

# 1 Microbiotyping the sinonasal microbiome

2 Ahmed Bassiouni<sup>1</sup>, Sathish Paramasivan<sup>1</sup>, Arron Shiffer<sup>2</sup>, Matthew R Dillon<sup>2</sup>, Emily Cope<sup>2</sup>,

3 Clare Cooksley<sup>1</sup>, Mohammad Javed Ali<sup>3</sup>, Benjamin Bleier<sup>4</sup>, Claudio Callejas<sup>5</sup>, Marjolein E

4 Cornet<sup>6</sup>, Richard G Douglas<sup>7</sup>, Daniel Dutra<sup>8</sup>, Christos Georgalas<sup>6</sup>, Richard J Harvey<sup>9,10</sup>, Peter H

5 Hwang<sup>11</sup>, Amber U Luong<sup>12</sup>, Rodney J Schlosser<sup>13</sup>, Pongsakorn Tantilipikorn<sup>14</sup>, Marc A

6 Tewfik<sup>15</sup>, Sarah Vreugde<sup>1</sup>, Peter-John Wormald<sup>1</sup>, J Gregory Caporaso<sup>2</sup>, and Alkis J Psaltis<sup>1</sup>

7 <sup>1</sup> Department of Otolaryngology, Head and Neck Surgery, University of Adelaide, Adelaide, Australia

8 <sup>2</sup> Pathogen and Microbiome Institute, Northern Arizona University, Arizona, USA

9 <sup>3</sup> Dacryology Service, LV Prasad Institute, Hyderabad, India

10 <sup>4</sup> Department of Otolaryngology, Massachusetts Eye and Ear Infirmary, Harvard Medical School, Boston, USA

11 <sup>5</sup> Department of Otolaryngology, Pontificia Universidad Catolica de Chile, Santiago, Chile

12 <sup>6</sup> Department of Otorhinolaryngology, Amsterdam UMC, Amsterdam, The Netherlands

13 <sup>7</sup> Department of Surgery, University of Auckland, Auckland, New Zealand

14 <sup>8</sup> Department of Otorhinolaryngology, University of Sao Paulo, Sao Paulo, Brazil

15 <sup>9</sup> Department of Otolaryngology, Rhinology and Skull base, University of New South Wales, Sydney, Australia

16 <sup>10</sup> Faculty of Medicine and Health sciences, Macquarie University, Sydney, Australia

17 <sup>11</sup> Department of Otolaryngology -Head and Neck Surgery, Stanford University, Stanford, California, USA

18 <sup>12</sup> Department of Otolaryngology -Head and Neck Surgery, University of Texas, Texas, USA

19 <sup>13</sup> Department of Otolaryngology, Medical University of South Carolina, Charleston, South Carolina, USA

20 <sup>14</sup> Department of Otorhinolaryngology, Faculty of Medicine, Siriraj Hospital, Mahidol University, Bangkok, Thailand

21 <sup>15</sup> Department of Otolaryngology - Head and Neck Surgery, McGill University, Montreal, Canada

## 22 **Corresponding author:**

23 Associate Professor Alkis J Psaltis

24 3C (Department of Otolaryngology, Head and Neck Surgery)

25 The Queen Elizabeth Hospital

26 28 Woodville Rd

27 Woodville South, SA 5011

28 Australia

29 Email: [alkis.psaltis@adelaide.edu.au](mailto:alkis.psaltis@adelaide.edu.au)

30 Phone: +61 08 8222 7158

31 Fax: +61 08 8222 7419

32 **Funding information, Disclosures and Conflicts of Interest (COI):**

33 Mohammad Javed Ali:

34 Receives royalties from Springer for his treatise “Principles and Practice of Lacrimal Surgery” and “Atlas  
35 of Lacrimal Drainage Disorders”.

36 No conflict of interest relevant to this study.

37 Ahmed Bassiouni, Clare Cooksley:

38 No conflict of interest to declare.

39 Benjamin Bleier:

40 Grant Funding: R01 NS108968-01 NIH/NINDS (Bleier PI) – This isn’t relevant to this study.

41 Consultant for: Gyrus ACMI Olympus, Canon, Karl Storz, Medtronic, and Sinopsys.

42 Equity: Cerebent, Inc, Arrinex.

43 COI: None relevant to this study.

44 Claudio Callejas:

45 No conflict of interest to declare.

46 J Gregory Caporaso, Matthew R Dillon, Arron Shiffer:

47 No conflicts of interest to declare. This work was funded in part by National Science Foundation Award  
48 1565100 to JGC.

49 Emily Cope:

50 Financial information: This work was partially funded under the State of Arizona Technology and  
51 Research Initiative Fund (TRIF), administered by the Arizona Board of Regents, through Northern  
52 Arizona University.

53 No relevant disclosures or COI.

54 Marjolein E Cornet:

55 No financial relationships or sponsors. No conflicts of interests.

56 Richard G Douglas:

57 Received consultancy fees from Lyra Therapeutics and is a consultant for Medtronic. These are not  
58 relevant to this study.

59 Daniel Dutra:

60 No conflict of interest to declare.

61 Richard J Harvey:

62 Consultant with Medtronic, Olympus and NeilMed pharmaceuticals. He has also been on the speakers'  
63 bureau for Glaxo-Smith-Kline, Seqiris and Astra-Zeneca.

64 No direct conflict of interest to declare.

65 Peter H Hwang:

66 Financial Relationships: Consultancies with Arrinex, Bioinspire, Canon, Lyra Therapeutics, Medtronic,  
67 Tivic.

68 Conflicts of Interest: None.

69 Amber U Luong:

70 Serves as a consultant for Aerin Medical (Sunnyvale, CA), Arrinex (Redwood City, CA), Lyra  
71 Therapeutics (Watertown, MA), and Stryker (Kalamazoo, MI) and is on the advisory board for

- 72 ENTvantage (Austin, TX).
- 73 Her department receives funding from Genentech/Roche (San Francisco, CA) and AstraZeneca  
74 (Cambridge, England).
- 75 No COI to declare related to this study.
- 76 Sathish Paramasivan:  
77 Supported by a Garnett Passe and Rodney Williams Memorial Foundation Academic Surgeon Scientist  
78 Research Scholarship.
- 79 No conflicts of interest to declare.
- 80 Alkis J Psaltis:  
81 Consultant for Aerin Devices and ENT technologies and is on the speakers' bureau for Smith and  
82 Nephew. Received consultancy fees from Lyra Therapeutics. These are not relevant to this study.
- 83 Rodney J Schlosser:  
84 Grant support from OptiNose, Entellus, and IntersectENT (not relevant to this study). Consultant for  
85 Olympus, Meda, and Arrinex (not relevant to this study).
- 86 Pongsakorn Tantilipikorn:  
87 No financial disclosures or conflict of interest.
- 88 Marc A Tewfik:  
89 Principal Investigator: Sanofi, Roche/Genentech, AstraZeneca.  
90 Speaker/Consultant: Stryker, Ondine Biomedical, Novartis, MEDA, Mylan.  
91 Royalties for book sales: Thieme.
- 92 Sarah Vreugde:  
93 No conflicts of interest relevant to this study.

94 Peter-John Wormald:

95 Receives royalties from Medtronic, Integra, and Scopis, and is a consultant for NeilMed. These are not  
96 relevant to this study.

97

## 98 **Abstract**

99 This study offers a novel description of the sinonasal microbiome, through an unsupervised machine  
100 learning approach combining dimensionality reduction and clustering. We apply our method to the  
101 International Sinonasal Microbiome Study (ISMS) dataset of 410 sinus swab samples. We propose three  
102 main sinonasal ‘microbiotypes’ or ‘states’: the first is *Corynebacterium*-dominated, the second is  
103 *Staphylococcus*-dominated, and the third dominated by the other core genera of the sinonasal microbiome  
104 (*Streptococcus*, *Haemophilus*, *Moraxella*, and *Pseudomonas*). The prevalence of the three microbiotypes  
105 studied did not differ between healthy and diseased sinuses, but differences in their distribution were  
106 evident based on geography. We also describe a potential reciprocal relationship between  
107 *Corynebacterium* species and *Staphylococcus aureus*, suggesting that a certain microbial equilibrium  
108 between various players is reached in the sinuses. We validate our approach by applying it to a separate  
109 16S rRNA gene sequence dataset of 97 sinus swabs from a different patient cohort. Sinonasal  
110 microbiotyping may prove useful in reducing the complexity of describing sinonasal microbiota. It may  
111 drive future studies aimed at modeling microbial interactions in the sinuses and in doing so may facilitate  
112 the development of a tailored patient-specific approach to the treatment of sinus disease in the future.

