

Analysis of head and neck carcinoma progression reveals novel and relevant stage-specific changes associated with immortalisation and malignancy.

Ratna Veeramachaneni^{1¶#a}, Thomas Walker^{1¶}, Antoine De Weck^{2&#b},
Timothée Revil^{3&}, Dunarel Badescu^{3&}, James O’Sullivan¹, Catherine
Higgins⁴, Louise Elliott⁴, Triantafillos Liloglou⁵, Janet M. Risk⁵, Richard
Shaw^{5,6}, Lynne Hampson¹, Ian Hampson¹, Simon Dearden⁷, Robert
Woodwards⁸, Stephen Prime⁹, Keith Hunter¹⁰, Eric Kenneth Parkinson⁹,
Ioannis Ragoussis³, Nalin Thakker^{1,4*}

1. Faculty of Biology, Medicine and Health, University of Manchester,
Manchester UK
2. Wellcome Trust Centre for Human Genetics, University of Oxford,
Oxford, UK
3. McGill University and Genome Quebec Innovation Centre, McGill
University, Montreal, Quebec, Canada
4. Department of Cellular Pathology, Manchester University NHS
Foundation Trust, Manchester, UK
5. Department of Molecular and Clinical Cancer Medicine, Institute of
Translational Medicine, University of Liverpool
6. Department of Head and Neck Surgery, Aintree University Hospitals
NHS Foundation Trust.

7. Precision Medicine and Genomics, IMED Biotech Unit, Astra Zeneca, Cambridge, UK
8. Department of Oral and Maxillofacial Surgery, Pennine Acute NHS Trust, Manchester, UK
9. Centre for Immunology and Regenerative Medicine, Institute of Dentistry, Barts and the London School of Medicine and Dentistry, Queen Mary University of London, Turner Street, London, UK
10. School of Clinical Dentistry, University of Sheffield, Sheffield, UK

* Corresponding Author

E-mail: nthakker@manchester.ac.uk

¶ These authors contributed equally to this work.

& These authors contributed equally to this work.

#a RV Bluestone Center for Clinical Research, New York University College of Dentistry, New York, New York 10010-4086, USA.

#b Novartis Institute for BioMedical Research, Basel, Switzerland.

1 **Abstract**

2 Head and neck squamous cell carcinoma (HNSCC) is a widely prevalent
3 cancer globally with high mortality and morbidity. We report here changes
4 in the genomic landscape in the development of HNSCC from potentially
5 premalignant lesions (PPOLS) to malignancy and lymph node metastases.
6 Frequent likely pathological mutations are restricted to a relatively small
7 set of genes including *TP53*, *CDKN2A*, *FBXW7*, *FAT1*, *NOTCH1* and
8 *KMT2D*; these arise early in tumour progression and are present in PPOLs
9 with *NOTCH1* mutations restricted to cell lines from lesions that
10 subsequently progressed to HNSCC. The most frequent genetic changes
11 are of consistent somatic copy number alterations (SCNA). The earliest
12 SCNAs involved deletions of *CSMD1* (8p23.2), *FHIT* (3p14.2) and *CDKN2A*
13 (9p21.3) together with gains of chromosome 20. *CSMD1* deletions or
14 promoter hypermethylation were present in all of the immortal PPOLs and
15 occurred at high frequency in the immortal HNSCC cell lines (promoter
16 hypermethylation ~63%, hemizygous deletions ~75%, homozygous
17 deletions ~18%). Forced expression of *CSMD1* in the HNSCC cell line
18 H103 showed significant suppression of proliferation ($p=0.0053$) and
19 invasion *in vitro* ($p=5.98 \times 10^{-5}$) supporting a role for *CSMD1* inactivation
20 in early head and neck carcinogenesis. In addition, knockdown of *CSMD1*
21 in the *CSMD1*-expressing BICR16 cell line showed significant stimulation
22 of invasion *in vitro* ($p=1.82 \times 10^{-5}$) but not cell proliferation ($p=0.239$).
23 HNSCC with and without nodal metastases showed some clear differences
24 including high copy number gains of *CCND1*, hsa-miR-548k and *TP63* in

25 the metastases group. GISTIC peak SCNA regions showed significant
26 enrichment (adj $P < 0.01$) of genes in multiple KEGG cancer pathways at
27 all stages with disruption of an increasing number of these involved in the
28 progression to lymph node metastases. Sixty-seven genes from regions
29 with statistically significant differences in SCNA/LOH frequency between
30 immortal PPOL and HNSCC cell lines showed correlation with expression
31 including 5 known cancer drivers.

32

33 **Lay Summary**

34 Cancers affecting the head and neck region are relatively common. A
35 large percentage of these are of one particular type; these are generally
36 detected late and are associated with poor prognosis. Early detection and
37 treatment dramatically improve survival and reduces the damage
38 associated with the cancer and its treatment. Cancers arise and progress
39 because of changes in the genetic material of the cells. This study focused
40 on identifying such changes in these cancers particularly in the early
41 stages of development, which are not fully known. Identification of these
42 changes is important in developing new treatments as well as markers of
43 behaviour of cancers and also the early or 'pre-malignant' lesions. We used
44 a well-characterised panel of cell lines generated from pre-malignant
45 lesions as well as cancers, to identify mutations in genes, and an increase
46 or decrease in number of copies of genes. We mapped new and previously
47 identified changes in these cancers to specific stages in the development

48 of these cancers and their spread. We additionally report here for the first
49 time, alterations in *CSMD1* gene in early premalignant lesions; we further
50 show that this is likely to result in increased ability of the cells to spread
51 and possibly, multiply faster as well.

52

53 **Introduction**

54 Globally, head and neck carcinomas account for over 550,000 new cases
55 per annum with a mortality of approximately 275,000 cases per year (1).
56 By far, the commonest site of cancer within this region is the oral cavity
57 and the commonest type of tumour is squamous cell carcinoma (SCC),
58 which accounts for over 90% of all malignant tumours at this site. HNSCC
59 is associated with high mortality having an overall 5-year survival rate of
60 less than 50%. Furthermore, both the disease and the multimodal
61 treatments options involved are associated with high morbidity (2).

62

63 The molecular pathology of head and neck squamous carcinoma has been
64 extensively studied previously and some of the common somatic genetic
65 changes have been variably characterised (3-7). There have been some
66 studies of the multistage evolution of these tumours (8) but this is less
67 well characterised. A small number of tumours arise from pre-existing
68 lesions (known as potentially malignant lesions or PPOLs) such as
69 leukoplakia or erythroplakia, which display variable epithelial dysplasia
70 (9). However, the vast majority are thought to arise *de novo* from

71 macroscopically normal appearing mucosa or possibly undiagnosed
72 PPOLs. Support for the latter comes from data showing that the
73 transcriptional signatures of PPOLs are retained in unrelated samples of
74 SCC both *in vivo* (10) and *in vitro* (11). Nevertheless, it is clear that
75 tumours arise from within a wide field bearing the relevant genetic
76 alterations and that there is a risk of synchronous or metachronous
77 tumours (12), (13). A fuller understanding of the events in evolution of
78 these cancers may permit the development of biomarkers or effective
79 therapeutic interventions possibly targeting not just tumours but also the
80 early changes in the field. The first multi-step model proposed for
81 carcinogenesis in HNSCC (14) suggested typical alterations associated
82 with progression from normal mucosa to invasive carcinoma, with
83 dysplasia reflecting an earlier stage of cancer progression.

84

85 We, and others, have previously shown that both SCCs and PPOLs yield
86 either mortal and immortal cells *in vitro* (15), (16) and, sometimes,
87 mixtures of the two (15), (16), (17). The status of these mortal cells is
88 unclear. Unlike the immortal cells, they lack inactivation of *TP53* and
89 *CDKNA*, but our limited previous investigations show that they are
90 genetically stable. Nevertheless, they often have extended replicative
91 lifespans (15, 16), possess neoplastic phenotypes, such as resistance to
92 suspension-induced terminal differentiation (15) and have expression
93 signatures which are distinct from both immortal cells and normal cells
94 (11). Furthermore, these characteristics are present in mortal cells from

95 both PPOLs and SCCs (11) suggesting the presence of distinct pathways
96 for the development of mortal and immortal SCC. Our preliminary work
97 established that the mortal PPOL were cytogenetically diploid and had low
98 levels of LOH (15) but the immortal PPOL and SCCs have never been
99 subjected to extensive genomic analysis. In addition, whilst extensive
100 genetic analysis of HNSCC has been carried out in recent years (Stransky
101 2011; Agrawal et al 2011; Pickering et al 2013, The Cancer Genome
102 Network 2015), including the identification of key driver mutations, the
103 stages in the cancer progression, at which they occur and the resulting
104 phenotypes are still unknown. There are numerous previous studies of
105 PPOLs limited to determining the frequency of alterations in specific gene
106 or specific genetic regions. Exceptions to this are the study by
107 Bhattacharya and colleagues (8) and Wood and colleagues (18), which
108 reported a comprehensive analysis of copy number variation in primary
109 PPOLs and HNSCC although the later study, the PPOL analyses was
110 confined to metachronous lesions. Here we extend their findings by a
111 combination of exome/targeted sequencing and SNP/CGH array analyses,
112 using our unique panel of mortal cultures and immortal cell lines derived
113 from both PPOLs and HNSCC, to show that mostly these genetic
114 alterations are mostly associated with cellular immortalisation and
115 increase with the stage of tumour progression in this class of SCC
116 keratinocyte.

117

118 **Results**

119 **Mutation analyses**

120 Several recent genomic sequencing studies have fully characterised
121 mutations in HNSCC (3),(4), (5), (6). In order to map these previously
122 identified mutations to progression of HNSCC, a small previously well-
123 characterised panel of samples consisting of 3 PPOL mortal cultures, 7
124 PPOL cell lines (from progressing and non-progressing lesions) 1 mortal
125 culture derived from HNSCC and 11 HNSCC cell lines (16, 19, 20) were
126 selected for exome-sequencing or targeted sequencing of the top 40
127 genes identified as altered in these cancers (3) using HaloPlex Target
128 Enrichment System (Agilent, Santa Clara, CA, USA). The sample details
129 are given in Fig. 1 and in S1 Tables.

130

131 For exome sequencing, approximately 6 gigabases of sequence mapped
132 to the human genome with an average of 65.7% (Range 33.8% to
133 86.1%) of the targeted exome covered at twenty-fold or higher (S2 Fig).
134 The lower coverage was sample specific and these samples are indicated
135 in Fig. 1. For HaloPlex sequencing, approximately 800 megabases of
136 sequence mapped to the human genome with an average of 94.8% of the
137 targeted exons covered at twenty-fold or higher (S2 Table).

138

139 For calling pathological mutations, we used the strategy described in
140 detail in the legend for Fig. 1; this was stringent and therefore, it is
141 possible that some genuine pathological mutations may have been
142 excluded. However, the full dataset with GEMINI framework annotation

143 (21) and raw data is available through Dryad Digital Repository (URL:
144 <https://datadryad.org/review?doi=doi:10.5061/dryad.314k5k5> Provisional
145 doi:10.5061/dryad.314k5k5)

146

147 Given the small numbers of samples examined in our study, we targeted
148 our analyses to 167 cancer drivers in head and neck cancers as identified
149 by IntOGen (Release 2014.12) The HaloPlex sequencing panel included 8
150 of the 10 most frequently mutated HNSCC driver genes. Limiting analyses
151 to known cancer drivers also effectively excluded possible false positives
152 that can arise due to DNA replication timing and low transcriptional
153 activity (22). Thus, our significant driver mutations (Fig. 1) largely mirror
154 those identified by Lawrence and colleagues, 2013 (22) following
155 correction for these factors. Significant mutations are shown
156 schematically in Fig. 1 and detailed in S 2C-2D Tables.

157

158 Mutations were rare in mortal cultures (Fig. 1). A single missense variant
159 of *FBXW7* predicted to be deleterious with low confidence was observed in
160 all 3 PPOL cultures (in addition to several immortal HNSCCs) and one high
161 impact *NOTCH1* mutation was observed in HNSCC culture BICR80.

162

163 As with previous studies (4), (3), (5), (6), the mutation analyses revealed
164 a small set of genes (*TP53*, *CDKN2A*, *FBXW7*, *FAT1*, *NOTCH1* and *KMT2D*)
165 as the most common targets for likely deleterious sequence mutations in
166 immortal PPOL and HNSCC cell lines (Fig. 1). Mutation of *TP53* and

167 *CDKN2A* as an early event in head and neck carcinogenesis is well
168 established (23), (24), (16), (25). In the present study, however, we
169 demonstrate for the first time that mutations in the other commonly
170 mutated genes in HNSCC also occur early and are not only present in
171 PPOLs derived from progressing lesions (D19, D20, D35) but also in non-
172 progressing lesions (D4, D34, D38).

173

174 As SCNAs may be an alternative mechanism for gain or loss of function of
175 a gene, we examined the frequency of SCNAs of genes predicted to be
176 cancer drivers by IntOgen but not showing sequence variation, in our
177 panel. Many cancer drivers showed a high frequency of SCNAs (S3
178 Tables). These included well-characterised tumour suppressors and
179 oncogenes implicated in HNSCC such as *CDKN2A*, *CCND1*, *EGFR*, *PIK3CB*
180 but also other cancer drivers that are less well characterised in HNSCC
181 (e.g., *CTTN*, *NDRG1*, *MLL3*, *ROBO2*). In addition, there were clear
182 differences between PPOLs and HNSCCs as well as between cell lines from
183 HNSCC with and without lymph node metastases. Although it can be
184 argued that some of these SCNAs may be bystander changes in
185 chromosomal gains or losses targeting other genes, consistent high
186 frequency changes are likely to be important as indicated by inclusion of
187 already well-characterised head and neck cancer drivers such as *MYC*,
188 *PIK3CA* and *CDKN2A*.

189

190 **Early changes in evolution of HNSCC**

191 Somatic copy number alterations in 7 PPOL cell lines, 11 mortal cell
192 cultures derived from PPOL, were identified using Illumina HumanHap550
193 Genotyping Beadchip and Infinium Assay II. The full dataset including raw
194 data and Nexus Copy Number v5.1 (BioDiscovery, Inc., CA, USA) data are
195 available at Dryad Digital Repository (URL:
196 <https://datadryad.org/review?doi=doi:10.5061/dryad.314k5k5> Provisional
197 doi:10.5061/dryad.314k5k5)

198

199 ***Mortal cultures are genetically stable.***

200 Mortal cultures derived from PPOL were genetically stable and showed
201 very few copy number changes, data that are consistent with their diploid
202 chromosome complement and previous limited loss of heterozygosity
203 analysis (15) (Fig 1 and S4 Fig.). There were no statistically significant
204 differences ($P > 0.05$) in SCNA between mortal PPOLS and matched
205 fibroblasts.

206

207 One mortal cell culture from PMOL (D17) that has an extended lifespan
208 and does not express CDKN2A but regulates telomerase normally and has
209 a functional *TP53* gene (11, 16, 19). This cell culture showed deletion of
210 one chromosome arm 9p with duplication of the homologous chromosome
211 9, 3 copies each of chromosome 2 and 7 and uniparental trisomy of
212 chromosome 5 (data not shown). It is possible that the changes involving
213 9p reflect loss of a normal *CDKN2A* allele and duplication of an allele with
214 hypomorphic mutation.

215

216 ***Immortal PPOL cell lines show progressive changes principally***
217 ***involving chromosome 3, 8, 9 and 20.***

218 Immortal cultures derived from PPOLs showed consistent SCNAs, thereby
219 defining the earliest changes in the development of HNSCC (Fig. 2a).
220 Overall, there were statistically significant ($P < 0.05$) losses on
221 chromosome 3p, 8p and 9p coupled with gains of chromosome 20
222 compared to normal fibroblasts (Fig. 2b). However, there was clear
223 heterogeneity in the nature of the changes observed in the 7 immortal
224 PPOL samples (Fig. 2c) with some specimens showing more frequent
225 changes than others. Meaningful testing for significant differences
226 between subsets of PPOLs was not possible due to the small sample size
227 and thus, hierarchical clustering was used to group samples by SCNAs
228 (Fig. 2c). This was remarkably similar to the clustering previously
229 observed for these cell lines by gene expression profiles (11). The
230 samples clustered into two main groups that suggested a correlation of
231 genomic changes with grade of dysplasia; the two cell lines derived from
232 PPOLs with mild dysplasia clustered separately from those derived from
233 PPOLs with higher-grade dysplasia. Furthermore, in the latter group, the
234 three cell lines (D19, D20, D35) from PPOLs that progressed to
235 carcinomas (progressive PPOLs or P-PPOL), clustered in a discrete
236 subgroup distinct from cell lines (D34, D4) derived from non-progressive
237 lesions (N-PPOL). Distinct gene expression signatures for normal tissues,
238 PPOLs and HNSCCs have been reported previously (10). In the present

239 study, we provide evidence that these signatures may at least in part
240 reflect the underlying somatic copy number changes.

241

242 Cell lines from progressive lesions (D19, D20, D35) were characterised by
243 consistent arm-level losses of chromosome arms 3p and 8p with
244 homozygous deletions of *FHIT* (3p14.1) and *CSMD1* (8p23.2) (S5 Fig). In
245 addition, two of the three P-PPOLs showed further arm-level SCNAs on
246 several other chromosomes (+3q, +5p, +7p +8q, -13p, -13q, -18p -18q,
247 +20). The pattern of SCNAs in cell lines derived from lesions that had not
248 progressed to date (NP-PPOLs, D34, D4, D9 and D38) reflected their
249 earlier stage of evolution (Fig. 2c). These cells harboured focal SCNAs
250 involving the *CSMD1* (3 of 4 cell lines) and *FHIT* (2 of 4 cell lines) on
251 chromosome 8p23.2 and 3p14.1 respectively, and showed arm-level
252 gains of chromosome 20 (3 of 4 cell lines). Additionally, other
253 chromosomal arms (+3q, +5p, +7p +8q, -13p, -13q, -18p -18q) also
254 showed variable and largely focal SCNA.

