

1 **Investigating the effect of sexual behaviour on oropharyngeal**  
2 **cancer risk: a methodological assessment of Mendelian**  
3 **randomization**

4 Mark Gormley, BDS\* <sup>1,2,3</sup>, Tom Dudding, PhD<sup>1,2,3</sup>, Linda Kachuri, PhD<sup>4</sup>, Kimberley Burrows,  
5 PhD<sup>1,3</sup>, Amanda HW Chong, PhD<sup>1,3</sup>, Richard M Martin, PhD<sup>1,3,5</sup>, Steven Thomas, PhD<sup>2,5</sup>,  
6 Jessica Tyrrell, PhD<sup>6</sup>, Andrew R Ness, PhD<sup>5</sup>, Paul Brennan, PhD<sup>7</sup>, Marcus R Munafò, PhD<sup>1,8</sup>,  
7 Miranda Pring, PhD<sup>2</sup>, Stefania Boccia, PhD<sup>9,10</sup>, Andrew F Olshan, PhD<sup>11</sup>, Brenda Diergaarde,  
8 PhD<sup>12</sup>, Rayjean J Hung, PhD<sup>13,14</sup>, Geoffrey Liu, MD<sup>14,15</sup>, Eloiza Tajara, PhD<sup>16</sup>, Patricia Severino,  
9 PhD<sup>17</sup>, Tatiana N Toporcov<sup>18</sup>, Martin Lacko, MD PhD<sup>19</sup>, Tim Waterboer, PhD<sup>20</sup>, Nicole  
10 Brenner, PhD<sup>20</sup>, George Davey Smith, FRS<sup>1,3</sup>, Emma E Vincent, PhD<sup>1,3,21</sup>, Rebecca C  
11 Richmond, PhD<sup>1,3</sup>

## 12 **Affiliations**

13 <sup>1</sup> MRC Integrative Epidemiology Unit, Population Health Sciences, Bristol Medical School, University of Bristol,  
14 Bristol, UK.

15 <sup>2</sup> Bristol Dental Hospital and School, University of Bristol, Bristol, UK.

16 <sup>3</sup> Department of Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK.

17 <sup>4</sup> Department of Epidemiology & Biostatistics, University of California San Francisco, San Francisco, USA.

18 <sup>5</sup> University Hospitals Bristol and Weston NHS Foundation Trust National Institute for Health Research Bristol  
19 Biomedical Research Centre, University of Bristol, Bristol, UK.

20 <sup>6</sup> University of Exeter Medical School, RILD Building, RD&E Hospital, Exeter, UK.

21 <sup>7</sup> Genetic Epidemiology Group, World Health Organization, International Agency for Research on Cancer, Lyon,  
22 France.

23 <sup>8</sup> School of Psychological Science, Faculty of Life Sciences, University of Bristol, Bristol, UK.

24 <sup>9</sup> Section of Hygiene, University Department of Life Sciences and Public Health, Università Cattolica del Sacro  
25 Cuore, Roma, Italia.

26 <sup>10</sup> Department of Woman and Child Health and Public Health - Public Health Area, Fondazione Policlinico  
27 Universitario A. Gemelli IRCCS, Roma, Italy.

28 <sup>11</sup> Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina, US.

29 <sup>12</sup> Department of Human Genetics, Graduate School of Public Health, University of Pittsburgh, and UPMC  
30 Hillman Cancer Center, Pittsburgh, US.

31 <sup>13</sup> Prosserman Centre for Population Health Research, Lunenfeld-Tanenbaum Research Institute, Sinai Health  
32 System, Toronto, Canada.

33

34 <sup>14</sup> Dalla Lana School of Public Health, University of Toronto, Toronto, Canada.

35 <sup>15</sup> Princess Margaret Cancer Centre, Toronto Canada

36 <sup>16</sup> Department of Molecular Biology, School of Medicine of São José do Rio Preto, São Paulo, Brazil.

37 <sup>17</sup> Albert Einstein Research and Education Institute, Hospital Israelita Albert Einstein, São Paulo, Brazil.

38 <sup>18</sup> Department of Epidemiology, School of Public Health, University of São Paulo, Brazil.

39 <sup>19</sup> Department of Otorhinolaryngology and Head and Neck Surgery, Research Institute GROW, Maastricht  
40 University Medical Center, Maastricht, The Netherlands.

41 <sup>20</sup> Infections and Cancer Epidemiology, Deutsches Krebsforschungszentrum, Heidelberg, Germany.

42 <sup>21</sup> School of Cellular and Molecular Medicine, University of Bristol, Bristol, UK.

43 \* Corresponding author: [mark.gormley@bristol.ac.uk](mailto:mark.gormley@bristol.ac.uk)

44

## 45 **Abstract**

46

47 Human papilloma virus infection is known to influence oropharyngeal cancer (OPC) risk,  
48 likely via sexual transmission. However, sexual behaviour has been correlated with other  
49 risk factors including smoking and alcohol, meaning independent effects are difficult to  
50 establish. Here we evaluate aspects of sexual behaviour in relation to the risk of OPC (2,641  
51 cases and 6,585 controls), using genetic variants associated with age at first sex (AFS) and  
52 number of sexual partners (NSP) to perform Mendelian randomization (MR) analyses. While  
53 univariable MR showed a causal effect of both later AFS and increasing NSP on OPC, results  
54 attenuated in the multivariable models (AFS IVW OR 0.7, 95%CI 0.4, 1.2,  $p=0.21$ ; NSP IVW  
55 OR 0.9, 95%CI 0.5 1.7,  $p=0.76$ ). We also found evidence for correlated pleiotropy in the  
56 genetic instruments for sexual behaviour, emphasising the need for multivariable  
57 approaches when performing MR of complex behavioural traits and the triangulation of  
58 evidence.

59

60

61

62

63

64

## 65 Introduction

66

67 Head and neck squamous cell carcinoma (HNSCC) is a heterogeneous disease <sup>1</sup>, which can  
68 originate from the mucosa of the oral cavity, oropharynx and larynx. Worldwide, there are  
69 over half a million incident cases each year, resulting in more than 200,000 deaths annually  
70 <sup>2</sup>. While using tobacco products and consuming alcohol are well-established risk factors  
71 across all HNSCC subsites, oral human papilloma virus (HPV) infection has been identified as  
72 another risk factor, particularly within the oropharyngeal subsite <sup>3-6</sup>. In developed countries  
73 such as the USA, 60-70% of oropharyngeal cancer (OPC) cases are reported to be HPV-  
74 positive <sup>7</sup>, compared to only around 5% of all oral cancer (OC) cases. Oncogenic HPV type-16  
75 (HPV16) is the most common type found in approximately 90% of HPV-positive  
76 oropharyngeal tumours <sup>8-10</sup>. Antibodies against HPV oncoproteins may be potential  
77 biomarkers, with case-control studies demonstrating seropositivity for late (L1) and early  
78 (E1, E2, E4, E6, E7) HPV16 proteins were associated with oropharyngeal cancer <sup>11-14</sup>.

79

80 HPV is thought to be sexually transmitted via oro-genital contact <sup>9,15-20</sup> and may enter the  
81 oropharyngeal mucosa via abrasions in the reticulated tonsillar epithelium <sup>21</sup>. One large  
82 pooled analysis investigating the role of sexual behaviour in HNSCC showed an increased  
83 risk of OPC with having a history of six or more lifetime sexual partners (OR 1.3, 95%  
84 confidence intervals (95%CI), 1.0, 1.5) and four or more oral sex partners (OR 2.3, 95%CI 1.4,  
85 3.6). A positive association was observed among men who had oral sex (OR 1.6, 95%CI 1.1,  
86 2.3) and those with an earlier age at sexual debut (OR 2.4, 95%CI 1.4, 5.1) <sup>15</sup>. Conversely,  
87 there was no association reported between oral sex practice and head and neck cancer in a  
88 more recent meta-analysis of 17 studies (OR 1.1, 95%CI: 0.9, 1.4), suggesting inconsistency

89 in these findings, although 12 of these 17 studies failed to stratify by oral and oropharyngeal  
90 subsite<sup>22</sup>. Furthermore, associations have typically been investigated using case-control  
91 studies<sup>5</sup>, with self-reported sexual behaviour which may be subject to recall bias and  
92 misreporting. Positive associations have also been found between sexual behaviour,  
93 sexually transmitted infections and other risk factors for HNSCC, including alcohol  
94 consumption, indicating the possibility of residual confounding exists<sup>23</sup>.

95

96 Mendelian randomization (MR) is an approach to causal analysis which attempts to  
97 overcome shortcomings of conventional observational studies by using single nucleotide  
98 polymorphisms (SNPs) which are randomly allocated at conception and known to be reliably  
99 associated with modifiable risk factors of interest. These genetic instruments can be used to  
100 estimate the effects of risk factors on disease outcomes, in this case sexual behaviours on  
101 oropharyngeal cancer<sup>24,25</sup>, which are less prone to unidentified confounding or reverse  
102 causation than conventional epidemiological analysis. Large-scale genome-wide association  
103 studies (GWAS) have been performed for sexual behaviour traits, including number of  
104 sexual partners (NSP)<sup>26,27</sup> and age at first sex (AFS)<sup>28</sup>. MR makes three key assumptions in  
105 that the genetic instrument (i) is robustly associated with the risk factor (i.e., 'relevance'),  
106 (ii) does not share a common cause with the outcome (i.e., 'exchangeability'), and (iii)  
107 affects the outcome only through the risk factor (i.e., 'exclusion restriction principle') to  
108 check for genetic pleiotropy<sup>24,25</sup>.

109

110 Here, we applied two-sample Mendelian randomization (MR) using summary-level genetic  
111 data from the largest available GWAS for each sexual behaviour (sample 1) and  
112 oropharyngeal cancer (sample 2). We first conducted univariable MR analysis to assess the

113 effects of NSP and AFS on oropharyngeal cancer risk. We next performed univariable MR  
114 analysis to explore the effect of sexual behaviours on HPV seropositivity. Genetic proxies for  
115 complex human behaviours are more likely to have broad pleiotropic effects and may  
116 influence multiple upstream pathways that indirectly impact on sexual behaviour. In  
117 particular, genetic variants associated with sexual behaviour may also influence the disease  
118 outcome via other head and neck cancer risk factors, such as smoking and alcohol  
119 consumption. For this reason we performed a number of sensitivity analyses: i) MR methods  
120 to account for horizontal pleiotropy, ii) MR of sexual behaviours on positive (cervical cancer  
121 and seropositivity for *Chlamydia trachomatis*) and negative control outcomes (lung and oral  
122 cancer), iii) Causal Analysis Using Summary Effect estimates (CAUSE), to account for  
123 correlated and uncorrelated horizontal pleiotropic effects<sup>29</sup>, iv) multivariable MR analysis to  
124 account for the effects of smoking, alcohol and risk tolerance.

125

126 Despite observing an association between genetically predicted AFS and NSP and risk of  
127 oropharyngeal cancer using univariable MR, further multivariable analysis indicated  
128 violation of the core MR assumptions, likely due to correlated pleiotropy. This highlights the  
129 importance of performing comprehensive multivariable MR and sensitivity analyses, in  
130 addition to testing for correlated pleiotropy when using genetic instruments to proxy  
131 complex human behaviours.

