

1 Integrated analysis of cervical squamous cell carcinoma
2 cohorts from three continents reveals conserved subtypes of
3 prognostic significance.

4

5

6 Ankur Chakravarthy^{1,*}, Ian Reddin^{2,*}, Stephen Henderson³, Cindy Dong⁴, Nerissa
7 Kirkwood⁴, Maxmilan Jeyakumar⁴, Daniela Rothschild Rodriguez⁴, Natalia Gonzalez
8 Martinez⁴, Jacqueline McDermott⁵, Xiaoping Su⁶, Nagayasau Egawa⁷, Christina S
9 Fjeldbo⁸, Vilde Eide Skingen⁸, Mari Kyylesø Halle¹⁰, Camilla Krakstad¹⁰, Afschin
10 Soleiman¹¹, Susanne Sprung¹², Peter Ellis⁴, Mark Wass⁴, Martin Michaelis⁴, Heidi
11 Lyng^{8, 9}, Heidi Fiegl¹³, Helga Salvesen¹⁰, Gareth Thomas², John Doorbar⁷, Kerry
12 Chester^{5, #}, Andrew Feber^{14, 15, #}, Tim R Fenton^{2, 4, #}.

13

14

15 ¹ Princess Margaret Cancer Centre, University Health Network. Toronto, Ontario,
16 Canada

17 ² School of Cancer Sciences, Cancer Research UK Centre, Faculty of Medicine,
18 University of Southampton, Southampton, UK

19 ³ UCL Cancer Institute, Bill Lyons Informatics Centre, University College London,
20 London, UK

21 ⁴ School of Biosciences, Division of Natural Sciences, University of Kent,
22 Canterbury, UK

23 ⁵ UCL Cancer Institute, University College London, London, UK

24 ⁶ MD Anderson Cancer Center, Houston, Texas, USA

25 ⁷ Department of Pathology, University of Cambridge, Cambridge, UK

26 ⁸ Department of Radiation Biology, Oslo University Hospital, Oslo, Norway

27 ⁹ Department of Physics, University of Oslo, Oslo, Norway

28 ¹⁰ Department of Obstetrics and Gynaecology, Haukeland University Hospital,
29 Bergen, Norway; Centre for Cancer Biomarkers, Department of Clinical Science,
30 University of Bergen, Norway

31 ¹¹ INNPATh, Institute of Pathology, Tirol Kliniken Innsbruck, Innsbruck, Austria.

32 ¹² Institute of Pathology, Medical University of Innsbruck, Innsbruck, Austria.

33 ¹³ Department of Obstetrics and Gynaecology, Medical University of Innsbruck,
34 Innsbruck, Austria.

35 ¹⁴ Centre for Molecular Pathology, Royal Marsden Hospital Trust, London, UK

36 ¹⁵ Division of Surgery and Interventional Science, University College London,
37 London, UK

38

39 *Equal contribution

40 #Correspondence to: k.chester@ucl.ac.uk; a.feber@ucl.ac.uk or

41 t.fenton@soton.ac.uk

42

43

44

45

46

47

48

49 Abstract

50

51 Human papillomavirus (HPV)-associated cervical cancer represents one of the
52 leading causes of cancer death worldwide. Although low-middle income countries
53 are disproportionately affected, our knowledge of the disease predominantly
54 originates from populations in high-income countries. Using the largest multi-omic
55 analysis of cervical squamous cell carcinoma (CSCC) to date, totalling 643 tumours
56 and representing patient populations from the USA, Europe and Sub-Saharan Africa,
57 we identify two CSCC subtypes (C1 and C2) with differing prognosis. C1 tumours
58 are largely HPV16-driven, display increased cytotoxic T-lymphocyte infiltration and
59 frequently harbour *PIK3CA* and *EP300* mutations. C2 tumours are associated with
60 shorter overall survival, are frequently driven by HPVs from the HPV18-containing
61 alpha-7 clade, harbour alterations in the Hippo signalling pathway and increased
62 expression of immune checkpoint genes, *B7-H3* (also known as *CD276*) and *NT5E*
63 (also known as *CD73*) and *PD-L2* (also known as *PDCD1LG2*). In conclusion, we
64 identify two novel, therapy-relevant CSCC subtypes that share the same defining
65 characteristics across three geographically diverse cohorts.

66

67 --

68

69 Despite screening and the introduction of prophylactic human papillomavirus (HPV)
70 vaccination in developed countries, cervical cancer continues to be one of the
71 leading worldwide causes of cancer-related deaths in women¹. Prognosis for
72 patients with metastatic disease remains poor, thus new treatments and effective

73 molecular markers for patient stratification are urgently required. Cervical cancer is
74 caused by at least 14 high-risk human papillomaviruses (hrHPVs), with HPV16 and
75 HPV18 together accounting for over 70% of cases worldwide, with some variation by
76 region^{1,2}. Cervical squamous cell carcinoma (CSCC) is the most common
77 histological subtype of cervical cancer, accounting for approximately 60-70% of
78 cases, again with some variation seen across different populations². Adeno- and
79 adenosquamous histology are both associated with poor prognosis³⁻⁶, while the
80 relationship, if any, between HPV type and cervical cancer prognosis remains
81 unclear⁷. HPV type is also associated with histology, with HPV16 more commonly
82 found in CSCC, while adenocarcinomas are more likely to harbour HPV18². Previous
83 landmark studies described the genomic landscape of cervical cancer in different
84 populations⁸⁻¹¹ and in some cases identified subtypes based on gene expression,
85 DNA methylation and/or proteomic profiles^{8,9}. The Cancer Genome Atlas (TCGA)
86 network identified clusters based on RNA, micro-RNA, protein/phospho-protein, DNA
87 copy number alterations and DNA methylation patterns and combined data from
88 multiple platforms to define integrated iClusters⁸. In their analysis, only clustering
89 based on the expression levels and/or phosphorylation state of 192 proteins as
90 measured by reverse-phase protein array (RPPA) was associated with outcome,
91 with significantly shorter overall survival (OS) observed for a cluster of cervical
92 cancers exhibiting increased expression of Yes-associated protein (YAP) and
93 features associated with epithelial-to-mesenchymal transition (EMT) and a reactive
94 tumour stroma. Since TCGA's RPPA analysis was restricted to 155 tumours
95 including SCCs, adeno- and adenosquamous carcinomas, we set out to test the
96 hypothesis that with data from more samples, we could identify a set of
97 transcriptional and epigenetic features associated with prognosis within CSCC and

98 to establish whether it is also present in independent patient cohorts representing
99 different geographical locations and ethnicities. To identify molecular subtypes and
100 prognostic correlates, we identified a set of 643 CSCCs (all HPV-positive), for which
101 clinico-pathological data and genome-wide DNA methylation profiles were either
102 publicly available or generated in this study, and for which in most cases, matched
103 gene expression and somatic mutation data were also available (Table 1).

104

105 Table 1: Summary of clinicopathological characteristics for five cervical cancer
106 cohorts.

107

	Discovery Cohort		Validation Cohorts			Total
	TCGA	Bergen	Innsbruck	Oslo	Uganda	
Cohort Numbers	236	37	28	248	94	643
Stage						
I	122	33	16	23	15	209
II	54	3	6	173	45	281
III	39	0	4	61	31	135
IV	14	1	2	13	2	32
NA	7	0	0	0	1	8
Age						
Median (Range)	47 (20-88)	42 (28-64)	49 (29-91)	54 (22-82)	45 (26-82)	
HPV Type						
16	136	22	14	168	39	379
18	26	2	5	30	14	77
45	19	4	0	9	14	46
Other	55	9	9	41	27	141
HPV Clade						
Alpha 7	57	6	6	40	32	141
Alpha 9	172	29	18	183	54	456
Other	7	2	4	25	8	46
Treatment						
Surgery alone	NA	18	5	0	NA	23
Surgery and radiotherapy	NA	14	17	0	NA	31
Surgery, radiotherapy and chemotherapy	NA	5	6	0	NA	11
Radiotherapy alone	NA	0	0	47	NA	47
Radiotherapy and chemotherapy	NA	0	0	201	NA	201
Overall Survival						
Median (Range)	1.9 (0-17.5)	8 (1.8-13.2)	9.8 (0.1-23.2)	4.2 (0.3-12.7)	1.1 (0-2.4)	
Survival Status						
Alive	181	32	16	187	42	458
Dead	55	5	12	61	52	185
Cluster Assignment						
C1	175	32	24	198	69	498
C2	61	5	4	50	25	145
% C2	25.8	13.5	14.3	20.2	26.6	22.6
HIV Status						
Positive	NA	NA	NA	NA	59	59
Negative	NA	NA	NA	NA	35	35
Available Data						
RNA-seq	236	37	NA	NA	94	367
Methylation	236	37	28	248	94	643
Mutation	236	37	NA	NA	94	367
RPPA	137	NA	NA	NA	NA	137
Microarray (Illumina HumanWG-6 v3)	NA	NA	NA	137	NA	137
Microarray (Illumina HumanHT-12 v4)	NA	NA	NA	109	NA	109
Copy number (methylation derived)	236	37	28	248	NA	549

109

110 **Identification of two gene expression-based clusters in cervical squamous cell**
111 **carcinoma**

112

113 Molecular and clinical differences between cervical adeno/adenosquamous and
114 CSCCs are well documented¹²⁻¹⁴ and gene expression differences were apparent in
115 multi-dimensional TSNE analysis based on the top 10% most variable genes of three
116 previously published cervical cancer cohorts^{8,9,11} with available RNA-seq data
117 (Supplementary Fig. S1a-d). To examine molecular and clinical heterogeneity
118 specifically within SCC we focused all subsequent analysis on a collection of
119 confirmed HPV-positive CSCCs from the USA, Europe and Uganda, as shown in
120 Table 1.

121

122 236 cervical SCCs profiled by TCGA were defined as our discovery cohort (Table 1,
123 Supplementary Table S1) and consensus clustering was performed using the top
124 10% most variable genes (n=1377 genes, Supplementary Table S2). Consensus
125 cluster membership heatmaps, delta area plot, consensus cumulative distribution
126 function (CDF) and proportion of ambiguous clusters (PAC) indicated the optimal
127 number of clusters was two (Fig. 1a, Supplementary Fig. 2), the larger of which
128 (n=175) was designated C1 while the smaller cluster (n=61) was designated C2
129 (Supplementary Table S1). Modelling transcriptomic differences between these two
130 clusters identified 938 differentially expressed genes (DEGs, FDR=0.01, FC > 2)
131 (Fig. 1b, Supplementary Table S3). Tumours in C1 predominantly harbour HPV
132 types from the HPV16-containing alpha-9 clade (150/175) while 38 of 61 C2 tumours
133 contained HPV types from the HPV18-containing alpha-7 clade. C2 tumours were

134 13.3 times more likely to harbour alpha 7 HPVs than C1 tumours ($p = 1.8 \times 10^{-14}$,
135 Fishers Exact Test) (Fig. 1b).

136

137 Univariate analysis of 5-year overall survival (OS) revealed worse outcomes for
138 patients with C2 tumours (HR = 2.54, $p = 0.001$; Fig. 1c) and in Cox regression
139 including age, tumour stage and HPV type as covariates along with cluster
140 membership, only membership of the C2 cluster (HR = 2.44, $p = 0.017$ 95% CI 1.18,
141 5.05) and a tumour stage of IV (HR versus stage I = 4.74, $p < 0.001$, 95% CI 2.1,
142 10.7) were independent predictors of five-year OS (Table 2). The relationship
143 between cluster and OS is also clear when restricting the analysis to HPV16-
144 containing tumours in each cluster in both univariate analysis (HR = 3.39, $p = 0.004$;
145 Fig. 1d) and multivariate analysis, including age and tumour stage as covariates (HR
146 = 3.89, $p = 0.003$, 95% CI 1.57, 9.67; Supplementary Table S4).

147

148

149 **Identification of C1 and C2 CSCCs and association with prognosis in** 150 **independent SCC cohorts**

151

152 To further investigate the association between C1/C2 cluster membership and OS,
153 we assembled a combined validation cohort consisting of 313 CSCC patients treated
154 at three centres in Europe (Bergen ($n = 37$), Oslo ($n = 248$) and Innsbruck ($n = 28$)),
155 for which detailed clinical information were available and for which genome-wide
156 DNA methylation profiles from Illumina Infinium 450k arrays (the same platform used
157 by TCGA) were either available or generated in this study (Table 1). Since RNA-seq
158 data were not available for all European samples, cluster membership was assigned

159 using a support vector machine (SVM) classification model based on 129 CpG sites
160 (methylation variable positions, MVPs) at which methylation differed significantly
161 between tumours in C1 vs C2 clusters in the discovery cohort (Fig. 2a, b; mean
162 delta-Beta > 0.25, FDR < 0.01, Supplementary Table S5), 18 of which were located
163 within 12 genes differentially expressed between the clusters (Supplementary Fig.
164 S3). MVP and DEG signatures were also used to assign cluster membership to 94
165 CSCCs from a Ugandan cohort originally profiled by the Cancer Genome
166 Characterization Initiative (CGCI)⁹, for which both DNA methylation and RNA-seq
167 data were available. C2 tumours from all cohorts clustered together using TSNE
168 analysis based on the MVP signature (Fig. 2c) and high concordance between DEG
169 and MVP-based cluster allocation was observed in all cohorts for which both gene
170 expression (RNA-seq for Uganda and Bergen or Illumina bead chip arrays for Oslo)
171 and DNA methylation data were available (Supplementary Fig. S4a, b). Single-
172 sample gene set enrichment analysis (ssGSEA) confirmed differential expression of
173 the signature genes in tumours classified as C1 or C2 using DNA methylation data
174 (Supplementary Fig. 4c). 59 of 313 (18.8%) tumours in the combined European
175 cohort (Fig. 2b, Supplementary Table S7) and 25 of 94 (26.6%) tumours in the
176 Ugandan cohort were classified as C2 (Supplementary Fig. S5a, Supplementary
177 Table S6). As in the discovery cohort, most C1 tumours from the European and
178 Ugandan cohorts harboured alpha-9 HPV types (260/325) while C2 tumours were
179 3.9 times more likely to harbour alpha-7 HPVs than C1 tumours ($p = 1.07 \times 10^{-6}$,
180 Fishers Exact Test) (Fig. 2b, Supplementary Fig. S5a). Interestingly 80% (20/25) of
181 Ugandan C2 patients were human immunodeficiency virus (HIV) positive, while only
182 56% (39/69) of C1 patients were HIV positive (Supplementary Fig. S5a).

