



## RESEARCH ARTICLE

# Integrated analysis of oral tongue squamous cell carcinoma identifies key variants and pathways linked to risk habits, HPV, clinical parameters and tumor recurrence [version 1; referees: 2 approved]

Neeraja Krishnan<sup>1</sup>, Saurabh Gupta<sup>1</sup>, Vinayak Palve<sup>1</sup>, Linu Varghese<sup>1</sup>, Swetansu Pattnaik<sup>1</sup>, Prach Jain<sup>1</sup>, Costerwell Khyriem<sup>1</sup>, Arun Hariharan<sup>1</sup>, Kunal Dhas<sup>1</sup>, Jayalakshmi Nair<sup>1</sup>, Manisha Pareek<sup>1</sup>, Venkatesh Prasad<sup>1</sup>, Gangotri Siddappa<sup>2</sup>, Amritha Suresh<sup>2</sup>, Vikram Kekatpure<sup>3</sup>, Moni Kuriakose<sup>2,3</sup>, Binay Panda<sup>1,4</sup>

<sup>1</sup>Ganit Labs, Bio-IT Centre, Institute of Bioinformatics and Applied Biotechnology, Bangalore, 560 100, India

<sup>2</sup>Integrated Head and Neck Oncology Program, Mazumdar Shaw Centre for Translational Research, Bangalore, 560 099, India

<sup>3</sup>Head and Neck Oncology, Mazumdar Shaw Medical Centre, Bangalore, 560 099, India

<sup>4</sup>Strand Life Sciences, Bangalore, 560 024, India

**v1** First published: 05 Nov 2015, 4:1215 (doi: [10.12688/f1000research.7302.1](https://doi.org/10.12688/f1000research.7302.1))  
Latest published: 05 Nov 2015, 4:1215 (doi: [10.12688/f1000research.7302.1](https://doi.org/10.12688/f1000research.7302.1))

## Abstract

Oral tongue squamous cell carcinomas (OTSCC) are a homogeneous group of tumors characterized by aggressive behavior, early spread to lymph nodes and a higher rate of regional failure. Additionally, the incidence of OTSCC among younger population (<50yrs) is on the rise; many of whom lack the typical associated risk factors of alcohol and/or tobacco exposure. We present data on single nucleotide variations (SNVs), indels, regions with loss of heterozygosity (LOH), and copy number variations (CNVs) from fifty-paired oral tongue primary tumors and link the significant somatic variants with clinical parameters, epidemiological factors including human papilloma virus (HPV) infection and tumor recurrence. Apart from the frequent somatic variants harbored in TP53, CASP8, RASA1, NOTCH and CDKN2A genes, significant amplifications and/or deletions were detected in chromosomes 6-9, and 11 in the tumors. Variants in CASP8 and CDKN2A were mutually exclusive. CDKN2A, PIK3CA, RASA1 and DMD variants were exclusively linked to smoking, chewing, HPV infection and tumor stage. We also performed a whole-genome gene expression study that identified matrix metalloproteases to be highly expressed in tumors and linked pathways involving arachidonic acid and NF-k-B to habits and distant metastasis, respectively. Functional knockdown studies in cell lines demonstrated the role of CASP8 in a HPV-negative OTSCC cell line. Finally, we identified a 38-gene minimal signature that predicts tumor recurrence using an ensemble machine-learning method. Taken together, this study links molecular signatures to various clinical and epidemiological factors in a homogeneous tumor population with a relatively high HPV prevalence.

## Open Peer Review

Referee Status:

|   | Invited Referees |            |
|---|------------------|------------|
|   | 1                | 2          |
| <b>version 1</b><br>published<br>05 Nov 2015                | <br>report       | <br>report |
| <b>1 Nishant Agrawal</b> , The Johns Hopkins University USA |                  |            |
| <b>2 Thomas Carey</b> , University of Michigan USA          |                  |            |
| <b>Discuss this article</b>                                 |                  |            |
| Comments (0)  |                  |            |

**Corresponding author:** Binay Panda ([binay@ganitlabs.in](mailto:binay@ganitlabs.in))

**How to cite this article:** Krishnan N, Gupta S, Palve V *et al.* **Integrated analysis of oral tongue squamous cell carcinoma identifies key variants and pathways linked to risk habits, HPV, clinical parameters and tumor recurrence [version 1; referees: 2 approved]** *F1000Research* 2015, 4:1215 (doi: [10.12688/f1000research.7302.1](https://doi.org/10.12688/f1000research.7302.1))

**Copyright:** © 2015 Krishnan N *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Grant information:** Research presented in this article is funded by Department of Electronics and Information Technology, Government of India (Ref No:18(4)/2010-E-Infra., 31-03-2010) and Department of IT, BT and ST, Government of Karnataka, India (Ref No:3451-00-090-2-22). *The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

**Competing interests:** None.

**First published:** 05 Nov 2015, 4:1215 (doi: [10.12688/f1000research.7302.1](https://doi.org/10.12688/f1000research.7302.1))

## Introduction

Squamous cell carcinomas of the head and neck (HNSCC) are the sixth leading cause of cancer worldwide<sup>1</sup>. Tumors of the head and neck region are heterogeneous in nature with different incidences, mortalities and prognoses for different subsites and accounts for almost 30% of all cancer cases in India<sup>2</sup>. Oral cancer is the most common subtype of head and neck cancer in humans, with a worldwide incidence in >300,000 cases. The disease is an important cause of death and morbidity, with a 5-year survival of less than 50%<sup>1,2</sup>. Recent studies have identified various genetic changes in many subsites of the head and neck using high-throughput sequencing assays and computational methods<sup>3-7</sup>. Such multi-tiered approaches using the exomes, genomes, transcriptomes and methylomes from different squamous cell carcinomas have generated data on key variants and in some cases, their biological significance, aiding our understanding of disease progression. Some of the above sequencing studies have identified key somatic variants and linked them with patient stratification and prognostication. This, along with the associated epidemiology, enables one to look beyond the discovery of driver mutations, and identify predictive signatures in HNSCC.

A previous study from the cancer genome atlas (TCGA) consortium with HNSCC patients (N = 279) identified somatic mutations in *TP53*, *CDKN2A*, *FAT1*, *PIK3CA*, *NOTCH1*, *KMT2D* and *NSD1* at a frequency greater than 10%<sup>7</sup>. Additionally, the TCGA study identified loss of *TRAF3* gene, amplification of *E2F1* in human papilloma virus (HPV)-positive oropharyngeal tumors, along with mutations in *PIK3CA*, *CASP8* and *HRAS*, and co-amplifications of the regions 11q13 (harboring *CCND1*, *FADD* and *CTTN*) and 11q22 (harboring *BIRC2* and *YAP1*), in HPV-negative tumors, described to play an important role in pathogenesis and tumor development<sup>7</sup>. Chromosomal losses at 3p and 8p, and gains at 3q, 5p and 8q were also observed in HNSCC<sup>7</sup>. Tumors originating in the anterior/oral part of the tongue, or oral tongue squamous cell carcinoma (OTSCC) tend to be different from those at other subsites as oral tongue tumors are associated more with younger patients<sup>8</sup> and spread early to lymph nodes<sup>9</sup>. Additionally oral tongue tumors have a higher regional failure compared to gingivo-buccal cases<sup>10</sup> in oral cavity. Tobacco (both chewing and smoking) and alcohol are common risk factors for this group of tumors among older patients<sup>8</sup>. The role of HPV, both as an etiological agent and/or risk factor along with its role as a good prognostic marker in OTSCC, unlike in oropharyngeal tumors, is currently uncertain. It remains to be explored whether HPV acts as an etiological agent in the development of oral tongue tumors or simply represents a superinfection in patients. Additionally, HPV infection status currently does not influence disease management in OTSCC.

Here, we present data towards a comprehensive molecular characterization of OTSCC. We performed exome sequencing, whole-genome gene expression, and genotyping arrays using fifty primary tumors along with their matched control samples, towards identification of somatic variants (mutations and indels), significantly up- and down-regulated genes, loss of heterozygosity (LOH) and copy number variations (CNVs). We integrated all the molecular data along with the clinical parameters and epidemiology such as tumor stage, nodal status, HPV infection, risk habits and tumor recurrence to interpret the effect of changes in the process of cancer

development in oral tongue. We identified significant somatic variations in *TP53* (38%), *RASA1* (8%), *CASP8* (8%), *CDKN2A* (6%), *NOTCH1* (4%), *NOTCH2* (4%), and *PIK3CA* (4%) from the exome sequencing study in OTSCC. The key variants were validated using an additional set of primary tumor samples. Variants in *TP53* and *NOTCH1* were found in mutually exclusive sets of tumors. Additionally, we found frequent aberrations in chromosomes 6-9, and 11 in tumor samples. We observed a strong association between somatic variations in some key genes with one or more risk habits; for example, *CDKN2A* and *PIK3CA* with smoking; *CASP8* with consuming alcohol and chewing tobacco; *RASA1* with chewing and tumor stage, and HPV infection, along with *DMD* and *PIK3CA*. From the gene expression analysis, we found matrix metalloproteases (MMPs) to be highly expressed in OTSCC. Pathway analysis identified Procaspase-8, Notch, Wnt, p53, extracellular matrix (ECM)-receptor interaction, JAK-STAT and PPAR to be some of the significantly altered pathways in OTSCC. We implemented an ensemble machine-learning method<sup>11</sup> and identified a minimal gene signature set that distinguished a group of tumors with loco-regional recurrence from the non-recurrent set. Finally, we performed functional analysis of *CASP8* gene in HPV-negative and HPV-positive OTSCC cell lines to establish its role in the process of tumor development.