## 113 **Keywords**

114 microbiome, sinus, next-generation sequencing, 16S rRNA gene, chronic rhinosinusitis, microbiotype

115

## 116 MAIN TEXT

117 Chronic rhinosinusitis (CRS) is a heterogenous, multi-factorial inflammatory disorder with a complex and  
118 incompletely understood aetiopathogenesis.<sup>1</sup> A potential role of the sinonasal microbiome and its  
119 “dysbiosis” in CRS pathophysiology has recently gained increased interest. The nature of the microbial  
120 dysbiosis and its role in disease causation and progression however remains unclear, with conflicting  
121 findings from the small sinonasal microbiome studies published thus far.

122 We recently reported the findings of our multi-national, multicenter “International Sinonasal Microbiome  
123 Study” or ISMS.<sup>2</sup> This study, the largest and most diverse of its kind to date, attempted to address many  
124 of the limitations of the smaller previous studies, by standardizing collection, processing and analysis of  
125 the samples. Furthermore, its large sample size and multinational recruitment, meant that it was more  
126 likely to capture geographical and centre-based differences if present. A recent meta-analysis of published  
127 sinonasal 16S rRNA sequences revealed that the largest proportion of variance was attributed to  
128 differences between studies,<sup>3</sup> highlighting a role for performing a large multi-centre study that employed  
129 a unified methodology.

130 Contrary to the findings of previous studies, our international cohort showed no significant differences in  
131 alpha or beta diversity between the three groups of patients analyzed: healthy control patients without  
132 CRS and the two phenotypes of CRS patients, those with polyps (CRS<sub>w</sub>NP) and those without  
133 (CRS<sub>s</sub>NP). The study however revealed a potential grouping of samples as demonstrated on beta diversity  
134 exploratory analysis.<sup>2</sup> Accordingly, we hypothesized that the bacteriology of the sinuses could be  
135 categorized into various clusters of similar compositions. We inquired whether these potential groups  
136 would aid in describing the sinonasal microbial composition of patients or associate with clinical features.  
137 Similar attempts performed on gut microbiota in healthy individuals were termed *enterotyping*.<sup>4</sup> The  
138 clinical relevance of gut enterotypes remain the topic of research, and sometimes controversy. A previous  
139 exploration of clusters of sinus microbiota in patients was performed by Cope et al.<sup>5</sup> in which the authors

140 reported four compositionally distinct sinonasal microbial community states; the largest group of patients  
141 were dominated by a continuum of Staphylococcaceae and Corynebacteriaceae demonstrating a  
142 reciprocal relationship.<sup>5</sup>

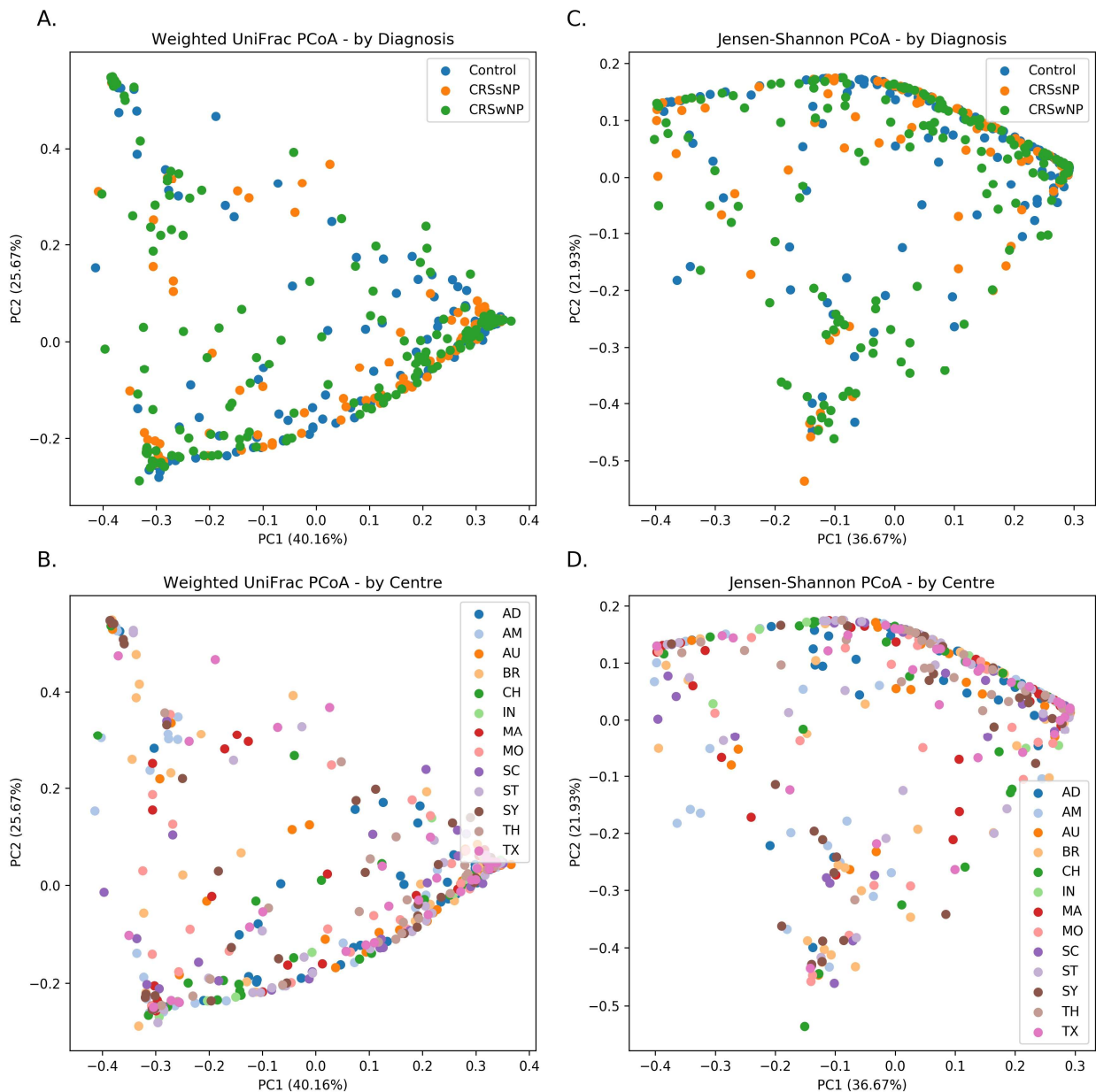
143 In this manuscript, we attempt “microbiotyping” to explain interpatient heterogeneity of the bacterial  
144 communities in the paranasal sinuses, and are the first to describe “sinonasal microbiotypes” across the  
145 first large, multi-centre cohort of individuals with and without CRS. We model our analysis on previous  
146 attempts of enterotyping the gut microbiome. We then describe the composition of these microbiotypes,  
147 explore potential clinical associations and validate microbiotyping on a separate sinus microbiome  
148 dataset.

149



## 150 RESULTS

### 151 Basic characteristics of the study cohort and beta diversity plots



152

153 **Figure 1: Beta diversity ordination plots.**

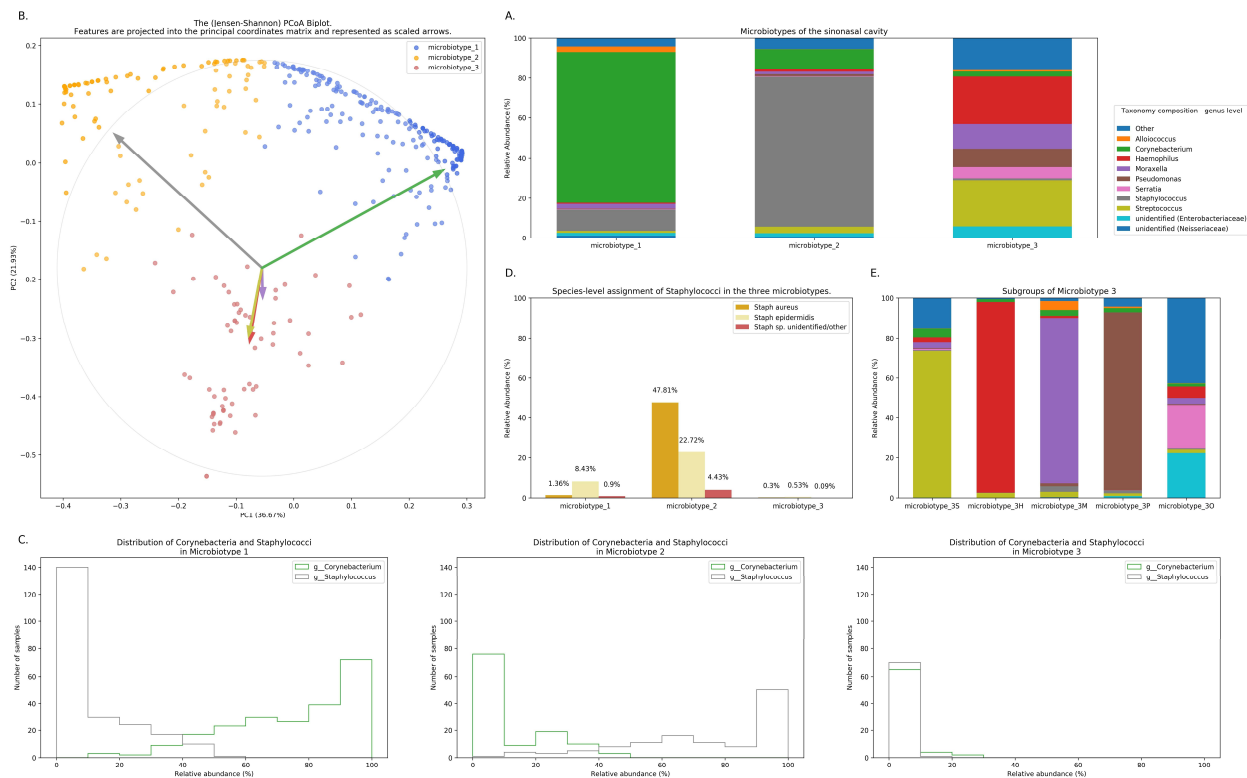
154 The main ISMS study cohort was described in our previous publication.<sup>2</sup> In brief, 410 samples were

