

1 **A comparison of epithelial cell content of oral samples estimated using cytology and**
2 **DNA methylation**

3

4 **Author list** (suggested)

5 Yen Ting Wong¹, Michael A Tayeb², Laurence B Lovat³, Andrew E Teschendorff^{4,5}, Rafal
6 Iwasiow², Jeffrey M Craig^{1,6,7}

7

8 **Affiliations**

9 ¹ IMPACT Strategic Research Centre, School of Medicine, Barwon Health, Geelong, Vic,
10 Australia

11 ² DNA Genotek Inc., Ottawa, ON, Canada

12 ³ Division of Surgery & Interventional Science, UCL, London WC1E 6BT, UK

13 ⁴ CAS Key Laboratory of Computational Biology, CAS-MPG Partner Institute for
14 Computational Biology, Shanghai Institutes for Biological Sciences, 320 Yue Yang Road,
15 Shanghai 200031, PR China.

16 ⁵ UCL Cancer Institute, 72 Huntley Street, University College London, London WC1E 6BT,
17 UK.

18 ⁶ Murdoch Children's Research Institute, Department of Paediatrics, The University of
19 Melbourne, Royal Children's Hospital, Melbourne, Vic., Australia.

20 ⁷ Corresponding author

21 jeffrey.craig@deakin.edu.au

22 **Abstract**

23 Saliva and buccal samples are popular for epigenome wide association studies (EWAS) due
24 to their ease of collection compared and their ability to sample a different cell lineage
25 compared to blood. As these samples contain a mix of white blood cells and buccal epithelial
26 cells that can vary within a population, this cellular heterogeneity may confound EWAS. This
27 has been addressed by including cellular heterogeneity obtained through cytology at the time
28 of collection or by using cellular deconvolution algorithms built on epigenetic data from
29 specific cell types. However, to our knowledge, the two methods have not yet been
30 compared. Here we show that the two methods are highly correlated in saliva and buccal
31 samples ($R = 0.84$, $P < 0.0001$) by comparing data generated from cytological staining and
32 Infinium MethylationEPIC arrays and the EpiDISH deconvolution algorithm from buccal and
33 saliva samples collected from twenty adults. In addition, by using an expanded dataset from
34 both sample types, we confirmed our previous finding that age has a significant negative
35 correlation with epithelial cell proportion in both sample types. However, children and adults
36 showed a large within-population variation in cellular heterogeneity. Our results validate the
37 use of the EpiDISH algorithm in estimating the effect of cellular heterogeneity in EWAS and
38 showed DNA methylation generally underestimates the epithelial cell content obtained from
39 cytology.

40

41 **Keywords** Buccal, Cytology, DNA methylation, cell-type heterogeneity, epithelial cell,
42 EWAS, saliva

43 **1. Introduction**

44 Cellular heterogeneity is a major potential confounder of epigenome-wide association studies
45 (EWAS) due to the cell type-specific state of DNA methylation. This is particularly the case in
46 oral samples, which are a mixture of epithelial cells from the ectoderm germ cell lineage and
47 immune cells from the mesoderm lineage (1-3). We and others have found that cellular
48 heterogeneity in oral samples is influenced by the method of sample collection, with buccal
49 swabs containing a much higher proportion of epithelial cells than saliva (2, 3). We have also
50 shown that epithelial cell proportion is also strongly influenced by age and oral health status
51 (3). Deconvolution of cellular heterogeneity can be achieved by measuring the proportion of
52 each cell type using cytology of collected cells or through algorithms that use DNA methylation
53 data from specific cell types to generate estimates (1, 2). Such measures can then be used in
54 EWAS models to correct for cellular heterogeneity. There are studies comparing cytology
55 estimates of tumor purity to DNA methylation based estimates and mRNA expression based
56 {Chakravarthy A #2018, Aran D # 2015} However, to our knowledge, no study has compared
57 these two methods for oral samples. We aimed to compare epithelial cell content of buccal
58 samples via ORAcollect•DNA kits and saliva obtained via passive drool collected in
59 Oragene•DNA kits, measured using cytology and estimated with the reference based EPIDISH
60 algorithm (2). We hypothesised that estimations of epithelial cell content would be highly
61 correlated between the two methods. In a sub-study, using customized ORAcollect•DNA
62 collection instructions, we compared two similar methods of collection differing in collection
63 site and duration.