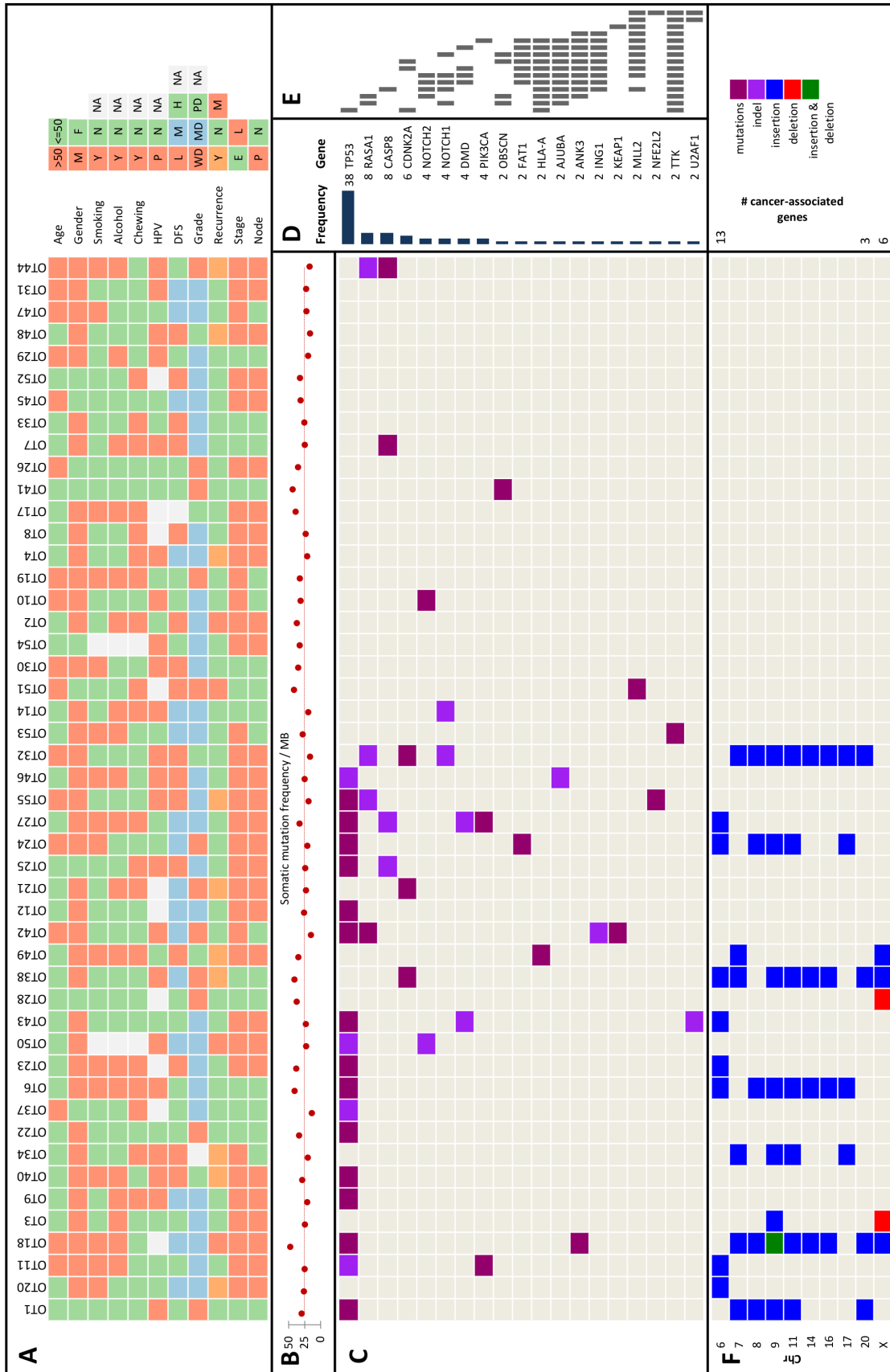
## Results

### Habits, clinical parameters and epidemiology

We collected tumor and matched control (adjacent normal and/or lymphocytes) samples from 50 patients diagnosed with OTSCC, with informed consent. Data from patient habits, epidemiology and clinical parameters are presented in [Figure 1A](#) and [Additional file 1A](#). About two-thirds of the patients (N = 31) included in our study were in the younger age group ( $\leq 50$  yrs), with 20% female patients in the total pool. Approximately, 70% of the patients were positive for at least one risk habit, namely, smoking, alcohol consumption or chewing tobacco (33% of patients smoked tobacco, 40% consumed alcohol and 42% chewed tobacco). HPV infection status in the primary tumors was established with type-specific qPCR or HPV16 digital PCR. Thirty-three percent of the patients were deceased at the time of completing the analysis. About 60% of the tumors were moderately differentiated, 25% well differentiated and the rest were poorly differentiated. Among the patients recruited, 60% were node-positive, 70% had no recurrence, 9% had distant metastasis and 24% had loco-regional recurrence at the time of completing the analysis. The mean and median follow-up durations for patients were nearly 30 months and 21 months, respectively. About 27% of the tumors were early stage tumors (T1N0M0 and T2N0M0) and the rest 73% were late stage tumors (tumors belonging to the rest of the TNM stage).

### Discovery and validation of significant somatic variants and their relationship with other parameters

We re-discovered variants, as described previously<sup>12</sup> 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



**Figure 1. Key variants in OTSCC and their relationship with habits, clinical and epidemiological parameters.** **A.** The OTSCC samples are represented in color-codes with their corresponding status on: node (P: positive, N: negative); stage (E: early, L: late), recurrence (Y: loco-regionally recurrent, N: non-recurrent and M: distant metastatic); grade (WD: well-differentiated, MD: moderately-differentiated and PD: poorly-differentiated); disease-free survival or DFS (L: low/ $\leq 12$ mo, M: mid/ $12-24$ mo and H: high/ $>24$ mo); HPV (P: positive and N: negative); and habits (chewing, alcohol and smoking, Y: yes and N: no). **B.** Somatic mutation frequency per megabase (MB) is represented as scatterplot with the median point as a line dotted line. **C.** Genes with significant somatic variants. **D.** Frequency histogram of nineteen cancer-associated genes bearing somatic missense and nonsense variants (mutations and indels). **E.** Columns representing mutually exclusive sets of genes. **F.** Significant copy number insertions and deletions (CNVs), alongside the chromosome cytobands (the numbers of cancer-associated genes within each cytoband are listed on the right).

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  $\geq 5\%$  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 *RASAI* 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 megabase (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  $\sim 2.5$  ([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 stage had no significant impact on the distribution of these nucleotides ([Additional file 5](#)). We used the workflow described in the *Methods* section to identify somatic mutations and indels in tumor samples following which we used three functional tools, IntOGen<sup>19</sup>, MutSigCV<sup>21</sup> and MuSiC2<sup>22</sup> for variant interpretations ([Additional file 6](#)). In order to identify genes harboring significant variants, we used the intersection of these tools, following the criteria that the somatic variants be callable in the matched control sample and present in a single sequencing read in the control sample. This resulted in a final list of 19 cancer-associated genes ([Figure 1C](#)), which were divided into three categories with varying mutation frequencies ([Figure 1D](#)). The three frequency tiers were  $\geq 30\%$  (*TP53*), 6–30% (*RASAI*, *CASP8* and *CDKN2A*) and 2–5% (*NOTCH1*, *NOTCH2*, *DMD* and *PIK3CA* were prominent among them).

Next, we looked for mutual exclusivity of finding somatic variants in the genes and found that many of these genes harbor variants in a mutually exclusive manner across samples ([Figure 1E](#)), suggesting the possibility that there might be some common pathway(s) involved in the development of OTSCC. We observed mutual exclusivity among somatic variants in *NOTCH1* and *NOTCH2* genes, and expanded this finding to identifying 15 such mutually exclusive sets ([Figure 1E](#)). Among them, *CDKN2A*, *HLA-A* and *TTK* form a mutually exclusive set with *TP53*; *RASAI*, *OBSCN*, *HLA-A*, *AJUBA* and *TTK* are mutually exclusive with either *NOTCH1* alone, or *NOTCH2* and *ANK3* together; *NOTCH1*, *NOTCH2*, *HLA-A*, *AJUBA*, *ANK3*, *TTK*, *MLL2*, *ING1* or *KEAP1*, are mutually exclusive with *CASP8* alone, or *FAT1* and *DMD* together; *FAT1*, *HLA-A*, *AJUBA*, *ANK3*, *TTK*, *MLL2*, *ING1* or *KEAP1*, are mutually exclusive with *PIK3CA* or *DMD* or *NOTCH1* and *OBSCN*, or *CDKN2A* and *OBSCN*; *U2AF1*, *MLL2* and *TTK* form a small mutually exclusive set. We

juxtaposed the positions of the somatic variants from final list of all 19 genes ([Additional file 7](#)) detected in OTSCC against those found in the TCGA data using the cBioPortal. We found that the somatic variants in OTSCC were in the same domains where mutations were observed earlier in many of the genes ([Additional file 7](#)).

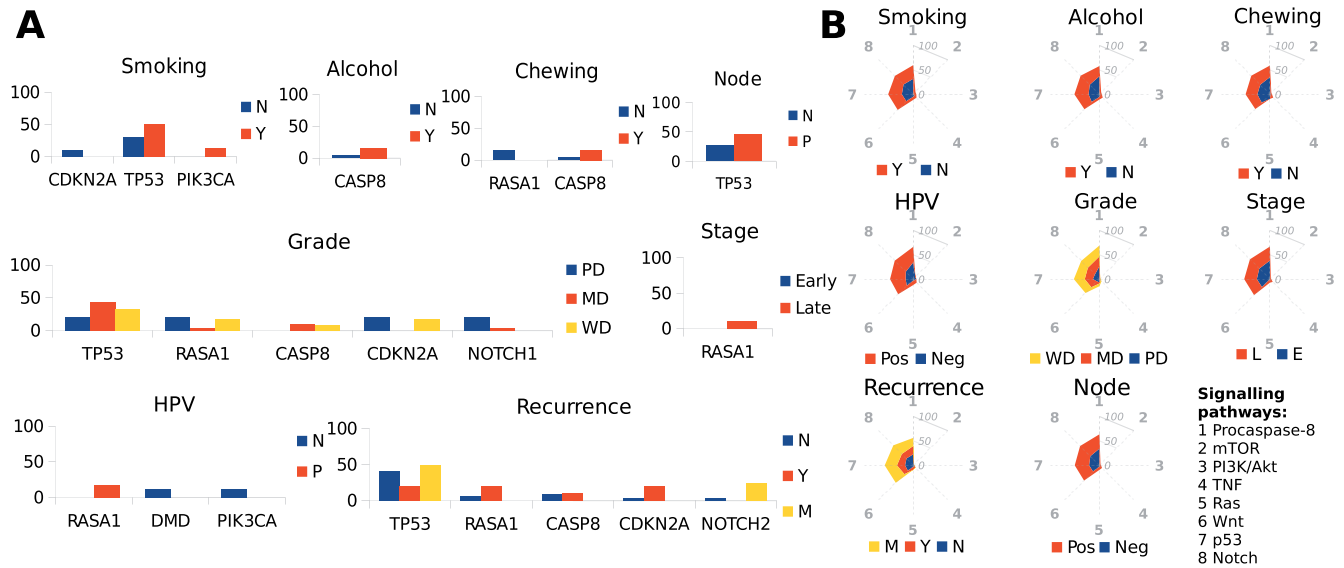
CNV analyses using data from the whole-genome single nucleotide polymorphism (SNP) genotyping arrays revealed a large chunk of chromosome 9, bearing cancer-associated genes like *CDKN2A*, *NF1* and *MRPL4*, to be affected in about 17% of the tumors ([Figure 1F](#) and [Additional file 8](#)). We found several CNVs of short stretches (in low kb range) within chromosomes 6–8, 11, 17 and X in many tumors.