132

## 133 **Results**

134

135 *Univariable Mendelian Randomization*

136 Using 139 SNPs robustly and independently associated with AFS (**Supplementary Data 1**),  
 137 there was evidence of a protective effect of later AFS on OPC (IVW OR 0.4, 95%CI 0.3, 0.7,  
 138 per standard deviation (SD),  $p = <0.001$ ) which was consistent across methods robust to  
 139 horizontal pleiotropy (MR-Egger, weighted median, and weighted mode) (**Table 1** &  
 140 **Supplementary Fig.1**). Using 117 SNPs (**Supplementary Data 1**) independently associated  
 141 with NSP, we found evidence to suggest an adverse effect of increased NSP on the risk of  
 142 OPC (IVW OR 2.2, 95%CI 1.3, 3.8 per SD,  $p < 0.001$ ). These results were consistent across the  
 143 other MR methods (**Table 1** & **Supplementary Fig.1**). The protective effect of later AFS was  
 144 consistent across all geographical regions, with the most precise effects seen in the  
 145 European (IVW OR 0.4, 95%CI 0.2, 0.8,  $p = <0.001$ ) and North American population (IVW OR  
 146 0.4, 95%CI 0.2, 0.8,  $p = 0.01$ ) (**Table 2**). There was also suggestive evidence for an adverse  
 147 effect of increasing NSP across regions, with the strongest effect again in the North  
 148 American population (IVW OR 3.0, 95%CI 1.4, 6.5,  $p = 0.01$ ) (**Table 3**).

149

150 **Table 1.** Univariable Mendelian randomization results for age at first sex and number of  
 151 sexual partners on risk of oropharyngeal cancer.

Outcome	Exposure/ Outcome datasets	Outcome N	Controls N	Method	Age at first sex (N SNPs 139)		Number of sexual partners (N SNPs 117)	
					OR (95%CI)	P	OR (95%CI)	P
OPC	UK Biobank/ GAME-ON	2,641	6,585	IVW	0.44 (0.28, 0.70)	<0.001	2.20 (1.27, 3.81)	<0.001
OPC	UK Biobank/ GAME-ON	2,641	6,585	Weighted median	0.41 (0.23, 0.75)	<0.001	2.57 (1.24, 5.29)	0.01
OPC	UK Biobank/ GAME-ON	2,641	6,585	Weighted mode	0.23 (0.04, 1.34)	0.10	3.57 (0.58, 21.69)	0.17
OPC	UK Biobank/ GAME-ON	2,641	6,585	MR-Egger	0.21 (0.03, 1.37)	0.10	1.88 (0.12, 29.49)	0.65

152

153

154 Abbreviations: OPC, oropharyngeal cancer; IVW, inverse variance weighted; OR, odds ratio; CI, confidence  
 155 intervals; P, *p*-value; NSP, number of sexual partners; AFS, age at first sex. AFS OR represents the exponential  
 156 change in odds of oropharyngeal squamous cell carcinoma per SD change (7.3-month delay) in age at first sex  
 157 NSP OR represents the exponential change in odds of oropharyngeal squamous cell carcinoma per SD increase  
 158 (0.94) in number of sexual partners.

159

160 **Table 2.** Inverse variance weighted univariable Mendelian randomization results for age at  
 161 first sex on risk of oropharyngeal cancer, by region.

162

Outcome	Region	N SNPs	Outcome N	Control N	Method	OR	CIL	CIU	P value
Oropharyngeal cancer	Europe	139	1,090	2,928	Inverse variance weighted	0.36	0.17	0.78	<0.001
Oropharyngeal cancer	North America	139	1,119	2,329	Inverse variance weighted	0.41	0.20	0.83	0.01
Oropharyngeal cancer	South America	139	205	727	Inverse variance weighted	0.38	0.07	1.95	0.24

163

164 Abbreviations: SE, standard error; OR, odds ratio; CIL, lower confidence interval; CIU, upper confidence  
 165 interval; P, *p*-value. OR represents the exponential change in odds of oropharyngeal squamous cell carcinoma  
 166 per SD change (7.3-month delay) in age at first sex.

167

168 **Table 3.** Inverse variance weighted univariable Mendelian randomization results for number  
 169 of sexual partners on risk of oropharyngeal cancer, by region.

170

Outcome	Region	N SNPs	Outcome N	Control N	Method	OR	CIL	CIU	P value
Oropharyngeal cancer	Europe	117	1,090	2,928	Inverse variance weighted	1.48	0.66	3.33	0.35
Oropharyngeal cancer	North America	117	1,119	2,329	Inverse variance weighted	2.99	1.37	6.51	0.01
Oropharyngeal cancer	South America	117	205	727	Inverse variance weighted	2.68	0.56	12.75	0.22

171



172 Abbreviations: SE, standard error; OR, odds ratio; CIL, lower confidence interval; CIU, upper confidence  
 173 interval; P, *p*-value. OR represents the exponential change in odds of oropharyngeal squamous cell carcinoma  
 174 per SD increase (0.94) in number of sexual partners.

175

176 *MR for effect of sexual behaviours on HPV seropositivity*

177 Using the NSP and AFS instruments, we next evaluated the effect of sexual behaviour on the  
 178 risk of HPV seropositivity in healthy individuals, using a GWAS of serological measures in UK  
 179 Biobank. There appeared to be some evidence for a protective effect of later AFS (IVW OR  
 180 0.5, 95%CI 0.2, 1.0, *p*=0.05) on HPV16 L1 seropositivity (**Table 4 & Supplementary Table 1**).  
 181 However, there was limited evidence for a similar protective effect on HPV18 L1, HPV16 E6  
 182 or E7 seropositivity. While there was some evidence that increasing NSP also increased the  
 183 likelihood of HPV16 E6 seropositivity (IVW OR 5.4, 95%CI 1.0, 28.3, *p*=0.05), this was  
 184 inconsistent among the other tested HPV antibodies (**Table 4 & Supplementary Table 2**).

185

186 **Table 4.** Inverse variance weighted univariable Mendelian randomization results of age at  
 187 first sex with HPV seropositivity.

Exposure	Outcome (serostatus)	N seropositive	N seronegative	Method	OR	CIL	CIU	P value
<b>Age at first sex</b>	HPV16 L1	344	7580	Inverse variance weighted	0.47	0.21	1.01	0.05
	HPV16 E6	133	7791	Inverse variance weighted	1.35	0.38	4.74	0.64
	HPV16 E7	256	7668	Inverse variance weighted	1.51	0.62	3.66	0.36
	HPV18 L1	191	7733	Inverse variance weighted	0.76	0.28	2.10	0.60
<b>Number of sexual partners</b>	HPV16 L1	344	7580	Inverse variance weighted	2.07	0.74	5.83	0.17
	HPV16 E6	133	7791	Inverse variance weighted	5.39	1.03	28.30	0.05
	HPV16 E7	256	7668	Inverse variance weighted	0.89	0.28	28.00	0.84
	HPV18 L1	191	7733	Inverse variance weighted	0.47	0.13	1.73	0.25

188

189

190 Abbreviations: SE, standard error; OR, odds ratio; CIL, lower confidence interval; CIU, upper confidence  
191 interval; P, *p*-value. Age at first sex OR represents the exponential change in odds of HPV seropositivity per SD  
192 change (7.3-month delay) in age at first sex. Number of sexual partners OR represents the exponential change  
193 in odds of HPV marker seropositivity per SD increase (0.94) in number of sexual partners.

194

### 195 *Sensitivity analyses*

196 There was limited evidence of weak instrument bias (F-statistic >10) and the proportion of  
197 variance in the phenotype ( $R^2$ ) explained by the genetic instruments ranged from 1 - 2%  
198 (**Supplementary Table 3**). There was limited evidence for heterogeneity in the SNP effect  
199 estimates for the AFS instrument (QIVW 159.4,  $p= 0.10$ ; Q MR-Egger 158.6,  $p= 0.10$ ), but  
200 clear evidence of heterogeneity in the NSP instrument (QIVW 155.6,  $p= 0.007$ ; Q MR-Egger  
201 155.6,  $p= 0.006$ ) (**Supplementary Table 4**). MR-Egger intercepts were not indicative of  
202 directional pleiotropy (**Supplementary Table 5**) but there were outliers present on visual  
203 inspection in both scatter and leave-one-out plots (**Supplementary Figs.2 & 3**). MR-PRESSO  
204 identified 8 outliers for AFS and 7 outliers for NSP, which when corrected for, yielded effects  
205 consistent with univariable MR for both instruments (**Supplementary Tables 6-8**). There was  
206 evidence of violation of the NOME assumption for both AFS and NSP genetic instruments  
207 (i.e.,  $I^2$  statistic <0.90) (**Supplementary Table 9**), so MR-Egger was performed with SIMEX  
208 correction. The effects were consistent with previous MR-Egger results for AFS, but there  
209 was significant attenuation of the NSP effect on oropharyngeal cancer (SIMEX corrected MR-  
210 Egger OR 3.6, 0.4, 32.1,  $p= 0.25$ ) (**Supplementary Table 10**). These estimates should  
211 however be interpreted with caution, given evidence of high dilution in the SNP-exposure  
212 effects<sup>30</sup>.

213

214 *Positive and negative control analyses*

215 Univariable MR analysis conducted within UK Biobank found a protective effect for later AFS  
216 on cervical cancer, which is known to be another HPV-driven cancer type (IVW OR 0.4,  
217 95%CI 0.3, 0.7,  $p < 0.001$ ) (**Supplementary Table 11**). A similar effect was found when  
218 assessing the effect of AFS on *C. trachomatis* seropositivity based on pGP3 antigen, another  
219 positive control (IVW OR 0.4, 95%CI 0.3, 0.6,  $p < 0.001$ ) (**Supplementary Table 11**). There  
220 was also evidence for an adverse effect of increasing NSP on cervical cancer risk (IVW OR  
221 1.9, 95%CI 1.0, 3.9,  $p = 0.06$ ) and a positive association between NSP and *C. trachomatis*  
222 serostatus (IVW OR 2.4, 95%CI 1.4, 4.1,  $p < 0.001$ ) (**Supplementary Table 12**).

223

224 Using lung cancer as a negative control, in univariable MR there was a strong protective  
225 effect of AFS (IVW OR 0.1 95%CI, 0.1, 0.3  $p < 0.001$ ) (**Supplementary Table 11**) and an  
226 adverse effect of increasing NSP (IVW OR 7.1 95%CI, 2.4, 21.6  $p < 0.001$ ) (**Supplementary**  
227 **Table 12**), indicating violation of the MR assumptions. A protective effect was also observed  
228 in relation to AFS with oral cancer, another negative control (IVW OR 0.6, 95%CI 0.4, 1.0,  $p =$   
229 0.03) (**Supplementary Table 11**); however, there was no effect for NSP on oral cancer (IVW  
230 OR 1.2, 95%CI 0.7, 2.0,  $p = 0.47$ ) (**Supplementary Table 12**).

231

232 While there was no strong evidence for directional pleiotropy (**Supplementary Table 13**),  
233 there was some evidence of heterogeneity (**Supplementary Table 14**) for both AFS and NSP  
234 in the lung and oral cancer analyses, suggesting that pleiotropy may be present<sup>31</sup>. While  
235 scatter and leave-one-out plots showed no obvious outliers (**Supplementary Figs.4-7**), MR-  
236 PRESSO identified outliers for AFS and for NSP across all positive and negative controls.

237 When corrected for outliers, the lung cancer results remained consistent with the  
238 univariable MR, suggesting further violation of the MR assumptions for the AFS and NSP  
239 instruments (**Supplementary Tables 15-17**).

