

1 **Description of Supplementary Data Files**

2

3

4 File Name: **Supplementary Data 1**

5 Legend: Clinical data for all sequenced tumours and matched controls.

6

7

8

9 File Name: **Supplementary Data 2**

10 Legend: Matrix with WES results including all mutated genes in mouse TSCCs. Table
11 containing the selected mutations (Missense, Nonsense, Silent, Ess_splice, Indel,
12 Deletions, Start Lost and Stop Lost) in each gene (row) in each sample (columns).
13 Comma separated mutations indicate multiple events present in the same gene in the
14 same sample.

15

16

17 File Name: **Supplementary Data 3**

18 Legend: Tongue recurrent mutated genes in human TSCC. List of 465 human genes
19 evaluated from target-exome sequencing.

20

21

22 File Name: **Supplementary Data 4**

23 Legend: *DNdScv* analysis output: annotation of mutations per sample. The first table
24 contains the list of driver genes identified with the algorithm. For each gene, the
25 following informations are indicated: number of mutations (synonymous, missense,
26 nonsense, essential splice sites), maximum-likelihood estimates (MLEs) of the dN/dS
27 ratios for each gene for each type of mutations, *p*-values and *q*-values. The second
28 table contains the annotation for each mutation in each sample. For each mutation,
29 the following info is indicated: sample and gene where the mutation occurs,
30 chromosome and position, aminoacid, nucleotide and codon change, impact.

31

32 File Name: **Supplementary Data 5**

33 Legend: Gene Ontology analysis of the significant genes from dndscv analysis. List
34 of 25 GO terms (biological process) for which a statistically significant enrichment
35 was detected in the driver genes. For each category, p -values, FDR and enrichment
36 score are indicated. The enrichment score is calculated as $(b/n)/(B/N)$, where: **b** is
37 the number of genes in the intersection; **n** is the number of genes in the top of the
38 user's input list or in the target set when appropriate; **B** is the total number of genes
39 associated with a specific GO term; **N** is the total number of genes. The statistical
40 test used for enrichment analysis is typically based on a hypergeometric model.

41

42

43 File Name: **Supplementary Data 6**

44 Legend: Linear regression analysis to establish the relationship between the total
45 number of mutations and the VAF value in each gene and the clinical parameters
46 (rows). The estimates for the model's coefficients, the standard deviation of the
47 sampling distribution of the estimate of the coefficient under the standard regression
48 assumptions, the *tvalue* and two-sided p -value ($Pr(>|t|)$) are shown for each
49 correlation (gene versus clinical parameter).