255

256 ***GISTIC analyses identifies key genes deleted early in progression*** 257 ***of HNSCC***

258 Significant peaks of copy number gains and losses were identified using
259 GISTIC (26). The genes in the peak regions at different levels of
260 stringency (Q=0.05-0.25) are shown in S6 Table. These included the
261 homozygously deleted genes described above (*FHIT*, *CSMD1*, *CDKN2A*,
262 *CDKN2B*). In addition, the peaks included homozygous deletions of *FAT1*,

263 thereby supporting loss-of-function mutations in HNSCC and identifying
264 inactivation of *FAT1* as an early change in HNSCC development.
265 Homozygous deletions of other genes (*NCKAP5*, *SORBS2*, *FAM190A*) not
266 previously implicated in HNSCC were also identified in peak regions. Some
267 genes such as *PTPRD*, *LRP1B* and *LINGO2*, which are target for
268 homozygous deletions in HNSCC cell lines, showed frequent hemizygous
269 deletions in the PPOL cell lines suggesting further selection in progression
270 to HNSCC. Deletion of *NCKAP5* was explored a little further (S7) as rare
271 SCNAs involving *NCKAP5* have been reported recently in HNSCC (3) (in
272 supplementary data) and in prostate cancer (27). Potentially deleterious
273 mutations have been reported in COSMIC (S7 Table). *NCKAP5*
274 homozygous deletions were present in both immortal PPOL and HNSCC
275 cell lines and the homozygous deletions eliminated one or more exons in
276 4 out of the 5 cell lines (S7 Fig) with the remaining cell line sustaining two
277 intronic deletions. However, deletions of *NCKAP5* are not common in
278 primary HNSCC and other tumour types (Tumorscape Release 1.6, our
279 tumour panel) and furthermore expression was not reduced significantly
280 in HNSCC and other tumour types (S7 Fig). Thus, *NCKAP5* is unlikely to
281 be frequent target for inactivation in HNSCC but the related pathways
282 may be more significant. *NCKAP5* has been shown to interact with *NCK1*
283 (STRING Release 9.1), an adaptor protein important in ligand-induced
284 activation of receptor tyrosine kinases (28), and also *APC* (BioGRID
285 Release 3.3). Our IntOGene (Release 2014.12) analyses showed a low
286 frequency of pan-cancer gain-of-function mutations in *NCK1*. In this

287 study, *NCK1* was within an extended GISTIC region and copy number
288 gains were seen in 60% of LN+ve cell lines with relatively low frequency
289 (~14%) gains in PPOL and LN-ve cell lines.

290

291 ***CSMD1 shows SCNAs early in development of HNSCCs and***
292 ***functional analyses suggest a tumour suppressor role***

293 Inactivation of *TP53* and *CDKN2A* are almost universal in both P-PPOL and
294 NP-PPOL immortal cell lines (16) and in HNSCC in vivo unless integrated
295 oncogenic HPV is present (7). Our data indicate that specific focal
296 deletions of *FHIT* and *CSMD1* are also common early changes. *FHIT* is
297 already well characterised and appears to have a tumour suppressor
298 function (29, 30) but relatively little is known about *CSMD1*. Thus, we
299 explored whether *CSMD1* functions as a tumour suppressor in HNSCC.
300 *CSMD1* is deleted in many tumour types (31) and also shows rare somatic
301 mutations (32). We observed both homozygous and hemizygous deletions
302 of *CSMD1* (5/28 and 21/28 respectively) in HNSCC cell lines (S8 Fig).
303 Analysis of expression in ProteomicsDB (33) suggests that the highest
304 expression in the body is in the oral epithelium although generally the
305 high expression is reported in brain and testis (Gene, NCBI).

306

307 Only one immortal PPOL (D4) sustained a nonsense mutation (G1579X) in
308 *CSMD1* and no likely pathogenic mutations were observed in HNSCC.
309 *CSMD1* promoter methylation analyses by pyrosequencing revealed
310 hypermethylation in 9 of 12 (75%) HNSCC cell lines with matching normal

311 samples (and several other HNSCC cell lines without matching normal
312 tissues), 3 of 7 (~43%) PPOL cell lines, and 15 of 24 (~63%) primary
313 HNSCCs (Figure 2). The highest level of promoter hypermethylation was
314 seen in immortal HNSCC cell lines with either hemizygous deletions
315 (BICR56, BICR22, BICR82, BICR10 and T4) or absence of deletions
316 (BICR63, BICR68 and H314), and in the only immortal PPOL cell line (D9)
317 that lacked *CSMD1* deletions (Fig. 3, and S5 and 8 Fig). Additionally, the
318 frequency of promoter methylation in primary tumours (~63%) was not
319 too dissimilar to the frequency of hemizygous deletions (~75%) in our
320 HNSCC cell lines. These findings suggest that frequent inactivation of
321 *CSMD1* in HNSCC occurs by deletion and/or promoter hypermethylation.

322

323 Stable transfection of full-length *CSMD1* cDNA into the H103 cell line,
324 which lacks endogenous expression of *CSMD1* expression, resulted in a
325 significant inhibition of proliferation ($p=0.0053$) and invasion ($p=5.98 \times 10^{-5}$
326 - Matrigel *in vitro* assay). Results from a representative clone are shown
327 in Fig. 4 (with data from all clones shown in S9 Fig). *CSMD1* expression
328 was also silenced by stable transfection with an shRNA vector in the
329 BICR16 cell line, which has a very low but detectable level of endogenous
330 *CSMD1* expression. This cell line has complete hemizygous deletion of the
331 *CSMD1* together with several small homozygous deletions of which one
332 involves an exon (S8 Fig) although the functional state of the protein is
333 unknown. The stable *CSMD1*-silenced clones showed a more variable
334 effect on both proliferation and invasion; clones displayed a significant

335 increase in invasion compared to parent cells ($p = 1.82 \times 10^{-5}$) but loss of
336 *CSMD1* did not have a significant effect on the rate of proliferation ($p =$
337 0.239). Data from representative clones are shown in Fig. 4 (with data
338 from all clones show in S9 Fig). It is possible that the inconsistent
339 proliferation results with gene silencing are due to the fact that these cell
340 lines have acquired the necessary cancer traits with some *CSMD1*
341 expression and that these traits are not significantly impacted by
342 additional knockdown of residual and possibly hypofunctional *CSMD1*.
343 Overall, however, these data support a role for *CSMD1* as a tumour
344 suppressor gene inactivated in the very early stages of HNSCC
345 development.

346

347 **Later changes in evolution of HNSCC**

348 SCNAs were analysed in two panels of tumour cell lines (S10 Fig) using
349 two different approaches (SNP array and array CGH). Since there was
350 little difference between the two panels in both high copy number
351 alterations (gains >2 and homozygous deletions) and low copy number
352 alterations (gains ≤ 2 and hemizygous deletions), the data were merged
353 for further analyses.

354

355 ***Progression to HNSCC is characterised by increased frequency of***
356 ***SCNAs of chromosomal regions involved in PPOLs as well further***
357 ***SCNAs of specific additional regions***

358 Overall, progression to HNSCC was characterised by an increased
359 frequency in the SCNAs observed in immortal PPOLs (Fig. 2) with
360 statistically significant increases in losses of proximal part of Chr3p,
361 Chr4q and Chr18q and gains of Chr20q. In addition, there was additional
362 loss of Chr10p coupled with gains of Chr5p, Chr9q, Chr14q and Chr11q.
363 Some genes (e.g., *CDKN2A*, *CDKN2B*) showing homozygous deletions in
364 PPOLs showed an increase in frequency of the same in HNSCCs. There
365 were in addition, homozygous deletions of further genes such as *PTPRD*,
366 *LRP1B* and *LINGO2* some of which showed high frequency hemizygous
367 losses in PPOLs. High copy number gains, which were rare in PPOLs, were
368 more frequent and principally centred on chromosome 11q. As with PPOL-
369 derived cell lines significant peaks of copy number gains and losses were
370 identified using GISTIC (26). The genes in the peak regions at different
371 levels of stringency ($Q=0.05-0.25$) are shown in S6 Tables.

372

373 ***Identification of key genes in progression from PPOLs to HNSCC***
374 ***by integrative analyses***

375 In addition to looking at differences in frequency of SCNAs between and
376 early and late lesions in HNSCC progression, we used integrative analyses
377 to further delineate key genes in transition from PPOLs to HNSCC. Gene
378 expression array data were available for 29 samples (11). We identified
379 genes that showed significant correlation of copy number with gene
380 expression, from genomic regions that showed statistically significant
381 difference in frequency of SCNA/LOH between PPOLs and HNSCC cell lines

382 (S11 Table). This identified 67 genes (Fig. 5) of which 50 have been
383 previously reported to be show some association with cancer (including
384 *NOTCH1* and *PIK3CA*) using PUBMED search. Nine of the 50 genes have
385 been previously associated with HNSCC including *DVL3*, and 5 genes
386 (*NOTCH1*, *PPP6C*, *RAC1*, *EIF4G1*, *PIK3CA*) were identified as mutational
387 cancer drivers in IntOGen (Release 2014.12).

388

389 ***Cell lines from HNSCC with and without lymph node metastases***
390 ***show specific differences in SCNAs.***

391 The aggregate HNSCC cell line data concealed notable differences
392 between cell lines from tumours with and without lymph nodal metastases
393 (LN+ve and LN-ve respectively). Lymph node metastases status of 17 of
394 28 HNSCC lines was known (10 LN+ve; 7 LN-ve). The LN+ve cell lines
395 showed higher frequency of high copy (>2) number gains of a 1.76 Mb
396 region at chromosome 11q13.2-q13.3 ($p < 0.05$) that encompasses 10
397 genes including *CCND1* and the microRNA hsa-miR-548k (Fig. 2).
398 Similarly, a higher frequency of high copy number gains on two regions
399 on chromosome arm 3q was observed involving *NAALADL2* (chromosome
400 band 3q26.32), and *TP63* and *CLDN1* (chromosome band 3q28).
401 Hierarchical clustering of all HNSCC cell lines by high-copy number
402 aberrations also revealed two major groupings defined by the presence or
403 absence of the high-copy number amplicon on chromosome band
404 11q13.2-q13.4 ($p < 0.01$, q bound < 0.1) (S12 Figure). In the group
405 lacking the amplicon, only 2 of 8 cell lines with known nodal staging, were

406 derived from LN+ve cancers and each of these involved a single node less
407 than 2 cm (TNM stage N1). By contrast, in the group with the amplicon, 7
408 of 9 cell lines were derived from LN+ve cancers and 4 of the 7 tumours
409 were graded as TNM stage N2 or higher. These results are consistent with
410 very recent findings linking the 11q13.3 amplicon with poor survival in
411 HNSCC patients (34). Comparison of all copy number changes also
412 revealed further important differences. LN+ve cell lines had a higher
413 frequency of copy number gains of chromosome 3q, 12q, 14q and 20,
414 together with a higher frequency of copy number losses on distal regions
415 on chromosome 3p and 11q, and chromosomes 4, and 18q (Fig. 2). Some
416 SCNAs (copy number gains in chromosome regions 7p12.2-21.3 and
417 9q31.3-32) were more frequent in LN-ve cell lines than in LN+ve cell
418 lines.

419

420 ***GISTIC peak regions of SCNA in PPOLs and HNSCCs show***
421 ***significant progressive enrichment of genes involved in cancer***
422 ***pathways***

423 Statistically significant (adj. $P < 0.01$) enrichment of genes in KEGG
424 'pathways in cancer' as well as other specific cancer pathways was
425 observed in GISTIC regions in both PPOLs and HNSCC cell lines (Fig. 6
426 and S13 Table). This enrichment was observed even if the analysis was
427 limited to 1466 genes in GISTIC regions that showed significant
428 correlation with expression after correction for multiple testing (adj.
429 $p < 0.05$) in immortal PPOL and HNSCC cell lines for which expression data

430 was available. Surprisingly, many genes in GISTIC regions that showed
431 high frequency of SCNAs in HNSCC including known HNSCC cancer drivers
432 (e.g., *CDKN2A*, *MYC*) did not show correlation with gene expression.
433 Therefore, we tested whether integrative analysis was a reliable method
434 to predict *in vivo* protein expression. We examined expression of protein
435 by immunohistochemistry of two genes (*BCL2L1* encoding an apoptosis
436 regulator and *CLDN1* encoding a component of tight junctions in epithelia)
437 that show copy number gains. *BCL2L1* gene showed high frequency/low
438 copy number gains in PPOLs (71%) and LN+ve HNSCC (90%) but not in
439 LN-ve HNSCC (28%). *CLDN1* showed low frequency/low copy number
440 gains in PPOLs (14%) and high frequency/low copy number gains in
441 HNSCC (73% overall, 53% LN-ve HNSCC, 90% LN+ve HNSCCs). *BCL2L1*
442 but not *CLDN1* shows significant correlation of copy number with gene
443 expression in this study (adjP=0.01 and 0.74 respectively). However, in
444 PPOLs and HNSCC biopsies (S14 Fig) both *CLDN1* and *BCL2L1* showed
445 significantly increased expression (p<0.0001) in HNSCC compared to
446 normal tissues and PPOL. This indicated that correlation of SCNA with
447 transcript expression in integrative analyses may not be a reliable
448 surrogate indicator of functional importance of a gene, and that protein
449 expression may better reflect the underlying SCNA.

450

451 The increase in both the size and the number of the SCNA regions in
452 HNSCC compared to PPOLs, and the differences in SCNAs between LN+ve
453 and LN-ve HNSCC cell lines, were reflected in the progressive increase

454 and/or differences in the enrichment for relevant KEGG pathways genes
455 (Fig. 6). Given that individual proteins participate in multiple pathways
456 and processes, many of the same genes mapped to multiple cancer-
457 related and other pathways. Some of the pathways identified such as
458 TP53 signalling and the cell cycle, have been well characterised in HNSCC,
459 but others such as axon guidance, actin cytoskeleton and motility, and
460 ubiquitin–proteosome pathway less so.

461

462 The examination in our cell line panel, of the frequencies of SCNAs
463 involving genes in individual pathways, allowed us to map multiple
464 pathway changes to stages in HNSCC progression (Fig. 7 and S15 Fig).
465 For example, the peak regions of SCNAs in LN+ve cell lines showed
466 enrichment of specific genes in the TGFB pathway that would predict
467 dysfunction of this signalling pathway (Fig. 7a). High frequency losses of
468 receptors (*TGFBRII*, *ACVR2B* and *BMPRI1B*), common *SMAD4*, R-*SMAD2*
469 and inhibitory *SMAD7* as well as other intracellular effectors such as
470 *PPP2R1B* and *RHOA*, were observed. This was coupled with high
471 frequency gains of *CHRD* and *TGFB3*. In addition, there were SCNAs of
472 downstream targets with losses of *CDKN2B* (normally induced by the
473 pathway) and gains of *MYC* (normally down-regulated by TGFB
474 signalling). Similarly, enrichment in SCNAs of genes in the *NOTCH1*
475 pathway was observed in both LN-ve and LN+ve cell lines (Figure 6b).
476 However, there were differences in the frequencies of SCNAs of the
477 individual genes (Figure 6b) in the pathway suggesting different

478 alterations of the pathway. LN+ve cell lines showed almost universal
479 amplification of *HES1* and *DVL3* and universal loss of *KAT2B*. LN+ve cell
480 lines also showed a relatively low frequency (20%) of high and low copy
481 number gains of *NOTCH1*. Gain of *DVL3* and loss of *KAT2B* together with a
482 relatively high frequency of gain of *NUMB* and loss of *MAML2* may be
483 expected to disrupt *NOTCH1* signalling but high frequency gains of the
484 downstream targets *HES1* and *HEY1* may negate these changes or
485 alternatively suggest complex amplification and inhibition of subsets of
486 NOTCH signalling elements. Clearly, these genes act in multiple pathways
487 and it is difficult to determine the effect of gain or loss of any single gene
488 without functional analyses. Thus, copy number gains of *DVL3* may be of
489 significance in WNT signalling pathway or in the cross-talk between the
490 two pathways. Amplification of *HES1* and *DVL3* and loss of *KAT2B* were
491 also observed in LN-ve cell lines as well as PPOLs albeit at a lower
492 frequency. Instead, the LN-ve cell lines were characterised by a high
493 frequency, low copy number gains of *NOTCH1* and *LFNG*. In our small
494 sample set, cell lines with mutations and amplifications of the *NOTCH1*
495 locus were mutually exclusive. Some genes such as *MAML1* and *MAML2*
496 showed both copy number gains and losses possibly indicating further
497 heterogeneity in NOTCH pathway aberration.

498

499 ***Multiple genes in SCNA regions may play a role in cancer***
500 ***progression.***

501 The prevalence and commonality of deletions and amplifications involving
502 specific chromosomal arms and regions in a wide range of cancers
503 coupled with the enrichment of known cancer-related genes in peak
504 regions of SCNAs and supporting functional evidence for many of the
505 genes in previous studies indicates selection for the SCNAs in cancer cells
506 is driven by the presence of several relevant genes on these chromosomal
507 arms/regions. Clearly, however, pathway analyses will not identify genes
508 such as *CSMD1* which may be cancer-relevant but whose functions are
509 not characterised or mapped to known pathways. We tested this further
510 by examining *ADAMTS9*, a gene not identified in pathway analyses but
511 showing frequent SCNAs.

512

513 *ADAMTS9* along with several other genes in this region at 3p14.2, shows
514 a single copy loss in just over 80% of the HNSCC and 29% of PPOL lines
515 suggesting possible further selection of deletions of this or neighbouring
516 genes in cancer progression. *ADAMTS9* is inactivated by promoter
517 hypermethylation in other tumour types including nasopharyngeal
518 carcinoma and oesophageal squamous cell carcinoma (35), (36),
519 Functional analyses suggest a role for *ADAMTS9* in inhibiting angiogenesis
520 (37). In our array CGH data of 347 tumours of multiple sites
521 (unpublished), *ADAMTS9* showed copy number losses in 23% of the
522 tumours (S16A Table) and in Tumorscape (Release 1.6) analyses, it was
523 focally deleted in epithelial tumours ($Q=8.38E-6$, frequency=0.369). In
524 the present study, we analysed mutations reported in COSMIC (v76,

525 cancer.sanger.ac.uk), (38). Of 40 mutations, 3 were nonsense mutations,
526 and 13 of 37 missense mutations (derived principally from lower
527 gastrointestinal tract) were predicted to be deleterious by Polyphen and
528 SIFT analyses (S16B Table). We failed to identify any likely pathological
529 mutations in our exome analysis and it was not part of our HaloPlex
530 sequencing panel..

