
348 learning for *nonlinear* dimensionality reduction and multivariate analysis, that are able to  
349 perform a *hierarchical embedding*.

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

354 Minimum Curvilinearity (MC)<sup>23</sup>, the principle behind MCE and ncMCE, was invented with  
355 the aim to reveal nonlinear data structures also, and especially, in the case of datasets with few  
356 samples and many features. MC principle suggests that curvilinear (nonlinear) distances  
357 between samples may be estimated as pairwise distances over their Minimum Spanning Tree  
358 (MST), constructed according to a selected distance (Euclidean, correlation-based, etc.) in a  
359 multidimensional feature space (here the metagenomic data space). In this study, we considered  
360 Pearson-correlation based distance (refer to<sup>23</sup> for details on the way to compute the distance  
361 for the MST). The collection of all nonlinear pairwise distances forms a distance matrix called  
362 the MC-distance matrix or MC-kernel, which can be used as an input in algorithms for  
363 dimensionality reduction, clustering, classification and generally in any type of machine  
364 learning. In MCE and ncMCE, the MC-kernel (which is non-centred for ncMCE) is followed  
365 by dimensionality reduction using singular value decomposition (SVD), and then by the  
366 projection of the samples onto a two-dimensional space for visualisation and analysis. Thus,  
367 MCE/ncMCE is a form of nonlinear and parameter-free kernel PCA<sup>30</sup>. In the rest of the article  
368 we will simply use the name MCE to indicate both MCE and ncMCE, since the centring

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

371

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

373 A previous study by Alanis-Lobato *et al.* <sup>39</sup> showed that MCE is automatically able to  
374 unsupervisedly infer and visualize phylogenetic (hierarchical) relations directly from individual  
375 SNP profiles in human population genetics. Precisely, ncMCE detected separation between  
376 ethnic groups and provided an ordering over the discriminating dimension that was related to  
377 the phylogenetic organization of these populations.

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

389

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

391 The performance of the mentioned dimension reduction algorithms is evaluated as the ability  
392 to separate the samples in the first two dimensions of embedding since, as discussed above, this  
393 is one of the preferred unsupervised strategies to investigate the presence of patterns in  
394 multidimensional datasets. In order to quantitatively evaluate the performance, we use a

395 recently proposed index for sample separation<sup>59</sup>. This index can be defined for any separation-  
396 measure and in this study we considered three well-known measures: p-value of Mann-Whitney  
397 U test, Area Under the ROC-Curve (AUC) and Area Under the Precision-Recall curve (AUPR),  
398 that are regularly used to quantitatively measure the performance of a binary predictor.  
399 More precisely, in the 2D space a line is drawn between the centroids of the two groups that are  
400 compared, subsequently all the points are projected on this line and then a p-value, AUC and  
401 AUPR are computed for the projected points. This new index is named *projection-based*  
402 *separability index* (PSI) and can actually be applied not only in a 2D space, but in any N  
403 dimensional space. For the calculation of the centroids we consider the 2D-median of each  
404 cluster/class's group. In case more than two groups are present in a dataset, all the p-values,  
405 AUC and AUPR between the possible pair-groups are computed, and the average values of all  
406 the pairwise p-values, AUC and AUPR are chosen as an overall estimator of separation between  
407 the groups in the 2D reduced space. This case applies only to the Paroni Sterbini dataset, which  
408 is composed of three or, possibly, four groups of samples. All the other datasets are instead  
409 composed of two groups.  
410 It is important to note that the PSI was also applied to the data in the original high-dimensional  
411 (HD) space, as a reference to see how good the unsupervised dimension reduction approaches  
412 are in preserving the original group separability of the HD space.  
413 All the algorithms were tested considering (when allowed by the dimension reduction method)  
414 data centring or non-centring. In addition, multiple normalization options were investigated and  
415 the datasets were considered under a certain type of normalization: division by the column -  
416 which reports the OTU - sum (indicated by DCS); division by the row - which reports the  
417 sample - sum (indicated by DRS); function  $\log_{10}(1+x)$  applied to the dataset (indicated by  
418 LOG).  
419



## 420 *From Markov Clustering (MCL) to Minimum Curvilinear Markov Clustering* 421 *(MC-MCL)*

422 MCL is an unsupervised algorithm for the clustering of weighted graphs based on simulations  
423 of (stochastic) flow in graphs <sup>60</sup> (<http://micans.org/mcl/>). By varying a single parameter called  
424 inflation (with values between 1.1 and 10), clustering patterns on different scales of granularity  
425 can be detected. For clustering samples of a multidimensional dataset, the workflow starts with  
426 the computation of correlations (generally Pearson correlations) between the samples, and  
427 creates an edge between each pair of samples, where the edge-weight assumes the value of the  
428 respective pairwise positive sample correlation, or values zeros in case of negative correlations.  
429 This generates a weighted correlation graph (network), which is used as a map to simulate  
430 stochastic flows and detect the structural organization of clusters in the graph.

431 With the purpose of creating and testing a nonlinear variant of the MCL algorithm, we adopt  
432 an innovative algorithm which was recently proposed and called MC-MCL <sup>61</sup>. The idea is the  
433 following. The MC-kernel – discussed above in the MCE section - is a nonlinear distance matrix  
434 (or kernel) that expresses the pairwise relations between samples as a value of distance: small  
435 samples distance indicates sample similarity, while large samples distance indicates sample  
436 dissimilarity. Here we reverse (using the following function:  $f(x) = 1 - x$ ) and after this we  
437 put to zero the negative values of the *MC-distance* kernel to get a *MC-similarity* kernel, where  
438 small values (close to zero) indicate low sample similarity and large values (close to one)  
439 indicate high sample similarity. A technical detail: for the computation of the MC-distance  
440 kernel, it is necessary to firstly square root the original distances (correlation-based) between  
441 the samples. As already investigated in <sup>23</sup>, this attenuates the estimation of large distances and  
442 amplifies the estimation of short distances; consequently it helps to regularize the nonlinear  
443 distances inferred over the MST in order to subsequently use them for message passing <sup>23</sup> (such  
444 as affinity propagation) or flow simulation (such as MCL) clustering algorithms.

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

453

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

455 The clustering algorithms MCL and MC-MCL were applied to the datasets, either raw, or after  
456 the same normalization procedures used before dimensionality reduction (DCS: division by  
457 column (OTU) sum; DRS: division by row (sample) sum; LOG: function  $\log_{10}(1+x)$  applied  
458 to the dataset) and their performance was evaluated by means of accuracy. The accuracy is  
459 computed as the ratio of the number of samples assigned to the correct clusters over the total  
460 number of samples. For both MCL and MC-MCL, we tested Pearson and Spearman correlations  
461 to build the similarity measure to feed into the clustering methods. The Spearman correlation  
462 can also detect a subclass of nonlinear associations (which have monotonic shape function) or  
463 correct for outliers. Differently from what suggested for large gene datasets with thousands of  
464 samples in <sup>60</sup> (<http://micans.org/mcl/>), in this study we had to consider the whole set of original  
465 positive correlations without applying any threshold (cut-off) to the values. This was  
466 compulsory, since we considered datasets with few samples. In our case, to keep the graph  
467 connected, with one unique connected component, we could not introduce any kind of threshold  
468 that would otherwise alter the real graph connectivity (dividing the graph in disconnected  
469 components) and hence the clustering result. Since the MCL algorithm needs a single input

470 parameter (inflation) to control the granularity of the output clustering, we ran it for different  
471 inflation values until we achieved the desired number of clusters. Finally, in the Paroni Sterbini  
472 *et al.* dataset<sup>22</sup> it was not clear in advance whether the correct number of clusters present in the  
473 multidimensional space was three or four. Hence, we tested the clustering algorithms  
474 considering as output both three and four clusters' configurations, and we identified as the best  
475 solution the one that offered the highest accuracy.

476

### 477 ***PC-corr network***

478 Furthermore, we investigated the effect of PPI on the microbiota of gastric fluid and gastric  
479 mucosa in dyspeptic patients, and the changes induced by *H. pylori* infection on the gastric  
480 mucosal microbiota, by means of the PC-corr approach<sup>62</sup>. PC-corr represents a simple  
481 algorithm that associates to any PCA segregation a discriminative network of features'  
482 interactions<sup>62</sup>. It is a method for linear multivariate-discriminative correlation network reverse  
483 engineering, that, thanks to its multivariate nature, can help to stress and squeeze out the  
484 underlying combinatorial and multifactorial mechanisms that generate the differences between  
485 the studied conditions<sup>62</sup>. Hence, for the studied datasets, it can be employed to point out the  
486 possible presence of bacterial alterations and their interplay, induced by a medical treatment  
487 (PPIs in dyspepsia) or infectious state (*H. pylori*).

488

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

490 Dimensionality reduction was performed in MATLAB on the abundance matrix of genus-level  
491 taxonomic assignments, with samples in rows and taxonomic assignments (OTUs) in columns.  
492 For MDSwUF, the computation of the weighted UniFrac distance was performed in R. We used  
493 the following MATLAB functions to calculate PCA, MDS and NMDS (Sammon Mapping)  
494 respectively: *svd*, *cmdscale* and *mdscale*. For the calculation of Bray-Curtis dissimilarity, we  
495 used the function MATLAB *f\_braycurtis* in the Fathom Toolbox<sup>63</sup>

496 (<http://www.marine.usf.edu/user/djones/matlab/matlab.html>). Instead, for the calculation of the  
497 weighted Unifrac distance for all sample pairs, we used the R function *UniFrac* in the phyloseq  
498 package (<https://bioconductor.org/packages/release/bioc/html/phyloseq.html>), after creating a  
499 phyloseq-class object (with R function *phyloseq* in the same package) that contains both the  
500 abundance table (OTU table) and the phylogenetic tree. The MATLAB code for MCE/ncMCE  
501 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)  
502 [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  
503 MCL-edge software (<http://micans.org/mcl/>) in a Windows environment, following the  
504 procedure suggested by the authors in the software website. To apply this algorithm, we created  
505 a MATLAB function that generates automatically the input for MCL (equivalent to the  
506 *mcxarray* function in the software) and then uses a system call to run MCL in a UNIX-like  
507 environment (Cygwin, <https://www.cygwin.com/>). PC-corr method was performed in  
508 MATLAB on the abundance matrix of the genus-level taxonomic assignments, with samples in  
509 rows and taxonomic assignments in columns. The PC-corr algorithm is available as MATLAB  
510 function (as well as R function) at: [https://github.com/biomedical-cybernetics/PC-corr\\_net](https://github.com/biomedical-cybernetics/PC-corr_net).  
511 Then the obtained PC-corr networks were displayed by Cytoscape (<http://www.cytoscape.org/>).

512

## 513 **Results**

514 To answer the five questions stated in the Background section, we analysed the abovementioned  
515 16S rRNA gene sequencing datasets with information on PPI consumption in dyspeptic  
516 patients, following the workflow shown in Fig. 2. It is important to underline that, in one of the  
517 three initially analysed datasets (in Paroni Sterbini *et al.*<sup>22</sup>), we have the additional information  
518 on positivity or negativity to *H. pylori* infection. A fourth dataset (Parsons *et al.*<sup>29</sup>) is used only  
519 for the validation of the PC-corr network results and it contains not only information on PPI  
520 consumption but also additional information on positivity or negativity to *H. pylori* infection.

521 Unsupervised approaches were chosen for dimension reduction, and clustering because  
522 supervised (constrained) methods have been shown to perform poorly on small datasets, as  
523 explained in the paper by Smialowski *et al.*<sup>31</sup> and the work by Zagar and colleagues<sup>56</sup>.  
524 Firstly, we performed unsupervised dimension reduction, both linear and nonlinear (described  
525 in the ‘*Methods- PCA, MDS (or PCoA) and LDA*’ and ‘*Methods- Minimum Curvilinear*  
526 *Embedding (MCE)*’) and we focused on the first two dimensions of embedding as it is common  
527 practice<sup>25</sup>. As we will show, linear techniques will fail to bring out the patterns in the microbial  
528 datasets, related to PPI-treatment. Instead, nonlinear dimension reduction will reveal the  
529 presence of hidden patterns related to PPI treatment. In particular, in the gastric biopsies dataset  
530 (Paroni Sterbini *et al.*<sup>22</sup>), nonlinear dimension reduction will point out the evidence of PPI  
531 perturbation. Secondly, clustering algorithms were applied to the studied datasets to confirm  
532 that the hidden patterns detected by nonlinear dimension reduction are well posed. Finally, the  
533 PC-corr algorithm<sup>62</sup> is used to find the bacteria community (features) that make the difference  
534 between the patterns or groups, allowing our understanding of the PPI-induced and *H. pylori*-  
535 induced microbial perturbations.

536

### 537 **Gastric tissue dataset unsupervised analysis**

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

541 In our first analysis, we focused on the Paroni Sterbini *et al.* dataset<sup>22</sup> and, to facilitate the  
542 visualization of the sample separations in the 2D reduced space, we assigned: red colour to  
543 untreated dyspeptic patients without *H. pylori* infection (HP-); green colour to untreated  
544 dyspeptic patients with *H. pylori* infection (HP+); and blue colour to patients treated with PPI  
545 regardless of their *H. pylori* infection (PPI). However, to help to detect also the effect of the *H.*  
546 *pylori* infection we reported the labels close to each sample, with a ‘+’ indicating the infection

547 (PPI+) or a ‘-’ indicating the absence of infection (PPI-). Finally, we also tested whether this  
548 separation into three main groups (HP-, HP+, PPI) is more truthful, from the metagenomics  
549 data standpoint, than the one in four groups (HP-, HP+, PPI-, PPI+).

550 Figure 3 shows the results of the multivariate techniques widely employed in metagenomic  
551 studies, PCA (Fig. 3A), MDSbc (Fig. 3B) and MDSwUF (Fig. 3C), and NMDS (with Sammon  
552 Mapping) (Fig. 3D) (for more detail see the corresponding method section; the plots represents  
553 the best results based on average p-value in Supplementary Table S1), which could only  
554 differentiate the group of untreated *H. pylori* positive samples (green dots) with respect to the  
555 group of untreated *H. pylori* negative samples (red dots), and no further separation is  
556 significantly detectable. Considering the PSI results, the p-values are significant (p-value<0.05,  
557 Table 1 and Fig. 3) (evaluated in the 2D embedding space, for details see ‘*Procedure to evaluate*  
558 *the performance of the dimension reduction algorithms*’). PCA and NMDS exhibit the lowest  
559 p-value (0.0090), while MDSwUF and MDSbc displays p-values higher than 0.01 (respectively  
560 0.011 and 0.021). This trend is also confirmed by their AUC and AUPR values, with highest  
561 values for PCA (AUC=0.924, AUPR=0.960) and NMDS (AUC=0.924, AUPR=0.954). Indeed,  
562 in all the plots there is a visible trend of separation between PPI-treated (blue dots) and untreated  
563 (red and green dots) samples, but this is not sufficient to declare the presence of the complete  
564 separation, and a manifest ‘crowding problem’<sup>30</sup> mixes the two cohorts together. According to  
565 this output, the dataset appears to be strongly influenced by the presence of *H. pylori*, which is  
566 the predominant taxon (abundance > 50%, Supplementary Table S2, percent abundance sheet)  
567 in four of the untreated *H. pylori* positive patients: where *H. pylori* is predominant, sample  
568 groups are quite close to one another and far from all the other samples in all four multivariate  
569 analyses (Fig. 3). Thus, PCA and MDS mainly show us that these metagenomes separate  
570 according to *H. pylori* abundance, and there is no treatment-related pattern.

571 Non-centred MCE (Figure 4A, DCS normalization) was the best performing technique, with a  
572 p-value of 0.004, AUC of 0.967 and AUPR of 0.987 (Table 1) (for details see Supplementary

573 Table S1). It even outperforms the nonlinear methods NMDS (Sammon Mapping) and  
574 MDSwUF, since it is automatically able to infer the (hierarchical) phylogenetic relationship  
575 among the bacteria directly from the bacterial abundance of the dataset by performing a  
576 hierarchical embedding, as already shown in the study of Alanis-Lobato *et al.*<sup>39</sup> (see '*Methods-*  
577 *MCE to unsupervisedly infer and visualize phylogenetic (hierarchical) relations*').  
578 Furthermore, the MCE performance does not depend on its centring/non-centring, in fact the  
579 centred MCE version resolves the nonlinearity in the data too. Whereas, PCA regardless of  
580 being centred or non-centred does not resolve the nonlinearity in the data.

