

Cellular states are coupled to genomic and viral heterogeneity in HPV-related oropharyngeal carcinoma

Received: 16 September 2021

Accepted: 27 February 2023

Published online: 3 April 2023

 Check for updates

Sidharth V. Puram^{1,2,3,15} , Michael Mints^{4,5,6,15}, Ananya Pal¹, Zongtai Qi¹, Ashley Reeb¹, Kyla Gelev⁷, Thomas F. Barrett¹, Sophie Gerndt¹, Ping Liu^{7,8}, Anuraag S. Parikh⁹, Salma Ramadan¹, Travis Law¹, Edmund A. Mroz¹⁰, James W. Rocco¹⁰, Doug Adkins^{3,11}, Wade L. Thorstad^{3,8}, Hiram A. Gay^{3,8}, Li Ding^{11,12}, Randal C. Paniello^{1,3}, Patrik Pipkorn^{1,3}, Ryan S. Jackson^{1,3}, Xiaowei Wang^{7,13}, Angela Mazul¹, Rebecca Chernock¹⁴, Jose P. Zevallos^{1,3}, Jessica Silva-Fisher^{3,7} & Itay Tirosh⁴ 

Head and neck squamous cell carcinoma (HNSCC) includes a subset of cancers driven by human papillomavirus (HPV). Here we use single-cell RNA-seq to profile both HPV-positive and HPV-negative oropharyngeal tumors, uncovering a high level of cellular diversity within and between tumors. First, we detect diverse chromosomal aberrations within individual tumors, suggesting genomic instability and enabling the identification of malignant cells even at pathologically negative margins. Second, we uncover diversity with respect to HNSCC subtypes and other cellular states such as the cell cycle, senescence and epithelial-mesenchymal transitions. Third, we find heterogeneity in viral gene expression within HPV-positive tumors. HPV expression is lost or repressed in a subset of cells, which are associated with a decrease in HPV-associated cell cycle phenotypes, decreased response to treatment, increased invasion and poor prognosis. These findings suggest that HPV expression diversity must be considered during diagnosis and treatment of HPV-positive tumors, with important prognostic ramifications.

Head and neck squamous cell carcinoma (HNSCC) tumors of the oral cavity and larynx are typically linked to alcohol and tobacco exposure, while oropharyngeal squamous cell carcinoma (OPSCC) is more commonly associated with infection by human papillomavirus (HPV)¹.

HPV-associated OPSCCs have better prognosis than other forms of HNSCC, calling for treatment de-escalation to reduce side effects while maintaining an excellent prognosis. However, a subset of HPV-related OPSCCs responds poorly to treatment and recur², underscoring the

¹Department of Otolaryngology—Head and Neck Surgery, Washington University School of Medicine, St. Louis, MO, USA. ²Department of Genetics, Washington University School of Medicine, St. Louis, MO, USA. ³Siteman Cancer Center, Washington University in St. Louis, St. Louis, MO, USA. ⁴Department of Molecular Cell Biology, Weizmann Institute of Science, Rehovot, Israel. ⁵Department of Surgical and Perioperative Sciences, Urology and Andrology, Umeå University, Umeå, Sweden. ⁶Department of Oncology—Pathology, Karolinska Institute, Stockholm, Sweden. ⁷Department of Internal Medicine, Washington University School of Medicine, St. Louis, MO, USA. ⁸Department of Radiation Oncology, Washington University School of Medicine, St. Louis, MO, USA. ⁹Department of Otolaryngology—Head and Neck Surgery, Columbia University Irving Medical Center, New York, NY, USA. ¹⁰Department of Otolaryngology—Head and Neck Surgery, Ohio State University, Columbus, OH, USA. ¹¹Department of Medicine, Washington University School of Medicine, St. Louis, MO, USA. ¹²McDonnell Genome Institute, Washington University in St. Louis, St. Louis, MO, USA. ¹³Department of Pharmacology and Regenerative Medicine, University of Illinois at Chicago, Chicago, IL, USA. ¹⁴Department of Pathology and Immunology, Washington University School of Medicine, St. Louis, MO, USA. ¹⁵These authors contributed equally: Sidharth V. Puram, Michael Mints. ✉ e-mail: sidpuram@wustl.edu; itay.tirosh@weizmann.ac.il

need for a deeper understanding of these tumors and the development of new therapeutic approaches.

HPV is a sexually transmitted DNA virus, accounting for more than 5% of all cancer cases worldwide including all cervical cancers, most OPSCCs and a majority of vaginal and anal cancers³. The HPV oncogenic proteins, E6 and E7, inhibit the tumor suppressor proteins p53 and Rb, respectively, thereby activating the proliferation of infected epithelial cells^{4,5}. While these initial effects of HPV are largely understood, the subsequent events leading to tumorigenesis⁶, the biology of the resulting tumors and their vulnerabilities still remain poorly characterized.

Previously, we used single-cell RNA sequencing (scRNA-seq) to interrogate patient samples of HPV-negative oral cavity tumors^{7–10}. Here we turn our focus to OPSCC, profiling both HPV-positive and HPV-negative tumors. We uncover unanticipated diversity of chromosomal aberrations and of HPV expression patterns. Strikingly, each HPV-positive tumor harbors a subset of malignant cells in which HPV expression is not detected, and HPV-related phenotypes are decreased. These cells may influence prognosis and therapy response, highlighting their significance and opportunities for new interventions.

Results

scRNA-seq analysis of OPSCC

We profiled 16 treatment-naïve OPSCC primary tumor samples using the 10x chromium platform (Fig. 1a and Supplementary Tables 1 and 2). After removing low-quality cells and potential doublets (see Methods), we retained 70,970 cells, which were used to describe four distinct layers of cellular diversity (Fig. 1a): cell types, genetic clones, cellular states and HPV expression patterns.

We first clustered all cells and annotated the clusters by differentially expressed marker genes (Fig. 1b,c, Extended Data Fig. 1a,b and Supplementary Tables 3 and 4). Epithelial cells clustered primarily by patient identity, while nonepithelial cell types clustered together regardless of patient identity. We defined 12 nonepithelial clusters, including typical tumor microenvironment components (for example, T cells and fibroblasts) as well as less common components (for example, myofibroblasts and lymphovascular cells), each of which contained cells from multiple patients and expressed characteristic marker genes (Fig. 1d).

Based on the standard p16 staining, 12 of the OPSCC tumors were clinically defined as HPV-positive and 4 as HPV-negative (Extended Data Fig. 1c). We mapped the scRNA-seq reads of all epithelial cells to the five most common high-risk HPV genotypes¹¹ and identified transcripts of the most common HPV genotype (HPV16) in 11 of the 12 tumors clinically defined as HPV-positive and in none of the tumors clinically defined as HPV-negative (Extended Data Fig. 1d,e). In these 11 HPV-positive tumors, HPV16 transcripts were identified in an average of 53% (20–78% range) of epithelial cells (Supplementary Table 4). In one exceptional tumor (OP8), we did not identify any HPV transcripts despite clinical diagnosis as HPV-positive. Further testing and sequence analysis failed to identify evidence for any other HPV genotypes (Methods) although we cannot formally exclude the possibility of a rare undetected genotype. These results suggest either a false-positive clinical diagnosis by p16, clearance of the virus or a limitation in detecting HPV transcripts, which seems unlikely based on the other tumors.

The HPV16 genome contains eight genes, and these were detected at variable frequencies, with E5 being the most commonly detected, followed by E1. The pattern of HPV gene expression varied between tumors, with seven tumors expressing E5 at particularly high levels and others with a relative enrichment of E7 (Extended Data Fig. 1e). These distinct patterns of HPV expression, rather than a uniform expression of viral proteins, are consistent with previous findings¹².

Chromosomal aberrations identify malignant cells and clonality

We classified the epithelial cells into malignant and nonmalignant cells based on the inference of chromosomal copy-number aberrations

(CNAs)^{8,13–15}. At each chromosomal locus, an estimated copy number was calculated by averaging the normalized expression levels of the hundred adjacent genes compared to their expression in a reference set of fibroblasts and endothelial cells (Methods). Most epithelial cells had multiple CNAs, including characteristic CNAs of OPSCC (for example, 3p loss, 3q and 8q gain; Fig. 2a and Extended Data Fig. 2a). We classified epithelial cells into malignant and nonmalignant cells, by the combined evidence for CNAs across all chromosomes (that is, CNA signal) and the similarity of the CNA pattern to that of other cells from the same tumor (that is, CNA correlation), thereby defining a robust separation (Fig. 2b).

Overall, we classified 20,323 (85%) epithelial cells as malignant cells, while 2,625 (11%) cells were classified as nonmalignant (normal) epithelial cells (Extended Data Fig. 2a). The remainder were defined as unresolved and excluded from further analysis, likely reflecting doublets or low-quality cells. To validate the CNA-based classification, we compared epithelial cells from HPV-positive tumor samples to those from normal adjacent tissue of the same patients to derive a gene expression signature of HPV-related malignancy (Supplementary Table 3). Scoring epithelial cells from HPV-positive patients by this signature showed remarkable congruence with the CNA-based classification (Extended Data Fig. 2b,c). In HPV-positive tumors, CNA classification was largely consistent with the identification of HPV transcripts although we also detected a small subset of nonmalignant cells with HPV transcripts (Extended Data Fig. 3f,g and Supplementary Note).

The CNA analyses also uncovered distinct genetic subclones within individual tumors. For example, in OP17, the malignant cells were separated into two genetic subclones with both shared and clone-specific gains and losses (Fig. 2a,b). Overall, multiple CNA subclones were identified within 14 of 16 tumors (Fig. 2c). OP9 displayed particularly extensive subclonal diversity, with six different genetic subclones (Fig. 2d), for which we inferred a phylogenetic tree (Fig. 2e). The tree defined two major branches (subclones A–B and subclones C–F) that also differed in expression, with one branch expressing a unique program with many mesenchymal genes (for example, collagens; Fig. 2d, right-most column).

The two branches of OP9 also differed by the frequencies of HPV detection (62% versus 33%, $P < 2.2 \times 10^{-16}$, chi-square test; Fig. 2d). Similarly, a significant difference in HPV detection between subclones was found for eight of the ten HPV-positive tumors that had multiple subclones (Fig. 2c). For example, in OP13, three subclones had HPV detected in 90%, 4% and 0% of cells. In some cases, subclones differed not only in the frequency of HPV detection but also in the relative detection of distinct HPV genes (Extended Data Fig. 2e). Thus, the overall abundance and the relative expression patterns of HPV genes appear to be modified during tumor evolution and to vary both between and within HPV-positive tumors.

Malignant cells found in the histologically negative tumor margin

In three cases, we were also able to profile histologically negative margin tissues (adjacent normal). Most epithelial cells from these samples were classified by CNA analysis as nonmalignant, as expected. However, in one negative margin sample (OP34; Extended Data Fig. 3a), 29 of 80 epithelial cells were classified as malignant by CNAs and by the HPV-related malignancy expression signature (Fig. 2f and Extended Data Fig. 3b). A subset of these cells expressed HPV genes, further supporting their malignant classification (Fig. 2f). These malignant cells harbored all of the CNAs shared across OP34 subclones as well as unique CNAs (in chromosomes 4q, 9 and 22), thus representing a separate genetic subclone (Extended Data Fig. 3c). These results suggest further evolution of an invasive subclone beyond the leading front (histological edge) of the tumor. Notably, OP34 had clear resection margins on frozen and permanent histopathologic analysis (see Methods), indicating that malignant cells were not expected in the margin sample by traditional pathologic techniques. However, in a subset of OPSCCs,

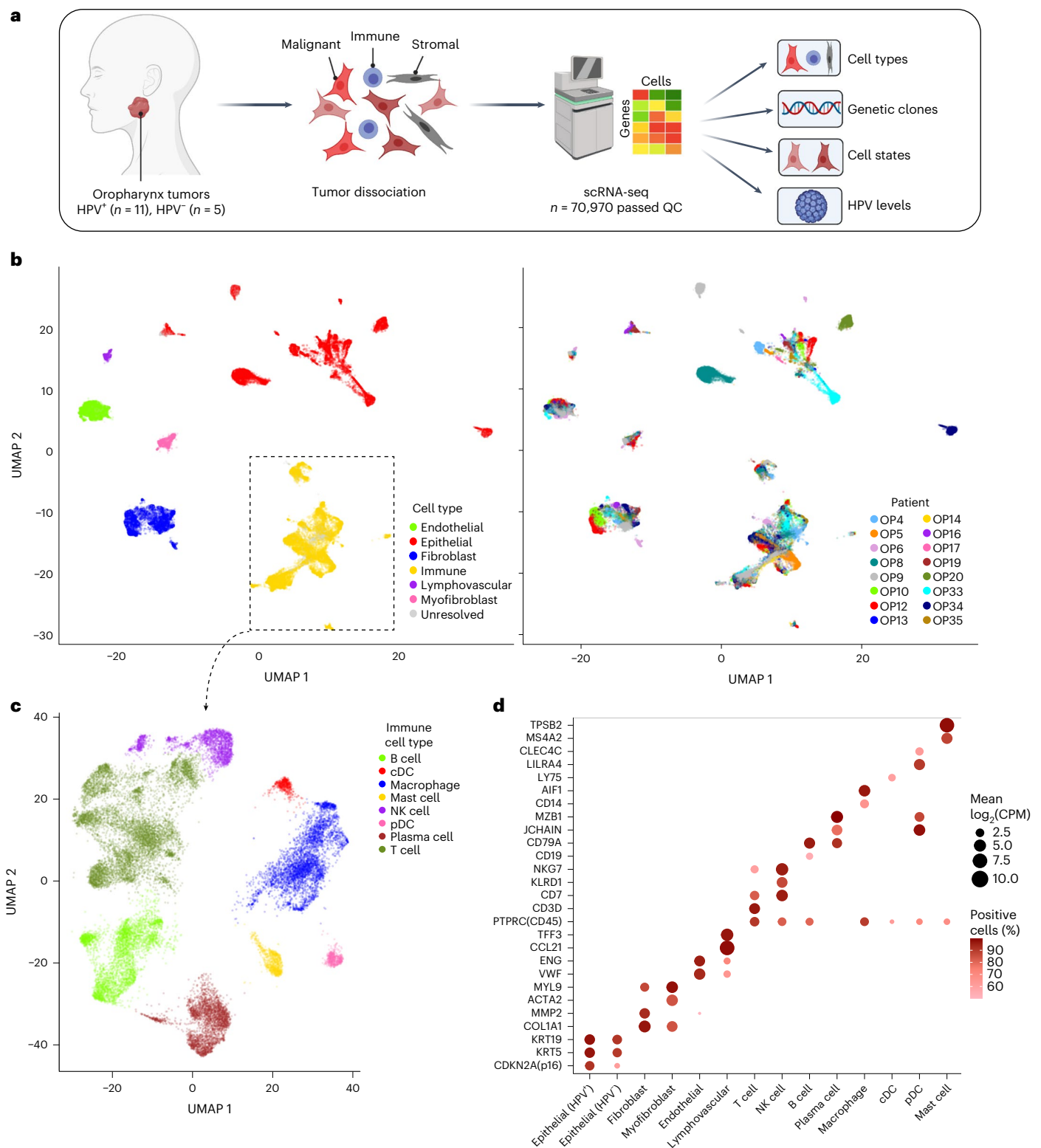


Fig. 1 | ScRNA-seq analysis of 16 OPSCC tumors. a, Scheme of the workflow for OPSCC profiling and subsequent analysis. **b**, UMAP plot of all cells that passed QC (n = 70,970), colored by cell type and patient. **c**, UMAP plot of all immune cells (n = 22,818), colored by immune cell type. **d**, Dot plot showing expression of

selected marker genes (y axis) by all cells assigned to each cell type (x axis). Dot size represents average expression, and dot color represents the fraction of cells with nonzero expression.