531

532 We analysed *ADAMTS9* promoter methylation by pyrosequencing in a
533 subset of the HNSCC cell lines where DNA from matching normal tissue
534 was available and also in the DNA from set of primary HNSCC with
535 matching normal tissues. Promoter hypermethylation was observed in
536 both immortal HNSCC cell lines (7/17, ~41%) and primary HNSCC (9/20,
537 ~45%) (Fig. 7A-C). By contrast, promoter hypermethylation was less
538 frequent in PPOL cell lines (1/7, ~14%) and was not observed at all in the
539 mortal PPOL or mortal HNSCC cell lines (Fig. 8) mirroring the pattern of
540 *ADAMTS9* SCNAs in PPOL and HNSCC. We also analysed the expression of
541 the *ADAMTS9* transcript in primary HNSCCs and also in other cancers
542 using Tissuescan panel (Origene, Maryland, USA). In primary HNSCCs, 17
543 of 23 samples (~74%) showed reduced expression compared to the
544 normal tissues (Fig. 8D-E). In other tumour types, reduced expression
545 was seen in carcinomas of the breast, colon, lung and thyroid (Fig. 8F).
546 These findings suggest a possible role for *ADAMTS9* in HNSCC and other
547 cancers. Similarly, several other genes not previously mapped to any
548 KEGG pathway but within GISTIC regions may be significant. For

549 example, *BOP1* on chromosome arm 8q24.3 is in GISTIC focal region of
550 deletions in both LN-ve and LN+ve HNSCC cell lines (S6 Table). PPOL and
551 HNSCC cell lines showed high frequency, low copy number gains (PPOL,
552 43%; LN-ve HNSCC, 53%; LN+ve HNSCC 90%) and low frequency, high
553 copy number gains (PPOL 0%; LN-ve HNSCC 14%; LN+ve HNSCC 10%)
554 of *BOP1* and SCNA correlated with gene expression (adj $p < 0.001$) in
555 integrative analyses. *BOP1* has recently been shown to be a downstream
556 target of Wnt signalling and promotes cell migration and metastases in
557 colorectal carcinomas (39) and epithelial–mesenchymal transition,
558 migration and invasion in hepatocellular carcinoma (40).

559

560 **DISCUSSION**

561 Many recent studies have reported genetic and epigenetic changes in
562 HNSCC (4), (3), (5), (6), (7). Here, we have extended these findings by
563 further detailed genomic analyses of a unique panel of cell lines from
564 premalignant lesions (PPOLs) as well as subsets of HNSCCs with and
565 without lymph node metastases. This approach has allowed us to map the
566 genetic changes to the stages of evolution of HNSCCs and to identify the
567 earliest abnormalities, which are significant in tumour progression. The
568 separate analyses of the HNSCC subgroups with respect to nodal
569 metastases has facilitated the identification of changes which may be
570 otherwise masked by looking at average changes in a heterogeneous
571 group.

572

573 The prevalence of the changes observed in this study must be regarded
574 with caution given the small sample size and the use of cell lines.
575 Changes observed in cell lines may not be fully representative of primary
576 tumour because of clonal selection and continuing evolution in culture.
577 However, the patterns of changes observed in our study are consistent
578 with studies of primary tumours. Our GISTIC analyses identified similar
579 peak and extended regions as those identified in 3131 primary tumours
580 by Beroukhim et al., 2010 (31). Balanced against this, one major
581 advantage of using cell lines is that extensive heterogeneity of primary
582 tumour samples can mask SCNAs (41), (42), whereas early passage cell
583 lines give cleaner results (43). Nevertheless, it is essential that the
584 findings are verified in larger sample sets of primary premalignant lesions
585 and tumours.

586

587 We have also demonstrated that techniques such as GISTIC, pathway and
588 integrative analyses for identifying the pertinent or significant changes
589 amongst the myriad changes observed in tumours have limitations.
590 *BCL2L1* and *CLDN1* both with high frequency copy number gains in
591 HNSCC, show increased protein expression in primary tumours despite
592 lack of correlation in integrative analyses for *CLDN1*. *ADAMTS9*, which is
593 in not in a GISTIC peak region in our sample set and does not map to a
594 KEGG pathway, showed frequent copy number loss and promoter
595 hypermethylation together with decreased gene expression in HNSCC but
596 not in PPOLs or normal tissues indicating potential role in HNSCC.

597

598 Given the relatively low frequency (~30% or less) of sequence mutations
599 of genes other than *TP53* both in the current and previous studies,
600 together with prevalence of consistent and frequent SCNA changes across
601 wide range of tumours (31), it is likely that the selection processes in the
602 clonal evolution of these tumours is driven by these SCNAs. Additionally,
603 we show here that the earliest changes are often characterised by focal
604 SCNAs (for example, deletions at chromosome 8p23 involving *CSMD1*)
605 with further selection for loss of whole arm or large region of arm during
606 progression or evolution. This suggests that regions of SCNAs harbour
607 multiple genes that collectively provide selective growth advantage. This
608 is further supported by the KEGG pathway analyses, which show
609 statistically significant enrichment of genes in the cancer relevant
610 pathways in peak regions of SCNAs identified in HNSCC using GISTIC.
611 Additionally, there are numerous functional studies in the literature of
612 different genes in the same chromosomal regions, which are involved in
613 development of tumours of same type. Although this suggests that there
614 are a large number of cancer drivers, it is possible that the genes do not
615 represent primary drivers but their alterations provide a further
616 cumulative selective advantage against a background of primary driver
617 loss-of-function or gain-of-function.

618

619 In the present study, we confirm previous observations (23), (24), (16),
620 (25) that loss-of-function of *TP53* (primarily through sequence mutations)

621 and *CDKN2A* (through SCNA, promoter hypermethylation and sequence
622 mutations) are early changes in HNSCC development. Furthermore, we
623 demonstrate for the first time that loss-of-function sequence mutations in
624 *NOTCH1*, *KMTD2* (*MLL2*) and *FBXW7* are present in PPOLs and represent
625 early but less common changes.

626

627 In accordance with previous observations (8), we demonstrate that a loss
628 of chromosome arms 3p, 8p and 9p and gains of chromosome 20 are the
629 earliest changes in HNSCC Using cell lines from progressive and non-
630 progressive PPOLs, we show that the earliest changes are characterised
631 by focal deletions and/or promoter hypermethylation of *CSMD1* on
632 chromosome arm 8p23.2. The role of *CSMD1* in cancer is relatively
633 unknown but deletions at this locus have been reported in many types of
634 cancers (31). Like many other commonly homozygously deleted genes in
635 HNSCC and other tumours, *CSMD1* (and *FHIT*) and are large genes that
636 are located in regions of low gene density or 'gene deserts' (S17 Table).
637 Deletions appear to occur more frequently in such regions and thus, at
638 least some of these SCNAs may be aetiologically unrelated to cancer
639 development and over represented through low selection pressure against
640 these changes (31). Furthermore, in this study, a few pseudogenes in
641 'gene deserts' were also found to have sustained homozygous deletions
642 (Supplementary Data 17) supporting the notion that at least of some of
643 these SCNAs are 'non-specific' changes. Nevertheless, there is convincing
644 functional evidence supporting a tumour suppressor role for at least some

645 of these genes including *FHIT* (30), (29) and other large genes such as
646 *DCC* that are located in gene deserts (44). In the present study, we
647 provide functional evidence for the first time for a tumour suppressor role
648 for *CSMD1* in head and neck squamous mucosa. Our findings support
649 findings in breast cancer reported by (45). *CSMD1* encodes for a predicted
650 transmembrane protein with a multidomain extracellular structure that is
651 likely to act as a multi-ligand receptor mediating endocytosis of ligands.
652 However, this remains to be characterised. One previous study has
653 reported functional analyses in melanoma cell lines demonstrating a role
654 for *CSMD1* in reducing proliferation and invasive potential possibly
655 through *SMAD* pathway (46).

656

657 In this study, whilst we did not detect mutations in *FAT1* (a gene known
658 to be mutated in HNSCC) in PPOL, we did identify hemizygous and
659 homozygous deletions of the gene confirming this gene as an early target
660 for inactivation in HNSCC development. In addition to *FAT1*, we also
661 identified novel homozygous deletions in *NCKAP5*, *SORBS2* and *FAM190A*
662 in PPOLs. *SORBS2* has been shown to induce cellular senescence (47).
663 *FAM190A* is a structural or regulatory component of mitosis and its loss
664 may contribute to chromosomal instability (48). Little is known of *NCKAP5*
665 and we demonstrated that this is not a frequent target of SCNAs or
666 reduced expression in HNSCC or other tumour types. However, its
667 product interacts with product of *NCK1*, a gene in the extended GISTIC
668 region on chromosome arm 3q; frequent amplifications of *NCK1* in LN+ve

669 HNSCC cell lines used in this study and pan-tumour low frequency gain-
670 of-function mutations in the IntOGen dataset suggest a possible novel
671 cancer-relevant pathway. Interestingly, NCK proteins are essential
672 signalling elements in cytoskeleton organisation and cellular motility (49)
673 and in the present study, there was a significant enrichment of the genes
674 in the cytoskeleton organisation and cell motility pathways in the
675 extended GISTIC regions in this study. NCK1 has been reported to be
676 necessary for EGFR-mediated migration and metastases in pancreatic
677 cancer (50) and depletion of NCK1 increases UV-induced TP53
678 phosphorylation and apoptosis (51).

679

680 In contrast to study by Bhattacharya and colleagues (8), this study did
681 not find evidence of significant subset of immortal PPOLs or HNSCCs that
682 didn't show loss at 8pter-p23.1 together with gains on 3q24-qter, 8q12-
683 q24.2, and chromosome 20. However, our sample size is relatively small
684 and there may be a selection in culture of the cells with these genetic
685 alterations. Interestingly, in the previous study (8) the subgroup lacking
686 these genetic alterations showed genetic stability, lack of *TP53* mutations
687 and a much-reduced predisposition to metastasis. It is interesting to
688 speculate whether these correspond to 'mortal' PPOLs and HNSCC
689 cultures, which are genetically stable and also lack *TP53* mutations.

690 Our findings with respect to mutations and SCNA are similar but not
691 identical to those reported by (18) in synchronous dysplasia and HNSCC.
692 In our study, we were able to show progressive changes with

693 transformation to malignancy and lymphovascular spread. We note that in
694 their study only a minority of low-grade dysplasias showed changes that
695 were present in high-grade dysplasia and HNSCC, and they suggested
696 that SCNAs were not necessary for the low-grade dysplasias to develop.
697 We interpret this with caution as unsurprisingly, low-grade dysplasia can
698 be very difficult to distinguish from normal tissues with absolute certainty
699 and agreement between histopathologists is generally weak in grading
700 low-grade dysplasia. Thus, some low-grade dysplasia lesions may
701 represent genuine dysplasia whilst others may present normal tissues.
702 Regardless of this, what is clear from our study is that the SCNAs arise
703 after the breakdown of cellular senescence. The mortal PPOLs do not
704 display SCNAs. Additionally, the loss of expression of *CDKN2A* is nearly
705 ubiquitous in our PPOL panel (16) and in PPOL tissues *in vivo* (24),(52),
706 (11). Furthermore, PPOL D17, which has lost *CDKN2A* expression whilst
707 remaining mortal and retaining normal *TP53* and telomerase status, has
708 only minimal chromosomal gains and no losses, suggesting that CNAs
709 follow breakdown of senescence.

710 The results of this study demonstrated that GISTIC extended SCNA
711 regions in the PPOLs and in HNSCC with and without nodal metastases,
712 harbour genes that are implicated in a number of cancer-relevant KEGG
713 pathways. Several of the genes that we identified have been reported
714 previously as being associated with HNSCC but the present study is the
715 first to systematically identify these genes and others, in the context of
716 cancer-relevant pathways and map them to cancer progression. We have

717 described this in context of TGFB and NOTCH signalling pathways and
718 provided supplementary data for other cancer-related pathways.

719

720 Our data show enrichment of genes of the TGFB signalling pathway in the
721 GISTIC regions and higher frequency of copy number loss associated in
722 LN+ve HNSCC compared to LN-ve HNSCC and PPOLs. Disruption of TGFB
723 pathway in HNSCC is well established (53). *Smad2*-null mice, for
724 example, develop spontaneous HNSCC and consistent with our findings in
725 this study, copy number losses of *SMAD2*, *SMAD4* and *TGFBRII* are
726 associated with an aggressive tumour phenotype and lymph node
727 metastases (54), (53). In the present study, we also report relatively
728 common SCNAs of other genes in TGFB pathway that have seldom been
729 reported and or not at all in HNSCC. These anomalies include copy
730 number losses of activin/BMP receptor *ACVR2B*, downstream and SMAD-
731 independent pathway effector *RHOA*, together with copy number gains of
732 BMP inhibitor chordin (*CHRD*) and *ID1*. *ID1* is a helix-loop-helix protein
733 that has been shown to induce immortalization in keratinocytes (55) and
734 overexpression of the protein has been reported recently in HNSCC (56).

735

736 *NOTCH1* mutations (4), (3) and *NOTCH1* pathway alterations have been
737 reported in HNSCC (5). We have extended these findings and
738 demonstrate that *NOTCH1* mutations are present in two of the three
739 progressive PPOLs but not in any of the non-progressive lesions. This
740 indicates that *NOTCH1* inactivation is a relatively early event but still

741 consistent with previous observations that suggest that *NOTCH1*
742 inactivation plays a key role in progression to invasive carcinoma of
743 already initiated cells (57). Our data are also consistent with that of
744 Agrawal and colleagues (4) who have shown no association or mutual
745 exclusivity of *NOTCH1* and *TP53* mutations. Interestingly, the results of
746 the present study demonstrate that SCNAs of several genes in the
747 *NOTCH1* signalling pathway including *NOTCH1*, are more frequent than
748 loss-of-function mutations of *NOTCH1*. Furthermore, many (but not all) of
749 these SCNAs are gain of function changes in the *NOTCH1* pathway. This is
750 consistent with recent observations in primary HNSCC demonstrating
751 over-expression of both ligands and receptors in this pathway (58), (59).
752 Although mainly loss-of-function mutations in *NOTCH1* have been
753 reported in HNSCC to date, activating mutations in HNSCC in a Chinese
754 population have also been described recently (60). Our data, therefore
755 add support to the emerging consensus for dual oncogenic and tumour
756 suppressive role for *NOTCH1* in HNSCC although further functional
757 analyses are necessary to confirm this proposal. Some SCNA's of *NOTCH1*
758 pathway genes such as copy number gains of *DVL3* may also be
759 significant in the context of cross-talk with the WNT signalling pathway.

760

761 In this study, we identified a potentially deleterious *NOTCH1* sequence
762 variant in two immortal poorly differentiated PPOL lines and in two
763 immortal and one mortal HNSCC lines. If the mutation in the mortal line
764 which has a wild type p53 (57) and no detectable CNVs (this study) or

765 LOH (15) is pathological as predicted, it would suggest that *NOTCH1*
766 mutations are independent of genomic instability. The presence of
767 *NOTCH1* mutations in this line also offers a plausible explanation for the
768 poor differentiation of this line in both suspension (15, 61) and surface
769 culture (15) and is also consistent with recent data showing that the
770 knockdown of *NOTCH1* expression in human keratinocytes recreates a
771 poorly differentiated epithelium reminiscent of dysplasia (62).
772 Furthermore, *NOTCH1* has been shown to mediate keratinocyte
773 stratification (61, 63) and stem cell maintenance (64). However, *NOTCH1*
774 deletion has also been shown to promote tumourigenesis and tumour
775 progression through paracrine effects (65), the former of which would
776 also be consistent with *NOTCH1* mutations in PPOLs . Moreover, this last
777 observation would be consistent with the reports of *NOTCH2* and *NOTCH3*
778 mutations in human SCC (3) because loss-of-function of these paralogues
779 promotes tumourigenesis in mouse skin in a paracrine fashion but does
780 not replicate the effect of *NOTCH1* deletion on keratinocyte differentiation
781 (65). It has also been reported that *NOTCH1* is a TP53 target gene (61)
782 and as most of the immortal PPOL and HNSCC lines have TP53 mutations
783 this could be an additional mechanism of its inactivation and is consistent
784 with the altered regulation of hairy enhancer of split 2 (*HES2*) in these
785 lines (11, 61).

786

787 Integrative analyses in a subset of samples for genes from genomic
788 regions showing a significant difference in frequency of SCNA between

789 PPOLs and HNSCC identified 67 genes that showed correlation with gene
790 expression including *NOTCH1*, *PPP6C*, *RAC1*, *EIF4G1*, *PIK3CA* and *DVL3*.
791 The role of *NOTCH1* and *PIK3CA* in HNSCC is well established. *PPP6C*,
792 *RAC1* and *EIF4G1* are also likely cancer drivers according to IntoGen.
793 *RAC1* activation has been reported previously in HNSCC and mediates
794 invasive properties of HNSCC (66); our observation of copy number gains
795 in tumour progression is consistent with this finding. *EIF4G1* is a
796 translation initiation factor that is part of the multi-subunit complex EIF4F
797 that facilitates recruitment of mRNA to the ribosome. This is a rate-
798 limiting step protein synthesis initiation phase. *EIF4G1* is amplified in
799 many tumour types in TCGA data sets (67). Over-expression of *EIF4G1*
800 promotes tumour cell survival and formation of tumour emboli through
801 increased translation of specific mRNAs in inflammatory breast cancer
802 (68). Components of EIF4F are also targets of C-MYC and initiate further
803 translation of specific targets including C-MYC (69). *PPP6C* loss-of-
804 function mutations have been reported in melanomas (70) but in the
805 present study, we observed copy number gains particularly in LN-ve
806 HNSCC; the significance of this observation is unclear. *DVL3* is a
807 transducer of both canonical and non-canonical Wnt signaling pathways
808 (71). Association with HNSCC has not been reported previously but it is
809 amplified in many tumour types in TCGA data sets (67) and inhibition of
810 Wnt signalling in HNSCC results in inhibition of growth and metastases
811 (72).