581 While MDS and PCA are confounded by the mixture of factors characterizing the samples and  
582 do not manage to resolve the differences between treated and untreated samples, non-centred  
583 MCE is the only technique that visibly separates samples by ordering them along the second  
584 dimension into three groups, detecting a treatment-related structure in the data (Fig. 4A). This  
585 is plausible, because in any non-centred embedding the first dimension points towards the  
586 centre of the manifold<sup>30</sup>, while the second dimension in the case of non-centred MCE represents  
587 the direction of higher topological nonlinear extension of the manifold. Interestingly, untreated  
588 *H. pylori* negative samples (red dots, HP-) gather in the upper tail of the samples' distribution,  
589 while treated samples (blue dots, PPI), both *H. pylori* test positive (PPI+) and negative (PPI-),  
590 are mixed and show no other internal discernible groups. Untreated *H. pylori* positive samples  
591 (green samples, HP+) gather at the bottom of the plot (Fig. 4A). Unlike the other approaches,  
592 non-centred MCE detects a treatment-related structure in the data and separates patients into  
593 three, not four, groups: PPI-treated, untreated *H. pylori* negative and untreated *H. pylori*  
594 positive. This last group appears as a subgroup marginally discriminating from the PPI-treated  
595 group and the topology of the samples seems to suggest that PPI treatment modifies the gastric  
596 microbiota of *H. pylori*-negative patients with dyspeptic symptoms and gastric mucosa  
597 inflammation, shifting their gastric ecosystem in the same direction of PPI-treated *H. pylori*-  
598 positive patients. We speculate that the fact that PPI treatment and *H. pylori* infection determine

599 the samples to gather in a similar position (i.e. out of the PPI-untreated/HP-negative group) in  
600 the non-centred MCE reduced space, indicates that both the PPI drugs and *H. pylori* induce an  
601 ecological change in the stomach, which might be driven by similar mechanisms. As a matter  
602 of fact, *H. pylori* can colonize the acidic lumen of the stomach thanks to its ability to hydrolyse  
603 urea into carbon dioxide (CO<sub>2</sub>) and ammonia (NH<sub>3</sub>)<sup>64</sup>, thus increasing the intragastric pH. On  
604 the other hand, PPIs obtain the same result through the inhibition of acid secretion in gastric  
605 parietal cells, which blocks H<sup>+</sup>/K<sup>+</sup>-ATPases. Both processes are therefore shifting the gastric  
606 environment towards an alkaline condition. Thus, MCE provides an ordering of the groups  
607 along the second dimension that is related to pH increment (from HP- to PPI+).

608 Similarly to the Paroni Sterbini *et al.* microbial dataset, the Tripartite-Swiss-roll dataset (that is  
609 a synthetic dataset containing nonlinear structures obtained by tri-partitioning a discrete Swiss-  
610 Roll manifold<sup>38</sup> in a three-dimensional space, for more details see the method section: The  
611 Tripartite-Swiss-Roll dataset'), presents a hierarchical-organized nonlinearity (Fig. 1A). And  
612 also in this case, similarly to the result of the Paroni Sterbini *et al.* analysis, non-centred MCE  
613 is able to perform a hierarchical embedding that orders the hidden subgroups of the dataset  
614 along the second dimension of embedding (Fig. 4B). On the contrary - as already commented  
615 in the method section - PCA, MDSbc and NMDS (Fig. 1B-D) were unable to resolve the  
616 nonlinearity of the Tripartite-Swiss-Roll: its three partitions are either superimposed (Fig. 1B,  
617 D) or twisted in a horseshoe shape (Fig. 1C). Indeed, the Tripartite-Swiss-Roll is purposely  
618 created to reproduce a manifold that is nonlinear and discontinuous (broken in three parts) such  
619 as the results of MCE analysis of Paroni Sterbini *et al.* seems to be.

620 For the Paroni Sterbini dataset, we also performed a supervised linear approach for dimension  
621 reduction, LDA (Supplementary Figure S1), yet the cross-validation test showed that this  
622 constrained technique could re-assign samples to their groups with 54% of error (ldaCVer in  
623 Supplementary Table S3), confirming its statistical invalidity for the small size dataset problem.



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

634 While MC-MCL represents the minimum curvilinear version of MCL, MCE is the minimum  
635 curvilinear version of PCA, particularly valuable for small sample size datasets. The principle  
636 behind them is MC<sup>23</sup>, that suggests that curvilinear (nonlinear) distances between samples may  
637 be estimated as pairwise distances over their Minimum Spanning Tree (MST) (constructed  
638 according to a selected distance). In fact, as explained in <sup>65</sup>, to approximate nonlinear  
639 (curvilinear) distances between the points of the manifold it is not necessary to reconstruct the  
640 nearest-neighbour graph. Indeed, a greedy routing process (that exploits a norm, for instance  
641 Euclidean) between the points in the multidimensional space is enough to efficiently navigate  
642 the hidden network that approximates the manifold in the multidimensional space. And a  
643 preferable greedy routing strategy, at the basis of MC-kernel, is the minimum spanning tree  
644 (MST).

645 Overall, we can conclude that both MCE in dimensionality reduction and MC-MCL in  
646 clustering perform better than the respective non-MC-based versions, and this result confirms  
647 the presence of nonlinear complexity in this dataset, generated by a three-body interaction  
648 (presence of three clusters). In addition, when considering correlation-based distances, they do  
649 not react to the presence of compositionality, since pairwise correlations are computed between

650 samples. Compositionality instead is a problem that arises when the correlations is computed  
651 between OTUs (features) from metagenomics abundance data (which are normalized by dividing  
652 each OTU count to the total sum of counts in the sample <sup>66,67</sup>), which yields unreliable results  
653 due to dependency of microbial relative abundances.

654 Moreover, because of the discovered major nonlinear complexity in the Paroni Sterbini gastric  
655 biopsy dataset, we wanted to verify whether it was generated by multi-grouping (three-body  
656 interaction problem associated to the presence of three hidden clusters). To do so, we applied  
657 PCA to three subsampled versions of the dataset (with the best normalization originally found  
658 for the complete dataset), each corresponding to the combination of two groups (Fig. 5A-C),  
659 and PCA could find significant separation (p-values <0.02 and AUC, AUPR > 0.80). To further  
660 confirm that the presence of multiple sample groups generates the data complexity, we did the  
661 same for the Tripartite Swiss-Roll (Fig. 5D-F), where we recovered the discrimination, even  
662 though two comparisons overlap to some extent (Fig. 5D and F). Furthermore, to have another  
663 confirmation that the PPI-treated samples are not separable for *H. pylori* infection, we analysed  
664 the dataset considering exclusively the PPI-treated samples. The result is that no internal  
665 separation related to *H. pylori* infection emerges within the PPI-treated patients, as shown by  
666 the best MCE result (Supplementary Figure S2).

667 In conclusion, the results confirm that linear techniques, even if supervised like LDA, are not  
668 able to resolve the differences in the data due to the presence of nonlinear complexity generated  
669 by the three-body interaction (HP-, HP+ and PPI). Once the complexity is reduced to a two-  
670 body interaction, the problem tends to vanish and PCA can detect significant differences  
671 between the groups, as shown by the PCA pairwise comparisons.

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

675

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

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

681 Table 1, top panel, shows the best results in performance of unsupervised dimension reduction  
682 (PCA, MDSwUF, MDSbc, NMDS, MCE, for details see *'Methods - PCA, MDS (or PCoA) and*  
683 *LDA'* and *'Methods - Minimum Curvilinear Embedding (MCE)'*) according to the PSI  
684 (projection-based separability index) in the space of the first two dimensions of embedding,  
685 based on the p-value of Mann-Whitney U test, AUC and the AUPR, on the three different  
686 datasets (for more details on the PSI see *'Methods - Procedure to evaluate the performance of*  
687 *the dimension reduction algorithms'*). The mean performance across all datasets is shown in  
688 the last column of the table for each method. The corresponding ranked performance for each  
689 method, based on p-value, AUC and AUPR, is presented instead in Table 2. For the Paroni  
690 Sterbini dataset, we show the results for three different labels (untreated HP-, untreated HP+  
691 and PPI-treated). For the Amir datasets, the p-values were computed for two groups, identified  
692 by the presence or absence of PPI treatment. The PSI was also applied to the data in the original  
693 high-dimensional (HD) space, as a reference to see how good the unsupervised dimension  
694 reduction approaches are in preserving the group separability in the HD. Moreover, the average  
695 p-value, AUC and AUPR best results with standard error on the original datasets, when  
696 applying leave-one-out-cross-validation (LOOCV), are shown in Supplementary Table S5.

697 For the Paroni Sterbini dataset, the PSI evaluation in the first two dimensions of embedding  
698 identifies MCE as the best dimension reduction technique that is able to preserve the group  
699 separability in the HD space. Surprisingly, MCE (presented in Fig. 4A, p-value= 0.0040, AUC  
700 = 0.967, AUPR=0.987) outdoes HD in sample separation in three groups (for HD, p-value=  
701 0.0056, AUC= 0.937, AUPR=0.967). Similarly, in Amir4, MCE (p-value=0.0047, AUC=0.906,

702 AUPR=0.920) succeeds in preserving the separability of the original HD space (in HD, p-  
703 value=0.0003, AUC=0.984, AUPR=0.985), better than the other dimension reduction methods.  
704 Finally, dimension reduction analysis on the Amir3 dataset shows that esophageal biopsies were  
705 significantly different before and after PPI treatment, as shown by MDSwUF results (p-value=  
706 0.0002, AUC=1=AUPR), that surpass the p-value, AUC and AUPR values in HD space (p-  
707 value=0.0011, AUC=0.953, AUPR=0.957). Markedly, MDSwUF reaches a value of AUPR and  
708 AUC of 1, meaning perfect classification of the samples.

709 Overall, when averaging across all datasets, the two metrics based on AUC and AUPR pointed  
710 out that MDSwUF (AUC=0.932, AUPR= 0.949) gave the best results of separability compared  
711 to HD (AUC=0.958, AUPR=0.970), followed by MCE with closer results (AUC=0.919,  
712 AUPR=0.933), while MCE gave the highest separability according to p-value (p-  
713 value=0.0055). Then PCA is the third best result (p-value=0.0095, AUC=0.896, AUPR=0.914),  
714 followed by NMDS and MDSbc. However, to conclude what is the best method, we considered  
715 an evaluation based on ranking (Table 2). It is important to note that MCE was the dimension  
716 reduction approach that ranked first in performance across all the datasets, followed by  
717 MDSwUF (Table 2). Hence, the results of sample separability suggest the presence of hidden  
718 patterns that emerge by applying nonlinear dimension reduction techniques like MCE and  
719 MDSwUF.

720 Then clustering algorithms, MCL and its Minimum Curvilinear version (for more information  
721 see '*Methods - From Markov Clustering (MCL) to Minimum Curvilinear Markov Clustering*  
722 (*MC-MCL*)'), were used to confirm the well-possedeness of the hidden patterns that were  
723 recognized by nonlinear dimension reduction. The best results as highest accuracies in each  
724 dataset and the mean performance across all the datasets are exhibited in Table 1, bottom panel.  
725 As already discussed in the previous section, the minimum curvilinear version of MCL (MC-  
726 MCL, acc=0.67) outperforms the MCL clustering algorithm (acc=0.58) in the Paroni Sterbini  
727 dataset, confirming the presence of underlying non-linear complexity in the data. However, the

728 accuracy doesn't reach high values, because of the difficulty in pattern recognition generated  
729 by the three-body problem in the HD space. Curiously, the accuracies for four clusters (HP-,  
730 HP+, PPI-, PPI+) drop to 0.58 for MC-MCL and to 0 for MCL, supporting the hypothesis that  
731 three clusters are more congruous than four clusters. Notably in Amir3, MC-MCL attains high  
732 clustering accuracy (acc=0.81), compared to MCL (acc=0.69). This is the dataset for which,  
733 surprisingly, Amir and collaborators did not find significant changes in the esophageal tissue  
734 microbiota following PPI-treatment, using classical MDS unsupervised multivariate method  
735 with unweighted UniFrac distance <sup>21</sup>. Instead, in the gastric fluid dataset (Amir 4), MC-MCL  
736 and MCL got the same accuracy of 0.75, where a significant separation of samples according  
737 to PPI consumption was already proved in the original article <sup>21</sup>.

738 However, we have to clarify that normalizations besides scaling (DRS and DCS) and log-  
739 transformation ( $\log(1+x)$ ) could potentially lead to different performance results of  
740 unsupervised analysis. Normalization is crucial to address uneven sampling depth and sparsity  
741 (high proportion of zeros) in microbiome data, like rarefying an OTU table, that is randomly  
742 sampling without replacement from each sample such that all samples have the same number  
743 of total counts (sequencing depth) <sup>68-71</sup> ([http://qiime.org/scripts/single\\_rarefaction.html](http://qiime.org/scripts/single_rarefaction.html)). This  
744 normalization is recommended to moderate the sensitivity of UniFrac distances to sequencing  
745 (sampling) depth <sup>50,72</sup>, especially differences in the presence of rare OTUs <sup>48</sup>, nonetheless it is  
746 also considered statistically improper due to the omission of data <sup>72</sup>.

747 Another normalization was introduced in 2010 by Anders and colleagues for general sequence  
748 count data (function *varianceStabilizingTransformation* implemented in the Bioconductor  
749 DESeq2 package), that uses a Variance-Stabilization Transformation (VST) by modelling  
750 microbiome count data with Negative Binomial (NB) distribution <sup>69,72</sup>.

751 We also provide the results with these two different normalizations, and we further confirm that  
752 the data are segregated in the HD space when pre-processed according to them, as shown in the  
753 p-value, AUC and AUPR tables in Additional file (for negative binomial, Supplementary

754 Tables S5-6; for rarefaction, Supplementary Table S11-12). Interestingly, across all the datasets  
755 MCE decreases its performance with these pre-processing techniques, remarkably with rarefied  
756 datasets, while the other linear techniques improve in performance (Supplementary Table S6  
757 for negative binomial; Supplementary Table S12 for rarefaction), suggesting that these  
758 adjustments linearize the datasets. Indeed, since MCE is a hierarchical technique, it needs the  
759 presence of nonlinearity to perform well. In a similar way, with these two normalizations the  
760 accuracy of MC-MCL drops down (less remarkably in the rarefaction datasets), while the  
761 performance of MCL does not increment (Supplementary Table S9 for negative binomial;  
762 Supplementary Table S14 for rarefaction). It is true that some pre-processing steps such as  
763 negative binomial tend to linearize the data but, in this manner, they can also remove important  
764 nonlinear discriminative information, as we show with the results of unsupervised analysis.  
765 Therefore, some pre-processing approaches can also cancel important nonlinear discriminant  
766 information present in the analysed data.

767

## 768 **Network analysis clarifies the effect of PPI-treatment on the gastric** 769 **microbiota**

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

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

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

778 Since we wanted to investigate the effect of PPI intake on gastric microbiota in dyspepsia, we  
779 analysed: Amir4 for gastric fluid microbiota <sup>21</sup> and Paroni Sterbini et al. dataset <sup>22</sup> for gastric

780 mucosal microflora, in the latter case restricting to PPI-treated *H. pylori*-negative (PPI-) and  
781 untreated *H. pylori* negative patients (HP-). In both studies, the samples from dyspeptic patients  
782 were analysed using the same next-generation sequencing technologies for direct sequencing  
783 of 16S rRNA gene amplicons, 454 Pyrosequencing.