183

184 Univariate analysis indicated lower 5-year OS in C2 tumours from the European
185 cohort (Fig. 2d) and Cox regression controlling for FIGO stage, age, HPV type and
186 treatment (surgery alone, surgery with radio-chemotherapy, surgery with
187 radiotherapy alone, radio-chemotherapy and radiotherapy alone) again identified C2
188 status but not HPV type to be an independent predictor of 5-year OS (HR = 2.54, p
189 =0.003, 95% CI 1.4, 4.7) along with tumour stage and inclusion of chemotherapy in
190 the treatment regimen (Table 2). As in the discovery cohort, a significant prognostic
191 difference was identified between the C1 and C2 subgroups when considering only
192 the HPV16-positive tumours (n = 204) in both univariate (Supplementary Fig. S5b)
193 and multivariate analyses (HR = 2.64, p = 0.02, 95% CI = 1.16, 6; Supplementary
194 Table S4). Interestingly the prognostic difference was even greater among 78
195 patients in the European cohort that did not receive chemotherapy (Supplementary
196 Fig. S5c; multivariate HR = 4.4, p = 0.005, 95% CI = 1.58, 12.3). At 94 patients, the
197 Ugandan cohort was underpowered for comparing survival between C1 and C2
198 tumours and survival rates in the Ugandan cohort were much lower than in the other
199 cohorts (Supplementary Fig. S5d), thus we did not attempt a combined survival
200 analysis including these patients. Taken together, the C1/C2 clusters identified in the
201 TCGA cohort (USA) are apparent in cohorts of CSCC patients from Europe and
202 Uganda and tumours can be accurately assigned to cluster using either gene
203 expression or DNA methylation profiles. C1/C2 cluster is an independent predictor of
204 5-year OS in both the TCGA (n = 236) and European (n = 313) cohorts and remains
205 so when only HPV16+ tumours are considered. There is no difference in the
206 breakdown of C1 and C2 tumours by stage (Supplementary Table S7).

207

208

209

210 Table 2 – Five-year survival analysis for all cohorts

211

	Univariate			Multivariate		
	Hazard Ratio	p Value	95% CI	Hazard Ratio	p Value	95% CI
TCGA	2.54	0.002	1.42, 4.56	2.44	0.02	1.18, 5.05
Bergen	5.28	0.07	0.87, 31.9	98.1	<0.001	8.41, 1145
Innsbruck	0	1	0, Inf	0	1	0, Inf
Oslo	1.74	0.07	0.96, 3.14	2.36	0.012	1.21, 4.62
Europe Combined	1.68	0.07	0.97, 2.90	2.54	0.003	1.40, 4.67
Uganda	NA	NA	NA	NA	NA	NA

212

213

214

215 Relationships between C1/C2 and clusters previously identified by TCGA

216

217 Of the 178 tumour samples that made up the core set in the TCGA's landmark study
 218 into cervical cancer genomics/epigenomics⁸, 140 CSCCs were present in our
 219 discovery cohort of 236 (Supplementary Table S8). This enabled comparisons
 220 between our gene expression-based cluster allocations and the subtypes defined by
 221 TCGA (Fig. 3). TCGA analysis included integrated clustering using multiomics data
 222 (three iClusters, two of which ('keratin-high' and 'keratin-low' were composed entirely
 223 of CSCCs) and clustering based on transcriptomic data (three mRNA clusters).
 224 There is considerable overlap between our C1 cluster and TCGA's mRNA C2 cluster
 225 (84/106) and keratin-high iCluster (80/106), and between our C2 cluster and TCGA's

226 mRNA C3 cluster (19/34) and keratin-low iCluster (27/34). Neither the mRNA C3 nor
227 the keratin-low iCluster were associated with poor prognosis in TCGA's analysis and
228 given the increased expression of a subset of keratin genes (including *KRT7*, *KRT8*
229 and *KRT18*) in C2 tumours (Fig. 3), we decided against adopting the keratin-high /
230 keratin-low nomenclature for our clusters. We also examined the relationship
231 between our subtypes and three clusters defined by TCGA based on reverse phase
232 protein array (RPPA) data. Notably, 57% of C2 TCGA tumours with RPPA data
233 available belong to the EMT cluster compared with only 25% of C1 tumours (Fig. 3)
234 and, consistent with the proteomic classification, C2 tumours display higher EMT
235 mRNA expression scores, as defined by TCGA⁸ than C1 tumours (Supplementary
236 Fig. S6). Although there is greater concordance between C2 and the TCGA EMT
237 cluster compared to C1, it is clearly distinct from the EMT cluster.

238

239 **Genomic analyses of prognostic clusters**

240

241 To investigate whether C1 and C2 tumours differ at the genomic level in addition to
242 the transcriptomic and epigenomic differences observed above, whole-exome data
243 was obtained for SCCs from three cohorts, TCGA⁸, Bergen¹¹ and Uganda⁹. This
244 amounted to 367 samples, 29 of which were classed as hypermutated by standards
245 set by TCGA⁸ (>600 mutations). The median tumour mutation burden (TMB) was
246 2.04/Mb for all tumour, 2.11/Mb for C1 tumours and 1.82/Mb for C2 tumours
247 (1.92/Mb, 1.94/Mb and 1.72/Mb respectively after removal of hypermutated
248 samples). We detected four mutation signatures for the combined cohorts
249 (Supplementary Fig. S7): as expected based on previous studies^{8,9,11}, COSMIC
250 signatures 2 and 13 (characterised by C>T transitions or C>G transversions

251 respectively at TpC sites attributed to cytosine deamination by APOBEC enzymes);
252 age-related COSMIC signature 1 (characterised by C>T transitions attributed to
253 spontaneous deamination of 5' methylated cytosine) and COSMIC signature 5, for
254 which the underlying mutational process is unknown¹⁵
255 (<https://cancer.sanger.ac.uk/signatures/>). The proportion of mutations attributable to
256 each signature did not vary between clusters (Fig. 4).

257

258 Having excluded the hypermutated samples, we next performed dNdScv analysis¹⁶
259 on each cohort, followed by p-value combination using sample size weighted
260 Fisher's method followed by FDR correction¹⁷ to permit identification of significantly
261 mutated genes (SMGs) across the entire dataset. This combined approach, followed
262 by analysis of individual samples by cluster identified 34 SMGs (Fig. 4,
263 Supplementary Table S9), 21 of which (highlighted by †) have not previously been
264 identified as SMGs in cervical cancer^{8,9,11}. Of the 34 SMGs, 21 were significantly
265 mutated in only C1 samples, two genes in only C2 samples, three genes in both C1
266 and C2 individual analysis, and eight genes were only significantly mutated when
267 both C1 and C2 clusters were analysed together (Fig. 4, Supplementary Table S9).
268 The frequency of mutations in SMGs that had been previously observed was
269 comparable between combined cohort and each respective SMG study
270 (Supplementary Table S10). Among the 21 genes that have not previously been
271 identified as significantly mutated in cervical cancer, six are SMGs in other SCCs,
272 including head and neck (*NOTCH1*, *JUB* (also known as *AJUBA*), *MLL2* (also known
273 as *KMT2D*), *RB1*, *PIK3R1*)¹⁸, oesophageal (*MLL2*, *NOTCH1*, *RB1*)¹⁹ and lung SCC
274 (*NOTCH1*, *RB1*, *MLL2*, *CREBBP* (also known as *KAT3A*))²⁰. Conversely, several
275 genes previously identified as SMGs in cervical cancer, including *TP53*, *ARID1A* and

276 *TGFBR2* are significantly mutated in adenocarcinoma but not in CSCC^{8,11}.
277 Comparing somatic mutation rates in SMGs between clusters using binomial
278 regression identified *PIK3CA* (FDR = 0.001) and *EP300* (FDR = 0.046) mutations as
279 disproportionally more common in C1 tumours and *STK11* (FDR = 0.005) and *NF2*
280 (FDR = 0.045) as enriched in C2 tumours (Fig. 4). *STK11* is also under-expressed in
281 C2 tumours compared with C1 tumours (Supplementary Table S3).

282

283 **C2 tumours display Hippo pathway alterations and increased YAP1 activity**

284

285 Two SMGs from our analysis (*LATS1* and *NF2*) are core members of the HIPPO
286 signalling pathway, while SMGs *FAT1*, *JUB* and *STK11* are known regulators of
287 HIPPO signalling²¹⁻²³. Mutations in *LATS1*, *FAT1*, *JUB*, *STK11* or *NF2* (the latter two
288 of which are significantly mutated specifically in C2 tumours, Fig. 4) result in aberrant
289 activation of the downstream transcription factor, yes1 associated transcriptional
290 regulator (YAP1)²⁴⁻²⁸, the expression of which is also elevated at the mRNA level in
291 C2 tumours (Table S3).

292

293 We generated segmented copy number data for all tumours (combining TCGA and
294 European validation cohort samples for which the necessary data were available for
295 maximum statistical power), which identified 211 focal candidate copy number
296 alterations (CNAs) at FDR < 0.1. Following binomial regression, we identified five
297 discrete CNAs that differed in frequency between C1 and C2 clusters (Fig. 5a; FDR
298 < 0.1, log₂ (Odds Ratio) > 1). All five were more prevalent in C2 tumours and
299 included 11q11 and 1q21.2 deletions and 6p22.1, 11q22.1 and 11q22.2 gains.
300 11q22.2 contains matrix metalloproteinase genes (MMPs) which are well known to

301 be involved in metastasis²⁹, but notably 11q22.1 contains the *YAP1* gene.
302 Furthermore, analysis of Reverse Phase Protein Assay (RPPA) data from TCGA
303 revealed significantly higher *YAP1* protein expression in C2 tumours (Fig. 5b). We
304 confirmed that of the 137 TCGA cases for which RPPA data were available, cases
305 with *YAP1* amplification (8/37 C2 tumours and 6/100 C1 tumours) also showed
306 increased *YAP1* mRNA and protein expression (Supplementary Fig. S8). In total 10
307 genes from a 22 gene signature that predicts HIPPO pathway activity in cancer³⁰ are
308 differentially expressed between C1 and C2 tumours (Supplementary Table S3).

309

310 **Differences in the tumour immune microenvironment between C1 and C2** 311 **tumours.**

312

313 The nature of the tumour immune microenvironment, particularly the abundance of
314 tumour infiltrating lymphocytes (TILs) is a strong prognostic factor in cervical
315 cancer³¹⁻³³. We used DNA methylation data to compare the cellular composition of
316 TCGA tumours³⁴, observing differences in the proportions of multiple cell types
317 between the subgroups (Fig. 6a); most notably decreased CD8+ (cytotoxic T
318 lymphocytes (CTL)), and a marked elevation of neutrophil and CD56+ natural killer
319 (NK)-cells in C2 tumours. Repeating this method with the validation cohorts
320 produced results that were remarkably similar (Fig. 6b). Differences in the
321 proportions of cell types between C1 and C2 in the validation cohort mirrored those
322 in the TCGA cohort, decreased CTL, and elevated neutrophil, NK-cell and
323 endothelial cell levels were observed in C2 tumours. Importantly, this was not driven
324 by any single validation cohort, as individual cohorts displayed consistent patterns of
325 differences in the proportion of cell types between C1 and C2 tumours, especially

326 with regards to CTLs, neutrophils and NK-cells (Supplementary Fig. S9a-d). C2
327 tumours also exhibit markedly higher neutrophil:CTL ratios (Supplementary Fig. S9e,
328 f) and neutrophil:lymphocyte (CTL, B-cell and Treg) ratios (NLR, Supplementary Fig.
329 S10); established adverse prognostic factors in cervical cancer³⁵⁻³⁷. At 0.7, the NLR
330 in C1 tumours across all cohorts was less than half that observed in C2 tumours
331 (1.85).

332

333 Validation of MethylCIBERSORT cell estimates was performed for a subset of
334 samples from the Innsbruck cohort using CD8 (CTLs) and myeloperoxidase (MPO,
335 neutrophils) immunohistochemistry (IHC)-based scores from a pathologist blinded to
336 cluster designation (Supplementary Fig. S11a-c) and for CTLs in the Oslo cohort
337 samples using comparison of MethylCIBERSORT estimates to CD8 IHC-based
338 digital pathology scores (Supplementary Fig. S11d).

339

340 Also of potential significance regarding the tumour immune microenvironment, is the
341 presence of two immune checkpoint genes, *CD276* (also known as *B7-H3*) and
342 *NT5E* (also known as *CD73*) in the set of 938 signature DEGs that separate the
343 clusters (Table S3). Both *B7-H3* and *NT5E*, along with a third immune checkpoint
344 gene (*PD-L2*) are expressed at higher levels in C2 tumours (Supplementary Fig.
345 S12) and hypomethylation of two CpGs in the *NT5E* promoter is evident in C2
346 tumours (Supplementary Table S5). All three suppress T-cell activity³⁸⁻⁴⁰ and *B7-H3*
347 expression has been linked to poor prognosis in cervical cancer^{41,42}.

348

349 **Evidence for differences in stromal fibroblast phenotype between C1 and C2**
350 **tumours**

351

352 Gene set enrichment analysis using Metascape⁴³ suggested increased EMT
353 (Supplementary Table S9) in C2 tumours, with 52 of 200 genes in in the EMT
354 Hallmark gene set upregulated. As noted above there is also greater overlap
355 between the C2 cluster and an EMT cluster defined by TCGA and based on RPPA
356 data (Figure 3). Single-cell RNA sequencing and xenografting studies strongly
357 suggest that rather than arising from the tumour cells (few of which have undergone
358 EMT at any given time^{44–46}), mesenchymal gene signatures in bulk tumour
359 expression data instead derive from stromal cells including fibroblasts, which can
360 adopt various phenotypes and play an important role in shaping the tumour immune
361 microenvironment^{47,48}. In addition to YAP1, which has been linked to the formation of
362 cancer-associated fibroblasts (CAFs)⁴⁹ (as well as EMT^{50–52} and angiogenesis⁵³), C2
363 tumours display increased expression of the CAF marker genes *FAP* and *SERPINE1*
364 (also known as *PAI-1*)⁵⁴; the latter evidenced at both mRNA and protein levels
365 (Supplementary Table S3, Fig. 5b). Overall fibroblast content as estimated by
366 MethylCIBERSORT is similar between C1 and C2 tumours (Fig. 6a, b) but given
367 recent findings regarding the extent and prognostic significance of CAF
368 heterogeneity in the tumour microenvironment^{55–59}, we hypothesized that CAF
369 phenotype rather than overall abundance, may differ between C1 and C2 tumours.
370 To examine this, hierarchical clustering was performed based on the expression of
371 eight gene sets (68 genes) curated by Qian et al⁵⁶, representing CAF-related
372 biological processes and which are differentially expressed across six CAF
373 phenotypes recently identified in a pan-cancer analysis⁵⁹. C2 tumours cluster
374 together, displaying increased expression of proinflammatory genes associated with
375 an inflammatory (pan-iCAF2) CAF phenotype, C1 tumours appear more

376 heterogenous with respect to expression of the signature genes used to define CAF
377 phenotypes; there is upregulation of assorted myofibroblastic (myoCAF) genes in a
378 subgroup C1 tumours, including various collagens, ECM genes and TGFb-
379 associated genes, as well as ‘contractile’ genes such as smooth muscle actin
380 (*ACTA2*, Fig. 6c). While *ACTA2* is commonly used to identify myoCAF, it is also
381 expressed by pericytes and smooth muscle cells, which share the contractile
382 phenotype (and express for example, *MYH11*)^{47,59,60}. Consistent with this, C2
383 tumours are 4.8x ($p = 1.78 \times 10^{-9}$, Fisher’s Exact Test) more likely to be classified as
384 ‘CAF-high’ than C1 tumours using a four-gene CAF index defined by Ko et al⁴⁸.
385 Indeed, three of the four CAF index genes (*TGFBI*, *TGFB2* and *FN1*) appear in the
386 938 DEG signature that separates C2 from C1 tumours (Supplementary Table S3).