812

813 In conclusion, we have further characterised specific genetic changes that
814 mark progression in head and neck squamous cell carcinogenesis.
815 Although genomic landscapes and progression models of SCCHN (14),
816 (4), (6), (5), (3) have been published previously, we have been able to
817 use our well characterised cell line panels to tentatively assign genetic
818 changes, including novel ones, to specific stages of progression in
819 transcriptionally distinct mortal and immortal classes (9), (11) of the
820 disease and also to cell function.

822

823 **Materials and Methods**

824 **Samples**

825 This study was approved by the UK National Research Ethics Service
826 Research Ethics Committee (08/H1006/21).

827

828 Details of the samples are shown in Supplementary Data 1. For SNP
829 array, the sample set consisted of 16 HNSCC cell lines, 7 PPOL cell lines
830 and 11 mortal cell cultures derived from PPOL. DNA from matching
831 fibroblasts was available for 6 HNSCC cell lines, 1 immortal PPOL cell line
832 and 2 mortal PPOL cultures. The sample culture conditions were as
833 described previously (15) and DNA was prepared from cell lines using
834 standard protocols. The sample set for array CGH consisted of 12 HNSCC
835 cell lines.

836

837 **SNP array analyses**

838 SNP genotyping of the primary HNSCC panel was performed using the
839 Illumina HumanHap550 Genotyping Beadchip and Infinium Assay II as per
840 standard protocols. DNA from the cell lines was quantitated with
841 NanoDrop (Thermo Scientific) and 750ng was used per assay.

842

843 **Array CGH**

844 ArrayCGH data were kindly provided by Dr Simon Deardon, AstraZeneca,
845 UK.

846

847 **Data analyses**

848 The average genotype call rate was 98.25%; genotype data from two
849 samples with a call rate of <95% in BeadStudio v3.1 (Illumina) were
850 excluded from analyses. Over 75% of samples had GenTrain score
851 (measure of reliability based on the total array of calls for a given SNP) of
852 ≥ 0.7 and none were below 0.4. Data were pre-processed in
853 GenomeStudio v2009.1 (Illumina) and imported into Nexus Copy Number
854 v5.1 (BioDiscovery, Inc., CA, USA) and OncoSNP v2.7 (73) for further
855 analyses. ArrayCGH data from a second HNSCC panel were also imported
856 into Nexus Copy Number v7.5 (BioDiscovery, Inc., CA, USA).

857

858 The robust variance sample QC calculation (a measure of probe to probe
859 variance after major outliers due to copy number breakpoints are
860 removed from the calculation) in Nexus Copy Number v7.5 (BioDiscovery,
861 Inc.), was used to assess the quality of the samples. Data for samples
862 with a score > 0.2 were excluded and the score for the remaining
863 samples was in the range 0.03-0.20. For identification of copy number
864 and copy neutral changes, the BioDiscovery's SNPRank Segmentation
865 Algorithm was used with significance level of 1×10^{-6} and a minimum
866 number of probes per segment of 5. Thresholds for determining copy
867 number variation were set at -1 for homozygous deletion, -0.18 for
868 hemizygous deletion, 0.18 for gain (single copy gain) and 0.6 for high
869 gain (2 or more copies). An area was considered to be showing LOH if

870 95% of the probes in the region had a B allele frequency of >0.8 or <0.2
871 (homozygous frequency and value thresholds of 0.95 and 0.8
872 respectively). Allelic imbalance was defined as 95% of the probes in the
873 region showing a B allele frequency of between 0.2 and 0.4 or 0.6 and 0.8
874 (i.e. heterozygous imbalance threshold of 0.4).

875

876 Areas of the genome with a statistically high frequency of aberration (Q-
877 bound value $\leq 0.5-0.25$ and G-score cut-off ≤ 1) after correction for
878 multiple testing using FDR correction (Benjamini & Hochberg), were
879 identified using the GISTIC (Genomic Identification of Significant Targets
880 in Cancer) approach (26).

881

882 Group comparisons were made in Nexus with differences in frequency of
883 specific events at any chromosomal location tested for significance by
884 two-tailed Fisher's Exact Probability Test with an accepted significance
885 level of $p < 0.01$ at a defined level of percentage difference. More stringent
886 Q-bound values based on two-tailed Fisher's Exact Probability Test
887 corrected for multiple testing using Benjamini-Hochberg FDR correction
888 (Benjamini & Hochberg) as well as a minimum of set percent difference in
889 frequencies between the two groups; significance was accepted at <0.25 .

890

891 **Promoter methylation analyses**

892 Genomic DNA (approximately 1mg each) from all the HNSCC and PPOL
893 cell lines and matched primary HNSCC samples used were subjected to

894 bisulfite treatment using the EZ DNA methylation™ kit (Zymo Research,
895 U.S.A) according to the manufacturer's protocol. 30ng each of the
896 bisulfite-treated DNA was used for the pyrosequencing reaction.

897

898 PCR and sequencing primers for the pyrosequencing methylation analyses
899 of the CpG rich promoter region were designed using the pyro-Q-CpG
900 software for the genes *ADAMTS9* and *CSMD1*. The forward primer was
901 biotinylated and was used in low concentration (5pmol) along with
902 amplification cycles of 45 to exhaust the primers. The forward and
903 reverse primer sequences used in the study were *ADAMTS9*- F-
904 5'agagatttttaaagttaaaagttgg3', R-5'tcctcctaccctcctta3' with the
905 sequencing primer 5'cctcctaccctcctta3' and *CSMD1*-F-5'
906 gtagtttagatagatagagtttagttt3', R-5' acaaatctcctttctcca3' with its
907 sequencing primer 5' aaatctcctttctccaacct3'. Optimized annealing
908 temperature for *ADAMTS9* and *CSMD1* PCR primer pairs is 54°C. Using
909 bisulfite-treated DNA as a template, regions of interest were amplified by
910 standard PCR cycling conditions in a 96-well plate using Qiagen's Hot start
911 Taq Polymerase to avoid nonspecific amplification. The specificity of the
912 PCR products was then verified by agarose gel electrophoresis. For the
913 pyrosequencing reaction, the PCR product was made single stranded by
914 immobilizing the incorporated biotinylated primer on streptavidin-coated
915 beads. The sequence run and analysis were done on the PyroMark™Q96
916 MD pyrosequencer (Qiagen) according to the manufacturer's instructions.
917 The sequence runs were analysed using the Pyro Q-CpG software. The

918 peak heights observed represented the quantitative proportion of the
919 alleles. The software generated methylation values for each CpG site and
920 also the mean methylation percentage for all the CpG sites analyzed.

921

922 **Generation of stable *CSMD1*-expression modulated clones**

923 A panel of nine OSCC cell lines was profiled for *CSMD1* transcript and
924 protein expression status (data not shown). This identified the *CSMD1*-
925 expressing cell line BICR16 and the *CSMD1* non-expressing cell line H103.

926

927 BICR16 cells were used to generate stable *CSMD1*-silenced monoclonal
928 and polyclonal lines with HuSH 29mer pRS shRNA (Origene, Rockville, MD,
929 USA). The degree of silenced *CSMD1* expression level was confirmed by
930 RT-qPCR and flow cytometry assays. Primer sensitivity assays determined
931 the limit of *CSMD1* transcript detection at five ORF copies per light-cycler
932 well.

933

934 Stable forced *CSMD1*-expressing monoclonal cells were generated by
935 transfection of 15.5kb pCMV6-*CSMD1* expression plasmid (Origene,
936 Rockville, MD, USA) into the *CSMD1*-deleted cell line H103. The plasmid
937 was linearized at the SexAI restriction site within a predetermined 14%
938 non-essential region and used to generate seven *CSMD1*-expressing
939 monoclonal cell lines using standard methods. The degree of forced
940 *CSMD1* expression was confirmed by RT-qPCR and flow cytometry assays.

941

942 Cell proliferation was determined with CellTiter 96® Aqueous Cell
943 Proliferation MTS Assay (Promega, Southampton, UK) as per the
944 manufacturer's protocols. Cell growth was calculated as a percentage
945 growth change from the 24-hour time point and population-doubling
946 times were determined. Gel-invasion was assayed using trans-well BD
947 Bio-coat Matrigel Invasion Chambers and control wells (BD Biosciences,
948 Oxford) as per manufacturer's protocols. Optimal seeding densities were
949 determined empirically. Triplicate Matrigel invasion chambers were used
950 for each clone from a minimum of two different Matrigel batches. Three
951 fields of view were captured for each Matrigel or control chamber (outer,
952 middle, centre areas, each at 120° rotation from one another).
953 Percentage invasion was calculated for each clone and expressed as an
954 invasion index (ratio of clone to parent percentage invasion). Statistical
955 analysis was performed in IBM SPSS 20 & 22 (Wilcoxon Sign-Rank Test)
956 with alpha levels set at 0.05.

957

958 **Integrative analysis**

959 Before correlating the SCNA and gene expression values corrections were
960 applied for polyploidy and heterogeneity.

961

962 When considered a by-product of instability rather than a response of
963 biological significance, the copy number (CN) values were altered to
964 remove the ubiquitous chromosomal amplification observed in polyploid
965 samples. To do so, all the CN values inferred for SNPs in chromosome

966 arms with mean CN larger than 2.5 were reduced by one unit. The
967 reduction was however rejected in the cases of heterozygous copy neutral
968 calls (CN2 LOH0), copy losses (CN1) and homozygous deletions (CN0).
969 These states were not altered.

970

971 Due to heterogeneity, the gene expression values obtained did not
972 represent the expression of the CN-altered cells solely. For each genomic
973 region, a non-negligible proportion of cells do not harbour any alteration.
974 The proportion of cells with normal heterozygous copy number in each
975 region was estimated. When investigating the correlation between SCNAs
976 and mRNA expression the weighted mean CN and LOH values of this
977 mixture were used.

978

979 **Exome sequencing**

980 Targeted enrichment and sequencing were performed on 1-3 µg of DNA
981 extracted from the cell lines. Enrichment was performed using the
982 SureSelect Human All Exon 50 MB v4 Kit (Agilent, Santa Clara, CA, USA)
983 for the Illumina system. Sequencing was carried out on a HiSeq 2500
984 sequencer (Illumina Inc, San Diego, CA, USA), following the
985 manufacturer's protocols.

986

987 **HaloPlex Sequencing**

988 Targeted enrichment and sequencing were performed on 225ng of DNA
989 extracted from the cell lines. Enrichment was performed using a custom

990 HaloPlex Kit (Agilent, Santa Clara, CA, USA) targeting 41 genes.
991 Sequencing was undertaken on a MiSeq sequencer (Illumina Inc, San
992 Diego, CA, USA), following the manufacturer's protocols.

993

994 **Sequence data analysis**

995 Raw paired-end reads were trimmed using Trimmomatic v0.33 to a
996 minimum length of 30 nucleotides. Illumina Truseq adapters were
997 removed in palindrome mode. A minimum Phred quality score of 30 was
998 required for the 3'end. Single end reads as well as paired end reads that
999 failed previous minimum quality controls were discarded. Individual read
1000 groups were aligned, using bwa v0.7.12 with default parameters, to the
1001 UCSC hg19 reference human genome from Illumina iGenomes web site.
1002 Trimming rates and insert length were controlled on each read group
1003 based on metrics reported by Trimmomatic, and Picard v1.128
1004 respectively.

1005

1006 Aligned reads from multiple read groups belonging to the same sample
1007 were indexed, sorted and merged using sambamba v0.5.4. Amplification
1008 duplicates were removed using Picard.

1009

1010 Various quality controls parameters were used including the obtained
1011 target coverage of the Nextera Rapid Capture exome library v1.2,
1012 mapping rates and duplication rates, based on metrics collected for each

1013 sample using Samtools v1.2, Picard v1.141, bedtools v2.25.0, and
1014 aggregated using custom Python v2.7.9 codes.

1015

1016 We applied GATK v3.5.0 base quality score recalibration and indel
1017 realignment [14] with standard parameters. We performed SNP and
1018 INDEL discovery and genotyping across each cohort of samples
1019 simultaneously using standard hard filtering parameters according to
1020 GATK Best Practices recommendations.

1021

1022 All variants were annotated with functional prediction using SnpEff v4.2.
1023 Additionally, functional annotation of variants found in two public
1024 databases (NCBI dbSNP v144 and dbNSFP v2.9) was added using SnpSift,
1025 part of the same software package. Multiallelic variants were decomposed
1026 and normalized using vt. A GEMINI v0.18.3; a database was created [15],
1027 and variants selected according to functional rules. Finally, they were
1028 manually validated against read alignments, using Integrative Genomics
1029 Viewer software (IGV) v2.3. Coverage metrics were calculated using
1030 bedtools.

1032

1033 **Figure Legends**

1034 **Fig 1.** Mutations of IntOgen (version 2014.12; <http://www.intogen.org>) -
1035 predicted head and neck cancer driver genes in mortal PPOL (M-PPOL)
1036 cultures, immortal PPOL (IM-PPOL) progressive and non –progressive (P
1037 and NP respectively) cell lines and mortal HNSCC cultures (M) and
1038 immortal (IM) cell line panels. The green shaded sample number indicate
1039 samples that were subject to HaloPlex target enrichment and sequence
1040 analyses of specific genes The remainder were subject to exome
1041 sequencing. The gene names shaded indicate genes in the HaloPlex panel
1042 that are also HNSCC cancer drivers as indicated by IntOgen. Samples with
1043 at least one mutation predicted to be high impact (nonsense mutations,
1044 indels, frameshift, splice site) are shown in red, samples with at least one
1045 missense mutation predicted to be deleterious by both Polyphen and SIFT
1046 are shown in yellow, and samples with at least one missense mutation
1047 predicted to be deleterious by Polyphen or SIFT (but not both) are shown
1048 in blue. Any variants previously reported as constitutional SNP was
1049 excluded regardless of the minor allele frequency unless it had been
1050 demonstrated to be pathogenic previously. The ranking is from IntOgen
1051 and indicates ranking by frequency of mutations in HNSCC. *Indicates a
1052 gene with the same variant in each sample with a variant. † Indicates
1053 samples with low coverage (<80% x20) in exome sequencing. Only genes
1054 showing significant mutations by criteria used are shown here.

1055

1056 **Fig 2. Somatic copy number changes in HNSCC A.** Density plots and
1057 karyograms showing copy number gains (blue) and losses (red) in normal
1058 fibroblasts, mortal PPOLs (M-PPOL), immortal PPOLs (IM-PPOL) and
1059 HNSCC cell lines with and without lymph node metastases (LN+ve and
1060 LN-ve respectively). **B.** Subtraction karyograms showing differences in
1061 copy number gains (blue) and losses (red) between (i) normal fibroblasts,
1062 and mortal PPOLs (M-PPOL), (ii) immortal PPOLs (IM-PPOL) and all
1063 HNSCC cell lines and (iii) HNSCC cell lines with and without lymph node
1064 metastases (LN+ve and LN-ve respectively). Both high and all copy
1065 number changes are shown for LN+ve and LN-ve HNSCC cell lines. Blue
1066 arrows on the right indicate lanes with regions of difference showing
1067 statistical significance ($p < 0.05$). **C.** Density plots and karyograms showing
1068 copy number gains (blue) and losses (red) in individual immortal PPOL
1069 cell lines showing hierarchical clustering and grade of dysplasia.

1070

1071 **Fig 3. Clustered column chart representing pyrosequencing**
1072 **analyses of the *CSMD1* promoter region.** Sample ID is shown on the
1073 X-axis and mean methylation percentage is represented on the Y-axis. For
1074 each sample, a mean methylation percentage greater than 5% was
1075 considered as significant promoter methylation. **A.** Primary HNSCCs and
1076 matching normal tissues. Significant promoter methylation was observed
1077 in 15 of 24 (~63%) primary HNSCC. **B.** Mortal and immortal HNSCC cell
1078 lines. Significant promoter methylation was observed in 12 of 17 (~70%)
1079 immortal HNSCC cell lines but not in any of the mortal HNSCC cell lines

1080 (indicated by *). The highest levels of promoter methylation were in cell
1081 lines with hemizygous deletions (BICR56, BICR22, BICR82, BICR10 and
1082 T4) or no deletions (BICR63, BICR68 and H314). **C.** PPOL cell
1083 lines/cultures. Significant promoter methylation was observed 3 of 7
1084 (~43%) immortal PPOL lines but not in any of the mortal PPOL cultures
1085 (indicated by *). The highest level of promoter methylation was in a
1086 single immortal cell line (D9) that had no deletions at the *CSMD1* locus.

1087

1088 **Fig 4. Phenotypic effects of *CSMD1* expression modulation in**
1089 **HNSCC.** Left panel: forced *CSMD1* expression in the *CSMD1* non-
1090 expressing H103 cell line. Right panel: silencing of *CSMD1* in the *CSMD1*-
1091 expressing BICR16 cell line. A. *CSMD1* mRNA transcript quantification by
1092 RT-qPCR and protein quantification by flow cytometry for generated
1093 clones, presented as fold change normalised to the reference cell line.
1094 Left: H103 *CSMD1*-expressing clones and H103 *CSMD1*-negative parent
1095 cells and *CSMD1*-disrupted clone H103-mcl-21, normalised to SCC116
1096 *CSMD1* expressing cells (red outline). Right: BICR16 *CSMD1*-silenced
1097 clones normalised to *CSMD1*-expressing BICR16 parent cell line. Boxplots
1098 represent RQ normalised to reference cells. Each box plot is the relative
1099 quantification (RQ) of two plates each of triplicate target and reference
1100 gene CT values plus and minus log-transformed standard deviations and
1101 so incorporates intra-plate variance. Standard boxes depict the first-third
1102 quartiles, whiskers depict ± 1.5 IQR. Median values are provided. Bar
1103 charts represent *CSMD1* protein fold change normalised to reference cells.