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

790 In Amir4 (gastric fluid), PCA revealed that gastric fluid samples were separated into two groups  
791 according to PPI treatment along PC2 and their difference is significant ( $p$ -value  $< 0.01$ )  
792 (Supplementary Figure S3). Hence, we built the PC-corr network<sup>62</sup> using the loadings of PC2  
793 at cut-off 0.5 (Supplementary Figure S4).

794 Similarly for the Paroni Sterbini dataset (gastric mucosa), PCA (Supplementary Figure S5)  
795 could (significantly or close to significance) separate PPI-treated *H. pylori*-negative patients  
796 from untreated *H. pylori*-negative patients along PC2 and PC15 ( $p$ -value along PC2 = 0.014,  $p$ -  
797 value along PC15=0.054). Therefore we built the PC-corr network for both PC2 and PC15  
798 discriminating dimension using 0.5 cut-off (Supplementary Figure S6, panel A and B).

799 Subsequently, to investigate how PPI is affecting the microbiota in the gastric environment, we  
800 considered the conserved network, which is obtained as the union of the two PC-corr networks  
801 (obtained for PC2 and PC15) derived from the Paroni Sterbini gastric mucosa dataset  
802 intersected with the PC-corr network derived from the Amir4 gastric fluid dataset. The resulting  
803 conserved network displays the bacteria with same trend in the two datasets, i.e. either increased  
804 or decreased with PPI-treatment, respectively in red and black colour, as emphasized by the  
805 violet circle at the centre of Figure 6. Figure 7 is the same as Figure 6 but here the nodes are

806 coloured according to phylum-level taxonomy. The conserved network which arises at the  
807 overlap between the two PC-corr networks (union of Paroni Sterbini networks intersected with  
808 the Amir4 network) is statistically significant ( $p$ -value=1.00e-04), as a result of the statistical  
809 test based on trying to obtain the same conserved network by random resampling the bacteria  
810 in the two networks (Supplementary Figure S7), implying the difficulty of generating this  
811 intersection simply at random (since this intersection lies to the right of the critical value at the  
812 0.05 level in the distribution of overlap). This is an important result because it confirms the  
813 robustness of the detected conserved network as a microbiota signature perturbed by PPI  
814 treatment. The top and bottom panels in Figure 6 and 7 show instead the remaining part of  
815 Amir4's network (top panel) and of Paroni Sterbini's network (bottom panel) that are not in the  
816 intersection, and therefore might be more specific for the gastric fluid and mucosa respectively.  
817 The PPI-perturbed conserved network is characterized by a main interconnected module with  
818 nine bacteria of four different phyla (*Bacteroidetes*, *Fusobacteria*, *Proteobacteria*, *Firmicutes*)  
819 that are positively associated (red edges) and by two single bacteria order without interactions  
820 (*Streptophyta*, *Clostridiales*), all being increased following PPI treatment, except *Streptophyta*  
821 that is instead decreased with PPI-treatment (Fig. 6 and 7). Note that a mix between genera,  
822 phyla and order of bacteria can be found in the networks. The reason behind it is the availability  
823 of detail information regarding different bacteria. Some of the spotted bacteria (*Veillonella*,  
824 *Clostridiales*, *Campylobacter*) were already observed in previous studies. The genus  
825 *Veillonella* was found increased in relation to PPI use <sup>16</sup> in the gut microbiome and has been  
826 associated with increased susceptibility to *Clostridium difficile* infection <sup>73</sup>. These Gram-  
827 negative anaerobic cocci with lactate fermenting abilities are abundant in the human  
828 microbiome and are normally found in the intestines and oral mucosa of humans <sup>74</sup>.  
829 Interestingly, they favour nitrite accumulation in the stomach during nitrate reduction,  
830 promoting a carcinogenic effect <sup>1</sup>. In addition, the order *Clostridiales*, that is associated to  
831 *Clostridium difficile* infection, was also seen significantly changed in the gastrointestinal tract,



832 however Freedberg *et al.*<sup>4</sup> found it significantly decreased during PPI use, in contrast to our  
833 results. PPIs use also increases the risk of other enteric infections, apart from *C. difficile*  
834 infection, such as campylobacteriosis, as reported in<sup>75,76</sup>. Moreover, half of the bacteria present  
835 in the network normally colonize the human oral cavity. Indeed, it is the main purpose of PPI  
836 treatment to increase the stomach pH, and the higher pH of treated patients is known to  
837 favour the growth of bacteria that usually reside in the mouth and esophagus and are not  
838 adapted to survive the normal gastric acidity<sup>6,20</sup>. Among genera usually reported as part of  
839 the normal flora of the gastrointestinal tract, only *Veillonella* is found regularly at other sites,  
840 like the mouth<sup>77</sup>. *Leptotrichia* species mostly colonize the oral cavity and they were isolated  
841 from various human infections, suggesting that they are emerging human pathogens<sup>78,79</sup>.  
842 *Oribacterium* also inhabits the mouth, besides the upper respiratory tract<sup>80</sup>. *Prevotella* is a  
843 genus of Gram-negative bacteria that tend to colonize the human gut, mouth and vagina, and  
844 may cause infections, mostly observed in the oral cavity (odontogenic infections)<sup>79</sup>.  
845 *Porphyromonas* has been found by<sup>81</sup> as part of the salivary microbiome. Both *Prevotella* and  
846 *Porphyromonas* contribute to the formation of abscesses and soft tissue infections in various  
847 part of the body and they can cause infections, including periodontal and endodontal diseases  
848<sup>82</sup>. *Capnocytophaga* are inhabitants of the oral cavity too, and these opportunistic pathogens  
849 can cause infections (both in immunocompromised and immunocompetent hosts), the severity  
850 of which depend on the immune status of the host<sup>83,84</sup>. As well, *Granulicatella* are Gram-  
851 positive cocci normally found in the oral flora and are uncommon causes of infections,  
852 nevertheless they can cause infections, including bloodstream infection and infective  
853 endocarditis<sup>85</sup>. Besides, the genus *Fusobacterium* inhabits the mucosal membranes of humans  
854 and all its species are parasites of humans<sup>86</sup>, and some species are found in the oral cavity. The  
855 remaining bacteria (*Campylobacter*, *Bulleidia*) do not belong to the oral microbiota<sup>82</sup>. The  
856 genus *Campylobacter* was increased in relation to PPI use and the increased abundance of these  
857 Gram-negative bacteria has the potential to cause diseases and infections in humans (most

858 commonly diarrhoea). Due to the induced increase of pH, PPI is hypothesised to facilitate  
859 gastrointestinal infections and a study by Brophy *et al.*<sup>87</sup> reported an increased risk of  
860 *Campylobacter* infection following PPI therapy. Moreover Campylobacteriosis, mostly caused  
861 by eating undercooked foods derived from poultry or other warm-blooded animals or contact  
862 with contaminated water or ice<sup>88</sup>, has been shown by the Dutch National Institute for Public  
863 Health and the Environment to noticeably increase in incidence when PPI use grows<sup>75</sup>.  
864 Altogether, PC-corr approach was applied on gastric fluid and gastric mucosal datasets (in the  
865 latter case, excluding the samples positive to *H. pylori* infection) to investigate how PPI is  
866 affecting the gastric microbiota (both gastric fluid and gastric mucosal microbiota), because of  
867 PC-corr's ability to pinpoint the combination of bacteria that play a major role in the  
868 discrimination of the samples, in this case according to PPI intake. The PC-corr conserved  
869 network identified eleven genera and order of bacteria, which belong to the phyla  
870 (*Bacteroidetes*, *Fusobacteria*, *Proteobacteria*, *Firmicutes*) commonly found in the stomach  
871 which, with exception of *Streptophyta*, demonstrated increased abundance following PPI  
872 treatment. Mostly all the found bacteria were not reported in previous studies, except  
873 *Veillonella*, *Clostridiales* and *Campylobacter*, but they were found as inhabitants of the oral  
874 cavity and/or possible cause of infections and diseases in humans. Hence, and in concordance  
875 to previous studies<sup>6,20</sup>, these results point out that PPI treatment, by increasing the intragastric  
876 pH, favours the growth of bacteria that usually reside in the mouth and survive through the  
877 harsh acidic conditions of the stomach. Furthermore, the results suggest that PPI-associated  
878 increases of some bacterial populations may lead to infections and diseases or increase  
879 susceptibility for other bacterial infections (like *Veillonella*) or promote a carcinogenic effect  
880 (like *Veillonella*). Previous studies have highlighted that PPI intake is associated with decreased  
881 bacterial richness<sup>16,18,89,90</sup>, increased risk of enteric and other infections (e.g. caused by  
882 *Salmonella*, *Clostridium difficile*, *Shigella*, *Listeria*)<sup>17,91</sup>, increase in the abundance of oral and  
883 upper GI tract commensals and potential pathogenic bacteria (e.g. *Enterococcus*,

884 *Streptococcus*, *Staphylococcus*, and *Escherichia coli* )<sup>16,17</sup> in the gut microbiota. Nevertheless,  
885 our analysis by means of PC-corr does not spot single bacteria perturbed in the gastric  
886 environment by PPI treatment, but a community of bacteria is altered in abundance by PPIs and  
887 their inter-specific bacterial interactions in the gastric niche.

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

892

### 893 **Network analysis clarifies the effect of *H. pylori* infection on gastric mucosal** 894 **microbiota**

895 The stomach was long thought sparsely colonized by bacteria due to the gastric microbicidal  
896 acidic barrier (pH<4.0)<sup>92</sup>. This view dramatically changed with the discovery of the Gram-  
897 negative bacterium *H. pylori* in the 1980's by Warren and Marshall<sup>93</sup>, that is a carcinogenic  
898 bacterial pathogen infecting the stomach of more than one-half of the world's  
899 human population. This human pathogen is able to survive in the highly acidic environment  
900 within the stomach by producing cytoplasmic urease that, by catalysing the hydrolysis of urea  
901 into CO<sub>2</sub> and NH<sub>4</sub>, produces a neutralizing ammonia cloud around it<sup>19,94,95</sup>. However, most *H.*  
902 *pylori* avoid the acidic environment of the gastric lumen by swimming towards the mucosal cell  
903 surface (using their polar flagella and chemotaxis mechanisms) and may adhere and invade the  
904 gastric mucosal epithelial cells<sup>96,97</sup>. Hence, it doesn't represent a dominant species in gastric  
905 fluid microbiota<sup>98</sup>, but was found to generally to reside in the gastric mucosae<sup>5,96,99</sup>.

906 Persistent (chronic) infection with this Gram-negative bacterium induces changes in gastric  
907 physiology and immunology, e.g. reduced gastric acidity and parietal cell mass, perturbed  
908 nutrient availability, local innate immune responses<sup>100,101</sup>, that most probably induces shift in  
909 gastric microbiota composition<sup>100</sup>. Although *H. pylori* colonization usually persists in the

910 human stomach for many decades without adverse effects, the infection of this bacteria is  
911 associated with increased risk for several diseases, including peptic ulcers, chronic gastritis,  
912 mucosa-associated lymphoid tissue lymphoma, gastric adenocarcinoma <sup>102,103</sup>, and dyspepsia  
913 <sup>104,105</sup>. The potential alterations induced by the *H. pylori* can in turn lead to dysbiosis and may  
914 cause aberrant proinflammatory immune responses <sup>106</sup>, susceptibility to bacterial pathogens and  
915 increased risk of gastric disease, including cancer <sup>1,107</sup>. However, the effect of *H. pylori*  
916 infection on overall composition of gastric microbiota at genus level and the bacterial interplay  
917 in presence of this widespread human infection remain unclear.

918 To investigate the influence of *H. pylori* infection on the gastric mucosal microbiota, we  
919 analysed: 1) Paroni Sterbini *et al.* <sup>22</sup> considering only PPI-untreated dyspeptic patients, either  
920 infected (HP+) or not by *H. pylori* (HP-); 2) Parsons *et al.* <sup>29</sup> restricting to PPI-untreated patients  
921 from: i) normal stomach group with no evidence of *H. pylori* infection; ii) *H. pylori* gastritis  
922 group with evidence of *H. pylori* infection. Even though the same technology is important for  
923 a comparative study, unfortunately in the literature there was no such data available like Paroni  
924 Sterbini's one, that is 16S rRNA gene pyrosequencing data (derived from gastric mucosal  
925 microflora in dyspeptic untreated patients either positive or negative for *H. pylori*). Despite  
926 this, the two studied datasets, obtained with two different next-generation sequencing  
927 technologies for direct sequencing of 16S rRNA gene amplicons (454 Pyrosequencing for  
928 Paroni Sterbini *et al.* and Illumina MiSeq for Parsons *et al.*) <sup>108</sup>, both contain community  
929 profiling of gastric mucosa-associated microbiota in PPI-untreated *H. pylori*-negative and -  
930 positive subjects. However, for the sake of clarity, we have to specify a difference: while in  
931 Paroni Sterbini's dataset the gastric mucosal biopsy specimens were collected from patients  
932 with dyspepsia, this is not the case for Parsons's data.

933 To enhance the understanding of the *H. pylori*-triggered microbial perturbation in this  
934 ecological niche, we employed again PC-corr algorithm, that is able to associate to any PCA  
935 analysis of an omic dataset, where a sample separation emerges, a network of discriminative

936 features (for details see '*Methods-PC-corr network*'). The analysis of the 16S rRNA sequencing  
937 data was restricted only the overlapping OTUs, excluding *Helicobacter* because our goal is to  
938 investigate its impact on the rest of the microbial network.

939 In Paroni Sterbini's dataset, since PCA could significantly separate gastric mucosal biopsy  
940 samples of PPI-untreated patients according to *H. pylori*-positivity (p-value=0.01) along PC2  
941 (Supplementary Figure S8), the PC-corr network was constructed from PC2 loadings at 0.5 cut-  
942 off (Supplementary Figure S9). Similarly, for Parsons' dataset, since PCA (Supplementary  
943 Figure S10) could significantly separate patients from the normal stomach group with no  
944 evidence of *H. pylori* infection and PPI-untreated (Control) from *H. pylori* gastritis group  
945 positive to *H. pylori* infection and not using PPIs (HPGas) along PC1 (p-value along PC1  
946 <0.01.), the PC-corr network was constructed from this discriminating dimension at 0.5 cut-off  
947 (Supplementary Figure 11). The obtained microbial differential networks (top panel for and  
948 bottom panel in Figure 8, coloured according to phylum level) pinpointed, from the system  
949 point of view, the bacteria affected by *H. pylori* infection in the gastric mucosa, that are  
950 precisely bacteria whose abundance is decreased in *H. pylori*-positive patients. A presumable  
951 explanation of this trend is already pointed out in literature, where the presence of *H. pylori*  
952 leads to a reduced gastric microbial diversity<sup>109-111</sup>. Nevertheless, in some cases the diversity  
953 increases again, because of diverse factors that allow survival and colonization of bacteria in  
954 the stomach<sup>1,112</sup>. Then, the preserved network of gastric mucosa microbiota was constructed  
955 by intersecting the two PC-corr networks obtained from Paroni Sterbini's and Parsons's dataset.

956 Figure 8, middle panel, shows the conserved network (violet circle), which presents the  
957 common bacteria coloured according to phylum level and their associations. The spotted  
958 bacteria display decreased abundance with *H. pylori* infection (i.e. increased in *H. pylori*-  
959 *negative* subjects) in both the two 16S rRNA gene sequencing data. By performing a statistical  
960 test based on random resampling of the bacteria in the two networks, we verified that the shown  
961 bacterial conserved network is statistically significant and difficult to be generated at random

962 (p-value=1.00e-04), because getting this intersection at random is very rare (Supplementary  
963 Figure S12). The top and bottom panels in Figure 8 show instead the remaining part of Paroni  
964 Sterbini's network (top panel) and of Parsons's network (bottom panel) that are not in the  
965 intersection. At the genus level, a study by Klymiuk *et al.*<sup>113</sup> identified *Actinomyces*,  
966 *Granulicatella*, *Veillonella*, *Fusobacterium*, *Neisseria*, *Helicobacter*, *Streptococcus*, and  
967 *Prevotella* as significantly different between the *H. pylori*-positive and *H. pylori*-negative  
968 gastric samples. These bacteria do not emerge in the conserved network, while they all (except  
969 *Neisseria*) appear altered (decreased) during *H. pylori* infection in the study by Parsons and  
970 colleagues (present in the bottom panel of Figure 8).