240

#### 241 *Investigating correlated pleiotropy using CAUSE*

242 We used GWAS summary statistics to evaluate evidence for an effect of AFS and NSP on  
243 oropharyngeal cancer, using the Causal Analysis using Summary Effect estimates (CAUSE)  
244 method to account for correlated pleiotropy<sup>32</sup>. For AFS, CAUSE suggested there was  
245 relatively similar evidence for sharing (correlated pleiotropy) ( $p= 0.02$ ) and causal models  
246 ( $p= 0.05$ ) compared to the null (no effect) model (**Supplementary Table 18 &**  
247 **Supplementary Fig.8**). Comparing both shared and causal models, there was limited  
248 evidence that the causal model fit the data better than the sharing model ( $p= 0.44$ ),  
249 indicating that correlated pleiotropy could not be discounted. When investigating the causal  
250 effect of NSP on oropharyngeal cancer, neither shared ( $p= 0.30$ ) nor causal ( $p= 0.27$ ) models  
251 appeared to fit in comparison to the null model, providing limited evidence for a causal  
252 effect of NSP (**Supplementary Table 19 & Supplementary Fig.9**).

253

#### 254 *Multivariable MR*

255 In total there were 21 overlapping SNPs identified between genetic instruments  
256 (**Supplementary Table 20**) and LD score regression highlighted strong genetic correlation  
257 between the traits ( $r_g = |0.20-0.63|$ ) **Supplementary Table 21 & Supplementary Fig.10**).  
258 Multivariable MR analysis was therefore carried out to investigate the direct causal effect of  
259 AFS and NSP on oropharyngeal cancer after accounting for the other sexual behaviour,  
260 smoking, alcohol, and risk tolerance. While the effect of NSP diminished (IVW OR 0.8, 95%CI

261 0.3, 2.0,  $p=0.60$ ), the AFS effect remained (IVW OR 0.4, 95%CI 0.2, 0.9,  $p=0.04$ ), after  
262 accounting for the other sexual behaviour in multivariable MR (**Tables 5 & 6; Fig.1**). When  
263 accounting for smoking and risk tolerance, the effect of AFS remained consistent within the  
264 oropharyngeal subsite (**Table 5 & Fig.2**). However, there was attenuation of the effect for  
265 AFS when controlling for drinks per week (IVW OR 0.7, 95%CI 0.4, 1.2,  $p=0.21$ ). The effect of  
266 NSP on oropharyngeal cancer attenuated when accounting for lifetime smoking (IVW OR  
267 0.9, 95%CI 0.5 1.72,  $p=0.76$ ), alcohol consumption (IVW OR 1.5, 95%CI 0.8, 2.8,  $p=0.27$ ) and  
268 risk tolerance (IVW OR 2.0, 95%CI 0.9, 4.4,  $p=0.07$ ) (**Table 6 & Fig.3**) These results suggest  
269 the NSP and AFS instruments may include pleiotropic variants related to smoking and  
270 drinking behaviours. Some of the multivariable models including smoking initiation and  
271 drinks per week showed high levels of heterogeneity and therefore further risk of invalid  
272 instruments (**Tables 5 & 6**). However, the MR-Egger intercepts in the multivariable analyses  
273 were consistent with the null, indicative of no further directional pleiotropy (**Supplementary**  
274 **Table 22**) and the effects estimated were also consistent across both IVW and MR-Egger  
275 models (**Tables 5 & 6**).

276

277

278

279

280

281

282

283

284

285

286 **Fig.1** Forest plot comparing univariable and multivariable Mendelian randomization effects

287 of age at first sex and number of sexual partners on oropharyngeal cancer risk.

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

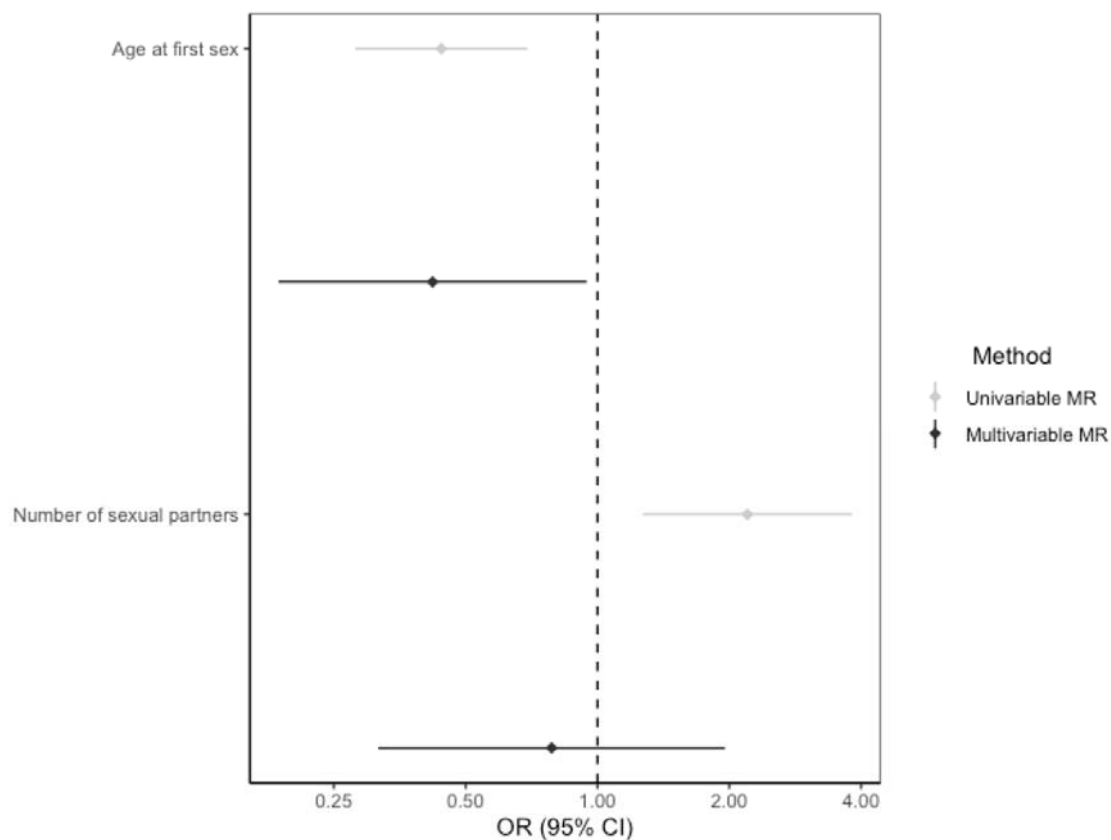
304

305

306

307

308



Effect estimates are reported on the log odds scale with 95% confidence intervals. Age at first sex point estimate represents the exponential change in odds of oropharyngeal squamous cell carcinoma per SD change (7.3-month delay) in age at first sex. Number of sexual partners point estimate represents the exponential change in odds of oropharyngeal squamous cell carcinoma per SD increase (0.94) in number of sexual partners.

309 **Table 5.** Multivariable Mendelian randomization for age at first sex with risk of  
 310 oropharyngeal cancer.

311

Exposure	Exposure dataset	N SNPs	F-stat	Q-stat	P-value for instrument validity	Method	AFS OR	95% CI	P
<b>Number of Sexual partners</b>	UK Biobank <sup>33</sup>	152	7.81	214.14	3.77x10 <sup>-4</sup>	IVW	0.42	0.19, 0.94	0.04
						MR Egger	0.24	0.06, 1.01	0.05
<b>Comprehensive Smoking Index</b>	UK Biobank <sup>34</sup>	174	8.87	191.26	0.14	IVW	0.48	0.25, 0.95	0.03
						MR-Egger	0.71	0.23, 2.19	0.56
<b>Smoking initiation</b>	GSCAN <sup>35</sup>	215	6.43	250.58	0.04	IVW	0.42	0.21, 0.83	0.01
						MR-Egger	0.61	0.21, 1.74	0.35
<b>Drinks per week</b>	GSCAN <sup>35</sup>	147	28.88	164.77	0.11	IVW	0.72	0.43, 1.20	0.21
						MR-Egger	0.43	0.14, 1.28	1.28
<b>Risk tolerance</b>	UK Biobank <sup>26</sup>	160	13.68	171.18	0.21	IVW	0.53	0.30, 0.93	0.03
						MR-Egger	0.24	0.08, 0.72	0.01

312

313 Abbreviations: IVW, inverse variance weighted; AFS, age at first sex; OR, odds ratio; CI, confidence intervals; P,

314 *p-value*; Q-stat, Cochran's Q statistic; F-stat, conditional F-statistic. AFS OR represents the odds ratio of

315 oropharyngeal squamous cell carcinoma per SD change (7.3-month delay) in age at first sex.

316

317

318

319

320 **Table 6.** Multivariable Mendelian randomization for number of sexual partners with risk of  
 321 oropharyngeal cancer.

322

Exposure	Exposure dataset	N SNPs	F-stat	Q-stat	P-value for instrument validity	Method	NSP OR	95% CI	P
Age at first sex	UK Biobank <sup>33</sup>	152	7.26	214.14	3.77x10 <sup>-4</sup>	IVW	0.79	0.32, 1.96	0.60
						MR Egger	0.97	0.35, 2.65	0.95
Comprehensive Smoking Index	UK Biobank <sup>34</sup>	157	12.18	168.74	0.20	IVW	0.91	0.48, 1.72	0.76
						MR-Egger	0.86	0.39, 1.88	0.70
Smoking initiation	GSCAN <sup>35</sup>	195	7.35	204.67	0.25	IVW	1.51	0.77, 2.97	0.23
						MR-Egger	1.66	0.66, 4.15	0.28
Drinks per week	GSCAN <sup>35</sup>	117	22.7	151.65	0.011	IVW	1.45	0.75, 2.79	0.27
						MR-Egger	1.61	0.76, 3.41	0.21
Risk tolerance	UK Biobank <sup>26</sup>	125	7.06	145.56	0.072	IVW	2.04	0.93, 4.44	0.07
						MR-Egger	2.12	0.66, 6.83	0.21

323

324 Abbreviations: IVW, inverse variance weighted; NSP, number of sexual partners; OR, odds ratio; CI, confidence  
 325 intervals; P, p-value; Q-stat, Cochran's Q statistic; F-stat, conditional F-statistic. NSP OR represents the odds  
 326 ratio of oropharyngeal squamous cell carcinoma per SD increase (0.94) in number of sexual partners.