1104 Error bars are \pm standard deviation. B. Effects of modulation of *CSMD1*
1105 expression on cell proliferation. Left: *CSMD1*-expressing H103
1106 monoclonal clones compared to *CSMD1*-negative H103 parent and control cells.
1107 *CSMD1* expression resulted in a reduced growth rate compared to
1108 *CSMD1*-negative parent (black line) and H103-mcl-21 cells (dotted black
1109 line) (shaded area) ($p = 0.0053$). Right: Cell proliferation of *CSMD1*-
1110 silenced BICR16 clones compared to *CSMD1*-expressing BICR16 parent
1111 cells. *CSMD1* silencing did not result in a significant change in growth rate
1112 compared to *CSMD1*-expressing parent cells (black line) ($p = 0.239$). This
1113 observation was further confirmed using an additional *CSMD1*-silenced
1114 OSCC cell line with 11 silenced monoclonal clones and 1 silenced polyclonal clone (data
1115 not shown). Plots represent triplicate points from duplicate 96AQ assays
1116 for 96 hours with growth rates normalised to achieve relative fold change
1117 values for 0 -72hrs. Error bars are ± 1 SD. Growth rate differences across
1118 the BICR16 clones and parent cells illustrate the proliferation variance
1119 inherent across the clone pool, rather than differences being specifically
1120 due to loss of expressed *CSMD1* underpinning a clonal growth variation
1121 effect. C. Effects of modulation of *CSMD1* expression on gel invasion.
1122 Three trans-well chambers of 2 representative clones and parent cells are
1123 displayed (for full dataset see S7). The invasion index of generated clones
1124 vs. parent is depicted as bar charts (white and black bars respectively).
1125 Left: *CSMD1*-expressing clones vs. *CSMD1*-negative H103 parent cells.
1126 *CSMD1* expression results in a marked decrease of gel invasion
1127 ($p=5.975 \times 10^{-5}$). Right: *CSMD1*-silenced clones vs. *CSMD1*-expressing

1128 BICR16 parent cells. *CSMD1* silencing results in a marked increase in gel
1129 invasion ($p=1.822 \times 10^{-05}$).

1130

1131 **Fig 5. Frequency of SCNA of genes from chromosomal regions**
1132 **showing significant difference in SCNAs between PPOLs and**

1133 **HNSCC and significant correlation with expression. A.** Copy number

1134 gains **B.** Copy number losses. Genes reported to be associated with any

1135 cancer previously by PUBMED search ('Gene name, Cancer'). **Genes

1136 reported to be associated with HNSCC previously by PUBMED search

1137 ('Gene name, HNSCC', oral cancer). Chromosomal regions showing

1138 differences between PPOLS and HNSCC are indicated by letters on the

1139 left; in top panel (A), the regions are: A, chr3:151842842-152984767; B,

1140 chr3:157160169-158933894; C, chr3:160599782-161161335; D,

1141 chr3:161705726-162566403; E, chr3:170691970-173911476; F,

1142 chr3:177707736-180451354; G, chr3:185035692-185529164; H,

1143 chr3:198578155-199298372; I, chr7:1631815-7317208; J,

1144 chr9:85367221-85868737; K, chr9:97596715-98774190; L,

1145 chr9:99427940-100373274; M, chr9:101504820-101852863; N,

1146 chr9:123376788-125083831; O, chr9:126825879-127177239; P,

1147 chr9:129518474-129927677; Q, chr9:132985780-136636113; R,

1148 chr9:137291321-139534231; S, Chr14:62865875-63862093; T,

1149 chr3:186575268-187080482. In bottom panel (B), the regions are: A,

1150 chr3:57677987-58154068; B, chr3:61422777-73764765; C,

1151 chr3:87035206-88461236; D, Chr3:78706269-79206160; E,

1152 chr10:16585949-17022250; F, chr10:26833288-28028304; G,
1153 chr3:161705726-162566403; chr17:8028078-9207567. The regions are
1154 ranked A onwards in order of decreasing statistical significance.

1155

1156 **Fig 6. Enrichment of genes in cancer-relevant pathways in GISTIC**
1157 **extended regions in PPOLs and HNSCCs.** Statistically significant
1158 enrichment (adjusted $P < 0.01$) for specific pathways is indicated by a red
1159 box and lack of enrichment by a green box. The data are shown in detail
1160 in Supplementary Data 9. The ranges for adjusted P values corrected for
1161 multiple testing were: 'All HNSCC', $1.97E-12 - 0.0098$; 'LN+ve HNSCC',
1162 $5.76E-14 - 0.0083$; 'LN-ve HNSCC', $7.52E-06 - 0.0074$; PPOL, $5.34E-09 -$
1163 0.0007 . In each section, the pathways were ranked from top to bottom in
1164 order of level of significance in the 'All HNSCC' group with highest level of
1165 significance at the top.

1166

1167 **Fig 7. Frequency of SCNAs of genes involved in selected cancer**
1168 **pathways that are significantly enriched in the GISTIC regions in**
1169 **PPOLs and HNSCC. A.** TGFB pathway **B.** NOTCH pathway. For each
1170 pathway, two charts are shown illustrating the frequency of copy number
1171 gains (top panel) and losses (bottom panel) in PPOLs, all HNSCC, and
1172 HNSCCs with and without nodal metastases (LN+ve and LN-ve
1173 respectively). **Genes showing significant correlation with expression in
1174 integrative analyses after correction for multiple testing (adj. $p < 0.05$).
1175 *Genes showing nominal significance ($p < 0.05$) only are indicated by a

1176 single asterisk. Only genes showing at least 40% frequency of SCNA in at
1177 least one subgroup, are shown.

1178

1179 **Fig 8. Promoter hypermethylation analyses of *ADAMTS9*.** Sample ID
1180 is shown on the X-axis and the mean methylation percentage is
1181 represented on the Y-axis. For a sample, a mean methylation percentage
1182 greater than 5% was considered as significant promoter methylation. **A.**
1183 Primary HNSCCs and matching normal samples. Promoter methylation
1184 was observed in 9 of 20 (~45%) primary HNSCCs (results from samples
1185 3232, 3241, 3242 and 3247 were excluded from analysis as either the
1186 normal or tumour reaction failed). **B.** Mortal and immortal HNSCC cell
1187 lines. Promoter methylation was observed in 7 of 17 (~41%) immortal
1188 HNSCC cell lines, but not in any of the mortal HNSCC cell lines (indicted
1189 by *). **C.** PPOL cell lines/cultures. Promoter methylation was observed in
1190 only 1 of 7 (~14%) immortal PPOL cell lines and none of the mortal PPOL
1191 cell lines (indicted by *).

1192

1193 **Supplementary Data**

1194 S1. Details of study samples. Table 1 Clinical data of the PPOLs and
1195 HNSCC cell lines/cultures; Table 2 Matching fibroblasts for PPOLs and
1196 HNSCC cell lines/cultures; Table 3 Primary HNSCC and matching normal
1197 mucosa samples used for pyrosequencing methylation analyses and gene
1198 expression analyses.

1199

1200 S2A. Cumulative distribution of coverage showing fraction of targeted
1201 bsd (Y-axis) that's were covered by at least certain depth (x-axis) in
1202 exome sequencing

1203

1204 S2B. Cumulative distribution of coverage showing fraction of targeted
1205 bsd (Y-axis) that's were covered by at least certain depth (x-axis) in
1206 HaloPlex sequencing

1207

1208 S2C. Significant variants in HNSCC driver genes from exome sequencing.

1209

1210 S2D. Significant variants in HNSCC driver genes from HaloPlex
1211 sequencing.

1212

1213 S3. Frequency of copy number alterations in IntOgen-derived cancer
1214 driver genes PPOLs and HNSCC cell lines. Table 1: Cancer drivers showing
1215 low copy number gains (≤ 2) ordered by frequency in LN+ve HNSCC cell
1216 lines with minimum frequency of 40%; Table 2: Cancer drivers showing
1217 low copy number gains (≤ 2) ordered by frequency in LN-ve HNSCC cell
1218 lines with minimum frequency of 40%; Table 3: Cancer drivers showing
1219 high copy number gains (>2) ordered by frequency in LN+ve HNSCC cell
1220 lines with minimum frequency of 10%; Table 4: Cancer drivers showing
1221 high copy number gains (>2) ordered by frequency in LN-ve HNSCC cell
1222 lines with minimum frequency of 10%; Table 5: Cancer drivers showing
1223 hemizygous loss ordered by frequency in LN+ve HNSCC cell lines with

1224 minimum frequency of 40%; Table 6: Cancer drivers showing hemizygous
1225 loss ordered by frequency in LN-ve HNSCC cell lines with minimum
1226 frequency of 40%; Table 7: Cancer drivers showing homozygous loss
1227 ordered by frequency in LN+ve HNSCC cell lines.

1228

1229 S4. Subtraction karyogram of mortal PPOLs and matched normal
1230 fibroblasts.

1231

1232 S5. Immortal PPOL cell lines SCNA density plots for chromosome 3, 8, 9
1233 and 20.

1234

1235 S6. GISTIC peak regions in PPOL and HNSCC cell lines. Table 1. GISTIC
1236 peak regions in PPOLs with varying significance thresholds (Q bound
1237 value); Table 2. GISTIC peak regions in HNSCCs with varying significance
1238 thresholds (Q bound value); Table 3. GISTIC peak regions in LN-ve
1239 HNSCCs with Q bound value<0.25; Table 4. GISTIC peak regions in
1240 LN+ve HNSCCs with Q bound value<0.25.

1241

1242 S7. Analysis of *NCKAP5*. Table 1. Analysis of *NCKAP5* mutations listed in
1243 COSMIC Fig 1. Hemizygous and homozygous deletions at *NCKAP5* locus;
1244 Fig.2. Expression analyses of *NCKAP5* in HNSCC; Figure 3. Expression
1245 analyses of *NCKAP5* in other tumour types.

1246

1247 S8. Deletions in HNSCC cell lines at the *CSMD1* locus.

1248

1249 S9. Modulation of the *CSMD1* expression in HNSCC cell lines.

1250

1251 S10. Comparison of SCNA in two HNSCC panels.

1252

1253 S11. Integrative analyses of somatic copy number changes and gene
1254 expression. Table 1. Copy number gains in PPOLs and HNSCC of genes
1255 that show correlation with expression in SCNA regions showing significant
1256 difference between PPOLs and HNSCCs; Table 2. Copy number losses in
1257 PPOLs and HNSCC of genes that show correlation with expression in SCNA
1258 regions showing significant difference between PPOLs and HNSCCs.

1259

1260 S12. Hierarchical cluster analyses of immortal HNSCC cell lines for high
1261 copy number SCNA – association with lymph node metastases.

1262

1263 S13. KEGG pathway genes enrichment in GISTIC peak regions

1264

1265 S14. Analysis of *CLDN1* and *BCL2L1* in primary HNSCC and PPOL. Fig 1.
1266 Expression of *CLDN1* in PPOLS and HNSC; Figure 2. Expression of *BCL-XL*
1267 in PPOLS and HNSCC

1268

1269 S15. SCNAs of genes in cancer-related KEGG pathway enriched in GISTIC
1270 regions. Fig 1 Pathway in cancer; Fig 2 Apoptosis pathway; Fig 3 Axon
1271 guidance pathway; Fig 4 Cell adhesion pathway; Fig 5 Cell cycle pathway;

1272 Fig 6 Endocytosis pathway; Fig 7. JAK-Stat pathway; Fig 8 MAPK
1273 pathway; Fig 9 Ubiquitin-proteasome pathway; Fig 10 WNT pathway.

1274

1275 S16A. Somatic copy number changes in *ADAMTS9* in a multitumour cell
1276 line panel.

1277

1278 S16B. Analyses of *ADAMTS9* sequence variants identified in this study
1279 reported in Stransky et al., 2011 and COSMIC database.

1280

1281 S17. Gene size and density for genes sustaining homozygous deletions in
1282 immortal HNSCC cell lines.

1284 **References**

1285

1286 1. Ferlay J, Soerjomataram I, Dikshit R, Eser S, Mathers C, Rebelo M,
1287 et al. Cancer incidence and mortality worldwide: sources, methods and
1288 major patterns in GLOBOCAN 2012. *Int J Cancer*. 2015;136(5):E359-86.

1289 2. Pace-Balzan A, Shaw RJ, Butterworth C. Oral rehabilitation following
1290 treatment for oral cancer. *Periodontol 2000*. 2011;57(1):102-17.

1291 3. Stransky N, Egloff AM, Tward AD, Kostic AD, Cibulskis K,
1292 Sivachenko A, et al. The mutational landscape of head and neck
1293 squamous cell carcinoma. *Science*. 2011;333(6046):1157-60.

1294 4. Agrawal N, Frederick MJ, Pickering CR, Bettegowda C, Chang K, Li
1295 RJ, et al. Exome sequencing of head and neck squamous cell carcinoma
1296 reveals inactivating mutations in NOTCH1. *Science*.
1297 2011;333(6046):1154-7.

1298 5. Pickering CR, Zhang J, Yoo SY, Bengtsson L, Moorthy S, Neskey DM,
1299 et al. Integrative genomic characterization of oral squamous cell
1300 carcinoma identifies frequent somatic drivers. *Cancer Discov*.
1301 2013;3(7):770-81.

1302 6. India Project Team of the International Cancer Genome C.
1303 Mutational landscape of gingivo-buccal oral squamous cell carcinoma
1304 reveals new recurrently-mutated genes and molecular subgroups. *Nat*
1305 *Commun*. 2013;4:2873.

- 1306 7. The Cancer Genome Atlas N. Comprehensive genomic
1307 characterization of head and neck squamous cell carcinomas. *Nature*.
1308 2015;517(7536):576-82.
- 1309 8. Bhattacharya A, Roy R, Snijders AM, Hamilton G, Paquette J,
1310 Tokuyasu T, et al. Two distinct routes to oral cancer differing in genome
1311 instability and risk for cervical node metastasis. *Clin Cancer Res*.
1312 2011;17(22):7024-34.
- 1313 9. Hunter KD, Parkinson EK, Harrison PR. Profiling early head and neck
1314 cancer. *Nat Rev Cancer*. 2005;5(2):127-35.
- 1315 10. Ha PK, Benoit NE, Yochem R, Sciubba J, Zahurak M, Sidransky D, et
1316 al. A transcriptional progression model for head and neck cancer. *Clin*
1317 *Cancer Res*. 2003;9(8):3058-64.
- 1318 11. Hunter KD, Thurlow JK, Fleming J, Drake PJ, Vass JK, Kalna G, et al.
1319 Divergent routes to oral cancer. *Cancer Res*. 2006;66(15):7405-13.
- 1320 12. Bedi GC, Westra WH, Gabrielson E, Koch W, Sidransky D. Multiple
1321 head and neck tumors: evidence for a common clonal origin. *Cancer Res*.
1322 1996;56(11):2484-7.
- 1323 13. Tabor MP, Brakenhoff RH, Ruijter-Schippers HJ, Kummer JA,
1324 Leemans CR, Braakhuis BJ. Genetically altered fields as origin of locally
1325 recurrent head and neck cancer: a retrospective study. *Clin Cancer Res*.
1326 2004;10(11):3607-13.
- 1327 14. Califano J, van der Riet P, Westra W, Nawroz H, Clayman G,
1328 Piantadosi S, et al. Genetic progression model for head and neck cancer:
1329 implications for field cancerization. *Cancer Res*. 1996;56(11):2488-92.

- 1330 15. Edington KG, Loughran OP, Berry IJ, Parkinson EK. Cellular
1331 immortality: a late event in the progression of human squamous cell
1332 carcinoma of the head and neck associated with p53 alteration and a high
1333 frequency of allele loss. *Mol Carcinog.* 1995;13(4):254-65.
- 1334 16. McGregor F, Muntoni A, Fleming J, Brown J, Felix DH, MacDonald
1335 DG, et al. Molecular changes associated with oral dysplasia progression
1336 and acquisition of immortality: potential for its reversal by 5-azacytidine.
1337 *Cancer Res.* 2002;62(16):4757-66.
- 1338 17. Rheinwald JG, Beckett MA. Tumorigenic keratinocyte lines requiring
1339 anchorage and fibroblast support cultured from human squamous cell
1340 carcinomas. *Cancer Res.* 1981;41(5):1657-63.
- 1341 18. Wood HM, Daly C, Chalkley R, Senguvan B, Ross L, Egan P, et al.
1342 The genomic road to invasion-examining the similarities and differences in
1343 the genomes of associated oral pre-cancer and cancer samples. *Genome*
1344 *Med.* 2017;9(1):53.
- 1345 19. Muntoni A, Fleming J, Gordon KE, Hunter K, McGregor F, Parkinson
1346 EK, et al. Senescing oral dysplasias are not immortalized by ectopic
1347 expression of hTERT alone without other molecular changes, such as loss
1348 of INK4A and/or retinoic acid receptor-beta: but p53 mutations are not
1349 necessarily required. *Oncogene.* 2003;22(49):7804-8.
- 1350 20. McGregor F, Wagner E, Felix D, Soutar D, Parkinson K, Harrison PR.
1351 Inappropriate Retinoic Acid Receptor- β Expression in Oral Dysplasias:
1352 Correlation with Acquisition of the Immortal Phenotype. *Cancer Research.*
1353 1997;57(18):3886-9.