971 Our analysis pinpoints a conserved network from two independent 16S rRNA gene sequencing  
972 data, that reveals microbial communities altered by *H. pylori* infection and their interactions in  
973 the gastric mucosa. It revealed a main core of six associated bacteria (with positive association,  
974 red edges) and two single nodes without any interaction with the main module, from three  
975 different phyla (*Proteobacteria*, *Firmicutes*, *Actinobacteria*) all resulting decreased in *H.*  
976 *pylori*-infected subjects (that is increased in non-infected subjects). The decreased abundance  
977 of the phyla *Firmicutes* and *Actinobacteria* in *H. pylori*-positive patients with respect to *H.*  
978 *pylori*-negative subjects was already shown in a previous study by Maldonado-Contreras *et al.*  
979 <sup>114</sup>. In addition, other studies have demonstrated an increased colonization of *Proteobacteria* in  
980 *H. pylori*-positive patients<sup>114,115</sup>, while the obtained conserved PC-corr network shows that the  
981 bacteria from this phylum are instead decreased in those individuals. Among the spotted  
982 bacteria, *Methylobacterium* is a genus of facultative methylotrophic bacteria that are commonly  
983 found in diverse natural environments (such as leaf surfaces, soil, dust, and fresh water) and in  
984 hospital environment due to contaminated tap water. *Methylobacterium* species can cause  
985 health care-associated infections (mainly catheter infection), especially in  
986 immunocompromised patients<sup>116</sup>. In addition, *Sphingomonas* plays a role in human health, as  
987 some of the sphingomonads (in particular *Sphingomonas paucimobilis*) are the cause of a range

988 of mostly nosocomial, non-life-threatening infections. *Sphingomonas* species are widely spread  
989 in nature, having been isolated from many sources, from water habitats to clinical settings <sup>117</sup>,  
990 *Pseudomonas*, due to its great metabolic versatility, can also colonize different types of niches  
991 <sup>118</sup>, including soil and water, in addition to plant and animal associations, and includes  
992 pathogenic species in humans <sup>119</sup>. *Acinetobacter* species are instead common, free-living  
993 saprophytes found in soil, water, sewage and foods and are ubiquitous organisms in hospitals.  
994 They have been increasingly identified as a key source of infection in debilitated patients in  
995 hospitals, due to their rapid development of resistance to antimicrobials <sup>120</sup>. In particular, one  
996 species, *Acinetobacter lwoffii*, can trigger gastritis, apart from *H. pylori* <sup>121</sup>. *Propionibacterium*,  
997 so named for their unique ability to synthesize propionic acid by using unusual transcarboxylase  
998 enzymes <sup>122</sup>, are primarily facultative pathogens and commensals of humans, living on the skin,  
999 while other members are widely employed for synthesizing vitamin B<sub>12</sub>, tetrapyrrole  
1000 compounds, and propionic acid, as well as used as probiotics <sup>123</sup>. *Catonella* is another node in  
1001 the network and this bacterial genus is obligative anaerobic, non-spore-forming and non-motile,  
1002 with one known species (*Catonella morbi*) from the human gingival crevice <sup>124,125</sup>, that has been  
1003 associated with periodontitis <sup>124</sup> and endocarditis <sup>126</sup>. Besides, the bacterial genus  
1004 *Enhydrobacter* so far contains a single species, *Enhydrobacter aerosaccus*, a Gram negative  
1005 non-motile bacterium that is both oxidase and catalase positive and shows gas vacuoles <sup>127,128</sup>.  
1006 *Bulleidia*, a Gram-positive, non-spore-forming, anaerobic and non-motile genus, has one  
1007 known species too (*Bulleidia extracta*)<sup>129</sup>.  
1008 In conclusion, by means of the PC-corr approach, we determined the combination of bacteria  
1009 responsible for the difference between *H. pylori*-positive and *H. pylori*-negative gastric mucosa  
1010 of untreated patients and their microbe-microbe interactions. All the bacteria, both in the  
1011 conserved network and not, were decreased in *H. pylori*-infected individuals (i.e. increased in  
1012 *H. pylori*-negative group). *H. pylori*, like acid suppressing medications (for the treatment of  
1013 dyspepsia), alters the population structure of the gastric and intestinal microbiota <sup>130</sup> and

1014 regularly, this bacterium constitutes most of the gastric microbiota <sup>112</sup>, literally depleting  
1015 bacterial biodiversity. Moreover, most of the identified bacteria represent bacteria of potential  
1016 health concern, as agents of diseases and infections.

1017

## 1018 **Discussion**

1019 This study indicates the necessity of including nonlinear multidimensional techniques into  
1020 clinical studies based on 16S metagenomic sequencing data, since drawing a study's  
1021 conclusions by solely relying on linear techniques, such as PCA and MDS, can lead to data  
1022 misinterpretation and impair the translational path from research to diagnostic. In the era of  
1023 post-genomics and systems approaches, nonlinear dimension reduction and clustering by MCE  
1024 and MC-MCL can offer new insights into complex clinical 16S metagenomics data, like the  
1025 ones studied in this article or the presence of clinical sub-types, and serve as a valuable tool in  
1026 the run towards precision medicine. Moreover, this study shows how it is possible to  
1027 complement multivariate analysis by means of network analysis employing PC-corr algorithm,  
1028 that accounts for the bacteria responsible for the sample discrimination and their co-occurrence  
1029 relationships. Precisely, from the system point of view the obtained microbial differential  
1030 networks pinpointed marked bacteria-bacteria interactions and modules affected by PPI  
1031 treatment in the gastric environment in dyspepsia and by *H. pylori* infection in the gastric  
1032 mucosa. We suggest that our findings can be an important starting point to design new therapies  
1033 that consider not only *H. pylori* infection but also the directly associated microbial alterations  
1034 as well as the indirect alterations due to the drugs used for *H. pylori* eradication such as PPI.

1035

## 1036 **List of abbreviations**

1037 LDA: Linear Discriminant Analysis

1038 MC: Minimum Curvilinearity



- 1039 MCE: Minimum Curvilinear Embedding
- 1040 MCL: Markov Clustering
- 1041 MC-MCL: Minimum Curvilinear Markov Clustering
- 1042 MDS: Multidimensional Scaling
- 1043 MDSbc: Multidimensional Scaling with Bray-Curtis dissimilarity
- 1044 MDSwUF: Multidimensional Scaling with weighted UniFrac distance
- 1045 MST: minimum spanning tree
- 1046 ncMCE: non-centred Minimum Curvilinear Embedding
- 1047 NMDS: non-metric (Sammon criterion) Multidimensional Scaling
- 1048 PC: Principal Component
- 1049 PCA: Principal Component Analysis
- 1050 PCoA: Principal Coordinate Analysis
- 1051 PPI: Proton Pump Inhibitor
- 1052 PSI: Projection-based separability index
- 1053 SVD: Singular Value Decomposition

1054

## 1055 **Declarations**

### 1056 **Ethics approval and consent to participate**

1057 Not applicable, because the used datasets have been generated by previous biomedical  
1058 studies, for which ethics approvals and consents were formerly collected.

1059

### 1060 **Consent for publication**

1061 Not applicable

1062

### 1063 **Availability of data and materials**

1064 Not applicable.

1065

1066 **Competing interests**

1067 The authors declare that they have no competing interests.

1068

1069 **Funding**

1070 This work was supported by the Dresden International Graduate School for Biomedicine and  
1071 Bioengineering (DIGS-BB), granted by the Deutsche Forschungsgemeinschaft (DFG) in the  
1072 context of the Excellence Initiative. PS is supported by Estonian Research Council Starting  
1073 Grant PUT1130.

1074

1075 **Authors' contributions**

1076 CVC developed Minimum Curvilinearity (MCE), Minimum Curvilinear Markov Clustering  
1077 (MC-MCL) and the Projection-based Separability Index (PSI). CVC conceived all the study  
1078 and the data analysis workflow with feedbacks from MiSc and SWG. SC, CD and AP  
1079 performed the computational analysis of the data and realized the figures under the CVC  
1080 guidance. SC, CD, AP together with CVC wrote the manuscript with valuable suggestions of  
1081 PS. FPS, LM, GC, GI, BP, MaSa, GG and AG provided data and knowledge about the Paroni  
1082 Sterbini *et al.* data cohort. BNP, UZI and MP provided data and knowledge about the Parsons  
1083 *et al.* data cohort. All authors discussed the results and revised the manuscript.

1084

1085 **Acknowledgements**

1086 Not applicable

1087

1088 **References**

1089 1. Nardone, G. & Compare, D. The human gastric microbiota: Is it time to rethink the  
1090 pathogenesis of stomach diseases? *United Eur. Gastroenterol. J.* **3**, 255–260 (2015).

- 1091 2. Quigley, E. M. M. Gut microbiome as a clinical tool in gastrointestinal disease  
1092 management: are we there yet? *Nat. Rev. Gastroenterol. Hepatol.* **14**, 315–320 (2017).
- 1093 3. Strand, D. S., Kim, D. & Peura, D. A. 25 years of proton pump inhibitors: A  
1094 comprehensive review. *Gut and Liver* **11**, 27–37 (2017).
- 1095 4. Freedberg, D. E., Lebwohl, B. & Abrams, J. A. The impact of proton pump inhibitors  
1096 on the human gastrointestinal microbiome. *Clinics in Laboratory Medicine* **34**, 771–  
1097 785 (2014).
- 1098 5. Wu, W. M., Yang, Y. S. & Peng, L. H. Microbiota in the stomach: new insights. *J. Dig.*  
1099 *Dis.* **15**, 54–61 (2014).
- 1100 6. Vesper, B. *et al.* The Effect of Proton Pump Inhibitors on the Human Microbiota. *Curr.*  
1101 *Drug Metab.* **10**, 84–89 (2009).
- 1102 7. Scarpignato, C. *et al.* Effective and safe proton pump inhibitor therapy in acid-related  
1103 diseases ? A position paper addressing benefits and potential harms of acid  
1104 suppression. *BMC Med.* **14**, 179 (2016).
- 1105 8. Yadlapati, R. & Kahrilas, P. J. When is proton pump inhibitor use appropriate? *BMC*  
1106 *Med.* **15**, 36 (2017).
- 1107 9. Harmon, R. C. & Peura, D. A. Evaluation and management of dyspepsia. *Therap. Adv.*  
1108 *Gastroenterol.* **3**, 87–98 (2010).
- 1109 10. Malfertheiner, P. *et al.* Management of *Helicobacter pylori* infection—the Maastricht  
1110 IV/ Florence Consensus Report. *Gut* **61**, 646–664 (2012).
- 1111 11. Rosen, R. *et al.* 16S community profiling identifies proton pump inhibitor related  
1112 differences in gastric, lung, and oropharyngeal microflora. *J. Pediatr.* **166**, 917–923  
1113 (2015).
- 1114 12. Lanas, A. We are using too many PPIs, and we need to stop: A European perspective.  
1115 *American Journal of Gastroenterology* **111**, 1085–1086 (2016).
- 1116 13. Vakil, N. Prescribing proton pump inhibitors: Is it time to pause and rethink? *Drugs* **72**,

- 1117 437–445 (2012).
- 1118 14. Tran-Duy, A., Spaetgens, B., Hoes, A. W., de Wit, N. J. & Stehouwer, C. D. A. Use of  
1119 Proton Pump Inhibitors and Risks of Fundic Gland Polyps and Gastric Cancer:  
1120 Systematic Review and Meta-analysis. *Clin. Gastroenterol. Hepatol.* **14**, 1706-1719.e5  
1121 (2016).
- 1122 15. Malfertheiner, P., Kandulski, A. & Venerito, M. Proton-pump inhibitors:  
1123 Understanding the complications and risks. *Nat. Rev. Gastroenterol. Hepatol.* **14**, 697–  
1124 710 (2017).
- 1125 16. Imhann, F. *et al.* Proton pump inhibitors affect the gut microbiome. *Gut* **65**, 740–748  
1126 (2016).
- 1127 17. Jackson, M. A. *et al.* Proton pump inhibitors alter the composition of the gut  
1128 microbiota. *Gut* **65**, 749–756 (2016).
- 1129 18. Tsuda, A. *et al.* Influence of proton-pump inhibitors on the luminal microbiota in the  
1130 gastrointestinal tract. *Clin. Transl. Gastroenterol.* **6**, e89 (2015).
- 1131 19. Williams, C. & McColl, K. E. L. Review article: proton pump inhibitors and bacterial  
1132 overgrowth. *Aliment. Pharmacol. Ther.* **23**, 3–10 (2006).
- 1133 20. Sanduleanu, S., Jonkers, D., De Bruine, A., Hameeteman, W. & Stockbrügger, R. W.  
1134 Non-Helicobacter pylori bacterial flora during acid-suppressive therapy: Differential  
1135 findings in gastric juice and gastric mucosa. *Aliment. Pharmacol. Ther.* **15**, 379–388  
1136 (2001).
- 1137 21. Amir, I., Konikoff, F. M., Oppenheim, M., Gophna, U. & Half, E. E. Gastric  
1138 microbiota is altered in oesophagitis and Barrett’s oesophagus and further modified by  
1139 proton pump inhibitors. *Environ. Microbiol.* **16**, 2905–2914 (2014).
- 1140 22. Paroni Sterbini, F. *et al.* Effects of Proton Pump Inhibitors on the Gastric Mucosa-  
1141 Associated Microbiota in Dyspeptic Patients. *Appl. Environ. Microbiol.* **82**, 6633–6644  
1142 (2016).

- 1143 23. Cannistraci, C. V., Ravasi, T., Montevecchi, F. M., Ideker, T. & Alessio, M. Nonlinear  
1144 dimension reduction and clustering by Minimum Curvilinearity unfold neuropathic  
1145 pain and tissue embryological classes. in *Bioinformatics* **27**, i531–i539 (2011).
- 1146 24. Kinross, J. M., Darzi, A. W. & Nicholson, J. K. Gut microbiome-host interactions in  
1147 health and disease. *Genome Med.* **3**, 14 (2011).
- 1148 25. Legendre, P. & Legendre, L. F. J. *Numerical ecology*. **24**, (Elsevier, 2012).
- 1149 26. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community  
1150 sequencing data. *Nat. Methods* **7**, 335–6 (2010).
- 1151 27. Wang, Q., Garrity, G. M., Tiedje, J. M. & Cole, J. R. Naïve Bayesian classifier for  
1152 rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl. Environ.*  
1153 *Microbiol.* **73**, 5261–5267 (2007).
- 1154 28. Caporaso, J. G. *et al.* PyNAST: A flexible tool for aligning sequences to a template  
1155 alignment. *Bioinformatics* **26**, 266–267 (2010).
- 1156 29. Parsons, B. N. *et al.* Comparison of the human gastric microbiota in hypochlorhydric  
1157 states arising as a result of. *PLOS Pathog.* **13**, 1–19 (2017).
- 1158 30. Cannistraci, C. V., Alanis-Lobato, G. & Ravasi, T. Minimum curvilinearity to enhance  
1159 topological prediction of protein interactions by network embedding. *Bioinformatics*  
1160 **29**, 199–209 (2013).
- 1161 31. Smialowski, P., Frishman, D. & Kramer, S. Pitfalls of supervised feature selection.  
1162 *Bioinformatics* **26**, 440–443 (2009).
- 1163 32. Ringnér. What is principal component analysis? *Nat. Biotechnol.* **26**, 303–304 (2008).
- 1164 33. Jolliffe, I. T. Principal Component Analysis. *Springer Ser. Stat.* **98**, 487 (2002).
- 1165 34. Dinsdale, E. A. *et al.* Multivariate analysis of functional metagenomes. *Front. Genet.* **4**,  
1166 41 (2013).
- 1167 35. Ramette, A. Multivariate analyses in microbial ecology. *FEMS Microbiol. Ecol.* **62**,  
1168 142–160 (2007).