327

328

329

330

331

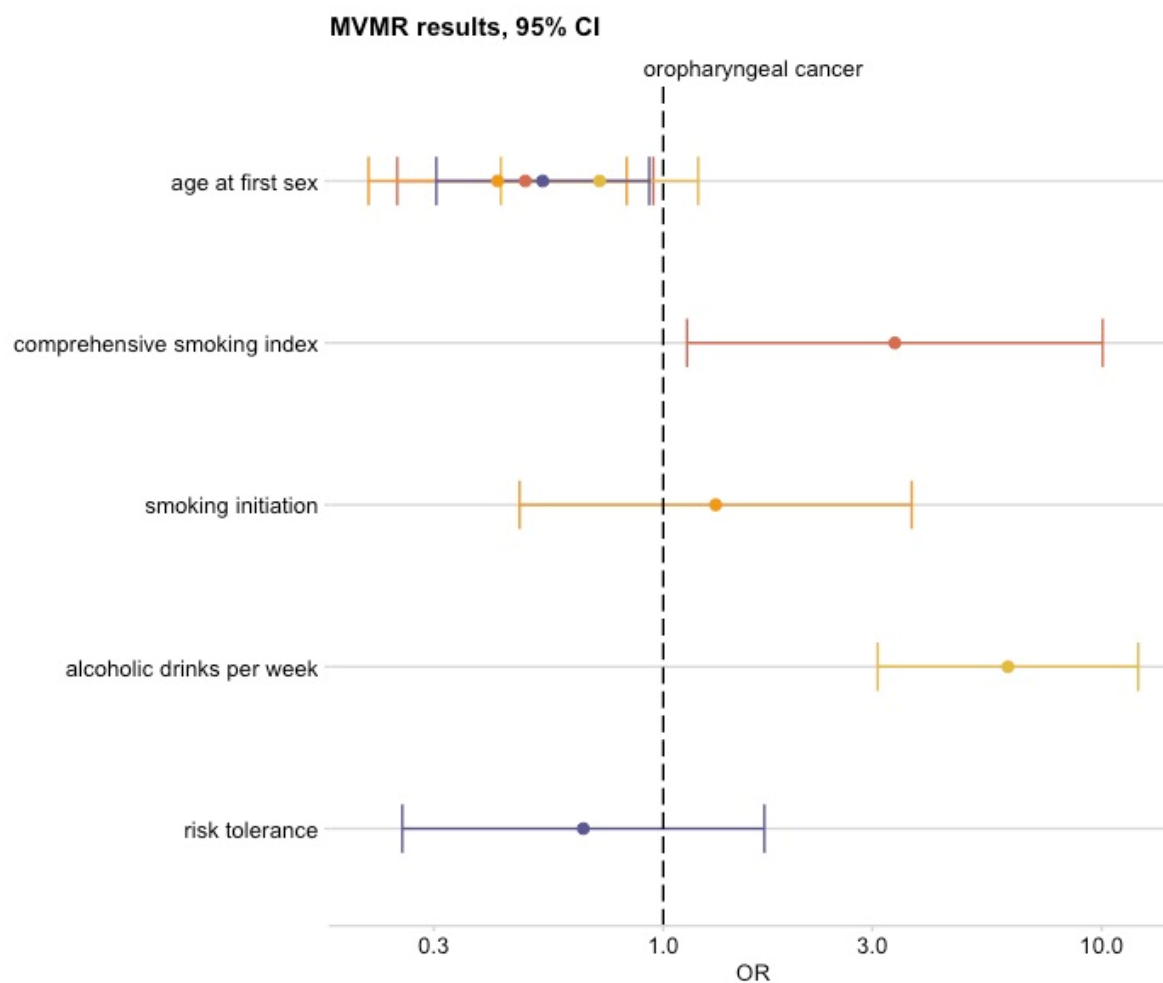
332



333 **Fig.2** Forest plot showing multivariable Mendelian randomization results for age at first sex  
334 single nucleotide polymorphisms with risk of oropharyngeal cancer.

335

336



368

369

370

371

372 Effect estimates on oropharyngeal cancer risk are reported on the log odds scale with 95% confidence

373 intervals. Age at first sex OR represents the change in odds of oropharyngeal squamous cell carcinoma per SD

374 change (7.3-month delay) in age at first sex. Comprehensive smoking index (dark orange), smoking initiation

375 (light orange), alcoholic drinks per week (yellow), risk tolerance (blue).

376

377

378

379 **Fig.3** Forest plot showing multivariable Mendelian randomization results for number of  
380 sexual partners single nucleotide polymorphisms with risk of oropharyngeal cancer.

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

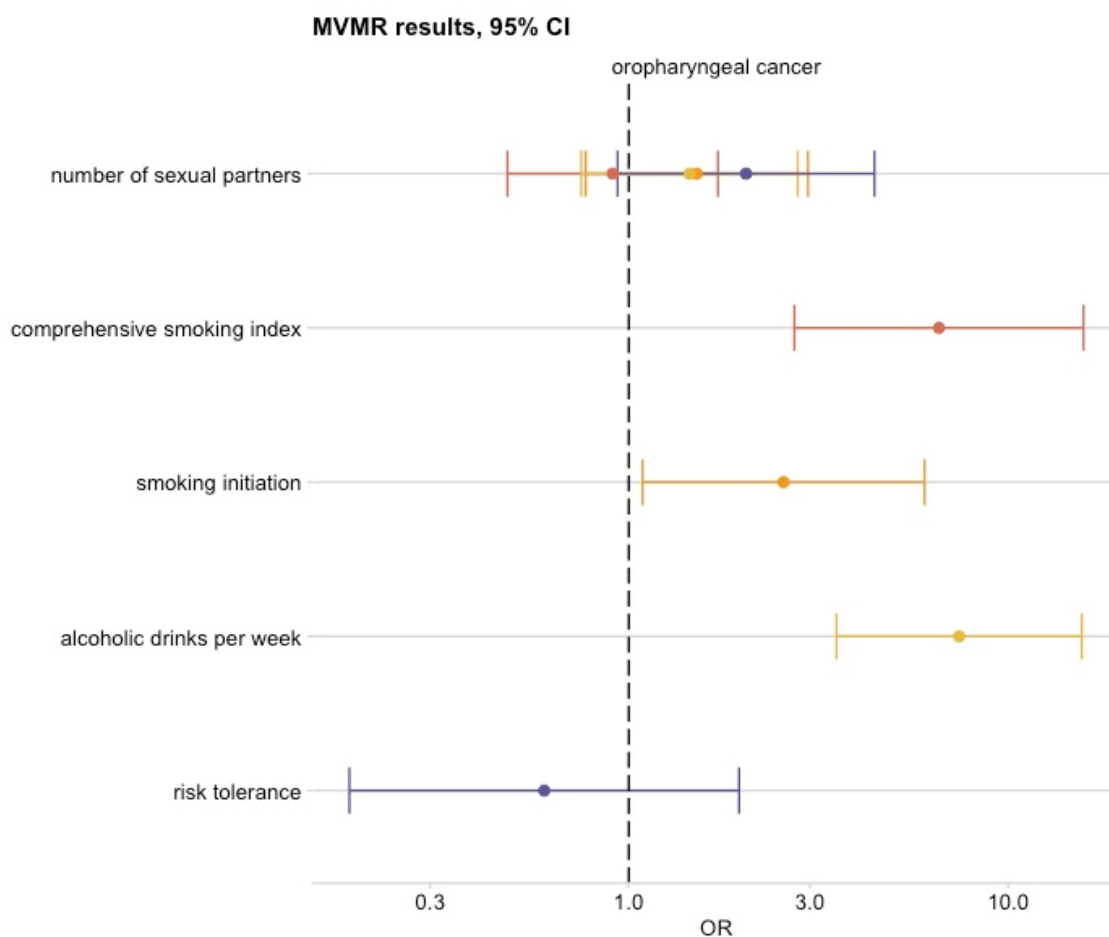
400

401

402

403

404



410

411

412 Effect estimates on oropharyngeal cancer risk are reported on the log odds scale with 95% confidence

413 intervals. Number of sexual partners OR represents the change in odds of oropharyngeal squamous cell

414 carcinoma per SD change (0.94) in number of sexual partners. Comprehensive smoking index (dark orange),

415 smoking initiation (light orange), alcoholic drinks per week (yellow), risk tolerance (blue).

416

417

418

419

420 In additional multivariable MR analysis of AFS and NSP on lung cancer, effects for both  
421 instruments were attenuated once smoking was included in the model. With AFS, this was  
422 clearly seen when controlling for smoking initiation (IVW OR 1.1, 95%CI 0.8, 1.6,  $p=0.57$ )  
423 and a change in direction of the effect of AFS was evident when controlling for the  
424 comprehensive smoking index (IVW OR 2.0, 95%CI 1.3, 3.0,  $p<0.001$ ) (**Supplementary Table**  
425 **23 & Supplementary Fig.11**). Similarly, there was limited evidence for an effect of NSP on  
426 lung cancer when controlling for the comprehensive smoking index (IVW OR 0.7, 95%CI 0.4,  
427 1.1,  $p=0.09$ ). The MR-Egger intercept deviated from the null in the multivariable models  
428 including smoking, suggestive of further directional pleiotropy in this analysis  
429 (**Supplementary Table 24**).

430

## 431 Discussion

432

433 In this study we applied Mendelian randomization to evaluate the effects of both later age  
434 at first sex and increased number of sexual partners on the risk of oropharyngeal cancer. We  
435 observed convergence between genetic pathways influencing sexual behaviors and  
436 susceptibility to oropharyngeal cancer, which may be partly mediated by HPV infection,  
437 however, we also uncovered complex correlated pleiotropy with other putative risk factors.  
438 Univariable MR results suggested a protective effect of later age at first sex and an adverse  
439 effect of increased number of sexual partners. However, these effects attenuated in the  
440 multivariable MR analyses that controlled for smoking behaviour and alcohol consumption.  
441 While there was suggestive evidence for an effect of sexual behaviours on some HPV16  
442 serology measures and in cervical cancer (supportive of a causal mechanism via HPV  
443 infection), the same direction of effect was observed in negative control analysis (lung and

444 oral cancer) indicating potential violation of the MR assumptions. Furthermore, CAUSE  
445 provided less support for a causal effect of AFS and NSP on oropharyngeal cancer risk,  
446 highlighting the risk of correlated pleiotropy in the genetic instruments for these complex  
447 behavioural traits.

448

#### 449 *Sexual behaviours and HPV transmission*

450 Over 90% of HPV positive OPC is caused by the high-risk genotype 16, with almost all oral  
451 infections thought to be sexually acquired<sup>36</sup>. HPV is a small non-enveloped DNA virus, with  
452 its genome encoding for both early oncoproteins E6/E7 and the late capsid proteins such as  
453 L1. The overexpression of these oncogenes is thought to stimulate proliferation and lateral  
454 expansion of epithelial basal cells, progressing to a malignant phenotype. HPV E6 forms a  
455 complex which leads to rapid degradation of tumour suppressor protein p53, resulting in  
456 deregulation of cell cycle checkpoints. E7 binds to a complex which ubiquitinates another  
457 tumour suppressor protein, retinoblastoma (pRb), again resulting in uncontrolled G1/S  
458 phase of the cell cycle<sup>37</sup>. While the transmission of HPV via sexual intercourse is well known  
459 and HPV in turn a major risk factor for cervical malignancies, the role of HPV in  
460 oropharyngeal cancer risk has only been acknowledged in recent decades<sup>8</sup>. Among OPC  
461 cases, HPV16 E6 serology is a good biomarker (~99% specificity, >90% sensitivity) and  
462 therefore both E6 and E7 are highly associated with this disease<sup>38</sup>. However, when studying  
463 these antibodies in the general population, E6 seroprevalence appears to be very low (0.5-  
464 1%), but in comparison with low incidence rates of HPV-positive OPC this figure is still high,  
465 suggesting that not all individuals in the general population who have HPV16 E6  
466 seropositivity will develop an oropharyngeal tumour or other HPV-associated cancer<sup>38</sup>.  
467 Consequently, we performed this analysis in UK Biobank and observed a strong and

468 consistent association with sexual behavior. In our univariable MR analysis, the effects of  
469 AFS and NSP instruments on risk of HPV16 and HPV18 seropositivity were not consistent,  
470 compared with recent observational studies which demonstrate an association between  
471 serology markers and sexual behaviour responses in UK Biobank<sup>38</sup>. This could be as a result  
472 of the small number of seropositive HPV16 (n= <450) and HPV18 (n= 265) cases within the  
473 UK Biobank pilot study used in our genetic analysis or that results from genetic proxies and  
474 questionnaire data are not directly comparable<sup>39</sup>. Using serology measures to predict HPV  
475 seropositivity or a HPV-positive OPC diagnosis is not straightforward, often requiring the use  
476 of multiple markers simultaneously<sup>40</sup>. Going forward, more reliable tests may emerge which  
477 could improve our prediction of both the infection and disease.