### Linking habits, HPV infection, nodal status, tumor grade and recurrence, with genes harboring somatic variants and the associated pathways

We further classified the 19 cancer-associated genes from the previous analyses and linked those with habits, clinical parameters and HPV infection. Among the genes harboring significant somatic variants, we found *CDKN2A* to be mutated only in the never-smokers and past smokers, *PIK3CA* to be mutated only in the smokers, and *TP53* to be mutated at a 20% greater frequency in the smokers, *CASP8* mutation has a 12% greater frequency in those that consumed alcohol or chewed tobacco. *RASAI* was exclusively mutated only in the non-chewers ([Figure 2A](#)). HPV-negative patients harbored somatic variants in *DMD* and *PIK3CA*, while HPV-positive patients alone had somatic variants in *RASAI*. Only the moderate- and well-differentiated tumor samples harbored variants in *CASP8*, while *NOTCH1* was mutated largely in the poorly-differentiated tumors. Node-positive tumors had a 19% greater occurrence of *TP53* variants. Somatic variants in *RASAI* occurred exclusively in the late stage tumors ([Figure 2A](#)). We further studied the association of affected cancer-related signaling pathways with habits and clinical parameters, and found that recurrence and HPV infection had the highest impact ([Figure 2B](#)). The Procaspase-8 activation, Notch, p53 and Wnt signaling pathways were linked most with many of the clinical parameters, HPV infection and habits ([Figure 2B](#)).

### Differentially expressed genes in OTSCC

Significant ( $q$  val  $\leq 0.05$ ) differentially expressed genes with a fold change of at least 1.5 revealed a consistent pattern of differential expression across the tumor samples (21 up- and 23 down-regulated genes, [Figure 3A](#) and [Additional file 9](#)). Genes involved in peroxisome proliferator-activated receptor (PPAR) signaling- (e.g., *MMP1*) and ECM-receptor interaction pathways (*LAMC2* and *SPP1*) were up-regulated and *CRNN*, *APOD*, *SCARA5* and *RERGL* were down-regulated in a majority of tumors ([Figure 3A](#)). Next, we studied the pathways involving genes with aberrant expression and their link with HPV infection and other clinical parameters. Genes in the arachidonic acid metabolism and Toll-like receptors were differentially expressed in patients with no smoking history (never smokers or past smokers) and alcohol habits ([Figure 3B](#)).



**Figure 2. Relationship between genes harboring somatic variants with clinical-, epidemiological parameters and signaling pathways.** **A.** Histograms showing relationship between genes with significant somatic variants and various clinical and epidemiological parameters. For genes solely mutated in one of the clinical or epidemiological categories, or those mutated at a  $\geq 5\%$  frequency between two categories. **B.** Stack net charts of relative patient fraction (%) for each of the eight cancer-associated signaling pathways and their relationship with various clinical and epidemiological parameters.

*SERPINE1* (a gene in HIF-1 signaling pathway) was differentially expressed in patients that are habits-negative. The NF- $\kappa$ -B signaling pathway was differentially expressed only in metastasized tumors.

### Functional studies with *CASP8* in OTSCC cell lines

*CASP8* is mutated in a significant number of oral tongue tumors [this study, 5, 7]. Caspase-8 is an important and versatile protein that plays a role in both apoptotic (extrinsic or death receptor-mediated) and non-apoptotic processes<sup>13,14</sup>. We studied the functional consequences of *CASP8* knockdown through a siRNA-mediated method in an HPV-positive UM:SCC-47<sup>15</sup> and an HPV-negative UPCI:SCC040<sup>16</sup> OTSCC cell lines. Prior to the functional assay, the concentration of siRNA required for silencing, extent of *CASP8* knockdown and cisplatin sensitivity ( $IC_{50}$ ) in both these cell lines was tested (Additional file 10). The invasion of cells was greater in both UM:SCC-47 and UPCI:SCC040 cell lines when *CASP8* was knocked down (Figure 4A). To analyze the effect of caspase-8 on the migration property of cells, scratches were made on the confluent monolayer of cells and the wound closure area was measured at different time points (0hr, 15hr, 23hr & 42hr, Figure 4B). The wound closure was faster in *CASP8* knockdown HPV-negative cells compared to the HPV-positive cells. At 15hr, 23hr and 48hrs, about 65%, 90% and 100% of the wound got closed respectively in the HPV-negative cell line compared to 50%, 70% and 85% respectively during the same time period in the HPV-positive cell lines (Figure 4B). siRNA knockdown of *CASP8* rescued the chemosensitivity caused by cisplatin treatment as evident by the MTT (3-(4,5-dimethylthiazol-2-yl)-2,5-diphenyltetrazolium bromide) survival assay (Figure 4C). Interestingly, we found that the extent

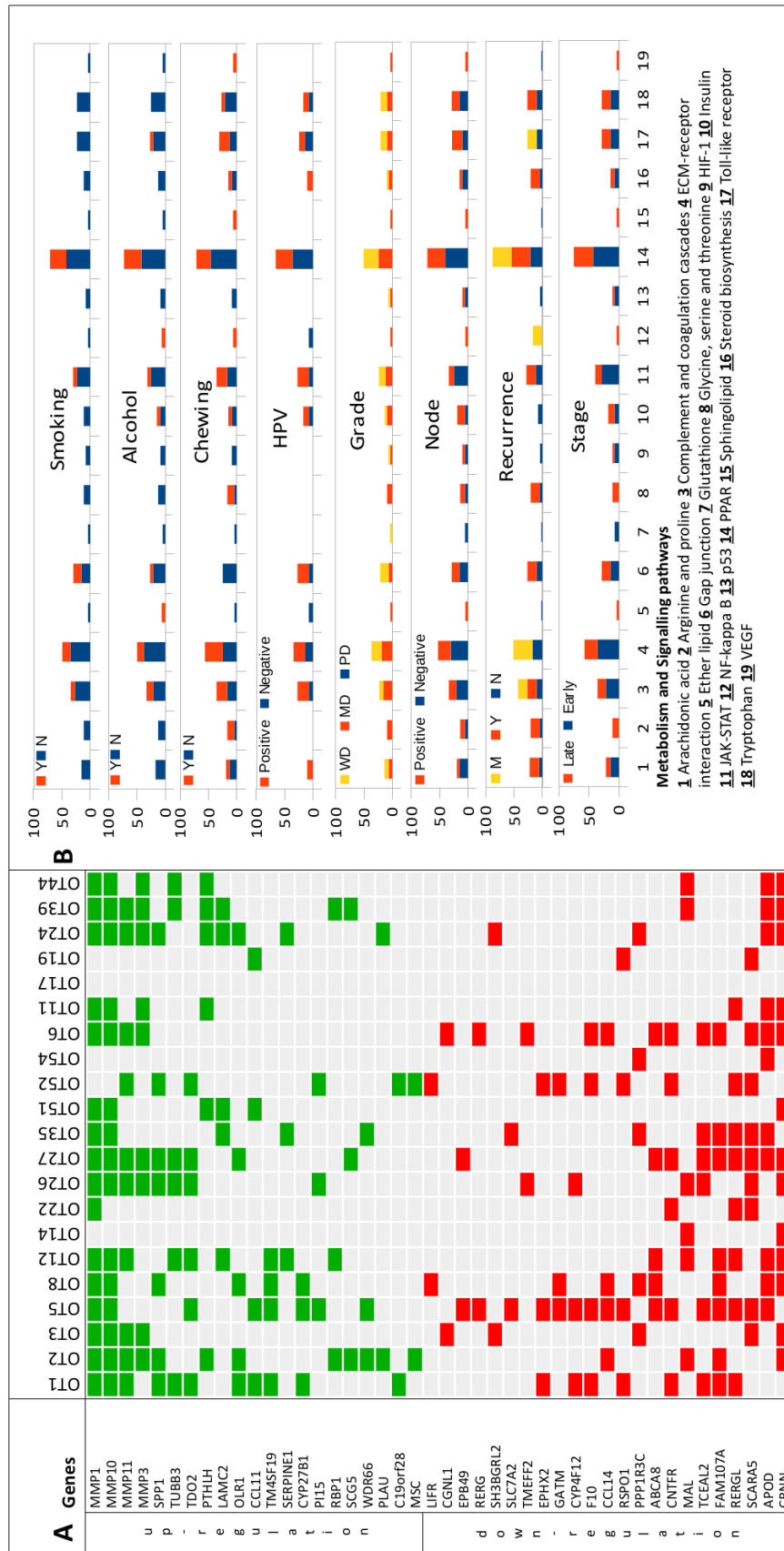
of rescue is greater in the HPV-negative cell line compared to the HPV16-positive one.