387

388 **Discussion**

389

390 In this study we hypothesized that by drawing upon several cervical cancer cohorts
391 for which ‘omics data, clinical information and HPV typing were either available or for
392 which we were able to profile samples ourselves, we would be able to gain further
393 insight into CSCC – the most common histological cervical cancer subtype.
394 Clustering of CSCCs according to the 10% most variable genes identified two
395 clusters (C1 and C2) that bear resemblance to the keratin-high and keratin-low
396 iClusters originally defined by TCGA⁸. Cluster membership is an independent
397 predictor of 5-year OS and CSCCs can be accurately assigned to cluster using either
398 a 938 gene expression signature or a 129 MVP DNA methylation signature,
399 providing a means by which to gain prognostic information for cervical cancer
400 patients. While HPV16 and the alpha-9 clade to which it belongs have been

401 associated with longer PFS and OS in several studies^{61–66}, the relationship between
402 HPV genotype and cervical cancer prognosis remains unclear, as highlighted by a
403 recent meta-analysis⁷. In our multivariate analyses, membership of the C2 cluster
404 but not HPV type was an independent predictor of poor prognosis in both the
405 discovery and validation cohorts and remained so when only HPV16-positive
406 tumours in either cohort were considered. Possibly, the reason that HPV16 and other
407 alpha-9 HPV types have been associated with more favourable outcomes in certain
408 studies is that these viruses are more likely to cause C1-type tumours.

409

410 Adeno- and adenosquamous carcinomas, which are thought to arise from the
411 columnar epithelium of the endocervix, have been linked to poor prognosis in
412 cervical cancer^{3–6} and to avoid differences due to histology, we focused our study
413 entirely on CSCC. Interestingly, of the 14 keratin genes that are differentially
414 expressed between C1 and C2 tumours, three (*KRT7*, *KRT8* and *KRT18*) that are
415 upregulated in C2 were classified as marker genes for columnar-like tumours with a
416 possible endocervical origin in a recent study that used single cell RNA-sequencing
417 and lineage tracing experiments to explore cell-of-origin for CSCC and
418 adenocarcinoma⁶⁷. In contrast, C1 tumours display increased expression of *KRT5*, a
419 marker of the squamous-like subtype with a proposed ectocervical origin identified
420 by Chumduri et al⁶⁷ (Fig. 3). Other signature genes (*TP63*, *CERS3*, *CSTA*, *CLCA2*,
421 *DSC3* and *DSG3*) upregulated in C1 tumours are also markers of the squamous-like
422 subtype, while further columnar-like marker genes (*MUC5B* and *RGL3*) are
423 upregulated in C2 tumours (Supplementary Table S3). Squamous-like tumours are
424 significantly enriched in the C1 sub-group, a C1 tumour is 4.9x more likely to be
425 squamous-like than columnar-like or unclassified (Fisher Exact Test, $p = 0.0003$).

426 This suggests that C2 tumours, although SCCs, harbour features associated with
427 adenocarcinoma; possibly even hinting at a different cell-of-origin for C1 versus C2
428 tumours. The greater frequency with which alpha 7 HPV types are found in C2 SCCs
429 is another feature shared with adenocarcinoma.

430

431 Our analysis suggests differences in the tumour immune microenvironment between
432 C1 and C2 CSCCs, that are highly reproducible across cohorts from the USA,
433 Europe and Uganda and that might explain the differential prognosis associated with
434 these clusters. In addition to the high neutrophil:lymphocyte ratio, the increased
435 expression of cytokines including IL-6, TGF- β and G-CSF and of the chemokines
436 CXCL1-3 in C2 tumours suggests pro-tumourigenic (N2) polarisation of these
437 neutrophils⁶⁸⁻⁷³, which is typical of tumours with a high NLR⁷⁴. The observation that
438 CSCCs occurring in HIV⁺ patients from the Ugandan/CGCI cohort are much more
439 likely to be of the C2 subtype than those in HIV⁻ patients hints at a possible
440 relationship between the immune competence of the patient and the likelihood of
441 developing a C2 tumour. This requires further investigation but is consistent with
442 greater evidence of existing anti-tumour immune responses in C1 tumours.

443

444 Finally, it is interesting to note that three targetable immune checkpoint proteins (B7-
445 H3, NT5E and PD-L2) are expressed at higher levels in C2 tumours. In addition to its
446 immune suppressive effects, B7-H3 has been linked to key processes that are
447 upregulated in these tumours including EMT and angiogenesis, through the
448 activation of NF- κ B signalling and the downregulation of E-cadherin expression^{75,76}.
449 Interestingly, the expression of B7-H3 and NT5E on CAFs has been linked to poor
450 prognosis in gastric and colorectal cancer, respectively^{40,77}. Also of relevance given

451 our observation of differing CAF phenotype between clusters is the report that a CAF
452 subtype (CAF-S1) identified in breast cancer that displays high levels of B7-H3 and
453 NT5E expression is seen in tumours with low levels of CTL infiltration⁷⁸. PD1/PD-L1
454 immune checkpoint blockade (pembrolizumab) was recently FDA-approved for first-
455 line treatment of metastatic cervical cancer in combination with chemotherapy in
456 patients whose tumours express PD-L1^{79,80}, while CTLA4 blockade (Ipilimumab) has
457 also shown promising activity, both as a single agent^{81,82} and in combination with
458 PD1 blockade (Nivolumab)⁸³. Efficacy of PD1 blockade in cervical cancer has been
459 linked to the presence of a CD8+FoxP3+CD25+ T-cell subset⁸⁴ and an important
460 limitation of our study is the inability to differentiate between CD8+ T-cell
461 phenotypes. Nonetheless, identification of alternative, targetable immune checkpoint
462 molecules in C2 tumours provides a potential therapeutic strategy for a subset of
463 cervical cancers that respond poorly to chemoradiotherapy and that, given their low
464 overall levels of T-cell infiltrates, are maybe less likely to respond to PD1 blockade
465 than C1 tumours.

466

467 In conclusion, we show that CSCCs can be categorised in two novel tumour types,
468 C1 and C2, among which C1 tumours have a more favourable outcome. Although
469 HPV16 is more likely to cause C1 tumours and HPV18 C2 tumours, HPV type is not
470 an independent predictor of prognosis, suggesting it is the tumour type rather than
471 the causative HPV type that is critical for the disease outcome. Notably, the key
472 molecular and cellular characteristics of C1 and C2 tumours are consistent among
473 cohorts from the US, Europe, and Sub-Saharan Africa. This suggests that the
474 findings and underlying principle: that CSCC can develop along two trajectories

475 associated with differing clinical behaviour that can be identified using defined gene
476 expression or DNA methylation signatures, are of broad relevance.

477

478

479 **Methods**

480

481 **Patient samples**

482 All patients gave written, informed consent before inclusion. Samples from Bergen
483 were collected in a population-based setting from patients treated at the Department
484 of Obstetrics and Gynaecology, Haukeland University Hospital, Bergen, Norway,
485 from May 2001 to May 2011. The study has been approved by the regional ethical
486 committee (REK 2009/2315, 2014/1907 and 2018/591). For more details on sample
487 collection see 11,79. Samples from Innsbruck were collected and processed at the
488 Department of Obstetrics and Gynaecology of the Medical University of Innsbruck.
489 The study was reviewed and approved by the Ethics committee of the Medical
490 University of Innsbruck (reference number: AN2016-0051 360/4.3; 374/5.4: 'Biobank
491 study: Validation of a DNA-methylation based signature in cervical cancer') and
492 conducted in accordance with the Declaration of Helsinki. Samples from Oslo (n =
493 268) were collected from patients participating in a previously published prospective
494 clinical study⁸⁰ approved by the Regional Committee for Medical Research Ethics in
495 southern Norway (REK no. S-01129). Limited quantities of patient tumour samples
496 and extracted DNA may remain and the distribution of these materials is subject to
497 ethical approval at the institutions from which they were collected. Note that the
498 cases in the Oslo cohort were not treated with surgery. The samples used for
499 molecular analysis were diagnostic biopsies from the primary tumour. In all other

500 cases, specimens were from resections of the primary tumour. Those interested in
501 working with these samples should contact the authors to discuss their requirements.

502

503

504 **Dataset assembly**

505 DNA methylation (Illumina Infinium 450k array) and RNAseq data were obtained for
506 CESC from the TCGA data portal. TCGA mutation data were obtained from the MC3
507 project on SAGE Synapse (syn7214402). RNAseq data for the Uganda cohort was
508 obtained from the TCGA data portal and DNA methylation (Illumina Infinium EPIC
509 array) and mutation data from National Cancer Institute's Genome Data Commons
510 Publication Page at [https://gdc.cancer.gov/about-data/publications/CGCI-HTMCP-](https://gdc.cancer.gov/about-data/publications/CGCI-HTMCP-CC-2020)
511 [CC-2020](#). DNA methylation (Illumina Infinium 450k array) and gene expression
512 (Illumina HumanHT-12 V4.0 expression beadchip) data from the Oslo cohort were
513 obtained from the Gene Expression Omnibus (GSE68339). RNAseq data were
514 obtained for the Bergen cohort from dbGaP (phs000600/DS-CA-MDS 'Genomic
515 Sequencing of Cervical Cancers') under the authorisation of project #14589
516 "Investigating the mechanisms by which viruses and carcinogens contribute to
517 cancer development" and were converted to fastq files using SRA-dump from the
518 SRA Toolkit (<http://ncbi.github.io/sra-tools/>). Kallisto⁸¹ was then used to quantify
519 expression of GENCODE GrCh37 transcripts, rebase repeats and transcripts from
520 20 different high-risk HPV types with bias correction. Where IDAT files for 450k data
521 were available, they were parsed using *minfi*⁸² and were subjected to Functional
522 Normalisation⁸³, followed by BMIQ-correction⁸⁴ for probe type distribution (which
523 was performed for all methylation data). For TCGA samples, viral type allocation was
524 performed using VirusSeq⁸⁵.

525

526 Only squamous cell carcinomas were considered in this study to avoid confounding
527 from histology. Multidimensional visualisation of the molecular differences in
528 histology was performed using Rtsne R package with parameters available in
529 Supplementary Table S12, and the top 10% most variable genes using mean
530 absolute deviation after pre filtering of low count genes (n = 1,385). Final cohort
531 numbers and summaries are shown in Table 1.

532

533 **Generation of 450k methylation profiles**

534 100ng DNA was bisulphite converted using the EZ DNA Methylation kit (Zymo
535 Research) as per manufacturer's instructions. Bisulphite converted DNA hybridised
536 to the Infinium 450K Human Methylation array and processed in accordance with the
537 manufacturer's recommendations.

538

539 **HPV typing**

540 HPV16 or 18 was detected in 208 samples from the Oslo cohort by PCR, using the
541 primers listed in86. The PCR products were detected by polyacrylamide gene
542 electrophoresis or the Agilent DNA 1000 kit (Agilent Technologies Inc, Germany).
543 Samples from the Innsbruck cohort and the remaining non-HPV16/18 samples from
544 the Oslo cohort (n=40) were HPV-typed by DDL Diagnostic Laboratory (Netherlands)
545 using the SPF10 assay, in which a PCR-based detection of over 50 HPV types is
546 followed by a genotyping assay (LIPA₂₅) that identifies 25 HPV types (HPV 6, 11, 16,
547 18, 31, 33, 34, 35, 39, 40, 42, 43, 44, 45, 51, 52, 53, 54, 56, 58, 59, 66, 68/73, 70
548 and 74). If more than one HPV type was identified in a sample (e.g. HPV16 and

549 HPV18), that sample was designated “Other” as HPV type in the study. HPV type
550 data for the remaining samples were published previously^{8,9,11}.

551

552 **Prognostic analyses and tumour clustering**

553 Unsupervised consensus clustering was performed on TCGA SCC samples using r
554 package ConsensusClusterPlus. After prefiltering of genes to remove those with low
555 read counts (75% samples read count < 1), only the top 10% most variable genes
556 using mean absolute deviation were considered for clustering (n = 1,385). 80% of
557 tumours were sampled over 1000 iterations using all genes. PAM clustering
558 algorithm was used and clustering distance was measured using Pearson’s
559 correlation. An optimum number of clusters (K) of 2 was obtained by using the
560 proportion of ambiguously clustered pairs (PAC) using thresholds of 0.1 and 0.9 to
561 define the intermediate sub-interval. PAC was used as it accurately infers K87.
562 Limma-voom on RNAseq data and limma on BMIQ and Functionally-normalised
563 450k and EPIC data were used to identify differentially expressed genes (DEGs,
564 FDR = 0.01, FC > 2) and methylation variable positions (MVPs, FDR = 0.01, mean
565 delta-Beta > 0.25) between the 2 clusters, C1 and C2. The 116 MVPs
566 (Supplementary Table S13) common to the 450k and EPIC arrays were used to
567 allocate clusters for the Ugandan cohort. The mean delta-Beta threshold for MVPs
568 was determined as it delivered the highest concordance between DEG and MVP
569 signature cluster allocation in the Bergen cohort (89.5%) and high concordance in
570 the Ugandan cohort (91.5%). The caret R package and limma were used to develop
571 an SVM using 5 iterations of 5-fold Cross-Validation using DEGs and MVPs to
572 allocate RNAseq samples in Ugandan and Bergen cohorts, 450k samples in Bergen,
573 Innsbruck and Oslo cohorts and EPIC samples in Ugandan cohort to these

574 subgroups. Multidimensional visualisation using R package Rtsne was performed on
575 the TCGA and European cohorts with available DNA methylation data combined
576 using the 129 MVPs and parameters as shown in Supplementary Table S12.