478

#### 479 *Regional differences in sexual behaviour and HPV prevalence*

480 Although the incidence of OPC in South America is similar to that in Western Europe and  
481 North America, the prevalence of HPV16 is reportedly low<sup>41</sup>. Latin America has an  
482 estimated overall HPV-positive head and neck cancer prevalence of between 3 – 4%,  
483 compared with 25% in European and North American populations<sup>41-43</sup>. This could partly be  
484 explained by differences in data collection and methods used to detect HPV. Despite Latin  
485 American countries having an average age of sexual debut between 18-19 years old<sup>44</sup>, the  
486 International Head and Neck Cancer Epidemiology (INHANCE) Consortium found that these  
487 countries reported higher mean numbers of sexual partners (e.g., Brazil n=22), compared  
488 with North American (e.g., USA, Atlanta n=10) or European (e.g., Warsaw n= 15) populations  
489<sup>15</sup>. Stratifying by region in our univariable MR analysis, we found a consistent protective  
490 effect for AFS and similarly, a consistent increased risk effect for NSP across all three regions  
491 (Europe, North America, and South America), with evidence for the most precise effects in

492 the North American population. In the largest pooled analysis, authors also report possible  
493 recall or reporting biases, given that some of the sexual behaviour interviews were carried  
494 out with family members nearby, in addition to small sample sizes (<150 cases)<sup>15</sup> which may  
495 have affected their results.

496

#### 497 *Confounding by other risk factors*

498 While transmission of HPV to the upper aerodigestive tract is thought to be through oral  
499 sexual contact<sup>9,15-20</sup>, a more recent meta-analysis reported no association between oral sex  
500 practices and head and neck cancer risk<sup>22</sup>. This could be explained by the inclusion of older  
501 studies<sup>22</sup>, which may not have captured the more recent rise in number of HPV-positive  
502 OPC cases which has been described by some as an ‘epidemic’ and predicted to overtake  
503 oral cancer within the next decade<sup>45</sup>. However, a study in the UK found that there was no  
504 change in the proportion of HPV-attributable cases from 2002-2011, although the incidence  
505 of OPC doubled over the same time period and national surveys have not described an  
506 increase in oral sex behaviour<sup>1,46</sup>. In one multi-national study of 1,626 men aged 18–73  
507 years with 4-year follow-up, no association was detected between oral sexual behaviours  
508 and incident HPV infection, but oral oncogenic HPV was found to be more prevalent in  
509 current smokers compared with non-smokers<sup>47</sup>. Furthermore, tobacco exposure induces  
510 proinflammatory and immunosuppressive effects, which could potentially increase the  
511 likelihood of HPV infection and persistence<sup>48,49</sup>. Since risk factors such as smoking and  
512 alcohol consumption are strongly associated with sexual behaviour and are well established  
513 in the aetiology of HNSCC, this may confound the relationship between sexual behaviours  
514 with HPV transmission and similarly oropharyngeal cancer in observational studies<sup>50,51</sup>.

515

516 Although Mendelian randomization analysis minimises the likelihood of confounding, since  
517 germline genetic variants should not theoretically be influenced by subsequent  
518 environmental confounders, pleiotropy is a major concern whereby genetic variants  
519 associated with the exposure (sexual behaviours) are related to the outcome  
520 (oropharyngeal cancer) through alternative, independent biological pathways. We used a  
521 series of analyses to evaluate the potential for pleiotropy. We first performed several  
522 methods (MR-Egger<sup>52</sup>, weighted median<sup>53</sup> and weighted mode<sup>54</sup>) which allow for the  
523 existence of horizontal pleiotropy and correct for this. We also identified and corrected for  
524 outlier SNPs most likely to exhibit pleiotropic effects. In univariable MR analyses, estimates  
525 were consistent with an effect of AFS and NSP on oropharyngeal cancer risk. However, in  
526 further MR analysis taking lung cancer as a negative control, we observed the same  
527 direction of effect for AFS and NSP which we did not expect, since there is no plausible  
528 biological mechanism directly linking sexual behaviour with lung cancer risk. Evidence of an  
529 effect here indicates potential violation of the MR assumptions.

530

531 Strong genetic correlation between sexual behaviours and other risk factors such as  
532 smoking, alcohol and general risk tolerance were found using LD score regression. The  
533 genetic instruments used in MR may therefore comprise variants which primarily influence  
534 other risk factors, which could induce correlated pleiotropy (**Fig.4**). We conducted two  
535 subsequent analyses to evaluate this. The CAUSE approach provided limited evidence for  
536 any effect of NSP on oropharyngeal cancer and was unable to distinguish an effect of AFS  
537 from the situation of correlated pleiotropy. We also performed multivariable MR to control  
538 for alcohol, smoking and risk tolerance, so as to determine the direct causal effect of sexual  
539 behaviours on oropharyngeal cancer. Effect estimates attenuated when alcohol and

540 smoking were taken into account in the multivariable MR models, again highlighting the role  
541 of potential pleiotropy in the genetic instruments for sexual behaviour.

542

#### 543 *Strengths and limitations*

544 MR was employed in this study in an attempt to overcome the drawbacks of conventional  
545 epidemiological studies. However, MR makes various assumptions which if violated may  
546 generate spurious conclusions. For example, sexual behaviours are difficult to instrument  
547 genetically due to measurement error (e.g., as a result of reporting bias) and because they  
548 are time-varying as well as context and culture-dependent. This could hamper the detection  
549 of genetic associations related to these traits which has implications for genetic instrument  
550 strength (the first assumption of MR), given the low percentage of variation explained ( $R^2$ ),  
551 as well as potential violation of the no measurement error (NOME) assumption, with  
552 relatively low  $I^2$  values.

553

554 Additionally, the available genetic instruments are not specifically for oral sex, which is the  
555 conceptually relevant exposure and mode of HPV transmission. However, other sexual  
556 behaviours are likely to be correlated and developing genetic instruments for specific sexual  
557 activities poses some methodological and ethical challenges. While the random inheritance  
558 of genetic variants from parents to offspring means genotypes are typically much less  
559 associated with many potential confounders than directly measured exposures (the second  
560 MR assumption), an obvious violation of this is created due to population stratification  
561 which can introduce confounding of genotype-outcome associations. Although the GWAS  
562 for both NSP and AFS were adjusted for genetic principal components, given that sexual  
563 behaviours are strongly socially patterned, residual population structure may reintroduce



564 confounding into MR analysis. Although a rare outcome, there is potential sample overlap  
565 present as head and neck cancer cases were not excluded from previously published AFS or  
566 NSP GWAS, however recent studies suggest the incurred bias is much less substantial than  
567 that due to weak instruments, or overestimation of the SNP-trait effect <sup>55,56</sup>.

568

569 The third major assumption of MR is the exclusion restriction principle (i.e., that the genetic  
570 variant affects the outcome exclusively through its effect on the exposure). We performed a  
571 series of comprehensive sensitivity analyses to evaluate potential violation of this  
572 assumption. While several pleiotropy-robust (MR-Egger, weighted median, and weighted  
573 mode) and outlier exclusion methods provided limited evidence for violation of this  
574 assumption, the results of the lung cancer negative control analysis, CAUSE method and  
575 multivariable MR all suggested violation of the exclusion restriction assumption in the  
576 univariable MR of sexual behaviours on oropharyngeal cancer risk. When multiple sources  
577 of evidence provide conflicting estimates, it is necessary to appraise the relative biases of  
578 the approaches in order to best “triangulate” evidence <sup>57,58</sup>. In this instance, it is possible  
579 that the primary phenotype for the genetic variants used to instrument the sexual  
580 behaviours has been mis-specified. For example, the genetic variants may be primarily  
581 associated with other traits (e.g., risk taking) and indirectly to sexual behaviours via the  
582 primary traits. Similarly, sexual behaviour instruments may be associated with traits which  
583 don’t have a direct negative connotation. In this instance, the Instrument Strength  
584 independent of Direct Effect (InSIDE) assumption of approaches such as MR-Egger is likely to  
585 be violated, whereas the CAUSE is less vulnerable to environmental confounders that are  
586 correlated with genetic variants than the other pleiotropy-robust methods. Multivariable  
587 MR was also used to directly model the potential indirect effects of the genetic variants via

588 other traits (smoking, alcohol, and risk tolerance) and supported the conclusions of the  
589 CAUSE method. Finally, we could not distinguish between HPV positive and negative  
590 oropharyngeal tumours in the GAME-ON summary data, which would require further  
591 analysis at an independent level or a GWAS of oropharyngeal cancer stratified by HPV  
592 status. The GWAS-by-subtraction approach<sup>59</sup> could be useful to account for latent factors of  
593 other behavioral traits to identify more specific genetic instruments for sexual behaviour, if  
594 valid instruments for these traits exist. More serological data may become available in the  
595 UK Biobank and other clinical genetic studies, which could enhance power to evaluate  
596 potential the extent to which any effect of sexual behaviour on cancer risk is mediated by  
597 HPV.

598

### 599 *Conclusions*

600 In conclusion, this study used a comprehensive series of MR analyses to investigate sexual  
601 behaviours in relation to oropharyngeal cancer. We initially observed an association  
602 between genetically predicted AFS and NSP and risk of oropharyngeal cancer using  
603 univariable MR. Despite using genetic variants strongly related to these traits in large-scale  
604 GWAS, further multivariate methods indicated violation of the core MR assumptions, likely  
605 due to correlated pleiotropy. Effect estimates attenuated when alcohol and smoking were  
606 taken into account in the multivariable MR models, highlighting the importance of  
607 performing these further analyses, particularly when using genetic instruments which proxy  
608 complex behavioural traits.

609

610

611

## 612 **Methods**

613

### 614 *Summary-level data for sexual behaviours*

615 Summary statistics for AFS were obtained from a GWAS conducted in the UK Biobank  
616 (n=397,338)<sup>28</sup>. AFS was treated as a continuous variable, with individuals considered as  
617 eligible if they had given a valid answer to the question “*What was your age when you first*  
618 *had sexual intercourse? (Sexual intercourse includes vaginal, oral or anal intercourse)*” and  
619 ages <12 years old were excluded. Since AFS had a non-normal distribution, a within-sex  
620 inverse rank normal transformation was applied<sup>28</sup>. Where possible the full 272 SNP AFS  
621 instrument was used, except in the primary analysis of OPC, whereby only 139 SNPs could  
622 be extracted from head and neck cancer data. We obtained summary statistics for the NSP  
623 instrument (117 SNPs) from a GWAS conducted in UK Biobank<sup>26</sup> (n=2370,711). NSP was  
624 treated as a continuous variable based on responses to the question: “*About how many*  
625 *sexual partners have you had in your lifetime?*”. Respondents who reported >99 lifetime  
626 sexual partners were asked to confirm their responses and a value of zero was assigned to  
627 participants who reported having never had sex, which was normalised separately for both  
628 males and females with an inverse rank normal transformation<sup>26</sup>. Both AFS and NSP GWAS  
629 adjusted for the top 10 principal components (accounting for population stratification), sex,  
630 and birth year. For AFS, those participants with family data were controlled with non-  
631 independence of family members or else one family member was included in the analysis<sup>28</sup>.

