1  **Integrated analysis of oral tongue squamous cell carcinoma identifies key**

2  **variants and pathways linked to risk habits, HPV, clinical parameters and tumor**

3  **recurrence**

4

5  Neeraja M Krishnan[1], Saurabh Gupta[1], Vinayak Palve[1], Linu Varghese[1], Swetansu

6  Pattnaik[1], Prachi Jain[1], Costerwell Khyriem[1], Arun K Hariharan[1], Kunal Dhas[1],

7  Jayalakshmi Nair[1], Manisha Pareek[1], Venkatesh K Prasad[1], Gangotri Siddappa[2],

8  Amritha Suresh[2], Vikram D Kekatpure[3], Moni Abraham Kuriakose[2,3] and Binay

9  Panda[1,4] [*]

10

11  [1]Ganit Labs, Bio-IT Centre, Institute of Bioinformatics and Applied Biotechnology,

12  Bangalore, India

13  [2]Integrated Head and Neck Oncology Program, Mazumdar Shaw Centre for

14  Translational Research, Bangalore, India

15  [3]Head and Neck Oncology, Mazumdar Shaw Medical Centre, Bangalore, India

16  [4]Strand Life Sciences, Bangalore, India

17

18  *corresponding author: binay@ganitlabs.in

19

20  **Abstract**

21  Oral tongue squamous cell carcinomas (OTSCC) are a homogenous group of

22  tumors characterized by aggressive behavior, early spread to lymph nodes and a

23  higher rate of regional failure. Additionally, the incidence of OTSCC among younger

24  population (<50yrs) is on a rise; many of who lack the typical associated risk factors

25  of alcohol and/or tobacco exposure. We present data on SNVs, indels, regions with

26  LOH, and CNVs from fifty-paired oral tongue primary tumors and link the significant

27  somatic variants with clinical parameters, epidemiological factors including HPV

28  infection and tumor recurrence. Apart from the frequent somatic variants harbored in

29  *TP53, CASP8, RASA1, NOTCH* and *CDKN2A* genes, significant amplifications and/or

30    deletions were detected in chromosomes 6-9, and 11 in the tumors. Variants in *CASP8*

31    and *CDKN2A* were mutually exclusive. *CDKN2A, PIK3CA, RASA1* and *DMD*

32    variants were exclusively linked to smoking, chewing, HPV infection and tumor

33    stage. We also performed whole-genome gene expression study that identified matrix

34    metalloproteases to be highly expressed in tumors and linked pathways involving

35    arachidonic acid and NF-κ-B to habits  and distant metastasis, respectively. Functional

36    knockdown studies in cell lines demonstrated the role of *CASP8* in HPV-negative

37    OTSCC cell line. Finally, we identified a 38-gene minimal signature that predicts

38    tumor recurrence using an ensemble machine learning method. Taken together, this

39    study links molecular signatures to various clinical and epidemiological factors in a

40    homogeneous tumor population with a relatively high HPV prevalence.

41

42    **Keywords:** oral tongue squamous cell carcinoma, somatic variants, gene expression,

43    CNV, LOH, tumor recurrence, *CASP8, TP53, RASA1*, HPV

44

45    **Introduction**

46          Squamous cell carcinomas of head and neck (HNSCC) are the sixth leading

47    cause of cancer worldwide [1]. Tumors of head and neck region are heterogeneous in

48    nature with different incidences, mortalities and prognosis for different subsites and

49    accounts for almost 30% of all cancer cases in India [2]. Oral cancer is the most

50    common subtype of head and neck cancers in humans, with a worldwide incidence in

51    >300,000 cases. The disease is an important cause of death and morbidity, with a 5-

52    year survival of less than 50% [1, 2]. Recent studies have identified various genetic

53    changes in many subsites of head and neck using high-throughput sequencing assays

54    and computational methods [3-7]. Such multi-tiered approaches using the exomes,

55    genomes, transcriptomes and methylomes from different squamous cell carcinomas

56    have generated data on key variants and in some cases, their biological significance,

57    aiding our understanding of disease progression. Some of the above sequencing

58    studies have identified key somatic variants and linked them with patient stratification

59    and prognostication. This, along with the associated epidemiology, enables one to

60    look beyond the discovery of driver mutations, and identify predictive signatures in

61    HNSCC.

62        A previous study from the cancer genome atlas (TCGA) consortium with

63    HNSCC patients (N = 279) identified somatic mutations in *TP53, CDKN2A, FAT1,*

64    *PIK3CA, NOTCH1, KMT2D* and *NSD1* at a frequency greater than 10% [7].

65    Additionally, the TCGA study identified loss of *TRAF3* gene, amplification of *E2F1*

66    in human papilloma virus (HPV)-positive oropharyngeal tumors, along with

67    mutations in *PIK3CA*, *CASP8* and *HRAS*, and co-amplifications of the regions 11q13

68    (harboring *CCND1*, *FADD* and *CTTN*) and 11q22 (harboring *BIRC2* and *YAP1*), in

69    HPV-negative tumors, described to play an important role in pathogenesis and tumor

70    development [7]. Chromosomal losses at 3p and 8p, and gains at 3q, 5p and 8q were

71    also observed in HNSCC [7]. Tumors originating in the anterior/oral part of tongue

72    or, oral tongue squamous cell carcinoma (OTSCC) tend to be different from those at

73    other subsites as oral tongue tumors are associated more with younger patients [8] and

74    spread early to lymph nodes [9]. Additionally oral tongue tumors have a higher

75    regional failure compared to gingivo-buccal cases [10] in oral cavity. Tobacco (both

76    chewing and smoking) and alcohol are common risk factors for this group of tumors

77    among older patients [8]. The role of HPV, both as an etiological agent and/or risk

78    factor along with its role as a good prognostic marker in OTSCC, unlike in

79    oropharyngeal tumors, is currently uncertain. It remains to be explored whether HPV

80    acts as an etiological agent in the development of oral tongue tumors or simply

81    represents a super infection in patients. Additionally, HPV infection status currently

82    does not influence disease management in OTSCC.

83         Here, we present data towards a comprehensive molecular characterization of

84    OTSCC. We performed exome sequencing, whole-genome gene expression, and

85    genotyping arrays using fifty primary tumors along with their matched control

86    samples, towards identification of somatic variants (mutations and indels),

87    significantly up- and down-regulated genes, loss of heterozygosity (LOH) and copy

88    number variations (CNVs). We integrated all the molecular data along with the

89    clinical parameters and epidemiology such as tumor stage, nodal status, HPV

90    infection, risk habits and tumor recurrence to interpret the effect of changes in the

91    process of cancer development in oral tongue. We identified significant somatic

92    variations in *TP53* (38%), *RASA1* (8%), *CASP8* (8%), *CDKN2A* (6%), *NOTCH1*

93    (4%), *NOTCH2* (4%), and *PIK3CA* (4%) from the exome sequencing study in

94    OTSCC. The key variants were validated using an additional set of primary tumor

95    samples. Variants in *TP53* and *NOTCH1* were found in mutually exclusive sets of

96    tumors. Additionally, we found frequent aberrations in chromosomes 6-9, and 11 in

97    tumor samples. We observed a strong association between somatic variations in some

98    key genes with one or more risk habits; for example, *CDKN2A* and *PIK3CA* with

99    smoking; *CASP8* with consuming alcohol and chewing tobacco; *RASA1* with chewing

100   and tumor stage, and HPV infection, along with *DMD* and *PIK3CA*. From the gene

101   expression analysis, we found matrix metalloproteases (MMPs) to be highly

102   expressed in OTSCC. Pathway analysis identified Procaspase-8, Notch, Wnt, p53,

103   extracellular matrix (ECM)-receptor interaction, JAK-STAT and PPAR to be some of

104   the significantly altered pathways in OTSCC. We implemented an ensemble machine

105    learning method [11] and identified a minimal gene signature set that distinguished a

106    group of tumors with loco-regional recurrence from the non-recurrent set. Finally, we

107    performed functional analysis of *CASP8* gene in HPV-negative and HPV-positive

108    OTSCC cell lines to establish its role in the process of tumor development.

109

110    **Results**

111    *Habits, clinical parameters and epidemiology*

112            We collected tumor and matched control (adjacent normal and/or

113    lymphocytes) samples from 50 patients diagnosed with OTSCC, with informed

114    consent. Data from patient habits, epidemiology and clinical parameters are presented

115    in Figure 1A and Additional file 1A. About two-thirds of the patients (N = 31)

116    included in our study were in the younger age group (≤50yrs), with 20% female

117    patients in the total pool. Approximately, 70% of the patients were positive for at least

118    one risk habit, namely, smoking, alcohol consumption or chewing tobacco (33% of

119    patients smoked tobacco, 40% consumed alcohol and 42% chewed tobacco). HPV

120    infection status in the primary tumors was established with type-specific qPCR or

121    HPV16 digital PCR. Thirty-three percentage of the patients were deceased at the time

122    of completing the analysis. About 60% of the tumors were moderately differentiated,

123    25% well differentiated and the rest were poorly differentiated. Among the patients

124    recruited, 60% were node-positive, 70% had no recurrence, 9% had distant metastasis

125    and 24% had loco-regional recurrence at the time of completing the analysis. The

126    mean and median follow up duration for patients were nearly 30 months and 21

127    months respectively. About 27% of the tumors were early stage tumors (T1N0M0 and

128    T2N0M0) and the rest 73% were late stage tumors (tumors belonging to the rest of the

129    TNM stage).

130

*Discovery and validation of significant somatic variants and their relationship with*

*other parameters*

We re-discovered variants, as described previously [12] using whole-genome

arrays, to validate the variant call accuracy as obtained from the exome sequencing

data. We validated ~99% of the SNPs discovered from Illumina sequencing in both

the tumor and matched control samples (Additional file 2). After filtering and

annotation, we identified 19 cancer-associated genes bearing significantly altered

somatic variants in OTSCC (Figure 1D). These were validated using Sanger

sequencing in two sets of samples, one using the same tumor-control pairs used in the

exome sequencing (the discovery set, Additional file 1A) and second, using an

additional 36-60 primary tumors (validation set, Additional file 1B) for genes altered

in ≥ 5% of the tumor samples. All the *TP53* variants were validated in the discovery

set. Three out of the four variants were validated for *CASP8*. The mutant alleles for

the heterozygous variants in *HLA-A*, *OBSCN*, *ING1*, *TTK* and *U2AF1* discovered by

exome sequencing were difficult to interpret from the results of the validation using

Sanger sequencing as they were present at a very low frequency (Additional file 3).