577

578 Samples from our validation cohort, comprise of cases from three European centres
579 (Bergen and Oslo in Norway and Innsbruck, Austria) and one African centre
580 (Uganda) were binned into these categories, and were used for subsequent
581 statistical analyses to identify genomic and microenvironmental correlates. Survival
582 analyses of epigenetic allocations were carried out using Cox Proportional Hazards
583 regression with age, tumour stage, HPV type, and with surgery, radiotherapy and
584 chemotherapy (given/not given) as covariates. R packages used were survival and
585 survminer. For all clinical analyses, stages were collapsed into Stages I, II, III and IV.

586

587 RNAseq data for Bergen and Ugandan samples, Illumina HumanWG-6 v3
588 microarray data for 137 of the Oslo samples and Illumina HumanHT-12 v4
589 microarray data for 109 of the Oslo samples were used to explore cluster allocation
590 concordance accuracy between DEG and MVP signature cluster allocation. ROC
591 curve and ssGSEA analysis were performed using R (scripts available at request).

592

593 **Previous study comparison**

594 140 TCGA samples from the core set analysis (TCGA, 2017) were present in our
595 TCGA SCC cohort. Previous cluster analysis by TCGA (2017) and Chumduri *et al.*
596 (2021) was compared with our C1 and C2 cluster allocation.

597

598 **Pathway analyses**

599 Pathway and gene sets were analysed with Metascape⁴³. Settings used were
600 minimum gene set overlap of 10, p value cutoff of 0.01 and minimum enrichment of
601 1.5. All functional set, pathway, structural complex and miscellaneous gene sets
602 were included in the analysis. Only hits with an FDR of less than 0.05 were included
603 in final results.

604

605 **Mutational analyses**

606 For TCGA data, mutation calls were obtained from SAGE synapse as called by the
607 MC3 project. Mutations for the Bergen cohort were obtained from¹¹. Ugandan
608 mutation calls were obtained from National Cancer Institute's Genome Data
609 Commons Publication Page at <https://gdc.cancer.gov/about-data/publications/CGCI->
610 [HTMCP-CC-2020](https://gdc.cancer.gov/about-data/publications/CGCI-). VCFs obtained for the Ugandan cohort samples were converted
611 to maf files using R package vcf2maf, filtered for whole-exome mutations only, and
612 combined. Significantly mutated genes (SMGs) were identified using dNdScv¹⁶
613 individually for the three cohorts. Hypermutated samples (>600 mutations⁸) were
614 excluded from this analysis. A weighted approach was used to combine p values for
615 each gene for the three cohorts. R package metapro¹⁷ function wFisher was used to
616 perform this task. Genes were considered SMGs if after FDR correction of combined
617 p values, $q < 0.1$. Analysis was repeated for only C1 and C2 samples individually.
618 Two genes were removed from our list. *MUC4* was removed due to the large size of
619 the gene and *GOLGA6L18* was removed as this gene and its aliases were not
620 recognised by R package maftools⁸⁸.

621

622 R package maftools was used to produce an oncoplot for SMGs, calculating tumour
623 mutational burden for individual samples, SMG mutation frequency and mutational

624 signatures for the combined cohorts. Binomial GLMs were used to estimate
625 associations between C1 and C2 clusters and SMG mutation frequencies.

626

627 The estimated exposures of each sample to the identified mutational signatures were
628 calculated using R package mutsignatures⁸⁹ and converted to proportion of
629 signature exposure per sample.

630

631

632 **Copy number analysis**

633 450k total intensities (Methylated and Unmethylated values) were used to generate
634 copy number profiles with normal blood samples from Renius et al⁹⁰ as the germline
635 reference. Functional normalisation⁸³ was used to regress out technical variation
636 across the reference and tumour datasets before merging and quantile normalisation
637 was used to normalise combined intensities followed by Circular Binary
638 Segmentation as previously described⁹¹. Median density peak correction was
639 performed to ensure centering before further analysis. GISTIC2.0⁹² was then used
640 to identify regions of significant copy number change at both arm and gene levels.
641 Candidate copy number changes were evaluated for association with cluster using
642 binomial GLMs. The parameters chosen were a noise threshold of 0.1 with arm-level
643 peel off and a confidence level of 0.95 was used to nominate genes targeted by copy
644 number changes. Binomial regression was finally used to estimate rates of
645 differential alteration.

646

647 **Reverse Phase Protein Assay analysis**

648 Reverse Phase Protein Assay (RPPA) data for the core TCGA CESC samples were
649 obtained from the NCI GDC Legacy Archive. Differentially expressed proteins
650 between C1 and C2 clusters were determined using R package limma (FDR = 0.05,
651 FC > 1.3).

652

653 **Tumour microenvironment analyses**

654 MethylCIBERSORT³⁴ was used to estimate tumour purity and abundances of nine
655 other microenvironmental cellular fractions using TCGA and validation cohort
656 methylation beta values. Fraction numbers were then normalised by cellular
657 abundance and differences between clusters C1 and C2 were estimated using
658 Wilcoxon's rank sum test with Benjamini Hochberg correction for multiple testing.
659 This analysis was performed separately on TCGA cohort and combined validation
660 cohort, as well as on each individual cohort.

661

662 Cancer associated fibroblast associated gene set lists were obtained from Qian et
663 al⁵⁶. TCGA, Bergen and Ugandan cohort sample RNAseq data was combined and
664 visualised for these gene set genes using R package NMF93.

665

666 **CAF Index calculation**

667 For cohorts that RNAseq data was available (TCGA, Bergen and Uganda), a CAF
668 index was calculated as described in Ko et al⁴⁸. The median CAF index value was
669 used as a threshold to allocate high or low CAF in tumour samples.

670

671 **Immunohistochemistry**

672 Immunohistochemical staining of samples from the Innsbruck cohort was conducted
673 by HSL-Advanced Diagnostics (London, UK) using the Leica Bond III platform with
674 Leica Bond Polymer Refine detection as per manufacturer's recommendations.
675 Sections from a series of 17 tumour samples from the validation cohort were stained
676 for CD8 (mouse monoclonal 4B11, Leica Biosystems PA0183, used as supplied for
677 15 minutes at room temperature. HIER was performed on-board using Leica ER2
678 solution (high pH) for 20 minutes), CD68 (mouse monoclonal PGM1, Agilent
679 M087601-2, used at a dilution of 1/50 for 15mins at room temperature. HIER was
680 performed on-board using Leica ER1 solution (low pH) for 20 minutes) or MPO
681 (rabbit polyclonal, Agilent A039829-2, used at a dilution of 1/4000 for 15 minutes at
682 room temperature without epitope retrieval. Scoring was performed blinded to cluster
683 membership by a histopathologist (JM) as follows: 0 = no positive cells / field (200X
684 magnification); 1 = 1 – 10 positive cells; 2 = 11 – 100 positive cells; 3 = 101 – 200
685 positive cells; 4 = 201 = 300 positive cells; 5 = over 300 positive cells.

686

687 For the Oslo cohort, manual CD8 staining was conducted using the Dako
688 EnVision™ Flex+ System (K8012, Dako). Deparaffinization and unmasking of epitopes
689 were performed using PT-Link (Dako) and EnVision™ Flex target retrieval solution at a
690 high pH. The sections were incubated with CD8 mouse monoclonal antibody (clone
691 4B11, 1:150, 0.2 µg IgG_{2b}/ml) from Novocastra (Leica Microsystems, Newcastle Upon
692 Tyne, UK) for 45 minutes. All CD8 series included positive controls. Negative controls
693 included substitution of the monoclonal antibody with mouse myeloma protein of the
694 same subclass and concentration as the monoclonal antibody. All controls gave
695 satisfactory results. CD8 pathology scores were given to each sample (blinded to
696 cluster membership) for connective tissue only, tumour only and both as follows: 0 = no

697 positive: 1 = <10% CD8 positive cells; 2 = 10-25% CD8 positive cells; 3 = 25-50% CD8
698 positive cells; 4 = >50% CD8 positive cells. For digital quantification scanned images
699 of all sections at a high resolution of 0.46 um/pixel (20x), which was reduced to 0.92
700 um/pixel for analysis, were used. Digital score was calculated by quantifying the area
701 fraction of stained CD8 cells in relation to the entire section in the digital assessment.

702

703

704

705 **Data availability**

706 Illumina Infinium 450k array DNA methylation data generated in-house from Bergen
707 and Innsbruck validation cohort samples have been deposited in the Gene
708 Expression Omnibus (accession number GSEXXXXX (to be deposited upon
709 publication)). For detailed information on all other datasets see 'Dataset Assembly'.

710

711 **Code availability**

712 All packages used have been published, are freely available and are referenced in
713 the methods. R markdowns used to run the analyses specific to this study are
714 available from the authors on request.

715

716 **Acknowledgements**

717 AC was supported by postgraduate research scholarships from UCL and received
718 additional research support from a Debbie Fund grant to KC and TRF. TRF was
719 supported by Rosetrees Trust (M229-CD1), Cancer Research UK (A25825), the
720 Biotechnology and Biosciences Research Council (Grant Ref: BB/V010271/1), the
721 Royal Society (IEC\R2\202256) and the Global Challenges Doctoral Centre at the
722 University of Kent. DNA methylation data were generated through funding provided

723 by the Debbie Fund and the results shown here are in part based upon data
724 generated by the TCGA Research Network: <https://www.cancer.gov/tcga> and the
725 Cancer Genome Characterization Initiative: <https://ocg.cancer.gov/programs/cgci>. AF
726 was supported by grants from the MRC (MR/M025411/1), PCUK(MA-TR15-009),
727 BBSRC (BB/R009295/1), TUF, Orchid and the UCLH BRC. The authors dedicate
728 this manuscript to the late Dr Helga Salvesen, a wonderful collaborator and
729 colleague who played a key role in the project.

730

731 **References**

732

- 733 1. de Martel, C., Plummer, M., Vignat, J. & Franceschi, S. Worldwide burden of cancer
734 attributable to HPV by site, country and HPV type. *International journal of cancer*
735 **141**, 664–670 (2017).
- 736 2. Li, N., Franceschi, S., Howell-Jones, R., Snijders, P. J. F. & Clifford, G. M. Human
737 papillomavirus type distribution in 30,848 invasive cervical cancers worldwide:
738 Variation by geographical region, histological type and year of publication.
739 *International journal of cancer* **128**, 927–935 (2011).
- 740 3. Jung, E. J. *et al.* Cervical adenocarcinoma has a poorer prognosis and a higher
741 propensity for distant recurrence than squamous cell carcinoma. *International Journal*
742 *of Gynecological Cancer* **27**, 1228–1236 (2017).
- 743 4. Huang, Y. T. *et al.* Clinical behaviors and outcomes for adenocarcinoma or
744 adenosquamous carcinoma of cervix treated by radical hysterectomy and adjuvant
745 radiotherapy or chemoradiotherapy. *International Journal of Radiation Oncology*
746 *Biology Physics* **84**, 420–427 (2012).
- 747 5. Zhou, J. *et al.* Comparison of clinical outcomes of squamous cell carcinoma,
748 adenocarcinoma, and adenosquamous carcinoma of the uterine cervix after definitive
749 radiotherapy: a population-based analysis. *Journal of cancer research and clinical*
750 *oncology* **143**, 115–122 (2017).
- 751 6. Galic, V. *et al.* Prognostic significance of adenocarcinoma histology in women with
752 cervical cancer. *Gynecologic Oncology* **125**, 287–291 (2012).
- 753 7. Chen, X. *et al.* Better or Worse? The Independent Prognostic Role of HPV-16 or HPV-
754 18 Positivity in Patients With Cervical Cancer: A Meta-Analysis and Systematic
755 Review. *Frontiers in oncology* vol. 10 1733 (2020).
- 756 8. Burk, R. D. *et al.* Integrated genomic and molecular characterization of cervical
757 cancer. *Nature* **543**, 378–384 (2017).
- 758 9. Gagliardi, A. *et al.* Analysis of Ugandan cervical carcinomas identifies human
759 papillomavirus clade-specific epigenome and transcriptome landscapes. *Nature*
760 *genetics* **52**, 800–810 (2020).

- 761 10. Huang, J. *et al.* Comprehensive genomic variation profiling of cervical intraepithelial
762 neoplasia and cervical cancer identifies potential targets for cervical cancer early
763 warning. *Journal of medical genetics* **56**, 186–194 (2019).
- 764 11. Ojesina, A. I. *et al.* Landscape of genomic alterations in cervical carcinomas. *Nature*
765 **506**, 371–375 (2014).
- 766 12. Chen, J. L. *et al.* Differential clinical characteristics, treatment response and prognosis
767 of locally advanced adenocarcinoma/adenosquamous carcinoma and squamous cell
768 carcinoma of cervix treated with definitive radiotherapy. *Acta Obstetrica et*
769 *Gynecologica Scandinavica* **93**, 661–668 (2014).
- 770 13. Williams, N. L., Werner, T. L., Jarboe, E. A. & Gaffney, D. K. Adenocarcinoma of the
771 cervix: should we treat it differently? *Current oncology reports* **17**, 16–17 (2015).
- 772 14. Hu, K., Wang, W., Liu, X., Meng, Q. & Zhang, F. Comparison of treatment outcomes
773 between squamous cell carcinoma and adenocarcinoma of cervix after definitive
774 radiotherapy or concurrent chemoradiotherapy. *Radiation oncology (London, England)*
775 **13**, 245–249 (2018).
- 776 15. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature*
777 **500**, 415–421 (2013).
- 778 16. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues.
779 *Cell* **171**, 1029-1041.e21 (2017).
- 780 17. Yoon, S., Baik, B., Park, T. & Nam, D. Powerful p-value combination methods to
781 detect incomplete association. *Scientific Reports* **11**, 6980 (2021).
- 782 18. Lawrence, M. S. *et al.* Comprehensive genomic characterization of head and neck
783 squamous cell carcinomas. *Nature* **517**, 576–582 (2015).
- 784 19. Lin, D.-C. *et al.* Genomic and molecular characterization of esophageal squamous cell
785 carcinoma. *Nature genetics* **46**, 467–473 (2014).
- 786 20. Hammerman, P. S. *et al.* Comprehensive genomic characterization of squamous cell
787 lung cancers. *Nature* **489**, 519–525 (2012).
- 788 21. Rauskolb, C., Sun, S., Sun, G., Pan, Y. & Irvine, K. D. Cytoskeletal tension inhibits
789 Hippo signaling through an Ajuba-Warts complex. *Cell* **158**, 143–156 (2014).
- 790 22. Nguyen, T. H., Ralbovska, A. & Kugler, J.-M. RhoBTB Proteins Regulate the Hippo
791 Pathway by Antagonizing Ubiquitination of LKB1. *G3 (Bethesda, Md.)* **10**, 1319–
792 1325 (2020).
- 793 23. Mohseni, M. *et al.* A genetic screen identifies an LKB1-MARK signalling axis
794 controlling the Hippo-YAP pathway. *Nature cell biology* **16**, 108–117 (2014).
- 795 24. Martin, D. *et al.* Assembly and activation of the Hippo signalome by FAT1 tumor
796 suppressor. *Nature communications* **9**, 2372 (2018).
- 797 25. Sourbier, C. *et al.* Targeting loss of the Hippo signaling pathway in NF2-deficient
798 papillary kidney cancers. *Oncotarget* **9**, 10723–10733 (2018).
- 799 26. Petrilli, A. M. & Fernández-Valle, C. Role of Merlin/NF2 inactivation in tumor
800 biology. *Oncogene* **35**, 537–548 (2016).
- 801 27. White, S. M. *et al.* YAP/TAZ Inhibition Induces Metabolic and Signaling Rewiring
802 Resulting in Targetable Vulnerabilities in NF2-Deficient Tumor Cells. *Developmental*
803 *cell* **49**, 425-443.e9 (2019).
- 804 28. Yang, H. *et al.* NF2 and Canonical Hippo-YAP Pathway Define Distinct Tumor
805 Subsets Characterized by Different Immune Deficiency and Treatment Implications in
806 Human Pleural Mesothelioma. *Cancers* **13**, 1561. doi: 10.3390/cancers13071561
807 (2021).
- 808 29. Gonzalez-Avila, G. *et al.* Matrix metalloproteinases participation in the metastatic
809 process and their diagnostic and therapeutic applications in cancer. *Critical reviews in*
810 *oncology/hematology* **137**, 57–83 (2019).