632

### 633 *Summary-level data for oropharyngeal cancer*

634 The largest available GWAS for oropharyngeal cancer was performed on 2,641 OPC cases  
635 and 6,585 matched controls from 12 studies which were part of the Genetic Associations

636 and Mechanisms in Oncology (GAME-ON) Network<sup>60</sup>. Cancer cases comprised the following  
637 ICD-10 codes: oropharynx (C01.9, C02.4 and C09.0–C10.9). Stratification was conducted by  
638 geographical region to evaluate potential heterogeneity in any effects given potential  
639 differences in the distribution of genetic variants for specific traits within populations. As  
640 GAME-ON included participants from Europe (45.3%), North America (43.9%) and South  
641 America (10.8%), this study was restricted to individuals of predominantly European  
642 ancestry to avoid the effect of population structure. Details of the studies included as well  
643 as the genotyping and imputation performed have been described previously<sup>60,61</sup>.

644

#### 645 *Univariable Mendelian randomization*

646 To assess effects of NSP and AFS, we used SNPs which reached genome-wide significance ( $p$   
647  $<5 \times 10^{-8}$ ) in the respective GWAS and for which pairwise  $r^2 < 0.1$  (with 250kb linkage  
648 disequilibrium (LD) windows), ensuring only independent SNPs were selected into the  
649 instrument. Two-sample MR analyses were conducted using the “TwoSampleMR” package  
650 (version 0.5.5) in R (version 4.0.2) to extract the SNPs instrumenting the risk factor from the  
651 oropharyngeal cancer GWAS. Harmonization of the direction of effects between exposure  
652 and outcome associations was performed and palindromic SNPs were aligned when minor  
653 allele frequencies (MAFs) were less than 0.3 or were otherwise excluded. SNP specific Wald  
654 estimates were calculated (SNP-outcome estimate divided by SNP-exposure estimate) and  
655 an inverse variance weighted (IVW) method applied to meta-analyse these in order to  
656 obtain an effect estimate of the risk factor on oropharyngeal cancer risk.

657

#### 658 *MR for sexual behaviours on HPV and C. trachomatis seropositivity*

659 Where there was evidence for an effect of sexual behaviour on oropharyngeal cancer risk,  
660 we also aimed to confirm the suspected aetiological link via HPV, by investigating the effects  
661 of NSP and AFS on a range of seropositivity measures against HPV16 L1, E6, E7 and HPV18  
662 L1 proteins. Here, seropositivity suggests previous HPV exposure, which can be a predictor  
663 of cancer. Generally, HPV16 L1 antibodies are considered cumulative exposure markers,  
664 while HPV16 E6 and E7 have been associated with HPV-driven cancers but not all those who  
665 test positive are expected to develop a HPV-driven cancer<sup>38</sup>. Summary-level genetic data  
666 for HPV16 and HPV18 serological measures were obtained from UK Biobank. We performed  
667 individual GWAS for each measure using a similar approach as described by Kachuri *et al.*<sup>62</sup>.  
668 Details on how these GWAS were conducted can be found in the **Supplementary**  
669 **information**.

670

#### 671 *Sensitivity analyses*

672 The strength of each genetic instrument was determined by the magnitude and precision of  
673 association with the sexual behaviour, which was considered to be sufficient if the  
674 corresponding F-statistic was >10. The fixed-effect IVW method provides an unbiased  
675 estimate in the absence of horizontal pleiotropy or when horizontal pleiotropy is balanced  
676<sup>31</sup>. To account for directional pleiotropy, we compared results with three other MR  
677 methods, which each make different assumptions about this: MR-Egger<sup>52</sup>, weighted median  
678<sup>53</sup> and weighted mode<sup>54</sup>. Further detail on these methods is provided in the **Supplementary**  
679 **information (see “Methods”)**.

680

#### 681 *Positive and negative control analyses*

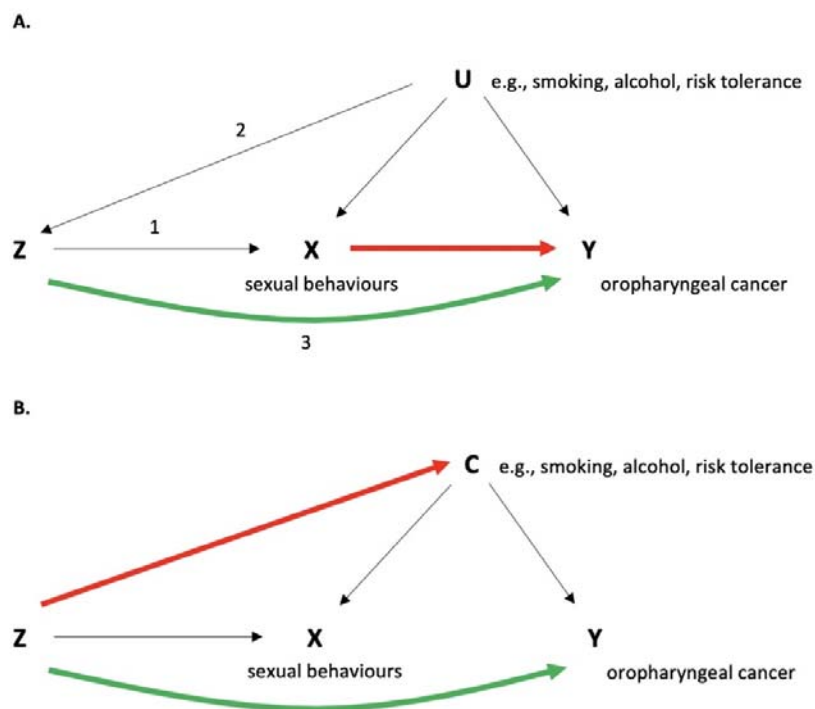
682 To further assess the specificity and sensitivity of the genetic instruments identified in  
683 relation to sexual behaviour, we conducted additional positive and negative control MR  
684 analyses. These aimed to appraise the role of AFS and NSP on a) cervical cancer and C.  
685 trachomatis seropositivity, as positive control outcomes where evidence of an effect would  
686 support the aetiological link via HPV; and b) lung cancer and oral cancer as negative  
687 controls, where a direct causal effect of sexual behaviour is unlikely and so where any  
688 evidence of an effect would indicate potential violation of the MR assumptions due to  
689 pleiotropy, population stratification or selection bias<sup>63</sup>. Details on the GWAS summary data  
690 used to conduct positive and negative control outcomes can be found in the **Supplementary**  
691 **information**.

692

#### 693 *Causal Analysis using Summary Effect estimates (CAUSE)*

694 While sensitivity analyses like MR-Egger, weighted median and weighted mode can detect  
695 horizontal or uncorrelated pleiotropy, whereby the genetic variant affects the exposure  
696 (sexual behaviours) and outcome (oropharyngeal cancer) through separate mechanisms,  
697 correlated pleiotropy is an alternative scenario which could generate spurious associations  
698 in MR. Here, the genetic variant affects the exposure and outcome via a shared heritable  
699 factor. Correlated pleiotropy may be present in the genetic instruments for AFS and NSP,  
700 which if undetected could lead to false positive results (**Fig.4**).

701



702

703

704 **Fig.4 A.** Genetic variants (Z) act as proxies or instruments to investigate if an exposure (X), is  
705 associated with a disease outcome (Y). Causal inference can be made between X and Y if the  
706 following conditions are upheld: (1) Z is a valid instrument, reliably associated with X  
707 ('relevance'); (2) Z is independent of any measured or unmeasured confounding factors (U)  
708 ('exchangeability') and (3) there is no independent association between Z and Y except  
709 through X ('exclusion restriction'). **B.** Directed acyclic graph (DAG) depicting correlated  
710 pleiotropy (C) whereby the genetic variant (Z) can affect the exposure (X) and the outcome  
711 (Y) via a shared heritable factor (C), for example here through smoking, alcohol, or general  
712 risk tolerance.

713

714 We used the CAUSE method in an attempt to identify potential correlated pleiotropy<sup>29</sup>.

715 CAUSE proposes that any causal effect of an exposure on the outcome leads to correlation

716 for all variants with a non-zero effect on the exposure, while a shared factor induces  
717 correlation for only a subset of exposure effect variants<sup>29</sup>. GWAS summary statistics were  
718 used to generate two models nested in a “null” effects model. The sharing model allows for  
719 horizontal pleiotropic effects but no causal effect ( $\gamma = 0$ ), whereas the causal model has  $\gamma$  as  
720 a free parameter. The Bayesian expected log pointwise posterior density (ELPD) is used to  
721 compare models, producing a one-sided  $p$  value which tests the best fitting model. In  
722 particular, if the hypothesis that the sharing model fits the data at least as well as the causal  
723 model is rejected, we can conclude that the data are consistent with a causal effect<sup>29</sup>.

724

#### 725 *LD Score regression*

726 Genetic correlation was calculated between the two sexual behaviour traits (AFS and NSP),  
727 smoking, alcohol, and general risk tolerance. Summary-level genome-wide association  
728 studies were obtained for alcohol consumption (drinks per week,  $n = 941,280$ ) and smoking  
729 initiation (a binary phenotype indicating whether an individual had ever smoked regularly)  
730 ( $n = 1,232,091$ ) from the GSCAN meta-analysis<sup>35</sup>. Summary statistics were also obtained  
731 from a GWAS of general risk tolerance ( $n = 939,908$ ), derived from a meta-analysis of UK  
732 Biobank ( $n = 431,126$ ) question “*Would you describe yourself as someone who takes risks?*”  
733 and 23andMe ( $n = 508,782$ ) question “*Overall, do you feel comfortable or uncomfortable*  
734 *taking risks?*”. The GWAS of risk tolerance was based on one’s tendency or willingness to  
735 take risks, making them more likely to engage in risk-taking behaviours more generally<sup>26</sup>.  
736 The regression was performed using pre-computed LD scores calculated based on  
737 individuals of European ancestry from 1000 Genomes European data and are appropriate  
738 for use with European-ancestry GWAS data<sup>64</sup>. This was filtered to HapMap3 SNPs as these  
739 are well-imputed in most studies<sup>65</sup>. SNPs found on chromosome 6 in the region 26MB to



740 34MB were excluded. GWAS summary statistics were converted for LD score regression  
741 using the munge\_sumstats.py command from the command line tool “ldsc”, and LD score  
742 regression was performed using the ldsc.py command.

743

#### 744 *Multivariable Mendelian randomization*

745 To account for the potential genetic overlap with other risk factors<sup>26</sup> for oropharyngeal  
746 cancer which may lead to correlated pleiotropy, we next conducted two-sample  
747 multivariable MR analysis. This accounted for the effects of the other sexual behaviour,  
748 smoking, alcohol consumption and risk tolerance in the MR of each sexual behaviour onto  
749 the cancer outcomes. First multivariable MR was carried out to assess the effect of genetic  
750 overlap between AFS and NSP using the genome-wide significant SNPs identified as  
751 instruments in the univariable analysis (272 SNPs for AFS and 117 SNPs for NSP). 196  
752 independent SNPs ( $p < 5 \times 10^{-8}$ ) were used in the analysis for smoking initiation, 60 SNPs for  
753 alcoholic drinks per week and 123 for risk tolerance, after excluding SNPs with a pairwise  $r^2$   
754  $> 0.001$ . To better capture lifetime smoking (duration, heaviness, and cessation), we used  
755 108 SNPs which make up the comprehensive smoking index, derived by Wootton et al in the  
756 UK Biobank ( $n = 462,690$ )<sup>34</sup>. SNP overlap was assessed between all instruments. We used  
757 generalised versions of Cochran’s Q statistical tests for both instrument strength and validity  
758<sup>66</sup>. Both the IVW and MR-Egger framework have been extended to estimate causal effects in  
759 multivariable MR analysis<sup>67,68</sup>, which was conducted using both the MVMR (version 0.2.0)  
760 and MendelianRandomization<sup>69</sup> (version 0.5.0) packages in R (version 4.0.2).