155 included in the analysis collected from 13 centres representing 5 continents. These samples are distributed

156 along three diagnosis groups as follows: 99 CRSsNP patients, 172 CRSwNP patients, and 139 (non-CRS)  
 157 controls. Beta diversity ordination plots (of weighted UniFrac and Jensen-Shannon distances) are shown  
 158 in Figure 1. The plots do not reveal any distinct grouping by disease state or by centre, but on visual  
 159 inspection show a triangular arrangement suggesting that samples lie on a continuum between three  
 160 distinct clusters, providing motivation for further analysis.

## 161 Composition of the three sinonasal microbiotypes

162 We applied our microbiotyping approach through the unsupervised dimensionality reduction and  
 163 clustering method described in the Methods. The composition of the resulting “sinonasal microbiotypes”  
 164 is found in Figure 2A.



165  
 166 **Figure 2: Microbiotyping the sinonasal microbiome.** (A) Taxonomic composition of the three  
 167 microbiotypes at the genus level. (B) Illustration of the assigned microbiotypes on the Jensen-Shannon  
 168 PCoA biplot. Arrows were used to depict the projection of the genera onto the PCoA matrix. Each arrow  
 169 is indicated by the color of the genus according to the Legend. (C) Histograms demonstrating the relative

170 *abundance of Corynebacterium and Staphylococcus. (D) Distribution of staphylococcal species (mean*  
171 *relative abundance). (E) Subgroups of microbiotype 3 (hierarchical density-based clustering).*

172 Microbiotype 1 is dominated by *Corynebacterium* (mean relative abundance of 75.29%). Microbiotype 2  
173 is dominated by *Staphylococcus* (mean relative abundance of 74.96%). Microbiotype 3 contained samples  
174 that were mostly constituted of *Streptococcus*, *Haemophilus*, *Moraxella*, *Pseudomonas* and other genera.

175 The Abundance/Prevalence tables for the microbiotypes is demonstrated in Supplementary Tables [S1A](#),  
176 [S1B](#) and [S1C](#).

177 We used a PCoA biplot to project features (genera) onto the PCoA matrix.<sup>6</sup> The 5 topmost abundant  
178 genera were overlaid on the PCoA plot as arrows, originating from the centre of the plot and pointing to  
179 the direction of the projected feature coordinates. (Figure 2B) Each arrow is indicated by the color of the  
180 genus according to the Legend in Figure 2A, and the length of each was normalized as a percentage of the  
181 longest arrow. The coloring of the samples in 2B in the PCoA scatter plot according to the microbiotype  
182 assignment is provided for additional illustration. (Figure 2B) We note that the biplot arrows show a  
183 quasi-orthogonal arrangement between the key genera that constitute the microbiome.

184 The distributions of the relative abundances of *Corynebacterium* and *Staphylococcus* in all three  
185 microbiotypes were plotted in histograms (Figure 2C). It was noted that in microbiotype 1, most samples  
186 have a high abundance of Corynebacteria (i.e. Corynebacteria dominate), while Staphylococci appeared  
187 to dominate in microbiotype 2 in most samples.

### 188 **Dissection of “sinonasal microbiotype 3”**

189 We observed that Microbiotype 3 included various genera that did not cluster into the major two  
190 microbiotypes. It was also evident that this microbiotype is more heterogeneous. Applying the K-Means  
191 algorithm we showed poor clustering on only the first two and three Principal Components, since this  
192 group included multiple signatures with various dominant organisms. Accordingly, we employed the

193 hierarchical density-based clustering algorithm “hdbscan”<sup>7</sup> on the full-dimensional OTU table. One  
194 advantage of this algorithm is that it can estimate the number of clusters, without *a priori* specification by  
195 the user. This algorithm also has the ability to detect “outliers” that fail to cluster with the rest of the  
196 groups and detaches them into a separate “Miscellaneous/Other” group. We ran this algorithm on samples  
197 in Microbiotype 3 and this revealed four clusters, each dominated by one of the genera of *Streptococcus*  
198 (21 samples), *Haemophilus* (16 samples), *Moraxella* (9 samples), and *Pseudomonas* (7 samples), with a  
199 mean relative abundance ranging from 73.49% to 95.5%. The fifth cluster was the assigned  
200 “Miscellaneous/Other” group (18 samples). We term these “sub-microbiotypes”: microbiotype 3S, 3H,  
201 3M, 3P, and 3O, respectively. (Figure 2E)

## 202 **Exploring microbiotypes at the species-level reveals potential antagonism between** 203 ***Corynebacterium* species and *Staphylococcus aureus***

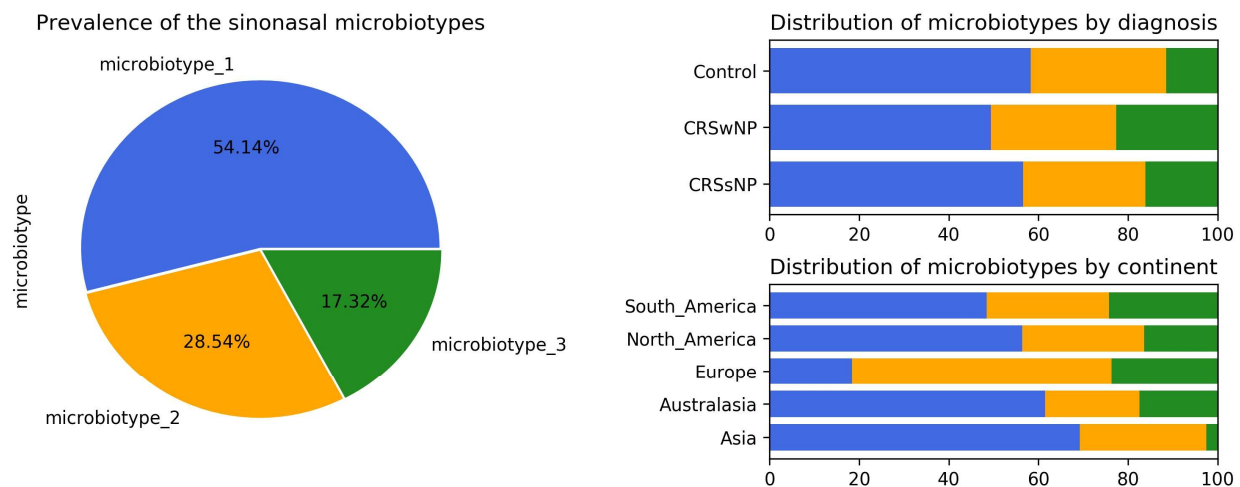
204 At present, species level assignment is limited by the current technology of 16S-surveys, the current state  
205 of microbial databases in general, and by our chosen short-read sequencing methodology. However,  
206 species level associations hold clinical significance for sinus health, since *Staphylococcus aureus* has  
207 been traditionally associated with biofilm formation and superantigen elaboration, both of which are  
208 associated with more severe sinus disease and poorer response to treatment. Furthermore nasal carriage of  
209 methicillin-resistant *Staphylococcus aureus* (MRSA) is a global health concern with implications that  
210 extend far beyond the sinuses. Moreover, our new QIIME 2-based pipeline<sup>8</sup> allows a higher “sub-OTU”  
211 resolution compared to older pipelines, offering an opportunity to resolve some taxa at species level when  
212 possible.<sup>9,10</sup>

213 We explored taxonomy assignment at the species level, with a focus on Staphylococcal species.  
214 Staphylococci were assigned to either *Staphylococcus aureus*, *Staphylococcus epidermidis* or unclassified  
215 *Staphylococcus*. We found that almost all of the assigned *Staphylococcus aureus* species were clustered in  
216 Microbiotype 2, forming 47.81% mean relative abundance of this Microbiotype, compared to 1.36% and  
217 0.3% in Microbiotype 1 and Microbiotype 3 respectively. (Figure 2E) Differential abundance of both

218 *Staphylococcus aureus* and *epidermidis* between the disease groups was confirmed as statistically  
219 significant using ANCOM.