Combining data from the validation set; the mutation frequencies for *RASA1* and

*CDKN2A* rose significantly to 10.71% and 16.47% in primary tumors respectively but

those for *TP53* and *CASP8* remained largely unchanged (Additional file 3).

The somatic mutation frequency per MB ranged from 10-45 with a median

around 25 (Figure 1B). The median value for transition to transversion (ti/tv) ratio for

both the tumor and its matched control samples was ~2.5 (Additional file 4). Overall,

*T->C* changes were most frequent, followed by *G->A* and then *T->G*. Habits

(smoking and alcohol consumption), nodal status, HPV infection, tumor grade and

155  stage had no significant impact on the distribution of these nucleotides (Additional

156  file 5). We used the workflow described in the *Methods* section to identify somatic

157  mutations and indels in tumor samples following which we used three functional

158  tools, IntOGen [19], MutSigCV [21] and MuSiC2 [22] for variant interpretations

159  (Additional file 6). In order to identify genes harboring significant variants, we used

160  the intersection of these tools, following the criteria that the somatic variants be

161  callable in the matched control sample and present in a single sequencing read in the

162  control sample. This resulted in a final list of 19 cancer-associated genes (Figure 1C),

163  which were divided into three categories with varying mutation frequencies (Figure

164  1D). The three frequency tiers were ≥ 30% (*TP53*), 6-30% (*RASA1*, *CASP8* and

165  *CDKN2A*) and 2-5% (*NOTCH1*, *NOTCH2*, *DMD* and *PIK3CA* were prominent

166  among them).

167      Next, we looked for mutual exclusivity of finding somatic variants in the

168  genes and found that many of these genes harbor variants in a mutually exclusive

169  manner across samples (Figure 1E), suggesting the possibility that there might be

170  some common pathway(s) involved in the development of OTSCC. We observed

171  mutual exclusivity among somatic variants in *NOTCH1* and *NOTCH2* genes, and

172  expanded this finding to identifying 15 such mutually exclusive sets (Figure 1E).

173  Among them, *CDKN2A*, *HLA-A* and *TTK* form a mutually exclusive set with *TP53;*

174  *RASA1*, *OBSCN*, *HLA-A*, *AJUBA* and *TTK* are mutually exclusive with either

175  *NOTCH1* alone, or *NOTCH2* and *ANK3* together; *NOTCH1*, *NOTCH2*, *HLA-A*,

176  *AJUBA*, *ANK3*, *TTK*, *MLL2*, *ING1* or *KEAP1*, are mutually exclusive with *CASP8*

177  alone, or *FAT1* and *DMD* together; *FAT1*, *HLA-A*, *AJUBA*, *ANK3*, *TTK*, *MLL2*, *ING1*

178  or *KEAP1*, are mutually exclusive with *PIK3CA* or *DMD* or *NOTCH1* and *OBSCN*, or

179  *CDKN2A* and *OBSCN*; *U2AF1*, *MLL2* and *TTK* form a small mutually exclusive set.

180 We juxtaposed the positions of the somatic variants from final list of all 19 genes

181 (Additional file 7) detected in OTSCC against those found in the TCGA data using

182 the cBioPortal. We found that the somatic variants in OTSCC were in the same

183 domains where mutations were observed earlier in many of the genes (Additional file

184 7).

185  Copy Number Variation (CNV) analyses using data from the whole-genome

186 SNP genotyping arrays revealed a large chunk of chromosome 9, bearing cancer-

187 associated genes like *CDKN2A, NF1* and *MRPL4*, to be affected in about 17% of the

188 tumors (Figure 1F and Additional file 8). We found several CNVs of short stretches

189 (in low kb range) within chromosomes 6-8, 11, 17 and X in many tumors.

190

191 *Linking habits, HPV infection, nodal status, tumor grade and recurrence, with genes*

192 *harboring somatic variants and the associated pathways*

193  We further classified the 19 cancer-associated genes from the previous

194 analyses and linked those with habits, clinical parameters and HPV infection. Among

195 the genes harboring significant somatic variants, we found *CDKN2A* to be mutated

196 only in the never smokers and past smokers, *PIK3CA* to be mutated only in the

197 smokers, and *TP53* to be mutated at a 20% greater frequency in the smokers, *CASP8*

198 has a 12% greater frequency in those that consumed alcohol or chewed tobacco.

199 *RASA1* was exclusively mutated only in the non-chewers. (Figure 2A). HPV-negative

200 patients harbored somatic variants in *DMD* and *PIK3CA*, while HPV-positive patients

201 alone had somatic variants in *RASA1*. Only the moderate- and well-differentiated

202 tumor samples harbored variants in *CASP8*, while *NOTCH1* was mutated largely in

203 the poorly-differentiated tumors. Node-positive tumors had a 19% greater occurrence

204 of *TP53* variants. Somatic variants in *RASA1* occurred exclusively in the late stage

205    tumors (Figure 2A). We further studied the association of affected cancer-related

206    signaling pathways with habits and clinical parameters, and found that recurrence and

207    HPV infection had the highest impact (Figure 2B). The Procaspase-8 activation,

208    Notch, p53 and Wnt signaling pathways were linked most with many of the clinical

209    parameters, HPV infection and habits (Figure 2B).

210

211    *Differentially expressed genes in OTSCC*

212         Significant ($q$ val ≤ 0.05) differentially expressed genes with a fold change of

213    at least 1.5 revealed a consistent pattern of differential expression across the tumor

214    samples (21 up- and 23 down-regulated genes, Figure 3A and Additional file 9).

215    Genes in PPAR signaling- (e.g., *MMP1*) and ECM-receptor interaction pathways

216    (*LAMC2* and *SPP1*) were up-regulated and *CRNN, APOD, SCARA5* and *RERGL* were

217    down-regulated in a majority of tumors (Figure 3A). Next, we studied the pathways

218    involving genes with aberrant expression and their link with HPV infection and other

219    clinical parameters. Genes in the arachidonic acid metabolism and Toll-like receptors

220    were differentially expressed in patients with no smoking history (never smokers or

221    past smokers) and alcohol habits (Figure 3B). *SERPINE1* (a gene in HIF-1 signaling

222    pathway) was differentially expressed in patients that are habits-negative. The NF-κ-B

223    signaling pathway was differentially expressed only in metastasized tumors.

224

225    *Functional studies with CASP8 in OTSCC cell lines*

226         *CASP8* is mutated in a significant number of oral tongue tumors [this study, 5,

227    7]. Caspase-8 is an important and versatile protein that plays a role in both apoptotic

228    (extrinsic or death receptor-mediated) and non-apoptotic processes [13, 14]. We

229    studied the functional consequences of *CASP8* knockdown through a siRNA-

230    mediated method in an HPV-positive UM:SCC-47 [15] and an HPV-negative

231    UPCI:SCC040 [16] OTSCC cell lines. Prior to the functional assay, the concentration

232    of siRNA required for silencing, extent of *CASP8* knockdown and cisplatin sensitivity

233    (IC$_{50}$) in both these cell lines was tested (Additional file 10). The invasion of cells was

234    greater in both UM:SCC-47 and UPCI:SCC040 cell lines when *CASP8* was knocked

235    down (Figure 4A). To analyze the effect of caspase-8 on the migration property of

236    cells, scratches were made on the confluent monolayer of cells and the wound closure

237    area was measured at different time points (0hr, 15hr, 23hr & 42hr, Figure 4B). The

238    wound closure was faster in *CASP8* knockdown HPV-negative cells compared to the

239    HPV-positive cells. At 15hr, 23hr and 48hrs, about 65%, 90% and 100% of the

240    wound got closed respectively in the HPV-negative cell line compared to 50%, 70%

241    and 85% respectively during the same time period in the HPV-positive cell lines

242    (Figure 4B). siRNA knockdown of *CASP8* rescued the chemo-sensitivity caused by

243    cisplatin treatment as evident by the MTT (3-(4,5-dimethylthiazol-2-yl)-2,5-

244    diphenyltetrazolium bromide) survival assay (Figure 4C). Interestingly, we found that

245    the extent of rescue is greater in the HPV-negative cell line compared to the HPV16-

246    positive one.

247    *Tumor recurrence prediction using random forests*

248         After cataloging the significantly altered genes in OTSCC, we wanted to see

249    whether there is a relationship between the altered genes and loco-regional recurrence

250    of tumors and metastasis. In order to do this, we used an ensemble machine learning

251    method implemented by variable elimination using random forests [11] (Figure 5).

252    We used multiple testing correction and the 0.632 bootstrapping method [17] to

253    estimate false positives. We discovered a 38-gene minimal signature that

254    discriminated between the non-recurring, loco-regionally recurring and distant

255    metastatic tumors (Figure 5). The .632+ bootstrap errors, indicative of prediction

256    specificity, varied across non-recurrent, recurrent and distant metastatic tumors. The

257    median error was low (0.03) and intermediate (0.3) for the non-recurrent and the loco-

258    regionally recurrent categories respectively but was relatively higher (1.0) for the

259    metastatic tumors. The errors were proportional to the number of representative

260    samples within each category.

261

262    *Major signaling pathways implicated in OTSCC*

263          We looked at significant pathways altered in OTSCC, taking into account all

264    the molecular changes in tumors and found apoptosis, HIF, Notch, mTOR, p53,

265    PI3K/Akt, Wnt and Ras to be some of the key signaling pathways affected in OTSCC

266    (Figure 6). In addition, histone methylation, cell cycle/immunity and mRNA splicing

267    processes were also affected. The complete list of pathways is provided in Additional

268    file 11.