761

762 Causal Analysis using Summary Effect Estimates, LD Score Regression and Multivariable  
763 Mendelian randomization approaches all require full GWAS summary data for the proposed

764 risk factors of interested. Where full data were available for the GWAS of NSP <sup>26</sup>, these have  
765 yet to be published for the GWAS of AFS, which we used to select genetic instruments .  
766 Therefore, for these approaches, we used another GWAS for AFS, also conducted using UK  
767 Biobank data (n=406,457), for which full summary data are publicly available  
768 (<https://gwas.mrcieu.ac.uk/datasets/ukb-b-6591/>). This GWAS was conducted using the  
769 MRC IEU UK Biobank GWAS pipeline, more details of which can be found in Elsworth *et al*,  
770 2019 <sup>70</sup>.

771

## 772 **Acknowledgements**

773 M.G. was a National Institute for Health Research (NIHR) academic clinical fellow and is  
774 currently supported by a Wellcome Trust GW4-Clinical Academic Training PhD Fellowship.  
775 This research was funded in part, by the Wellcome Trust [Grant number 220530/Z/20/Z].  
776 For the purpose of open access, the author has applied a CC BY public copyright licence to  
777 any Author Accepted Manuscript version arising from this submission. R.C.R. is a de Pass VC  
778 research fellow at the University of Bristol. J.T. is supported by an Academy of Medical  
779 Sciences (AMS) Springboard award, which is supported by the AMS, the Wellcome Trust,  
780 Global Challenges Research Fund (GCRF), the Government Department of Business, Energy  
781 and Industrial strategy, the British Heart Foundation and Diabetes UK (SBF004\1079).  
782 R.M.M. was supported by a Cancer Research UK (C18281/A20919) programme grant (the  
783 Integrative Cancer Epidemiology Programme). R.M.M. and A.R.N. are supported by the  
784 National Institute for Health Research (NIHR) Bristol Biomedical Research Centre which is  
785 funded by the National Institute for Health Research (NIHR) and is a partnership between  
786 University Hospitals Bristol NHS Foundation Trust and the University of Bristol. Department

787 of Health and Social Care disclaimer: The views expressed are those of the authors and not  
788 necessarily those of the NHS, the NIHR or the Department of Health and Social Care. This  
789 publication presents data from the Head and Neck 5000 study. The study was a component  
790 of independent research funded by the National Institute for Health Research (NIHR) under  
791 its Programme Grants for Applied Research scheme (RP-PG-0707-10034). The views  
792 expressed in this publication are those of the author(s) and not necessarily those of the  
793 NHS, the NIHR or the Department of Health. Core funding was also provided through awards  
794 from Above and Beyond, University Hospitals Bristol and Weston Research Capability  
795 Funding and the NIHR Senior Investigator award to A.R.N. Human papillomavirus (HPV)  
796 serology was supported by a Cancer Research UK Programme Grant, the Integrative Cancer  
797 Epidemiology Programme (C18281/A20919). B.D. and the University of Pittsburgh head and  
798 neck cancer case-control study are supported by US National Institutes of Health (NIH)  
799 grants: P50 CA097190, P30 CA047904 and R01 DE025712. The genotyping of the HNSCC  
800 cases and controls was performed at the Center for Inherited Disease Research (CIDR) and  
801 funded by the US National Institute of Dental and Craniofacial Research (NIDCR;  
802 1X01HG007780-0). The University of North Carolina (UNC) CHANCE study was supported in  
803 part by the National Cancer Institute (R01-CA90731). E.E.V is supported by Diabetes UK  
804 (17/0005587). E.E.V is also supported by the World Cancer Research Fund (WCRF UK), as  
805 part of the World Cancer Research Fund International grant programme (IIG\_2019\_2009).  
806 E.H.T and P.S. were supported by FAPESP grant 10/51168-0 (GENCAPO/Head and Neck  
807 Genome project). M.G., T.D., K.B., A.C., R.M.M., M.M., G.D.S, E.E.V. and R.C.R are part of the  
808 Medical Research Council Integrative Epidemiology Unit at the University of Bristol  
809 supported by the Medical Research Council (MC\_UU\_00011/1, MC\_UU\_00011/5,  
810 MC\_UU\_00011/6, MC\_UU\_00011/7).

811

## 812 **Competing Interests**

813           The authors declare no conflicts of interest. The funders had no role in the design of  
814 the study; the collection, analysis, and interpretation of the data; the writing of the  
815 manuscript; and the decision to submit the manuscript for publication. The authors alone  
816 are responsible for the views expressed in this article.

817

## 818 **Author Contributions**

819 M.G. and R.C.R. conceived the study and M.G. carried out data curation and analysis,  
820 validating the results separately. L.K. completed both the HPV and cervical cancer GWAS  
821 and helped with interpretation of these data. T.W. and N.B. produced serology data for HPV  
822 in the UK Biobank pilot and provided expertise on interpretation of these data. Head and  
823 neck cancer summary genetic data was obtained through multiple collaborations from  
824 studies lead by A.R.N., S.T., A.F.O., R.J.H., G.L., B.D., S.B., E.T., P.S., T.N.T., M.L. and P.B. The  
825 initial manuscript was drafted by M.G., L.K., G.D.S. and R.C.R. Expert guidance on MR  
826 methodology was provided by L.K., T.D., K.B., B.D., R.C.R., G.D.S, R.M.M. All authors M.G.,  
827 T.D., L.K., K.B., A.C., R.M.M., S.T., J.T., A.R.N., P.B., M.M., M.P., S.B., A.F.O., B.D., R.J.H., G.L.,  
828 E.T., P.S., T.N.T., M.L., T.W., N.B., G.D.S., E.V. and R.C.R. contributed to the interpretation of  
829 the results and critical revision of the manuscript. M.G. supervisory team includes R.C.R.,  
830 E.V., J.T., A.R.N and G.D.S.

831

832 **Data availability**

833

834 Summary-level analysis was conducted using publicly available GWAS data. Full summary  
835 statistics for the GAME-ON GWAS can be accessed via dbGAP (OncoArray: Oral and Pharynx  
836 Cancer; study accession number: phs001202.v1.p1) and via the IEU OpenGWAS  
837 project <https://gwas.mrcieu.ac.uk/>. Access to UK Biobank (<https://www.ukbiobank.ac.uk/>)  
838 data is available to researchers through application. For the purpose of open access, the  
839 author has applied a CC BY public copyright licence to any Author Accepted Manuscript  
840 version arising from this submission. UK Biobank approval was given for this project (ID  
841 40644 “Investigating aetiology, associations and causality in diseases of the head and neck”)  
842 and UK Biobank GWAS data was also accessed under the application (ID 15825 “MR-Base:  
843 an online resource for Mendelian randomization using summary data”- Dr Philip Haycock).”

844

845 **Code availability statement**

846

847 MR analyses were conducted using the “TwoSampleMR” package in R (version 3.5.3). A copy  
848 of the code and all files used in this analysis is available at:

849 [https://github.com/rcrichmond/sexual\\_behaviours\\_opc](https://github.com/rcrichmond/sexual_behaviours_opc)

850

851

852

853

854

855

856 **References**

857

- 858 1 Thomas, S. J., Penfold, C. M., Waylen, A. & Ness, A. R. The changing aetiology of head  
859 and neck squamous cell cancer: A tale of three cancers? *Clin Otolaryngol* **43**, 999-  
860 1003 (2018).
- 861 2 Sung, H. *et al.* Global cancer statistics 2020: GLOBOCAN estimates of incidence and  
862 mortality worldwide for 36 cancers in 185 countries. *CA: A Cancer Journal for*  
863 *Clinicians* (2021).
- 864 3 Syrjanen, K., Syrjanen, S., Lamberg, M., Pyrhonen, S. & Nuutinen, J. Morphological  
865 and immunohistochemical evidence suggesting human papillomavirus (HPV)  
866 involvement in oral squamous cell carcinogenesis. *Int J Oral Surg* **12**, 418-424 (1983).
- 867 4 Smith, E. M. *et al.* Human papillomavirus seropositivity and risks of head and neck  
868 cancer. *Int J Cancer* **120**, 825-832 (2007).
- 869 5 D'Souza, G. *et al.* Case-control study of human papillomavirus and oropharyngeal  
870 cancer. *New Engl J Med* **356**, 1944-1956 (2007).
- 871 6 Pan, C., Issaeva, N. & Yarbrough, W. G. HPV-driven oropharyngeal cancer: current  
872 knowledge of molecular biology and mechanisms of carcinogenesis. *Cancers Head*  
873 *Neck* **3**, 12 (2018).
- 874 7 Chaturvedi, A. K. *et al.* Worldwide Trends in Incidence Rates for Oral Cavity and  
875 Oropharyngeal Cancers. *Journal of Clinical Oncology* **31**, 4550-4559 (2013).
- 876 8 Gillison, M. L. *et al.* Evidence for a causal association between human papillomavirus  
877 and a subset of head and neck cancers. *J Natl Cancer Inst* **92**, 709-720 (2000).

- 878 9 Gillison, M. L., Chaturvedi, A. K., Anderson, W. F. & Fakhry, C. Epidemiology of  
879 Human Papillomavirus-Positive Head and Neck Squamous Cell Carcinoma. *J Clin*  
880 *Oncol* **33**, 3235-3242 (2015).
- 881 10 Castellsagué, X. *et al.* HPV Involvement in Head and Neck Cancers: Comprehensive  
882 Assessment of Biomarkers in 3680 Patients. *J Natl Cancer Inst* **108**, djv403 (2016).
- 883 11 Kreimer, A. R. *et al.* Evaluation of human papillomavirus antibodies and risk of  
884 subsequent head and neck cancer. *Journal of clinical oncology : official journal of the*  
885 *American Society of Clinical Oncology* **31**, 2708-2715 (2013).
- 886 12 Kreimer, A. R. *et al.* Kinetics of the Human Papillomavirus Type 16 E6 Antibody  
887 Response Prior to Oropharyngeal Cancer. *JNCI: Journal of the National Cancer*  
888 *Institute* **109** (2017).
- 889 13 Anantharaman, D. *et al.* Human papillomavirus infections and upper aero-digestive  
890 tract cancers: the ARCAGE study. *J Natl Cancer Inst* **105**, 536-545 (2013).
- 891 14 Ribeiro, K. B. *et al.* Low human papillomavirus prevalence in head and neck cancer:  
892 results from two large case-control studies in high-incidence regions. *Int J Epidemiol*  
893 **40**, 489-502 (2011).
- 894 15 Heck, J. E. *et al.* Sexual behaviours and the risk of head and neck cancers: a pooled  
895 analysis in the International Head and Neck Cancer Epidemiology (INHANCE)  
896 consortium. *Int J Epidemiol* **39**, 166-181 (2010).
- 897 16 Herrero, R. *et al.* Human papillomavirus and oral cancer: the International Agency for  
898 Research on Cancer multicenter study. *J Natl Cancer Inst* **95**, 1772-1783 (2003).
- 899 17 Schwartz, S. M. *et al.* Oral cancer risk in relation to sexual history and evidence of  
900 human papillomavirus infection. *J Natl Cancer Inst* **90**, 1626-1636 (1998).