tumor recurrence occurs despite surgery with widely clear margins, suggesting that individual malignant cells likely remain undetected in these cases, as might be the case in OP34 (ref. 16). We compared the expression of malignant cells from the margin to both malignant cells

from the core of the tumor and to nonmalignant epithelial cells in the margin sample (Fig. 2g and Supplementary Table 2). Fifty-seven genes were significantly ($P < 0.05$, t-test) upregulated in the malignant cells from the margin, including cytokeratin, EMT-related genes, APOBEC

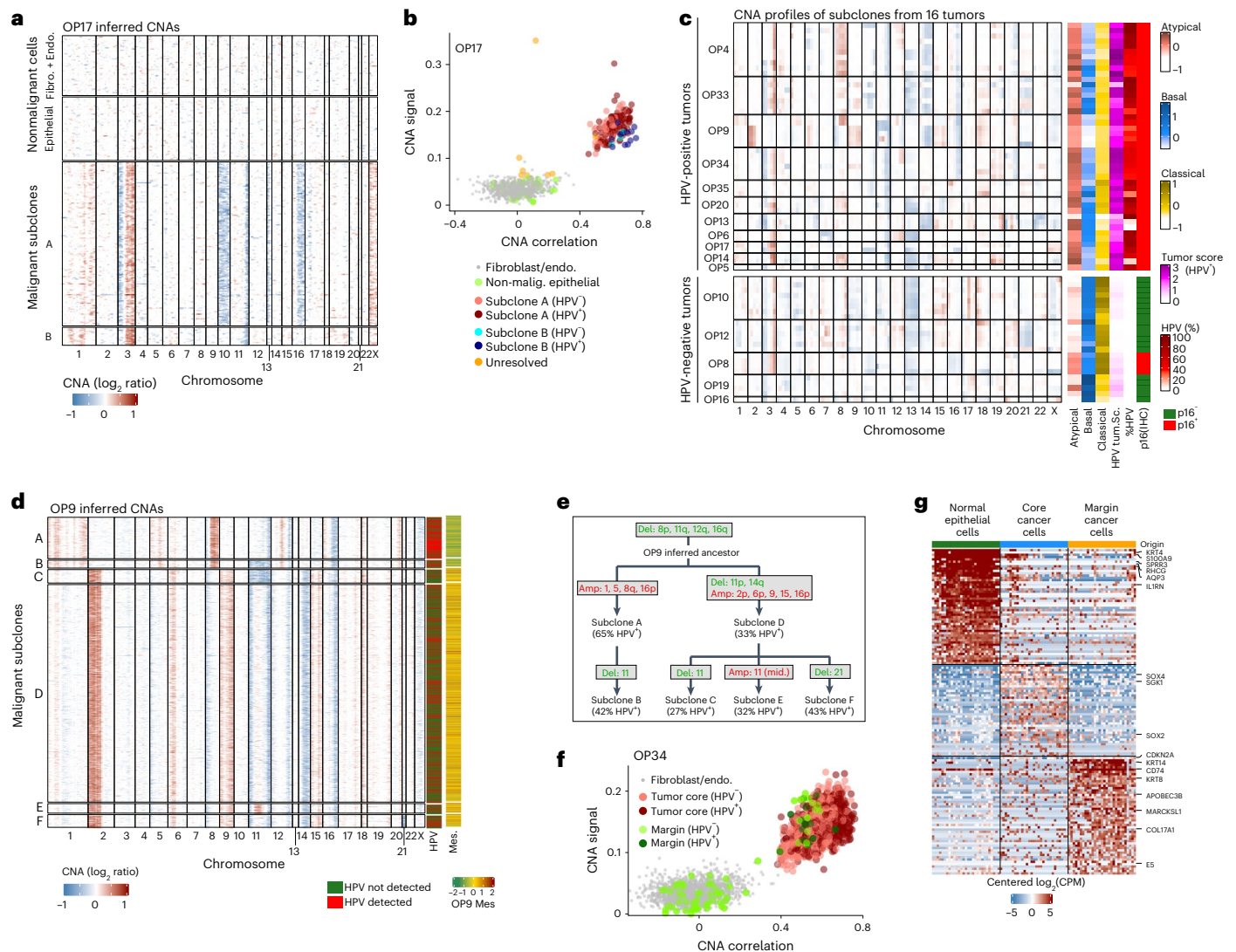


Fig. 2 | Inference of chromosomal aberrations for identification of malignant cells, genetic subclones and invasive cells. **a**, CNA plot of OP17, inferred by taking a 100-gene moving average of relative expression values across the transcriptome (Methods). Rows represent cells, arranged by genetic subclones and columns genes, arranged by chromosomal position. Fibroblasts and endothelial cells, used as a reference for CNA inference, as well as cells classified as nonmalignant epithelial cells, are shown above the malignant cells. **b**, Scatter plot of two CNA metrics used for classification of cells as malignant. CNA signal (y-axis) and CNA correlation (x-axis). All epithelial and stromal cells of OP17 are shown, colored by their cell type, subclone assignment and HPV expression. **c**, Left: average CNA profiles for all identified genetic subclones; rows represent subclones, ordered by patient and columns represent chromosomal positions (with five bins per chromosome). Right: scores of subclones (arranged as in left panel) for the TCGA subtypes and the HPV⁺ tumor signatures, the percentage of cells with HPV reads and the HPV clinical classification of the corresponding

tumor based on p16 staining. Subclone scores reflect the average scores of the cells in each subclone. **d**, CNA plot of malignant cells in OP9 as in **a**. Columns on the right show detection of HPV reads and the average expression of a mesenchymal signature found in OP9. **e**, Inferred phylogenetic tree of genetic subclones in OP9. The percentage of cells with detection of HPV reads is noted for each observed subclone; chromosomal deletions (green) and amplifications (red) are noted for each observed subclone as well as for the inferred ancestral cell. **f**, CNA signal and correlation scatter plot for OP34 as in **b**. Cells are colored by their origin (tumor core or margin sample) and by HPV expression. **g**, Heatmap of differentially expressed genes between the three subsets of epithelial cells in OP34—normal epithelial cells, invasive malignant cells and malignant cells from the tumor core. Rows represent genes, and columns represent cells. An equal number of cells is shown from each subset (to that end, cells from the normal and tumor core subsets were randomly sampled).

genes, immune-related genes and the HPV *E5* gene, clearly distinguishing this invasive population.

To explore the generalizability of identifying malignant cells within histologically negative margins, we searched for other scRNA-seq datasets containing matched tumor and negative margin samples. We identified two lung adenocarcinoma samples¹⁷ meeting these criteria and classified the cells from these samples by inferred CNAs (Extended Data Fig. 3d,e). One lung tumor did not contain malignant cells in the margin, while another tumor had malignant cells in the negative margin sample, representing a distinct subclone by CNAs and upregulating

47 genes (Extended Data Fig. 3d and Supplementary Table 3). These data highlight the presence of malignant cells in histologically negative margin biopsies that may drive adverse clinical outcomes.

OPSCC tumor diversity highlights three cellular subtypes

Overall, the diversity among malignant cells is linked to patient identity, to HPV status and to three TCGA subtypes—*atypical*, *basal* and *classical*¹⁸ (Fig. 3a, Extended Data Fig. 4 and Supplementary Table 5). Each tumor had a dominant subtype that, on average per tumor, covered 96% of the malignant cells that were confidently assigned to

any subtype (Fig. 3a,b). All HPV-positive tumors had a dominant *atypical* subtype. In contrast, HPV-negative tumors included two with a dominant *basal* subtype, two with a dominant *classical* subtype and OP8, with a majority of intermediate cells (not confidently assigned to any TCGA subtype), possibly reflecting its unique pattern as a p16-positive but HPV-negative tumor.

Despite the dominant subtype of each tumor, subsets of cells from five tumors were confidently assigned to a secondary subtype. In all of these cases, subtype heterogeneity was linked to genetic subclones (Fig. 3b). For example, while OP19 is dominated by the *basal* subtype, 4% of its malignant cells are classified as *atypical* and these primarily derive from subclone C. Interestingly, while OP19 is HPV-negative, it has high mRNA expression of *CDKN2A* (p16) (Extended Data Fig. 4c), perhaps relating to its secondary *atypical* subtype.

Intratumoral heterogeneity and epithelial senescence

To systematically search for additional patterns of intratumoral heterogeneity, malignant cells from each tumor were analyzed by nonnegative matrix factorization (NMF). The expression programs identified as variable within tumors were compared across tumors to define eight groups of recurrent expression programs (Fig. 3c). For each of the eight groups, we defined a consensus ‘meta-program’, annotated them by functional enrichments and scored all malignant cells for these meta-programs (Fig. 3d and Supplementary Table 6). A similar analysis was also performed for six common nonmalignant cell types (Extended Data Fig. 5 and Supplementary Table 7).

The malignant meta-programs included cell cycle (G1/S and G2/M phases), stress and hypoxia responses, oxidative phosphorylation, interferon response, partial EMT (p-EMT) and an epithelial senescence-associated (EpiSen) program. Notably, the latter two meta-programs appear to be enriched in HNSCC and associated with metastasis and drug responses, respectively^{19,20}. EpiSen was the most common pattern of heterogeneity in OPSCC, detected in 14 tumors (Fig. 3c), with high similarity to previously identified programs in oral cavity tumors⁸, cell lines¹⁹ and other squamous cancers^{21,22}. The fraction of EpiSen_{high} cells varied markedly between tumors, from less than 5% to more than 50% of malignant cells. Consequently, pseudobulk tumor profiles segregated the HPV-positive tumors primarily by the frequency of EpiSen_{high} cells (Extended Data Fig. 4b).

Undetectable HPV expression in a subset of malignant cells from HPV-driven tumors

Except for the G2/M meta-program, all other meta-programs differed significantly ($P < 0.05$, hypergeometric test) in their abundance between HPV-positive and HPV-negative tumors (Fig. 3e). Interestingly, we also noticed differences in meta-program abundances within the HPV-positive tumors when distinguishing between 66% of the malignant cells in which we detected HPV reads (denoted as *HPV_{on}* cells) and the remaining 34% in which we did not detect any HPV reads (denoted as *HPV_{off}* cells) (Fig. 3e). For example, expression of the G1/S meta-program was enriched in *HPV_{on}* cells, not only relative to

cells from HPV-negative tumors (denoted in Fig. 3e by asterisks within *HPV_{neg}*) but also relative to the *HPV_{off}* cells from HPV-positive tumors (denoted in Fig. 3e by asterisks within *HPV_{on}*). Conversely, the EpiSen meta-program was specifically enriched in *HPV_{off}* relative to *HPV_{on}* cells among the HPV-positive tumors.

These results raise the possibility that lack of HPV detection in *HPV_{off}* cells may not merely reflect limited scRNA-seq sensitivity, but may also correspond to unique cellular states associated with genetic or epigenetic repression of HPV genes. Such repression is consistent with the variability in HPV expression profiles (Extended Data Figs. 1e and 2e) and in the fraction of cells with detected HPV that we observed between tumors and subclones (Fig. 2c,d). Thus, HPV status may define not only two types of tumors (HPV-positive and HPV-negative), but at least three types of malignant cells (*HPV_{neg}*, *HPV_{on}* and *HPV_{off}*). Notably, no tumors were entirely *HPV_{off}*, rather just subsets of malignant cells within each of the HPV-positive tumors.

To examine if HPV genes are repressed in *HPV_{off}* cells or whether their expression is not detected due to technical limitations, we compared the frequency of *HPV_{off}* cells to the frequency in which other sets of genes are not detected. The fraction of *HPV_{off}* was significantly higher than expected based on the detection of other sets of control genes sampled so that each gene in the control gene set had similar average expression levels to one HPV gene (34% versus 9%, $P < 2.2 \times 10^{-16}$, z-test; Fig. 4a). RNAscope in situ hybridization (ISH; Fig. 4b and Extended Data Fig. 6a) as well as immunohistochemistry (IHC; Extended Data Fig. 6b) further supported the presence of *HPV_{off}* cells by demonstrating the absence of *E6* and *E7* RNA and *E6* protein in several tumor areas marked by p16. Together, these data uncover a subset of malignant cells (*HPV_{off}*) in which HPV expression is lost or reduced.

HPV_{off} cells are associated with HPV-negative phenotypes

We next identified genes that were differentially expressed between *HPV_{on}* and *HPV_{off}* cells across multiple tumors (Supplementary Table 8). EpiSen genes were enriched in *HPV_{off}* cells, while G1/S cell cycle genes were enriched in *HPV_{on}* cells (Fig. 4c and Extended Data Fig. 6c,d). To further characterize cell cycle differences between *HPV_{on}* and *HPV_{off}* cells, we divided all malignant cells into ten bins by expression of the G1/S program. In all HPV-positive tumors, *HPV_{on}* and *HPV_{off}* cells were significantly enriched ($P < 0.01$ for every patient, chi-square test) in higher and lower G1/S bins, respectively (Fig. 4d and Extended Data Fig. 6e). This consistent association of *HPV_{on}* with cell cycle is also seen across subclones (Extended Data Fig. 6f).