220 In light of this finding, we hypothesized a reciprocal or antagonistic relationship between  
221 *Corynebacterium sp.* and *Staphylococcus aureus* and investigated this using SparCC. This confirmed a  
222 significant negative correlation between *Corynebacterium* genus and the species *Staphylococcus aureus*  
223 (SparCC correlation coefficient = -0.339,  $p = 0.001$ ). Interestingly, *Staphylococcus epidermidis* positively  
224 correlated with *Corynebacterium* (SparCC correlation coefficient = 0.271,  $p = 0.001$ ). These results  
225 should be interpreted cautiously in light of 16S-sequencing limitations. Nevertheless, they do appear to  
226 correlate to previous findings in the literature, including *in vitro* experiments<sup>11</sup>, a murine nasal bacterial  
227 interaction model<sup>12</sup>, and a survey of nasal vestibule swabs in healthy individuals<sup>13</sup>. These results suggest  
228 that a benign or probiotic role is played by both *Corynebacterium spp.* and *Staphylococcus epidermidis*  
229 when interacting with *Staphylococcus aureus*.

### 230 Prevalence and distribution of the microbiotypes in different diagnoses and centres



231  
232 **Figure 3: Prevalence and distribution of the microbiotypes.**

233 Microbiotype 1 was assigned to 222 samples (54.1%), microbiotype 2 to 117 samples (28.5%), and  
234 microbiotype 3 to 71 samples (17.3%). The prevalence distribution of the sinonasal microbiotypes did not

235 appear to significantly differ by the disease state of the sinuses. (Figure 3) However, a Chi-Squared test  
236 on the contingency table by centre showed significantly different distributions by centre (FDR-corrected p  
237 < 0.001): there was a higher prevalence of microbiotype 2 in our European centre (Amsterdam), and a  
238 higher prevalence of microbiotype 1 in Asian and Australasian centres, with a much lower prevalence of  
239 microbiotype 3 in Asia. (Figure 3 and Table 1)

240 *Table 1: Distribution of microbiotypes by diagnosis and continent.*

variable	value	microbiotype_1	microbiotype_2	microbiotype_3	p value
Diagnosis	CRSsNP	56 (56.6%)	27 (27.3%)	16 (16.2%)	0.507
	CRSwNP	85 (49.4%)	48 (27.9%)	39 (22.7%)	
	Control	81 (58.3%)	42 (30.2%)	16 (11.5%)	
Continent	Asia	27 (69.2%)	11 (28.2%)	1 (2.6%)	< 0.001
	Australasia	67 (61.5%)	23 (21.1%)	19 (17.4%)	
	Europe	7 (18.4%)	22 (57.9%)	9 (23.7%)	
	North_America	89 (56.3%)	43 (27.2%)	26 (16.5%)	
	South_America	32 (48.5%)	18 (27.3%)	16 (24.2%)	

241

## 242 **Associations of microbiotypes with clinical variables**

243 We then explore the distribution of the three microbiotypes among multiple clinical variables in Table 2.  
244 This shows no significant difference for some variables including asthma, aspirin sensitivity, GORD,  
245 diabetes mellitus, and current smoking status, (FDR-corrected  $p > 0.05$ ; Chi-squared test). The cross  
246 tabulation however revealed a statistically significant association with “aspirin sensitivity” or aspirin-  
247 exacerbated respiratory disease (AERD) ( $p = 0.02$ ), although this did not persist after a Benjamini-  
248 Hochberg correction (corrected  $p = 0.077$ ). Patients who were aspirin-sensitive (or suffering from AERD)  
249 showed less prevalence of microbiotypes 1, 2 and a higher prevalence of microbiotype 3, compared to  
250 those who were not aspirin-sensitive. On the other hand, patients who were undergoing their “primary

251 surgery”, had a higher prevalence of microbiotype 1 and a lower prevalence of microbiotype 3, compared  
252 to those patients who had had previous surgeries, but these results were not statistically significant.

253 *Table 2: Distribution of microbiotypes by various clinical variables.*

variable	value	microbiotype_1	microbiotype_2	microbiotype_3	p value
Asthma	No	162 (56.4%)	81 (28.2%)	44 (15.3%)	0.906
	Yes	55 (51.4%)	31 (29.0%)	21 (19.6%)	
Aspirin sensitivity	No	202 (55.3%)	106 (29.0%)	57 (15.6%)	0.077
	Yes	12 (48.0%)	5 (20.0%)	8 (32.0%)	
Diabetes	No	189 (54.9%)	98 (28.5%)	57 (16.6%)	0.979
	Yes	22 (55.0%)	11 (27.5%)	7 (17.5%)	
GORD	No	177 (55.3%)	93 (29.1%)	50 (15.6%)	0.979
	Yes	35 (55.6%)	17 (27.0%)	11 (17.5%)	
Current Smoker	No	204 (54.4%)	110 (29.3%)	61 (16.3%)	0.077
	Yes	15 (57.7%)	4 (15.4%)	7 (26.9%)	
Primary surgery	No	92 (47.2%)	57 (29.2%)	46 (23.6%)	0.114
	Yes	130 (60.5%)	60 (27.9%)	25 (11.6%)	

254

### 255 **Validation of sinonasal microbiotyping on a separate dataset**

256 We validated our approach on a separate 16S dataset we called Dataset Two. As described in the Methods  
257 section, we validated this using an independent unsupervised approach and a semi-supervised approach  
258 guided by the Main Dataset.

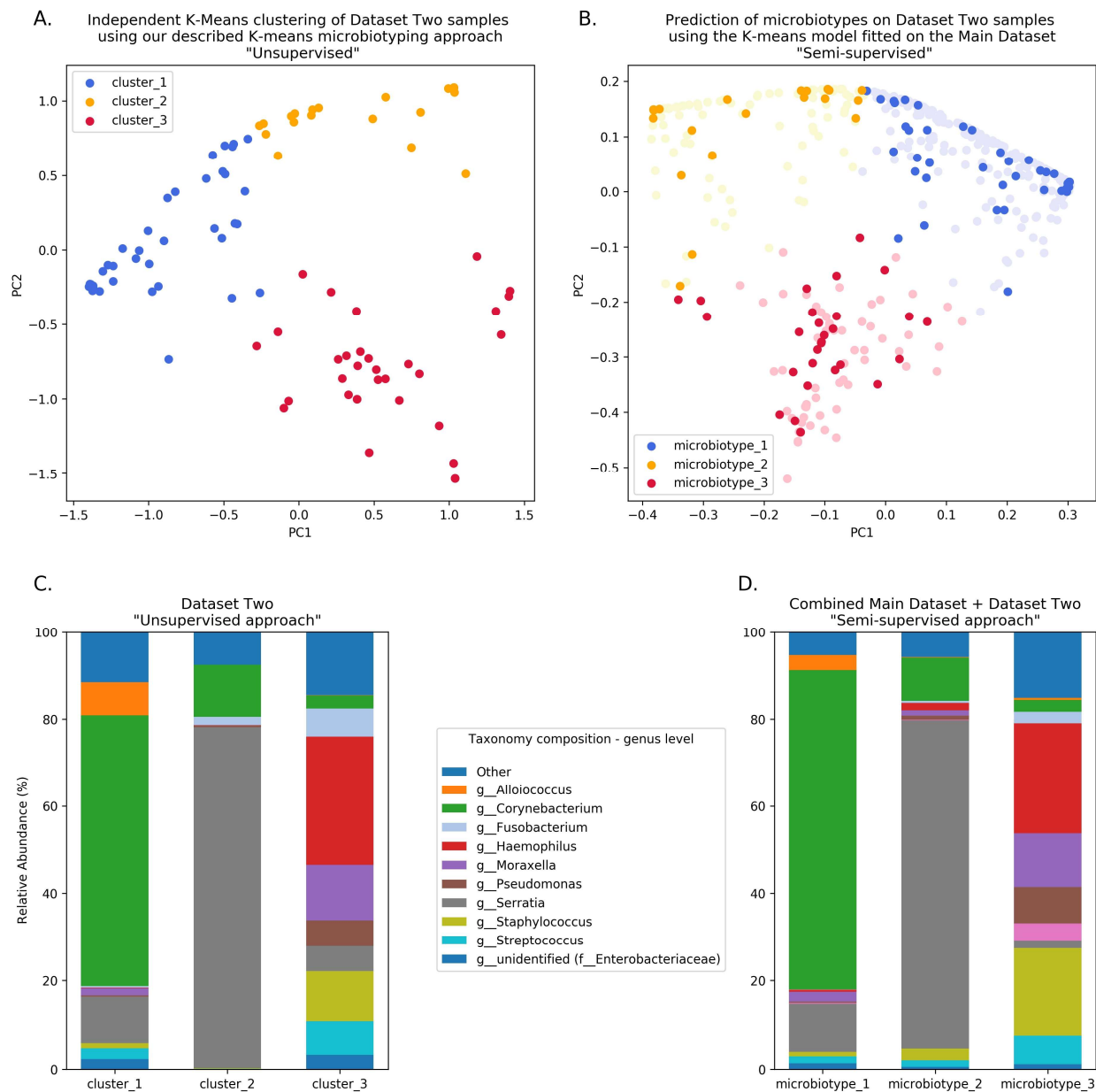
259 The first unsupervised approach yielded three clusters similar to the microbiotypes described on the Main  
260 Dataset, with one cluster exhibiting high mean relative abundance of Corynebacteria, a second cluster  
261 exhibiting high mean relative abundance of Staphylococcus, and a third cluster with other dominant  
262 genera. Plotting the first two Principal Components (Figure 4A) resulting from PCoA on the JSD matrix  
263 revealed the same triangular distribution of samples observed in Figure 1.