64

65 **2 Materials and methods**

66 **2.1 Participants**

67 Twenty adult volunteers from Deakin University provided informed consent to collect one
68 saliva sample and two buccal samples. Ethics approval was granted by the Human Research
69 Ethics Committee of the Royal Children's Hospital, Melbourne (#33174) and Deakin
70 University (2018-368). All methods were performed according to relevant protocols and
71 regulations. Participants also completed an oral health questionnaire, which included questions
72 about mouth injuries, oral infections, medications and smoking status (Supplementary
73 Methods).

74 **2.2 Oral sampling**

75 Oral samples were obtained from participants under supervision of the research team.
76 Participants were advised not to smoke, chew gum, or consume anything apart from water for
77 the 30 minutes prior to providing samples. Ten minutes prior to sample collection, they were
78 asked to rinse their mouth with water. Saliva samples were collected unstimulated via passive
79 drool for three to five minutes to allow sufficient time to collect to the fill line (2mL) of
80 Oragene•DNA collection devices (OG; DNA Genotek Inc, Ottawa, Canada). One hundred
81 microliters of saliva were then smeared onto a microscope slide and immediately fixed with
82 95% ethanol for 10 minutes and left to dry at room temperature. Oragene DNA-stabilising
83 chemistry contained within the device was then released into the remaining sample. Following
84 collection of saliva, two samples were collected from participants using ORAcollect•DNA
85 (OC, DNA Genotek Inc, Ottawa, Canada), a sponge-tipped oral sample collection kit,
86 sequentially using two collection methods. In the first (OCA) participants gently rubbed the
87 sponge ten times in a back-and-forth motion in the furrow between their lower teeth and inner
88 cheek on one side of their mouth. In the second (OCB), the sponge was rubbed up and down
89 against the inside of the cheek twenty times then rubbed ten seconds in a back-and-forth motion
90 in the furrows between their upper and lower right teeth and inner cheek on the opposite side
91 of their mouth. Each sponge was wiped along the length of a standard size microscope slide

92 and fixed as outlined for saliva. The sponge was then inserted into the ORACollect•DNA tube
93 containing DNA stabilising chemistry, capped tightly and mixed by inversion 15 times.

94 **2.3 Slide staining and microscopy**

95 Slides were stained using Diff-Quik as detailed elsewhere (4). All slides were deidentified and
96 analysed by two observers. Cell types were counted via bright field microscopy at 100x
97 magnification in regions with adequate cell density. Counts were used if the discrepancy
98 between observers was less than 10% of the total count of count for each cell type. For counts
99 which discrepancy between observers in more than 10% will be discarded and re-count again.
100 A minimum of 50 epithelial cells and a 100 cells total was counted. Cells were scored as
101 epithelial cells or immune cells, the latter including segmented cells, lymphocytes and
102 monocytes (3).

103 **2.4 DNA extraction**

104 Genomic DNA was extracted from 0.5 mL of each oral sample using ethanol precipitation via
105 prepIT•L2P kits (DNA Genotek Inc, Ottawa, Canada) following the manufacturer's protocol.
106 DNA concentration was measured using PicoGreen (Thermo Fisher Scientific, Canada) in a
107 SpectraMax M2 plate-based fluorimeter (Molecular Devices, CA, USA). DNA quality was
108 measured using a TapeStation (Agilent, Santa Clara, United States).

109 **2.5 DNA Methylation arrays**

110 Following genomic DNA extraction from all the samples, these genomic DNA samples were
111 treated with bisulphite to convert unmethylated cytosine into uracil and transformation of uracil
112 into thymine by amplification. Genome-wide analysis of DNA methylation was assessed using
113 Infinium MethylationEPIC arrays (Illumina, CA, USA) with probes of over 850,000
114 methylation sites at the GenoFIND Genomic Service Lab (DNA Genotek Inc, Ottawa, Canada).
115 Hybridization and scanning were performed according to manufacturer's instructions.