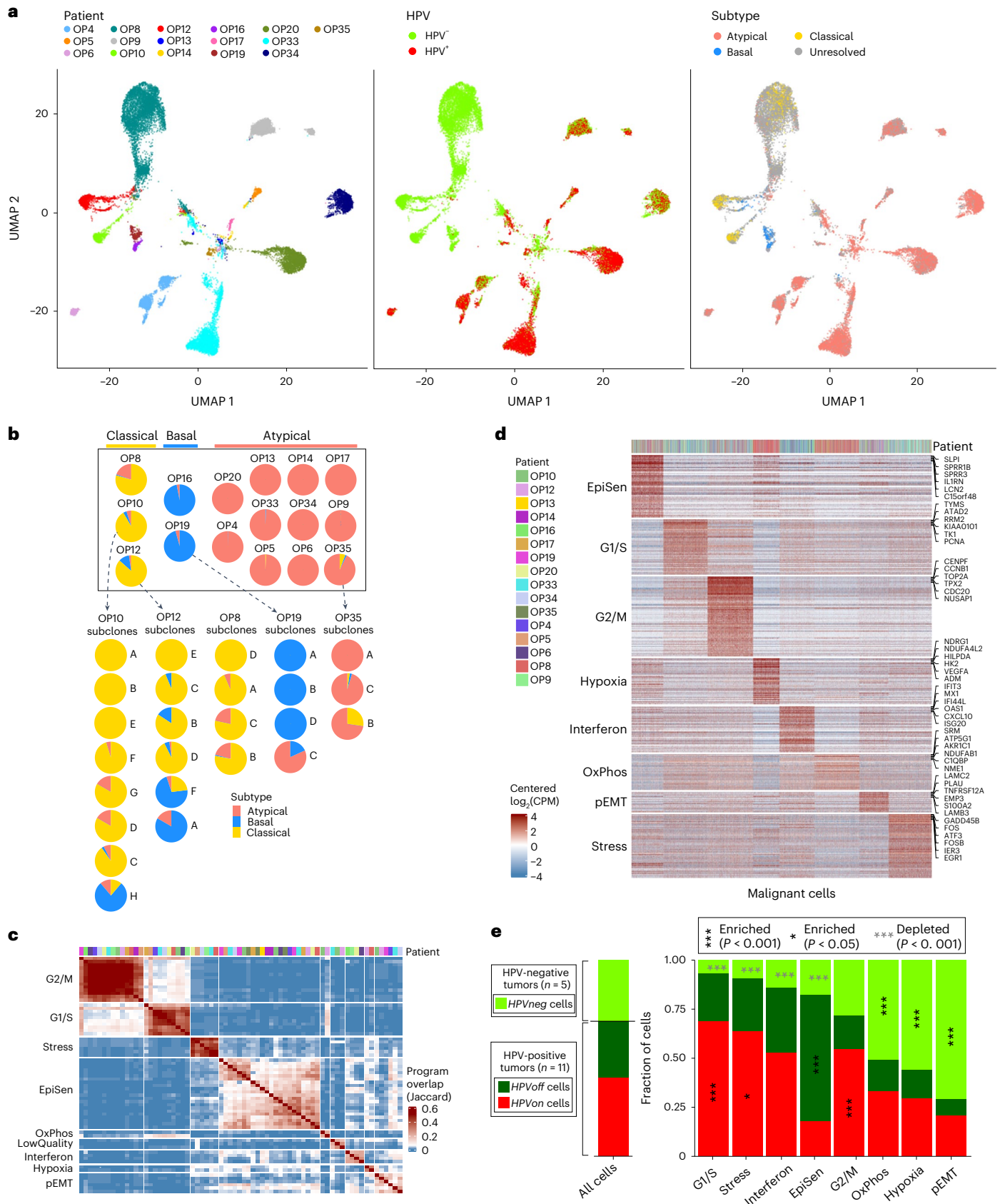
While all cancers are associated with increased proliferation, *HPV_{on}* cells have higher expression of the G1/S meta-program than other cancer types, based on reanalysis of multiple 10x scRNA-seq datasets. Most *HPV_{on}* cells (54%) were among the G1/S-high cells, compared to substantially lower fractions (defined in a similar manner) for the HPV-negative OPSCCs (from this work) as well as for nine other tumor cohorts (Fig. 4e). For *HPV_{off}* cell, 36% were G1/S-high, which is comparable to the other tumor cohorts although still higher than most of them. Thus, HPV is associated with aberrant activation of the G1/S

Fig. 3 | Diversity of OPSCC malignant cells. **a**, UMAPs of all malignant cells ($n = 20,323$) colored by the patient (left panel), HPV expression (middle panel) and TCGA subtype (right panel). Cells with smaller than 1.5-fold change between the top and the second highest subtype scores were defined as unresolved and marked in gray. **b**, Pie charts representing the fraction of cells assigned to each TCGA subtype (excluding unresolved cells), per patient (above) and per subclone for patients with multiple subtypes and multiple subclones (below). **c**, Hierarchical clustering of 69 NMF-derived program signatures from 16 patients (Methods). Signatures are clustered by Jaccard overlap. Groups of signatures, from which meta-programs are derived, are annotated on the left. Top panel shows the patient origin for each program using the same color map as in **d**. **d**, Expression of meta-program genes (rows) in all malignant cells (columns). Top panel indicates the patient origin for every cell. **e**, For each meta-program,

bar plot shows the fraction of cells, out of those assigned to that meta-program, in three HPV-related classes: cells from HPV-negative tumors (*HPV_{neg}*, light green) and cells from HPV-positive tumors in which HPV reads are detected (*HPV_{on}*, red) or undetected (*HPV_{off}*, dark green). Asterisks denote enrichment (black and vertical) or depletion (gray and horizontal); asterisks within the *HPV_{neg}* area denote enrichment/depletion in *HPV_{neg}* versus HPV-positive tumors (*HPV_{on}* and *HPV_{off}*), and asterisks within the *HPV_{on}* or *HPV_{off}* area denote enrichment in comparison between those two classes. The significance of enrichment/depletion was calculated using a hypergeometric test, corrected for multiple testing. Bar plot at the left shows the same analysis for all malignant cells. When calculating all fractions, 100 cells per patient and subset were randomly sampled 100 times to avoid patients with more cells skewing the results.

expression program, which is reduced in *HPVoff* cells, consistent with the possibility that HPV expression is suppressed in those cells. Loss of pRb repression by the HPV-E7 oncogene^{4,5} may thus decrease the proportion of *HPVoff* cells in the G1/S phase.

We further speculated that the second major difference between *HPVon* and *HPVoff* cells—the enrichment of EpiSen in *HPVoff* cells (Fig. 4c)—reflects the inactivation of p53 by HPV-E6, as the absence of HPV-E6 may enable cells to induce senescence. Multiple observations



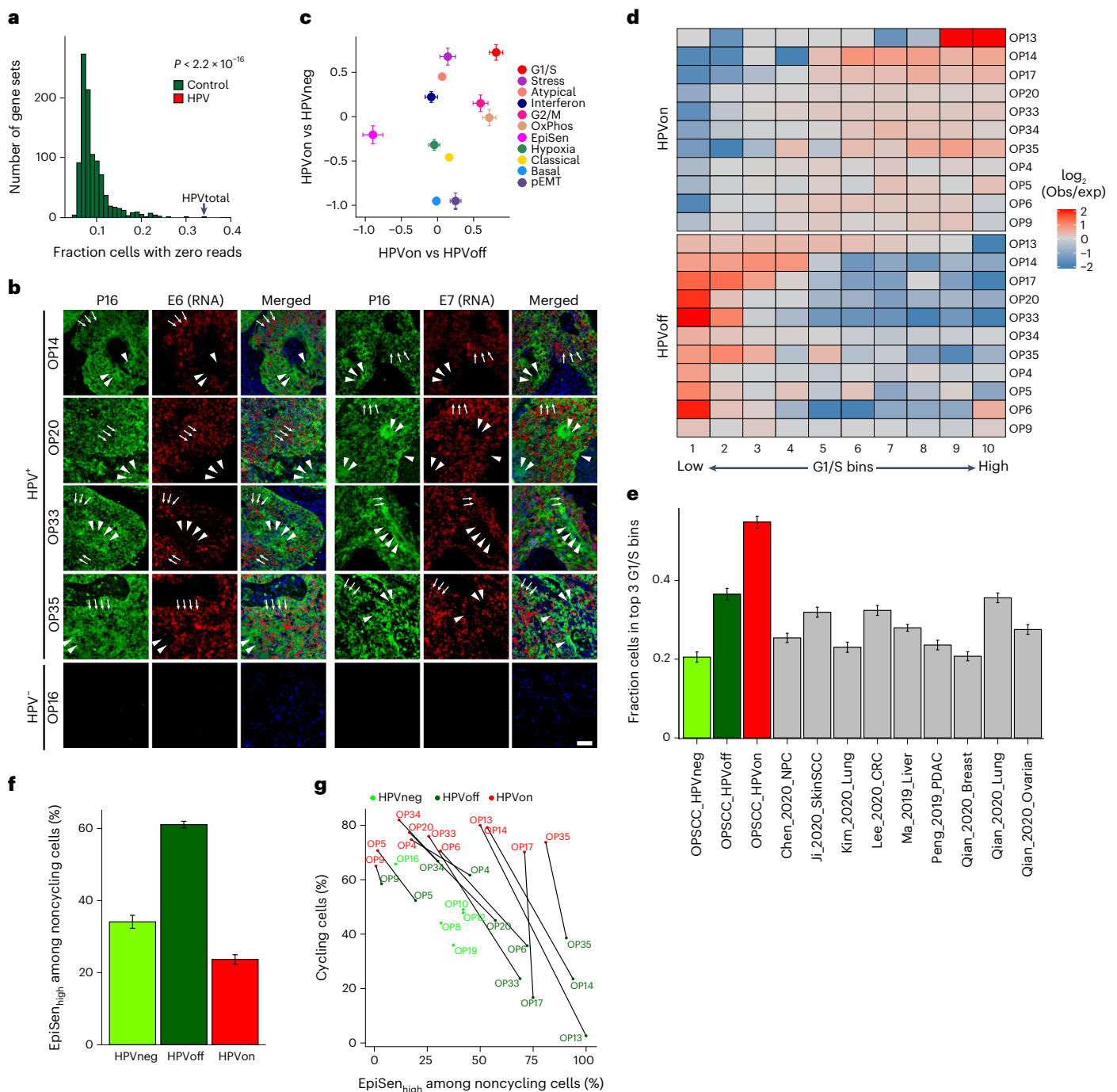


Fig. 4 | HPVoff cells and their association with cell cycle and senescence.

a, Fraction of cells with zero reads for the set of five detected HPV genes (*E1*, *E2*, *E5*, *E6* and *E7*) is significantly higher than that for 1,000 control gene sets ($P < 2.2 \times 10^{-16}$, z test). Each control gene set includes one non-HPV gene as the matched control for each of the five HPV genes. Control genes were randomly sampled among the 100 genes closest to the respective HPV gene, based on the average expression across all cancer cells from HPV-positive patients ($P < 2.2 \times 10^{-16}$, z test). **b**, RNA ISH (RNAScope) of representative HPV-positive (OP14, OP20, OP33 and OP35) and HPV-negative (OP16) tumors, for viral *E6* (left) and *E7* (right) RNA (red) with immunofluorescence costaining for regions of tumor as marked by p16 protein (green) and nuclei by DAPI (blue). HPV-positive tumors display regions of p16 positivity with the absence of *E6* and *E7* RNA signal, consistent with an HPVoff state (arrowheads), while other regions have p16 along with *E6* and *E7* expression (HPVon; arrows). HPV-negative tumors do not have signal for p16 protein or *E6* or *E7* RNA. Scale bar = 1,000 μm . **c**, Scatter plot of differences in meta-program expression between cells from different HPV classes. The x-axis shows mean difference, for all genes in each meta-program,

between HPVon and HPVoff cells within the same patient ($n = 11$ patients). The y-axis shows mean difference between HPVon cells, averaged across all HPV-positive patients ($n = 11$), and HPVneg cells, averaged across HPV-negative patients ($n = 5$). Error bars represent the standard error of the mean for each meta-program. **d**, \log_2 ratio of observed to the expected number of cells in each bin of G1/S scores ranked from low (left) to high (right), for each HPV-positive tumor (rows). Top and bottom rows correspond to the HPVon and HPVoff cells, respectively. **e**, Mean fraction of malignant cells with high G1/S expression, as defined by the top three bins of G1/S scores, in each HPV class from this work and in multiple external datasets ($n = 9$). Error bars show the standard error of the mean fraction from 100 repetitions with a sampling of 1,000 cells per dataset. **f**, Mean proportions of EpiSen_{high} cells among all noncycling cells, across HPV subsets in $n = 5$ (HPVneg) and $n = 11$ (HPVon, HPVoff) patients. The top 20% of all malignant cells by average expression of the EpiSen program genes were defined as EpiSen_{high}. Error bars represent standard error after resampling 100 times, each time sampling 200 cells per patient and subset. **g**, Proportions of cycling cells and of EpiSen_{high} cells among noncycling cells, for each patient and HPV subset.

link the EpiSen meta-program to senescence—it is induced in senescent keratinocytes and bronchial cells¹⁹ and is enriched in noncycling HNSCC cells, both in the oral cavity⁸ and in oropharynx tumors (Extended Data Fig. 6g). Notably, the enrichment of EpiSen in *HPVoff* cells remained significant when restricting the analysis to noncycling cells, to decouple the differences in induction of a senescence program from the differences in proliferation (Fig. 4f). In summary, *HPVon* and *HPVoff* cells from the same tumor differ in the fraction of cycling cells and in the induction of EpiSen among the noncycling cells, presumably reflecting the reduced activity of the two major HPV oncogenes (E6 and E7) in *HPVoff* cells. Notably, this observation is consistent across all 11 HPV-positive tumors (Fig. 4g).

TCGA data and cell lines support lower proliferation in *HPVoff* cells

To explore the functional significance of *HPVoff* cells, we examined the TCGA dataset of HPV-positive OPSCC tumors. We reasoned that bulk expression levels of HPV transcripts (normalized for tumor purity) could serve as an approximation for the fraction of *HPVon* versus *HPVoff* malignant cells. Consistent with our scRNA-seq analysis, normalized HPV expression correlates with G1/S scores across HPV-positive TCGA tumors (Fig. 5a). Similar results were obtained in the analysis of TCGA specimens for cervical cancer, suggesting that HPV expression levels are associated with G1/S induction across distinct contexts (Extended Data Fig. 7a).