269

270    **Discussion**

271          Squamous cell carcinomas of oral tongue are an aggressive group of tumors

272    with a higher incidence in the younger population (≤50yrs), which spread early to

273    lymph nodes and have a higher regional failure compared to gingivo-buccal cases [8-

274    10]. Previous sequencing studies [3-5, 7] grouped oral tongue tumors with tumors

275    from oral cavity but a rise in the incidence of oral tongue tumors, especially among

276    younger people who never smoked, consumed alcohol or chewed tobacco, warrants

277    further investigation of this subgroup of oral tumors. Additionally, the role of HPV in

278    oral tongue tumors, unlike in oropharyngeal cases [18-20], is not well understood both

279    in terms of incidence and prognosis. A meta-analysis of HPV-positive HNSCC

280  tumors from multiple studies conducted at multiple locations concluded that HPV-

281  positive patients, especially in oropharynx, have improved overall and disease-

282  specific survival [21]. A past study has presented data that the HPV incidence in oral

283  tongue is low [22] and some argue against any link between HPV infection and

284  aggressive oral tongue tumors [23]. Although there is no consensus on rate of HPV

285  incidence among oral tongue patients, it is generally believed that it is low compared

286  to oropharyngeal tumors. However, some studies in the past [24], albeit from a

287  different geography, established a much higher rate of HPV infection in oral tongue

288  tumors.

289         We applied stringent filtering steps and used multiple annotation tools to come

290  up with a list of 19 cancer-associated genes that harbored somatic variants in OTSCC.

291  Most of these genes were also found in other studies, including the recent TCGA

292  HNSCC study [7], with some notable differences. A comparison of somatic variants

293  discovered in all HSNCC studies, including the current study, is provided in

294  Additional file 12. The frequency for somatic changes in *CASP8*, *NOTCH1*, *CDKN2A*

295  and *FAT1* genes in previous studies [3-7] were, 4-34%, 13-18%, 2-16% and 13-50%,

296  respectively. This is different from what we found in the current study (8%, 4%, 6%

297  and 2% respectively for the same genes). This may partly be attributed to the total

298  number of tumor samples used in different studies but may also be due to a unique

299  pattern of mutations specific to oral tongue subsite. It appears from our study that the

300  latter is the case. For example, in one of the earlier studies [6] involving similar

301  number of patients as in the current study, *CASP8* and *FAT1* were mutated in 34%

302  and 50% of the patients but we find the frequency to be 8% and 2%, respectively. In

303  some earlier studies, it was not possible to categorize and identify oral tongue-specific

304  variants as the sites were classified under oral cavity [3, 4].

305    Although the somatic variants discovered from our study appear to be

306    distributed uniformly across the genome, the significant copy number variation events

307    are more concentrated in chromosomes 6-9 and 11 (Figure 1F and Additional file 8).

308    One of the most important genes harboring somatic mutations discovered in our study

309    is *CASP8*, the product for which derived from the precursor Procaspase-8. Caspase-

310    8,is an important protein implicated in both apoptotic and non-apoptotic pathways

311    [14]. Recent analysis from the TCGA study [7] suggests that mutations in *CASP8* co-

312    occur with mutations in *HRAS*, and are mutually exclusive with amplifications in

313    *FADD* gene . In our functional studies, the most important observation was that

314    caspase-8 shows different effects in HPV-positive and HPV-negative cells, the effect

315    being more pronounced in HPV-negative cells (Figure 4). Therefore, it is possible that

316    HPV-negative tumors activate a completely different set(s) of pathways and/or may

317    have different chemosensitivity towards drugs than the HPV-positive tumors. It was

318    shown previously that HPV-positive HNSCC cell lines are resistant to TRAIL (tumor

319    necrosis factor-related apoptosis-inducing ligand) and treatment of cells with the

320    proteasome inhibitor bortezomib sensitizes HPV-positive cells towards TRAIL-

321    induced cell death mediated by caspase-8 [25]. The E6 protein of HPV interacts with

322    the DED domain of caspase-8 and induces its activation by recruiting it to the nucleus

323    [26]. Our observation on the role of caspase-8-mediated apoptosis being more

324    pronounced in the HPV-negative OTSCC cell line is similar to the observation on the

325    role of *CASP8* in HPV-negative patients made earlier in TCGA study [7]. Taken

326    together, genes including *CASP*8 regulate key pathways (Figure 6) that might play

327    important role in the development of tumors in oral tongue.

328    In the past, several large sequencing studies have been undertaken in HNSCC

329    [3-5, 7] that contained very few HPV-positive oral tongue patients. Our study is based

330    on a unique patient cohort and attempts to link molecular signature with different

331    clinical and epidemiological parameters. The prevalence of HPV is very high in oral

332    tongue tumors from India, including in our cohort, compared with studies using

333    cohorts elsewhere. We don't see the same high prevalence of HPV in non-oral tongue

334    tumors in oral cavity, for example in buccal tumors, in one of our study (Palve et al.,

335    *unpublished observation*). The exact reason for this high prevalence is not know.

336    Additionally, the HPV positive patients that harbor mutations in *TP53* in some of the

337    patient samples is also counter-intuitive owing to the fact that E6 is known to block

338    p53. Although we don't know the reason behind this, there is a possibility that HPV

339    positive tumors harboring TP53 mutations represent a unique class of tumors and it

340    will be interesting to see if those tumors recur early or late compared to the HPV

341    positive tumors that have wild type p53 function. Therefore, this study is unique in

342    that respect.

343        Identifying signature for tumor recurrence prospectively in primary tumors

344    may add significant advantage to disease management. In order to do this, we used a

345    machine learning method using the molecular changes identified in this study, in three

346    batches of primary tumors; non-recurring, loco-regionally recurring and tumors with

347    distant metastasis. We identified a 38-gene signature to be significantly distinguishing

348    the three groups. The bootstrapping error for the non-recurring and the loco-regionally

349    recurring groups were low (N = 34, *.632 error* = 0.03 and N = 10, *.632 error* = 0.3

350    respectively) but not in the metastatic tumor group (N = 4, *.632 error = 1*). This was

351    due to the small sample numbers (N = 4) in the metastatic category, justifying the

352    need for a larger sample set to validate the signature. The 38 gene signature identified

353    in out study, however, needs to be validated in a much larger cohort in the future to

354    achieve its true potential as a prognostic panel in OTSCC.

355        Finally, we were keen to see if the current study leads to finding novel drug

356    candidates in OTSCC. We based our assumption on the fact that genome-wide

357    somatic variant discovery in tumors may give rise to possibilities of finding novel

358    drug targets/candidates or may led us to use existing drugs prescribed for other

359    indications. In an attempt to identify if any of the significantly altered genes found in

360    the current study could potentially act as drug targets, we screened for available drugs

361    against them. We found drugs against three targets out of which two have undergone

362    at least one clinical trial (Additional file 13).

363

364    **Methods**

365    *Informed consent, ethics approval and patient samples used in the study*

366        Informed consent was obtained voluntarily from each patient enrolled in the

367    study. Ethics approval was obtained from the Institutional Ethics Committees of the

368    Mazumdar Shaw Medical Centre. Matched control (blood and/or adjacent normal

369    tissue) and tumor specimens were collected and used in the study. Only those patients,

370    where the histological sections confirmed the presence of squamous cell carcinoma

371    with at least 70% tumor cells in the specimen, were used in the current study. At the

372    time of admission, patients were asked about the habits (chewing, smoking and/or

373    alcohol consumption). Fifty treatment-naïve patients who underwent staging

374    according to AJCC criteria, and curative intent treatment as per NCCN guideline

375    involving surgery with or without post-operative adjuvant radiation or chemo-

376    radiation at the Mazumdar Shaw Medical Centre were accrued for the study

377    (Additional file 1). Post-treatment surveillance was carried out by clinical and

378    radiographic examinations as per the NCCN guidelines.

379

380    *HPV detection*

381         HPV was detected by using q-PCR using HPV16- and HPV18-specific

382    TaqMan probes and primers and digital PCR using TaqMan probes and primers to

383    detect HPV in primary tumor samples.

384    *Exome Sequencing, read QC, alignment, variant discovery and post-processing filters*

385         Exome libraries were prepared using Agilent SureSelect, Illumina TruSeq and

386    Nextera exome capture kits (Additional file 14) following manufacturers'

387    specifications. Paired end sequencing was performed using HiSeq 2500 or GAIIx and

388    raw reads were generated using standard Illumina base caller. Read pairs were filtered

389    using *in house* scripts (Additional files 15 and 16) and only those reads having ≥75%

390    bases with ≥ 20 phred score and ≤ 15 Ns were used for sequence alignment against

391    human hg19 reference genome using NovoAlign [27] (v3.00.05). The aligned files

392    (*.sam) were processed using Samtools [28](v0.1.12a) and only uniquely mapped

393    reads from NovoAlign were considered for variant calling. The alignments were pre-

394    processed using GATK [29] (v1.2-62) in three steps before variant calling. First, the

395    indels were realigned using the known indels from 1000G (phase1) data. Second,

396    duplicates were removed using Picard (v1.39). Third, base quality recalibration was

397    done using CountCovariates and TableRecalibration from GATK (v1.2-62), taking

398    into account known SNPs and indels from dbSNP (build 138). Finally,

399    UnifiedGenotyper from GATK (v2.5-2) was used for variant calling, using known

400    SNPs and indels from dbSNP (build 138). Raw variants from GATK were filtered to

401    only include the PASS variants (standard call confidence ≥ 50) within the merged

402    exomic bait boundaries. Two out of 50 tumor samples did not confirm to the QC

403    standards, therefore excluded from all further analyses. Therefore, all the downstream

404    analyses were restricted to 48 primary tumors. The variants were further flagged as

405     novel or present in either dbSNP138 or COSMIC (v67) databases, based on their

406     overlap. In addition to GATK, we also used Dindel [30] to call indels. Both GATK

407     and Dindel calls were filtered for microsatellite repeats (flagged as STR). The raw

408     variant calls were used to estimate frequencies of nucleotide changes and

409     transition:transversion (ti/tv) ratios. Exome-filtered PASS variants specific to the

410     tumor samples, with respect to both location and actual call, were retained as somatic

411     variants, which were further filtered to exclude variants where the region bearing the

412     variant was not callable in the matched control sample, and those where the matched

413     control sample had even one read covering the variant allele.