264 Prevalence of the microbiotypes in this dataset (using the unsupervised approach) was as follows:  
265 microbiotype 1 assigned 39.2% of samples, microbiotype 2 with 26.8% of samples, and microbiotype 3  
266 with 34.0%.

267 The second semi-supervised approach yielded similar results (Figure 4; Supplementary Table), differing  
268 in the classification of only 3 samples (out of 97 samples; i.e. 3.09%). (See Supplementary Jupyter  
269 notebook) Two of these samples show *Staphylococcus* dominating the samples in combination with  
270 *Haemophilus*, with no overt dominance of one taxon over the other, making them more-or-less  
271 transitional samples between the signatures of microbiotypes 2 and 3. The third sample was dominated by  
272 *Staphylococcus* and *Corynebacterium*, making it a transitional sample between microbiotype 1 and  
273 microbiotype 2, with Staphylococcal species assigned to *epidermidis*, making this more appropriately  
274 assigned to microbiotype 1. (see Supplementary Jupyter notebook)

275 These results validate the microbiotyping approach and suggest that our approach and dataset could be  
276 used to guide classification of sinonasal samples sequenced in future separate studies. (Figure 4)  
277 Moreover, it points towards a potential clinical relevance of performing sinonasal microbiotyping.





278

279 **Figure 4: Validation of microbiotyping approach on Dataset Two.**

280

## 281 DISCUSSION

282 We demonstrate that the microbiota of most sinus swab samples could be classified into distinct  
283 signatures or archetypes, which we have termed “sinonasal microbiotypes”. We observed three main  
284 microbiotypes: the most prevalent being a *Corynebacterium*-dominated microbiotype (microbiotype 1),  
285 then a *Staphylococcus*-dominated microbiotype (microbiotype 2), and microbiotype 3 which includes  
286 samples dominated by *Streptococcus*, *Haemophilus*, *Moraxella*, *Pseudomonas*, and other genera (3S, 3H,  
287 3M, 3P, and 3O respectively).

288 As we have previously reported,<sup>2</sup> the sinus microbiota are dominated by the genera *Corynebacterium* and  
289 *Staphylococcus* (microbiotypes 1 and 2). A similar clustering approach to the sinus microbiome was  
290 applied by Cope and colleagues, who utilized Dirichlet multinomial mixture models (DMMs),<sup>5</sup> and  
291 reported that most samples in their study were occupied by a continuum of Staphylococcaceae and  
292 Corynebacteriaceae.<sup>5</sup> It appears that, regardless the statistical or clustering methodology utilized, it is  
293 most likely that the sinonasal microbiome consists of core organisms<sup>2</sup> that have a distinct co-occurrence  
294 pattern. This could be explored through a network analysis approach and should be a future area of study.

295 *Staphylococcus aureus* has been perceived to be an important pathogen in sinus inflammatory disease.  
296 *Staphylococcus aureus* biofilms may act as a nidus for recurrent infections<sup>14,15</sup> and as a “nemesis” of  
297 otherwise-successful sinus surgery.<sup>16–18</sup> *Staphylococcus aureus* is also a producer of exotoxins, which in  
298 some cases can serve as superantigens, and these have been previously described as playing an important  
299 role in the pathogenesis of CRSwNP.<sup>19</sup> *Pseudomonas aeruginosa* biofilms are also virulent organisms that  
300 are difficult to eradicate from the sinuses, and have been associated with worse clinical outcomes.<sup>20</sup> Both  
301 these organisms are important pathogens in the chronic mucociliary dysfunction exhibited in cystic  
302 fibrosis. However, methicillin-resistant *Staphylococcus aureus* (MRSA) is an important nasal colonizer  
303 that could asymptotically colonize the nose. What determines the clinical course, between  
304 asymptomatic colonization versus symptomatic pathogenicity, remains an interesting topic of research. In

305 this study, we identified a potential reciprocal relationship between *Staphylococcus aureus* and  
306 *Corynebacterium*. Being aware of the challenges of compositional data analysis, we utilized for this  
307 purpose the specialized SparCC algorithm which infers correlations from compositional data.<sup>21</sup> This  
308 finding needs to be supported by future co-culture experiments, but suggests that *Corynebacterium sp.*  
309 may be a “cornerstone” of sinus microbial health. It is important to note that our bioinformatic  
310 methodology has been intentionally designed to utilize state-of-the-art software methods at every step of  
311 the analysis pipeline, in order to maximise the resolution of taxonomy assignment.<sup>8,9,22</sup> Nevertheless, our  
312 approach is still confined within the limitations of current 16S sequencing methodologies, and the  
313 confidence of assignment is reduced beyond the genus level. Our analysis pipeline could not delineate  
314 between different *Corynebacterium* at the species level and *Staphylococcus aureus* at the strain level.  
315 Hence functional difference between samples with same species remain to be determined using a  
316 functional metagenomics approach. A recent study suggest that by incorporating location information or  
317 “sample-level metadata” into species-level assignment accuracy could be improved.<sup>23</sup> In our study, the  
318 differential relationships of both *Staphylococcus aureus* and *epidermidis* towards *Corynebacteria*  
319 (negative and positive associations, respectively) could be of clinical significance and is worthy of future  
320 investigation. We performed a post-hoc inspection of species-level assignment in Dataset Two, to  
321 investigate whether this finding will be reproducible in a separate dataset. This confirmed clustering of  
322 almost all *Staphylococcus aureus* species in microbiotype 2. (Supplementary Results in Jupyter  
323 Notebook)

324 Interestingly, we found that the distribution of the sinonasal microbiotypes was not significantly dis-  
325 similar amongst healthy controls and CRS patients. There appeared to be no significant associations with  
326 other clinical variables such as asthma and aspirin-sensitivity after controlling for multiple comparisons.  
327 (Table 2) The distribution of the microbiotypes however differed according to centre/location of  
328 collection. (Figure 3) As such, we cannot conclude based on our study that microbiotypes could function  
329 independently as a disease biomarker. Although not reaching statistical significance (chi squared  $p >$

330 0.05) the prevalence of microbiotype 3 was higher in CRSsNP and CRSwNP, compared to controls. It  
331 could be the case that chronicity of inflammation -on its own- is not a determinant of a dysbiotic  
332 microbiome, but whether there is a clinically-evident “sinus infection” current at the time of sample  
333 collection. In this theory, stable chronic sinuses with no overt signs of acute or chronic infection, may  
334 remain similar to a “healthy sinus microbiome”. Only when the sinuses are clinically infected (as evident  
335 on clinical symptoms and endoscopic findings), the microbiota become disrupted and the dysbiosis  
336 exaggerated. It is important to note that *Streptococcus*, *Haemophilus* and *Moraxella* (represented here in  
337 microbiotype 3) have been traditionally implicated in acute infections of the upper respiratory tract  
338 including acute rhinosinusitis and acute otitis media. Unfortunately, information regarding acute  
339 exacerbations was not explored within this study.