116 **2.6 Pre-processing of Illumina Infinium array data**

117 MethylationEPIC array analysis was performed using the R statistical programming language
118 (www.R-project.org) and Bioconductor packages (5). Raw intensity data (IDAT) files were
119 imported into R (3.6.3; <http://cran.r-project.org/>). Data quality was assessed using the *minfi*
120 (v1.34.0) Bioconductor package (5). The MethylationEPIC probes were filtered by removing
121 those with poor signal to noise ratio (mean detection p-value of >0.01), cross-reactivity to
122 multiple genomic locations, containing a single nucleotide polymorphism at the CpG site, or
123 map to sex chromosomes (6). Data was then subjected to subset-quantile within array
124 normalisation (SWAN), (7)) and between-array normalisation (SQN) (8). The
125 HEpiDISH/EpiDISH and Robust partial correlation (RPC) algorithms were applied to estimate
126 proportions of epithelial, fibroblast and immune cells from MethylationEPIC array data (9).

127 **2.7 Data analysis**

128 Descriptive statistical analyses were conducted on the age of the participant and proportions of
129 epithelial and immune cells. The assumption of normality of the independent and dependent
130 variables for each cell type was tested using the Shapiro-Wilk test. A Kruskal-Wallis ANOVA
131 analysis was conducted to test for statistically significant differences in cell proportion and
132 DNA yield between collection methods OCA, OCB and OG. In a post-hoc test, the Dunnett's
133 test with Bonferroni correction was applied to identify the relatively small but significant
134 differences among collection methods. Variables collected with insufficient number and
135 information will not be included in the statistical analysis.

136 Percentage of epithelial cells and estimated cell-type fractions from EpiDISH were graphed
137 using box and whisker plots, which included information on interquartile range (boxes, 25th to
138 75th percentiles, boxes), median (horizontal lines), data within 5th-95th percentiles (whiskers),
139 outliers (circles), and mean (crosses). The proportion of epithelial cells in oral samples

140 estimated from cytology and DNA methylation was tested using Pearson correlation
141 coefficient. To investigate the age effect on epithelial cell content estimated using DNA
142 methylation, the buccal and saliva sample data from this study was analysed along with seven
143 of our other studies, three published (3, 10, 11) and four unpublished. These cohorts' details
144 included to investigate the age effect on epithelial cell content was described in supplementary
145 method.

146

147 **3 Results**

148 **3.1 Determination of epithelial and immune cell proportions using cytology**

149 Slides from all twenty adults (mean age 26.9 years, range 21 to 48 years, 60% female) were
150 analysable i.e. had sufficient cells for analysis. Seven individuals reported recent gum bleeding
151 within the seven days preceding their collection day. Examples of microscopic fields of view
152 are shown in **Figure 1**. Epithelial cells were large, with low nuclear to cytoplasmic ratio and
153 immune cells were much smaller with a high nuclear to cytoplasmic ratio. Immune cells
154 included granulocytes with segmented nuclei, lymphocytes with round, dense nuclei
155 surrounded by cytoplasm, monocytes with kidney-shaped nuclei. Between two and twenty
156 fields of view at 100x magnification were required to score the minimum number of cells.

157

158 *Figure 1 around here*

159

160 Results for estimations of epithelial cell proportions determined by cytology and DNA
161 methylation analysis are shown in **Supplementary Table 1** and **Figure 2**. The mean proportion
162 of epithelial cells in saliva (58%, SD 17.1%), was significantly lower with than sponge
163 collection methods OCA (86.0%, SD 9.9%) and OCB (87.0%, SD 11.2%), $p < 0.0001$. A

164 28.5% mean difference with SD 6.5% in compared saliva to cheek swab methods. There was
165 no evidence for a difference in epithelial cell proportions between OCA and OCB ($p = 0.6$).
166 There was also no evidence of an influence of recent gum bleeding (p value =0.5) and sex (p
167 value = 0.9) on epithelial cell proportion across all methods of sampling; results for individual
168 oral collection methods were similar.