- 1169 36. Moitinho-Silva, L. *et al.* Specificity and transcriptional activity of microbiota  
1170 associated with low and high microbial abundance sponges from the Red Sea. *Mol.*  
1171 *Ecol.* **23**, 1348–1363 (2014).
- 1172 37. Bayer, K. *et al.* GeoChip-based insights into the microbial functional gene repertoire of  
1173 marine sponges (high microbial abundance, low microbial abundance) and seawater.  
1174 *FEMS Microbiol. Ecol.* **90**, 832–843 (2014).
- 1175 38. Tenenbaum, J. B., de Silva, V. & Langford, J. C. A global geometric framework for  
1176 nonlinear dimensionality reduction. *Science* **290**, 2319–23 (2000).
- 1177 39. Alanis-Lobato, G., Cannistraci, C. V., Eriksson, A., Manica, A. & Ravasi, T.  
1178 Highlighting nonlinear patterns in population genetics datasets. *Sci. Rep.* **5**, 8140  
1179 (2015).
- 1180 40. Legendre, P. & De Cáceres, M. Beta diversity as the variance of community data:  
1181 Dissimilarity coefficients and partitioning. *Ecol. Lett.* **16**, 951–963 (2013).
- 1182 41. Paliy, O. & Shankar, V. Application of multivariate statistical techniques in microbial  
1183 ecology. *Mol. Ecol.* **25**, 1032–1057 (2016).
- 1184 42. Zand, M. S., Wang, J. & Hilchey, S. Graphical Representation of Proximity Measures  
1185 for Multidimensional Data: Classical and Metric Multidimensional Scaling. *Math. J.*  
1186 **17**, (2015).
- 1187 43. Cox, M. A. A. & Cox, T. F. Multidimensional Scaling. *Handb. Data Vis.* (2008).  
1188 doi:10.1007/978-3-540-33037-0\_14
- 1189 44. Sammon, J. W. A Nonlinear Mapping for Data Structure Analysis. *IEEE Trans.*  
1190 *Comput.* **C18**, 401–409 (1969).
- 1191 45. Beals, E. W. Bray-curtis ordination: An effective strategy for analysis of multivariate  
1192 ecological data. in *Advances in Ecological Research* **14**, 1–55 (1984).
- 1193 46. Bray, J. R. & Curtis, J. T. An Ordination of the Upland Forest Communities of  
1194 Southern Wisconsin. *Ecol. Monogr.* **27**, 325–349 (1957).

- 1195 47. Whittaker, R. H. Vegetation of the Siskiyou Mountains, Oregon and California. *Ecol.*  
1196 *Monogr.* **30**, 279–338 (1960).
- 1197 48. Lozupone, C., Lladser, M. E., Knights, D., Stombaugh, J. & Knight, R. UniFrac: An  
1198 effective distance metric for microbial community comparison. *ISME J.* **5**, 169–172  
1199 (2011).
- 1200 49. Lozupone, C. A., Hamady, M., Kelley, S. T. & Knight, R. Quantitative and qualitative  
1201 beta diversity measures lead to different insights into factors that structure microbial  
1202 communities. *Appl. Environ. Microbiol.* **73**, 1576–85 (2007).
- 1203 50. Lozupone, C. & Knight, R. UniFrac: a new phylogenetic method for comparing  
1204 microbial communities. *Appl. Environ. Microbiol.* **71**, 8228–35 (2005).
- 1205 51. Chen, J. *et al.* Associating microbiome composition with environmental covariates  
1206 using generalized UniFrac distances. *Bioinformatics* **28**, 2106–13 (2012).
- 1207 52. Podani, J. & Miklós, I. Resemblance Coefficients and the Horseshoe Effect in Principal  
1208 Coordinates Analysis. *Ecology* **83**, 3331–3343 (2002).
- 1209 53. Papadopoulos, F., Psomas, C. & Krioukov, D. Network mapping by replaying  
1210 hyperbolic growth. *IEEE/ACM Trans. Netw.* **23**, 198–211 (2015).
- 1211 54. Muscoloni, A., Thomas, J. M., Ciucci, S., Bianconi, G. & Cannistraci, C. V. Machine  
1212 learning meets complex networks via coalescent embedding in the hyperbolic space.  
1213 *Nat. Commun.* **8**, 1615 (2017).
- 1214 55. Muscoloni, A. & Cannistraci, C. V. Minimum curvilinear automata with similarity  
1215 attachment for network embedding and link prediction in the hyperbolic space. (2018).
- 1216 56. Zagar, L. *et al.* Stage prediction of embryonic stem cell differentiation from genome-  
1217 wide expression data. **27**, 2546–2553 (2011).
- 1218 57. Ryu, T., Seridi, L. & Ravasi, T. The evolution of ultraconserved elements with  
1219 different phylogenetic origins. *BMC Evol. Biol.* **12**, 236 (2012).
- 1220 58. Sales, S. *et al.* Gender, Contraceptives and Individual Metabolic Predisposition Shape a

- 1221 Healthy Plasma Lipidome. *Sci. Rep.* **6**, 27710 (2016).
- 1222 59. Acevedo, A., Ciucci, S., Kuo, M. J., Durán, C. & Cannistraci, C. V. Measuring group-  
1223 separability in geometrical space for evaluation of pattern recognition and embedding  
1224 algorithms. *ArXiv:1912.12418* 1–20 (2019).
- 1225 60. van Dongen, S. Graph clustering by flow simulation. *Graph Stimul. by flow Clust.*  
1226 (2000). doi:10.1016/j.cosrev.2007.05.001
- 1227 61. Duran, C., Acevedo, A., Ciucci, S., Muscoloni, A. & Cannistraci, C. Nonlinear Markov  
1228 Clustering by Minimum Curvilinear Sparse Similarity. *ArXiv:1912.12211* 1–17 (2019).
- 1229 62. Ciucci, S. *et al.* Enlightening discriminative network functional modules behind  
1230 Principal Component Analysis separation in differential-omic science studies. 1–24  
1231 (2017). doi:10.1038/srep43946
- 1232 63. Jones, D. L. The Fathom Toolbox for Matlab: multivariate ecological and  
1233 oceanographic data analysis. *Coll. Mar. Sci. Univ. South Florida, St. Petersburg, FL,*  
1234 *USA* (2014).
- 1235 64. Montecucco, C. & Rappuoli, R. Living dangerously: how *Helicobacter pylori* survives  
1236 in the human stomach. *Nat. Rev. Mol. Cell Biol.* **2**, 457–466 (2001).
- 1237 65. Boguñá, M., Krioukov, D. & Claffy, K. C. Navigability of complex networks. *Nat.*  
1238 *Phys.* **5**, 74–80 (2008).
- 1239 66. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data.  
1240 *PLoS Comput. Biol.* **8**, (2012).
- 1241 67. Kurtz, Z. D. *et al.* Sparse and Compositionally Robust Inference of Microbial  
1242 Ecological Networks. *PLoS Comput. Biol.* **11**, e1004226 (2015).
- 1243 68. Wong, R. G., Wu, J. R. & Gloor, G. B. Expanding the UniFrac toolbox. *PLoS One* **11**,  
1244 e0161196 (2016).
- 1245 69. Weiss, S. *et al.* Normalization and microbial differential abundance strategies depend  
1246 upon data characteristics. *Microbiome* **5**, 27 (2017).



- 1247 70. Navas-Molina, J. A. *et al.* Advancing our understanding of the human microbiome  
1248 using QIIME. in *Methods in Enzymology* **531**, 371–444 (2013).
- 1249 71. Hughes, J. B. & Hellmann, J. J. The application of rarefaction techniques to molecular  
1250 inventories of microbial diversity. in *Methods in Enzymology* **397**, 292–308 (2005).
- 1251 72. McMurdie, P. J., Holmes, S., Hoffmann, C., Bittinger, K. & Chen, Y. Waste Not, Want  
1252 Not: Why Rarefying Microbiome Data Is Inadmissible. *PLoS Comput. Biol.* **10**,  
1253 e1003531 (2014).
- 1254 73. Antharam, V. C. *et al.* Intestinal dysbiosis and depletion of butyrogenic bacteria in  
1255 *Clostridium difficile* infection and nosocomial diarrhea. *J. Clin. Microbiol.* **51**, 2884–  
1256 2892 (2013).
- 1257 74. Vesth, T. *et al.* Veillonella, Firmicutes: Microbes disguised as Gram negatives. *Stand.*  
1258 *Genomic Sci.* **9**, (2013).
- 1259 75. Bouwknegt, M., van Pelt, W., Kubbinga, M., Weda, M. & Havelaar, A. Potential  
1260 association between the recent increase in campylobacteriosis incidence in the  
1261 Netherlands and proton-pump inhibitor use – an ecological study. *Eurosurveillance* **19**,  
1262 20873 (2014).
- 1263 76. Leonard, J., Marshall, J. K. & Moayyedi, P. Systematic review of the risk of enteric  
1264 infection in patients taking acid suppression. *Am. J. Gastroenterol.* **102**, 2047–2056  
1265 (2007).
- 1266 77. Allaker, R. P. Non-sporing anaerobes: Wound infection; periodontal disease; abscess;  
1267 normal flora. *Med. Microbiol. Eighteenth Ed.* 359–364 (2012). doi:10.1016/B978-0-  
1268 7020-4089-4.00051-2
- 1269 78. Eribe, E. R. K. & Olsen, I. Leptotrichia species in human infections II. *J. Oral*  
1270 *Microbiol.* **9**, 1368848 (2017).
- 1271 79. Liu, D. *Molecular detection of human bacterial pathogens.* (CRC press, 2011).
- 1272 80. Carlier, J.-P. *Oribacterium.* in *Bergey's Manual of Systematics of Archaea and*

- 1273 *Bacteria* 1–5 (John Wiley & Sons, Ltd, 2015). doi:10.1002/9781118960608.gbm00649
- 1274 81. Wang, K. *et al.* Preliminary analysis of salivary microbiome and their potential roles in  
1275 oral lichen planus. *Sci. Rep.* **6**, 22943 (2016).
- 1276 82. Torok, E., Moran, E. & Cooke, F. *Oxford Handbook of Infectious Diseases and*  
1277 *Microbiology*. (Oxford University Press, 2009).  
1278 doi:10.1093/med/9780198569251.001.0001
- 1279 83. Jolivet-Gougeon, A., Sixou, J.-L., Tamanai-Shacoori, Z. & Bonnaure-Mallet, M.  
1280 Antimicrobial treatment of Capnocytophaga infections. *Int. J. Antimicrob. Agents* **29**,  
1281 367–373 (2007).
- 1282 84. Piau, C., Arvieux, C., Bonnaure-Mallet, M. & Jolivet-Gougeon, A. Capnocytophaga  
1283 spp. involvement in bone infections: a review. *Int. J. Antimicrob. Agents* **41**, 509–515  
1284 (2013).
- 1285 85. Cargill, J. S., Scott, K. S., Gascoyne-Binzi, D. & Sandoe, J. A. T. Granulicatella  
1286 infection: Diagnosis and management. *J. Med. Microbiol.* **61**, 755–761 (2012).
- 1287 86. Hofstad, T. The Genus Fusobacterium. in *The Prokaryotes* 1016–1027 (Springer New  
1288 York, 2006). doi:10.1007/0-387-30747-8
- 1289 87. Brophy, S. *et al.* Incidence of Campylobacter and Salmonella Infections Following  
1290 First Prescription for PPI: A Cohort Study Using Routine Data. *Am. J. Gastroenterol.*  
1291 **108**, 1094–1100 (2013).
- 1292 88. Allos, B. M. Campylobacter infections. in *Bacterial Infections of Humans:*  
1293 *Epidemiology and Control* 189–211 (Springer US, 2009). doi:10.1007/978-0-387-  
1294 09843-2\_9
- 1295 89. Lee, C. & Hong, S. N. Does long-term proton pump inhibitor therapy affect the health  
1296 of gut microbiota? *Gut and Liver* **10**, 865–866 (2016).
- 1297 90. Seto, C. T., Jeraldo, P., Orenstein, R., Chia, N. & DiBaise, J. K. Prolonged use of a  
1298 proton pump inhibitor reduces microbial diversity: Implications for Clostridium

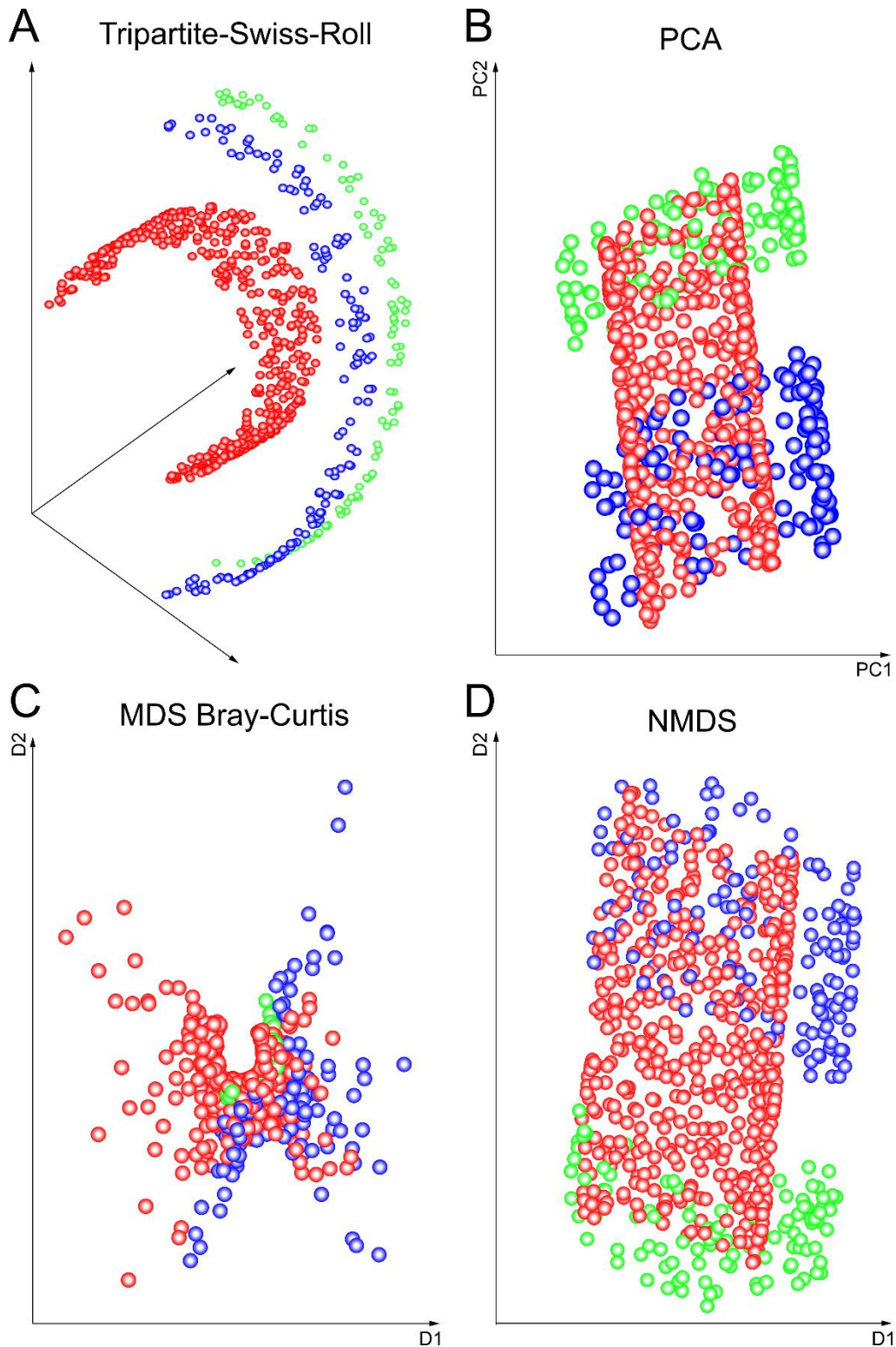
- 1299 difficile susceptibility. *Microbiome* **2**, (2014).
- 1300 91. Bavishi, C. & DuPont, H. L. Systematic review: The use of proton pump inhibitors and  
1301 increased susceptibility to enteric infection. *Alimentary Pharmacology and*  
1302 *Therapeutics* **34**, 1269–1281 (2011).
- 1303 92. Olbe, L. *Proton pump inhibitors*. (Birkhäuser, 2012).
- 1304 93. Warren, J. R. & Marshall, B. Unidentified curved bacilli on gastric epithelium in active  
1305 chronic gastritis. *Lancet* **321**, 1273–1275 (1983).
- 1306 94. Ha, N. *et al.* Supramolecular assembly and acid resistance of *Helicobacter pylori*  
1307 urease. *Nat. Struct. Biol.* **8**, 505–509 (2001).
- 1308 95. Berger, A. Scientists discover how helicobacter survives gastric acid. *Br. Med. J.* **29**,  
1309 268 (2000).
- 1310 96. Amieva, M. R. & El-Omar, E. M. Host-Bacterial Interactions in *Helicobacter pylori*  
1311 Infection. *Gastroenterology* **134**, 306–323 (2008).
- 1312 97. Scott Merrell, D. *et al.* Adhesion and Invasion of Gastric Mucosa Epithelial Cells by  
1313 *Helicobacter pylori*. *Front. Cell. Infect. Microbiol* **6**, 1593389–159 (2016).
- 1314 98. von Rosenvinge, E. C. *et al.* Immune status, antibiotic medication and pH are  
1315 associated with changes in the stomach fluid microbiota. *ISME J.* **7**, 1354–1366 (2013).
- 1316 99. Eun, C. S. o. *et al.* Differences in gastric mucosal microbiota profiling in patients with  
1317 chronic gastritis, intestinal metaplasia, and gastric cancer using pyrosequencing  
1318 methods. *Helicobacter* **19**, 407–416 (2014).
- 1319 100. Cao, L. & Yu, J. Effect of *Helicobacter pylori* Infection on the Composition of Gastric  
1320 Microbiota in the Development of Gastric Cancer. *Gastrointest. tumors* **2**, 14–25  
1321 (2015).
- 1322 101. Brawner, K. M., Morrow, C. D. & Smith, P. D. Gastric microbiome and gastric cancer.  
1323 *Cancer J.* **20**, 211–6 (2014).
- 1324 102. Cover, T. L. & Blaser, M. J. *Helicobacter pylori* in health and disease.