- 1354 21. Paila U, Chapman BA, Kirchner R, Quinlan AR. GEMINI: Integrative
1355 Exploration of Genetic Variation and Genome Annotations. PLoS Comput
1356 Biol. 2013;9(7):e1003153.
- 1357 22. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K,
1358 Sivachenko A, et al. Mutational heterogeneity in cancer and the search for
1359 new cancer-associated genes. Nature. 2013;499(7457):214-8.
- 1360 23. El-Naggar AK, Lai S, Luna MA, Zhou X-D, Weber RS, Goepfert H, et
1361 al. Sequential p53 mutation analysis of pre-invasive and invasive head
1362 and neck squamous carcinoma. International Journal of Cancer.
1363 1995;64(3):196-201.
- 1364 24. Papadimitrakopoulou Vali IJ, Lippman M Scott, Lee Soo Jin, Fan
1365 Hong You, Clayman Gary, Ro Y Jay, Hittelman N Walter, Lotan Reuben,
1366 Hong K Waun and Mao Li Frequent inactivation of p16INK4a in oral
1367 premalignant lesions. Oncogene. 1997;14:1799-803.
- 1368 25. Qin G-Z, Park JY, Chen S-Y, Lazarus P. A high prevalence of p53
1369 mutations in pre-malignant oral erythroplakia. International Journal of
1370 Cancer. 1999;80(3):345-8.
- 1371 26. Beroukhim R, Getz G, Nghiemphu L, Barretina J, Hsueh T, Linhart
1372 D, et al. Assessing the significance of chromosomal aberrations in cancer:
1373 Methodology and application to glioma. Proceedings of the National
1374 Academy of Sciences. 2007;104(50):20007-12.
- 1375 27. Kluth M, Galal R, Krohn A, Weischenfeldt J, Tsourlakis C, Paustian L,
1376 et al. Prevalence of chromosomal rearrangements involving non-ETS

- 1377 genes in prostate cancer. *International journal of oncology*.
1378 2015;46(4):1637-42.
- 1379 28. Buday L, Wunderlich L, Tamas P. The Nck family of adapter
1380 proteins: regulators of actin cytoskeleton. *Cell Signal*. 2002;14(9):723-
1381 31.
- 1382 29. Joannes A, Bonnomet A, Bindels S, Polette M, Gilles C, Burlet H, et
1383 al. Fhit regulates invasion of lung tumor cells. *Oncogene*.
1384 2010;29(8):1203-13.
- 1385 30. Roz L, Gramegna M, Ishii H, Croce CM, Sozzi G. Restoration of
1386 fragile histidine triad (FHIT) expression induces apoptosis and suppresses
1387 tumorigenicity in lung and cervical cancer cell lines. *Proc Natl Acad Sci U*
1388 *S A*. 2002;99(6):3615-20.
- 1389 31. Beroukhim R, Mermel CH, Porter D, Wei G, Raychaudhuri S,
1390 Donovan J, et al. The landscape of somatic copy-number alteration across
1391 human cancers. *Nature*. 2010;463(7283):899-905.
- 1392 32. Shull AY, Clendenning ML, Ghoshal-Gupta S, Farrell CL, Vangapandu
1393 HV, Dudas L, et al. Somatic mutations, allele loss, and DNA methylation
1394 of the Cub and Sushi Multiple Domains 1 (CSMD1) gene reveals
1395 association with early age of diagnosis in colorectal cancer patients. *PLoS*
1396 *one*. 2013;8(3):e58731.
- 1397 33. Wilhelm M, Schlegl J, Hahne H, Gholami AM, Lieberenz M, Savitski
1398 MM, et al. Mass-spectrometry-based draft of the human proteome.
1399 *Nature*. 2014;509(7502):582-7.

- 1400 34. Gross AM, Orosco RK, Shen JP, Egloff AM, Carter H, Hofree M, et al.
1401 Multi-tiered genomic analysis of head and neck cancer ties TP53 mutation
1402 to 3p loss. *Nat Genet.* 2014;46(9):939-43.
- 1403 35. Lo PH, Leung AC, Kwok CY, Cheung WS, Ko JM, Yang LC, et al.
1404 Identification of a tumor suppressive critical region mapping to 3p14.2 in
1405 esophageal squamous cell carcinoma and studies of a candidate tumor
1406 suppressor gene, ADAMTS9. *Oncogene.* 2007;26(1):148-57.
- 1407 36. Lung HL, Lo PH, Xie D, Apte SS, Cheung AK, Cheng Y, et al.
1408 Characterization of a novel epigenetically-silenced, growth-suppressive
1409 gene, ADAMTS9, and its association with lymph node metastases in
1410 nasopharyngeal carcinoma. *Int J Cancer.* 2008;123(2):401-8.
- 1411 37. Lo PH, Lung HL, Cheung AK, Apte SS, Chan KW, Kwong FM, et al.
1412 Extracellular protease ADAMTS9 suppresses esophageal and
1413 nasopharyngeal carcinoma tumor formation by inhibiting angiogenesis.
1414 *Cancer Res.* 2010;70(13):5567-76.
- 1415 38. Forbes SA, Beare D, Gunasekaran P, Leung K, Bindal N, Boutselakis
1416 H, et al. COSMIC: exploring the world's knowledge of somatic mutations
1417 in human cancer. *Nucleic Acids Research.* 2015;43(D1):D805-D11.
- 1418 39. Qi J, Yu Y, Akilli Öztürk Ö, Holland JD, Besser D, Fritzmann J, et al.
1419 New Wnt/ β -catenin target genes promote experimental metastasis and
1420 migration of colorectal cancer cells through different signals. *Gut.* 2015.
- 1421 40. Chung K-Y, Cheng IKC, Ching AKK, Chu J-H, Lai PBS, Wong N.
1422 Block of proliferation 1 (BOP1) plays an oncogenic role in hepatocellular

- 1423 carcinoma by promoting epithelial-to-mesenchymal transition.
1424 Hepatology. 2011;54(1):307-18.
- 1425 41. Reed AL, Califano J, Cairns P, Westra WH, Jones RM, Koch W, et al.
1426 High Frequency of p16 (CDKN2/MTS-1/INK4A) Inactivation in Head and
1427 Neck Squamous Cell Carcinoma. Cancer Research. 1996;56(16):3630-3.
- 1428 42. Cairns P, Polascik TJ, Eby Y, Tokino K, Califano J, Merlo A, et al.
1429 Frequency of homozygous deletion at p16/CDKN2 in primary human
1430 tumours. Nat Genet. 1995;11(2):210-2.
- 1431 43. Munro J, Stott FJ, Vousden KH, Peters G, Parkinson EK. Role of the
1432 alternative INK4A proteins in human keratinocyte senescence: evidence
1433 for the specific inactivation of p16INK4A upon immortalization. Cancer
1434 Res. 1999;59(11):2516-21.
- 1435 44. Castets M, Broutier L, Molin Y, Brevet M, Chazot G, Gadot N, et al.
1436 DCC constrains tumour progression via its dependence receptor activity.
1437 Nature. 2012;482(7386):534-7.
- 1438 45. Escudero-Esparza A, Bartoschek M, Gialeli C, Okroj M, Owen S,
1439 Jirstrom K, et al. Complement inhibitor CSMD1 acts as tumor suppressor
1440 in human breast cancer. Oncotarget. 2016;7(47):76920-33.
- 1441 46. Tang MR, Wang YX, Guo S, Han SY, Wang D. CSMD1 exhibits
1442 antitumor activity in A375 melanoma cells through activation of the Smad
1443 pathway. Apoptosis. 2012;17(9):927-37.
- 1444 47. Liesenfeld M, Mosig S, Funke H, Jansen L, Runnebaum IB, Dürst M,
1445 et al. SORBS2 and TLR3 induce premature senescence in primary human
1446 fibroblasts and keratinocytes. BMC Cancer. 2013;13(1):1-11.

- 1447 48. Patel K, Scrimieri F, Ghosh S, Zhong J, Kim M-S, Ren YR, et al.
1448 FAM190A Deficiency Creates a Cell Division Defect. *The American Journal*
1449 *of Pathology*. 2013;183(1):296-303.
- 1450 49. Bashaw GJ, Klein R. Signaling from Axon Guidance Receptors. *Cold*
1451 *Spring Harbor Perspectives in Biology*. 2010;2(5).
- 1452 50. Huang M, Anand S, Murphy EA, Desgrosellier JS, Stupack DG,
1453 Shattil SJ, et al. EGFR-dependent pancreatic carcinoma cell metastasis
1454 through Rap1 activation. *Oncogene*. 2012;31(22):2783-93.
- 1455 51. Errington TM, Macara IG. Depletion of the adaptor protein NCK
1456 increases UV-induced p53 phosphorylation and promotes apoptosis. *PloS*
1457 *one*. 2013;8(9):e76204.
- 1458 52. Kresty LA, Mallery SR, Knobloch TJ, Song H, Lloyd M, Casto BC, et
1459 al. Alterations of p16(INK4a) and p14(ARF) in patients with severe oral
1460 epithelial dysplasia. *Cancer Res*. 2002;62(18):5295-300.
- 1461 53. Malkoski SP, Wang X-J. Two sides of the story? Smad4 loss in
1462 pancreatic cancer versus head-and-neck cancer. *FEBS Letters*.
1463 2012;586(14):1984-92.
- 1464 54. Huntley SP, Davies M, Matthews JB, Thomas G, Marshall J, Robinson
1465 CM, et al. Attenuated type II TGF- β receptor signalling in human
1466 malignant oral keratinocytes induces a less differentiated and more
1467 aggressive phenotype that is associated with metastatic dissemination.
1468 *International Journal of Cancer*. 2004;110(2):170-6.
- 1469 55. Alani RM, Hasskarl J, Grace M, Hernandez M-C, Israel MA, Munger
1470 K. Immortalization of primary human keratinocytes by the helix-loop-

- 1471 helix protein, Id-1. Proceedings of the National Academy of Sciences.
1472 1999;96(17):9637-41.
- 1473 56. Lin J, Guan Z, Wang C, Feng L, Zheng Y, Caicedo E, et al. Inhibitor
1474 of Differentiation 1 Contributes to Head and Neck Squamous Cell
1475 Carcinoma Survival via the NF- κ B/Survivin and Phosphoinositide 3-
1476 Kinase/Akt Signaling Pathways. American Association for Cancer
1477 Research. 2010;16(1):77-87.
- 1478 57. Burns JE, Clark LJ, Yeudall WA, Mitchell R, Mackenzie K, Chang SE,
1479 et al. The p53 status of cultured human premalignant oral keratinocytes.
1480 Br J Cancer. 1994;70(4):591-5.
- 1481 58. Sun Q, Wang R, Luo J, Wang P, Xiong S, Liu M, et al. Notch1
1482 promotes hepatitis B virus X protein-induced hepatocarcinogenesis via
1483 Wnt/beta-catenin pathway. Int J Oncol. 2014;45(4):1638-48.
- 1484 59. Lin JT, Chen MK, Yeh KT, Chang CS, Chang TH, Lin CY, et al.
1485 Association of high levels of Jagged-1 and Notch-1 expression with poor
1486 prognosis in head and neck cancer. Ann Surg Oncol. 2010;17(11):2976-
1487 83.
- 1488 60. Song X, Xia R, Li J, Long Z, Ren H, Chen W, et al. Common and
1489 complex Notch1 mutations in Chinese oral squamous cell carcinoma. Clin
1490 Cancer Res. 2014;20(3):701-10.
- 1491 61. Lefort K, Mandinova A, Ostano P, Kolev V, Calpini V, Kolfschoten I,
1492 et al. Notch1 is a p53 target gene involved in human keratinocyte tumor
1493 suppression through negative regulation of ROCK1/2 and MRCKalpha
1494 kinases. Genes Dev. 2007;21(5):562-77.

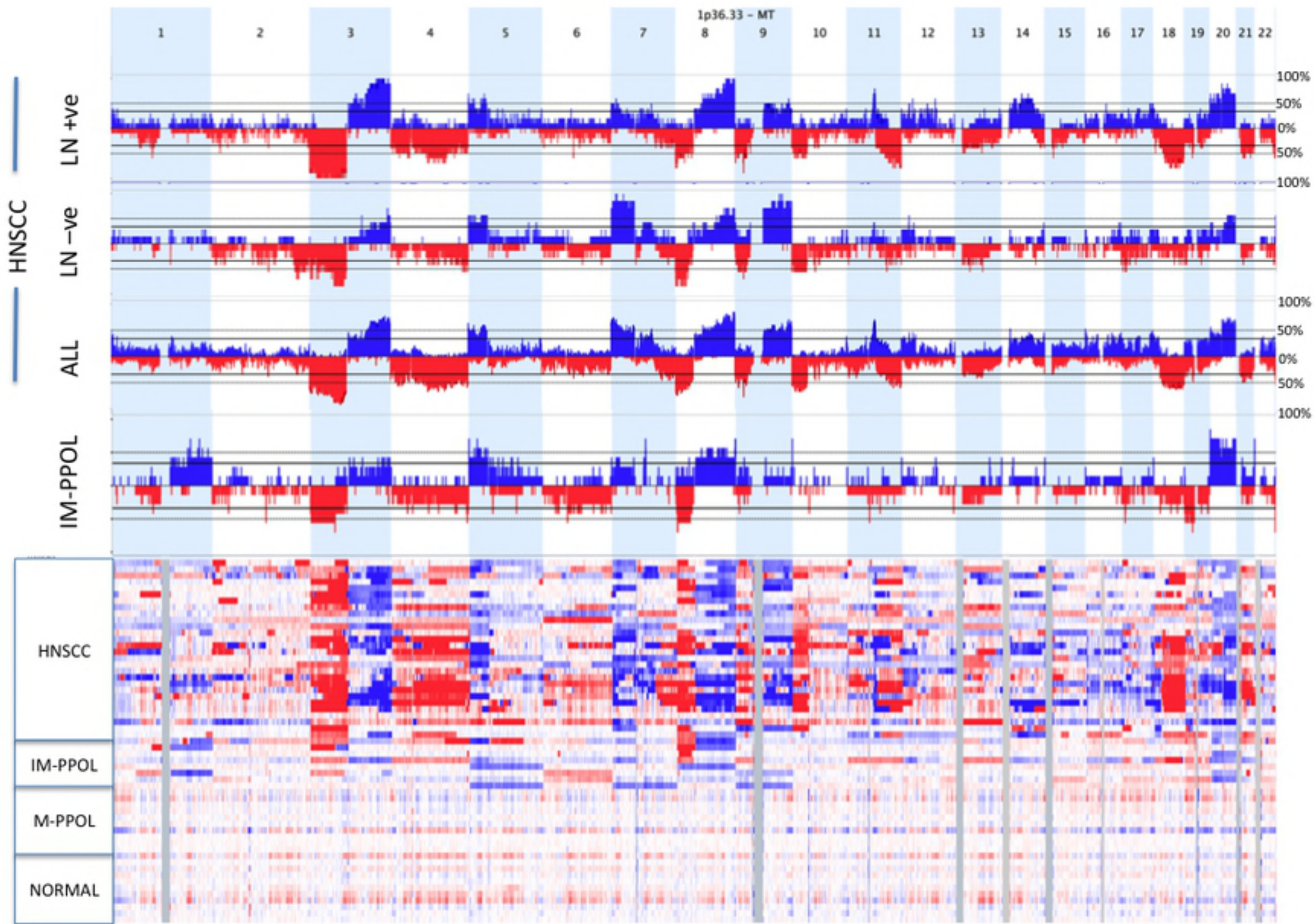
- 1495 62. Sakamoto K, Fujii T, Kawachi H, Miki Y, Omura K, Morita K, et al.
1496 Reduction of NOTCH1 expression pertains to maturation abnormalities of
1497 keratinocytes in squamous neoplasms. *Lab Invest.* 2012;92(5):688-702.
- 1498 63. Nickoloff BJ, Qin JZ, Chaturvedi V, Denning MF, Bonish B, Miele L.
1499 Jagged-1 mediated activation of notch signaling induces complete
1500 maturation of human keratinocytes through NF-kappaB and PPARgamma.
1501 *Cell Death Differ.* 2002;9(8):842-55.
- 1502 64. Lowell S, Jones P, Le Roux I, Dunne J, Watt FM. Stimulation of
1503 human epidermal differentiation by delta-notch signalling at the
1504 boundaries of stem-cell clusters. *Curr Biol.* 2000;10(9):491-500.
- 1505 65. Demehri S, Turkoz A, Kopan R. Epidermal Notch1 loss promotes
1506 skin tumorigenesis by impacting the stromal microenvironment. *Cancer*
1507 *Cell.* 2009;16(1):55-66.
- 1508 66. Patel V, Rosenfeldt HM, Lyons R, Servitja J-M, Bustelo XR, Siroff M,
1509 et al. Persistent activation of Rac1 in squamous carcinomas of the head
1510 and neck: evidence for an EGFR/Vav2 signaling axis involved in cell
1511 invasion. *Carcinogenesis.* 2007;28(6):1145-52.
- 1512 67. Cerami E, Gao J, Dogrusoz U, Gross BE, Sumer SO, Aksoy BA, et al.
1513 The cBio Cancer Genomics Portal: An Open Platform for Exploring
1514 Multidimensional Cancer Genomics Data. *Cancer Discovery.*
1515 2012;2(5):401-4.
- 1516 68. Silvera D, Arju R, Darvishian F, Levine PH, Zolfaghari L, Goldberg J,
1517 et al. Essential role for eIF4GI overexpression in the pathogenesis of
1518 inflammatory breast cancer. *Nat Cell Biol.* 2009;11(7):903-8.