169

170 *Figure 2 around here*

171

172 **3.2 Determination of epithelial cell proportions using DNA methylation analysis**

173 Saliva samples showed a significantly higher mean of total DNA yield per mL (33.7 μ g, SD
174 24.2 μ g) compared to oral sponge collection methods OCA (4.1 μ g, SD 1.57 μ g) and OCB
175 (5.9 μ g, SD 2.71 μ g), $p < 0.0001$ for both comparisons (**Figure 3**). Although DNA yield was
176 approximately 1.7x higher in OCB compared to OCA, this difference was not significant ($p =$
177 0.083).

178

179 *Figure 3 around here*

180

181 We next used the EpiDISH and robust partial correlation (RPC) algorithm on Infinium
182 MethylationEPIC data to estimate cell type proportions. Although this method calculates
183 proportion of epithelial, immune and fibroblast cell types, we found that the proportion of
184 fibroblasts was negligible (mean = 0.4%) (**Supplementary Table 1**). As this meant that the
185 proportion of immune and epithelial cells had a correlation of -1.0, we limited our analysis to
186 the latter. As with cytology, the mean proportion of buccal epithelial cells determined by DNA

187 methylation in saliva (25.4%, SD 17.1%), was significantly lower than cheek swab methods
188 OCA (69.5%, SD 18.8%) and OCB (75.5%, SD 17.0%), $p < 0.0001$ (**Figure 2**). A 47.1% mean
189 difference with SD 0.7% in compared saliva to cheek swab methods. There was no evidence
190 for a difference in epithelial cell proportion between OCA and OCB ($p = 0.11$).

191

192 **3.3 Comparison between epithelial cell proportions estimated using cytology and** 193 **DNA methylation**

194 To address our hypothesis that proportions of epithelial cells present in oral samples estimated
195 using DNA methylation analysis represented the cell proportions as measured by cytology, we
196 pooled all samples and compared both methods (**Figure 4**). The two methods were strongly
197 correlated ($R = 0.84$, $P < 0.0001$). However, the intercept of the line of best fit (methylation %
198 = $[1.32 \times \text{cytology \%}] - 45\%$) on the x axis was 34%. A 20% mean difference of DNA
199 methylation (SD 17.7%, IQR 63.7%) compared to cytology (SD 12.7%, IQR 50.8%)
200 (**Supplementary Table 1**). Methods correlated similarly in buccals ($R = 0.75$, $P < 0.0001$) and
201 saliva ($R = 0.72$, $P < 0.0001$) (**Supplementary Figure 1**).

202

203 *Figure 4 around here*

204

205 **3.4 An age effect on epithelial cell content in saliva**

206 In our previous paper, we observed that epithelial cell content of buccal swabs and saliva was
207 lower in adults compared to children (3). To investigate a possible age effect using epithelial
208 cell content estimated using DNA methylation, we combined buccal swab and saliva data from
209 this study with seven of our other Infinium array studies, including three published (3, 10, 11)
210 and four unpublished (**Figure 5, Supplementary Table 2**). We found a moderate negative

211 correlation between age and epithelial cell content estimated by DNA methylation ($R = -0.72$,
212 $p < 0.0001$, 0.59% of epithelial cell content estimated from DNA methylation), with age
213 accounting for 14% of the variation in epithelial cell content. We found a stronger relationship
214 in buccals ($R = -0.85$, $P < 0.0001$) compared to saliva ($R = -0.28$, $P < 0.0001$) (**Supplementary**
215 **Figure 2**).

216

217 **4 Discussion**

218 **4.1 Influences on epithelial content of oral samples**

219 Buccal and saliva samples have a proven utility for epigenomics (12, 13) and other cell-based
220 omics (14, 15). As these samples are mixtures of epithelial and immune cells, deconvolution
221 of these cellular mixtures is of utmost importance. Although cellular deconvolution
222 algorithms based on reference sample types have been applied to epigenomic studies (2), to
223 our knowledge, the validity of such algorithms has not yet been tested using cytology of
224 primary samples. We aimed to address this issue.

225 Our cytological analysis of adults with a mean age of 26 years showed that the epithelial
226 content of ORAcollect•DNA (OC) samples was 86.5%, similar to the 83.4% we previously
227 obtained using Copan flocced swabs in adults 16 years older (3). In the present study,
228 epithelial cell content of saliva, but not buccal samples, was significantly higher than in our
229 earlier study, which agrees with our previous finding that age has a much greater effect on
230 saliva than on buccal samples. We also found that the epithelial content of ORAcollect•DNA
231 collected samples was around 47% higher than that of saliva. This difference was 11% larger
232 than that of our previous study, which may also reflect an age effect.