- 811 30. Wang, Y. *et al.* Comprehensive Molecular Characterization of the Hippo Signaling
812 Pathway in Cancer. *Cell reports* **25**, 1304-1317.e5 (2018).
- 813 31. Gooden, M. J., de Bock, G. H., Leffers, N., Daemen, T. & Nijman, H. W. The
814 prognostic influence of tumour-infiltrating lymphocytes in cancer: a systematic review
815 with meta-analysis. *British journal of cancer* **105**, 93–103 (2011).
- 816 32. Jordanova, E. S. *et al.* Human leukocyte antigen class I, MHC class I chain-related
817 molecule A, and CD8+/regulatory T-cell ratio: which variable determines survival of
818 cervical cancer patients? *Clinical cancer research* □: *an official journal of the*
819 *American Association for Cancer Research* **14**, 2028–2035 (2008).
- 820 33. Nedergaard, B. S., Ladekarl, M., Thomsen, H. F., Nyengaard, J. R. & Nielsen, K. Low
821 density of CD3+, CD4+ and CD8+ cells is associated with increased risk of relapse in
822 squamous cell cervical cancer. *British journal of cancer* **97**, 1135–1138 (2007).
- 823 34. Chakravarthy, A. *et al.* Pan-cancer deconvolution of tumour composition using DNA
824 methylation. *Nature communications* **9**, 3220–3221 (2018).
- 825 35. Mizunuma, M. *et al.* The pretreatment neutrophil-to-lymphocyte ratio predicts
826 therapeutic response to radiation therapy and concurrent chemoradiation therapy in
827 uterine cervical cancer. *International journal of clinical oncology* **20**, 989–996 (2015).
- 828 36. Lee, Y. Y. *et al.* Pretreatment neutrophil:lymphocyte ratio as a prognostic factor in
829 cervical carcinoma. *Anticancer Research* **32**, 1555–1561 (2012).
- 830 37. Huang, Q. T. *et al.* Prognostic significance of neutrophil-to-lymphocyte ratio in
831 cervical cancer: A systematic review and meta-analysis of observational studies.
832 *Oncotarget* **8**, 16755–16764 (2017).
- 833 38. Prasad, D. V. R. *et al.* Murine B7-H3 Is a Negative Regulator of T Cells. *The Journal*
834 *of Immunology* **173**, 2500 (2004).
- 835 39. Zang, X. *et al.* B7x: a widely expressed B7 family member that inhibits T cell
836 activation. *Proceedings of the National Academy of Sciences of the United States of*
837 *America* **100**, 10388–10392 (2003).
- 838 40. Yu, M. *et al.* CD73 on cancer-associated fibroblasts enhanced by the A2B-mediated
839 feedforward circuit enforces an immune checkpoint. *Nature Communications* **11**, 515
840 (2020).
- 841 41. Han, S. *et al.* Roles of B7-H3 in Cervical Cancer and Its Prognostic Value. *Journal of*
842 *Cancer* **9**, 2612–2624 (2018).
- 843 42. Huang, C. *et al.* B7-H3, B7-H4, Foxp3 and IL-2 expression in cervical cancer:
844 Associations with patient outcome and clinical significance. *Oncology reports* **35**,
845 2183–2190 (2016).
- 846 43. Zhou, Y. *et al.* Metascape provides a biologist-oriented resource for the analysis of
847 systems-level datasets. *Nature communications* **10**, 1523–1526 (2019).
- 848 44. Ruscetti, M., Quach, B., Dadashian, E. L., Mulholland, D. J. & Wu, H. Tracking and
849 Functional Characterization of Epithelial-Mesenchymal Transition and Mesenchymal
850 Tumor Cells during Prostate Cancer Metastasis. *Cancer research* **75**, 2749–2759
851 (2015).
- 852 45. Fischer, K. R. *et al.* Epithelial-to-mesenchymal transition is not required for lung
853 metastasis but contributes to chemoresistance. *Nature* **527**, 472–476 (2015).
- 854 46. Zheng, X. *et al.* Epithelial-to-mesenchymal transition is dispensable for metastasis but
855 induces chemoresistance in pancreatic cancer. *Nature* **527**, 525–530 (2015).
- 856 47. Puram, S. v *et al.* Single-Cell Transcriptomic Analysis of Primary and Metastatic
857 Tumor Ecosystems in Head and Neck Cancer. *Cell* **171**, 1611-1624.e24 (2017).
- 858 48. Ko, Y.-C. *et al.* Index of Cancer-Associated Fibroblasts Is Superior to the Epithelial-
859 Mesenchymal Transition Score in Prognosis Prediction. *Cancers* **12**, 1718 (2020).

- 860 49. Shen, T. *et al.* YAP1 plays a key role of the conversion of normal fibroblasts into
861 cancer-associated fibroblasts that contribute to prostate cancer progression. *Journal of*
862 *Experimental & Clinical Cancer Research* **39**, 36 (2020).
- 863 50. Zanonato, F. *et al.* Genome-wide association between YAP/TAZ/TEAD and AP-1 at
864 enhancers drives oncogenic growth. *Nature cell biology* **17**, 1218–1227 (2015).
- 865 51. Shao, D. D. *et al.* KRAS and YAP1 converge to regulate EMT and tumor survival.
866 *Cell* **158**, 171–184 (2014).
- 867 52. Schlegelmilch, K. *et al.* Yap1 acts downstream of α -catenin to control epidermal
868 proliferation. *Cell* **144**, 782–795 (2011).
- 869 53. Kim, J. *et al.* YAP/TAZ regulates sprouting angiogenesis and vascular barrier
870 maturation. *The Journal of clinical investigation* **127**, 3441–3461 (2017).
- 871 54. Sakamoto, H. *et al.* PAI-1 derived from cancer-associated fibroblasts in esophageal
872 squamous cell carcinoma promotes the invasion of cancer cells and the migration of
873 macrophages. *Laboratory Investigation* **101**, 353–368 (2021).
- 874 55. Neuzillet, C. *et al.* Inter- and intra-tumoural heterogeneity in cancer-associated
875 fibroblasts of human pancreatic ductal adenocarcinoma. *The Journal of pathology* **248**,
876 51–65 (2019).
- 877 56. Qian, J. *et al.* A pan-cancer blueprint of the heterogeneous tumor microenvironment
878 revealed by single-cell profiling. *Cell research* **30**, 745–762 (2020).
- 879 57. Mhaidly, R. & Mechta-Grigoriou, F. Fibroblast heterogeneity in tumor micro-
880 environment: Role in immunosuppression and new therapies. *Seminars in Immunology*
881 **48**, 101417 (2020).
- 882 58. Hutton, C. *et al.* Single-cell analysis defines a pancreatic fibroblast lineage that
883 supports anti-tumor immunity. *Cancer cell* **39**, 1227–1244.e20 (2021).
- 884 59. Galbo, P. M., Zang, X. & Zheng, D. Molecular Features of Cancer-associated
885 Fibroblast Subtypes and their Implication on Cancer Pathogenesis, Prognosis, and
886 Immunotherapy Resistance. *Clinical Cancer Research* **27**, 2636 (2021).
- 887 60. Chen, Z. *et al.* Single-cell RNA sequencing highlights the role of inflammatory cancer-
888 associated fibroblasts in bladder urothelial carcinoma. *Nature Communications* **11**,
889 5077 (2020).
- 890 61. Rader, J. S. *et al.* Genetic variations in human papillomavirus and cervical cancer
891 outcomes. *International journal of cancer* **144**, 2206–2214 (2019).
- 892 62. Wright, J. D. *et al.* Human papillomavirus type and tobacco use as predictors of
893 survival in early stage cervical carcinoma. *Gynecologic oncology* **98**, 84–91 (2005).
- 894 63. Yang, S. H., Kong, S. K., Lee, S. H., Lim, S. Y. & Park, C. Y. Human papillomavirus
895 18 as a poor prognostic factor in stage I-IIA cervical cancer following primary surgical
896 treatment. *Obstetrics & gynecology science* **57**, 492–500 (2014).
- 897 64. Burger, R. A. *et al.* Human papillomavirus type 18: association with poor prognosis in
898 early stage cervical cancer. *Journal of the National Cancer Institute* **88**, 1361–1368
899 (1996).
- 900 65. Schwartz, S. M. *et al.* Human papillomavirus and prognosis of invasive cervical
901 cancer: a population-based study. *Journal of clinical oncology* □: *official journal of the*
902 *American Society of Clinical Oncology* **19**, 1906–1915 (2001).
- 903 66. Hang, D. *et al.* Independent prognostic role of human papillomavirus genotype in
904 cervical cancer. *BMC Infectious Diseases* **17**, 391 (2017).
- 905 67. Chumduri, C. *et al.* Opposing Wnt signals regulate cervical squamocolumnar
906 homeostasis and emergence of metaplasia. *Nature cell biology* **23**, 184–197 (2021).
- 907 68. Fridlender, Z. G. *et al.* Polarization of tumor-associated neutrophil phenotype by TGF-
908 beta: “N1” versus “N2” TAN. *Cancer cell* **16**, 183–194 (2009).

- 909 69. Zhu, Q. *et al.* The IL-6–STAT3 axis mediates a reciprocal crosstalk between cancer-
910 derived mesenchymal stem cells and neutrophils to synergistically prompt gastric
911 cancer progression. *Cell Death & Disease* **5**, e1295–e1295 (2014).
- 912 70. Ohms, M., Möller, S. & Laskay, T. An Attempt to Polarize Human Neutrophils
913 Toward N1 and N2 Phenotypes in vitro. *Frontiers in immunology* **11**, 532 (2020).
- 914 71. SenGupta, S. *et al.* Triple-Negative Breast Cancer Cells Recruit Neutrophils by
915 Secreting TGF- β and CXCR2 Ligands. *Frontiers in immunology* **12**, 659996 (2021).
- 916 72. Casbon, A. J. *et al.* Invasive breast cancer reprograms early myeloid differentiation in
917 the bone marrow to generate immunosuppressive neutrophils. *Proceedings of the*
918 *National Academy of Sciences of the United States of America* **112**, E566-75 (2015).
- 919 73. Shaul, M. E. *et al.* Tumor-associated neutrophils display a distinct N1 profile
920 following TGF β modulation: A transcriptomics analysis of pro- vs. antitumor TANs.
921 *Oncoimmunology* **5**, e1232221 (2016).
- 922 74. Kim, Y., Lee, D., Lee, J., Lee, S. & Lawler, S. Role of tumor-associated neutrophils in
923 regulation of tumor growth in lung cancer development: A mathematical model. *PloS*
924 *one* **14**, e0211041 (2019).
- 925 75. Xie, C. *et al.* Soluble B7-H3 promotes the invasion and metastasis of pancreatic
926 carcinoma cells through the TLR4/NF- κ B pathway. *Scientific reports* **6**, 27528 (2016).
- 927 76. MacGregor, H. L. *et al.* High expression of B7-H3 on stromal cells defines tumor and
928 stromal compartments in epithelial ovarian cancer and is associated with limited
929 immune activation. *Journal for ImmunoTherapy of Cancer* **7**, 357 (2019).
- 930 77. Zhan, S. *et al.* Overexpression of B7-H3 in α -SMA-Positive Fibroblasts Is Associated
931 With Cancer Progression and Survival in Gastric Adenocarcinomas. *Frontiers in*
932 *Oncology* **9**, 1466 (2020).
- 933 78. Costa, A. *et al.* Fibroblast Heterogeneity and Immunosuppressive Environment in
934 Human Breast Cancer. *Cancer cell* **33**, 463-479.e10 (2018).
- 935 79. Chung, H. C. *et al.* Efficacy and Safety of Pembrolizumab in Previously Treated
936 Advanced Cervical Cancer: Results From the Phase II KEYNOTE-158 Study. *Journal*
937 *of Clinical Oncology* **37**, 1470–1478 (2019).
- 938 80. Colombo, N. *et al.* Pembrolizumab for Persistent, Recurrent, or Metastatic Cervical
939 Cancer. *New England Journal of Medicine* (2021) doi:10.1056/nejmoa2112435.
- 940 81. Lheureux, S. *et al.* Association of Ipilimumab With Safety and Antitumor Activity in
941 Women With Metastatic or Recurrent Human Papillomavirus-Related Cervical
942 Carcinoma. *JAMA oncology* **4**, e173776–e173776 (2018).
- 943 82. Mayadev, J. S. *et al.* Sequential Ipilimumab After Chemoradiotherapy in Curative-
944 Intent Treatment of Patients With Node-Positive Cervical Cancer. *JAMA oncology* **6**,
945 92–99 (2020).
- 946 83. Naumann, R. W. *et al.* Efficacy and safety of nivolumab (Nivo) + ipilimumab (Ipi) in
947 patients (pts) with recurrent/metastatic (R/M) cervical cancer: Results from
948 CheckMate 358. in *Annals of Oncology* v898–v899 (2019).
949 doi:10.1093/annonc/mdz394.
- 950 84. Heeren, A. M. *et al.* Efficacy of PD-1 blockade in cervical cancer is related to a
951 CD8+FoxP3+CD25+ T-cell subset with operational effector functions despite high
952 immune checkpoint levels. *Journal for ImmunoTherapy of Cancer* **7**, 43 (2019).
- 953 85. Halle, M. K. *et al.* Clinicopathologic and molecular markers in cervical carcinoma: a
954 prospective cohort study. *American journal of obstetrics and gynecology* **217**, 432.e1-
955 432.e17 (2017).
- 956 86. Lando, M. *et al.* Identification of eight candidate target genes of the recurrent 3p12-
957 p14 loss in cervical cancer by integrative genomic profiling. *The Journal of pathology*
958 **230**, 59–69 (2013).