- 1325 *Gastroenterology* **136**, 1863–73 (2009).
- 1326 103. Sanders, M. K. & Peura, D. A. Helicobacter pylori-Associated Diseases. *Curr.*  
1327 *Gastroenterol. Rep.* **4**, 448–54 (2002).
- 1328 104. Talley, N. J. Helicobacter pylori and dyspepsia. *Yale J. Biol. Med.* **72**, 145–51 (1999).
- 1329 105. Shadwell, J. Helicobacter pylori–associated dyspepsia. *2016*
- 1330 106. Noto, J. M. & Peek, R. M. The gastric microbiome, its interaction with Helicobacter  
1331 pylori, and its potential role in the progression to stomach cancer. *PLoS Pathogens* **13**,  
1332 (2017).
- 1333 107. Schwabe, R. F. & Jobin, C. The microbiome and cancer. *Nature Reviews Cancer* **13**,  
1334 800–812 (2013).
- 1335 108. Fraher, M. H., O’Toole, P. W. & Quigley, E. M. M. Techniques used to characterize  
1336 the gut microbiota: a guide for the clinician. *Nat. Rev. Gastroenterol. Hepatol.* **9**, 312–  
1337 322 (2012).
- 1338 109. Andersson, A. F. *et al.* Comparative Analysis of Human Gut Microbiota by Barcoded  
1339 Pyrosequencing. *PLoS One* **3**, e2836 (2008).
- 1340 110. Bik, E. M. Molecular analysis of the bacterial microbiota in the human stomach. *Proc.*  
1341 *Natl. Acad. Sci. USA* **103**, 732–737 (2006).
- 1342 111. Llorca, L. *et al.* Characterization of the gastric microbiota in a pediatric population  
1343 according to Helicobacter pylori status. in *Pediatric Infectious Disease Journal* **36**,  
1344 173–178 (2017).
- 1345 112. Jo, H. J. The effect of H. pylori infection on the gastric microbiota. in *Helicobacter*  
1346 *pylori* (ed. Kim, N.) 529–533 (Springer Singapore, 2016). doi:10.1007/978-981-287-  
1347 706-2\_54
- 1348 113. Klymiuk, I. *et al.* The Human Gastric Microbiome Is Predicated upon Infection with  
1349 Helicobacter pylori. *Front. Microbiol.* **8**, 2508 (2017).
- 1350 114. Maldonado-Contreras, A. *et al.* Structure of the human gastric bacterial community in

- 1351 relation to *Helicobacter pylori* status. *ISME J.* **5**, 574–579 (2011).
- 1352 115. Aviles-Jimenez, F., Vazquez-Jimenez, F., Medrano-Guzman, R., Mantilla, A. &  
1353 Torres, J. Stomach microbiota composition varies between patients with non-atrophic  
1354 gastritis and patients with intestinal type of gastric cancer. *Sci. Rep.* **4**, 4202 (2015).
- 1355 116. Kovaleva, J., Degener, J. E. & van der Mei, H. C. Methylobacterium and its role in  
1356 health care-associated infection. *J. Clin. Microbiol.* **52**, 1317–21 (2014).
- 1357 117. White, D. C., Sutton, S. D. & Ringelberg, D. B. The genus *Sphingomonas*: physiology  
1358 and ecology. *Curr. Opin. Biotechnol.* **7**, 301–306 (1996).
- 1359 118. Madigan, M., Martinko, J., Stahl, D. and Clark, D. Brock Biology of Microorganisms.  
1360 321 (2012).
- 1361 119. Özen, A. I. & Ussery, D. W. Defining the *Pseudomonas* genus: where do we draw the  
1362 line with *Azotobacter*? *Microb. Ecol.* **63**, 239–48 (2012).
- 1363 120. Towner, K. The genus *Acinetobacter*. in *The Prokaryotes* 545–577 (Springer New  
1364 York, 2006). doi:10.1007/978-3-642-30194-0
- 1365 121. Rathinavelu, S., Zavros, Y. & Merchant, J. L. *Acinetobacter lwoffii* infection and  
1366 gastritis. *Microbes Infect.* **5**, 651–657 (2003).
- 1367 122. Cheung, Y. F., Walsh, C. & Fung, C. H. Stereochemistry of Propionyl-Coenzyme A  
1368 and Pyruvate Carboxylations Catalyzed by Transcarboxylase. *Biochemistry* **14**, 2981–  
1369 2986 (1975).
- 1370 123. Piwowarek, K., Lipińska, E., Hać-Szymańczuk, E., Kieliszek, M. & Ścibisz, I.  
1371 *Propionibacterium* spp.—source of propionic acid, vitamin B12, and other metabolites  
1372 important for the industry. *Applied Microbiology and Biotechnology* **102**, 515–538  
1373 (2018).
- 1374 124. Moore, L. V. H. & Moore, W. E. C. *Oribaculum catoniae* gen. nov., sp. nov.; *Catonella*  
1375 *morbi* gen. nov., sp. nov.; *Hallella seregens* gen. nov., sp. nov.; *Johnsonella ignava* gen.  
1376 nov., sp. nov.; and *Dialister pneumosintes* gen. nov., comb. nov., nom. rev., Anaerobic

- 1377 Gram-Negative Bacilli from. *Int. J. Syst. Bacteriol.* **44**, 187–192 (1994).
- 1378 125. Willems, A. & Collins, M. D. *Catonella*. in *Bergey's Manual of Systematics of*  
1379 *Archaea and Bacteria* 1–7 (John Wiley & Sons, Ltd, 2015).  
1380 doi:10.1002/9781118960608.gbm00641
- 1381 126. Menon, T. & Kumar, V. N. *Catonella morbi* as a cause of native valve endocarditis in  
1382 Chennai, India. *Infection* **40**, 581–582 (2012).
- 1383 127. Balows, A., Truper, H., Dvorkin, M., Harder, W. & Schleifer, K. *The Prokaryotes. A*  
1384 *Handbook on the Biology of Bacteria: Proteobacteria: Gamma subclass. The*  
1385 *prokaryotes* (Springer, 1991). doi:10.1007/0-387-30745-1
- 1386 128. Staley, J. T., Irgens, R. L. & Brenner, D. J. *Enhydrobacter aerosaccus* gen. nov., sp.  
1387 nov., a Gas-Vacuolated, Facultatively Anaerobic, Heterotrophic Rod. *Int. J. Syst.*  
1388 *Bacteriol.* **37**, 289–291 (1987).
- 1389 129. Wade, W. G. & Downes, J. *Bulleidia*. *Bergey's Manual of Systematics of Archaea and*  
1390 *Bacteria* (2015). doi:doi:10.1002/9781118960608.gbm00760
- 1391 130. Kienesberger, S. *et al.* Gastric *Helicobacter pylori* Infection Affects Local and Distant  
1392 Microbial Populations and Host Responses. *Cell Rep.* **14**, 1395–1407 (2016).  
1393  
1394

1395 **Figures and tables**



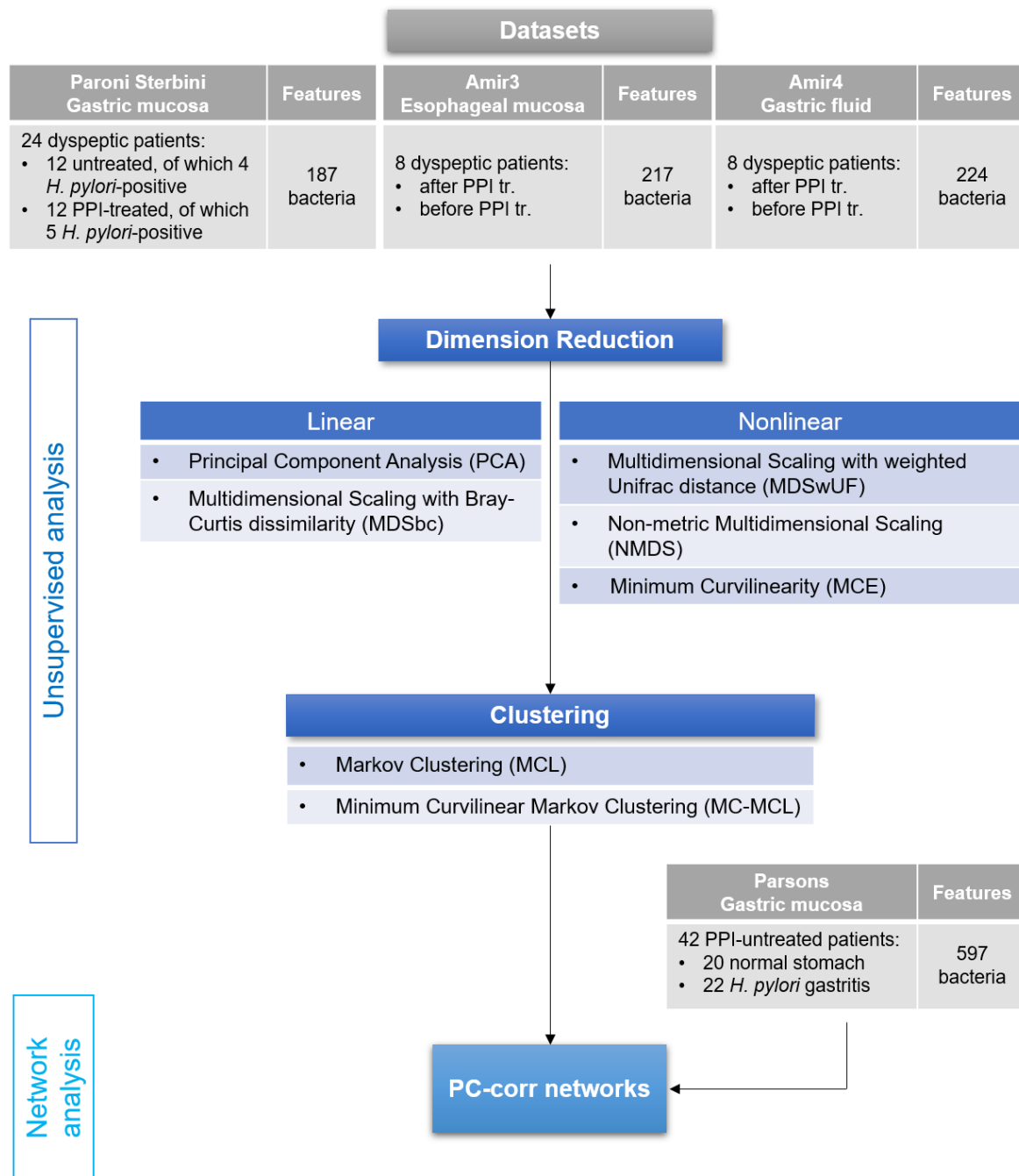
1396

1397 **Figure 1. The Tripartite-Swiss-Roll as an example of data nonlinear organization.**

1398 A) Tripartite-Swiss-Roll; B) PCA; C) MDS (Bray-Curtis dissimilarity); D) NMDS (Sammon Mapping).

1399 The three different colours (red, blue and green) represent the three partitions of the Swiss-roll manifold.

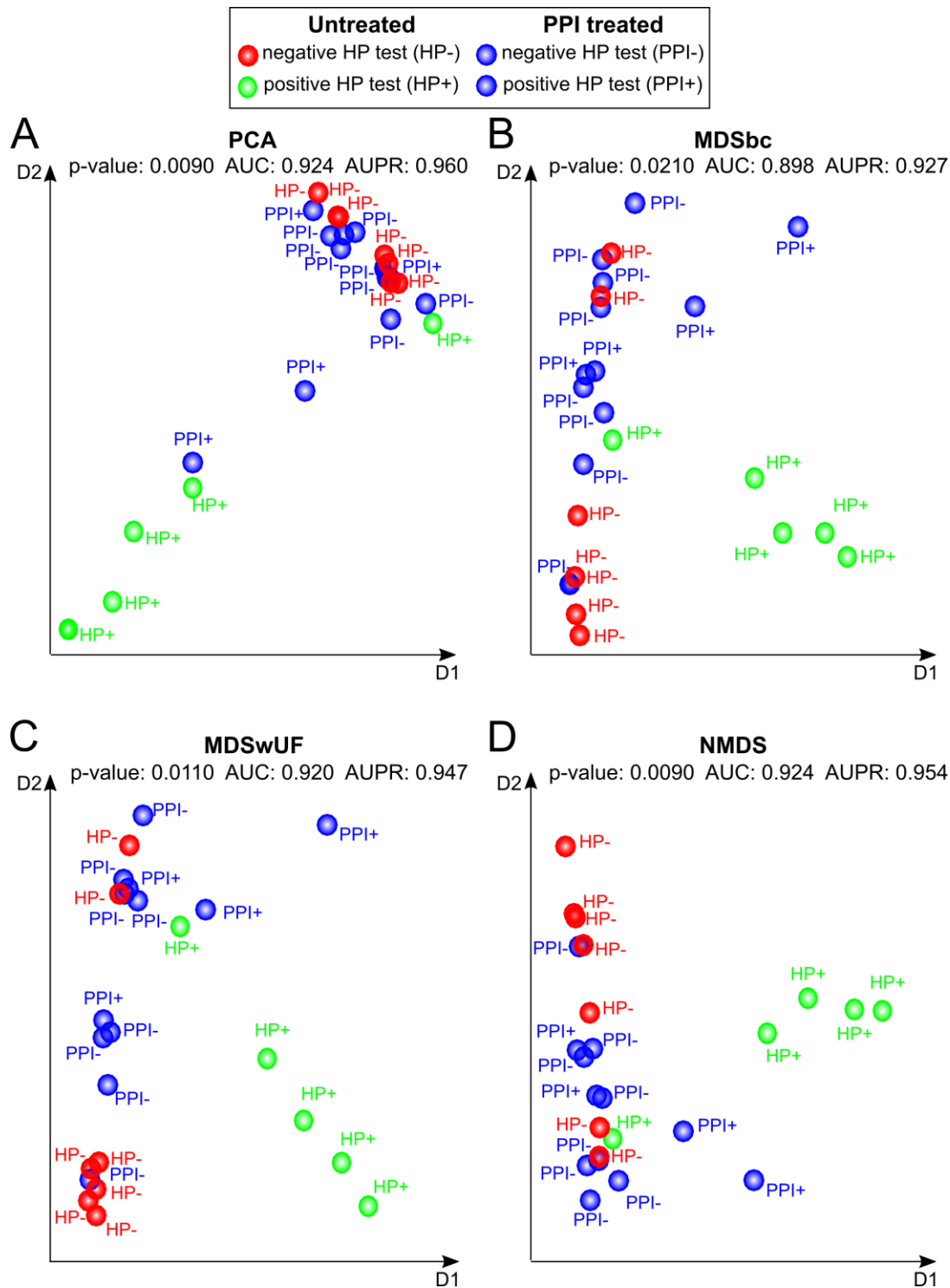
1400 This figure shows the inability of PCA, MDS and NMDS to reveal the inner nonlinear structure of the  
 1401 Tripartite-Swiss-Roll, which appears collapsed (B, D) or with a horseshoe shape (C).  
 1402