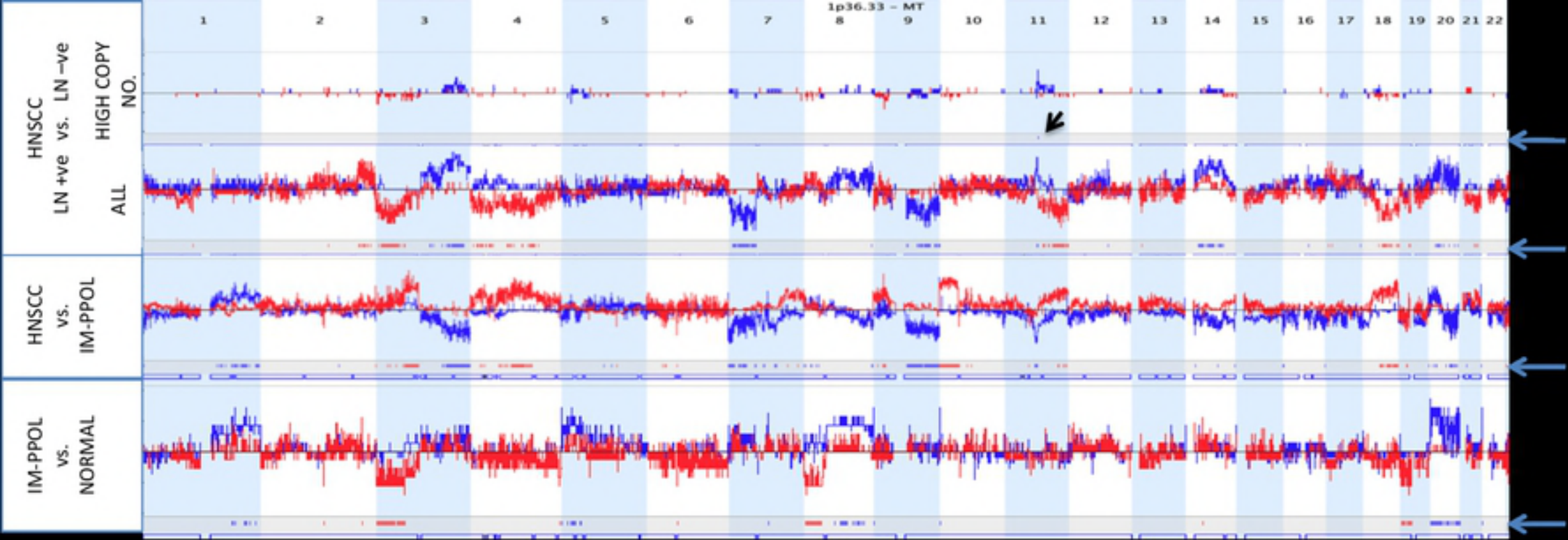
- 1519 69. Lin C-J, Cencic R, Mills JR, Robert F, Pelletier J. c-Myc and eIF4F Are
1520 Components of a Feedforward Loop that Links Transcription and
1521 Translation. *Cancer Research*. 2008;68(13):5326-34.
- 1522 70. Krauthammer M, Kong Y, Ha BH, Evans P, Bacchiocchi A, McCusker
1523 JP, et al. Exome sequencing identifies recurrent somatic RAC1 mutations
1524 in melanoma. *Nat Genet*. 2012;44(9):1006-14.
- 1525 71. Wallingford JB, Habas R. The developmental biology of Dishevelled:
1526 an enigmatic protein governing cell fate and cell polarity. *Development*.
1527 2005;132(20):4421-36.
- 1528 72. Rudy SF, Brenner JC, Harris JL, Liu J, Che J, Scott MV, et al. In vivo
1529 Wnt pathway inhibition of human squamous cell carcinoma growth and
1530 metastasis in the chick chorioallantoic model. *Journal of Otolaryngology -
1531 Head & Neck Surgery*. 2016;45(1):1-8.
- 1532 73. Yau C, Mouradov D, Jorissen RN, Colella S, Mirza G, Steers G, et al.
1533 A statistical approach for detecting genomic aberrations in heterogeneous
1534 tumor samples from single nucleotide polymorphism genotyping data.
1535 *Genome Biol*. 2010;11(9):R92.

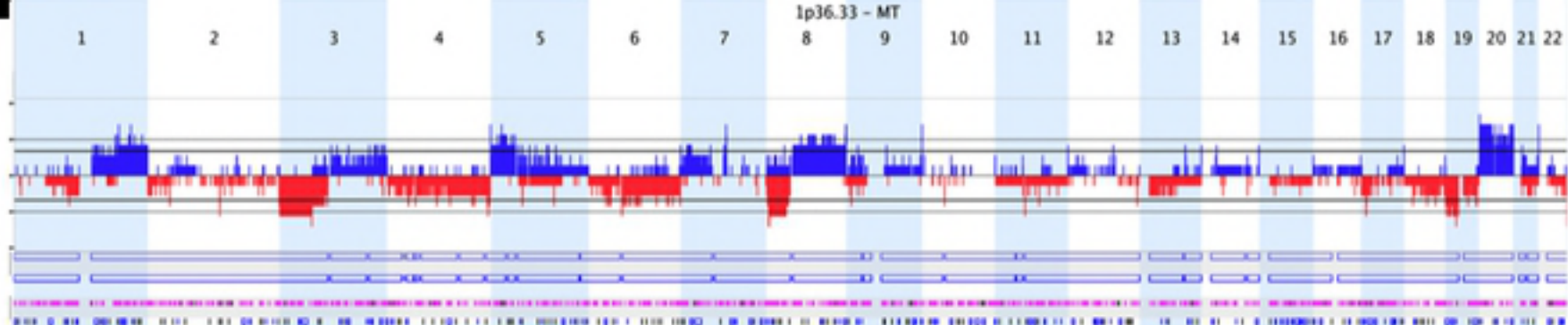
1536

RANK	GENE	M-PPOL			IM-PPOL						HNSCC											
		E2	E4	E5	NP			P			M	IM										
					D9	D38 ⁺	D4 ⁺	D34 ⁺	D20 ⁺	D35	D19	BICR80	BICR16	BICR10	BICR68	BICR82	H413	BICR63	BICR78	BICR22 ⁺	BICR37	BICR18
1	TP53						Red	Yellow		Red			Red		Yellow		Red		Yellow	Red		Red
2	CDKN2A								Red				Red				Red					Yellow
3	PIK3CA						Yellow															
4	CASP8							Red							Red							
5	FBXW7	Blue					Yellow								Blue				Blue			
6	NFE2L2					Yellow			Yellow													
8	FAT1													Blue	Red				Red			
9	NOTCH1								Red	Yellow		Red					Blue					
10	KMT2D (MLL2)					Red			Red						Blue							Red
13	ZNF750																					Red
18	ARID2																					Red
20	KDM6A*						Blue		Blue											Blue		Blue
21	FUBP1																			Blue		
25	MACF1																					Red
28	KALRN	Blue																				
31	NCOR1					Blue																
33	SMARCA4											Red										Blue
35	HLA-A									Red												Red
37	WHSC1								Yellow													
57	ATRX						Red		Blue													
61	NF1					Blue		Blue		Blue						Blue						Blue
62	ATM						Red															
68	KDM5C																			Red		
83	ATR						Yellow															
85	TRIO								Blue													
88	BAZ2B								Yellow													
89	BRCA1																					Blue
95	GNAS					Red														Yellow	Red	
99	CCAR1						Red													Red		
100	NCKAP1						Blue															
108	CDK12																					Blue
110	BNC2					Red			Red													Blue
116	DDX3X										Yellow											
122	ELF1								Blue													
125	CASP1																			Red		Yellow
128	MGA*					Blue			Blue													Blue
137	ATIC																			Yellow		
141	BRWD1					Yellow																
150	GPSM2								Yellow													
158	NEDD4L																			Blue		
160	NUP107*					Red	Red	Red	Red	Red										Red		Red
164	NR4A2*					Yellow		Yellow														Yellow
167	TCF4																					Blue

bioRxiv preprint doi: <https://doi.org/10.1101/364205>; this version posted July 9, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.