### Tumor recurrence prediction using random forests

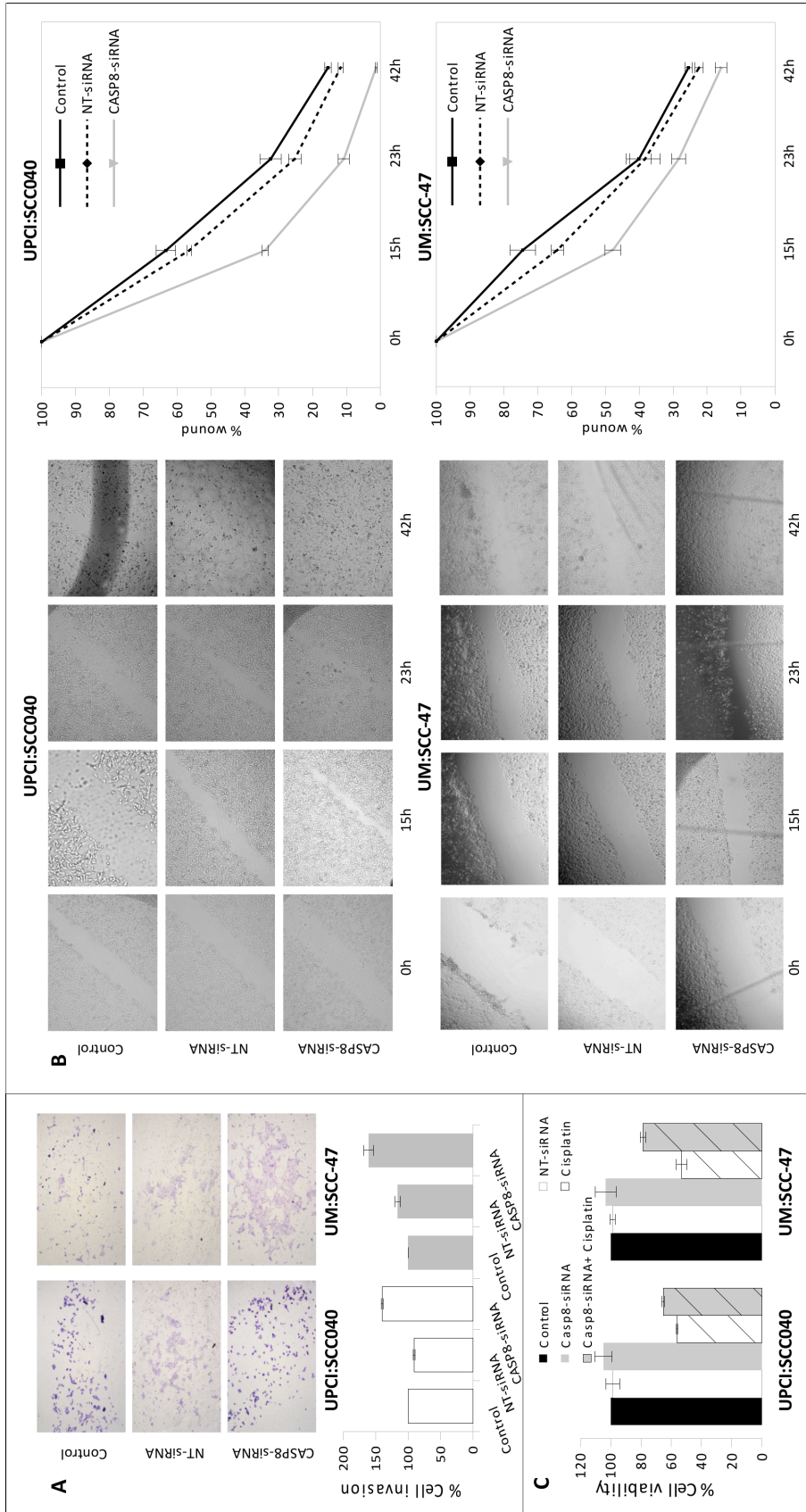
After cataloging the significantly altered genes in OTSCC, we wanted to see whether there is a relationship between the altered genes and loco-regional recurrence of tumors and metastasis. In order to do this, we used an ensemble machine-learning method implemented by variable elimination using random forests<sup>11</sup> (Figure 5). We used multiple testing correction and the 0.632 bootstrapping method<sup>17</sup> to estimate false positives. We discovered a 38-gene minimal signature that discriminated between the non-recurring, loco-regionally recurring and distant metastatic tumors (Figure 5). The .632+ bootstrap errors, indicative of prediction specificity, varied across non-recurrent, recurrent and distant metastatic tumors. The median error was low (0.03) and intermediate (0.3) for the non-recurrent and the loco-regionally recurrent categories respectively but was relatively higher (1.0) for the metastatic tumors. The errors were proportional to the number of representative samples within each category.

### Major signaling pathways implicated in OTSCC

We looked at significant pathways altered in OTSCC, taking into account all the molecular changes in tumors and found apoptosis, HIF, Notch, mTOR, p53, PI3K/Akt, Wnt and Ras to be some of the key signaling pathways affected in OTSCC (Figure 6). In addition, histone methylation, cell cycle/immunity and mRNA splicing processes were also affected. The complete list of pathways is provided in Additional file 11.

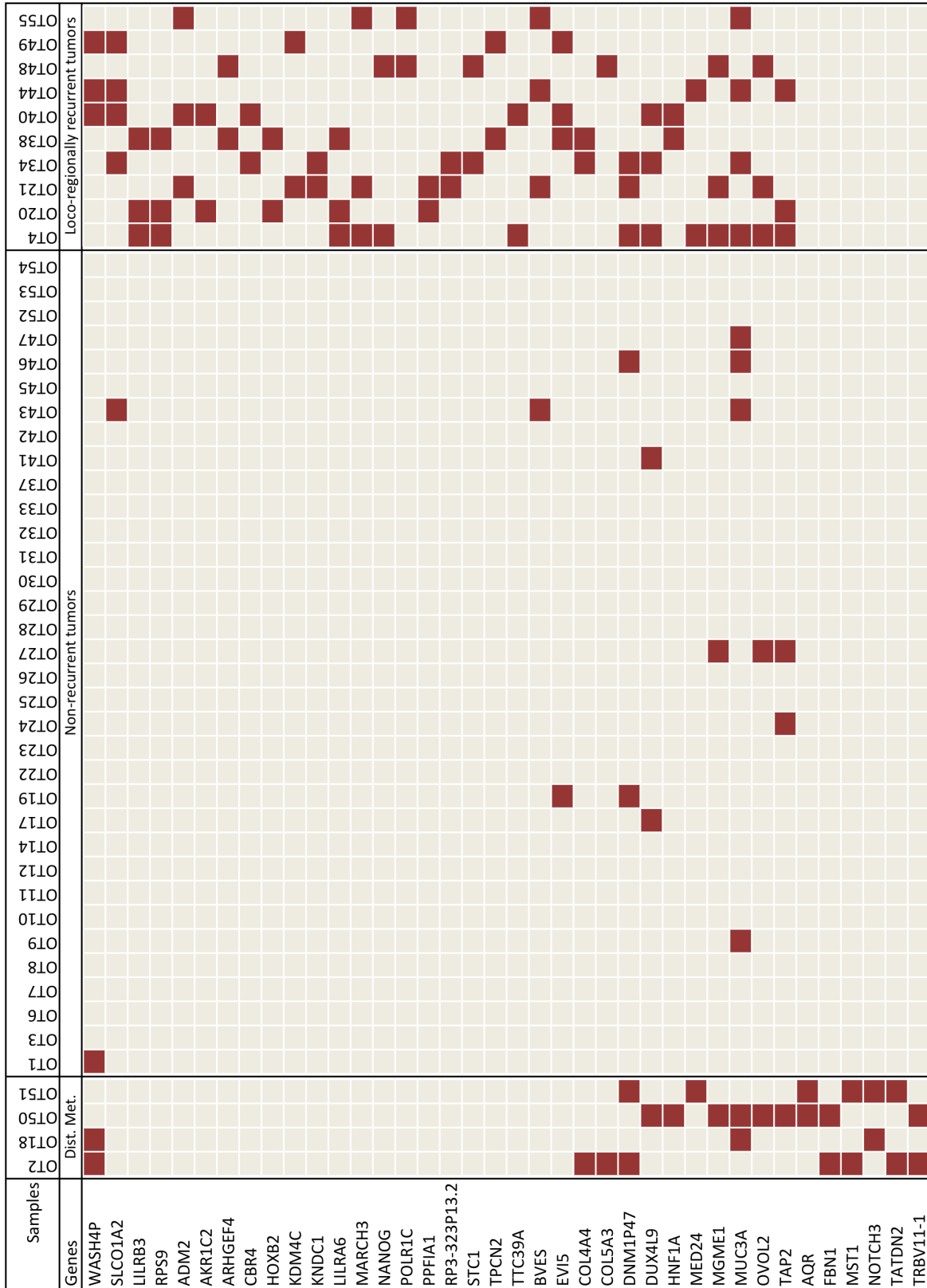


**Figure 3. Differentially expressed genes, affected pathways and their relationship with clinical and epidemiological parameters.** A. Expression changes (green – up- and red – down-regulation) representing significantly differentially expressed genes in tumors. B. Stacked histograms representing relative patient fraction (%) for each of the 19 cancer-associated pathways and their relationship with clinical and epidemiological parameters.

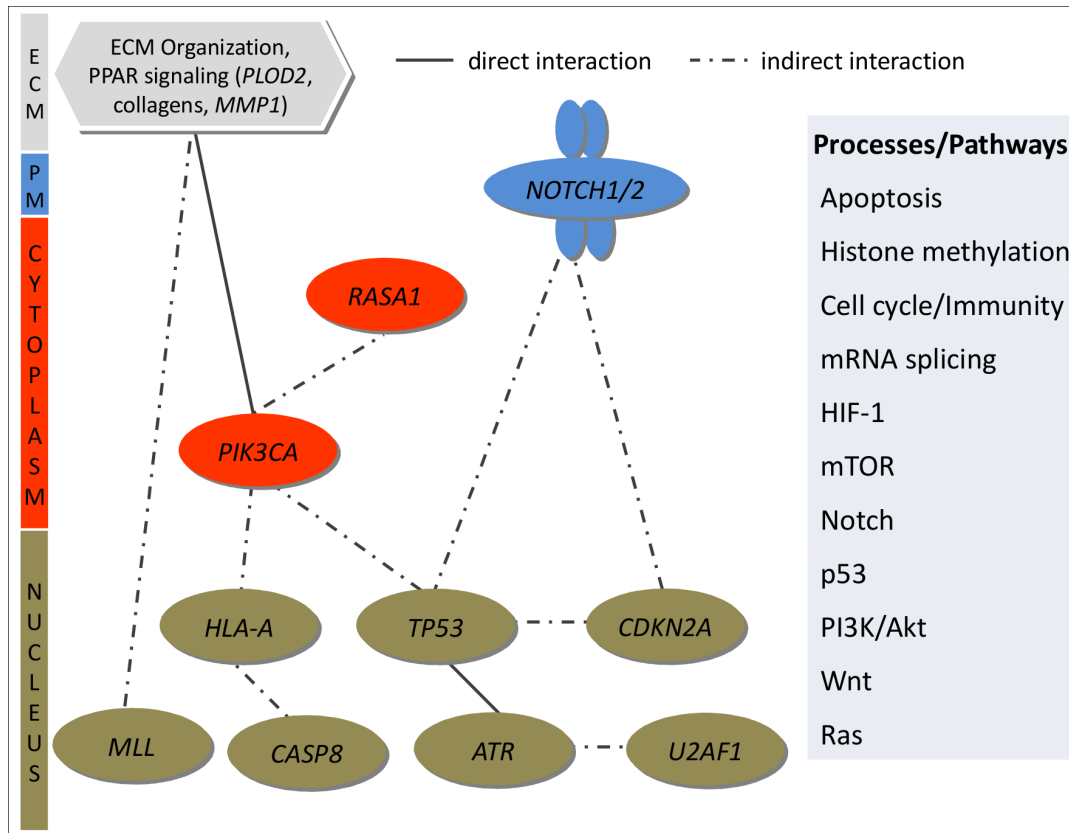


**Figure 4. Role of CASP8 in HPV-positive and HPV-negative OTSCC cell lines.** Results from **A**. Matrigel cell invasion assay (plotted with respect to the control cells), **B**. Wound healing assay, and **C**. MTT cell survival assay (plotted with respect to the control cells) in UPCI:SCC040 (HPV-negative) and UM:SCC-47 (HPV-positive) cell lines.





**Figure 5. A minimal gene signature for tumor recurrence.** Genes harboring somatic variants (in color) that are a part of the minimal signature set for tumor recurrence derived from random forest analyses are used.



**Figure 6. Significantly affected pathways in OTSCC.** Genes harboring significant somatic variants and with expression changes in tumors were used in Cytoscape to derive a set of important signaling pathways implicated in OTSCC.