340 Regarding geographical differences: Asia and Australasia showed an over-representation of microbiotype  
341 1. Europe had a higher prevalence of microbiotype 2. Unfortunately, the study only included one  
342 European centre (Amsterdam) so it is difficult to be certain whether this finding generalizes to other  
343 locations in Europe. The driving factors for these geographical differences could be multiple, including  
344 but not limited to clinical practices such as local antibiotic prescriptions for CRS and timing of  
345 recruitment of patients for sinus surgery, as discussed previously.<sup>2</sup>

346 We have adapted our methodology from the enterotyping approach taken by Arumugam et al.<sup>4</sup> for  
347 classifying bacterial signatures of the gut microbiome. In their original manuscript, they described three  
348 different enterotypes in the gut dominated by *Prevotella*, *Bacteroidetes*, and *Ruminococcus* respectively.<sup>4</sup>  
349 Several papers have correlated gut enterotypes with various clinical variables.<sup>24,25</sup> Despite this,  
350 enterotyping as an approach to population stratification has not been without its controversies. Several  
351 authors have criticized the definition of distinct clusters, since it neglects intra-cluster variation and  
352 gradients between clusters.<sup>26–29</sup> We provide answers to previous critique<sup>28</sup> to enterotyping as it applies to  
353 our study in Supplementary Table S2. It is important to note these valid criticisms to any community  
354 typing approach. In our experiment, the clusters or types lie on a continuum, with some samples falling in

355 the gradients between two, or perhaps even all three microbiotypes (see ordination plots). The histograms  
356 in Figure 2 also suggest this, but they do show most samples in each microbiotype feature a high relative  
357 abundance of a dominating genus in many samples. We investigated a simple dominance measure, the  
358 Berger-Parker (BP) alpha diversity index,<sup>30</sup> in the combined datasets' 507 samples. The Berger-Parker  
359 index simply reports the relative abundance of the most dominant taxon in a sample. This found that only  
360 24.9% of samples had a dominating taxon that only had a relative abundance of 50% or less. On the other  
361 hand, 51.9% of samples had the dominant taxon exhibiting a relative abundance of greater than 70% of  
362 the sample.(Supplementary Results in Jupyter notebook; Supplementary Figure S1) This shows that in  
363 most samples, there is one dominating organism. Based on these results, the microbiotyping approach is  
364 therefore proposed to reduce complexity about modeling bacterial interactions in the sinuses, and not to  
365 suggest that each type is a walled-off discrete cluster. Further investigations into the local substructures of  
366 each type will be required to further explore the roles and interactions of its constituent taxa. Another  
367 limitation of our description of microbiotypes is that they may as well describe different community  
368 “states” rather than community “types”, since we do not have longitudinal data to describe how these  
369 clusters behave with the passage of time and treatments. Hence, we could not confirm whether these are  
370 stable, consistent communities across time. It may well be that intermediate samples lying between two or  
371 more microbiotypes are representing a transitional state. An important future avenue of research is to  
372 conduct a longitudinal study to investigate the temporal stability of these clusters.

373 We predict that ongoing sinonasal microbiome research and consequent large meta-analyses of  
374 microbiota studies, with novel tools (such as QIITA<sup>31</sup>) enabling such large-scale studies, will allow the  
375 refinement of these types and further clarify their clinical/microbiological utility. Our methodological  
376 approach to describe the microbiotypes is not exclusive, as alternative statistical or machine-learning  
377 approaches could be employed to investigate them. In light of this, we expect that international multi-  
378 centre standardization and rationalization of the sinonasal microbiotypes would be possible in the future,  
379 similar to the recent proposed effort to standardize enterotyping of the gut microbiota by Costea et al.<sup>29</sup>

## 380 CONCLUSION

381 We investigated the ISMS dataset through an approach modeled on human gut microbiome enterotyping  
382 and we found three major microbial community types or “microbiotypes” as clusters that lie on a  
383 continuum, based on an unsupervised machine learning approach that involved dimensionality reduction  
384 and clustering. Microbiotypes did not show an association with disease state or clinical variable,  
385 suggesting that they could not function as independent disease biomarkers. The description of these  
386 microbiotypes has also unveiled a potential reciprocal relationship between *Staphylococcus aureus* and  
387 *Corynebacterium spp.* in the sinuses that requires further investigation in future studies. The findings  
388 were validated on a separate previously unpublished sinus bacterial 16S gene dataset. Microbiotypes are  
389 therefore proposed to reduce the complexity of modeling bacterial interactions in the sinuses, and in this  
390 sense hold microbiological and clinical relevance that could potentially influence medical and surgical  
391 treatment of CRS patients.

392

## 393 **METHODS**

### 394 **The “International Sinonasal Microbiome Study (ISMS)” dataset**

395 We perform the primary analysis on the dataset obtained from the “International Sinonasal Microbiome  
396 Study (ISMS)” project.<sup>2</sup> In summary, this dataset is a multi-centre 16S-amplicon dataset which includes  
397 endoscopically-guided, guarded swabs collected from the sinuses (in particular the middle meatus /  
398 anterior ethmoid region) of 532 participants in 13 centres representing 5 continents. Details of sample  
399 collection, DNA extraction and sequencing methodologies are described in the original report.<sup>2</sup> The 16S  
400 gene region sequenced was the V3–V4 hypervariable region, utilizing primers  
401 (CCTAYGGGRBGCASCAG forward primer) and (GGACTACNNGGGTATCTAAT reverse primer)  
402 according to protocols at the sequencing facility (the Australian Genome Research Facility; AGRF).  
403 Sequencing was done on the Illumina MiSeq platform (Illumina Inc., San Diego, CA) with 300-base-pairs  
404 paired-end Illumina chemistry

### 405 **Bioinformatics pipeline**

406 Details of the bioinformatic pipeline is detailed in the original report.<sup>2</sup> In summary, we utilized a QIIME  
407 2-based pipeline.<sup>8</sup> Forward and reverse fastq reads were joined<sup>32</sup>, quality-filtered,<sup>33</sup> abundance-filtered<sup>34</sup>,  
408 then denoised using deblur<sup>9</sup> through QIIME 2-based plugins. This yielded a final feature table of high-  
409 quality, high-resolution Amplicon Sequence Variants (ASVs). Taxonomy assignment and phylogenetic  
410 tree generation<sup>35</sup> was done against the Greengenes<sup>36</sup> database; and taxonomy was assigned using the  
411 QIIME 2 BLAST assigner.<sup>22</sup> A rarefaction minimum depth cut-off was chosen at 400 and this yielded 410  
412 samples out of the original 532 for downstream analysis. The same pipeline was then applied on DataSet  
413 Two for purposes of validation of microbiodotyping. We chose to reproduce exactly all the original pipeline  
414 steps on DataSet Two, despite being a completely separate dataset, to reduce bias.

## 415 **Delineating the microbiotypes of the sinonasal microbiome**

416 Our approach was guided by the “enterotyping” method described by Arumugam et al.<sup>4</sup> with adaptations.  
417 We constructed a sample distance matrix using the Jensen-Shannon distance (JSD) metric, as used in the  
418 original “enterotypes” paper.<sup>4</sup> The Jensen-Shannon distances were calculated between samples in the  
419 genus-level-assigned table in a pairwise fashion using the JSD function in the R package “philentropy”  
420 with a log (log<sub>10</sub>) base. Following this, Principal Coordinate analysis (PCoA) was done on the distance  
421 matrix for dimensionality reduction and visualization. Clustering was then performed using a standard K-  
422 means clustering algorithm, as implemented in the machine learning Python package scikit-learn (version  
423 0.20.1);<sup>37</sup>) on the first two principal components (PCs) obtained from the PCoA, with the number of  
424 clusters (k) chosen at 3 based on visual inspection of the beta diversity PCoA plots. Average silhouette  
425 scores, as implemented in scikit-learn, for the range (k = 2 - 8) were calculated to assess clustering  
426 quality, and this revealed the highest silhouette scores: 0.61 and 0.6 for [k=4] and [k=3] respectively. The  
427 three resulting clusters were defined as the three sinonasal microbiotypes. For further exploration of the  
428 subgroups that constitute microbiotype 3, we used the hierarchical density-based clustering algorithm  
429 “hdbscan”<sup>7</sup> on the full-dimensional feature table. Genera were projected onto the PCoA matrix using a  
430 biplot approach<sup>6</sup>, as implemented in scikit-bio’s function “*pcoa\_biplot*”. Genera were represented in the  
431 biplot figure as arrows, originating from the centre of the plot pointing to the direction of the projected  
432 feature coordinates, and the lengths normalized as a percentage of the longest arrow. We utilized  
433 “Analysis of Compositions of Microbiomes (ANCOM)”<sup>38</sup> for identifying differentially-abundant taxa.  
434 Taxa genus level and Staphylococcus species level co-occurrence/correlation analysis were done after  
435 taxonomy assignment using SparCC algorithm,<sup>21</sup> in the fast implementation in FastSpar.<sup>39</sup>

## 436 **Validating microbiotypes on a second sinonasal microbiome dataset**

437 To infer whether our classification could be generalizable to other sinonasal microbiome samples not  
438 included in this study, we sought to validate our microbiotyping approach on a separate, previously-  
439 unpublished, 16S dataset. This dataset includes sinonasal microbiome swabs collected from private and