As a complementary approach, we analyzed scRNA-seq data from three HPV-positive cell lines¹⁹. Although HPV expression was identified in most cells, we found *HPVoff* subpopulations in each of the cell lines (Fig. 5b and Extended Data Fig. 7b). In two cell lines, *HPVoff* cells were also associated with decreased G1/S scores (Fig. 5c and Extended Data Fig. 7c). Immunocytochemistry confirmed that expression of HPV proteins E6 and E7 (but not of p16) correlates with proliferation (Fig. 5d,e and Extended Data Fig. 7d,e). Moreover, the knockdown of E6 and E7 in these lines did not affect p16 expression but reduced proliferation (Extended Data Fig. 7f,g and Supplementary Fig. 1).

Single-cell clones from these two cell lines showed a spectrum of HPV expression, which we denoted as *HPVon*, *HPVoff* and intermediate clones (Fig. 5f,g). *HPVon* and *HPVoff* clones largely maintained their relative HPV expression levels over multiple passages (Extended Data Fig. 7h,i), demonstrating the heritability of these states. *HPVon* clones were enriched with cycling cells and were more proliferative (Fig. 5h,i

and Extended Data Fig. 7j). Notably, serum starvation of *HPVon* clones suppressed their proliferation with little to no effect on HPV expression (Extended Data Fig. 7k,l), suggesting that HPV expression does not merely reflect that a cell is cycling but rather directly promotes the cell cycle through the function of E6 and E7 (refs. 4,5). Taken together, these results highlight an association between heterogeneity of HPV expression levels and G1/S cell cycle activity.

HPVoff cells are epigenetically regulated and may be associated with invasion and drug resistance

In contrast to the observed expression differences of HPV genes between *HPVon* and *HPVoff* clones, the genomic copy numbers of HPV genes were comparable between *HPVon* and *HPVoff* clones (Extended Data Fig. 8a). Moreover, DNAScope ISH experiments did not show substantial differences in E6 and E7 between and within different tumors (Extended Data Fig. 8b,c). These observations support the possibility of epigenetic regulation of HPV expression. We, therefore, treated HPV-positive cell lines with inhibitors of two epigenetic regulators, EZH2 and DNA methyltransferases (DNMT). Inhibition of EZH2 substantially reduced HPV expression in 93VU147T with limited effect on SCC47, while DNMT inhibition reduced HPV expression in SCC47 but not in 93VU147T (Fig. 5j and Extended Data Fig. 8d,e). These effects were largely specific to *HPVon* clones (Extended Data Fig. 8f). Thus, epigenetic regulators may direct HPV expression and heterogeneity.

Next, we turned to examine the impact of heterogeneity in HPV expression on cancer phenotypes. Aberrant cell cycle activity of *HPVon* cells might render them susceptible to standard cancer treatments, while the less proliferative *HPVoff* cells may have reduced susceptibility to such treatments. Indeed, treatment of SCC47 cells with cisplatin and of 93VU147T cells with radiation, reduced the expression of HPV genes (Fig. 5k, left) without affecting HPV genomic copy numbers (Fig. 5k, right), consistent with the possibility that *HPVon* cells were preferentially eliminated. As expected, treatment of individual *HPVon* and *HPVoff* clones revealed that *HPVon* clones are more susceptible to these cytotoxic agents (Extended Data Fig. 8g). The aberrant cell cycle of *HPVon* cells might also diminish their migration and invasive capacity, due to potential migration-proliferation tradeoffs^{23,24}. Indeed, we found increased invasiveness of *HPVoff* clones compared to *HPVon* clones in both cell lines (Extended Data Fig. 8h,i).

The association of *HPVoff* cells with increased invasion and resistance to treatments suggests that the fraction of *HPVoff* cells in

Fig. 5 | Regulation and function of *HPVoff* cells. **a**, Scatter plot of all HPV-positive OPSCC samples in the TCGA cohort, showing correlation (two-sided Pearson correlation test) between the relative expression of HPV genes and of the genes in the G1/S meta-program. Relative expression values reflect residuals, after normalizing each sample for malignant cell purity (using the epithelial signature from Supplementary Table 3). **b**, UMAP of 1,422 cells from three HPV-positive cell lines colored by HPV expression. Cells with at least one read from an HPV16 gene were considered HPV+. **c**, Differences in expression of the G1/S meta-program between HPV subsets in the HPV-positive cell line 93VU147T. Cells were divided into five bins of equal size, ranked by G1/S expression. The Y-axis shows the mean ratio of cells belonging to an HPV subset in a bin versus the expected number of cells, assuming random distribution across bins. Error bars are standard error after 100 resampling runs, where 100 cells per subset were randomly selected. *P* value is based on the chi-square test, comparing the distributions of cells per bin between the groups. **d**, Immunocytochemistry images of 93VU147T cells probed with Ki67 (red) and E6 (green, top) or E7 (green, bottom). Nuclei were stained and visualized with DAPI (blue). Scale bar = 100 μ m. **e**, Bar plot (mean \pm s.e.m.) shows percentage of Ki67 positive cells among E6 and E7 positive cells (HPV detected; red) and E6 and E7 negative cells (HPV not detected; green). 50 cells were counted across four fields (*P* < 0.00001, chi-square). **f,g**, Bar plot (mean \pm s.e.m.) shows relative expression of E6 and E7 among single clones (*n* = 10) derived from 93VU147T (**f**) and SCC47 (**g**) revealing diversity in HPV expression (*P* < 0.00001, ANOVA). **h,i**, Line graph

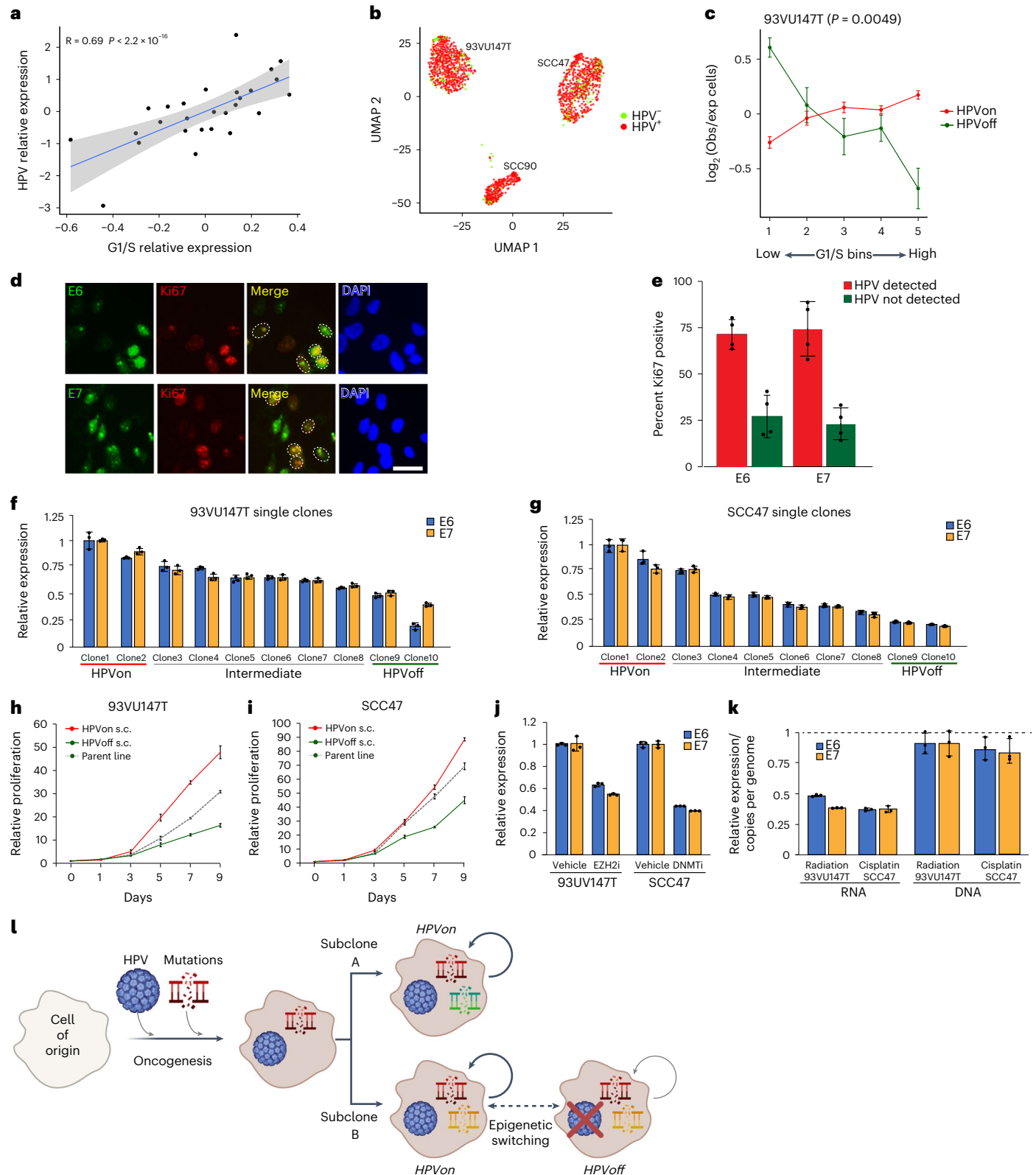
(mean \pm s.e.m.) shows relative proliferation of *HPVon* and *HPVoff* single clones derived from 93VU147T (**h**) and SCC47 (**i**) compared to parent line. *HPVon* single clones displayed substantially more relative proliferation than *HPVoff* single clones (*n* = 3; *P* < 0.00001, two-sided *t*-test). **j**, Left: bar plot (mean \pm s.e.m.) shows relative expression of E6 and E7 in 93VU147T cells treated with vehicle or tazemetostat (EZH2 inhibitor) (left). Right: bar plot (mean \pm s.e.m.) shows relative expression of E6 and E7 in SCC47 cells treated with vehicle or decitabine (DNMT inhibitor). Tazemetostat and decitabine substantially reduced relative E6 and E7 expression compared to vehicle in 93VU147T cells and SCC47 cells, respectively (*n* = 3; *P* < 0.001 and *P* < 0.00001, ANOVA). **k**, Left: Bar plot (mean \pm s.e.m.) shows relative expression of E6 and E7 in 93VU147T and SCC47 cells treated with radiation or cisplatin, respectively, normalized to control cells (dashed line) (*n* = 3; *P* < 0.00005, two-sided *t*-test). Right: Bar plot (mean \pm s.e.m.) depicts HPV copies per genome of E6 and E7 (normalized to albumin) for 93VU147T and SCC47 cells treated with radiation or cisplatin, respectively, normalized to control cells (dashed line). There were no significant differences in HPV copies in genomic DNA in radiation or cisplatin-treated cells compared to control (*n* = 3). **l**, Model of genomic and viral heterogeneity in HPV-related OPSCC. A combination of HPV infection and associated genetic mutations trigger oncogenesis. Some genetic subclones continue to express HPV (*HPVon*), while others may undergo epigenetic switching with repression of HPV expression (*HPVoff*) and an associated decrease in cell cycle (circled arrows).

a tumor might have clinical implications. Accordingly, HPV-positive tumors with low normalized HPV expression tend to have reduced recurrence-free survival compared to those with higher HPV expression (Extended Data Fig. 8j). This analysis was hindered by a small sample size, and the effect on survival had borderline statistical significance ($P = 0.05$), highlighting the need for further analysis with a larger patient cohort, while raising the intriguing possibility that loss

or reduction of HPV expression in subsets of cells may have a negative effect on patient survival.

Discussion

Our comprehensive scRNA-seq analysis reveals unappreciated diversity, both in genomic CNA profiles (Fig. 2) and in HPV gene expression (Fig. 4), within individual OPSCC tumors. The observed genomic



diversity may reflect an HPV-driven genomic instability, consistent with previous studies^{25–27}. This instability allowed us to robustly detect invasive malignant cells in pathologically normal tissue, which needs to be examined further in larger cohorts but may ultimately guide improved analyses of tumor margins.

The diversity of HPV expression is observed at three levels. First, different HPV genes are expressed at distinct levels in each tumor and cell line examined. Overall, *E5* is the most highly expressed HPV gene in tumors, but not in cell lines, highlighting the need to better understand its regulation and function. Second, these HPV expression patterns vary among tumors, among cell lines and even among genetic subclones of the same tumor. Thus, HPV integration and/or expression patterns may be modulated during tumor initiation and clonal evolution. Third, we do not detect HPV expression in a subset of cells, and the number of such cells is substantially higher than would be expected by the technical limitations of scRNA-seq. While we cannot distinguish between partial and complete repression of HPV genes, cells with undetected HPV mRNA (*HPVoff*) are associated with a decrease in the phenotypes that are driven by HPV oncogenes, namely aberrant cell cycle (through *E7*) and avoidance of senescence (through *E6*). These results suggest that reduced HPV levels persist for a sufficient degree and time to invoke phenotypes that partially resemble HPV-negative cells. *HPVoff* cells present a paradigm of heterogeneity in HPV expression within each HPV-positive tumor that is associated with unique cell states and clinical implications.

Given the strong evolutionary pressures to suppress viruses and the multitude of viral-protective mechanisms, it is tempting to speculate that even in a successful viral infection and a resulting tumor, the virus may still be suppressed in a subset of cells, thereby leading to the observed *HPVoff* cells. In cell lines, such suppression appears to be driven by epigenetic mechanisms: reduced HPV expression in *HPVoff* clones is not mirrored by reduced copy numbers at the DNA level and can be achieved by inhibition of epigenetic regulators. However, we also found significant variability in the fraction of *HPVoff* cells between tumor subclones within patients, suggesting that genetic evolution further modulates the transition towards *HPVoff* cells. We, therefore, speculate that multiple mechanisms, both genetic and epigenetic, regulate HPV expression levels and the emergence of *HPVoff* cells (see the model in Fig. 5I).

The potential clinical significance of *HPVoff* cells is hinted by their decreased response to treatments, increased invasion in vitro, and by the trend to worsen disease-free survival in HPV-positive patients with a larger proportion of *HPVoff* cells. We speculate that aberrant HPV-driven cell cycle activity facilitates responses to chemotherapy and radiation, partially accounting for the improved prognosis of HPV-positive tumors. *HPVoff* cells could resume their growth after treatment, possibly even switching their HPV expression back on, and may provide the basis for recurrent HPV-positive tumors.