- 959 87. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-
960 seq quantification. *Nature Biotechnology* **34**, 525–527 (2016).
- 961 88. Aryee, M. J. *et al.* Minfi: a flexible and comprehensive Bioconductor package for the
962 analysis of Infinium DNA methylation microarrays. *Bioinformatics (Oxford, England)*
963 **30**, 1363–1369 (2014).
- 964 89. Fortin, J.-P. *et al.* Functional normalization of 450k methylation array data improves
965 replication in large cancer studies. *Genome biology* **15**, 503 (2014).
- 966 90. Teschendorff, A. E. *et al.* A beta-mixture quantile normalization method for correcting
967 probe design bias in Illumina Infinium 450 k DNA methylation data. *Bioinformatics*
968 *(Oxford, England)* **29**, 189–196 (2013).
- 969 91. Chen, Y. *et al.* VirusSeq: software to identify viruses and their integration sites using
970 next-generation sequencing of human cancer tissue. *Bioinformatics (Oxford, England)*
971 **29**, 266–267 (2013).
- 972 92. Lyng, H. *et al.* Intratumor chromosomal heterogeneity in advanced carcinomas of the
973 uterine cervix. *International Journal of Cancer* **111**, 358–366 (2004).
- 974 93. \square enbabaoğlu, Y., Michailidis, G. & Li, J. Z. Critical limitations of consensus
975 clustering in class discovery. *Scientific reports* **4**, 6207 (2014).
- 976 94. Mayakonda, A., Lin, D.-C., Assenov, Y., Plass, C. & Koeffler, H. P. Maftools:
977 efficient and comprehensive analysis of somatic variants in cancer. *Genome research*
978 **28**, 1747–1756 (2018).
- 979 95. Fantini, D., Vidimar, V., Yu, Y., Condello, S. & Meeks, J. J. MutSignatures: an R
980 package for extraction and analysis of cancer mutational signatures. *Scientific Reports*
981 **10**, 18217 (2020).
- 982 96. Reinius, L. E. *et al.* Differential DNA methylation in purified human blood cells:
983 implications for cell lineage and studies on disease susceptibility. *PloS one* **7**, e41361–
984 e41361 (2012).
- 985 97. Feber, A. *et al.* Using high-density DNA methylation arrays to profile copy number
986 alterations. *Genome biology* **15**, R30–R30 (2014).
- 987 98. Mermel, C. H. *et al.* GISTIC2.0 facilitates sensitive and confident localization of the
988 targets of focal somatic copy-number alteration in human cancers. *Genome Biology* **12**,
989 R41 (2011).
- 990 99. Gaujoux, R. & Seoighe, C. A flexible R package for nonnegative matrix factorization.
991 *BMC Bioinformatics* **11**, 367 (2010).
- 992

993

994 **Figure Legends**

995

996 **Figure 1 Consensus clustering produces two prognostic clusters in TCGA**

997 **SCC cohort. a)** Consensus clustering of 236 TCGA HPV+ SCC patients. **b)** There

998 were 938 differentially expressed genes between the two clusters. **c)** 5 year survival

999 between the 2 SCC subgroups. **d)** 5 year survival between the 2 SCC subgroups

1000 considering only HPV16+ tumours. Statistics from univariate Cox regression.

1001

1002 **Figure 2 Cluster allocation of validation cohorts using methylation signature.**

1003 **a)** A signature of DNA methylation ($dB > 0.25$, $FDR < 0.01$) separates C1 and C2
1004 SCC subgroups in the TCGA cohort. **b)** The methylation patterns are reproduced in
1005 a validation dataset from three European centres ($n = 313$). **c)** C2 tumours from
1006 TCGA and European validation cohorts cluster together based on the 129 MVP
1007 signature. **d)** 5 year survival curve for combined European validation cohorts.
1008 Statistics from univariate Cox regression.

1009

1010 **Figure 3 Comparison of SCC subgroups with previous studies.** Cluster analysis
1011 had previously been performed on 140 TCGA SCC tumours in two studies – one
1012 determined clusters based on cell of origin markers (Chumduri *et al*, 2021, red), one
1013 determined clusters based on integrated omics data (TCGA Network, 2017, orange).
1014 The heatmap at the bottom of plot represents expression levels of cytokeratin genes
1015 present in our C2 gene signature.

1016

1017 **Figure 4 Genomic summary of significantly mutated genes (SMGs) in SCC**
1018 **cohorts.** Main plot shows mutation type and frequencies for 34 SMGs identified
1019 using dNdSCV on TCGA, Bergen and Ugandan cohorts (367 total patients). Grey
1020 bars at top of plot represent TMB per sample. Grey bars to left of plot represent
1021 significance of SMG, larger bar is more significant. Barchart to the right shows
1022 proportion of a genes mutations in by cluster (blue = C1, red = C2). Black box
1023 around bar represents a significant difference in mutation frequency between the
1024 clusters ($p < 0.05$) while a gold box means no significant difference between the

1025 clusters. The plot at the bottom of figure represents the mutational signatures that
1026 contribute towards each individuals tumour mutational burden.

1027 [Gene name key – blue – unique to C1 analysis, red = unique to C2 analysis, black =
1028 both in C1 and C2 individual analyses, black* = only significant when combining both
1029 clusters for analysis, † = novel SMG in cervical cancer, ‡ = not significant in
1030 combined cluster analysis but significant in C1 only analysis].

1031

1032 **Figure 5 Copy number and protein level differences between SCC subgroups.**

1033 **a)** Volcano plot showing differences in GISTIC copy number peak frequencies
1034 between C1 and C2 tumours, with $-\log_{10}(\text{FDR})$ on the y axis and the odds ratio on
1035 the x axis. **b)** Volcano plot showing differentially abundant proteins and phospho-
1036 proteins (FDR < 0.05, FC > 1.3, represented by yellow dots) between C1 and C2
1037 TCGA tumours, as measured by Reverse Phase Protein Array.

1038

1039 **Figure 6 Differences in the tumour microenvironment between cervical cancer**
1040 **subgroups.** Plot showing median abundances (x-axis) and median differences (%,
1041 y-axis) for different cell types estimated using MethylCIBERSORT, with significant
1042 differences in orange, for **a)** TCGA discovery cohort and **b)** combined validation
1043 cohorts. **c)** C2 tumours cluster together using CAF geneset genes.

1044

1045 **Supplementary Figure Legends**

1046

1047 **Supplementary Figure S1 – tSNE clustering by histology in cervical cancer**
1048 **cohorts.** Unsupervised tSNE analysis using top 10% most variable genes for
1049 cervical cancer cohorts **a)** TCGA (1385 most variable genes), **b)** Ugandan (1371)

1050 and **c)** Bergen (1430). Concordance of most variable genes was high amongst the 3
1051 cohorts **(d)**.

1052

1053 **Supplementary Figure S2 Consensus clustering using ConsensusClusterPlus.**

1054 **a)** Consensus CDF plot. PAC score = CDF at 0.9 consensus index – CDF at 0.1
1055 consensus index for each curve. **b)** Delta area plot used in decision of optimum
1056 number of clusters.

1057

1058 **Supplementary Figure S3 Genes that are both differentially expressed and**

1059 **differentially methylated between C1 and C2 subgroups.** Datapoints represent
1060 methylated variable positions (in either the 3'UTR, body of gene, intergenic region or
1061 gene promoter) in genes that are also differentially expressed between C1 and C2
1062 subgroups. Datapoints in the top left quadrant are MVPs that are hypomethylated in
1063 genes that are also upregulated in C2 tumours. Those in the bottom right quadrant
1064 are hypermethylated in genes that are downregulated in C2 tumours.

1065

1066 **Supplementary Figure S4 Concordance between gene expression and DNA**

1067 **methylation-derived cluster membership. a)** The percentage of samples that are
1068 designated the same cluster allocation by gene expression signature and
1069 methylation signatures based on varying delta Beta thresholds. **b)** ROC curves
1070 showing the accuracy with which C1 or C2 cluster membership can be predicted
1071 using DNA methylation differences (MVPs) in samples from the validation cohorts for
1072 which either RNA-seq (Bergen, n=37, and Uganda, n=94, HPV+ SCC cases),
1073 Illumina HumanHT-12 V4.0 expression beadchip array (Oslo SCC cases, n=109) or
1074 Illumina HumanWG-6 v3.0 expression beadchip array (Oslo SCC cases, n=139)

1075 gene expression data were available. **c)** Single sample gene set enrichment analysis
1076 (ssGSEA) for validation cohorts used in panel B. The y-axis represents the ssGSEA
1077 score for each sample, compared with the genes from the C2 gene expression
1078 signature. P-values from Wilcoxon rank-sum test.

1079

1080 **Supplementary Figure S5 Validation SCC cohorts. a)** Ugandan validation cohort
1081 clustering based on 116 MVP signature. Kaplan-meier curves for **b)** HPV16+
1082 European validation cohort SCC patients; **c)** European validation cohort SCC
1083 patients without chemotherapy treatment and **d)** 5 year survival for the 5 individual
1084 cohorts in this study.

1085

1086 **Supplementary Figure S6 Elevation of epithelial mesenchymal transition (EMT)**
1087 **score is evident in C2 tumours. a)** EMT score derived by TCGA for 140 HPV+
1088 squamous TCGA cervical cancer tumours in our study. EMT score is higher in C2
1089 tumours.

1090

1091 **Supplementary Figure S7 Mutational signatures of combined HPV+ squamous**
1092 **cervical cancer cohorts.** COSMIC mutational signatures identified in combined
1093 HPV+ squamous cervical cancer cohort including genomic data from TCGA, Bergen
1094 and Ugandan cohorts.

1095

1096

1097 **Supplementary Figure S8 Increased levels of YAP in tumours with YAP1**
1098 **amplification.** YAP1 expression **(a)**, and YAP protein levels **(b)** unphosphorylated,
1099 **c)** phosphorylated) are higher in tumours that contain YAP1 amplifications.

1100

1101 **Supplementary Figure S9 Differences in immune microenvironment between**
1102 **SCC subgroups in individual cohorts.** Median abundances (x-axis) and median
1103 differences (% , y-axis) for different cell types estimated using MethylCIBERSORT,
1104 with significant differences in orange for cohorts from **a)** Bergen, **b)** Innsbruck, **c)**
1105 Oslo and **d)** Uganda. C2 tumours display increased neutrophil:CTL ratios as
1106 estimated using MethylCIBERSORT for **e)** TCGA discovery cohort and **f)** combined
1107 validation cohorts.

1108

1109 **Supplementary Figure S10 Immune cell ratios by cluster using**
1110 **MethylCIBERSORT estimates.** **a)** Neutrophil:CD19 estimate ratios for combined
1111 cohorts. **b)** Neutrophil:Treg estimate ratios for combined cohorts.

1112

1113 **Supplementary Figure S11 Comparison of MethylCIBERSORT estimates and**
1114 **immunohistochemistry(IHC)-based scoring.** Correlations between
1115 MethylCIBERSORT estimates and IHC-based scoring for **a)** CD8+ T-cells, **b)**
1116 neutrophils (MPO+), **c)** CD8+ T-cell:neutrophil ratio in 14 SCCs from the Innsbruck
1117 validation cohort and **d)** CD8+ T-cells for 229 SCCs from the Oslo validation cohort.
1118 Trendlines are derived from linear modelling, shaded areas represent 95% CI of
1119 trendlines.

1120

1121 **Supplementary Figure S12 Upregulation of immune checkpoint genes in C2**
1122 **SCCs.** Upregulation of **a)** *B7-H3* (*CD276*), **b)** *NT5E* (*CD73*) and **c)** *PD-L2*
1123 (*PDCD1LG2*) was observed in poor prognosis C2 tumours. Analysis performed with
1124 RNA-seq data from TCGA, Bergen and Ugandan cohorts.

1125

1126 **Tables**

1127

1128 Table 1: Summary of clinicopathological characteristics for five cervical cancer
1129 cohorts.

1130

1131 Table 2 – Five-year survival analysis for all cohorts

1132

1133 **Supplementary Tables**

1134

1135 Table S1 - Clinical and pathologic characteristics of TCGA squamous cervical
1136 cancer cohort samples

1137

1138 Table S2 - Top 10% most variable genes in TCGA squamous cervical cancer cohort

1139

1140 Table S3 - 938 Differentially expressed genes between TCGA squamous cervical
1141 cancer clusters C1 and C2

1142

1143 Table S4 - 5 year survival uni- and multivariate analysis for HPV16+ patients in
1144 squamous cervical cancer cohorts