414       Scripts used to perform various filtering steps are provided in Additional file

415     16. The numbers of raw reads, after QC, alignment statistics, numbers of variants pre-

416     and post-filters are provided in Additional file 2.

417

418     *Detection of cross-contamination and identification of significant somatic variants*

419       We estimated cross-contamination using ContEst [31] in the tumor samples

420     (Additional file 16). Locus-wise and gene-wise driver scores were estimated by

421     CRAVAT [32] using the head and neck cancer database with the CHASM [33]

422     analysis option. Genes with a CHASM score of at least 0.35 were considered

423     significant for comparison with other functional analyses (Additional file 16).

424     Somatic mutations were normalized with respect to the exome bait size (MB) to

425     calculate the somatic mutation frequency per MB.

426     *Annotation and functional analyses of variants*

427       Annotation and functional analyses of somatic variants was performed using

428     IntoGen [34, 35](web version 2.4), MutSigCV [36, 37] and MuSiC2 [37]. Somatic

429     variants, filtered to contain only those callable in the matched normal but not covered

430     by any read in the control samples (VCF), were used for IntoGen with the 'cohort

431     analyses' option. We also ran MutsigCV1.3 with these variants using coverage from

432     un-filtered variants of all tumor samples (Additional file 16). Pooled alignments for

433     all normal and tumor samples (bam), each, along with pooled variants for all normal

434     samples (MAF) were analyzed using MuSiC2 to calculate the background mutation

435     rates (bmrs) for all genes, and identify a list of significantly mutated genes ($p$-value of

436     convolution test $\leq 0.05$; Additional file 16). A condensed list of 19 genes, common

437     between at least two analyses was compiled (Figure 1D).

438

439     *SNP genotyping and validation using Illumina whole-genome Omni LCG arrays*

440          High quality DNA (200ng), quantified by Qubit (Invitrogen), was used as the

441     starting material for whole-genome genotyping experiments following the

442     manufacturer's specifications. Briefly, the genomic DNA was denatured at room

443     temperature (RT) for 10 mins using 0.1N NaOH, neutralized and used for whole

444     genome amplification (WGA) under isothermal conditions, at $37^0$C for 20 hrs. Post

445     WGA, the DNA was enzymatically fragmented at $37^0$C for 1hr. The fragmented DNA

446     was precipitated with isopropanol at $4^0$C and resuspended in hybridization buffer. The

447     samples were then denatured at $95^0$C for 20 mins, cooled at RT for 30 mins and $35\mu$l

448     of each sample was loaded onto the Illumina HumanOmni 2.5-8 beadchip for

449     hybridization (20hrs at $48^0$C) in a hybridization chamber. The unhybridized probes

450     were washed away and the Chips (Human Omni2.5-8 v1.0 and v1.1, Additional file 2)

451     were prepared for staining, single base extension and scanning using Illumina's

452     HiScan system.

453          We filtered the SNP locations to retain only those, called without any error,

454     contained within the exome boundaries as per the sequencing baits, and which were

455     callable (covered by at least five sequencing reads). At these locations, we estimated

456     the overlap for individual SNP calls, i.e., chr/pos/ref/alt and for no calls; i.e.,

457     chr/pos/ref/ref; between sequencing and array platforms (Additional file 16).

458

459     *Discovering Copy number Variations (CNVs) and Loss of Heterozygosity (LOH)*

460            CNVs and LOHs were identified using cnvPartition 3.1.6 plugin in Illumina

461     GenomeStudio v2011.1, with default settings except for a minimum coverage of at

462     least 10 probes per CNV/LOH with a confidence score threshhold of at least 100

463     (Additional file 17). Somatic CNVs and LOHs were extracted by filtering out any

464     region common to CNVs and LOHs detected in its matched control. Somatic CNVs

465     and LOHs were further filtered with respect to common and disease-related CNVs

466     and LOHs using CNVAnnotator [38]. Overlaps with common CNVs and LOHs were

467     discarded, reporting only the overlaps with disease-related, and novel CNVs and

468     LOHs. We categorized the CNVs and LOHs within each cytoband and reported those

469     with an occurrence in at least 10% of the patient samples.

470

471     *Gene Expression Assay*

472            Gene expression profiling was carried out using Illumina HumanHT-12 v4

473     expression BeadChip (Illumina, San Diego, CA) in tumor and matched normal tissues

474     (Additional file 9) following manufacturer's specifications. Total RNA was extracted

475     from 20mg of tissue using PureLink RNA (Invitrogen) and RNeasy (Qiagen) Mini

476     kits. RNA quality was checked using Agilent Bioanalyzer using RNA Nano6000 chip.

477     Samples with poor RIN numbers, indicating partial degradation of RNA, were

478     processed using Illumina WGDASL assay as per manufacturer's recommendations.

479     The RNA samples with no degradation were labelled using Illumina TotalPrep RNA

480    Amplification kit (Ambion) and processed according to the array manufacturer's

481    recommendations. Gene expression data was collected using Illumina's HiScan and

482    analyzed with the GenomeStudio (v2011.1 Gene Expression module 1.9.0) and all

483    assay controls were checked to ensure quality of the assay and chip scanning. Raw

484    signal intensities were exported from GenomeStudio for pre-processing and analyzed

485    using R further.

486    Gene-wise expression intensities for tumor and matched control samples from

487    GenomeStudio were transformed and normalized using VST (Variance Stabilizing

488    Transformation) and LOESS methods, respectively, using the R package lumi [39].

489    The data was further batch-corrected using ComBat [40] (Additional file 16). The pre-

490    processed intensities for tumor and matched control samples were subjected to

491    differential expression analyses using the R package, limma [41] (Additional file 16).

492    Genes with significant expression changes (adjusted $P$ value $<= 0.05$) and fold change

493    of at least 1.5 were followed up with further functional analyses.

494

495    *Recurrence prediction using random forests*

496    We used presence or absence of somatic mutations/indels data in the entire set

497    of genes for all the OTSCC patients, along with their recurrence patterns as training

498    set for the random forests [11] analyses using the varSelRF package in R. This

499    method performs both backward elimination of variables and selection based on their

500    importance spectrum, and predicts recurrence patterns in the same set by iteratively

501    eliminating 2% of the least important predictive variables until the current OOB (out-

502    of-bag) error rate becomes larger than the initial or previous OOB error rates. In order

503    to understand the specificity of the best minimalistic predictors of tumor recurrence,

504    we estimated the 0.632+ error rate [17] over 50 bootstrap replicates. We used the

505     varSelRFBoot function from the varSelRF Bioconductor package to perform

506     bootstrapping. The .632+ method is described by the following formula:

507
$$Err^{0.632'} = Err^{0.632} + \frac{\left(Err^1 - err\right)\left(.368 \cdot .632 \cdot R'\right)}{1 - .368 \cdot R'}$$

508     where $Err^{(.632')}, Err^{(.632)}, Err^{(1)}$ and $err$ are errors estimated by the .632+ method, the

509     original .632 method, *leave-one-out* bootstrap method and *err* represents the error. *R'*

510     represents a value between 0 and 1. Another popular error correction method used is

511     *leave-one-out* bootstrap method. The .632+ method was designed to correct the

512     upward bias in the *leave-one-out* and the downward bias in the original .632 bootstrap

513     methods.

514         For all iterations of all random forest analyses, we confirmed that the variable

515     importance remained the same before and after correcting for multiple hypotheses

516     comparisons using pre- and post- Benjamin-Hochberg FDR-corrected *P* values (data

517     not shown).  R commands for variable elimination using random forests, 0.632+

518     bootstrapping and re-computing importance values after multiple comparisons testing

519     are provided in Additional file 16.

520

521     *Pathway analyses*

522         Consensus list of genes from analysis, filtering and annotation of variant calls

523     and from differential expression analysis using whole genome micro-arrays, were

524     mapped to pathways using the web version of Graphite Web [42] employing KEGG

525     and Reactome databases. The network of interactions between genes was drawn

526     originally using CytoScape [43] (v3.1.1) using the .sif file created by Graphite Web

527     (Additional file 16).

528

529     *Data Visualization*

530    We used Circos [44] (v0.66) (Additional files 18 and 19) for multi-

531    dimensional data visualization. Additionally, we used the cbioportal protal

532    (http://www.cbioportal.org/) to visualize variants within the 19 genes harboring

533    significant variants. All of the mandatory fields accepted by Mutation Mapper were

534    provided for select genes from our study to create structural representations for each

535    gene including domains. Such diagrams from our study, the HNSCC study and all

536    cancer studies from TCGA were collated using the image-editing tool, GIMP

537    (www.gimp.org). SNPs and indels were visualized for each individual tumor sample

538    using IGV [45] (v1.5.54), along with the reads supporting variants (Additional file

539    20).

540    *Validation of somatic variants using Sanger sequencing*

541    Primers were designed using the NCBI primer designing tool

542    (http://www.ncbi.nlm.nih.gov/tools/primer-blast/index.cgi?LINK_LOC=BlastHome)

543    and used in Sanger sequencing for validation. The sequences of all primers (IDT)

544    used for validation is provided in Additional file 21. We tested the specificity of the

545    designed primers using UCSC's tool, *In Silico* PCR. The variant-bearing region was

546    amplified by using specific primers and used in Sanger sequencing (Additional file

547    14). The somatic variants were confirmed by sequencing in the entire tumor and

548    matched control DNA set used for the exome sequencing followed by further

549    validation in 60 additional tumor samples (Additional file 1B).