233 Our findings also suggest that the type of buccal collector has minimal influence on the
234 proportion of epithelial cells collected and this may also be one reason why increasing

235 collection time for OC sponge did not increase the proportion of epithelial cells collected, nor
236 did it significantly increase DNA yield. However, future, larger studies are needed to further
237 test our hypotheses. However, there may be a danger that longer collection times penetrate
238 blood capillaries within the inner cheek, which would increase the proportion of immune
239 cells, which may negate any rise in epithelial cell numbers.

240 We found no evidence that recent gum bleeding influenced the proportions of epithelial cells
241 with either mode of sample collection. This disproved our hypothesis that gum bleeding
242 would decrease proportions of epithelial cells, possibly because the severity and temporary
243 nature of bleeding may be insufficient to cause a significant impact on immune cell numbers
244 and possibly because of our relatively small sample size.

245 **4.2 Comparison of epithelial cell proportions using cytology and DNA methylation**

246 Using the EpiDISH algorithm (2) on DNA methylation data generated by Infinium
247 MethylePIC arrays, we estimated that epithelial cell proportion was lower in saliva compared
248 to OC-collected samples by an average of 47%, a larger magnitude than that shown using
249 cytology. Across all samples, the correlation between the two methods of epithelial cell
250 estimation was very high ($R=0.84$). Taken together, these findings prove our primary
251 hypothesis and imply that post hoc deconvolution accurately estimates cellular heterogeneity
252 in oral samples. In a study comparing proportions of various blood cell types estimated using
253 flow cytology and a blood-specific DNA methylation-based algorithm, a wide range of
254 correlations, between 0.51 and 0.97 were observed (16). Our line of best fit showed that DNA
255 methylation underestimated the epithelial cell content determined using cytology by 34% at
256 0% epithelial cells and by 13% at 87% epithelial cells (**Figure 4**). This is larger than the over-
257 or under-estimations of up to 10% observed in the previous study of blood (16). This
258 discrepancy could be for a number of reasons. The reference dataset was derived from
259 Illumina InfiniumHM450 array data from 11 different epithelial cell lines (2) which may not

260 accurately represent buccal epithelial cells. We also cannot rule out the possibility that buccal
261 and immune cell types may have been differentially applied to slides prior to cytological
262 examination. Nevertheless, the high correlation between epithelial cell proportions based on
263 cytology and DNA methylation should still be sufficient to use the latter to generate
264 estimations across a set of biosamples for adjustment within EWAS.

265 **4.3 Investigating an age effect on the proportion of immune cells in buccal swabs**

266 We found a significant negative correlation of epithelial content in buccal swabs and saliva
267 with age (**Figure 5, Supplementary Figure 2**). This agrees with our previous study that
268 showed an effect in the same direction with buccal swabs and saliva in children and adults
269 (3). In our earlier study, we showed that epithelial cell proportion was significantly lower in
270 children with gingivitis. As gingivitis and other oral inflammatory pathologies such as
271 periodontitis increase in prevalence with age, this may result in an increase in immune cell
272 content of oral samples and a corresponding decrease in epithelial cell content.

273 **4.4 Strengths, limitations and future studies**

274 To our knowledge, this study is the first to analyse the correlation between cell proportions
275 in oral samples estimated using cytology and DNA methylation. Another strength is our
276 longitudinal analysis showing a decline of epithelial content of buccal swabs and saliva with
277 age. However, our sample size (n=20 for both buccal samples and saliva; all studies, n=579)
278 is relatively small, although our sample size for the study of age effects (n=579) was much
279 larger. Future, larger-scale studies that compare estimates of cell proportion using both
280 cytology and DNA methylation are required to validate our findings. Such studies should
281 include a wider age-group and measures of oral health.