D19

CIS

D35

CIS

D20

Moderate

D34

Moderate

D4

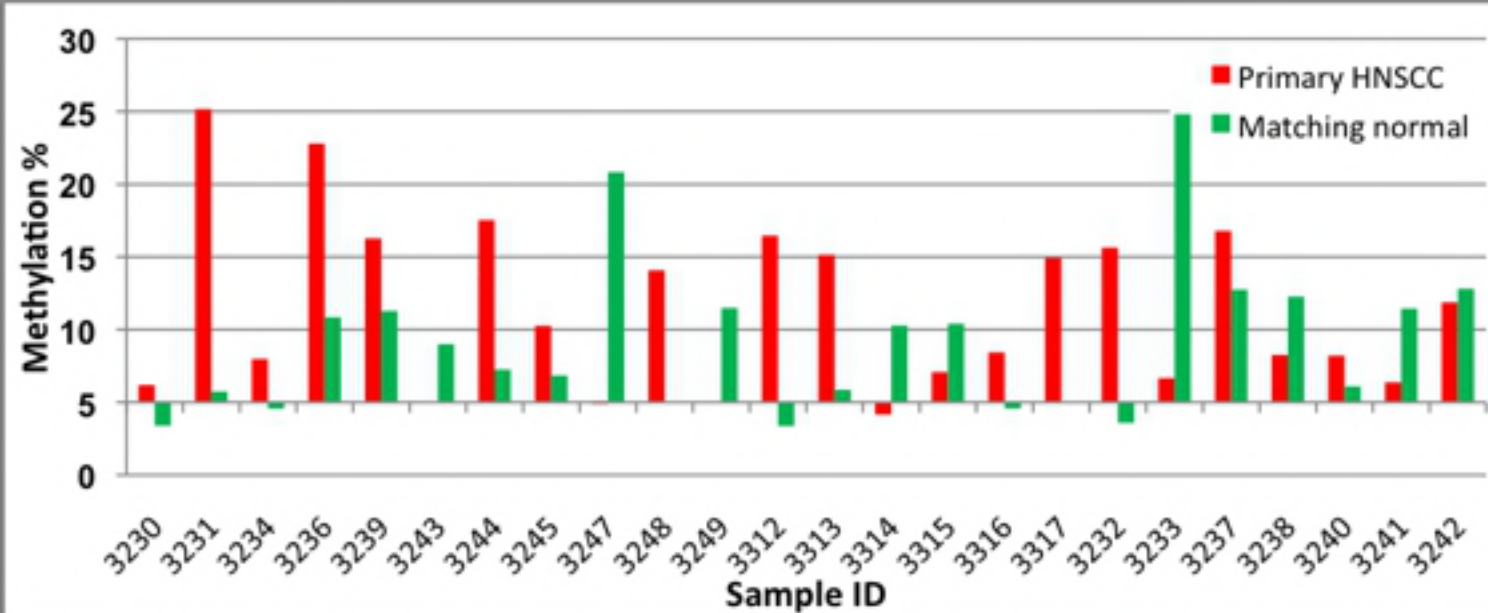
CIS

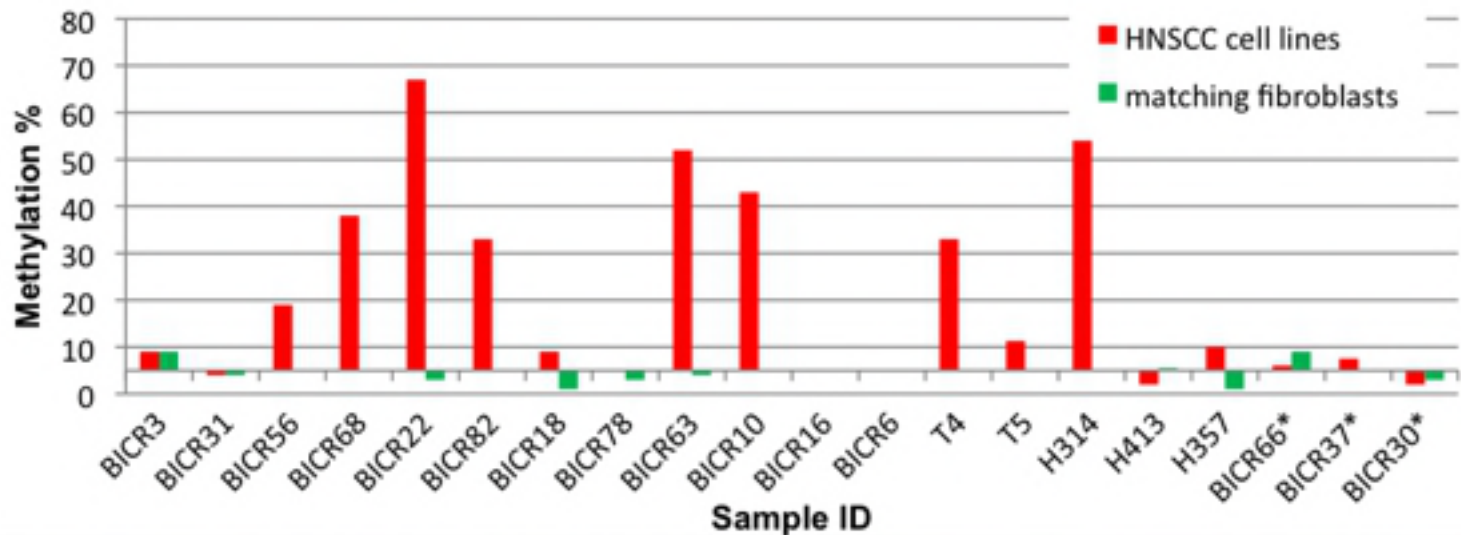
D9

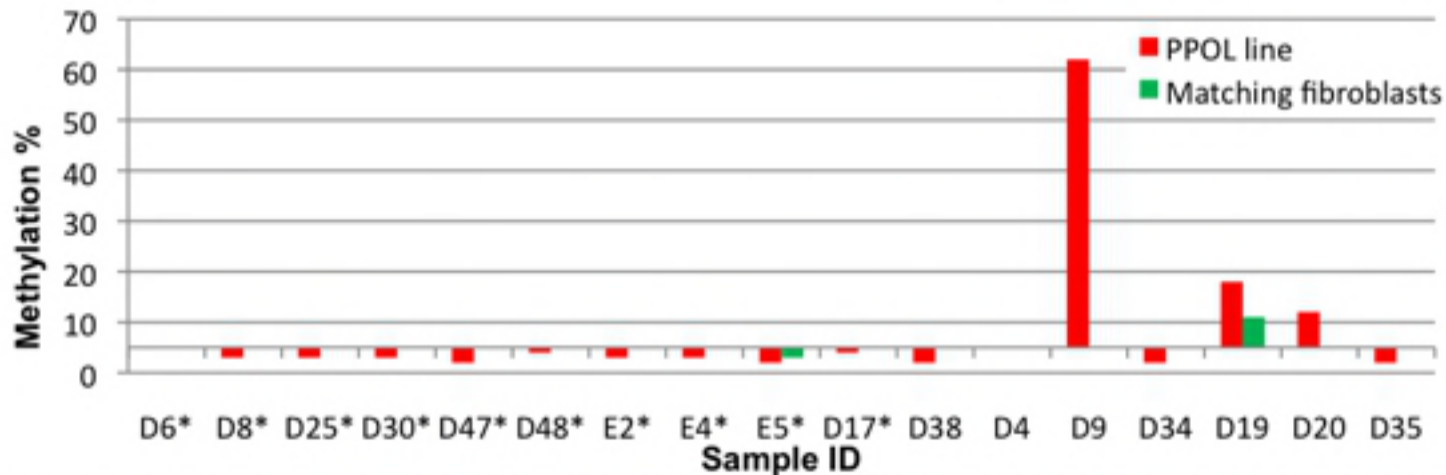
Mild

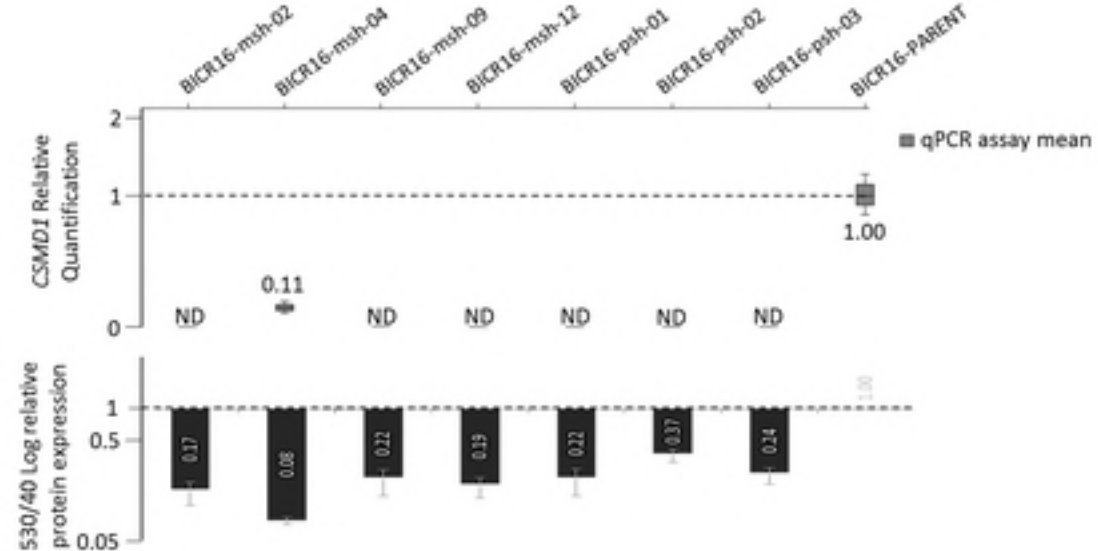
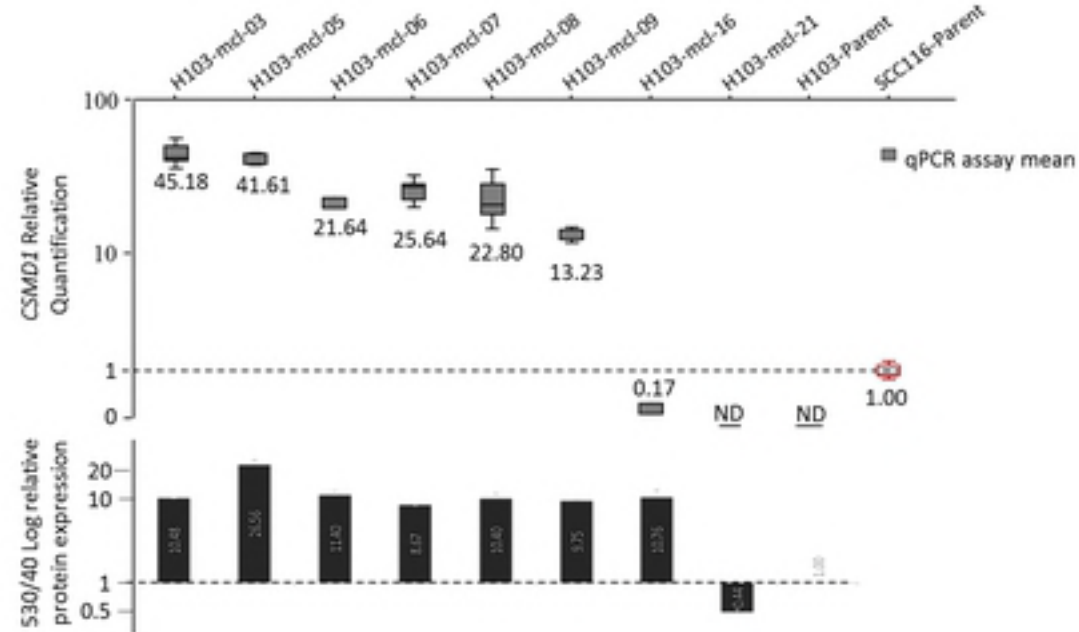
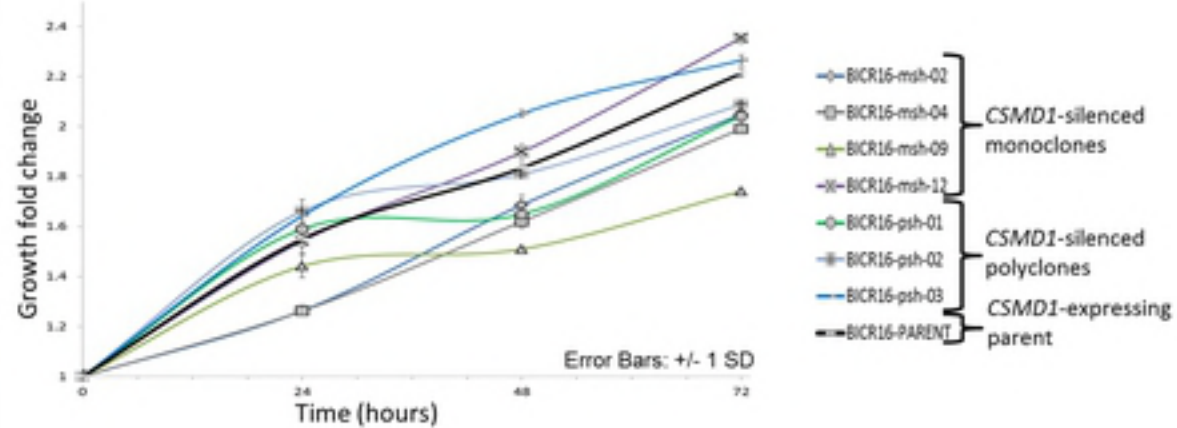
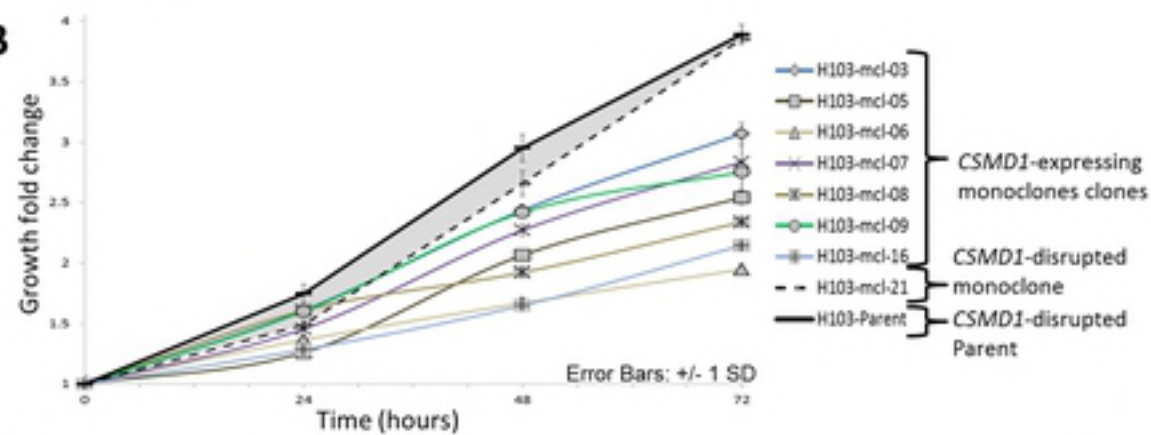
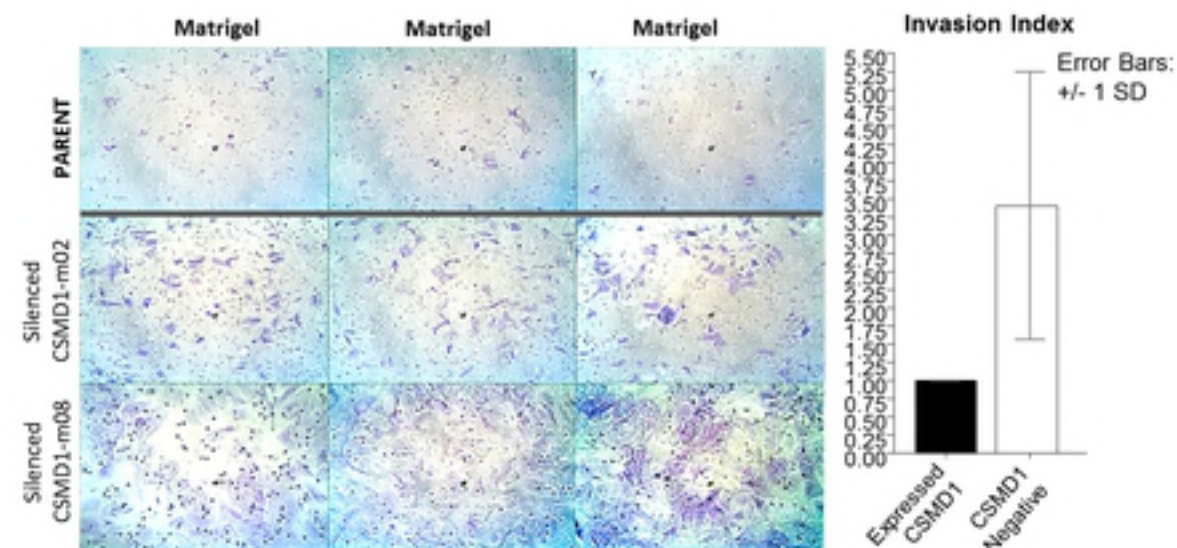
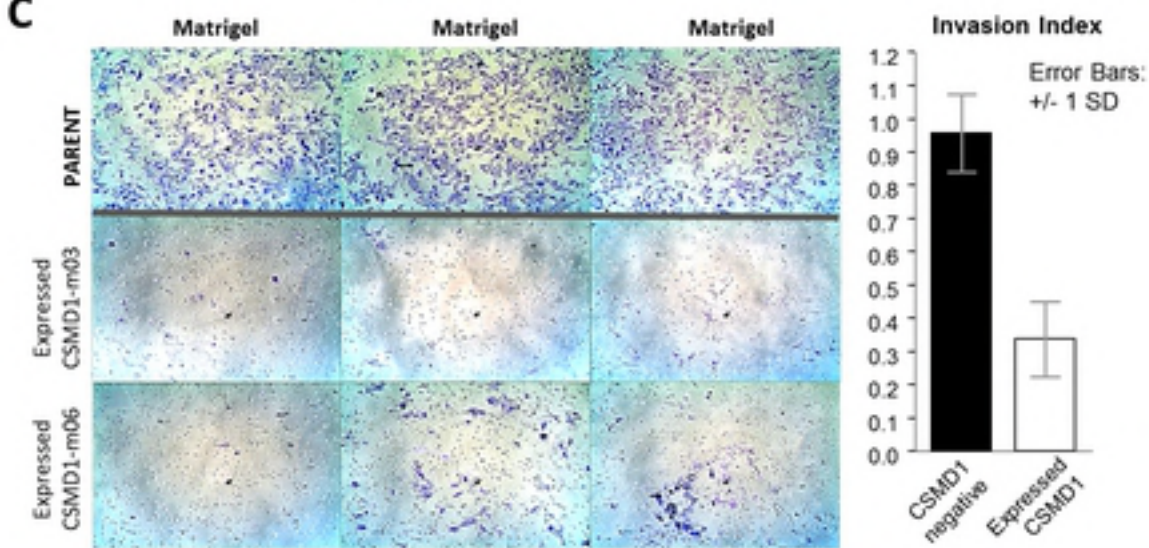
D38

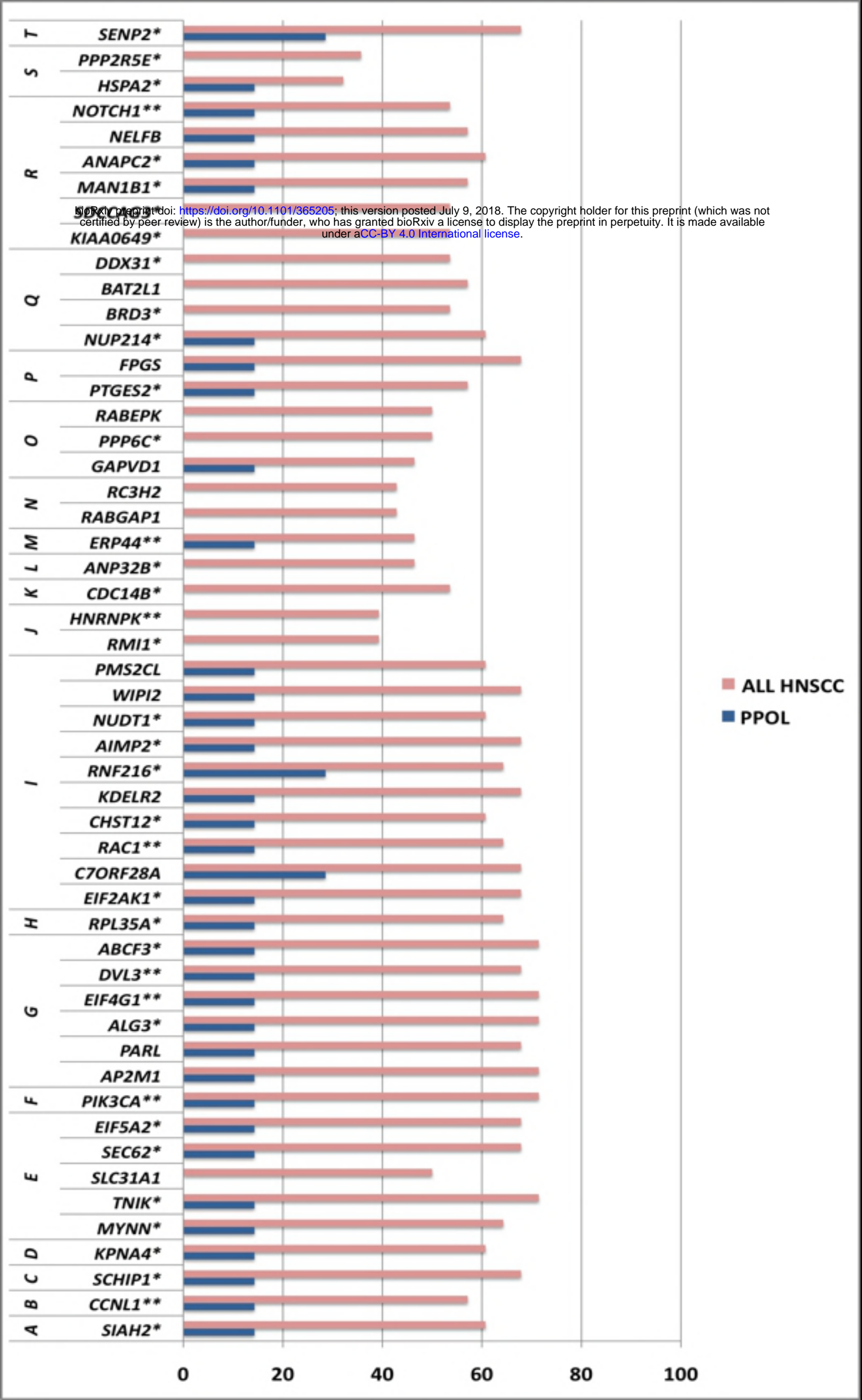
Mild

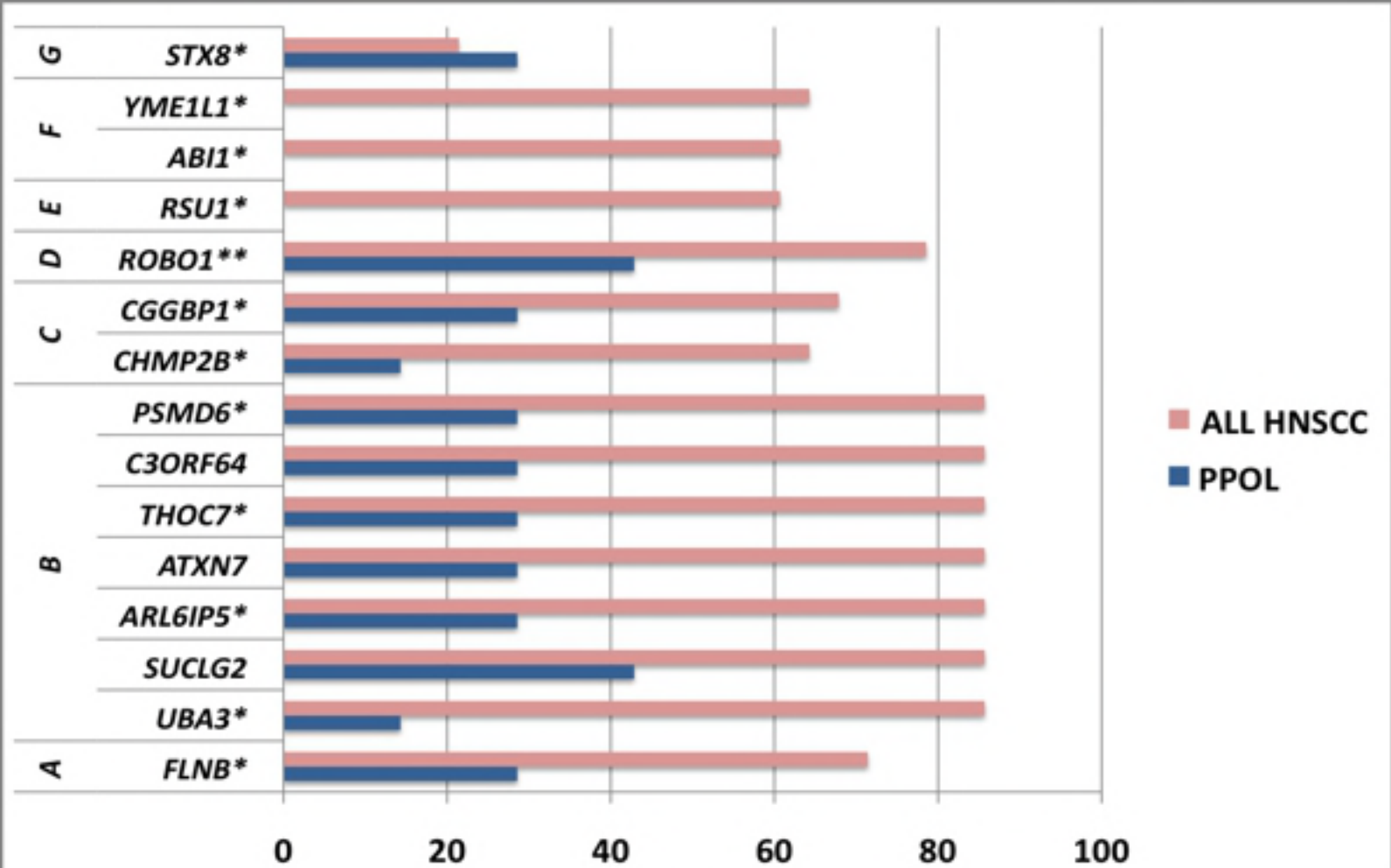






A**B****C**





KEGG PATHWAY

PPOLS

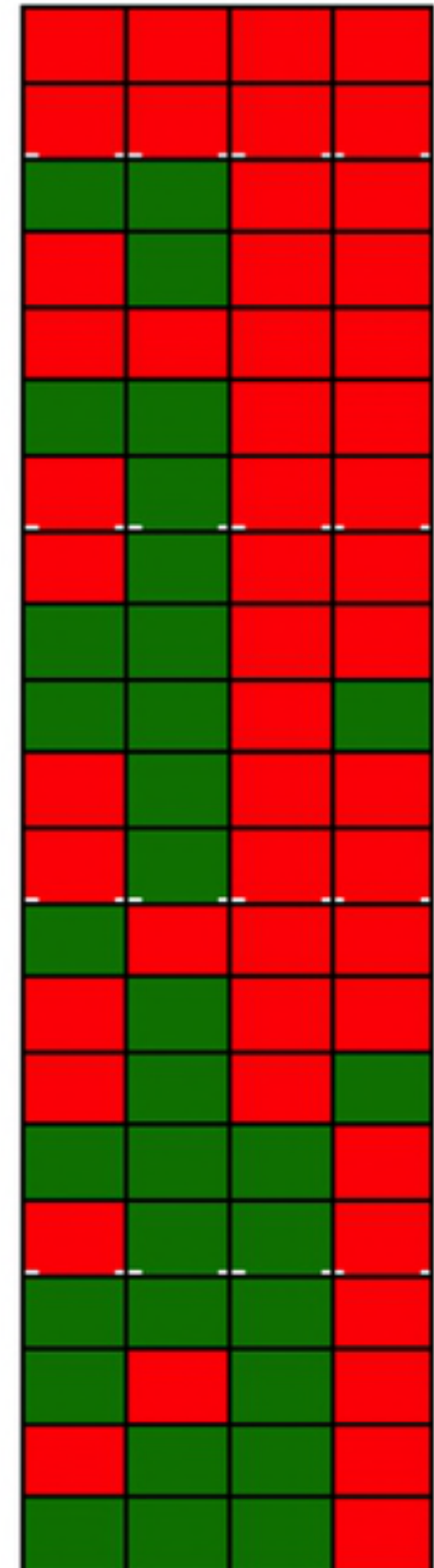
LN-VE HNSCC

ALL HNSCC

LN+VE HNSCC

Pathways in cancer
 Endocytosis
 Ubiquitin mediated proteolysis
 Wnt signaling pathway
 Focal adhesion
 Apoptosis
 Regulation of actin cytoskeleton
 MAPK signaling pathway
 Cell adhesion molecules (CAMs)
 mTOR signalling pathway
 Axon guidance
 p53 signaling pathway
 Cell cycle
 TGF-beta signaling pathway
 Erb signalling pathway
 Jak-STAT signaling pathway
 Toll-like receptor signaling pathway
 Regulation of autophagy
 Notch signaling pathway
 PPAR signaling pathway
 Mismatch repair

bioRxiv preprint doi: <https://doi.org/10.1101/365205>; this version posted July 9, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



KEGG CANCER-SPECIFIC PATHWAY

Small cell lung cancer
 Colorectal cancer
 Pancreatic cancer
 Melanoma
 Non-small cell lung cancer
 Endometrial cancer
 Glioma
 Bladder cancer
 Prostate cancer