550

551    *Cell culture and knockdown of CASP8 gene*

552    The human OTSCC cell lines UPCI:SCC040 (gift from Dr. Susan Gollin,

553    University of Pittsburgh, PA, USA) [16] and UM-SCC47 (gift from Dr. Thomas

554    Carey, University of Michigan, MI, USA) [15] were used in the study. All the cells

555    were maintained in Dulbecco's Modified Eagles' Media (DMEM) supplemented with

556    10% FBS, 1% MEM nonessential amino acids solution & 1% penicillin/streptomycin

557    mixture (Gibco) at 37°C with 5% $CO_2$ incubator.

558        We performed the siRNA-based knockdown using UPCI:SCC040 and

559    UM:SCC47 cell lines for *CASP8* gene. The expression of Caspase-8 was transiently

560    knocked down using ON-TARGETplus Human *CASP8* smart pool siRNA (L-

561    003466-00-0010; Dharmacon) along with an ON-TARGET plus Non-targeting

562    siRNA (D-001810-01-20; Dharmacon). The transfection efficiency for the two cell

563    lines (UPCI:SCC040 and UM:SCC47) were optimized using siGLO Red Transfection

564    Indicator (D-001630; Dharmacon). The siRNA duplexes were transfected using

565    Lipofectamine-2000 according to the manufacturer's instructions (Invitrogen). The

566    siRNA-oligo complexes medium was changed 8 hrs post transfection. The efficiency

567    of transfection along with the mRNA expression was analyzed at 24 and 48 hrs post

568    transfection by qRT-PCR. The specific down-regulation of *CASP8* was confirmed by

569    three independent experiments.

570

571    *RNA isolation and quantitative real-time PCR*

572        RNA was extracted from cell pellets and tissues using RNeasy Mini kit spin

573    columns (Qiagen) following manufacturer's protocol. Genomic DNA contamination

574    was removed by RNase-Free DNase Set (Qiagen) and the total RNA was eluted in

575    nuclease free water (Ambion). The RNA samples were estimated using Qubit 2.0

576    fluorometer (Invitrogen) and the integrity was checked by gel electrophoresis. The

577    RNA samples were stored at $-80^0$C until further used. The cDNA was synthesized

578    with 400ng total RNA, using a SuperScript-III first strand cDNA synthesis kit, and

579    following the manufacturer's instructions (Invitrogen). The cDNA was then subjected

580    for quantitative real-time PCR (q-RT-PCR) using KAPA SYBR FAST qPCR Master

581    Mix (KK4601, KAPA). The primer pairs used for testing the expression of caspase-8

582    in q-RT-PCR were, forward 5'-ATGATGACATGAACCTGCTGGA-3' and reverse

583    5'-CAGGCTCTTGTTGATTTGGGC-3'. The amplification was done on Stratagene

584    MX300P real time machine. To normalize inter-sample variation in RNA input, the

585    expression values were normalized with GAPDH. All amplification reactions were

586    done in triplicates, using nuclease free water as negative controls. The differential

587    gene expression was calculated by using the comparative $C_T$ method of relative

588    quantification [46].

589

590    *Assessment of cell viability (confirm)*

591         MTT cell proliferation assay was performed as per manufacturer's instructions

592    (Sigma) to assess cell viability. Briefly, cells were seeded on 96-well plates

593    containing DMEM with 10% FBS & incubated overnight. After treatment with 0.1%

594    DMSO (vehicle control), or Cisplatin for 48 hrs, medium was changed and 100 $\mu$l of

595    MTT solution (1mg/ml) was added to each well. The cells were further incubated for

596    4hrs at 37°C. The formazan crystals in viable cells were dissolved by adding 100$\mu$l of

597    dimethyl sulfoxide (DMSO) (Merck). The absorbance was recorded at 540 nm using

598    reference wavelength of 690 nm on micro plate reader (Tecan Systems). Data were

599    normalized to vehicle treatment, and the cell viability was calculated using GraphPad

600    Prism software (version 4.03; La Jolla, CA). All the experiments were performed in

601    triplicates.

602

603    *Wound healing assay*

604         Cells were cultured up to 80% confluency in 12 well plates; serum-starved for

605    24 hrs and then wounded using a $200\mu$l pipette tip. The wound was washed with 1x

606    PBS and the cells were grown in DMEM containing 10% FBS. Cells were imaged at

607    10x magnification at 0 hr, 15 hrs, 23 hrs and 42 hrs. For each well, three wounds were

608    made and the migration distance was photographed and measured using Carl Zeiss

609    software (Zeiss). Each experiment was performed in triplicates.

610

611    *Matrigel invasion assay*

612    The ECM gel (E1270, Sigma) was thawed overnight at 4ºC and plated at

613    requisite concentrations (for UPCI:SCC040: 1.5mg/ml and UMSCC047: 2mg/ml)

614    onto the transwell inserts and incubated overnight in the $CO_2$ incubator at 37ºC with

615    5% $CO_2$. Cells were serum-starved for overnight, harvested, counted and seeded

616    (UPCI:SCC040: 50,000 cells and UMSCC047: 20,000 cells per well) on top of the

617    matrigel transwell-inserts (2 mg/ml) in serum free medium as per manufacturer's

618    specifications (Sigma). D-MEM containing 10% FBS and 1% NEAA was added to

619    the lower chamber. The 24-well plates containing matrigel inserts with cells were

620    incubated in 37°C incubator for 48 hrs. At the end of incubation time, cells in the

621    upper chamber were removed with cotton swabs and cells that invaded the Matrigel to

622    the lower surface of the insert were fixed with 4% paraformaldehyde (Merk

623    Milipore), permeabilized with 100% methanol, stained with Giemsa (Sigma),

624    mounted on glass slides with DPX mounting agent and counted under a light

625    microscope (Zeiss). Each experiment was performed in triplicates.

626

627    **Conclusions**

628    We have catalogued genetic variants (somatic mutations, indels, CNVs and

629    LOHs) and transcriptomic (significantly up- and down-regulated genes) changes in

630 oral tongue squamous cell carcinoma (OTSCC) and used those in an integrated

631 approach linking genes harboring somatic variants with common risk factors like

632 tobacco and alcohol; clinical, epidemiological factors like tumor grade and HPV; and

633 tumor recurrence. We found *CASP8* gene to be significantly altered and play an

634 important role in apoptosis-mediated cell death in an HPV-negative OTSCC cell line.

635 Finally, we present data towards a minimal gene signature that can predict tumor

636 recurrence.

637

638 **Availability of supporting data**

639 All supporting data is available through the journal's website and on figshare.

640

641 **List of abbreviations:** Oral tongue squamous cell carcinoma (OTSCC), head and

642 neck squamous cell carcinoma (HNSCC), single nucleotide variant (SNV), insertion

643 and deletion (indel), loss of heterozygosity (LOH), copy number variation (CNV)

644

645 **Competing interests**

646 None

647 **Authors' contributions**

648 BP: conceived, designed and supervised the study, wrote the manuscript; NMK:

649 analyzed the data and wrote the manuscript; SG: analyzed the data and critically read

650 the manuscript; SP, PJ, CK, VKP: analyzed the data; VP, GS, AS: produced data on

651 *CASP8* functional analysis; LV, AKH, MP: produced sequencing data; KD and JN:

652 produced array data; GS, AS, VK and MAK: provided clinical data, clinical input and

653 associated clinical information.

654

## Acknowledgements

## References

1.      Ferlay J, Shin HR, Bray F, Forman D, Mathers C, Parkin DM: **Estimates of worldwide burden of cancer in 2008: GLOBOCAN 2008.** *Int J Cancer* 2010, **127:**2893-2917.

2.      Mishra A, Meherotra R: **Head and neck cancer: global burden and regional trends in India.** *Asian Pac J Cancer Prev* 2014, **15:**537-550.

3.      Stransky N, Egloff AM, Tward AD, Kostic AD, Cibulskis K, Sivachenko A, Kryukov GV, Lawrence MS, Sougnez C, McKenna A, et al: **The mutational landscape of head and neck squamous cell carcinoma.** *Science* 2011, **333:**1157-1160.

4.      Agrawal N, Frederick MJ, Pickering CR, Bettegowda C, Chang K, Li RJ, Fakhry C, Xie TX, Zhang J, Wang J, et al: **Exome sequencing of head and neck squamous cell carcinoma reveals inactivating mutations in NOTCH1.** *Science* 2011, **333:**1154-1157.

5.      Pickering CR, Zhang J, Yoo SY, Bengtsson L, Moorthy S, Neskey DM, Zhao M, Ortega Alves MV, Chang K, Drummond J, et al: **Integrative genomic characterization of oral squamous cell carcinoma identifies frequent somatic drivers.** *Cancer Discov* 2013, **3:**770-781.

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

7.      Cancer Genome Atlas N: **Comprehensive genomic characterization of head and neck squamous cell carcinomas.** *Nature* 2015, **517:**576-582.

8.      Llewellyn CD, Johnson NW, Warnakulasuriya KA: **Risk factors for squamous cell carcinoma of the oral cavity in young people--a comprehensive literature review.** *Oral Oncol* 2001, **37:**401-418.

684  9.    Kuriakose M, Sankaranarayanan M, Nair MK, Cherian T, Sugar AW, Scully C, Prime SS:

685        **Comparison of oral squamous cell carcinoma in younger and older patients in India.** *Eur*

686        *J Cancer B Oral Oncol* 1992, **28B:**113-120.

687  10.   Pathak KA, Das AK, Agarwal R, Talole S, Deshpande MS, Chaturvedi P, Pai PS, Chaukar

688        DA, D'Cruz AK: **Selective neck dissection (I-III) for node negative and node positive**

689        **necks.** *Oral Oncol* 2006, **42:**837-841.

690  11.   Breiman L: *Random Forests*. The Netherlands: Kliwer Academic Publishers; 2001.

691  12.   Pattnaik S, Vaidyanathan S, Pooja DG, Deepak S, Panda B: **Customisation of the exome**

692        **data analysis pipeline using a combinatorial approach.** *PLoS One* 2012, **7:**e30080.

693  13.   Oberst A, Green DR: **It cuts both ways: reconciling the dual roles of caspase 8 in cell**

694        **death and survival.** *Nat Rev Mol Cell Biol* 2011, **12:**757-763.