440 public patients attending the Otolaryngology Department (University of Adelaide) to have surgery done  
441 by the authors P.J.W., A.J.P. or the Otorhinolaryngology Service at the Queen Elizabeth Hospital in  
442 Adelaide, South Australia. Similar to the main dataset, these included CRS patients who underwent  
443 endoscopic sinus surgery for this sinus disease, and non-CRS control patients who underwent other  
444 otolaryngological procedures, such as tonsillectomy, septoplasty or skullbase tumour resection. Sample  
445 collection, and processing were done in a standardized fashion similar to that has been described in the  
446 ISMS main dataset, except that DNA extraction was carried out using the PowerLyzer Power-Soil DNA  
447 kit (MoBio Laboratories, Salona Beach, CA) as previously described<sup>40</sup>, rather than the Qiagen DNeasy kit  
448 (Qiagen, Hilden, Germany). Similar to the ISMS samples, library preparation and 16S sequencing were  
449 done at the Australian Genome Research Facility (AGRF) on the Illumina MiSeq platform (Illumina Inc.,  
450 San Diego, CA, USA) with the 300-base-pairs paired-end chemistry. Libraries were generated by  
451 amplifying (341F–806R) primers against the V3–V4 hypervariable region of the 16S gene  
452 (CCTAYGGGRBGCASCAG forward primer; GGACTACNNGGGTATCTAAT reverse primer).<sup>41</sup> PCR  
453 was done using AmpliTaq Gold 360 master mix (Life Technologies, Mulgrave, Australia) following a  
454 two-stage PCR protocol (29 cycles for the first stage; and 8 cycles for the second, indexing stage).  
455 Sequencing was done over two MiSeq runs in January 2015. We termed this dataset in this manuscript  
456 “Dataset Two”. This dataset comprises samples collected from 129 participants. Rarefaction at a cutoff of  
457 400 reads was performed, to match what was performed for the main dataset, and samples with read  
458 number less than 400 were excluded; this yielded a final feature table containing 97 samples, representing  
459 33 CRSsNP patients, 35 CRSwNP patients, and 29 controls.

460 We took two separate approaches to validation. The first approach is to replicate the previously-described  
461 unsupervised K-means microbiotyping methodology independently on samples in Dataset Two. We call  
462 this first approach the “unsupervised approach”. The second approach is to use the K-means model that  
463 was fitted on the samples from the Main Dataset to predict labels (i.e. microbiotypes) of the samples in

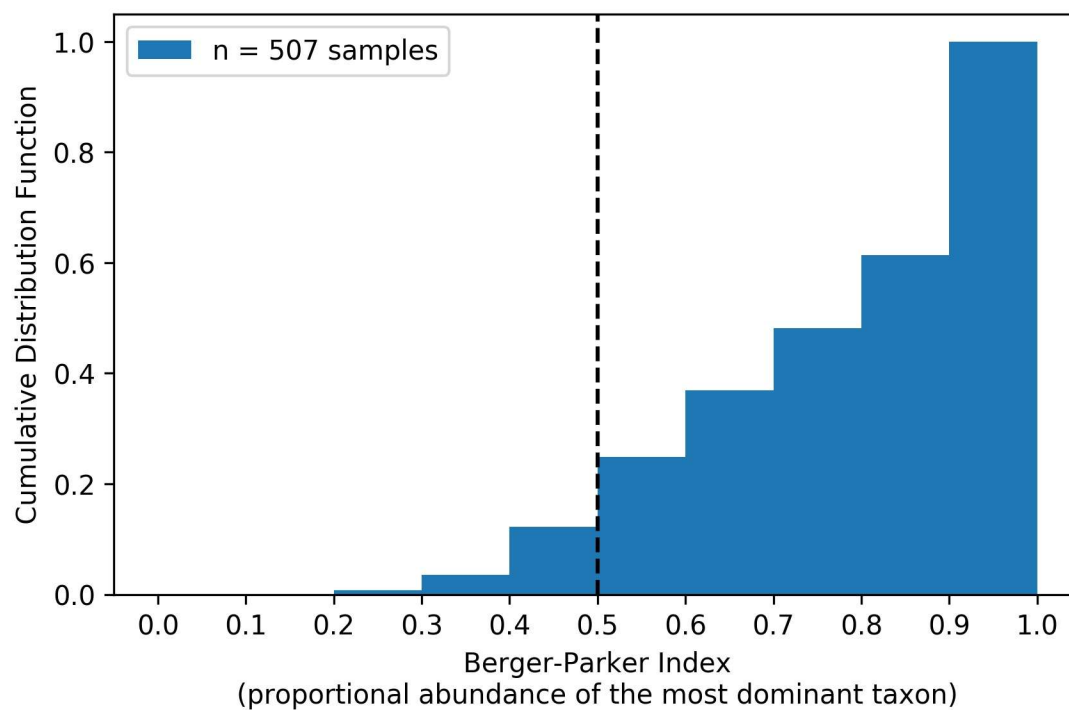
464 Dataset Two. As such, the Main Dataset is used as a “training dataset” in the language of machine  
465 learning. We called the second approach the “semi-supervised approach”.

#### 466 **Statistical Analysis**

467 All frontend analyses were done using the Jupyter notebook frontend<sup>42</sup> and utilizing the assistance of  
468 packages from the Scientific Python<sup>43</sup> stack (numpy, scipy, pandas, statsmodels), scikit-learn<sup>37</sup>, scikit-bio  
469 (<https://github.com/biocore/scikit-bio>) and omicexperiment  
470 (<https://www.github.com/bassio/omicexperiment>).

471

472 **Supplementary Figures**



473

474 **Figure S1: Cumulative distribution function of the Berger-Parker Index in the combined datasets.**

475

476 **Supplementary Tables**

477 *Table S1A: Predominant taxa of microbiotype 1.*

genus	Mean Relative Abundance (%)	Prevalence (%)
Corynebacterium	75.29	100
Staphylococcus	10.69	76.58
Alloiococcus	2.79	28.83
Moraxella	2.31	9.91
unidentified (Enterobacteriaceae)	1.41	15.32
unidentified (Neisseriaceae)	1.18	20.72
Streptococcus	1	21.62
Haemophilus	0.56	9.91
unidentified (Moraxellaceae)	0.44	2.7
Ralstonia	0.34	10.36

478

479 *Table S1B: Predominant taxa of microbiotype 2.*

genus	Mean Relative Abundance (%)	Prevalence (%)
Staphylococcus	74.96	100
Corynebacterium	9.87	64.1
Streptococcus	3.22	25.64
unidentified (Enterobacteriaceae)	1.82	15.38
Haemophilus	1.41	10.26
Moraxella	1.27	5.13
Ralstonia	1.19	11.97
Pseudomonas	1.05	6.84
Parvimonas	0.72	0.85
unidentified (Neisseriaceae)	0.61	7.69

480

481 *Table SIC: Predominant taxa of microbiotype 3.*

genus	Mean Relative Abundance (%)	Prevalence (%)
Haemophilus	23.78	40.85
Streptococcus	23.22	46.48
Moraxella	12.11	19.72
Pseudomonas	9.17	15.49
unidentified (Enterobacteriaceae)	5.74	9.86
Serratia	5.7	8.45
Klebsiella	2.75	4.23
Corynebacterium	2.56	46.48
Prevotella	1.44	12.68
Acinetobacter	1.38	1.41

482

483 *Table S2: Addressing previous criticism to gut enterotyping.*

Critique	Answer
Discrete clusters or a multi-dimensional gradient?	We acknowledge the a proportion of samples fall in the gradient between the proposed microbiotypes. Berger-Parker index investigation showed that most samples had one dominating taxon.
Do discrete clusters link to human disease?	No. We report that we could not find an association between the microbiotype and chronic sinusitis disease status.
Is sampling frame or selection bias affecting results?	No; Multi-centre international study with consecutive sampling methodology. We also validate on a separate dataset.
Use inappropriate visualization such as “star-burst plots”?	We did not use inappropriate visualizations.
Use a supervised approach “between-class analysis”?	We use an unsupervised clustering and dimensionality reduction approach.
Is an individual’s microbiotype stable over time?	Answer unknown; Future longitudinal studies required.