1145

1146 Table S5 - 129 MVP signature probes (European validation cohorts)

1147

1148 Table S6 - Combined validation cohort cluster allocation

1149

1150 Table S7 - Breakdown of tumour stage in C1 and C2 cluster by percentage

1151

1152 Table S8 - Clusters and EMT scores for TCGA squamous cervical cancer samples

1153

1154 Table S9 - Significantly mutated genes using dNdSCV analysis and combining
1155 cohorts

1156

1157 Table S10 - Mutation frequency in SMGs observed in previous studies

1158

1159 Table 11 - Gene set enrichment analysis of C2 gene expression signature genes
1160 using Metascape

1161

1162 Table S12 - Parameters for TSNE multidimensional visualisation analyses

1163

1164 Table S13 - 116 MVP signature probes (Ugandan validation cohort)

1165

1166

1167

1168

1169

1170

1171

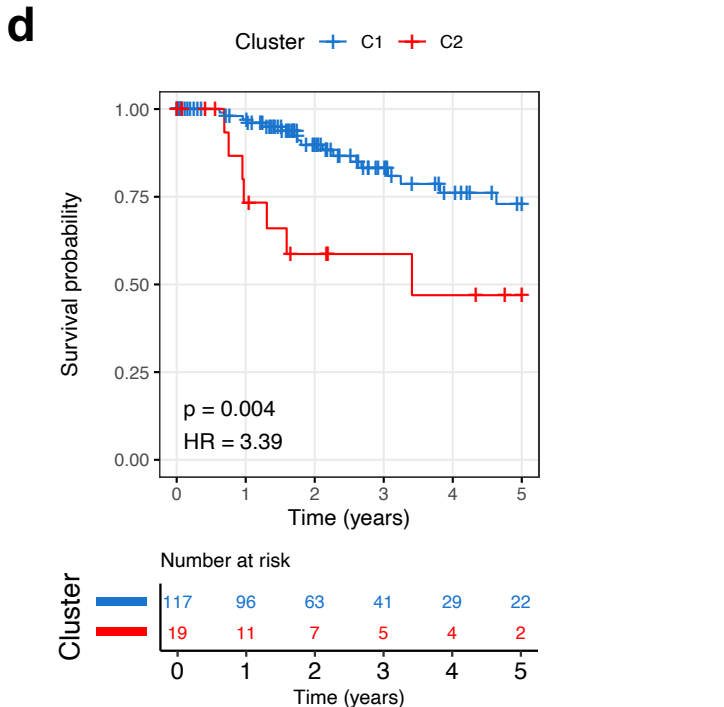
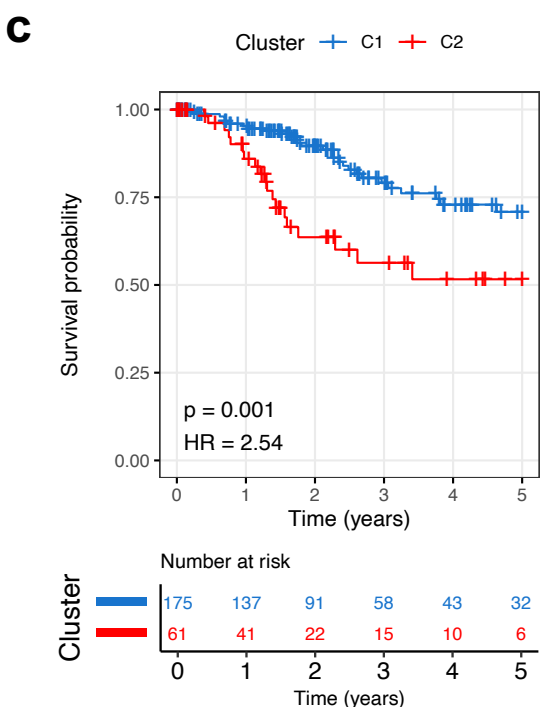
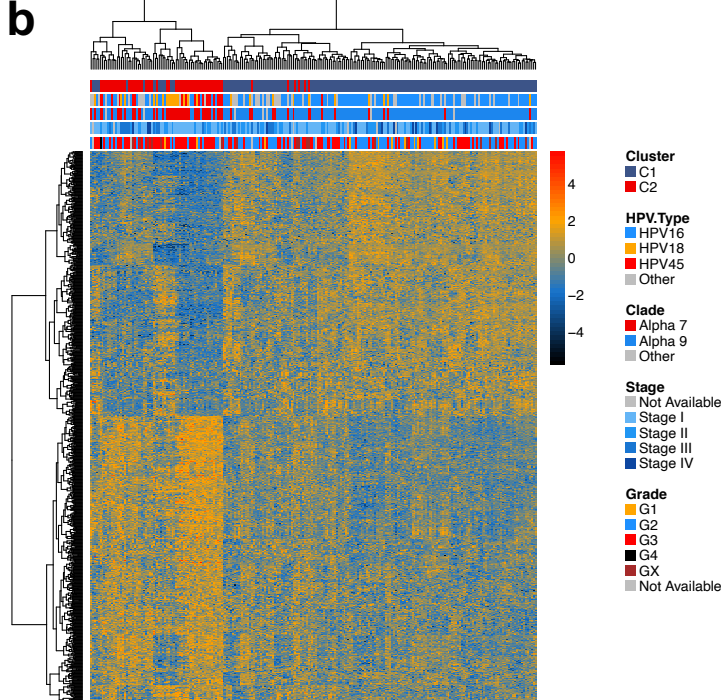
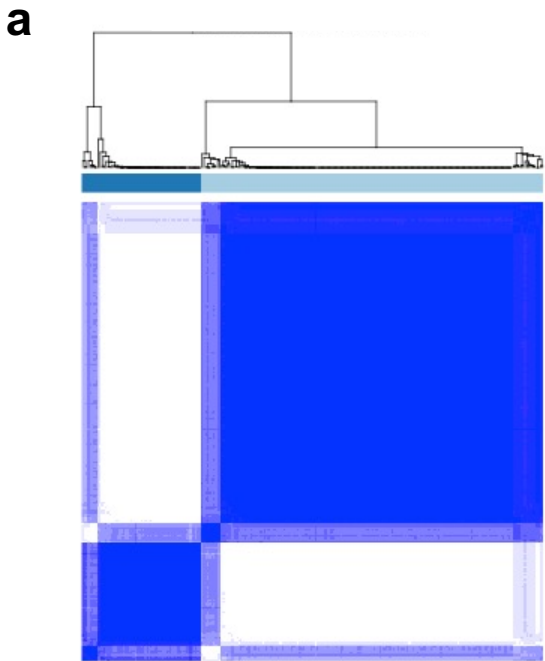


Figure 1 Consensus clustering produces two prognostic clusters in TCGA SCC cohort. a) Consensus clustering of 236 TCGA HPV+ SCC patients. **b)** There were 938 differentially expressed genes between the two clusters. **c)** 5 year survival between the 2 SCC subgroups. **d)** 5 year survival between the 2 SCC subgroups considering only HPV16+ tumours. Statistics from univariate Cox regression.

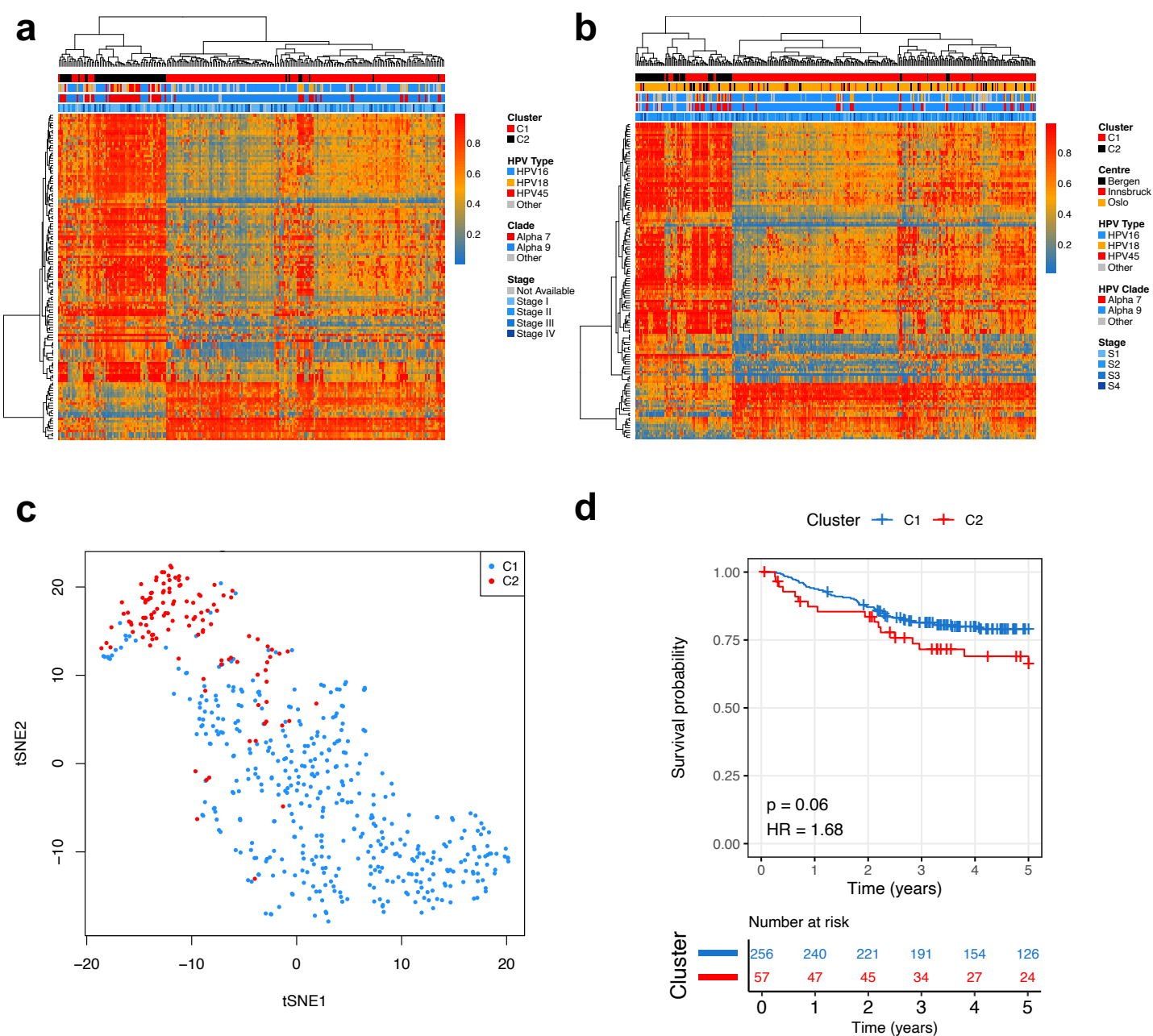


Figure 2 Cluster allocation of validation cohorts using methylation signature. **a)** A signature of DNA methylation ($dB > 0.25$, $FDR < 0.01$) separates C1 and C2 SCC subgroups in the TCGA cohort. **b)** The methylation patterns are reproduced in a validation dataset from three European centres ($n = 313$). **c)** C2 tumours from TCGA and European validation cohorts cluster together based on the 129 MVP signature. **d)** 5 year survival curve for combined European validation cohorts. Statistics from univariate Cox regression.

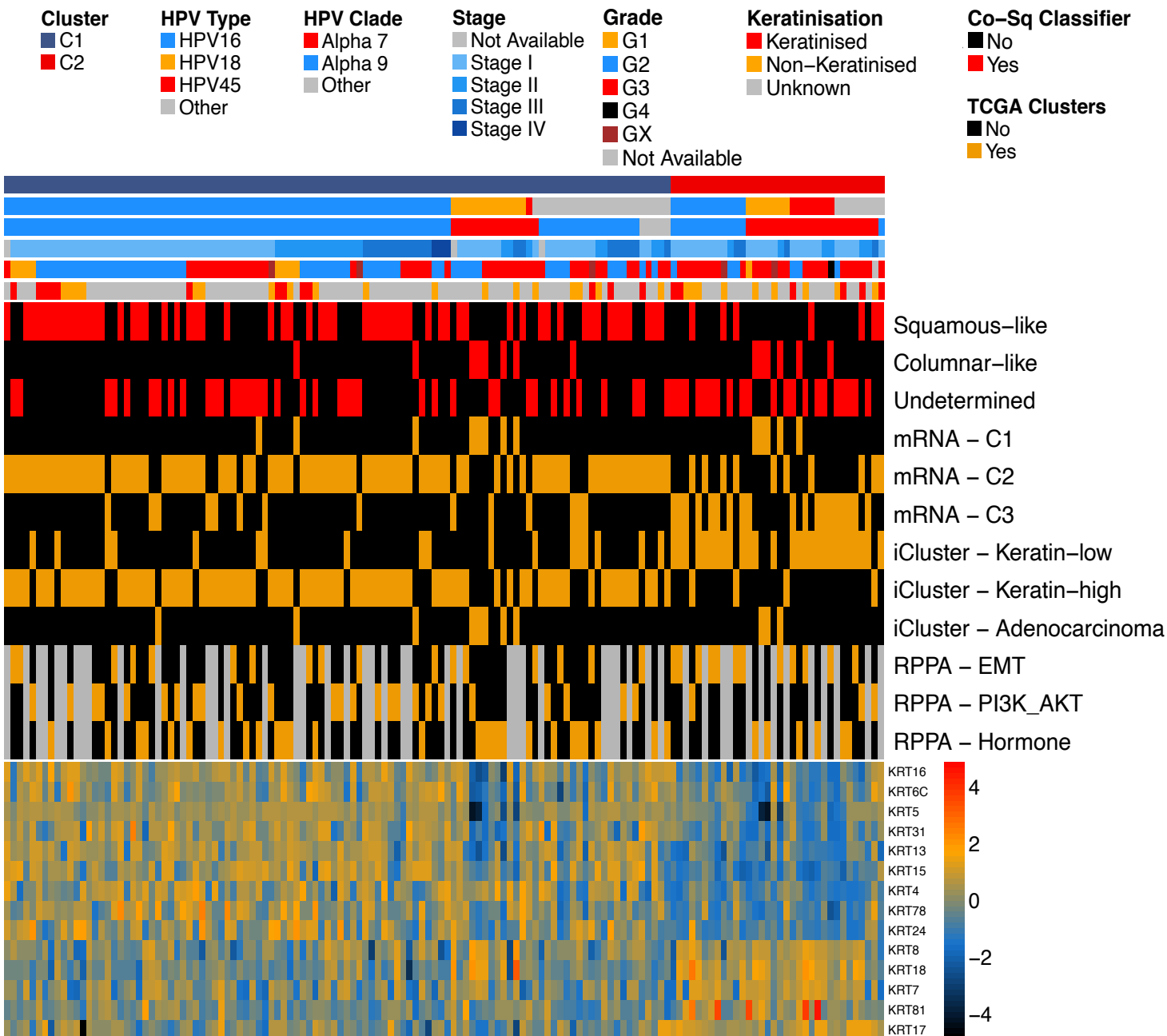


Figure 3 Comparison of SCC subgroups with previous studies. Cluster analysis had previously been performed on 140 TCGA SCC tumours in two studies – one determined clusters based on cell of origin markers (Chumduri *et al*, 2021, red), one determined clusters based on integrated omics data (TCGA Network, 2017, orange). The heatmap at the bottom of plot represents expression levels of cytokeratin genes present in our C2 gene signature.