695  14.   Fulda S: **Caspase-8 in cancer biology and therapy.** *Cancer Lett* 2009, **281:**128-133.

696  15.   Lansford CD, Grenman R, Bier H, Somers KD, Kim SY, Whiteside TL, Clayman GL,

697        Welkoborsky H-J, Carey TE: *Head and neck cancers*. New York: Kluwer Academic

698        Publishers; 2002.

699  16.   Telmer CA, An J, Malehorn DE, Zeng X, Gollin SM, Ishwad CS, Jarvik JW: **Detection and**

700        **assignment of TP53 mutations in tumor DNA using peptide mass signature genotyping.**

701        *Hum Mutat* 2003, **22:**158-165.

702  17.   Tibshirani R, Efron B: **Improvements on Cross-Validation: The .632 + Bootstrap Method.**

703        *J Am Stat Assoc* 1997, **92:**13.

704  18.   Kumar B, Cordell KG, Lee JS, Prince ME, Tran HH, Wolf GT, Urba SG, Worden FP,

705        Chepeha DB, Teknos TN, et al: **Response to therapy and outcomes in oropharyngeal**

706        **cancer are associated with biomarkers including human papillomavirus, epidermal**

707        **growth factor receptor, gender, and smoking.** *Int J Radiat Oncol Biol Phys* 2007, **69:**S109-

708        111.

709  19.   Worden FP, Kumar B, Lee JS, Wolf GT, Cordell KG, Taylor JM, Urba SG, Eisbruch A,

710        Teknos TN, Chepeha DB, et al: **Chemoselection as a strategy for organ preservation in**

711        **advanced oropharynx cancer: response and survival positively associated with HPV16**

712        **copy number.** *J Clin Oncol* 2008, **26:**3138-3146.

713  20.   Fakhry C, Zhang Q, Nguyen-Tan PF, Rosenthal D, El-Naggar A, Garden AS, Soulieres D,

714        Trotti A, Avizonis V, Ridge JA, et al: **Human papillomavirus and overall survival after**

715        **progression of oropharyngeal squamous cell carcinoma.** *J Clin Oncol* 2014, **32:**3365-3373.

716   21.  Dayyani F, Etzel CJ, Liu M, Ho CH, Lippman SM, Tsao AS: **Meta-analysis of the impact of**

717        **human papillomavirus (HPV) on cancer risk and overall survival in head and neck**

718        **squamous cell carcinomas (HNSCC).** *Head Neck Oncol* 2010, **2:**15.

719   22.  Dahlgren L, Dahlstrand HM, Lindquist D, Hogmo A, Bjornestal L, Lindholm J, Lundberg B,

720        Dalianis T, Munck-Wikland E: **Human papillomavirus is more common in base of tongue**

721        **than in mobile tongue cancer and is a favorable prognostic factor in base of tongue**

722        **cancer patients.** *Int J Cancer* 2004, **112:**1015-1019.

723   23.  Salem A: **Dismissing links between HPV and aggressive tongue cancer in young patients.**

724        *Ann Oncol* 2010, **21:**13-17.

725   24.  Elango KJ, Suresh A, Erode EM, Subhadradevi L, Ravindran HK, Iyer SK, Iyer SK,

726        Kuriakose MA: **Role of human papilloma virus in oral tongue squamous cell carcinoma.**

727        *Asian Pac J Cancer Prev* 2011, **12:**889-896.

728   25.  Bullenkamp J, Raulf N, Ayaz B, Walczak H, Kulms D, Odell E, Thavaraj S, Tavassoli M:

729        **Bortezomib sensitises TRAIL-resistant HPV-positive head and neck cancer cells to**

730        **TRAIL through a caspase-dependent, E6-independent mechanism.** *Cell Death Dis* 2014,

731        **5:**e1489.

732   26.  Manzo-Merino J, Massimi P, Lizano M, Banks L: **The human papillomavirus (HPV) E6**

733        **oncoproteins promotes nuclear localization of active caspase 8.** *Virology* 2014, **450-**

734        **451:**146-152.

735   27.  Novocraft: **Novoalign (www.novocraft.com).** 2011.

736   28.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin

737        R: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25:**2078-

738        2079.

739   29.  McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,

740        Altshuler D, Gabriel S, Daly M, DePristo MA: **The Genome Analysis Toolkit: a**

741        **MapReduce framework for analyzing next-generation DNA sequencing data.** *Genome*

742        *Res* 2010, **20:**1297-1303.

743   30.  Albers CA, Lunter G, MacArthur DG, McVean G, Ouwehand WH, Durbin R: **Dindel:**

744        **accurate indel calls from short-read data.** *Genome Res* 2011, **21:**961-973.

745   31.    Cibulskis K, McKenna A, Fennell T, Banks E, DePristo M, Getz G: **ContEst: estimating**

746        **cross-contamination of human samples in next-generation sequencing data.**

747        *Bioinformatics* 2011, **27:**2601-2602.

748   32.    Carter H, Chen S, Isik L, Tyekucheva S, Velculescu VE, Kinzler KW, Vogelstein B, Karchin

749        R: **Cancer-specific high-throughput annotation of somatic mutations: computational**

750        **prediction of driver missense mutations.** *Cancer Res* 2009, **69:**6660-6667.

751   33.    Douville C, Carter H, Kim R, Niknafs N, Diekhans M, Stenson PD, Cooper DN, Ryan M,

752        Karchin R: **CRAVAT: cancer-related analysis of variants toolkit.** *Bioinformatics* 2013,

753        **29:**647-648.

754   34.    Gundem G, Perez-Llamas C, Jene-Sanz A, Kedzierska A, Islam A, Deu-Pons J, Furney SJ,

755        Lopez-Bigas N: **IntOGen: integration and data mining of multidimensional oncogenomic**

756        **data.** *Nat Methods* 2010, **7:**92-93.

757   35.    Schroeder MP, Gonzalez-Perez A, Lopez-Bigas N: **Visualizing multidimensional cancer**

758        **genomics data.** *Genome Med* 2013, **5:**9.

759   36.    Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, Meyerson M,

760        Gabriel SB, Lander ES, Getz G: **Discovery and saturation analysis of cancer genes across**

761        **21 tumour types.** *Nature* 2014, **505:**495-501.

762   37.    Dees ND: **MuSiC2.** 2015.

763   38.    Zhao M, Zhao Z: **CNVannotator: a comprehensive annotation server for copy number**

764        **variation in the human genome.** *PLoS One* 2013, **8:**e80170.

765   39.    Du P, Kibbe WA, Lin SM: **lumi: a pipeline for processing Illumina microarray.**

766        *Bioinformatics* 2008, **24:**1547-1548.

767   40.    Johnson WE, Li C, Rabinovic A: **Adjusting batch effects in microarray expression data**

768        **using empirical Bayes methods.** *Biostatistics* 2007, **8:**118-127.

769   41.    Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK: **limma powers**

770        **differential expression analyses for RNA-sequencing and microarray studies.** *Nucleic*

771        *Acids Res* 2015, **43:**e47.

772   42.    Sales G, Calura E, Martini P, Romualdi C: **Graphite Web: Web tool for gene set analysis**

773        **exploiting pathway topology.** *Nucleic Acids Res* 2013, **41:**W89-97.

774  43.  Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B,

775      Ideker T: **Cytoscape: a software environment for integrated models of biomolecular**

776      **interaction networks.** *Genome Res* 2003, **13:**2498-2504.

777  44.  Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA:

778      **Circos: an information aesthetic for comparative genomics.** *Genome Res* 2009, **19:**1639-

779      1645.

780  45.  Thorvaldsdottir H, Robinson JT, Mesirov JP: **Integrative Genomics Viewer (IGV): high-**

781      **performance genomics data visualization and exploration.** *Brief Bioinform* 2013, **14:**178-

782      192.

783  46.  Schmittgen TD, Livak KJ: **Analyzing real-time PCR data by the comparative C(T)**

784      **method.** *Nat Protoc* 2008, **3:**1101-1108.

785

786  **Figure legends**

787

788  **Figure 1. Key variants in OTSCC and their relationship with habits, clinical and**

789  **epidemiological parameters.**

790      A. The OTSCC samples are represented in color-codes with their

791      corresponding status on; node (P: positive, N: negative); stage (E: early, L:

792      late), recurrence (Y: loco-regionally recurrent, N: non-recurrent and M: distant

793      metastatic); grade (WD: well-differentiated, MD: moderately-differentiated

794      and PD: poorly-differentiated); disease-free survival or DFS (L: low/≤12mo,

795      M: mid/12-24mo and H: high/>24mo); HPV (P: positive and N: negative); and

796      habits (chewing, alcohol and smoking, Y: yes and N: no). B. Somatic mutation

797      frequency per megabase (MB) is represented as scatterplot with the median

798      point as a fine dotted line. C. Genes with significant somatic variants. D.

799      Frequency histogram of nineteen cancer-associated genes bearing somatic

800      missense and nonsense variants (mutations and indels). E. Columns

801     representing mutually exclusive sets of genes. F. Significant copy number

802     insertions and deletions (CNVs), alongside the chromosome cytobands (the

803     numbers of cancer-associated genes within each cytoband are listed on the

804     right).

805

806     **Figure 2. Relationship between genes harboring somatic variants with clinical-,**

807     **epidemiological parameters and signaling pathways.**

808     A. Histograms showing relationship between genes with significant somatic

809     variants and various clinical and epidemiological parameters,.for genes solely

810     mutated in one of the clinical or epidemiological categories, or those mutated

811     at a >= 5% frequency between two categories. B. Stack net charts of relative

812     patient fraction (%) for each of the eight cancer-associated signaling pathways

813     and their relationship with various clinical and epidemiological parameters.

814

815     **Figure 3. Differentially expressed genes, affected pathways and their relationship**

816     **with clinical and epidemiological parameters.**

817     A. Expression changes (green – up- and red – down-regulation) representing

818     significantly differentially expressed genes in tumors. B. Stacked histograms

819     representing relative patient fraction (%) for each of the 19 cancer-associated

820     pathways and their relationship with clinical and epidemiological parameters.