1403 **Figure 2. Flowchart of the data analysis.** To answer the five questions under investigation in our study,  
 1404 we implemented a workflow based on machine learning tools. Following the flowchart shown in the  
 1405 figure, we analysed three 16S rRNA gene sequencing datasets with information on PPI use in dyspeptic  
 1406 patients; for one of the datasets (Paroni Sterbini *et al.* <sup>22</sup>), patients were also determined to be positive  
 1407 or negative to *H. pylori* infection.  
 1408



1409 Firstly, we performed unsupervised dimension reduction, both linear and nonlinear, in the first two  
1410 dimensions of embedding. Nonlinear dimension reduction will show the presence of hidden patterns,  
1411 in the form of sample groups. Secondly, nonlinear clustering was applied to confirm the well-  
1412 possessedness of the hidden patterns found by nonlinear dimension reduction. Lastly, our workflow ends  
1413 with the PC-corr algorithm, that reveals which combination of bacteria (features) are responsible for the  
1414 identified differences between the groups of samples. A fourth dataset (Parsons *et al*<sup>29</sup>.) is used only for  
1415 the validation of the PC-corr network results and it contains information of PPI treatment and *H. pylori*  
1416 infection.



1417

1418 **Figure 3. Dimension reduction techniques usually employed in metagenomic data analysis and**

1419 **applied to the Paroni Sterbini dataset.** The plots represent the best PCA and MDS results based on

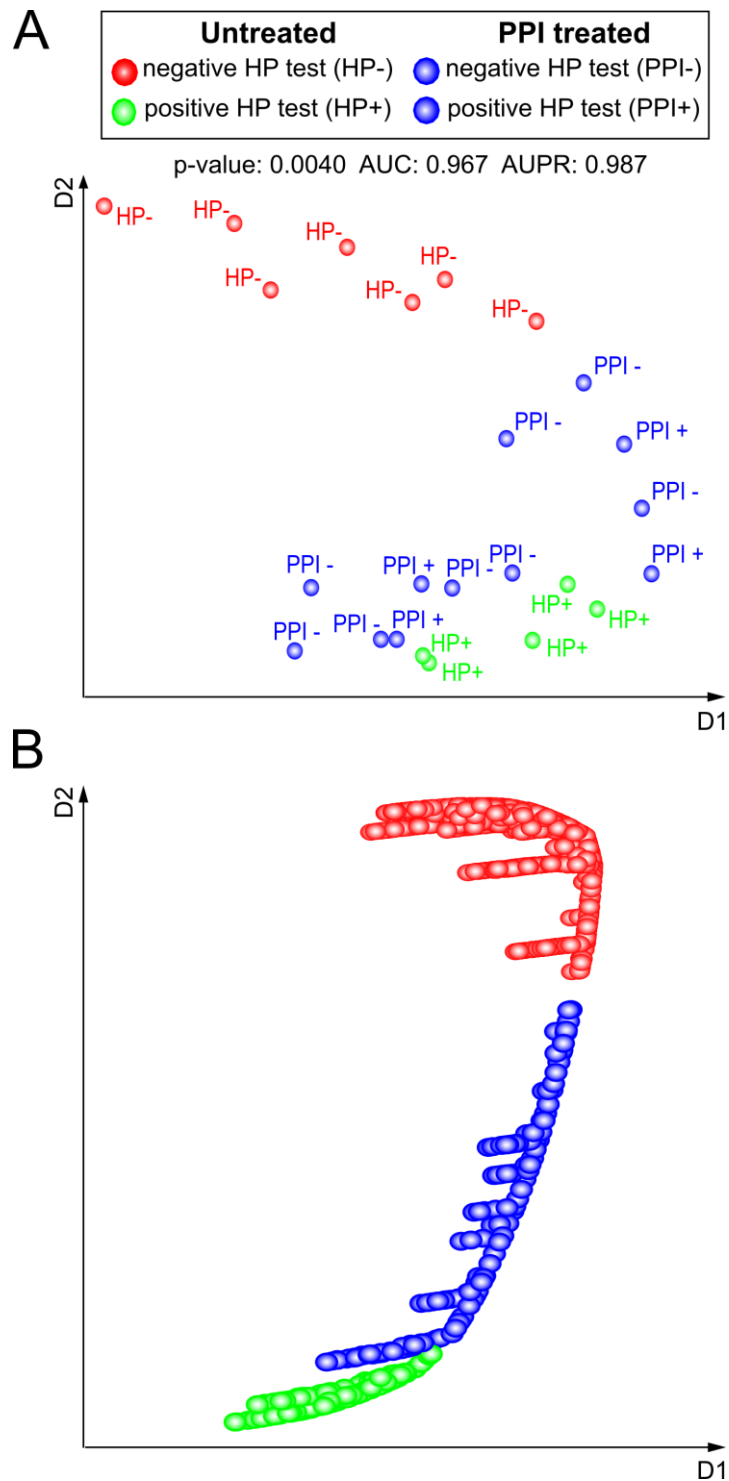
1420 (average) p-value projection-based separability index (PSI) for the three different labels (PPI-treated,

1421 untreated HP+ and untreated HP-), evaluated in the 2D embedding space. Moreover, also the average

1422 values of all pairwise AUC and AUPR PSI are reported as overall estimators of separation between the

1423 groups in the 2D reduced space. A) PCA; B) MDS with Bray-Curtis dissimilarity (MDSbc); C) MDS

1424 with weighted UniFrac distance (MDSwUF); D) non-metric MDS with Sammon Mapping (NMDS).  
1425 Blue dots represent PPI-treated samples, while red and green dots are the untreated samples which  
1426 resulted either negative (red) or positive (green) to the *H. pylori* test (histological observation and urease  
1427 test).



**Figure 4. MCE, a topological machine learning for nonlinear and hierarchical dimension reduction.** (A) Results on the Paroni Sterbini *et al.*<sup>22</sup> dataset. The shown best MCE result is based on (average) p-value projection-based separability index (PSI) for the three different labels (PPI-treated, untreated HP+ and untreated HP-), evaluated in the 2D embedding space under the DCS normalization. The average values of all pairwise AUC and AUPR PSI are reported as well as overall estimators of separation between the groups in the 2D reduced space. Blue dots represent PPI-treated samples, while

red and green dots are the untreated samples which resulted either negative (red) or positive (green) to the *H. pylori* test (histological observation and urease test). **(B)** Results on the Tripartite-Swiss Roll. The three different colours (red, blue and green) represent the three partitions of the Swiss-roll manifold.

**Table 1. Results of unsupervised analysis on the original 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 index 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 HP+ and untreated HP-). For the Amir datasets, the p-values were computed for two groups, identified by the presence or absence of PPI treatment.

**(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 HP+ and untreated HP-) and in brackets the results for four clusters (PPI-treated HP+, PPI-treated HP-, untreated HP+ and untreated HP-). Instead, for Amir datasets, the accuracies were computed for two groups,

identified according to the presence or absence of PPI treatment.

Dimension Reduction	P-value				
	Method	Paroni Sterbini	Amir3	Amir4	mean
	HD	0.0056	0.0011	0.0003	0.0023
	MCE	0.0040	0.0078	0.0047	0.0055
	MDSwUF	0.0110	0.0002	0.0104	0.0072
	PCA	0.0090	0.0047	0.0148	0.0095
	NMDS	0.0090	0.0148	0.0207	0.0148
MDSbc	0.0210	0.0148	0.0207	0.0188	

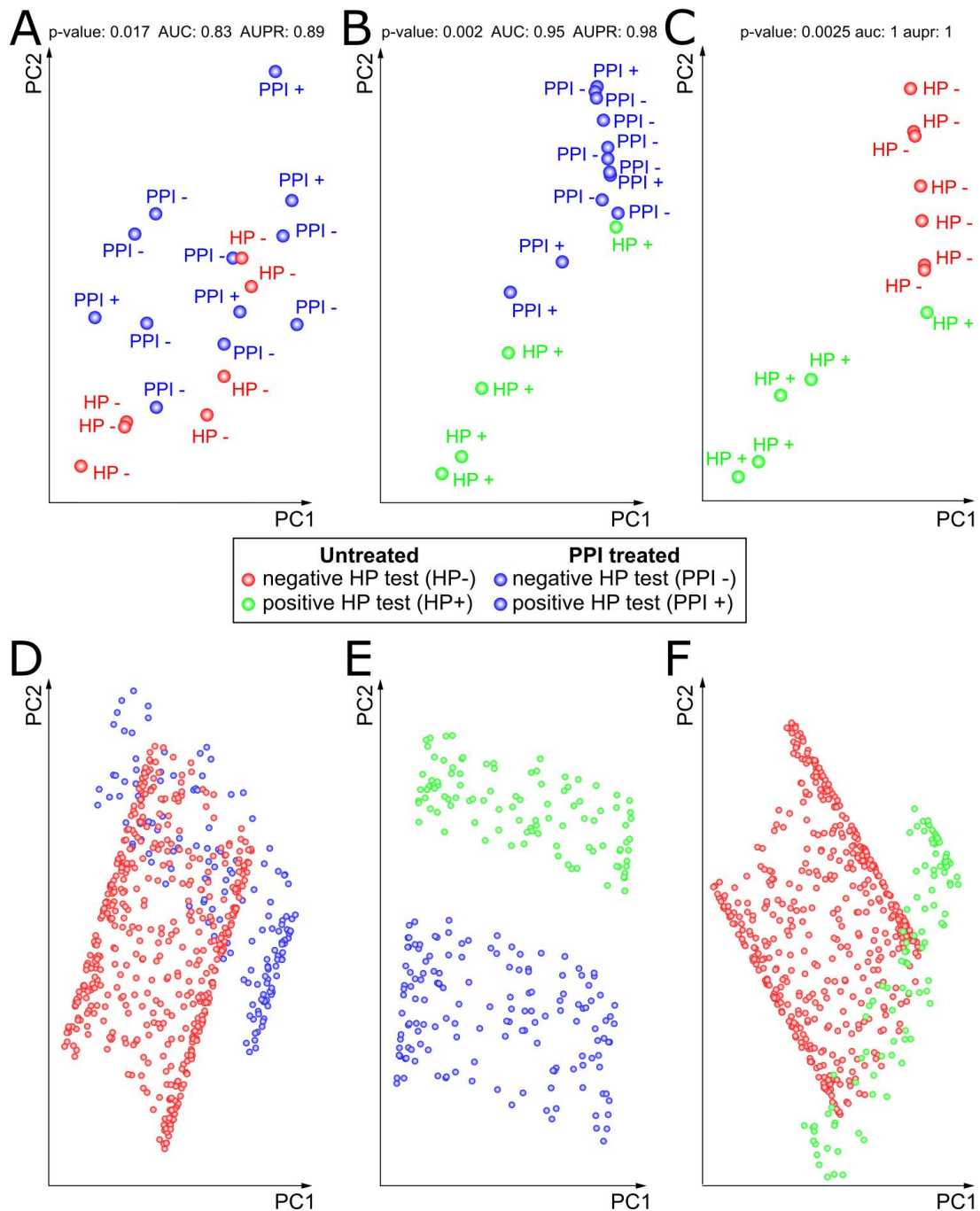
Dimension Reduction	AUC				AUPR					
	Method	Paroni Sterbini	Amir3	Amir4	mean	Method	Paroni Sterbini	Amir3	Amir4	mean
	HD	0.937	0.953	0.984	0.958	HD	0.967	0.957	0.985	0.970
	MDSwUF	0.920	1.000	0.875	0.932	MDSwUF	0.946	1.000	0.901	0.949
	MCE	0.967	0.883	0.906	0.919	MCE	0.987	0.891	0.920	0.933
	PCA	0.924	0.906	0.859	0.896	PCA	0.959	0.902	0.880	0.914
	NMDS	0.924	0.859	0.844	0.876	MDSbc	0.927	0.891	0.900	0.906
MDSbc	0.898	0.859	0.844	0.867	NMDS	0.954	0.871	0.873	0.899	

Clustering	Accuracy	Paroni Sterbini	Amir3	Amir4	Mean performance
	MC-MCL	0.67 (0.58)	0.81	0.75	0.74
	MCL	0.58 (0)	0.69	0.75	0.67

Note: all the P-values, AUC and AUPR can be found in Supplementary Table S1, while all the accuracies can be found in Supplementary Table S4.

1428 **Abbreviations:** HD: High Dimension; MCE: Minimum Curvilinear Embedding; MDSbc:  
 1429 Multidimensional Scaling with Bray-Curtis dissimilarity; MDSwUF: Multidimensional Scaling with  
 1430 weighted UniFrac distance; NMDS: Non-metric Multidimensional Scaling; PCA: Principal Component  
 1431 Analysis; MCL: Markov Clustering; MC-MCL: Minimum Curvilinear Markov Clustering; p-value:  
 1432 Mann-Whitney p-value; AUC: Area Under the Curve; AUPR: Area Under the Precision Recall.



**Figure 5. Pairwise PCA of Paroni Sterbini's gastric samples and of the Tripartite-Swiss-Roll. (A-C)** 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 p-value, AUC and AUPR PSI are reported as well as overall estimators of separation between the groups in the 2D reduced space (**D-F**) In a similar manner, PCA was applied to the three datasets obtained by subsetting

the Swiss-roll dataset, each one corresponding to a combination of two groups: D) red vs blue groups, E) blue vs green groups, F) red vs green groups.

**Table 2. Ranked performance of unsupervised dimension reduction techniques on the original datasets.** The table shows the ranked performance of unsupervised dimension reduction techniques according to the index for sample separation (based on Mann-Whitney P-value, AUC and AUPR) 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.