While the cell cycle behavior of *HPVoff* cells is reminiscent of HPV-negative cells, the levels of *CDKN2A* (encoding p16) are indistinguishable between *HPVoff* and *HPVon* cells based on the scRNA-seq, RNAscope ISH and IHC (Fig. 4b, Extended Data Figs. 6d and 8k). Moreover, the knockdown of *E6* and *E7* decreased the proliferation of cells but not their levels of p16 (Extended Data Figs. 1 and 7f,g). *CDKN2A* was, however, substantially upregulated in the few nonmalignant epithelial cells in which we detected HPV reads (Extended Data Fig. 3g), highlighting the robust and early effect of HPV infection on *CDKN2A*. These observations raise the possibility that *CDKN2A* activation is more stable than other HPV-driven effects and may persist for a long time through unknown mechanisms.

If HPV expression varies after infection while *CDKN2A* (p16) expression remains constitutively high long after infection, then *CDKN2A* could theoretically become a more sensitive readout for latent or past HPV infection than HPV itself, potentially explaining the common and reliable use of p16 as a clinical marker of HPV infection. Interestingly,

CDKN2A is highly expressed in several OPSCC tumors and cell lines in which we do not detect any HPV reads (Extended Data Fig. 4c). It is conceivable that such tumors had initially been driven by HPV, but that during tumor progression one/few of the clones lost HPV or its expression (as appears to be the case in OP13; Fig. 2c), and these clones could have taken over the tumor via clonal evolution. Such a scenario could explain our observation that OP8 is p16-positive, yet HPV-negative, and has a mixture of transcriptional TCGA subtypes associated with HPV-negative and HPV-positive tumors (Fig. 3a,b). Similarly, the subset of *atypical* cells in OP19 (a tumor with undetected HPV transcripts but high *CDKN2A* expression) may be a remnant of latent or past HPV infection, although *CDKN2A* expression could also reflect nonviral mechanisms of regulation.

In summary, our single-cell atlas of OPSCC shows that genes encoded by an oncovirus (in this case, HPV) may cease to be expressed in a subset of cells, thereby reducing the oncogenic properties of cells, but also relieving their associated vulnerabilities, which may allow malignant cells to survive antitumor treatments and then potentially re-express the virally encoded oncogene and resume growth. This model is conceptually similar to that of reversible drug-tolerant persister cells²⁸, except that the source of drug tolerance is directly connected to the repression of the oncogene. This model may be particularly relevant for virally induced cancers, in which antiviral mechanisms may drive such oncogenic diversity. Future studies will determine if this model is relevant in additional contexts and might reveal new opportunities to eradicate the elusive persister cells in OPSCC and other cancers.

Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41588-023-01357-3>.

References

- Gillison, M. L. et al. Distinct risk factor profiles for human papillomavirus type 16-positive and human papillomavirus type 16-negative head and neck cancers. *J. Natl Cancer Inst.* **100**, 407–420 (2008).
- Ang, K. K. et al. Human papillomavirus and survival of patients with oropharyngeal cancer. *N. Engl. J. Med.* **363**, 24–35 (2010).
- Brianti, P., De Flaminio, E. & Mercuri, S. R. Review of HPV-related diseases and cancers. *Nat. Microbiol.* **40**, 80–85 (2017).
- Doorbar, J., Egawa, N., Griffin, H., Kranjec, C. & Murakami, I. Human papillomavirus molecular biology and disease association. *Rev. Med. Virol.* **25**, 2–23 (2015).
- Graham, S. V. The human papillomavirus replication cycle, and its links to cancer progression: a comprehensive review. *Clin. Sci. Lond. Engl.* **131**, 2201–2221 (2017).
- Litwin, T. R., Clarke, M. A., Dean, M. & Wentzensen, N. Somatic host cell alterations in HPV carcinogenesis. *Viruses* **9**, E206 (2017).
- Parikh, A. et al. Malignant cell-specific CXCL14 promotes tumor lymphocyte infiltration in oral cavity squamous cell carcinoma. *J. Immunother. Cancer* **8**, e001048 (2020).
- Puram, S. V. et al. Single-cell transcriptomic analysis of primary and metastatic tumor ecosystems in head and neck cancer. *Cell* **171**, 1611–1624 (2017).
- Qi, Z., Barrett, T., Parikh, A. S., Tirosh, I. & Puram, S. V. Single-cell sequencing and its applications in head and neck cancer. *Oral Oncol.* **99**, 104441 (2019).
- Qi, Z. et al. Single-cell deconvolution of head and neck squamous cell carcinoma. *Cancers* **13**, 1230 (2021).

11. Castellsagué, X. et al. HPV involvement in head and neck cancers: comprehensive assessment of biomarkers in 3680 patients. *J. Natl Cancer Inst.* **108**, djv403 (2016).
 12. Ramqvist, T. et al. Studies on human papillomavirus (HPV) 16 E2, E5 and E7 mRNA in HPV-positive tonsillar and base of tongue cancer in relation to clinical outcome and immunological parameters. *Oral Oncol.* **51**, 1126–1131 (2015).
 13. Neftel, C. et al. An integrative model of cellular states, plasticity, and genetics for glioblastoma. *Cell* **178**, 835–849 (2019).
 14. Tirosh, I. et al. Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. *Science* **352**, 189–196 (2016).
 15. Tirosh, I. et al. Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma. *Nature* **539**, 309–313 (2016).
 16. Slootweg, P. J., Hordijk, G. J., Schade, Y., van Es, R. J. J. & Koole, R. Treatment failure and margin status in head and neck cancer. A critical view on the potential value of molecular pathology. *Oral Oncol.* **38**, 500–503 (2002).
 17. Maynard, A. et al. Therapy-induced evolution of human lung cancer revealed by single-cell RNA sequencing. *Cell* **182**, 1232–1251 (2020).
 18. Cancer Genome Atlas Network. Comprehensive genomic characterization of head and neck squamous cell carcinomas. *Nature* **517**, 576–582 (2015).
 19. Kinker, G. S. et al. Pan-cancer single-cell RNA-seq identifies recurring programs of cellular heterogeneity. *Nat. Genet.* **52**, 1208–1218 (2020).
 20. Parikh, A. S. et al. Immunohistochemical quantification of partial-EMT in oral cavity squamous cell carcinoma primary tumors is associated with nodal metastasis. *Oral Oncol.* **99**, 104458 (2019).
 21. Ji, A. L. et al. Multimodal analysis of composition and spatial architecture in human squamous cell carcinoma. *Cell* **182**, 497–514 (2020).
 22. Yao, J. et al. Single-cell transcriptomic analysis in a mouse model deciphers cell transition states in the multistep development of esophageal cancer. *Nat. Commun.* **11**, 3715 (2020).
 23. Gerlee, P. & Nelander, S. The impact of phenotypic switching on glioblastoma growth and invasion. *PLoS Comput. Biol.* **8**, e1002556 (2012).
 24. Giese, A. et al. Dichotomy of astrocytoma migration and proliferation. *Int. J. Cancer* **67**, 275–282 (1996).
 25. Akagi, K. et al. Genome-wide analysis of HPV integration in human cancers reveals recurrent, focal genomic instability. *Genome Res.* **24**, 185–199 (2014).
 26. Duensing, S. & Münger, K. The human papillomavirus type 16 E6 and E7 oncoproteins independently induce numerical and structural chromosome instability. *Cancer Res.* **62**, 7075–7082 (2002).
 27. Korzeniewski, N., Spardy, N., Duensing, A. & Duensing, S. Genomic instability and cancer: lessons learned from human papillomaviruses. *Cancer Lett.* **305**, 113–122 (2011).
 28. Shen, S., Vagner, S. & Robert, C. Persistent cancer cells: the deadly survivors. *Cell* **183**, 860–874 (2020).
- Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.
- Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.
- © The Author(s), under exclusive licence to Springer Nature America, Inc. 2023

Methods

Human tumor specimens

Patients with OPSCC at the Washington University School of Medicine gave informed consent preoperatively to take part in the study following institutional review board approval (201911095 and 201102323), complying with all relevant ethical regulations. A total of 16 patients were included in the study and received no compensation for providing tissue samples. Twelve patients were initially clinically classified as HPV-positive based on p16 staining performed in a CLIA-certified clinical laboratory and interpreted by a dedicated head and neck pathologist, while four were classified as HPV-negative. After further analysis of HPV-related reads, one of the HPV-positive samples (OP8) was reclassified as HPV-negative due to the absence of any detectable HPV reads (Detection of rare HPV genotypes). Age, sex and other demographic characteristics of human subjects providing samples are summarized in Supplementary Table 1, while pathologic features are summarized in Supplementary Table 2.

Sample processing and sequencing

Fresh biopsies of OPSCC were collected from the primary tumor at the time of surgical resection, and in some cases, additional tissue was obtained from metastatic lymph node tissue. For histologically negative margins, the surgeon thoroughly irrigated the primary tumor defect site and then obtained a separate biopsy just beyond the intraoperative, frozen section margin sent to pathology. In all cases, the intraoperative, frozen margin analysis returned clear and final permanent margin status was also confirmed to be negative. A small fragment was snap-frozen for bulk whole exome sequencing, and the remainder of the provided tissue was processed for scRNA-seq. Fresh samples were minced, washed with PBS (Thermo Fisher Scientific) and dissociated using a Human Tumor Dissociation Kit (Miltenyi Biotec) per manufacturer guidelines. Red blood cell lysis was performed with Ammonium-Chloride-Potassium lysis buffer per manufacturer protocol (Thermo Fisher Scientific), followed by dead cell removal using a dead cell removal kit to improve the viability if needed (Miltenyi Biotec). Viability was confirmed to be >80% in all samples based on trypan blue analysis (Thermo Fisher Scientific). Cell suspensions were filtered using a 40 μ m filter (Thermo Fisher Scientific) and dissociated cells were pelleted and resuspended in AutoMACS Rinsing Solution with 0.5% BSA (Miltenyi Biotec). The single-cell suspension was sorted using human CD45 magnetic MicroBeads (Miltenyi Biotec) to enrich CD45⁺ cells. Briefly, 20 μ l of CD45 MicroBeads per 10⁷ total cells was added to the cell suspension and incubated for 15 min at 2–8 °C. After incubation, CD45⁺ and CD45[−] cells were collected after the cells passed through the magnetic column. The CD45⁺ and CD45[−] cell pellets were then obtained after centrifuging at 450g for 5 min at 4 °C. Samples were processed using the Chromium Single Cell 3' (v2 Chemistry), and in two cases 5', platform with the target of ~10,000 cells (10x Genomics) following the manufacturer's instructions. Briefly, cells were added onto a chip to form Gel Bead-in-Emulsion in the Chromium instrument followed by cell lysis, barcoding, fragmentation, adapter ligation and addition of sample index to the libraries before sequencing. scRNA-seq libraries were sequenced on Illumina NovaSeq 6,000 machines with a minimal target read count of 0.5 billion per sample. In four patients (OP4, OP6, OP9 and OP14), including the two (OP9 and OP14) who underwent 5' sequencing, CD45⁺ and CD45[−] cell fractions were sequenced separately. In the remaining cases, cells from the two fractions were mixed at a CD45⁺:CD45[−] ratio of 1:2 (twice as many CD45[−] cells) before single-cell barcoding. After sequencing, the resulting FASTQ files were aligned to a custom genome, combining the human genome (grch38) with genomes of the main high-risk HPV genotypes—HPV16, 18, 31, 33 and 35—using Cell Ranger v4.0. Assemblies for the high-risk HPV types were downloaded from NCBI with the following GenBank accession numbers: HPV16 (GCA_000863945.2), HPV18 (GCA_000865665.1), HPV31 (GCA_003179095.1), HPV33 (GCA_003179955.1) and HPV35 (GCA_003180695.1). All viral FASTA files were concatenated onto the

grch38 FASTA. All viral gtf files were adapted for Cell Ranger usage with the mkgtf function in Cell Ranger and concatenated onto the grch38 gtf. Thereafter, cellranger mkref was used on the new FASTA and gtf to create the custom reference. All reads that aligned to HPV genes were aligned exclusively to HPV16.

Detection of rare HPV genotypes

One sample classified as HPV-positive by p16 staining turned out not to have any reads for the genotypes used in our alignment (16, 18, 31, 33 and 35). For this particular patient, single-cell alignment was also done against HPV45, without any aligned reads. To make sure that this sample actually was HPV-negative, we also tested this patient for 13 HPV genotypes (16, 18, 31, 33, 35, 39, 45, 52, 56, 58, 59, 66 and 68), all of which were negative (Supplementary Table 9), using a previously validated RT-PCR method²⁹.

Furthermore, to exclude the possibility that a rare HPV genotype might still be present, we used the HPV-EM tool³⁰, which detects all known HPV genotypes in human sequencing data, to reanalyze a number of known HPV-positive samples as positive controls in addition to OP12 (known HPV-negative, negative control) and OP8 (patient in question). The FASTQ files from each patient were treated as a pseudobulk sample and aligned to the human genome, whereafter unmapped reads were mapped, allowing for mismatches, to all known HPV genotypes. In the HPV-positive patients, *HPV_{on}* and *HPV_{off}* cells were analyzed separately. While we did find HPV16 reads in *HPV_{on}* cells from all analyzed HPV-positive patients in this cohort, as expected, we did not have any reads mapped to any HPV genotype in OP12 or in OP8 or in the *HPV_{off}* cells. No HPV genotypes other than HPV16 were detected in any sample (Supplementary Fig. 2).

Cell lines

HNSCC HPV-positive cell lines SCC47, 93VU147T and SCC90 were cultured in 3:1 Ham's F12 (Thermo Fisher Scientific):DMEM (Thermo Fisher Scientific) supplemented with 10% fetal bovine serum (Peak Serum) and 1% penicillin-streptomycin-glutamine (Thermo Fisher Scientific). Cells were maintained at or below confluency of 90% to optimize growth conditions.

Filtering and preprocessing

For each sample, cells with fewer than 1,000 detected genes were removed as low-quality. Doublets were identified by combining the results of three alternative methods that were published recently and implemented in R packages—scdBlFinder³¹, the hybrid score from scds³² and the doubletCells algorithm from scan³³. For each method, we set the expected doublet rate at 0.6%, per 500 cells per sample. Cells classified as doublets by at least two methods, totaling 2,744 of all 73,714 cells (3.7%), were removed as probable doublets.