## Discussion

Squamous cell carcinomas of the oral tongue are an aggressive group of tumors with a higher incidence in the younger population ( $\leq 50$  yrs), which spread early to lymph nodes and have a higher regional failure compared to gingivo-buccal cases<sup>8–10</sup>. Previous sequencing studies<sup>3–5,7</sup> grouped oral tongue tumors with tumors from the oral cavity, but a rise in the incidence of oral tongue tumors, especially among younger people who never smoked, consumed alcohol or chewed tobacco warrants further investigation of this subgroup of oral tumors. Additionally, the role of HPV in oral tongue tumors, unlike in oropharyngeal cases<sup>18–20</sup>, is not well understood both in terms of incidence and prognosis. A meta-analysis of HPV-positive HNSCC tumors from multiple studies conducted at multiple locations concluded that HPV-positive patients, especially in oropharynx, have improved overall and disease-specific survival<sup>21</sup>. A past study has presented data that the HPV incidence in the oral tongue is low<sup>22</sup> and some argue against any link between HPV infection and aggressive oral tongue tumors<sup>23</sup>. Although there is no consensus on rate of HPV incidence among oral tongue tumor patients, it is generally believed that it is low compared to

oropharyngeal tumors. However, some studies in the past<sup>24</sup>, albeit from a different geography, established a much higher rate of HPV infection in oral tongue tumors.

We applied stringent filtering steps and used multiple annotation tools to come up with a list of 19 cancer-associated genes that harbored somatic variants in OTSCC. Most of these genes were also found in other studies, including the recent TCGA HNSCC study<sup>7</sup>, with some notable differences. A comparison of somatic variants discovered in all HSNCC studies, including the current study, is provided in [Additional file 12](#). The frequency for somatic changes in *CASP8*, *NOTCH1*, *CDKN2A* and *FAT1* genes in previous studies<sup>3–7</sup> were, 4–34%, 13–18%, 2–16% and 13–50%, respectively. This is different from what we found in the current study (8%, 4%, 6% and 2% respectively for the same genes). This may partly be attributed to the total number of tumor samples used in different studies but may also be due to a unique pattern of mutations specific to the oral tongue subsite. It appears from our study that the latter is the case. For example, in one of the earlier studies<sup>6</sup> involving similar number of patients as in the current study, *CASP8* and *FAT1* were mutated

in 34% and 50% of the patients but we find the frequency to be 8% and 2%, respectively. In some earlier studies, it was not possible to categorize and identify oral tongue-specific variants as the sites were classified under oral cavity<sup>3,4</sup>.

Although the somatic variants discovered from our study appear to be distributed uniformly across the genome, the significant copy number variation events are more concentrated in chromosomes 6–9 and 11 (Figure 1F and Additional file 8). One of the most important genes harboring somatic mutations discovered in our study is *CASP8*, the product for which derived from the precursor Procaspase-8. Caspase-8 is an important protein implicated in both apoptotic and non-apoptotic pathways<sup>14</sup>. Recent analysis from the TCGA study<sup>7</sup> suggests that mutations in *CASP8* co-occur with mutations in *HRAS*, and are mutually exclusive with amplifications in the *FADD* gene. In our functional studies, the most important observation was that caspase-8 shows different effects in HPV-positive and HPV-negative cells, the effect being more pronounced in HPV-negative cells (Figure 4). Therefore, it is possible that HPV-negative tumors activate a completely different set(s) of pathways and/or may have different chemosensitivity towards drugs than the HPV-positive tumors. It was shown previously that HPV-positive HNSCC cell lines are resistant to TRAIL (tumor necrosis factor-related apoptosis-inducing ligand) and treatment of cells with the proteasome inhibitor bortezomib sensitizes HPV-positive cells towards TRAIL-induced cell death mediated by caspase-8<sup>25</sup>. The E6 protein of HPV interacts with the DED domain of caspase-8 and induces its activation by recruiting it to the nucleus<sup>26</sup>. Our observation on the role of caspase-8-mediated apoptosis being more pronounced in the HPV-negative OTSCC cell line is similar to the observation on the role of *CASP8* in HPV-negative patients made earlier in TCGA study<sup>7</sup>. Taken together, genes including *CASP8* regulate key pathways (Figure 6) that might play important role in the development of tumors in oral tongue.

In the past, several large sequencing studies have been undertaken in HNSCC<sup>3–5,7</sup> that contained very few HPV-positive oral tongue patients. Our study is based on a unique patient cohort and attempts to link molecular signature with different clinical and epidemiological parameters. The prevalence of HPV is very high in oral tongue tumors from India, including in our cohort, compared with studies using cohorts elsewhere. Currently we are completing a larger study on HPV prevalence in different head and neck subsites and we don't see the same high prevalence of HPV in non-oral tongue tumors in the oral cavity, for example in buccal tumors, in one of our studies (Palve *et al.*, unpublished observation from upcoming publication). The exact reason for this high prevalence is not known. Additionally, our observation that some HPV-positive patients harbored *TP53* mutations is counter-intuitive, owing to the fact that E6 is known to block p53. Although we don't know the reason behind this, there is a possibility that HPV-positive tumors harboring *TP53* mutations represent a unique class of tumors and it will be interesting to see if those tumors recur early or late compared to the HPV-positive tumors that have wild type p53 function. Therefore, this study is unique in that respect.

Identifying a signature for tumor recurrence prospectively in primary tumors may add significant advantage to disease management.

In order to do this, we used a machine-learning method using the molecular changes identified in this study, in three batches of primary tumors; non-recurring, loco-regionally recurring and tumors with distant metastasis. We identified a 38-gene signature to be significantly distinguishing the three groups. The bootstrapping error for the non-recurring and the loco-regionally recurring groups were low (N = 34, .632 error = 0.03 and N = 10, .632 error = 0.3 respectively) but not in the metastatic tumor group (N = 4, .632 error = 1). This was due to the small sample numbers (N = 4) in the metastatic category, justifying the need for a larger sample set to validate the signature. The 38 gene signature identified in our study, however, needs to be validated in a much larger cohort in the future to achieve its true potential as a prognostic panel in OTSCC.

Finally, we were keen to see if the current study leads to finding novel drug candidates in OTSCC. We based our assumption on the fact that genome-wide somatic variant discovery in tumors may give rise to possibilities of finding novel drug targets/candidates or may lead us to use existing drugs prescribed for other indications. In an attempt to identify if any of the significantly altered genes found in the current study could potentially act as drug targets, we screened for available drugs against them. We found drugs against three targets out of which two have undergone at least one clinical trial (Additional file 13).

## Methods

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

Informed consent was obtained voluntarily from each patient enrolled in the study. Ethical approval (NHH/MEC-CL/2014/197) was obtained from the Institutional Ethics Committees of the Mazumdar Shaw Medical Centre. Matched control (blood and/or adjacent normal tissue) and tumor specimens were collected and used in the study. Patients diagnosed and treated at the cancer clinic of the Mazumdar Shaw Medical Centre for oral tongue tumors were subjected to a screening procedure before being enrolled in the study. Only those patients, where the histological sections confirmed the presence of squamous cell carcinoma with at least 70% tumor cells in the specimen, were used in the current study. At the time of admission, patients were asked about the habits (chewing, smoking and/or alcohol consumption). Fifty treatment-naïve patients who underwent staging according to AJCC criteria, and curative intent treatment as per NCCN guideline involving surgery with or without post-operative adjuvant radiation or chemo-radiation at the Mazumdar Shaw Medical Centre were accrued for the study (Additional file 1). Post-treatment surveillance was carried out by clinical and radiographic examinations as per the NCCN guidelines.

### HPV detection

HPV was detected by using q-PCR (Applied Biosystems 9700) using HPV16- and HPV18-specific TaqMan probes and primers, and digital PCR (BioRad QX100) using TaqMan probes and primers to detect HPV in primary tumor samples. The primers, probes and cycling conditions for q-PCR and ddPCR were as follows. For q-PCR: 5' GCA CAG AGC TGC AAA CAA CT 3'; 3' GCA TAA ATC CCG AAA AGC AA 5'; probe-ATTAGAATGTGTACTGCAAGCA-FAM-BHQ and 5' TGA CAC TGT GCC TCA ATC CT 3'; 3' AGA GCC ACT TGG AGA GGG AG 5';

Probe-TGCCTGCTTCACCTGGCAGC-VIC-BHQ for HPV16 and HPV18 respectively. The cycling conditions for q-PCR were: 95°C : 3 min, 95°C : 30 sec, 55°C for HPV16 and 60°C for HPV18 : 30 sec, 72°C : 30 sec for 40 cycles. For ddPCR: 5' ACT GTC AAA AGC CAC TGT GT 3'; 3' GCT GGG TTT CTC TAC GTG TT 5' and Probe-AGGGGTCGGTGGACCGGTTCGATGT-FAM-BHQ for HPV16. The cycling conditions for ddPCR were: 95°C: 10 min, 95°C : 315 sec, 55°C : 20 sec for 40 cycles.

### Exome sequencing, read QC, alignment, variant discovery and post-processing filters

Exome libraries were prepared using Agilent SureSelect, Illumina TruSeq and Nextera exome capture kits (Additional file 14) following manufacturers' specifications. Paired end sequencing was performed using HiSeq 2500 or GAIIX and raw reads were generated using standard Illumina base caller (HCS 2.0). Read pairs were filtered using *in house* scripts (Additional file 15 and Additional file 16) and only those reads having  $\geq 75\%$  bases with  $\geq 20$  phred score and  $\leq 15$  Ns were used for sequence alignment against human hg19 reference genome using NovoAlign (v3.00.05)<sup>27</sup>. The aligned files (\*.sam) were processed using Samtools (v0.1.12a)<sup>28</sup> and only uniquely mapped reads from NovoAlign were considered for variant calling. The alignments were pre-processed using GATK (v1.2-62)<sup>29</sup> in three steps before variant calling. First, the indels were realigned using the known indels from 1000G (phase1) data. Second, duplicates were removed using Picard (v1.39). Third, base quality recalibration was done using CountCovariates and TableRecalibration from GATK (v1.2-62), taking into account known SNPs and indels from dbSNP (build 138). Finally, UnifiedGenotyper from GATK (v2.5-2) was used for variant calling, using known SNPs and indels from dbSNP (build 138). Raw variants from GATK were filtered to only include the PASS variants (standard call confidence  $\geq 50$ ) within the merged exomic bait boundaries. Two out of 50 tumor samples did not confirm to the QC standards, therefore excluded from all further analyses. Therefore, all the downstream analyses were restricted to 48 primary tumors. The variants were further flagged as novel or present in either dbSNP138 or COSMIC (v67) databases, based on their overlap. In addition to GATK, we also used Dindel<sup>30</sup> to call indels. Both GATK and Dindel calls were filtered for microsatellite repeats (flagged as STR). The raw variant calls were used to estimate frequencies of nucleotide changes and transition:transversion (ti/tv) ratios. Exome-filtered PASS variants specific to the tumor samples, with respect to both location and actual call, were retained as somatic variants, which were further filtered to exclude variants where the region bearing the variant was not callable in the matched control sample, and those where the matched control sample had even one read covering the variant allele.