282

283 **5 Acknowledgements**

284 The authors thank Ziad Marroushi and Christina Dillane for technical assistance and
285 Elizabeth Firth for her comments on the manuscript. We thank Evie Muggli, Jane Halliday,
286 Lata Vadlamudi and Tim Silk for permission to include unpublished data.

287

288 **6 Disclosure of interest**

289 The authors report no conflict of interest.

290 **7 References**

- 291 1. Teschendorff AE, Zheng SC. Cell-type deconvolution in epigenome-wide association
292 studies: a review and recommendations. *Epigenomics*. 2017;9(5):757-68.
- 293 2. Zheng SC, Breeze CE, Beck S, Dong D, Zhu T, Ma L, et al. EpiDISH web server:
294 Epigenetic Dissection of Intra-Sample-Heterogeneity with online GUI. *Bioinformatics*. 2019.
- 295 3. Theda C, Hwang SH, Czajko A, Loke YJ, Leong P, Craig JM. Quantitation of the
296 cellular content of saliva and buccal swab samples. *Sci Rep*. 2018;8(1):6944.
- 297 4. Yang GC, Papellas J, Wu HC, Waisman J. Application of Ultrafast Papanicolaou stain
298 to body fluid cytology. *Acta cytologica*. 2001;45(2):180-5.
- 299 5. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, et
300 al. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium
301 DNA methylation microarrays. *Bioinformatics*. 2014;30(10):1363-9.
- 302 6. Chen Y-a, Lemire M, Choufani S, Butcher DT, Grafodatskaya D, Zanke BW, et al.
303 Discovery of cross-reactive probes and polymorphic CpGs in the Illumina Infinium
304 HumanMethylation450 microarray. *Epigenetics*. 2013;8(2):203-9.

- 305 7. Maksimovic J, Gordon L, Oshlack A. SWAN: Subset-quantile within array
306 normalization for illumina infinium HumanMethylation450 BeadChips. *Genome biology*.
307 2012;13(6):R44.
- 308 8. Touleimat N, Tost J. Complete pipeline for Infinium® Human Methylation 450K
309 BeadChip data processing using subset quantile normalization for accurate DNA methylation
310 estimation. *Epigenomics*. 2012;4(3):325-41.
- 311 9. Teschendorff AE, Breeze CE, Zheng SC, Beck S. A comparison of reference-based
312 algorithms for correcting cell-type heterogeneity in Epigenome-Wide Association Studies.
313 *BMC Bioinformatics*. 2017;18(1):105.
- 314 10. Martino D, Loke YJ, Gordon L, Ollikainen M, Cruickshank MN, Saffery R, et al.
315 Longitudinal, genome-scale analysis of DNA methylation in twins from birth to 18 months of
316 age reveals rapid epigenetic change in early life and pair-specific effects of discordance.
317 *Genome Biol*. 2013;14(5):R42.
- 318 11. Mohandas N, Loke YJ, Hopkins S, Mackenzie L, Bennett C, Berkovic SF, et al.
319 Evidence for type-specific DNA methylation patterns in epilepsy: a discordant monozygotic
320 twin approach. *Epigenomics*. 2019;11(8):951-68.
- 321 12. Langie SA, Moisse M, Declerck K, Koppen G, Godderis L, Vanden Berghe W, et al.
322 Salivary DNA Methylation Profiling: Aspects to Consider for Biomarker Identification. *Basic*
323 *Clin Pharmacol Toxicol*. 2016.
- 324 13. Farah R, Haraty H, Salame Z, Fares Y, Ojcius DM, Said Sadier N. Salivary
325 biomarkers for the diagnosis and monitoring of neurological diseases. *Biomed J*.
326 2018;41(2):63-87.
- 327 14. Hassaneen M, Maron JL. Salivary Diagnostics in Pediatrics: Applicability,
328 Translatability, and Limitations. *Front Public Health*. 2017;5:83.

- 329 15. Shah S. Salivaomics: The current scenario. J Oral Maxillofac Pathol. 2018;22(3):375-
330 81.
- 331 16. Gervin K, Page CM, Aass HC, Jansen MA, Fjeldstad HE, Andreassen BK, et al. Cell
332 type specific DNA methylation in cord blood: A 450K-reference data set and cell count-
333 based validation of estimated cell type composition. Epigenetics. 2016;11(9):690-8.