## 485 REFERENCES

- 486 1. Fokkens, W. J. *et al.* EPOS 2012: European position paper on rhinosinusitis and nasal polyps 2012. A  
487 summary for otorhinolaryngologists. *Rhinology* **50**, 1–12 (2012).
- 488 2. Paramasivan, S. *et al.* The international sinonasal microbiome study (ISMS): A multi centre,  
489 international characterization of sinonasal bacterial ecology. *bioRxiv* 548743 (2019). doi:[10.1101/548743](https://doi.org/10.1101/548743)
- 490 3. Wagner Mackenzie, B. *et al.* Bacterial community collapse: A meta-analysis of the sinonasal  
491 microbiota in chronic rhinosinusitis. *Environmental Microbiology* **19**, 381–392 (2017).
- 492 4. Arumugam, M. *et al.* Enterotypes of the human gut microbiome. *Nature* **473**, 174–180 (2011).
- 493 5. Cope, E. K., Goldberg, A. N., Pletcher, S. D. & Lynch, S. V. Compositionally and functionally distinct  
494 sinus microbiota in chronic rhinosinusitis patients have immunological and clinically divergent  
495 consequences. *Microbiome* **5**, 53 (2017).
- 496 6. Legendre, P. & Legendre, L. *Numerical ecology*. (Elsevier, 2012).
- 497 7. McInnes, L., Healy, J. & Astels, S. HdbSCAN: Hierarchical density based clustering. *The Journal of*  
498 *Open Source Software* **2**, 205 (2017).
- 499 8. Bolyen, E. *et al.* *QIIME 2: Reproducible, interactive, scalable, and extensible microbiome data*  
500 *science*. (PeerJ Inc., 2018). doi:[10.7287/peerj.preprints.27295v1](https://doi.org/10.7287/peerj.preprints.27295v1)
- 501 9. Amir, A. *et al.* Deblur Rapidly Resolves Single-Nucleotide Community Sequence Patterns. *mSystems*  
502 **2**,
- 503 10. Thompson, L. R. *et al.* A communal catalogue reveals Earth’s multiscale microbial diversity. *Nature*  
504 **551**, 457–463 (2017).
- 505 11. Barrow, G. I. Microbial Antagonism by *Staphylococcus aureus*. *Microbiology* **31**, 471–481 (1963).
- 506 12. Cleland, E. J. *et al.* Probiotic manipulation of the chronic rhinosinusitis microbiome. *International*  
507 *Forum of Allergy & Rhinology* **4**, 309–314 (2014).
- 508 13. Lina, G. *et al.* Bacterial competition for human nasal cavity colonization: Role of *Staphylococcal* agr  
509 alleles. *Applied and Environmental Microbiology* **69**, 18–23 (2003).
- 510 14. Jervis-Bardy, J., Foreman, A., Boase, S., Valentine, R. & Wormald, P.-J. What is the origin of  
511 *Staphylococcus aureus* in the early postoperative sinonasal cavity? *International Forum of Allergy &*  
512 *Rhinology* **1**, 308–312
- 513 15. Drilling, A. *et al.* Cousins, siblings, or copies: The genomics of recurrent *Staphylococcus aureus*  
514 infections in chronic rhinosinusitis. *International Forum of Allergy & Rhinology* **4**, 953–960 (2014).
- 515 16. Psaltis, A. J., Weitzel, E. K., Ha, K. R. & Wormald, P.-J. The effect of bacterial biofilms on post-  
516 sinus surgical outcomes. *American Journal of Rhinology* **22**, 1–6
- 517 17. Foreman, A. & Wormald, P.-J. Different biofilms, different disease? A clinical outcomes study. *The*  
518 *Laryngoscope* **120**, 1701–1706 (2010).



- 519 18. Singhal, D., Foreman, A., Bardy, J.-J. & Wormald, P.-J. Staphylococcus aureus biofilms: Nemesis of  
520 endoscopic sinus surgery. *The Laryngoscope* **121**, 1578–1583 (2011).
- 521 19. Bachert, C., Zhang, N., Patou, J., van Zele, T. & Gevaert, P. Role of staphylococcal superantigens in  
522 upper airway disease. *Current Opinion in Allergy and Clinical Immunology* **8**, 34–38 (2008).
- 523 20. Bendouah, Z., Barbeau, J., Hamad, W. A. & Desrosiers, M. Biofilm formation by Staphylococcus  
524 aureus and Pseudomonas aeruginosa is associated with an unfavorable evolution after surgery for chronic  
525 sinusitis and nasal polyposis. *Otolaryngology–Head and Neck Surgery: Official Journal of American  
526 Academy of Otolaryngology-Head and Neck Surgery* **134**, 991–996 (2006).
- 527 21. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data. *PLOS*  
528 *Computational Biology* **8**, e1002687 (2012).
- 529 22. Bokulich, N. A. *et al.* Optimizing taxonomic classification of marker-gene amplicon sequences with  
530 QIIME 2's q2-feature-classifier plugin. *Microbiome* **6**, 90 (2018).
- 531 23. Kaehler, B. D., Bokulich, N., Caporaso, J. G. & Huttley, G. A. Species-level microbial sequence  
532 classification is improved by source-environment information. *bioRxiv* 406611 (2018).  
533 doi:[10.1101/406611](https://doi.org/10.1101/406611)
- 534 24. Wu, G. D. *et al.* Linking long-term dietary patterns with gut microbial enterotypes. *Science (New  
535 York, N.Y.)* **334**, 105–108 (2011).
- 536 25. Vandeputte, D. *et al.* Stool consistency is strongly associated with gut microbiota richness and  
537 composition, enterotypes and bacterial growth rates. *Gut* **65**, 57–62 (2016).
- 538 26. Jeffery, I. B., Claesson, M. J., O'Toole, P. W. & Shanahan, F. Categorization of the gut microbiota:  
539 Enterotypes or gradients? *Nature Reviews. Microbiology* **10**, 591–592 (2012).
- 540 27. Koren, O. *et al.* A Guide to Enterotypes across the Human Body: Meta-Analysis of Microbial  
541 Community Structures in Human Microbiome Datasets. *PLOS Computational Biology* **9**, e1002863  
542 (2013).
- 543 28. Knights, D. *et al.* Rethinking 'Enterotypes'. *Cell host & microbe* **16**, 433–437 (2014).
- 544 29. Costea, P. I. *et al.* Enterotypes in the landscape of gut microbial community composition. *Nature  
545 Microbiology* **3**, 8–16 (2018).
- 546 30. Berger, W. H. & Parker, F. L. Diversity of planktonic foraminifera in deep-sea sediments. *Science  
547 (New York, N.Y.)* **168**, 1345–1347 (1970).
- 548 31. Gonzalez, A. *et al.* Qiita: Rapid, web-enabled microbiome meta-analysis. *Nature Methods* **15**, 796–  
549 798 (2018).
- 550 32. Zhang, J., Kobert, K., Flouri, T. & Stamatakis, A. PEAR: A fast and accurate Illumina Paired-End  
551 reAd mergeR. *Bioinformatics* **30**, 614–620 (2014).
- 552 33. Bokulich, N. A. *et al.* Quality-filtering vastly improves diversity estimates from Illumina amplicon  
553 sequencing. *Nature Methods* **10**, 57–59 (2013).
- 554 34. Wang, J. *et al.* Minimizing spurious features in 16S rRNA gene amplicon sequencing. (PeerJ Inc.,  
555 2018). doi:[10.7287/peerj.preprints.26872v1](https://doi.org/10.7287/peerj.preprints.26872v1)

- 556 35. Janssen, S. *et al.* Phylogenetic Placement of Exact Amplicon Sequences Improves Associations with  
557 Clinical Information. *mSystems* **3**,
- 558 36. DeSantis, T. Z. *et al.* Greengenes, a Chimera-Checked 16S rRNA Gene Database and Workbench  
559 Compatible with ARB. *Applied and Environmental Microbiology* **72**, 5069–5072 (2006).
- 560 37. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*  
561 **12**, 2825–2830 (2011).
- 562 38. Mandal, S. *et al.* Analysis of composition of microbiomes: A novel method for studying microbial  
563 composition. *Microbial Ecology in Health and Disease* **26**, 27663 (2015).
- 564 39. Watts, S. C., Ritchie, S. C., Inouye, M. & Holt, K. E. FastSpar: Rapid and scalable correlation  
565 estimation for compositional data. *bioRxiv* 272583 (2018). doi:[10.1101/272583](https://doi.org/10.1101/272583)
- 566 40. Chan, C. L. *et al.* The microbiome of otitis media with effusion. *The Laryngoscope* **126**, 2844–2851  
567 (2016).
- 568 41. Yu, Y., Lee, C., Kim, J. & Hwang, S. Group-specific primer and probe sets to detect methanogenic  
569 communities using quantitative real-time polymerase chain reaction. *Biotechnology and Bioengineering*  
570 **89**, 670–679 (2005).
- 571 42. Kluyver, T. *et al.* Jupyter Notebooks a publishing format for reproducible computational workflows.  
572 in *Positioning and Power in Academic Publishing: Players, Agents and Agendas* (eds. Loizides, F. &  
573 Schmidt, B.) 87–90 (IOS Press, 2016). doi:[10.3233/978-1-61499-649-1-87](https://doi.org/10.3233/978-1-61499-649-1-87)
- 574 43. Oliphant, T. E. Python for Scientific Computing. *Computing in Science & Engineering* **9**, 10–20  
575 (2007).