Scripts used to perform various filtering steps are provided in Additional file 16. The numbers of raw reads, after QC, alignment statistics, numbers of variants pre- and post-filters are provided in Additional file 2.

### Detection of cross-contamination and identification of significant somatic variants

We estimated cross-contamination using ContEst (June 2013)<sup>31</sup> in the tumor samples (Additional file 16). Locus-wise and gene-wise

driver scores were estimated by CRAVAT<sup>32</sup> using the head and neck cancer database with the CHASM<sup>33</sup> analysis option. Genes with a CHASM score of at least 0.35 were considered significant for comparison with other functional analyses (Additional file 16). Somatic mutations were normalized with respect to the exome bait size (MB) to calculate the somatic mutation frequency per MB.

### Annotation and functional analyses of variants

Annotation and functional analyses of somatic variants was performed using IntoGen (web version 2.4)<sup>34,35</sup>, MutSigCV (v1.3.01)<sup>36,37</sup> and MuSiC2 (v0.1)<sup>37</sup>. Somatic variants, filtered to contain only those callable in the matched normal but not covered by any read in the control samples (VCF), were used for IntoGen with the 'cohort analyses' option. We also ran MutsigCV1.3 with these variants using coverage from un-filtered variants of all tumor samples (Additional file 16). Pooled alignments for all normal and tumor samples (BAM), each, along with pooled variants for all normal samples (MAF) were analyzed using MuSiC2 to calculate the background mutation rates (bmrs) for all genes, and identify a list of significantly mutated genes (*p*-value of convolution test  $\leq 0.05$ ; Additional file 16). A condensed list of 19 genes, common between at least two analyses was compiled (Figure 1D).

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

High quality DNA (200ng), quantified by Qubit 2.0 (Invitrogen), was used as the starting material for whole-genome genotyping experiments following the manufacturer's specifications. Briefly, the genomic DNA was denatured at room temperature (RT) for 10 mins using 0.1N NaOH, neutralized and used for whole genome amplification (WGA) under isothermal conditions, at 37°C for 20 hrs. Post-WGA, the DNA was enzymatically fragmented at 37°C for 1hr. The fragmented DNA was precipitated with isopropanol at 4°C and resuspended in hybridization buffer. The samples were then denatured at 95°C for 20 mins, cooled at RT for 30 mins and 35 $\mu$ l of each sample was loaded onto the Illumina HumanOmni 2.5-8 beadchip for hybridization (20hrs at 48°C) in a hybridization chamber. The unhybridized probes were washed away and the Chips (HumanOmni 2.5-8 v1.0 and v1.1, Additional file 2) were prepared for staining, single base extension and scanning using Illumina's HiScan system.

We filtered the SNP locations to retain only those, called without any error, contained within the exome boundaries as per the sequencing baits, and which were callable (covered by at least five sequencing reads). At these locations, we estimated the overlap for individual SNP calls, i.e., chr/pos/ref/alt and for no calls; i.e., chr/pos/ref/ref; between sequencing and array platforms (Additional file 16).

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

CNVs and LOHs were identified using cnvPartition 3.1.6 plugin in Illumina GenomeStudio v2011.1, with default settings except for a minimum coverage of at least 10 probes per CNV/LOH with a confidence score threshold of at least 100 (Additional file 17). Somatic CNVs and LOHs were extracted by filtering out any region common to CNVs and LOHs detected in its matched control. Somatic CNVs and LOHs were further filtered with respect to common and

disease-related CNVs and LOHs using CNVAnnotator<sup>38</sup>. Overlaps with common CNVs and LOHs were discarded, reporting only the overlaps with disease-related, and novel CNVs and LOHs. We categorized the CNVs and LOHs within each cytoband and reported those with an occurrence in at least 10% of the patient samples.

### Gene expression assay

Gene expression profiling was carried out using Illumina HumanHT-12 v4 expression BeadChip (Illumina, San Diego, CA) in tumor and matched normal tissues (Additional file 9) following manufacturer's specifications. Total RNA was extracted from 20mg of tissue using PureLink RNA (Invitrogen) and RNeasy (Qiagen) Mini kits. RNA quality was checked using Agilent Bioanalyzer 2100 using RNA Nano6000 chip. Samples with poor RNA integrity numbers (RIN) (<7), indicating partial degradation of RNA, were processed using Illumina WGDSL assay as per manufacturer's recommendations. The RNA samples with no degradation were labelled using Illumina TotalPrep RNA Amplification kit (Ambion) and processed according to the array manufacturer's recommendations. Gene expression data was collected using Illumina's HiScan and analyzed with the GenomeStudio (v2011.1 Gene Expression module 1.9.0) and all assay controls were checked to ensure quality of the assay and chip scanning. Raw signal intensities were exported from GenomeStudio for pre-processing and analyzed using R further.

Gene-wise expression intensities for tumor and matched control samples from GenomeStudio were transformed and normalized using VST (Variance Stabilizing Transformation) and LOESS methods, respectively, using the R package lumi<sup>39</sup>. The data was further batch-corrected using ComBat<sup>40</sup> (Additional file 16). The pre-processed intensities for tumor and matched control samples were subjected to differential expression analyses using the R package, limma<sup>41</sup> (Additional file 16). Genes with significant expression changes (adjusted *P* value ≤ 0.05) and fold change of at least 1.5 were followed up with further functional analyses.

### Recurrence prediction using random forests

We used presence or absence of somatic mutations/indels data in the entire set of genes for all the OTSCC patients, along with their recurrence patterns as training set for the random forests<sup>11</sup> analyses using the varSelRF package in R. This method performs both backward elimination of variables and selection based on their importance spectrum, and predicts recurrence patterns in the same set by iteratively eliminating 2% of the least important predictive variables until the current OOB (out-of-bag) error rate becomes larger than the initial or previous OOB error rates. In order to understand the specificity of the best minimalistic predictors of tumor recurrence, we estimated the 0.632+ error rate<sup>17</sup> over 50 bootstrap replicates. We used the varSelRFBoot function from the varSelRF Bioconductor package to perform bootstrapping. The .632+ method is described by the following formula:

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

where  $Err^{(.632')}$ ,  $Err^{(.632)}$ ,  $Err^{(1)}$  and  $err$  are errors estimated by the .632+ method, the original .632 method, *leave-one-out* bootstrap

method and  $err$  represents the error.  $R'$  represents a value between 0 and 1. Another popular error correction method used is *leave-one-out* bootstrap method. The .632+ method was designed to correct the upward bias in the *leave-one-out* and the downward bias in the original .632 bootstrap methods.

For all iterations of all random forest analyses, we confirmed that the variable importance remained the same before and after correcting for multiple hypotheses comparisons using pre- and post-Benjamin-Hochberg FDR-corrected *P* values. R commands for variable elimination using random forests, 0.632+ bootstrapping and re-computing importance values after multiple comparisons testing are provided in Additional file 16.

### Pathway analyses

Consensus list of genes from analysis, filtering and annotation of variant calls and from differential expression analysis using whole genome micro-arrays, were mapped to pathways using the web version of Graphite Web<sup>42</sup> employing KEGG and Reactome databases. The network of interactions between genes was drawn originally using CytoScape (v3.1.1)<sup>43</sup> using the .sif file created by Graphite Web (Additional file 16).

### Data visualization

We used Circos (v0.66)<sup>44</sup> (Additional file 18 and Additional file 19) for multi-dimensional data visualization. Additionally, we used the cbiportal protal (<http://www.cbiportal.org/>) to visualize variants within the 19 genes harboring significant variants. All of the mandatory fields accepted by Mutation Mapper were provided for select genes from our study to create structural representations for each gene including domains. Such diagrams from our study, the HNSCC study and all cancer studies from TCGA were collated using the image-editing tool, GIMP (v2.8.0) ([www.gimp.org](http://www.gimp.org)). SNPs and indels were visualized for each individual tumor sample using IGV (v1.5.54)<sup>45</sup>, along with the reads supporting variants (Additional file 20).

### Validation of somatic variants using Sanger sequencing

Primers were designed using the NCBI primer designing tool ([http://www.ncbi.nlm.nih.gov/tools/primer-blast/index.cgi?LINK\\_LOC=BlastHome](http://www.ncbi.nlm.nih.gov/tools/primer-blast/index.cgi?LINK_LOC=BlastHome)) and used in Sanger sequencing for validation. The sequences of all primers (IDT) used for validation is provided in Additional file 21. We tested the specificity of the designed primers using UCSC's tool, *In Silico* PCR. The variant-bearing region was amplified by using specific primers and used in Sanger sequencing (Additional file 14). The somatic variants were confirmed by sequencing in the entire tumor and matched control DNA set used for the exome sequencing followed by further validation in 60 additional tumor samples (Additional file 1B).

### Cell culture and knockdown of CASP8 gene

The human OTSCC cell lines UPCI:SCC040 (gift from Dr. Susan Gollin, University of Pittsburgh, PA, USA)<sup>16</sup> and UM-SCC47 (gift from Dr. Thomas Carey, University of Michigan, MI, USA)<sup>15</sup> were used in the study. All the cells were maintained in Dulbecco's Modified Eagles' Media (DMEM) supplemented with 10% FBS, 1% MEM nonessential amino acids solution & 1% penicillin/streptomycin mixture (Gibco) at 37°C with 5% CO<sub>2</sub> incubator.

We performed the siRNA-based knockdown using UPCI:SCC040 and UM:SCC47 cell lines for *CASP8* gene. The expression of Caspase-8 was transiently knocked down using ON-TARGETplus Human *CASP8* smart pool siRNA (L-003466-00-0010; Dharmacon) along with an ON-TARGETplus Non-targeting siRNA (D-001810-01-20; Dharmacon). The transfection efficiency for the two cell lines (UPCI:SCC040 and UM:SCC47) were optimized using siGLO Red Transfection Indicator (D-001630; Dharmacon). The siRNA duplexes were transfected using Lipofectamine-2000 according to the manufacturer's instructions (Invitrogen). The siRNA-oligo complexes medium was changed 8 hrs post transfection. The efficiency of transfection along with the mRNA expression was analyzed at 24 and 48 hrs post transfection by qRT-PCR. The specific down-regulation of *CASP8* was confirmed by three independent experiments.