821

822     **Figure 4: Role of *CASP8* in HPV-positive and HPV-negative OTSCC cell lines.**

823     Results from A. Matrigel cell invasion assay (plotted with respect to the

824     control cells), B. Wound healing assay, and C. MTT cell survival assay

825     (plotted with respect to the control cells) in UPCI:SCC040 (HPV-negative)

826    and UM:SCC-47 (HPV-positive) cell lines.

827

828    **Figure 5. A minimal gene signature for tumor recurrence.**

829    Genes harboring somatic variants (in color) that are a part of the minimal

830    signature set for tumor recurrence derived from random forest analyses are

831    used.

832

833    **Figure 6. Significantly affected pathways in OTSCC.**

834    Genes harboring significant somatic variants and with expression changes in

835    tumors were used in Cytoscape to derive a set of important signaling pathways

836    implicated in OTSCC.

837    **Additional files**

838    All the additional files can be downloaded from figshare

839    (http://figshare.com/s/c928faa66f2a11e586d506ec4b8d1f61).

840    Additional file 1:

841    Title: Patient details used in the study.

842    Additional file 2:

843    Title: Sequencing, read QC, alignment and variant calls, OMNI SNP

844    genotyping array validation.

845    Additional file 3:

846    Title: Validation using capillary gel electrophoresis based on Sanger

847    sequencing in A. discovery set and B. validation set.

848    Additional file 4:

849    Title: The ratio of transitions to transversions (ti/tv) was estimated using the

850    exome-filtered GATK PASS variants for tumor and matched control samples.

851          The dotted lines depict the respective median ti/tv ratios.

852     Additional file 5:

853          Title: Effect of habits, clinical parameters and HPV infection on individual

854          nucleotide change.

855     Additional file 6:

856          Title: Functional annotation of somatic variants using IntOGen, MuSiC2,

857          MutSigCV.

858     Additional file 7:

859          Title: Position and frequency of somatic variants in protein domains found in

860          this study, TCGA HNSCC, and in studies involving all cancer types using

861          mutation mapper in the cBioPortal.

862     Additional file 8:

863          Title: Cytoband-wise representation of CNVs found in all 48 samples along

864          with clinical parameters and patient epidemiology.

865     Additional file 9:

866          Title: Transformed, normalized and batch-corrected intensities following

867          expression assay and results from differential expression analyses.

868     Additional file 10:

869          Title: Functional validation for the role of *CASP8* in OTSCC cell lines.

870     Additional file 11:

871          Title: List of all pathways affected by somatic mutations/indels, copy number

872          variations and expression changes ($log_2FC \geq =0.6$).

873     Additional file 12:

874          Title: Comparative sample frequency of important variants found in this and

875          other HNSCC studies.

<antchor file="boilerplate">bioRxiv preprint doi: https://doi.org/10.1101/028845; this version posted October 11, 2015. 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-NC 4.0 International license.</antchor>

876    Additional file 13:

877        Title: Drug candidates and their targets in head and neck cancer.

878    Additional file 14:

879        Title: Supplementary Methods.

880    Additional file 15:

881        Title: Read QC filter scripts' executable.

882    Additional file 16:

883        Title: Scripts used in the study.

884    Additional file 17:

885        Title: GenomeStudio output of all LOHs and CNVs found using cnvPartition

886        plugin in GenomeStudio.

887    Additional file 18:

888        Title: Circos data and config files.

889    Additional file 19:

890        Title: Circular genomic representation using Circos (v0.66) of LOHs, somatic

891        variants, CNVs with >= 10% frequency of patients bearing them, and genes

892        with significant expression changes ($|\log_2 FC| >= 0.6$).

893    Additional file 20:

894        Title: IGV snapshots of all significantly mutated somatic variants in this study.

895    Additional file 21:

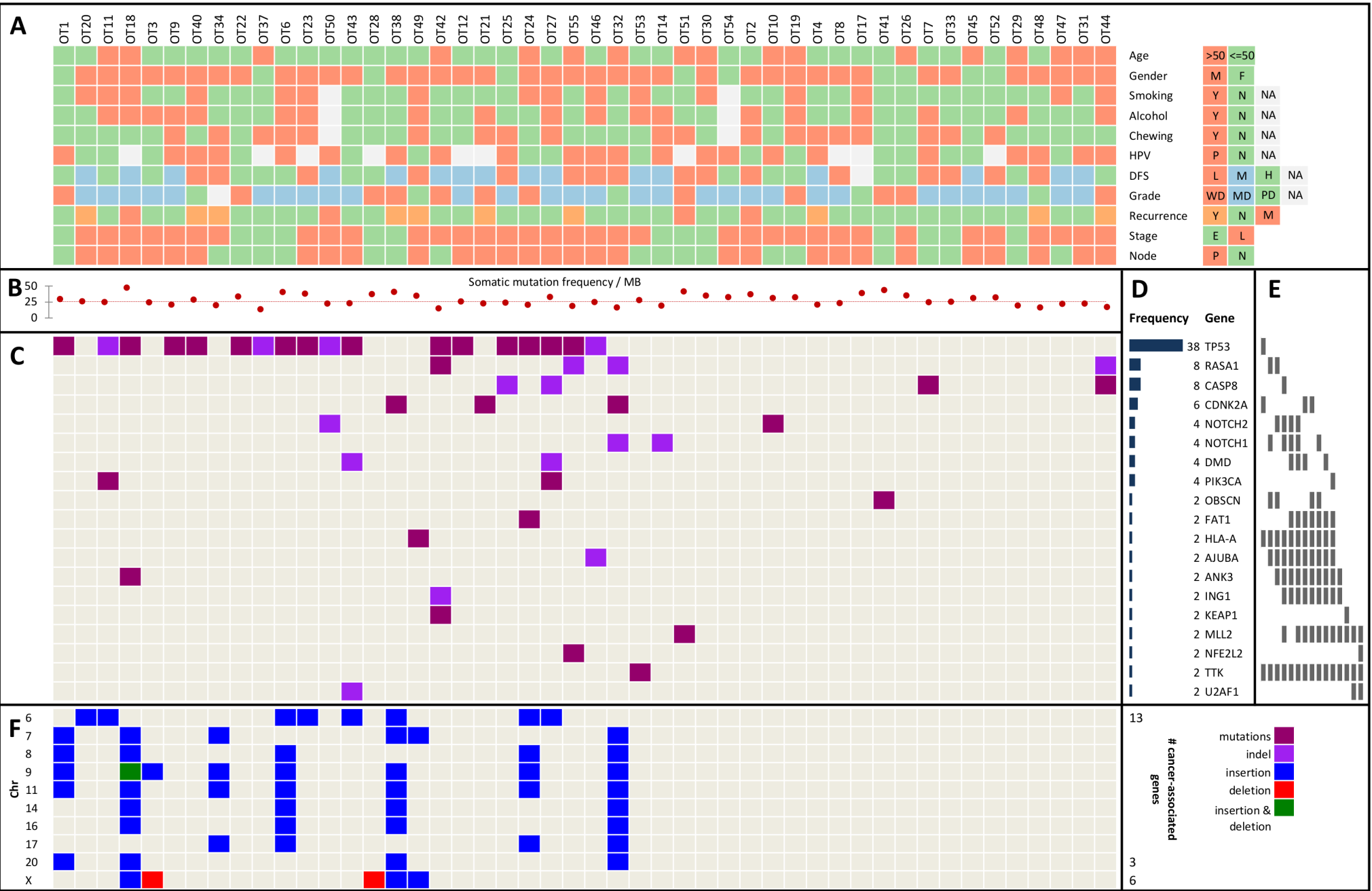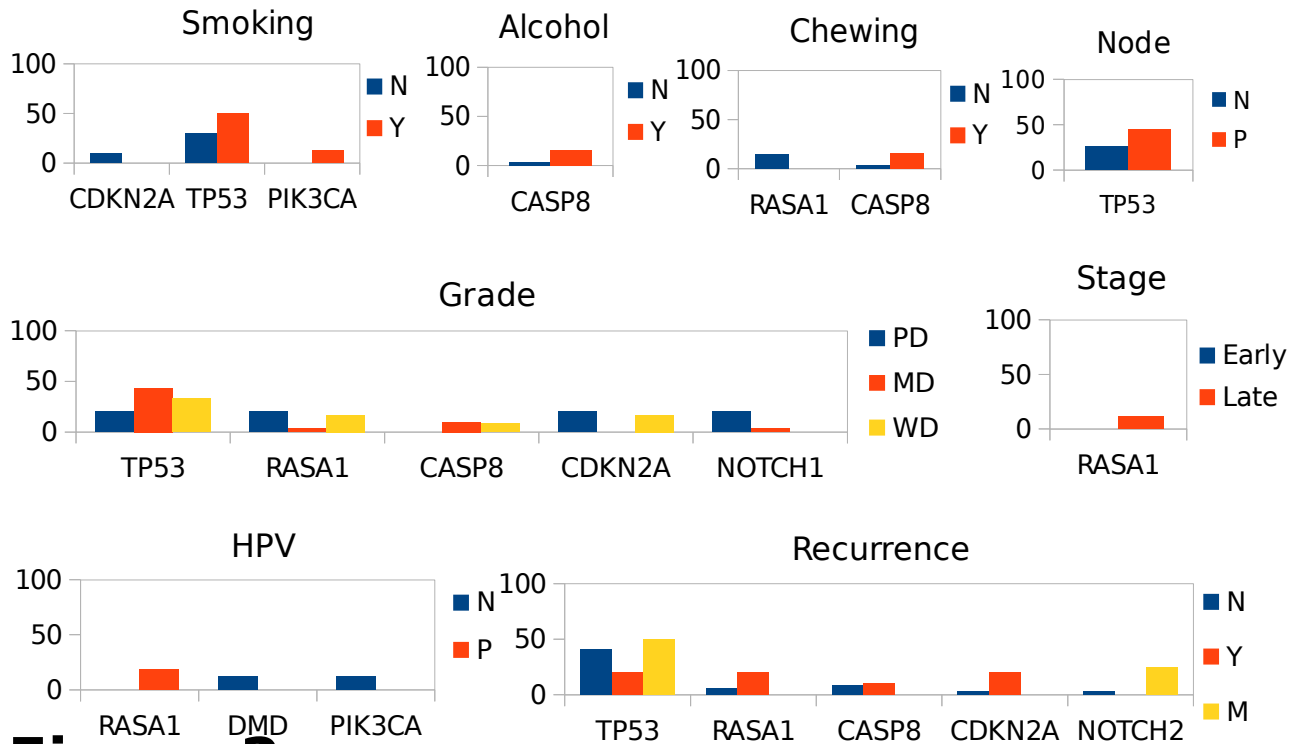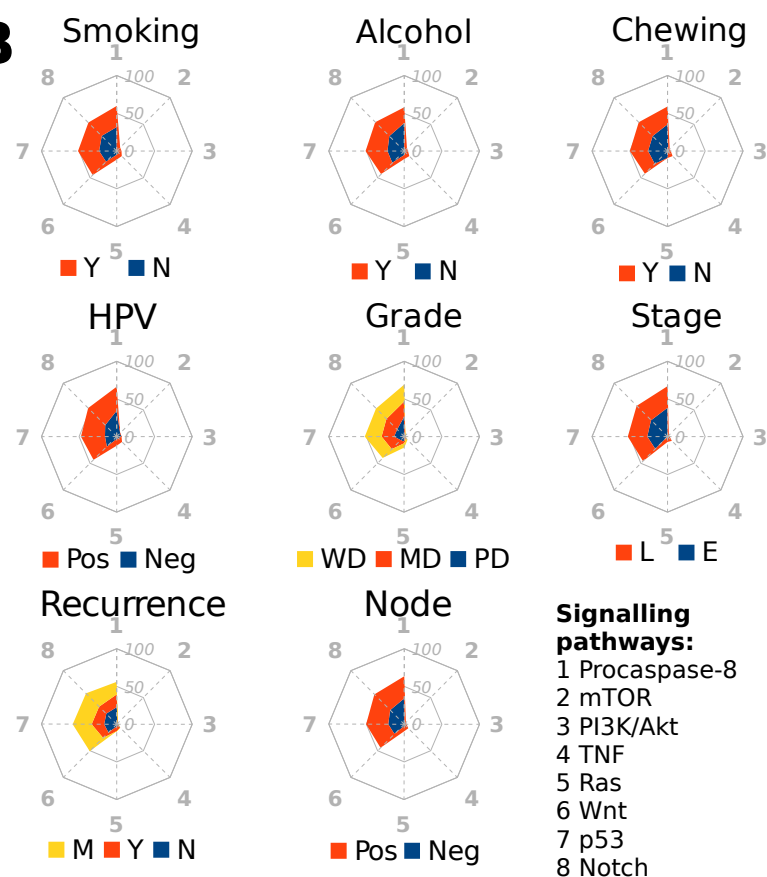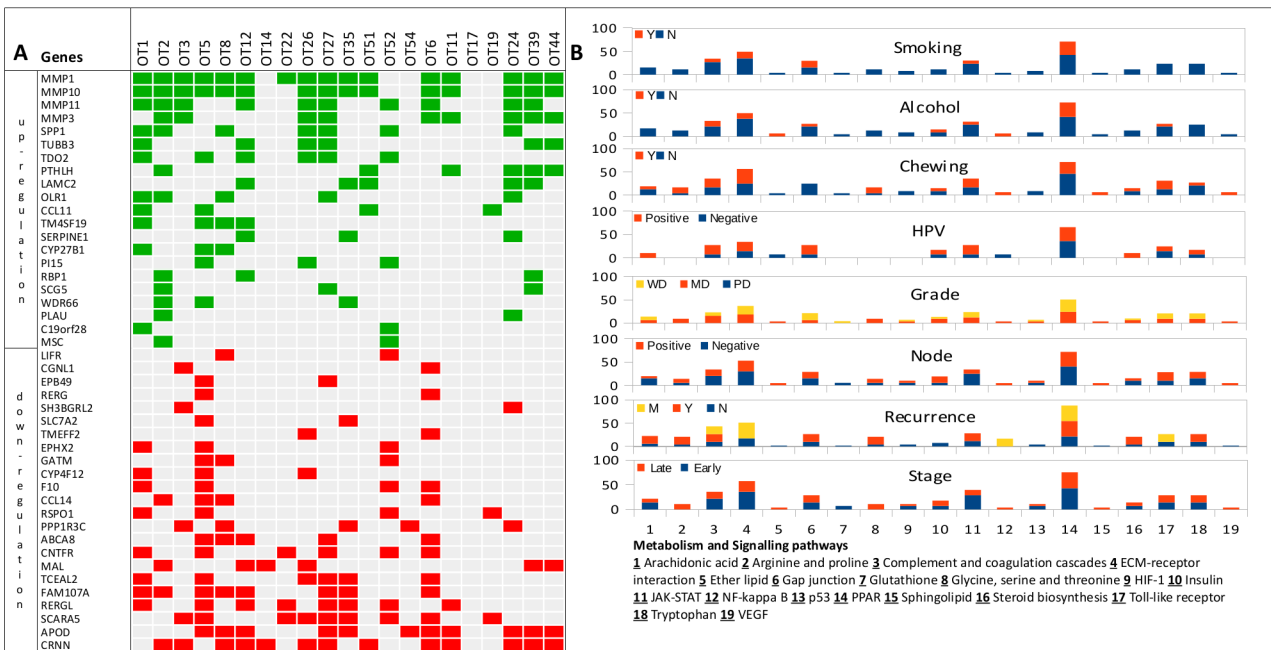896        Title: Primer sequences used in the Sanger validation study.