334

335 **Figure Legends**

336

337 **Figure 1. Examples of cellular morphology in oral samples.** Representative fields of view
338 from Diff-Quik staining of (A) saliva, 100x magnification and (B) OCA buccal sample, 400x
339 magnification. Both samples contain large epithelial cells (Epi) with dense nuclei, and smaller
340 immune cells, exemplified by lymphocytes (Lym), segmented cells (Seg) and monocytes
341 (Mono).

342

343 **Figure 2. Comparison of the percentage proportion of epithelial cells in oral samples,**
344 **estimated using cytology and DNA methylation arrays, collected using three different**
345 **methods (OCA, OCB and OG) estimated.** Means are indicated with crosses. The p value of
346 percentage of epithelial cell between OCA and OCB is $p > 0.05$. The p value of between buccal
347 sample collection (OCA and OCB) compared to saliva (OG) is $p < 0.0001$.

348

349 **Figure 3: Range of DNA yields for each oral sample type.** Box and whisker plots from saliva
350 (OG) and the two methods of buccal sample collection (OCA and OCB). Means are indicated
351 with an X.

352

353 **Figure 4: Comparison of the proportion of epithelial cells in oral samples estimated from**
354 **cytology and DNA methylation arrays.**

355

356 **Figure 5: Epithelial cell content of oral samples as a function of age in six studies.**

357

358 **Supplementary Figure 1: Comparison of the proportion of epithelial cells in oral samples**
359 **estimated from cytology and DNA methylation arrays. (A) Data from saliva epithelial cells;**
360 **(B) Data from buccal epithelial cells.**

361

362 **Supplementary Figure 2: Epithelial cell content of buccal and saliva samples as a function**
363 **of age. (A) Buccal data is from five studies (n = 344); (B) saliva data is from three studies (n**
364 **= 234).**

365

366 **Supplementary Table 1: Estimation of cell proportions and for each collection method.**

367

368 **Supplementary Table 2: Epithelial proportions of the buccal and saliva sample data from**
369 **this study and five of our other studies (n= 753).**

370

371 **8 Supplementary methods and results**

372 **8.1 Questionnaire**

373 Before sample collection, each adult participant was asked to complete an oral health
374 questionnaire. The participant was given a unique ID number, other details included birth year,
375 collection date and time, and sex were recorded. The oral health of each participant was
376 recorded via a questionnaire which asked about whether they had bleeding gums when brushing
377 their teeth, mouth ulcers, other mouth lesions, a cold, a sore throat, or other mouth infection
378 during the past week. Participants were also asked whether they used an inhaler, took
379 antibiotics, anti-inflammatories or blood thinners and whether they smoked.

380 **8.2 Study cohorts**

381 As stated in the method, to investigate the age effect on epithelial cell content estimated using
382 DNA methylation, the buccal and saliva sample data from this study was analysed along with
383 three published (3, 10, 11) and four unpublished of our other studies. First, the child
384 participants are recruited from part of the Peri/postnatal Epigenetic Twins Study (PETS)
385 cohort, an Australian twin birth research study based in Melbourne. These participants are
386 involved in a longitudinal study of DNA methylation at birth (n= 29, age 0 year) and age 18
387 months (n= 24, age 1.5 years) from buccal swabs (10). Second, Theda et al. study, their saliva
388 and buccal samples were collected from ten pairs of twins (n= 20, age range 6.4-7.1 years)
389 from PETS cohort and adult volunteers (n = 23, age range 23-59 years) (3). Third, an
390 epilepsy cohort consisted of monozygotic twin pairs (n= 28) age range of 14-67, who were
391 discordant for epilepsy without a known acquired cause (11).

392

393 The four unpublished studies were an Australian longitudinal study of community-based
394 children with ADHD cohort with match age (n=175, 10.4 years). A population subset of the
395 AQUA (Asking QUestions about Alcohol in pregnancy) cohort study (n=187 of neonatal
396 cheek swabs, match age at 0 year). Two data sets collected from buccal samples in the year

397 2015 (n= 63, age range 28-87 years) and year 2019 (n= 111, age range 35-90 years) provided
398 by a collaborator from University College London.