The UMI matrix was transformed to CPM (counts per million) by normalizing every gene by the total number of UMIs per sample. The CPM-matrix was then log₂-transformed, as log₂(CPM/10+1). Then, the data were mean-centered by subtracting the average expression of each gene from all values of that gene. We further filtered out the data, keeping only genes with either an average expression of >4 log₂ (CPM) across all cells or genes with >5 UMI counts in >20 cells. After filtering and doublet removal, the resulting matrix had 10,034 genes and 70,970 cells.

Scoring cells for gene expression signatures

Cells were scored for expression of various gene expression signatures, following our previously used approach¹³. For each cell, a relative expression score was defined by subtracting the average expression of the gene signature in a cell by that of a control gene set. The control gene set was defined by dividing all analyzed genes into 30 bins by average expression level and for each gene in the gene signature randomly sampling 100 genes from the same bin.

Cell-type assignment

The gene-cell matrix underwent dimension reduction using UMAP and Louvain clustering ($k = 200$), with each cluster assigned as either epithelial, stromal or immune based on the clusters' top 50 differentially expressed genes. All cells were also individually assigned to cell types by scoring for the expression of cell-type signature genes (Supplementary Table 3). Cells fulfilling either of the three following conditions: (1) HPV-positive cell in nonepithelial cluster, (2) highest cell signature score discordant with cluster assignment or (3) highest cell signature score less than $1.15 \times$ second highest signature score AND second highest signature scoring cell-type discordant with cluster assignment were set as unresolved and removed from further analysis; 3,342 cells were filtered out using this approach. Assignment to any nonimmune subtype in an immune cluster was defined as discordant and likewise for stromal and epithelial clusters. Cells classified as fibroblasts in epithelial clusters were kept, so as not to miss malignant cells undergoing EMT.

Cells assigned to any of the immune cell types were reclassified separately. A matrix was formed from just the immune cells, and batch correction was applied to samples from four patients—OP4, OP6, OP9 and OP14, because these samples had their CD45⁺ fractions sequenced separately and showed batch effects. Batch correction was performed by assigning every cell to an immune cell type by individual scoring as described above, followed by centering the expression values of cells from the four above-mentioned patients to the expression values of all other cells from the same cell type. Final assignments were achieved through dimension reduction and clustering of the corrected matrix.

Differential gene expression analysis

Whenever differential gene expression analysis was performed, a new UMI matrix was created containing only the relevant cells. It was then \log_2 (CPM/10+1)-transformed, filtered to only keep highly expressed genes and mean-centered as described above. P values were corrected using the Benjamini–Hochberg method.

Nonnegative matrix factorization

Diversity within cell types was studied through NMF. For every cell type, patients with at least 30 cells belonging to that cell type were selected. For each patient and the cells of that cell type, a new matrix was created by gene filtering and centering as described above. All negative values were set to zero, and NMF was performed using the *snmf/r* factorization algorithm from the NMF R package³⁴. For each sample and cell type, the algorithm was run 100 times and the factorization yielding the lowest approximation error was kept.

Every matrix was split into ten factors, each represented by 100 factor genes. The cells were then assigned to the factor with the highest average factor-genes expression. Factors for which fewer than ten cells were assigned were removed. The factor-gene lists from all patients were then compared, and Jaccard similarities were calculated between every pair of gene lists. Factors that did not have Jaccard similarities ≥ 0.2 with any other factor were removed because these did not represent recurrent programs. For the malignant cells, a total of 69 factors were kept for downstream analysis. The remaining factors were then hierarchically clustered, using Euclidean (1-Jaccard similarity) distance as a distance metric, with average linkage.

Clusters of factors (meta-clusters) were then used to define subtypes. For each of the meta-clusters representing at least two patients, all genes present in $>50\%$ of patients included in that meta-cluster were defined as a subtype signature. Cells from every cell type were then assigned to subtypes by creating matrices consisting just of cells from that cell type and scoring every cell for the subtype signatures as described in the section 'Scoring Cells for Gene Expression Signatures'. Annotations for fibroblast and endothelial cell subtypes were aided by subtype data^{35,36}.

HPV-positive tumor signature

An HPV-positive malignant signature score was defined from the three patients where adjacent as well as tumor tissue was provided. First, all HPV-positive epithelial cells in adjacent tissue samples were excluded. Thereafter, differential expression analysis between epithelial cells from tumor samples and epithelial cells from adjacent tissue samples was performed for each patient separately. Genes that ranked among the top 50 overexpressed genes in either the tumor sample ('up in cancer' genes) or the adjacent sample ('down in cancer' genes), consistent with all three comparisons, were kept. The signature score was then defined as the average normalized expression of the former genes minus the average normalized expression of the latter genes.

CNA inference

To define malignant cells, we first inferred CNAs from single-cell data using the method earlier published¹⁵. For each patient, a matrix was created from all epithelial and stromal (endothelial/fibroblasts) cells from that patient. The matrix was filtered, normalized and centered as described above. The genes were then reordered according to their chromosomal position, and extreme values of normalized expression were limited by setting the extremes at -3 and 3 . For each chromosome separately, a moving average was calculated at every chromosomal position by using a 100-gene window. As a baseline reference for normalization, an average CNA value was calculated for each stromal cell type, defining multiple potential reference profiles that represent cells with normal karyotype. Then, for each positive CNA value, we subtracted the maximum value of these potential references and, for each negative CNA value, we subtracted the minimum value of these potential references. Finally, all values between -0.15 and 0.15 , which we consider as likely reflecting noise rather than a genuine CNA signal, were set to zero. This resulted in a final matrix of CNA signal by cell.

Subclone assignments

The epithelial cells in the matrix of CNA values, derived as described above, were clustered to identify genetic subclones. Given the difficulty in selecting optimal clustering parameters, our approach was to initially use parameters that define a relatively large number of clusters (overclustering) and, subsequently, merge clusters that have the same set of inferred aberrations. Specifically, the matrix was first filtered, keeping only the top 2/3 of genes by the absolute value of their CNA signal. The matrix was then subjected to dimension reduction through UMAP, followed by the overclustering of the UMAP coordinate matrix through Louvain clustering with k set at 15. Clusters containing fewer than ten cells were merged into the most similar larger cluster by KNN distance, with $k = \ln$ (number of cells). For each cluster, an average CNA value for all cells in that cluster was calculated for every chromosome arm. Clusters were assigned as deleted or amplified at a chromosome arm if the average CNA value across the cells in the cluster over the genes in the chromosome arm was ≤ -0.15 or > 0.15 , respectively. Clusters were then merged if they satisfied two requirements as follows: (1) equal assignments across all chromosome arms and (2) maximum difference between clusters (across all chromosome arm) smaller than 0.15. This merging process was repeated, with new average values calculated, until all remaining clusters differed by at least one chromosome arm.

Malignant cell definitions

To separate malignant cells from nonmalignant epithelial cells, two metrics—CNA signal and CNA correlation—were calculated for each epithelial cell. CNA signal was defined as the average of absolute CNA values across the top 2/3 of genes by the CNA value. CNA correlation was defined as the correlation between the CNA values of every cell and the average CNA profile of the top 25% epithelial cells by CNA signal. For every subclone analyzed, cutoffs were set for both CNA signal and CNA correlation so that $<1\%$ of cells passing each threshold were stromal reference cells.

Cells passing both thresholds were classified as malignant cells, cells passing neither were classified as nonmalignant epithelial cells and those passing only one threshold were classified as unresolved.

Signature score comparisons by binning

The expression of gene set signatures between HPV-related subsets of malignant cells was compared through a binning approach. First, all cells received a signature score as described above. Cells were then ranked by their expression score and divided into ten bins of equal size. For each subset of cells, 100 cells per patient within the subset were randomly sampled to control for an uneven number of cells across patients. Then, for each bin and subset, the actual number of cells in that bin was compared to the expected number and provided an even distribution across all bins. This process was repeated 100 times, and a mean value of observed/expected was calculated across all runs.

Comparison of G1/S scores across datasets

Nine external scRNA-seq cancer datasets, acquired through 10x sequencing and comprising a total of 60,056 cancer cells,^{21,37–42} were used for this analysis. For each dataset, the cells annotated as malignant by the authors were selected and the UMI matrix \log_2 (CPM/10+1)-transformed as described above. Neither were genes filtered nor were expression values normalized by centering. Each dataset was then separately scored for the genes in our G1/S signature as described above. The *HPVon*, *HPVoff* and *HPVneg* cells from our study were considered as three separate datasets for this analysis and processed in the same way.

From each dataset, 1,000 cells were randomly sampled and all cells were divided into ten bins of equal size, ranked by G1/S score. We then looked at what proportion of cells per dataset were in the top three bins, representing a high G1/S score. This sampling was repeated 100 times, and the mean fraction of cells in the top bins, as well as standard error across the 100 runs, was used for the comparison of datasets.

TCGA analyses

An earlier study⁴³ provided data on HPV reads in ppm for each sample in a subset of the TCGA HNSCC cohort. Forty patients with available HPV data, of which 28 were HPV-positive, were oropharyngeal (ICD codes C01, C01.9, C09.9 and C10.9) and were used for our analysis. HPV read counts in ppm were \log_2 (ppm)-transformed, and mRNA expression data were \log_2 (TPM)-transformed. To account for differences in sample composition, we scored each sample for expression of epithelial genes (Supplementary Table 3) and created a linear model where we regressed all gene expression scores of interest, including HPV expression, against the epithelial score. Model residuals were used as final scores for downstream analysis. Grouping of patients into HPVhigh and HPVlow groups for survival analysis was done using the maxstat algorithm⁴⁴.

Cell cycle assignment

To set cutoffs for defining which epithelial cells express the cell cycle programs of the G1/S or G2/M phases and can thus be defined as cycling, we reasoned that (1) the vast majority of nonmalignant epithelial cells would not be cycling; and (2) by permuting the expression of each cell cycle gene across the nonmalignant epithelial cells we could reduce the signal of the potential few cycling cells, thereby defining reference profiles representing noncycling epithelial cells. We thus scored all the epithelial cells, as well as the permuted nonmalignant cells, for the G1/S and G2/M signatures. We then used the maximal observed scores of the permuted nonmalignant cells as cutoffs for the G1/S and G2/M signature scores.

Flow cytometry and single-cell clone expansion

Upon reaching 80% confluence, SCC47 and 93VU147T cells were trypsinized from the plate and filtered through 40 μ m filters (Thermo Fisher

Scientific). Filtered cells were washed with PBS and suspended in Hank's Balanced Salt Solution supplemented with 2 mM EDTA. To obtain a single cell from the bulk suspension, we performed cell sorting at the Siteman Flow Cytometry Core (Washington University School of Medicine). Briefly, standard forward scatter height versus area criteria were used to discard doublets and capture singlets. Cells were sorted into a 96-well plate containing 100 μ l of complete growth medium with 1X penicillin-streptomycin (Thermo Fisher Scientific). Plates containing sorted single cells were spun down at 200g for 5 min and incubated at 37 °C and 5% CO₂. Plates were scanned for single-cell colonies via microscope as soon as small aggregates of cells were visible by microscope, with single clone-derived colonies usually appreciable 2 weeks later. The confirmed single clones were transferred to 12-well plates and incubated for an additional 2 weeks to expand the clonal populations. To characterize the HPV genetic and transcriptomic heterogeneity in SCC47 and 93VU147T cell lines, we selected ten single-cell clones from each cell line.

For flow cytometric cell cycle analyses, cells were fixed with 70% ice-cold ethanol and stained with propidium iodide (30 μ g ml⁻¹ of PI (Sigma) with 200 μ g ml⁻¹ of RNase (Sigma) in 0.1% of Triton X-100 (Sigma) in PBS) for 1 h at room temperature. Cell cycle analysis was completed with at least 10,000 cells using CytoFLEX Flow Cytometer, and data were analyzed using FlowJo v9.0 software. The gating strategy is shown in Supplementary Fig. 3.

Cell treatment

SCC47 was treated with 0.2 μ M of cisplatin (Sigma) or equal volume of DMSO (vehicle) for 72 h. HPV-positive 93VU147T cells were irradiated with a dose of 8 Gy, and the controls were untreated and then collected 24 h later. SCC47 and 93VU147T cells were treated with varying dosages of the DNMT inhibitor decitabine (generously provided by Dr Ting Wang) and the H3K27 histone methyltransferase EZH2 inhibitor tazemetostat (MedKoo), respectively, with DMSO (vehicle) treatment as a control. Cells were treated for 72 h before they were collected for expression analysis of HPV genes.

CRISPRi knockdown

Golden-Gate cloning protocol was used to clone the HPV16 E6 and E7 targeting sgRNA oligos (Supplementary Table 10) into the sgOpti lentiviral vector backbone (Addgene). Lentivirus was generated via cotransfection of CRISPR plasmids into 293T cells with psPAX2 packaging plasmid (containing the *GAG/POL genes*) and pMD2.G plasmid supplying the *VSVG* envelope gene (Addgene) using PEI (2 μ g ml⁻¹). dCAS9-KRAB (lenti-dCAS9-dKRAB-blast, Addgene) SCC47 and 93VU147T cells expressing catalytically inactive dead Cas9 fused to transcriptional repressor KRAB were transduced with viral supernatant in the presence of 5 μ g/ml of polbrene (Santa Cruz Biotech) and then selected with puromycin and blastocidin (Life Technologies). Knockdown of genes was confirmed by qPCR.

Nucleic acid extraction and reverse transcription

The genomic DNA and total RNA of the HPV-positive cell lines (SCC47 and 93VU147T) were extracted by DNeasy Blood & Tissue Kit (QIAGEN) and RNeasy Plus Mini Kit (QIAGEN), respectively, according to the manufacturer's instructions. We performed first-strand synthesis with Maxima H Minus Reverse Transcriptase (Thermo Fisher Scientific) per manufacturer protocols using a dT18VN primer (Supplementary Table 10). Both genomic DNA and cDNA were stored at -20 °C.

qPCR analysis

qPCR was used to quantify relative HPV gene copies and expression. Primer sequences are listed in Supplementary Table 10. For the quantification of relative viral gene copies, we used the single-copy gene (albumin) as an internal reference. The amplification efficiency of the viral genes (E6 and E7) was compared to that of albumin to estimate the

relative viral gene copies⁴⁵. We also generated a standard curve with a 10-fold dilution series of plasmids containing the albumin gene (spanning 5 logs from 10⁵ to 10⁹ copies) to estimate the albumin gene copies. To quantify the relative HPV expression, we used GAPDH as an internal reference. We performed the qPCR with either 50 ng of the genomic DNA for gene copy analysis or the cDNA for gene expression analysis per manufacturer protocols using the ABI QuantStudio Q3 (Applied Biosystems). We generated melting curves after each PCR and all samples yielded a single peak.