**Figure 1**

**Figure 2**

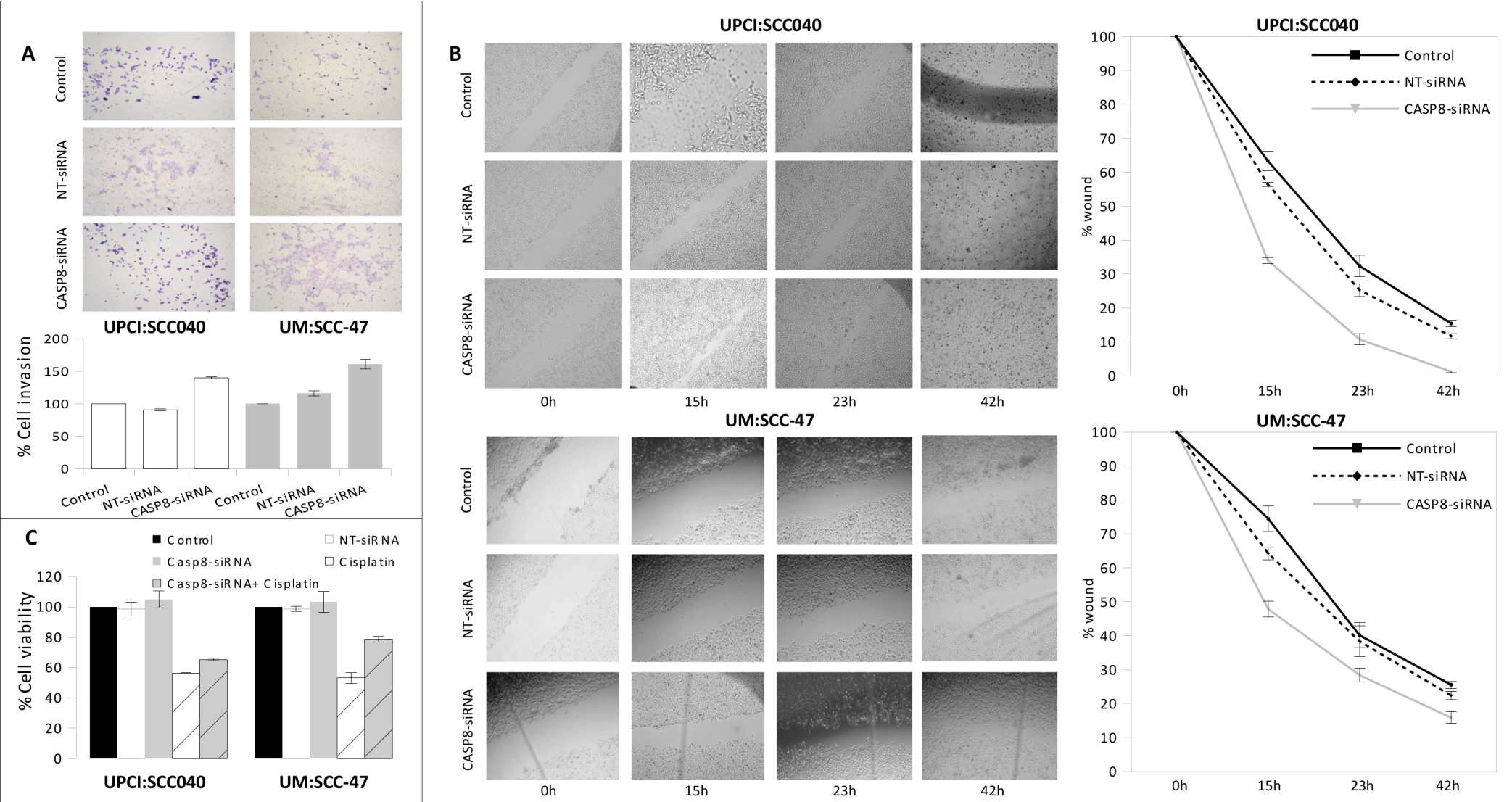Signalling pathways:
1 Procaspase-8
2 mTOR
3 PI3K/Akt
4 TNF
5 Ras
6 Wnt
7 p53
8 Notch

**Figure 3**

Figure 4

Figure 5

Figure 6