### RNA isolation and quantitative real-time PCR

RNA was extracted from cell pellets and tissues using RNeasy Mini kit spin columns (Qiagen) following manufacturer's protocol. Genomic DNA contamination was removed by RNase-Free DNase Set (Qiagen) and the total RNA was eluted in nuclease free water (Ambion). The RNA samples were estimated using Qubit 2.0 fluorometer (Invitrogen) and the integrity was checked by gel electrophoresis. The RNA samples were stored at  $-80^{\circ}\text{C}$  until further used. The cDNA was synthesized with 400ng total RNA, using a SuperScript-III first strand cDNA synthesis kit, and following the manufacturer's instructions (Invitrogen). The cDNA was then subjected for quantitative real-time PCR (q-RT-PCR) using KAPA SYBR FAST qPCR Master Mix (KK4601, KAPA). The primer pairs used for testing the expression of caspase-8 in q-RT-PCR were, forward 5'-ATGATGACATGAACCTGCTGGA-3' and reverse 5'-CAGGCTCTTGTTGATTTGGGC-3'. The amplification was done on Stratagene MX300P real time machine. The cycling conditions were: step-1  $95^{\circ}\text{C}$ -3min, step-2  $95^{\circ}\text{C}$ -3sec, step-3  $60^{\circ}\text{C}$ -60sec then repeat steps 2–3 for 40 cycles following dissociation curve at  $60^{\circ}\text{C}$ -60sec,  $95^{\circ}\text{C}$ -1min,  $60^{\circ}\text{C}$ -60sec.

To normalize inter-sample variation in RNA input, the expression values were normalized with GAPDH. All amplification reactions were done in triplicates, using nuclease free water as negative controls. The differential gene expression was calculated by using the comparative  $C_{\text{T}}$  method of relative quantification<sup>46</sup>.

### Assessment of cell viability (confirm)

MTT cell proliferation assay was performed as per manufacturer's instructions (Sigma) to assess cell viability. Briefly, cells were seeded on 96-well plates containing DMEM with 10% FBS & incubated overnight. After treatment with 0.1% DMSO (vehicle control), or Cisplatin for 48 hrs, medium was changed and 100  $\mu\text{l}$  of MTT solution (1mg/ml) was added to each well. The cells were further incubated for 4hrs at  $37^{\circ}\text{C}$ . The formazan crystals in viable cells were dissolved by adding 100 $\mu\text{l}$  of dimethyl sulfoxide (DMSO) (Merck). The absorbance was recorded at 540 nm using reference wavelength of 690 nm on micro plate reader (Tecan Systems). Data were normalized to vehicle treatment, and the cell viability was

calculated using GraphPad Prism software (version 4.03; La Jolla, CA). All the experiments were performed in triplicates.

### Wound healing assay

Cells were cultured up to 80% confluency in 12 well plates; serum-starved for 24 hrs and then wounded using a 200 $\mu\text{l}$  pipette tip. The wound was washed with 1 $\times$  PBS and the cells were grown in DMEM containing 10% FBS. Cells were imaged at 10 $\times$  magnification at 0 hr, 15 hrs, 23 hrs and 42 hrs. For each well, three wounds were made and the migration distance was photographed and measured using Carl Zeiss software (Zeiss). Each experiment was performed in triplicates.

### Matrigel invasion assay

The ECM gel (E1270, Sigma) was thawed overnight at  $4^{\circ}\text{C}$  and plated at requisite concentrations (for UPCI:SCC040: 1.5mg/ml and UM:SCC047: 2mg/ml) onto the transwell inserts and incubated overnight in the  $\text{CO}_2$  incubator at  $37^{\circ}\text{C}$  with 5%  $\text{CO}_2$ . Cells were serum-starved for overnight, harvested, counted and seeded (UPCI:SCC040: 50,000 cells and UM:SCC047: 20,000 cells per well) on top of the matrigel transwell-inserts (2 mg/ml) in serum-free medium as per manufacturer's specifications (Sigma). D-MEM containing 10% FBS and 1% NEAA was added to the lower chamber. The 24-well plates containing matrigel inserts with cells were incubated in  $37^{\circ}\text{C}$  incubator for 48 hrs. At the end of incubation time, cells in the upper chamber were removed with cotton swabs and cells that invaded the Matrigel to the lower surface of the insert were fixed with 4% paraformaldehyde (Merk Milipore), permeabilized with 100% methanol, stained with Giemsa (Sigma), mounted on glass slides with DPX mounting agent and counted under a light microscope (Zeiss). Each experiment was performed in triplicates.

### Conclusions

We have catalogued genetic variants (somatic mutations, indels, CNVs and LOHs) and transcriptomic (significantly up- and down-regulated genes) changes in oral tongue squamous cell carcinoma (OTSCC) and used those in an integrated approach linking genes harboring somatic variants with common risk factors like tobacco and alcohol; clinical, epidemiological factors like tumor grade and HPV; and tumor recurrence. We found *CASP8* gene to be significantly altered and play an important role in apoptosis-mediated cell death in an HPV-negative OTSCC cell line. Finally, we present data towards a minimal gene signature that can predict tumor recurrence.

### Author contributions

BP: conceived, designed and supervised the study, wrote the manuscript; NMK: analyzed the data and wrote the manuscript; SG: analyzed the data and critically read the manuscript; SP, PJ, CK, VKP: analyzed the data; VP, GS, AS: produced data on *CASP8* functional analysis; LV, AKH, MP: produced sequencing data; KD and JN: produced array data; GS, AS, VK and MAK: provided clinical data, clinical input and associated clinical information. All authors have seen and agreed to the final content of this manuscript.

## Competing interests

None.

(Ref No:18(4)/2010-E-Infra., 31-03-2010) and Department of IT, BT and ST, Government of Karnataka, India (Ref No:3451-00-090-2-22).

## Grant information

Research presented in this article is funded by Department of Electronics and Information Technology, Government of India

*I confirm that the funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

## Additional files

**Additional files for Krishnan *et al.*, 2015 'Integrated analysis of oral tongue squamous cell carcinoma identifies key variants and pathways linked to risk habits, HPV, clinical parameters and tumor recurrence'.**

All the additional files can be downloaded from Figshare (<http://figshare.com/s/c928faa66f2a11e586d506ec4b8d1f61>).

**Additional file 1.** Patient details used in the study.

**Additional file 2.** Sequencing, read QC, alignment and variant calls, OMNI SNP genotyping array validation.

**Additional file 3.** Validation using capillary gel electrophoresis based on Sanger sequencing in A. discovery set and B. validation set.

**Additional file 4.** The ratio of transitions to transversions (ti/tv) was estimated using the exome-filtered GATK PASS variants for tumor and matched control samples. The dotted lines depict the respective median ti/tv ratios.

**Additional file 5.** Effect of habits, clinical parameters and HPV infection on individual nucleotide change.

**Additional file 6.** Functional annotation of somatic variants using IntOGen, MuSiC2, MutSigCV.

**Additional file 7.** Position and frequency of somatic variants in protein domains found in this study, TCGA HNSCC, and in studies involving all cancer types using mutation mapper in the cBioPortal.

**Additional file 8.** Cytoband-wise representation of CNVs found in all 48 samples along with clinical parameters and patient epidemiology.

**Additional file 9.** Transformed, normalized and batch-corrected intensities following expression assay and results from differential expression analyses.

**Additional file 10.** Functional validation for the role of *CASP8* in OTSCC cell lines.

**Additional file 11.** List of all pathways affected by somatic mutations/indels, copy number variations and expression changes ( $\log_2FC \geq 0.6$ ).

**Additional file 12.** Comparative sample frequency of important variants found in this and other HNSCC studies.

**Additional file 13.** Drug candidates and their targets in head and neck cancer.

**Additional file 14.** Supplementary Methods.

**Additional file 15.** Read QC filter scripts' executable.

**Additional file 16.** Scripts used in the study.

**Additional file 17.** GenomeStudio output of all LOHs and CNVs found using cnvPartition plugin in GenomeStudio.

**Additional file 18.** Circos data and config files.

**Additional file 19.** Circular genomic representation using Circos (v0.66) of LOHs, somatic variants, CNVs with  $\geq 10\%$  frequency of patients bearing them, and genes with significant expression changes ( $|\log_2FC| \geq 0.6$ ).

**Additional file 20.** IGV snapshots of all significantly mutated somatic variants in this study.

**Additional file 21.** Primer sequences used in the Sanger validation study.

[Click here to access the data.](#)