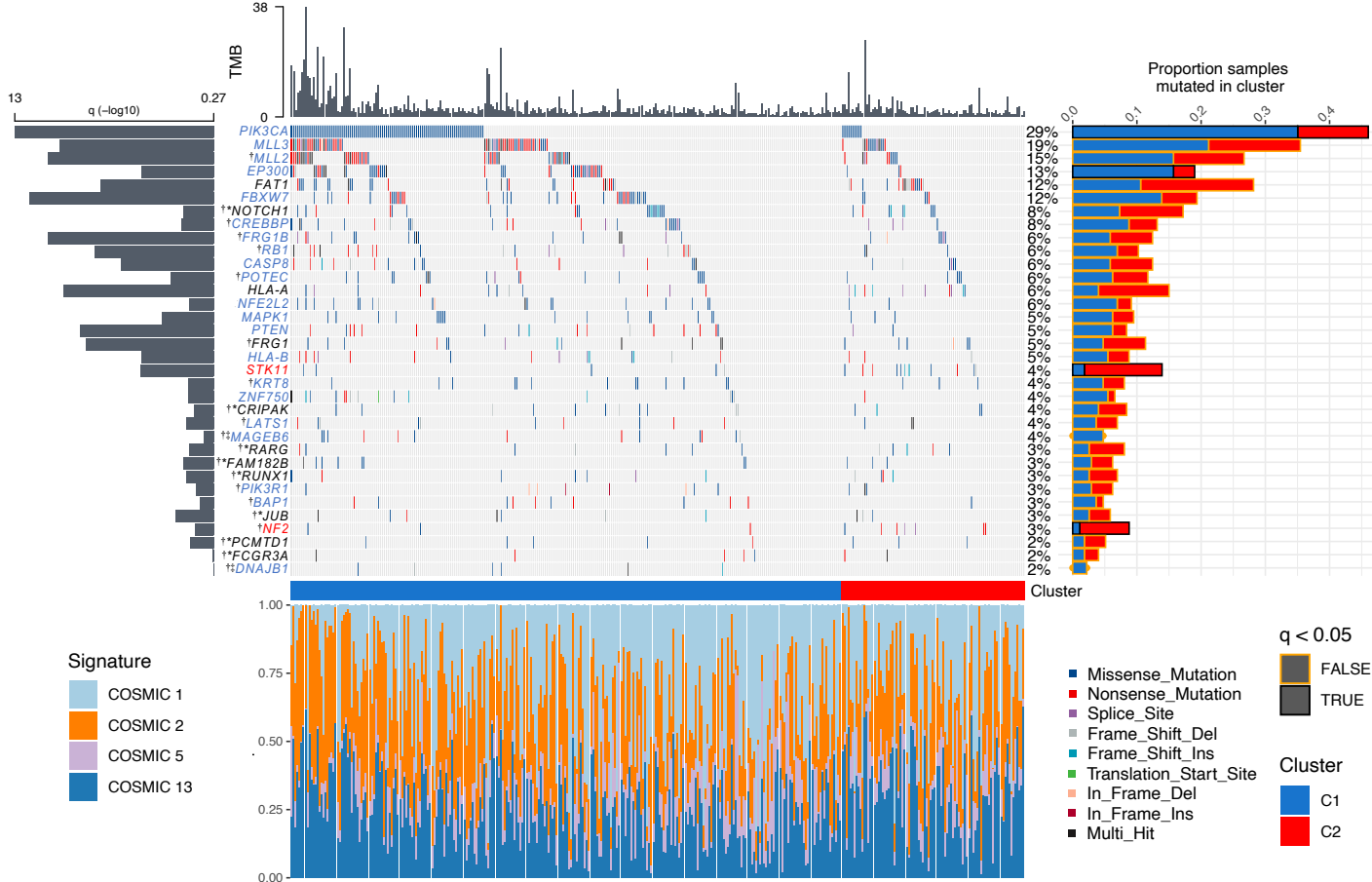


Figure 4 Genomic summary of significantly mutated genes (SMGs) in SCC cohorts. Main plot shows mutation type and frequencies for 34 SMGs identified using dNdSCV on TCGA, Bergen and Ugandan cohorts (367 total patients). Grey bars at top of plot represent TMB per sample. Grey bars to left of plot represent significance of SMG, larger bar is more significant. Barchart to the right shows proportion of a genes mutations in by cluster (blue = C1, red = C2). Black box around bar represents a significant difference in mutation frequency between the clusters ($p < 0.05$) while a gold box means no significant difference between the clusters. The plot at the bottom of figure represents the mutational signatures that contribute towards each individuals tumour mutational burden.

[Gene name key – blue – unique to C1 analysis, red = unique to C2 analysis, black = both in C1 and C2 individual analyses, black* = only significant when combining both clusters for analysis, † = novel SMG in cervical cancer, ‡ = not significant in combined cluster analysis but significant in C1 only analysis]

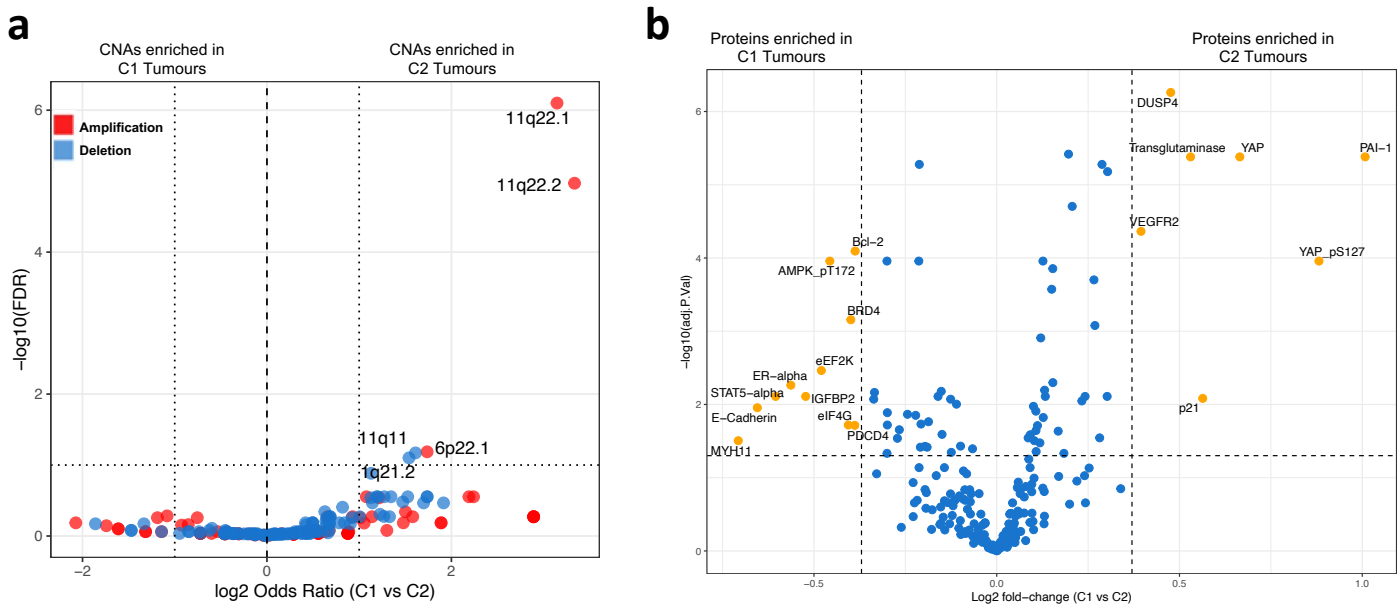


Figure 5 Copy number and protein level differences between SCC subgroups. **a)** Volcano plot showing differences in GISTIC copy number peak frequencies between C1 and C2 tumours, with $-\log_{10}(\text{FDR})$ on the y axis and the odds ratio on the x axis. **b)** Volcano plot showing differentially abundant proteins and phospho-proteins (FDR < 0.05, FC > 1.3, represented by yellow dots) between C1 and C2 TCGA tumours, as measured by Reverse Phase Protein Array.

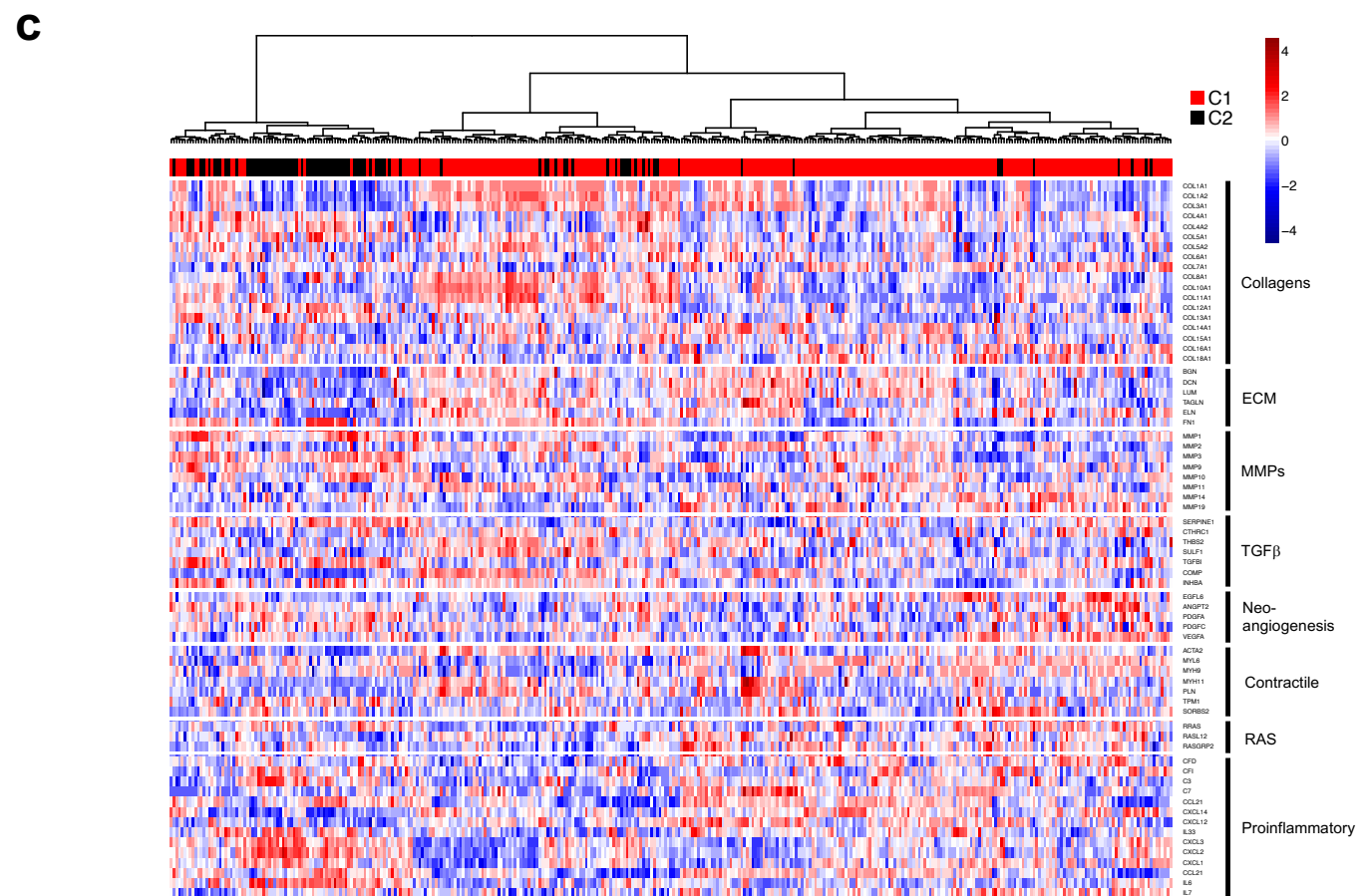
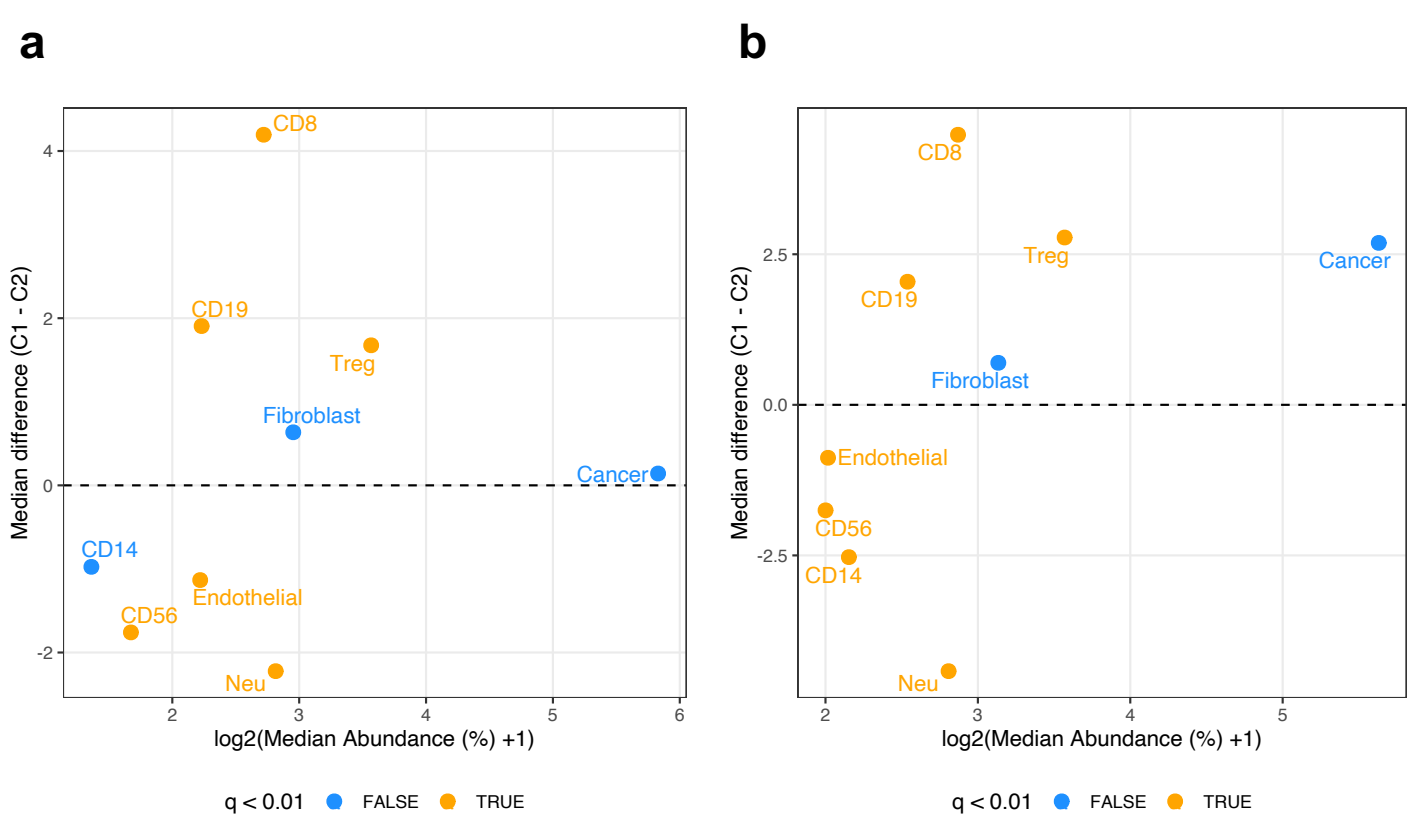


Figure 6 Differences in the tumour microenvironment between cervical cancer subgroups. Plot showing median abundances (x-axis) and median differences (% , y-axis) for different cell types estimated using MethylCIBERSORT, with significant differences in orange, for **a)** TCGA discovery cohort and **b)** combined validation cohorts. **c)** C2 tumours cluster together using CAF geneset genes.