- 901 18 Rajkumar, T. *et al.* Oral cancer in Southern India: the influence of body size, diet,  
902 infections and sexual practices. *Eur J Cancer Prev* **12**, 135-143 (2003).
- 903 19 Smith, E. M. *et al.* Age, sexual behavior and human papillomavirus infection in oral  
904 cavity and oropharyngeal cancers. *Int J Cancer* **108**, 766-772 (2004).
- 905 20 Shah, A. *et al.* Oral sex and human papilloma virus-related head and neck squamous  
906 cell cancer: a review of the literature. *Postgrad Med J* **93**, 704-709 (2017).
- 907 21 Doorbar, J. & Griffin, H. Refining our understanding of cervical neoplasia and its  
908 cellular origins. *Papillomavirus Res* **7**, 176-179 (2019).
- 909 22 Farsi, N. J. *et al.* Sexual behaviours and head and neck cancer: A systematic review  
910 and meta-analysis. *Cancer Epidemiol* **39**, 1036-1046 (2015).
- 911 23 Khadr, S. N. *et al.* Investigating the relationship between substance use and sexual  
912 behaviour in young people in Britain: findings from a national probability survey.  
913 *BMJ Open* **6**, e011961 (2016).
- 914 24 Davey Smith, G. & Hemani, G. Mendelian randomization: genetic anchors for causal  
915 inference in epidemiological studies. *Hum Mol Genet* **23**, R89-98 (2014).
- 916 25 Smith, G. D. & Ebrahim, S. 'Mendelian randomization': can genetic epidemiology  
917 contribute to understanding environmental determinants of disease? *Int J Epidemiol*  
918 **32**, 1-22 (2003).
- 919 26 Karlsson Linner, R. *et al.* Genome-wide association analyses of risk tolerance and  
920 risky behaviors in over 1 million individuals identify hundreds of loci and shared  
921 genetic influences. *Nat Genet* **51**, 245-257 (2019).
- 922 27 Ganna, A. *et al.* Large-scale GWAS reveals insights into the genetic architecture of  
923 same-sex sexual behavior. *Science* **365**, eaat7693 (2019).



- 924 28 Mills, M. C. *et al.* Identification of 370 loci for age at onset of sexual and reproductive  
925 behaviour, highlighting common aetiology with reproductive biology, externalizing  
926 behaviour and longevity. *bioRxiv*, 2020.2005.2006.081273 (2020).
- 927 29 Morrison, J., Knoblach, N., Marcus, J. H., Stephens, M. & He, X. Mendelian  
928 randomization accounting for correlated and uncorrelated pleiotropic effects using  
929 genome-wide summary statistics. *Nat Genet* (2020).
- 930 30 Bowden, J. *et al.* Assessing the suitability of summary data for two-sample  
931 Mendelian randomization analyses using MR-Egger regression: the role of the I<sup>2</sup>  
932 statistic. *Int J Epidemiol* **45**, 1961-1974 (2016).
- 933 31 Hemani, G., Bowden, J. & Smith, G. D. Evaluating the potential role of pleiotropy in  
934 Mendelian randomization studies. *Hum Mol Genet* **27**, R195-R208 (2018).
- 935 32 Morrison, J., Knoblach, N., Marcus, J. H., Stephens, M. & He, X. Mendelian  
936 randomization accounting for correlated and uncorrelated pleiotropic effects using  
937 genome-wide summary statistics. *Nat Genet* **52**, 740-747 (2020).
- 938 33 Elsworth, B. *et al.* The MRC IEU OpenGWAS data infrastructure. *bioRxiv*,  
939 2020.2008.2010.244293 (2020).
- 940 34 Wootton, R. E. *et al.* Evidence for causal effects of lifetime smoking on risk for  
941 depression and schizophrenia: a Mendelian randomisation study. *Psychol Med*, 1-9  
942 (2019).
- 943 35 Liu, M. Z. *et al.* Association studies of up to 1.2 million individuals yield new insights  
944 into the genetic etiology of tobacco and alcohol use. *Nat Genet* **51**, 237-+ (2019).
- 945 36 Pan, C., Issaeva, N. & Yarbrough, W. G. HPV-driven oropharyngeal cancer: current  
946 knowledge of molecular biology and mechanisms of carcinogenesis. *Cancers of the*  
947 *Head & Neck* **3**, 12 (2018).

- 948 37 Chung, C. H. & Gillison, M. L. Human Papillomavirus in Head and Neck Cancer: Its  
949 Role in Pathogenesis and Clinical Implications. *Clinical Cancer Research* **15**, 6758-  
950 6762 (2009).
- 951 38 Brenner, N. *et al.* Characterization of human papillomavirus (HPV) 16 E6 seropositive  
952 individuals without HPV-associated malignancies after 10 years of follow-up in the  
953 UK Biobank. *EBioMedicine* **62**, 103123 (2020).
- 954 39 Mentzer, A. J. *et al.* Identification of host-pathogen-disease relationships using a  
955 scalable Multiplex Serology platform in UK Biobank. *medRxiv*, 19004960 (2019).
- 956 40 Dahlstrom, K. R. *et al.* HPV Serum Antibodies as Predictors of Survival and Disease  
957 Progression in Patients with HPV-Positive Squamous Cell Carcinoma of the  
958 Oropharynx. *Clinical cancer research : an official journal of the American Association*  
959 *for Cancer Research* **21**, 2861-2869 (2015).
- 960 41 Perdomo, S., Martin Roa, G., Brennan, P., Forman, D. & Sierra, M. S. Head and neck  
961 cancer burden and preventive measures in Central and South America. *Cancer*  
962 *Epidemiology* **44**, S43-S52 (2016).
- 963 42 Kreimer, A. R., Clifford, G. M., Boyle, P. & Franceschi, S. Human papillomavirus types  
964 in head and neck squamous cell carcinomas worldwide: a systematic review. *Cancer*  
965 *Epidemiol Biomarkers Prev* **14**, 467-475 (2005).
- 966 43 Dayyani, F. *et al.* Meta-analysis of the impact of human papillomavirus (HPV) on  
967 cancer risk and overall survival in head and neck squamous cell carcinomas (HNSCC).  
968 *Head Neck Oncol* **2**, 15 (2010).
- 969 44 Gayet, C., Juarez, F. & Bozon, M. Vol. 5 67-90 (2013).
- 970 45 Bosetti, C. *et al.* Global trends in oral and pharyngeal cancer incidence and mortality.  
971 *Int J Cancer* **147**, 1040-1049 (2020).

- 972 46 Schache, A. G. *et al.* HPV-Related Oropharynx Cancer in the United Kingdom: An  
973 Evolution in the Understanding of Disease Etiology. *Cancer Res* **76**, 6598-6606  
974 (2016).
- 975 47 Kreimer, A. R. *et al.* Incidence and clearance of oral human papillomavirus infection  
976 in men: the HIM cohort study. *Lancet (London, England)* **382**, 877-887 (2013).
- 977 48 Castle, P. E. How does tobacco smoke contribute to cervical carcinogenesis? *J Virol*  
978 **82**, 6084-6086 (2008).
- 979 49 Arnson, Y., Shoenfeld, Y. & Amital, H. Effects of tobacco smoke on immunity,  
980 inflammation and autoimmunity. *J Autoimmun* **34**, J258-265 (2010).
- 981 50 Hashibe, M. *et al.* Interaction between Tobacco and Alcohol Use and the Risk of  
982 Head and Neck Cancer: Pooled Analysis in the International Head and Neck Cancer  
983 Epidemiology Consortium. *Cancer Epidem Biomar* **18**, 541-550 (2009).
- 984 51 Gormley, M. *et al.* A multivariable Mendelian randomization analysis investigating  
985 smoking and alcohol consumption in oral and oropharyngeal cancer. *Nat Commun*  
986 **11**, 6071 (2020).
- 987 52 Bowden, J., Davey Smith, G. & Burgess, S. Mendelian randomization with invalid  
988 instruments: effect estimation and bias detection through Egger regression. *Int J*  
989 *Epidemiol* **44**, 512-525 (2015).
- 990 53 Bowden, J., Smith, G. D., Haycock, P. C. & Burgess, S. Consistent Estimation in  
991 Mendelian Randomization with Some Invalid Instruments Using a Weighted Median  
992 Estimator. *Genet Epidemiol* **40**, 304-314 (2016).
- 993 54 Hartwig, F. P., Smith, G. D. & Bowden, J. Robust inference in summary data  
994 Mendelian randomization via the zero modal pleiotropy assumption. *Int J Epidemiol*  
995 **46**, 1985-1998 (2017).

- 996 55 Mounier, N. & Kutalik, Z. Correction for sample overlap, winner's curse and weak  
997 instrument bias in two-sample Mendelian Randomization. *bioRxiv*,  
998 2021.2003.2026.437168 (2021).
- 999 56 Minelli, C. *et al.* The use of two-sample methods for Mendelian randomization  
1000 analyses on single large datasets. *bioRxiv*, 2020.2005.2007.082206 (2020).
- 1001 57 Lawlor, D. A., Tilling, K. & Davey Smith, G. Triangulation in aetiological epidemiology.  
1002 *Int J Epidemiol* **45**, 1866-1886 (2017).
- 1003 58 Munafò, M. R. & Davey Smith, G. Robust research needs many lines of evidence.  
1004 *Nature* **553**, 399-401 (2018).
- 1005 59 Demange, P. A. *et al.* Investigating the genetic architecture of noncognitive skills  
1006 using GWAS-by-subtraction. *Nat Genet* **53**, 35-44 (2021).
- 1007 60 Lesseur, C. *et al.* Genome-wide association analyses identify new susceptibility loci  
1008 for oral cavity and pharyngeal cancer. *Nat Genet* **48**, 1544-1550 (2016).
- 1009 61 Dudding, T. *et al.* Assessing the causal association between 25-hydroxyvitamin D and  
1010 the risk of oral and oropharyngeal cancer using Mendelian randomization. *Int J*  
1011 *Cancer* **143**, 1029-1036 (2018).
- 1012 62 Kachuri, L. *et al.* The landscape of host genetic factors involved in immune response  
1013 to common viral infections. *Genome Medicine* **12**, 93 (2020).
- 1014 63 Sanderson, E., Richardson, T. G., Hemani, G. & Davey Smith, G. The use of negative  
1015 control outcomes in Mendelian randomization to detect potential population  
1016 stratification. *Int J Epidemiol* (2021).
- 1017 64 Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74  
1018 (2015).

- 1019 65 Altshuler, D. M. *et al.* Integrating common and rare genetic variation in diverse  
1020 human populations. *Nature* **467**, 52-58 (2010).
- 1021 66 Sanderson, E., Davey Smith, G., Windmeijer, F. & Bowden, J. An examination of  
1022 multivariable Mendelian randomization in the single-sample and two-sample  
1023 summary data settings. *Int J Epidemiol* (2018).
- 1024 67 Rees, J. M. B., Wood, A. M. & Burgess, S. Extending the MR-Egger method for  
1025 multivariable Mendelian randomization to correct for both measured and  
1026 unmeasured pleiotropy. *Stat Med* **36**, 4705-4718 (2017).
- 1027 68 Burgess, S., Dudbridge, F. & Thompson, S. G. Re: "Multivariable Mendelian  
1028 randomization: the use of pleiotropic genetic variants to estimate causal effects".  
1029 *Am J Epidemiol* **181**, 290-291 (2015).
- 1030 69 Yavorska, O. O. & Burgess, S. MendelianRandomization: an R package for performing  
1031 Mendelian randomization analyses using summarized data. *Int J Epidemiol* **46**, 1734-  
1032 1739 (2017).
- 1033 70 Mitchell, R., Elsworth, BL, Mitchell, R, Raistrick, CA, Paternoster, L, Hemani, G, Gaunt,  
1034 TR. (2019).  
1035