Matrigel invasion assay

Matrigel invasion assay was performed following an established protocol⁴⁶. Briefly, preformed matrigel invasion chambers (Corning) were prepared per manufacturer protocol. Serum-containing media was placed below the invasion chambers, and 1 × 10⁵ cells suspended in 200 µl serum-free media were placed above the invasion chambers and incubated for 24 h. Cells on the lower surface of the membrane were fixed with methanol, stained with crystal violet and counted in a blinded manner. Cells in serum-containing media were used as a negative control.

Immunocytochemistry

Cells were fixed in freshly prepared 4% paraformaldehyde for 20 min at room temperature, washed with PBS and subsequently blocked and permeabilized with 0.1% Triton X in 10% goat serum containing PBS at room temperature for 1 h. Cells were then probed with primary antibodies, Ki67 1:500 dilution (D2H0 rabbit mAb, Cell Signaling), E6 1:100 dilution (mouse antiviral, clone C1P5, Invitrogen), E7 1:100 dilution (mouse antiviral, clone TVG701Y, Invitrogen) diluted in 10% goat serum PBS and incubated overnight at 4 °C. Cells were washed with PBS and then probed with secondary antibody, goat antirabbit IgG (H+L), Fab2 Alexa Fluor 594, goat antimouse IgG (H+L) or Fab2 Alexa Fluor 488 (Cell Signaling) at 1:400 dilution in 2% goat serum containing PBS for 1 h at room temperature followed by PBS washes and mounted with DAPI (Fluoroshield, Sigma). Imaging was completed using Fluorescence Eclipse Ti2 Inverted microscope (Nikon).

Molecular fluorescent in situ hybridization

Viral RNA ISH (RNAScope) and DNA ISH (DNAScope) were performed using RNAScope 2.5 HD Reagent Kit Red assay combined with Immunohistochemistry (Advanced Cell Diagnostics; ACD) according to the manufacturer's instructions. Briefly, slides were baked in a dry air oven for 1 hour at 60 °C, deparaffinized (Xylene for five minutes twice followed by 100% ethanol for 2 min twice), hydrogen peroxide was applied for 10 min at room temperature, and codetection target retrieval was done using Steamer (BELL) for twenty minutes and washed with PBS-T. Tissue slides were then incubated overnight with p16-INK4a antibody (LSBio) in HybEZ Slide Rack in the Humidity Control Tray with damp humidifying paper and incubated overnight at 4 °C. The next day, postprimary fixation was done by washing slides with PBS-T and submerging slides in 10% NBF for 30 min at room temperature. Slides were washed with PBS-T, and Protease Plus was added to each slide for 30 min at 40 °C and then washed with distilled water. RNAScope antisense probes were utilized to target RNA of specified viral genes, while DNAScope sense probes were utilized to target DNA of specified viral genes⁴⁷. Selected probes were warmed at 40 °C and hybridized with specific oligonucleotide probes for 2 h at 40 °C in the HybEZ Humidifying System. Details of antibodies, probes and sequences are in Supplementary Table 11. RNA/DNA was then serially amplified and stained with Fast Red solution. Slides were blocked with a codetection blocker for 15 min at 40 °C and washed in PBS-T. Secondary Alexa Fluor 488 antibody (Abcam) was applied for 1 h at room temperature in the dark. Finally, slides were washed with PBS-T and counter-stained with DAPI (Sigma) and mounted with ProLong Gold Antifade Reagent (Invitrogen). RNAScope was optimized with a PPIB probe as a positive

control, while a DapB probe and no secondary antibody served as negative controls. All slides were imaged on the EVOS M5000 Imaging System (Invitrogen) for scoring.

Quantification was performed with CellProfiler using the ISH pipeline⁴⁸. Adjustments were made to accommodate cell size as well as green versus red staining. Dot staining was identified based on intensity and distinct pixel ranges for DAPI (nucleus, 20–50 pixels), green (p16, 10–30 pixels) and red (E6 or E7, 3–12 pixels). Cell size was identified using a 5-pixel radius from the nucleus, and images were overlaid to count dots per cell. A positive stain scoring of p16 was determined as greater than 1, while a negative stain score was determined as less than 1. Red (E6 and E7) RNA ISH signal was detected within individual cells and scored using ACD scoring bins. Bin scoring ranged from <1 dot/cell designated as bin 0, 1–3 dots per cell designated as bin 1, 4–9 dots per cell designated as bin 2, 10–15 dots per cell designated as bin 3, and >15 dots per cell designated as bin 4. Stain scoring remained the same for all genes of interest with positive staining determined as >0 dots per cell and negative staining as 0 dots per cell per ACD guidelines. DNA ISH signal was detected and scored in a similar fashion. RNase treatment was used to confirm that signal was specific to DNA. All RNA and DNA ISH images were reviewed with a dedicated head and neck pathologist, who confirmed the heterogeneous pattern of HPV RNA expression compared to a more homogeneous pattern of HPV DNA detected.

Dual-stain immunohistochemistry

FFPE blocks of the patient tumors were sectioned onto slides at 4 µm. Slides were baked at 60 degrees for 30 min followed by deparaffinization with xylene and graded ethanol. Diva Decloaker (Biocare Medical) was used for heat-mediated antigen retrieval for all stains. Blocking was performed with Dako Dual Endogenous Enzyme Block (5 min). HPV type 16/18 E6 Mouse Monoclonal antibody (1:50 dilution; Thermo Fisher Scientific, MA1-46057) was applied first and incubated for 30 min. Secondary antibody incubation was performed with the Dako EnVision+ Dual Link System-HRP for 30 min, followed by DAB staining for 5 min. Blocking with Dako Dual Endogenous Enzyme Block was then repeated for 5 min. Staining with P16-INK4A polyclonal antibody (1:75 dilution; Thermo Fisher Scientific, 10883-1-AP) was then performed with a 30 min incubation time. Dako PowerVision Poly-AP was used for secondary antibody staining (30 min), followed by incubation with AP Red substrate for 5 min. Sections were then mounted with a coverslip with Glycergel (Dako).

Cell proliferation

CellTiter-Glo (CTG) proliferation assays (Promega) were completed according to the manufacturer's protocols. Briefly, 2,000 cells were seeded per 96 wells in technical replicates of five. Cells were lysed on day 0 (1 h after seeding of cells), 1, 3, 5, 7 and 9 for HPV single clones and day 0, day 2, day 4, day 6, day 8, day 10 and day 12 for E6 and E7 HNSCC knockdown cell lines by the addition of the CTG reagent followed by the measurement of luminescence using the Biotek Cytation 5 (BioTek). Background luminescence was removed. Luminescence values were adjusted based on 2 µM Adenosine triphosphate luminescence measured on the same plate for each day.

Statistics

All functional experiments were performed with at least three independent biological replicates, with the number of replicates indicated in figure legends. Statistical analyses for functional experiments were performed with GraphPad Prism 4.0. All histograms are presented as mean + s.e.m. Student's *t*-test was used for comparisons in experiments with two sample groups. In experiments with more than two sample groups, ANOVA was performed followed by Bonferroni's post hoc test.

Software packages

Data analysis was performed in R (version 4.1.0) with the following packages used: caTools version 1.18.2 for calculating moving

averages, circlize version 0.4.13 for creating color palettes, class version 7.3–19 for classifying cells by kNN, clusterProfiler version 4.0.2 for enrichment analysis, ComplexHeatmap version 2.8.0 for plotting heatmaps, dplyr version 1.0.7 for data handling, FNN version 1.1.3 for creating kNN graphs, ggplot2 version 3.3.4 for creating plots, ggrepel version 0.9.1 for separating text labels in plots, gtools version 3.9.2 for random permutation of values, igraph version 1.2.6 for Louvain clustering, Matrix version 1.3–4 and Matrix.utils version 0.9.8 for working with sparse matrices, msigdb version 7.4.1 for enrichment analysis, NMF version 0.23.0 for performing NMF, parallel version 4.1.0 for parallelizing computation, reshape2 version 1.4.4 for data handling, scDBIFinder version 1.7.1, scds version 1.8.0 and scan version 1.22.1 for doublet detection, SingleCellExperiment version 1.14.1 for data handling, stringdist version 0.9.6.3 for calculating similarity between character strings and uwot version 0.1.10 for creating UMAP plots.

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

All scRNA-seq data produced by this study are available through the Gene Expression Omnibus with GEO accession [GSE182227](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE182227). TCGA bulk RNA-seq and clinical data for head and neck and cervical cancer are available through the Broad Genome Data Analysis Center Firehose (<https://gdac.broadinstitute.org/>). Single-cell datasets reanalyzed to compare proliferation rates are available through the Gene Expression Omnibus with accession numbers [GSE150430](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE150430) (nasopharyngeal carcinoma), [GSE131907](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE131907) (lung carcinoma), [GSE132465](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE132465), [GSE132257](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE132257), [GSE144735](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE144735) (CRC), [GSE125449](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE125449) (HCC), through the Chinese National Centre for Bioinformation Genome Sequence Archive (CNCB-GSA) with accession: [CRA001160](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=CRA001160) (PDAC) and through EMBL-EBI ArrayExpress with accession numbers [E-MTAB-8107](https://www.ebi.ac.uk/ena/arrayexpress/studies/E-MTAB-8107) (breast, ovarian, colorectal cancer), [E-MTAB-6149](https://www.ebi.ac.uk/ena/arrayexpress/studies/E-MTAB-6149) (lung) and [E-MTAB-6653](https://www.ebi.ac.uk/ena/arrayexpress/studies/E-MTAB-6653) (lung). Cell line data used for validation analysis are available through GEO with accession number [GSE157220](https://www.ncbi.nlm.nih.gov/geo/accession.cgi?acc=GSE157220). The NSCLC dataset used to validate finding malignant cells in normal samples is deposited as an NCBI BioProject with accession number [PRJNA591860](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA591860). Source data are provided with this paper.

Code availability

R code and functions for the analysis are available at https://github.com/micmin3/HPV_OPSCC_Analysis.

References

29. Gao, G. et al. A novel RT-PCR method for quantification of human papillomavirus transcripts in archived tissues and its application in oropharyngeal cancer prognosis. *Int. J. Cancer* **132**, 882–890 (2013).
30. Inkman, M. J. et al. HPV-EM: an accurate HPV detection and genotyping EM algorithm. *Sci. Rep.* **10**, 14340 (2020).
31. Germain, P.-L. scDBlFinder. R package version 1.6.0 <https://github.com/plger/scDBlFinder> (2021).
32. Bais, A. S. & Kostka, D. scds: computational annotation of doublets in single-cell RNA sequencing data. *Bioinform. Oxf. Engl.* **36**, 1150–1158 (2020).
33. Lun, A. T. L., McCarthy, D. J. & Marioni, J. C. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Res.* **5**, 2122 (2016).
34. Gaujoux, R. & Seoighe, C. A flexible R package for nonnegative matrix factorization. *BMC Bioinformatics.* **11**, 367 (2010).
35. Goveia, J. et al. An integrated gene expression landscape profiling approach to identify lung tumor endothelial cell heterogeneity and angiogenic candidates. *Cancer Cell* **37**, 21–36 (2020).
36. Buechler, M. B. et al. Cross-tissue organization of the fibroblast lineage. *Nature* **593**, 575–579 (2021).
37. Chen, Y.-P. et al. Single-cell transcriptomics reveals regulators underlying immune cell diversity and immune subtypes associated with prognosis in nasopharyngeal carcinoma. *Cell Res.* <https://doi.org/10.1038/s41422-020-0374-x> (2020).
38. Kim, N. et al. Single-cell RNA sequencing demonstrates the molecular and cellular reprogramming of metastatic lung adenocarcinoma. *Nat. Commun.* **11**, 2285 (2020).
39. Lee, H.-O. et al. Lineage-dependent gene expression programs influence the immune landscape of colorectal cancer. *Nat. Genet.* **52**, 594–603 (2020).
40. Ma, L. et al. Tumor cell biodiversity drives microenvironmental reprogramming in liver cancer. *Cancer Cell* **36**, 418–430 (2019).
41. Peng, J. et al. Single-cell RNA-seq highlights intra-tumoral heterogeneity and malignant progression in pancreatic ductal adenocarcinoma. *Cell Res.* **29**, 725–738 (2019).
42. Qian, J. et al. A pan-cancer blueprint of the heterogeneous tumor microenvironment revealed by single-cell profiling. *Cell Res.* **30**, 745–762 (2020).
43. Tang, K.-W., Alaei-Mahabadi, B., Samuelsson, T., Lindh, M. & Larsson, E. The landscape of viral expression and host gene fusion and adaptation in human cancer. *Nat. Commun.* **4**, 2513 (2013).
44. Ogłuszka, M., Orzechowska, M., Jędraszka, D., Witas, P. & Bednarek, A. K. Evaluate cutpoints: adaptable continuous data distribution system for determining survival in Kaplan–Meier estimator. *Comput. Methods Prog. Biomed.* **177**, 133–139 (2019).
45. Barczak, W., Suchorska, W., Rubiś, B. & Kulcenty, K. Universal real-time PCR-based assay for lentiviral titration. *Mol. Biotechnol.* **57**, 195–200 (2015).
46. Puram, S. V. et al. STAT3-iNOS signaling mediates EGFRvIII-induced glial proliferation and transformation. *J. Neurosci.* **32**, 7806–7818 (2012).
47. Deleage, C. et al. Defining HIV and SIV reservoirs in lymphoid tissues. *Pathog. Immun.* **1**, 68–106 (2016).
48. Erben, L., He, M.-X., Laeremans, A., Park, E. & Buonanno, A. A novel ultrasensitive in situ hybridization approach to detect short sequences and splice variants with cellular resolution. *Mol. Neurobiol.* **55**, 6169–6181 (2018).

Acknowledgements

This work was supported by the V Foundation (S.V.P.), Cancer Research Foundation (S.V.P.), Emerson Collective Cancer Research Fund (S.V.P.), Barnes Jewish Hospital Foundation (S.V.P.), NCI 1K08CA237732 (S.V.P.), Doris Duke Fund to Retain Clinician Scientists (S.V.P.), Doris Duke Clinician Scientist Development Award (S.V.P.), NIDCD T32DC000022 (T.F.B.), the Swedish Society of Medicine (M.M.), Washington University Department of Medicine Faculty Diversity Award (J.S.F.), Israel Science Foundation (I.T.), Rising Tide Foundation (S.V.P. and I.T.), Mexican Friends New Generation Grant (I.T.), Mauricio Schwarz (I.T.), Zuckerman STEM Leadership Program (I.T.) and The Dr Celia Zwillenberg-Fridman and Dr Lutz Zwillenberg Career Development Chair (I.T.). The funding sources had no involvement in the design, conduct and reporting of the research.