P-value

Method	Paroni Sterbini	Amir3	Amir4	mean
HD	2	2	1	1.67
MCE	1	4	2	2.33
MDSwUF	4	1	3	2.67
PCA	3	3	4	3.33
NMDS	3	5	5	4.33
MDSbc	5	5	5	5

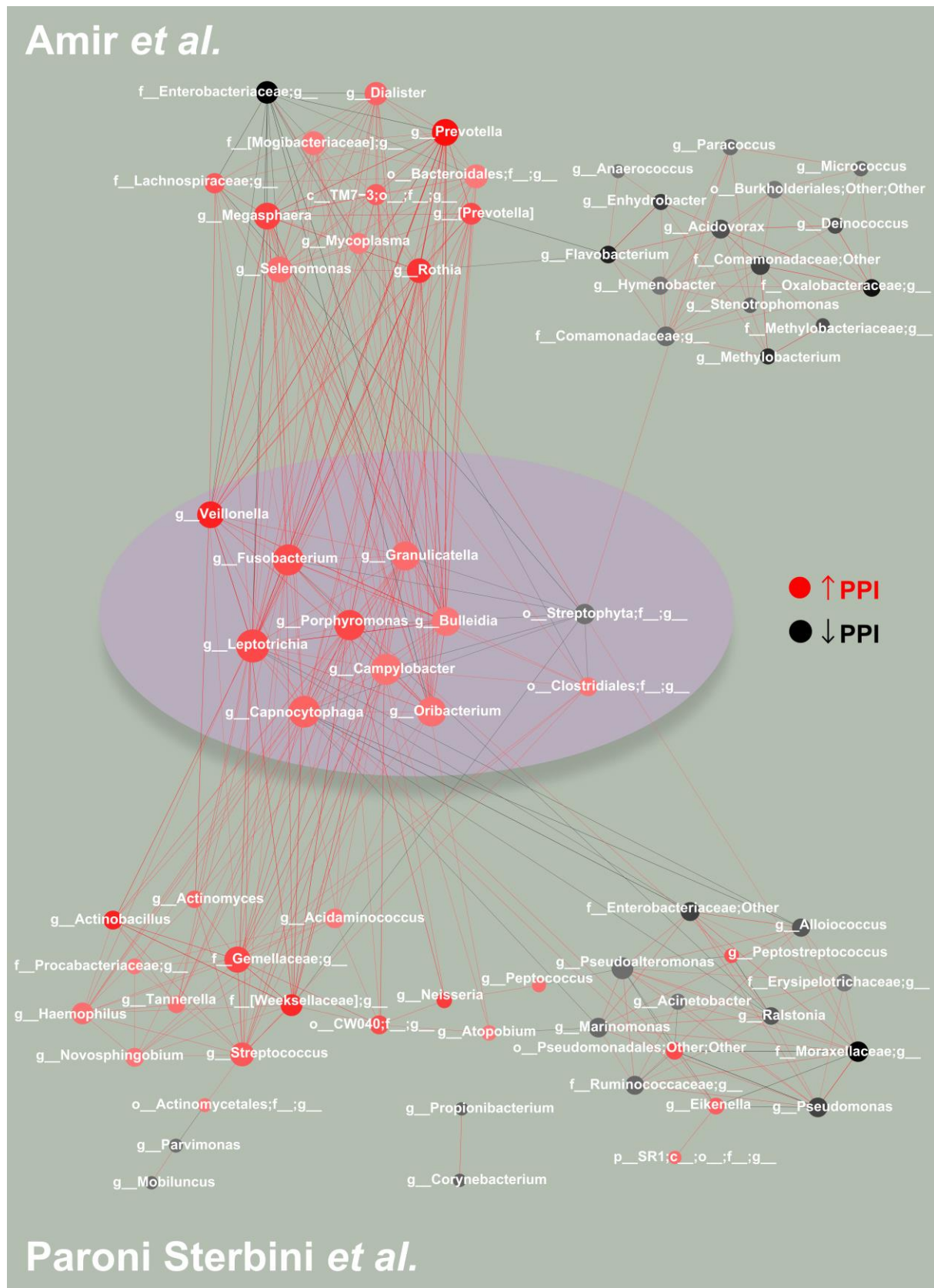
AUC

AUPR

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	4	2	2.33
MDSwUF	4	1	3	2.67	MDSwUF	5	1	3	3
PCA	3	3	4	3.33	PCA	3	3	5	3.67
NMDS	3	5	5	4.33	MDSbc	6	4	4	4.67
MDSbc	5	5	5	5	NMDS	4	5	6	5

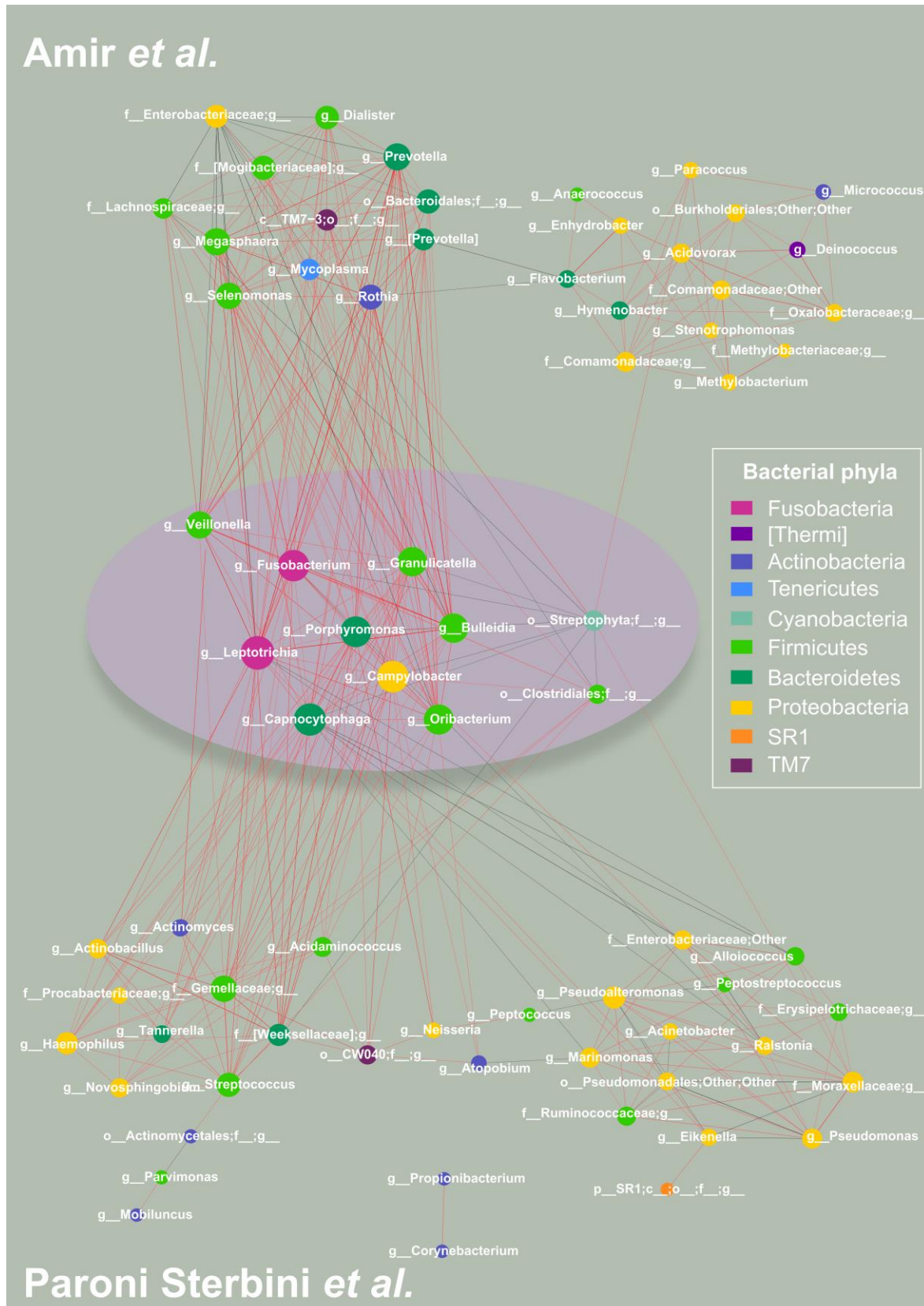
1433 **Abbreviations:** HD: High Dimension; MCE: Minimum Curvilinear Embedding; MDSbc:  
 1434 Multidimensional Scaling with Bray-Curtis dissimilarity; MDSwUF: Multidimensional Scaling with  
 1435 weighted UniFrac distance; NMDS: Non-metric Multidimensional Scaling; PCA: Principal Component  
 1436 Analysis; p-value: Mann-Whitney p-value; AUC: Area Under the Curve; AUPR: Area Under the  
 1437 Precision Recall.





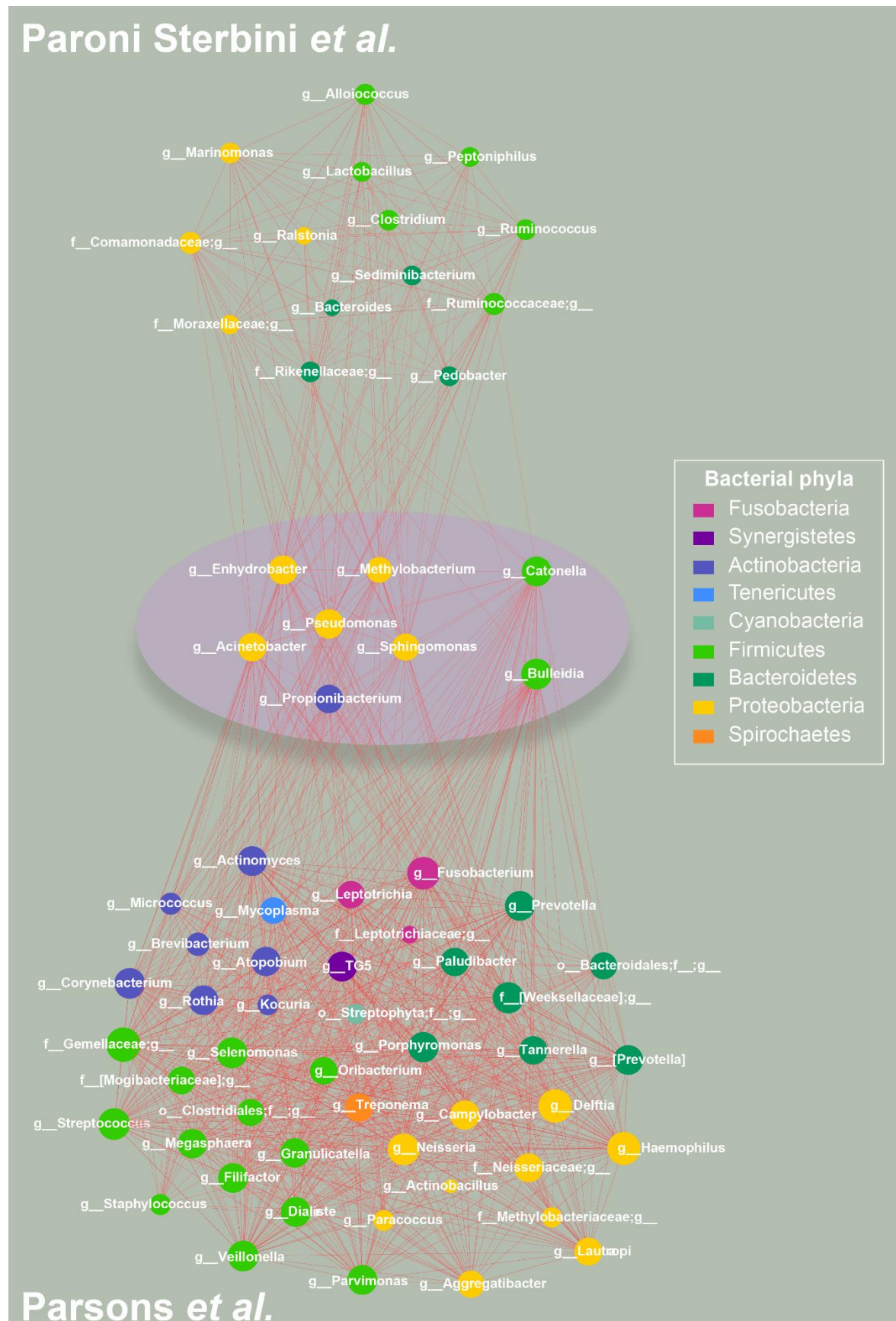
1438 **Figure 6. PC-corr method to unveil how PPI is affecting the microbiota in gastric environment in**  
 1439 **dyspeptic patients. (Middle panel)** To investigate the effect of PPIs on the gastric microbiota in  
 1440 dyspeptic patients, we constructed the conserved PC-corr network at 0.5 cut-off, by merging the PC-

1441 corr networks obtained from the gastric mucosa (Paroni Sterbini *et al.*<sup>22</sup>) and the gastric fluid (Amir *et*  
1442 *al.*<sup>21</sup>). To do so, we firstly considered the union of the two PC-corr networks obtained from the gastric  
1443 tissue dataset and then we intersected it with the PC-corr network from the gastric fluid dataset. All the  
1444 bacteria spotted in the conserved PC-corr network (violet circle) were found increased with PPI use. In  
1445 both the two studied datasets, red nodes indicate bacteria whose abundance is increased with PPI-  
1446 treatment, while black nodes indicate bacteria with lower abundance following treatment with this acid  
1447 suppressing medication. The common bacteria that showed an opposite trend in the two datasets, i.e.  
1448 microbial abundance increased in one dataset and decreased in the other dataset, were removed from the  
1449 network. (**Top panel**) The top panel shows the obtained Amir4's network, not in common with the  
1450 Paroni Sterbini's network. The module on the left side (except *Enterobacteriaceae*) include bacteria  
1451 more abundant following PPI-treatment in Amir4's data, while the module on the right (and  
1452 *Enterobacteriaceae*) is composed of decreased bacteria in abundance under PPI therapy in Amir4's data.  
1453 (**Bottom panel**) The bottom panel represents the part of Paroni Sterbini's network (union of the two  
1454 PC-corr network), that is not shared with Amir4's one. As in the top and middle panels, the colour of  
1455 the nodes represents if the bacteria display higher (red nodes) or lower abundance (black nodes) in PPI-  
1456 treated samples of Paroni Sterbini's dataset.



1457 **Figure 7. PC-corr networks to unveil how PPI is affecting the microbiota in gastric environment**  
 1458 **in dyspeptic patients, coloured according to phylum-level taxonomy.** To investigate the effect of  
 1459 PPIs on the gastric microbiota in dyspeptic patients, we constructed the conserved PC-corr network at

1460 0.5 cut-off, by merging the PC-corr networks obtained from the gastric mucosa (Paroni Sterbini *et al.*  
1461 <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  
1462 networks obtained from the gastric tissue dataset and then we intersected it with the PC-corr network  
1463 from the gastric fluid dataset. All the bacteria spotted in the conserved PC-corr network (violet circle)  
1464 were found increased with PPI use. (**Top panel**) The top panel shows the obtained Amir4's network,  
1465 not in common with the Paroni Sterbini's network. The module on the left side (except  
1466 *Enterobacteriaceae*) include bacteria more abundant following PPI-treatment in Amir4's data, while the  
1467 module on the right (and *Enterobacteriaceae*) is composed of decreased bacteria in abundance under PPI  
1468 therapy in Amir4's data. (**Bottom panel**) The bottom panel represents the part of Paroni Sterbini's  
1469 network (union of the two PC-corr network), that is not shared with Amir4's one. As in the top and  
1470 middle panels, nodes are coloured according to bacterial phylum level.



1471

1472 **Figure 8. PC-corr network to investigate the effect of *H. pylori* infection on the gastric mucosal**  
 1473 **microbiota, coloured according to phylum-level taxonomy. (Middle panel) To investigate the effect**

1474 of *H. pylori* infection on the gastric mucosal microbiota, we constructed the conserved PC-corr network  
1475 at 0.5 cut-off, by intersecting the PC-corr networks obtained from Paroni Sterbini *et al.*<sup>22</sup> and Parsons  
1476 *et al.*<sup>29</sup> dataset. All the bacteria spotted in the conserved PC-corr network (violet circle) were found  
1477 decreased in abundance with *H. pylori* infection. The common bacteria that showed an opposite trend  
1478 in the two datasets, i.e. microbial abundance increased in one dataset and decreased in the other dataset,  
1479 were removed from the network. (**Top panel**) The top panel show the obtained Paroni Sterbini's  
1480 network, not in common with the Parsons's network. It contains all bacteria whose abundance is  
1481 decreased in *H. pylori*-positive patients in Paroni Sterbini *et al.* dataset. (**Bottom panel**) The bottom  
1482 panel represent the part of Parsons's network that is not shared with Paroni Sterbini's one. As in the top  
1483 and middle panels, it includes bacterial communities decreased in *H. pylori*-infected patients.