## References

1. Ferlay J, Shin HR, Bray F, *et al.*: **Estimates of worldwide burden of cancer in 2008: GLOBOCAN 2008.** *Int J Cancer.* 2010; **127**(12): 2893–2917.  
[PubMed Abstract](#) | [Publisher Full Text](#)
2. Mishra A, Meherotra R: **Head and neck cancer: global burden and regional trends in India.** *Asian Pac J Cancer Prev.* 2014; **15**(2): 537–550.  
[PubMed Abstract](#) | [Publisher Full Text](#)
3. Stransky N, Egloff AM, Tward AD, *et al.*: **The mutational landscape of head and neck squamous cell carcinoma.** *Science.* 2011; **333**(6046): 1157–1160.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
4. Agrawal N, Frederick MJ, Pickering CR, *et al.*: **Exome sequencing of head and neck squamous cell carcinoma reveals inactivating mutations in *NOTCH1*.** *Science.* 2011; **333**(6046): 1154–1157.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
5. Pickering CR, Zhang J, Yoo SY, *et al.*: **Integrative genomic characterization of oral squamous cell carcinoma identifies frequent somatic drivers.** *Cancer Discov.* 2013; **3**(7): 770–781.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
6. India Project Team of the International Cancer Genome Consortium: **Mutational landscape of gingivo-buccal oral squamous cell carcinoma reveals new recurrently-mutated genes and molecular subgroups.** *Nat Commun.* 2013; **4**: 2873.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
7. Cancer Genome Atlas Network: **Comprehensive genomic characterization of head and neck squamous cell carcinomas.** *Nature.* 2015; **517**(7536): 576–582.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
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**(5): 401–418.  
[PubMed Abstract](#) | [Publisher Full Text](#)
9. Kuriakose M, Sankaranarayanan M, Nair MK, *et al.*: **Comparison of oral squamous cell carcinoma in younger and older patients in India.** *Eur J Cancer B Oral Oncol.* 1992; **28B**(2): 113–120.  
[PubMed Abstract](#) | [Publisher Full Text](#)
10. Pathak KA, Das AK, Agarwal R, *et al.*: **Selective neck dissection (I-III) for node negative and node positive necks.** *Oral Oncol.* 2006; **42**(8): 837–841.  
[PubMed Abstract](#) | [Publisher Full Text](#)
11. Breiman L: **Random Forests.** The Netherlands: Kluwer Academic Publishers; 2001; **45**(1): 5–32.  
[Publisher Full Text](#)
12. Pattnaik S, Vaidyanathan S, Pooja DG, *et al.*: **Customisation of the exome data analysis pipeline using a combinatorial approach.** *PLoS One.* 2012; **7**(1): e30080.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
13. Oberst A, Green DR: **It cuts both ways: reconciling the dual roles of caspase 8 in cell death and survival.** *Nat Rev Mol Cell Biol.* 2011; **12**(11): 757–763.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
14. Fulda S: **Caspase-8 in cancer biology and therapy.** *Cancer Lett.* 2009; **281**(2): 128–133.  
[PubMed Abstract](#) | [Publisher Full Text](#)
15. Lansford CD, Grenman R, Bier H, *et al.*: **Head and neck cancers.** New York: Kluwer Academic Publishers, 2002.  
[Publisher Full Text](#)
16. Telmer CA, An J, Malehorn DE, *et al.*: **Detection and assignment of *TP53* mutations in tumor DNA using peptide mass signature genotyping.** *Hum Mutat.* 2003; **22**(2): 158–165.  
[PubMed Abstract](#) | [Publisher Full Text](#)
17. Tibshirani R, Efron B: **Improvements on Cross-Validation: The .632+ Bootstrap Method.** *J Am Stat Assoc.* 1997; **92**(438): 548–560.  
[Publisher Full Text](#)
18. Kumar B, Cordell KG, Lee JS, *et al.*: **Response to therapy and outcomes in oropharyngeal cancer are associated with biomarkers including human papillomavirus, epidermal growth factor receptor, gender, and smoking.** *Int J Radiat Oncol Biol Phys.* 2007; **69**(2 Suppl): S109–111.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
19. Worden FP, Kumar B, Lee JS, *et al.*: **Chemoselection as a strategy for organ preservation in advanced oropharynx cancer: response and survival positively associated with HPV16 copy number.** *J Clin Oncol.* 2008; **26**(19): 3138–3146.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
20. Fakhry C, Zhang Q, Nguyen-Tan PF, *et al.*: **Human papillomavirus and overall survival after progression of oropharyngeal squamous cell carcinoma.** *J Clin Oncol.* 2014; **32**(30): 3365–3373.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
21. Dayyani F, Etzel CJ, Liu M, *et al.*: **Meta-analysis of the impact of human papillomavirus (HPV) on cancer risk and overall survival in head and neck squamous cell carcinomas (HNSCC).** *Head Neck Oncol.* 2010; **2**: 15.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
22. Dahlgren L, Dahlstrand HM, Lindquist D, *et al.*: **Human papillomavirus is more common in base of tongue than in mobile tongue cancer and is a favorable prognostic factor in base of tongue cancer patients.** *Int J Cancer.* 2004; **112**(6): 1015–1019.  
[PubMed Abstract](#) | [Publisher Full Text](#)
23. Salem A: **Dismissing links between HPV and aggressive tongue cancer in young patients.** *Ann Oncol.* 2010; **21**(1): 13–17.  
[PubMed Abstract](#) | [Publisher Full Text](#)
24. Elango KJ, Suresh A, Erode EM, *et al.*: **Role of human papilloma virus in oral tongue squamous cell carcinoma.** *Asian Pac J Cancer Prev.* 2011; **12**(4): 889–896.  
[PubMed Abstract](#)
25. Bullenkamp J, Rauff N, Ayaz B, *et al.*: **Bortezomib sensitises TRAIL-resistant HPV-positive head and neck cancer cells to TRAIL through a caspase-dependent, E6-independent mechanism.** *Cell Death Dis.* 2014; **5**: e1489.  
[PubMed Abstract](#) | [Publisher Full Text](#)
26. Manzo-Merino J, Massimi P, Lizano M, *et al.*: **The human papillomavirus (HPV) E6 oncoproteins promotes nuclear localization of active caspase 8.** *Virology.* 2014; **450–451**: 146–152.  
[PubMed Abstract](#) | [Publisher Full Text](#)
27. Novocraft: **Novoalign.** 2011.  
[Reference Source](#)
28. Li H, Handsaker B, Wysoker A, *et al.*: **The Sequence Alignment/Map format and SAMtools.** *Bioinformatics.* 2009; **25**(16): 2078–2079.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
29. McKenna A, Hanna M, Banks E, *et al.*: **The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.** *Genome Res.* 2010; **20**(9): 1297–1303.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
30. Albers CA, Lunter G, MacArthur DG, *et al.*: **Dindel: accurate indel calls from short-read data.** *Genome Res.* 2011; **21**(6): 961–973.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
31. Cibulskis K, McKenna A, Fennell T, *et al.*: **ContEst: estimating cross-contamination of human samples in next-generation sequencing data.** *Bioinformatics.* 2011; **27**(18): 2601–2602.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
32. Carter H, Chen S, Isik L, *et al.*: **Cancer-specific high-throughput annotation of somatic mutations: computational prediction of driver missense mutations.** *Cancer Res.* 2009; **69**(16): 6660–6667.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
33. Douville C, Carter H, Kim R, *et al.*: **CRAVAT: cancer-related analysis of variants toolkit.** *Bioinformatics.* 2013; **29**(5): 647–648.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
34. Gundem G, Perez-Llamas C, Jene-Sanz A, *et al.*: **IntOGen: integration and data mining of multidimensional oncogenic data.** *Nat Methods.* 2010; **7**(2): 92–93.  
[PubMed Abstract](#) | [Publisher Full Text](#)
35. Schroeder MP, Gonzalez-Perez A, Lopez-Bigas N: **Visualizing multidimensional cancer genomics data.** *Genome Med.* 2013; **5**(1): 9.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
36. Lawrence MS, Stojanov P, Mermel CH, *et al.*: **Discovery and saturation analysis of cancer genes across 21 tumour types.** *Nature.* 2014; **505**(7484): 495–501.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
37. Dees ND: **MuSIC2.** 2015.
38. Zhao M, Zhao Z: **CNVannotator: a comprehensive annotation server for copy number variation in the human genome.** *PLoS One.* 2013; **8**(11): e80170.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
39. Du P, Kibbe WA, Lin SM: **lumi: a pipeline for processing Illumina microarray.** *Bioinformatics.* 2008; **24**(13): 1547–1548.  
[PubMed Abstract](#) | [Publisher Full Text](#)
40. Johnson WE, Li C, Rabinovic A: **Adjusting batch effects in microarray expression data using empirical Bayes methods.** *Biostat.* 2007; **8**(1): 118–127.  
[PubMed Abstract](#) | [Publisher Full Text](#)
41. Ritchie ME, Phipson B, Wu D, *et al.*: **limma powers differential expression analyses for RNA-sequencing and microarray studies.** *Nucleic Acids Res.* 2015; **43**(7): e47.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
42. Sales G, Calura E, Martini P, *et al.*: **Graphite Web: Web tool for gene set analysis exploiting pathway topology.** *Nucleic Acids Res.* 2013; **41**(Web Server issue): W89–97.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
43. Shannon P, Markiel A, Ozier O, *et al.*: **Cytoscape: a software environment for integrated models of biomolecular interaction networks.** *Genome Res.* 2003; **13**(11): 2498–2504.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
44. Krzywinski M, Schein J, Birol I, *et al.*: **Circos: an information aesthetic for comparative genomics.** *Genome Res.* 2009; **19**(9): 1639–1645.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
45. Thorvaldsdóttir H, Robinson JT, Mesirov JP: **Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration.** *Brief Bioinform.* 2013; **14**(2): 178–192.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
46. Schmittgen TD, Livak KJ: **Analyzing real-time PCR data by the comparative C<sub>t</sub> method.** *Nat Protoc.* 2008; **3**(6): 1101–1108.  
[PubMed Abstract](#) | [Publisher Full Text](#)



# Open Peer Review

Current Referee Status:



Version 1

Referee Report 17 December 2015

doi:10.5256/f1000research.7870.r11116



**Thomas Carey**

Laboratory of Head and Neck Center Biology, Department of Otolaryngology-Head and Neck Surgery, University of Michigan, Ann Arbor, MI, USA

This article reports the the genetic analysis of 50 oral tongue primary tumor cancers treated at a single center, listed as paired sets. The study appears to be well done and thorough with carefully performed informatics.

In the abstract it is stated that 50 paired primary oral tongue cancers were studied, but doesn't state whether the tumors were paired to normal blood or other tissue.

The analysis includes single nucleotide variations, copy number variations, indels, regions with loss of heterozygosity. These somatic variations are linked to clinical parameters. In addition HPV was assessed by q-PCR and found to have a very high rate of positivity, which is surprising given the reported site of oral tongue. It was not clear if all HPV positive cases also had p16 expression. It would be valuable for the reader to know what portions of the HPV transcriptome was assessed. Did the authors examine E6 and E7 oncogene expression? It is questionable whether the HPV is a driver in 22 of 50 tumors. It would be important to verify the activity of HPV by assessing viral oncogene expression and level of expression. In tumors with HPV and mutant p53 it is unlikely that those are driven by HPV, given that p53 mutations in other HPV positive head and neck tumors is extremely rare. It would be helpful to have a table that shows how often mutant p53 was found in the HPV positive tumors and whether those patients had excessive use of carcinogenic substances.

There are some novel findings in this study and some very surprising correlations, for example, CDKN2A mutations were found only in non-smokers. This is highly surprising given that in the US abnormalities of this locus is common in head and neck cancers in smokers. Thus, grouping never smokers and former smokers may not be a fair grouping. What does former smoker mean in this population?

Only smoking and oral tobacco and alcohol use are discussed as etiologic factors. Are there no other factors in this part of India? Is betel nut or oral tobacco mixed with other substances included in the oral tobacco use category? Since these are all oral tongue cancers the information about etiology should be made more clear.

The mutual exclusivity of several genetic variants is interesting but needs validation.

The followup period is fairly short for the clinical correlates and relatively few patients recurred. Inclusion of a table of the clinical characteristics showing a breakdown of the T-class, N-class stage and

outcome would be helpful to evaluate the clinical outcome and the association with the minimal gene signature. However, the non-recurrent tumors have minimal genetic changes, whereas the recurrent and metastatic tumors are far more complex. So the low genetic complexity alone may be a marker of good outcome rather than the minimal gene signature described for poor outcome. Addition of the primary tumor size, nodal status and stage would be informative on figure 5, which shows the minimal signature set for tumor recurrence.

**I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

**Competing Interests:** No competing interests were disclosed.

Referee Report 09 November 2015

doi:[10.5256/f1000research.7870.r11114](https://doi.org/10.5256/f1000research.7870.r11114)



**Nishant Agrawal**

Ludwig Center for Cancer Genetics and Therapeutics, The Johns Hopkins University, Baltimore, MD, USA

The authors report an integrated genetic, epigenetic, and expression analysis from 50 oral tongue SCC. The findings were confirmed using a data set from TCGA. The manuscript presents surprising data which is interesting and has potential clinical implications.

1. It is surprising that 23 of 50 patients with oral tongue SCC are HPV/p16 positive. This is not entirely consistent with previously published literature. Although the authors comment on this, this finding is rather “unique.”
2. It is very possible that the improved survival is due to HPV status and not mutation status as these results maybe confounded.
3. Although 50 patients is a relatively large sample size for such a study, a much larger study is necessary to really have immediate clinical impact.

**I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.**

**Competing Interests:** No competing interests were disclosed.

---