Author contributions

S.V.P., M.M. and I.T. conceived and designed the study, interpreted the results and wrote the paper. Biological validation experiments were performed by A.P., Z.Q., A.R., T.F.B. and T.L. and overseen by S.V.P., with the exception of ISH experiments which were performed by K.G. and overseen by J.S.F. and HPV typing which was performed by P.L. and overseen by X.W. The computational analysis was performed by M.M. and overseen by I.T. Tumor samples and associated information were

provided by S.G. and S.R. A.S.P., E.A.M., J.W.R., D.A., W.L.T., H.A.G, L.D., R.C.P., P.P., R.S.J., A.M., R.C. and J.P.Z provided essential input on study design, interpretation of results and critical paper review. All authors approved the paper.

Competing interests

I.T. is a member of the Scientific Advisory Board (SAB) of Immunitas Therapeutics. All other authors report no competing interests.

Additional information

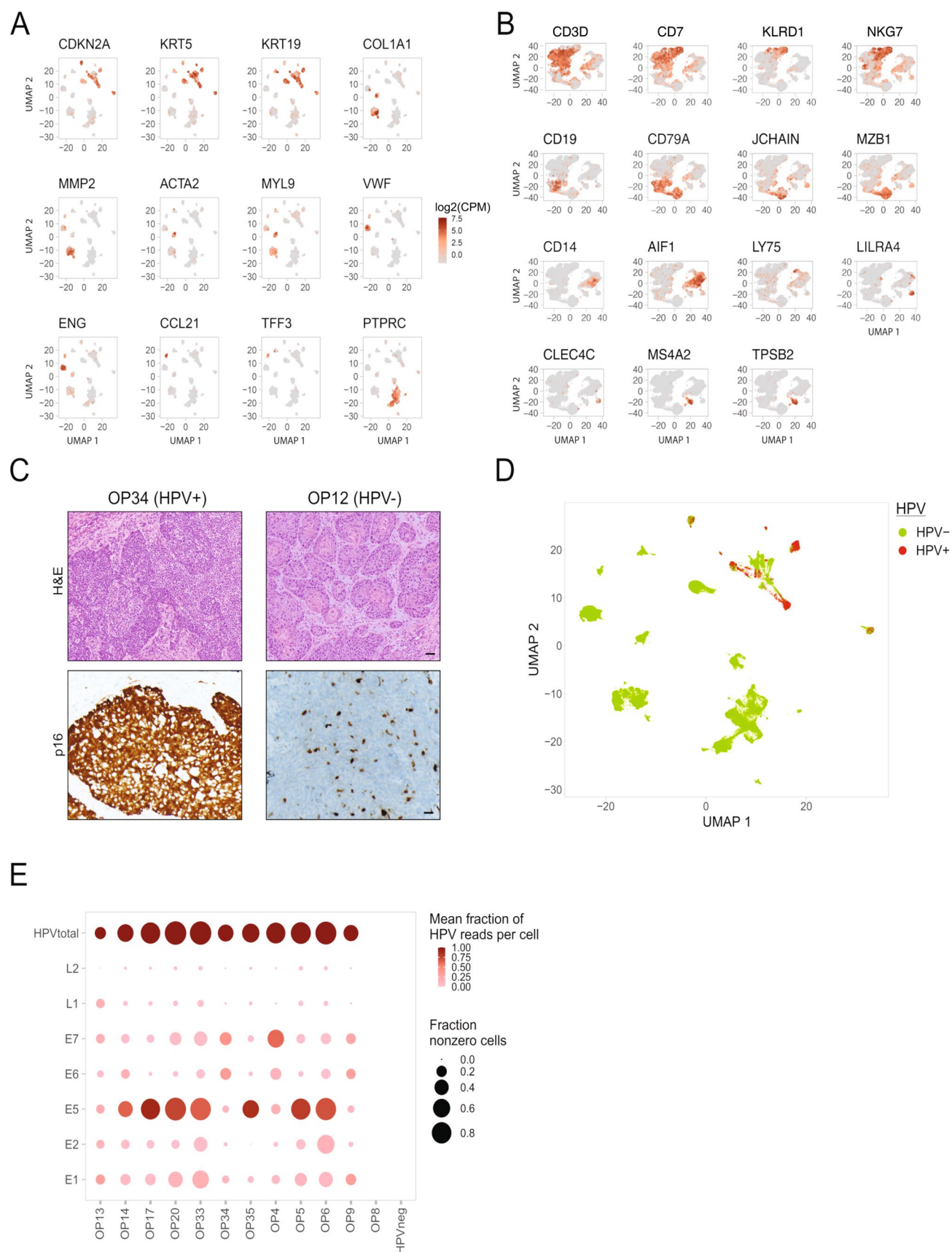
Extended data is available for this paper at <https://doi.org/10.1038/s41588-023-01357-3>.

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41588-023-01357-3>.

Correspondence and requests for materials should be addressed to Sidharth V. Puram or Itay Tirosh.

Peer review information *Nature Genetics* thanks David Sidransky, Lisa Mirabello, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

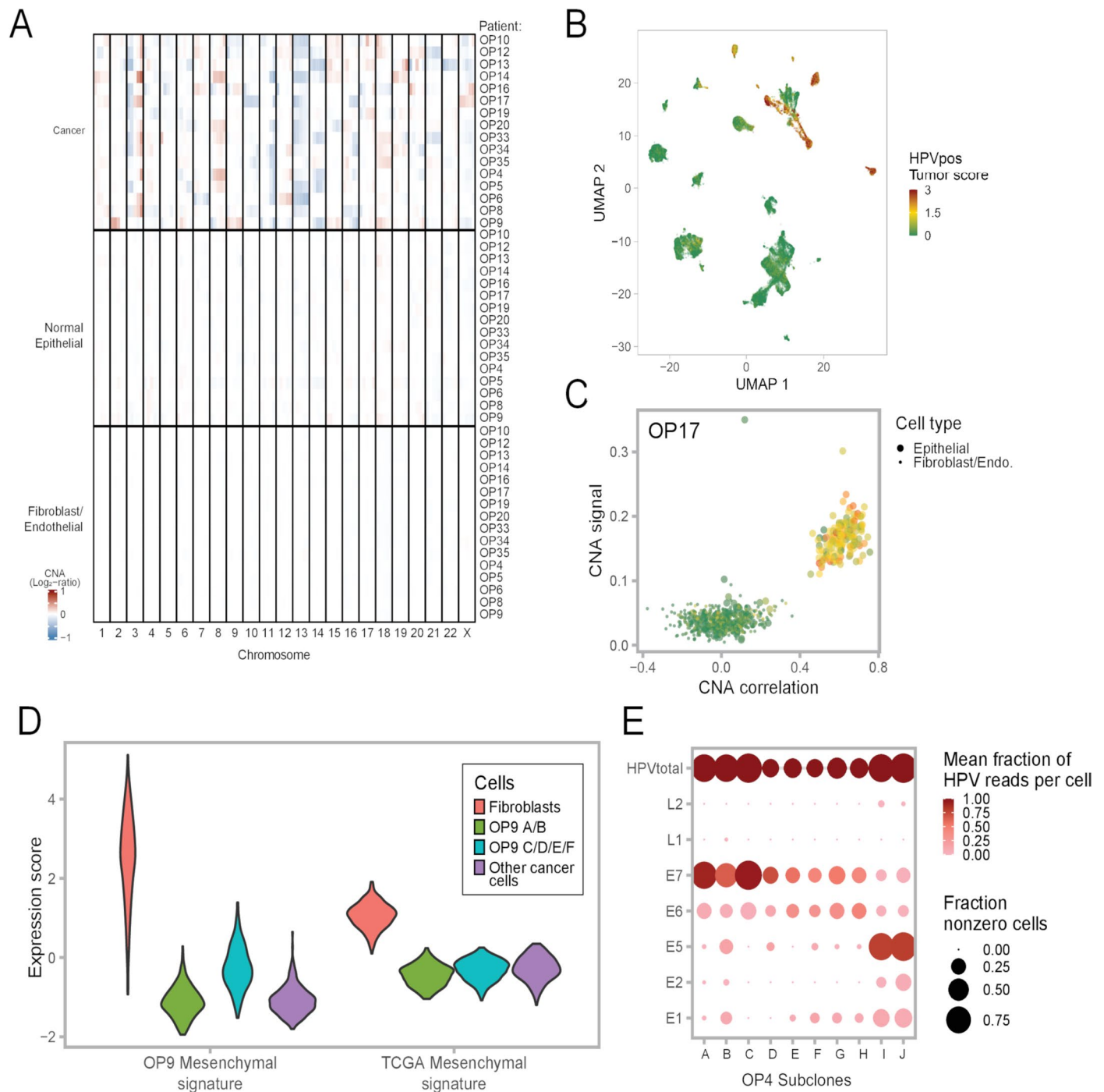
Reprints and permissions information is available at www.nature.com/reprints.



Extended Data Fig. 1 | See next page for caption.

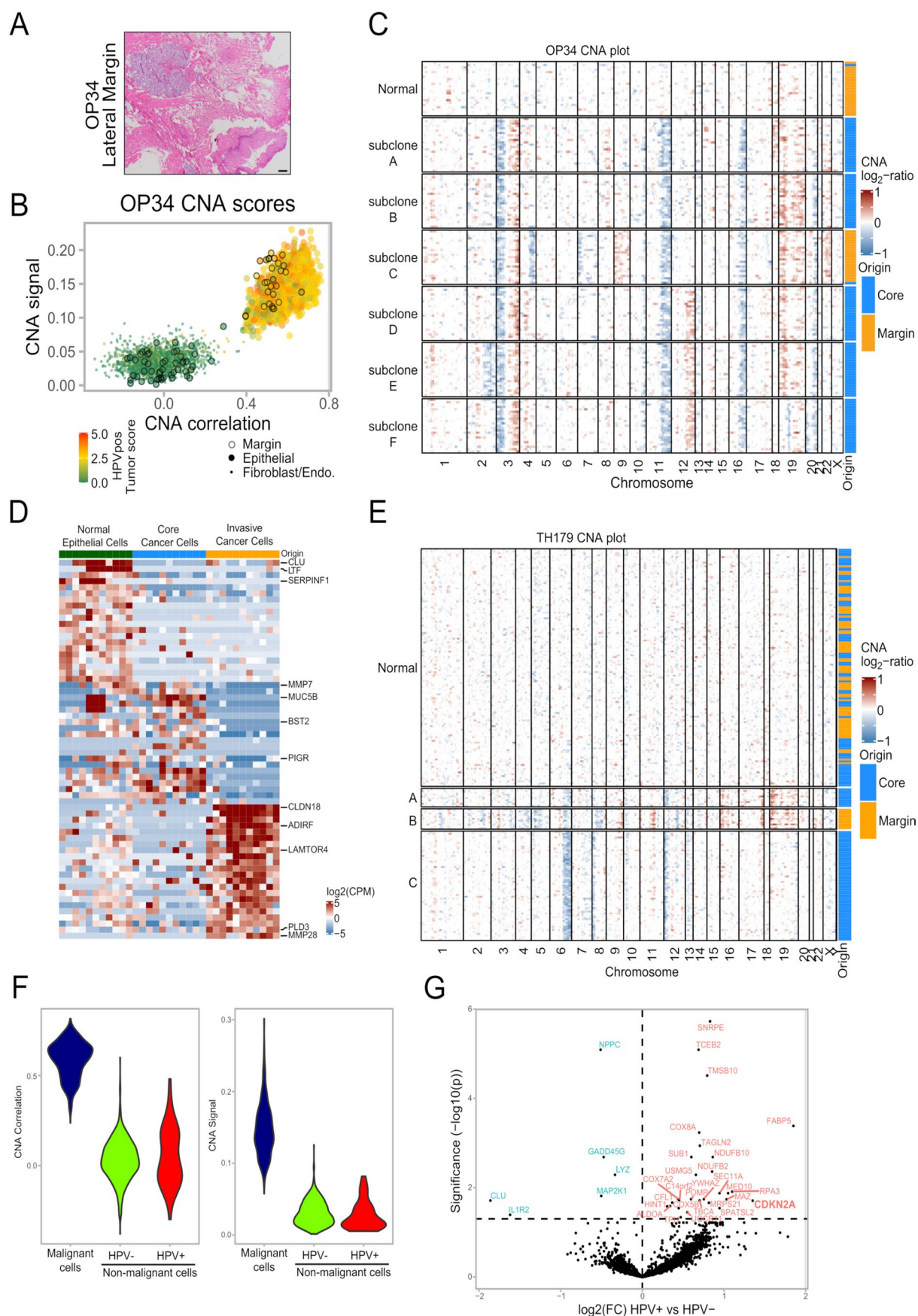
Extended Data Fig. 1 | Expression of marker genes and HPV genes, related to Fig. 1. (a) UMAP of all cells ($n = 70,970$) colored by expression of selected marker genes. (b) UMAP of immune cells ($n = 22,818$) colored by expression of selected marker genes. (c) Histologic sections of two representative HPV+ (p16+) and HPV- (p16-) oropharynx tumors (OP34 and OP12), stained by H&E (top) and p16 (bottom). Staining was repeated three independent times with similar results. Scale bar = 100 μm . (d) UMAP of all cells colored by detection of at least one

read from HPV16 genes. (e) Dot plot showing variability in expression of HPV genes (rows) across patients (columns). The last column summarizes all HPV-negative tumors. The top row shows the sum of HPV gene expression per patient (HPVtotal). The size of each dot represents the fraction of cells with at least one read for that gene in each patient, while the color represents the fraction of HPV reads in one patient that reflect the corresponding gene. For the latter metric, HPVtotal is set to 1.



Extended Data Fig. 2 | CNA patterns and controls, related to Fig. 2. (a) Average CNA profiles of malignant cells, normal epithelial cells and fibroblasts/endothelial cells used as reference for each patient. Each row is a cell subset within a patient. Rows are ordered by cell subset and patient ID. Columns are chromosomal positions. For each row and chromosome, the chromosome was split into five bins. **(b)** UMAP of all cells colored by HPV-positive tumor score. **(c)** CNA signal and correlation scatter plot of OP17. Cells are colored by their expression of the HPV-positive tumor score. **(d)** Violin plots showing

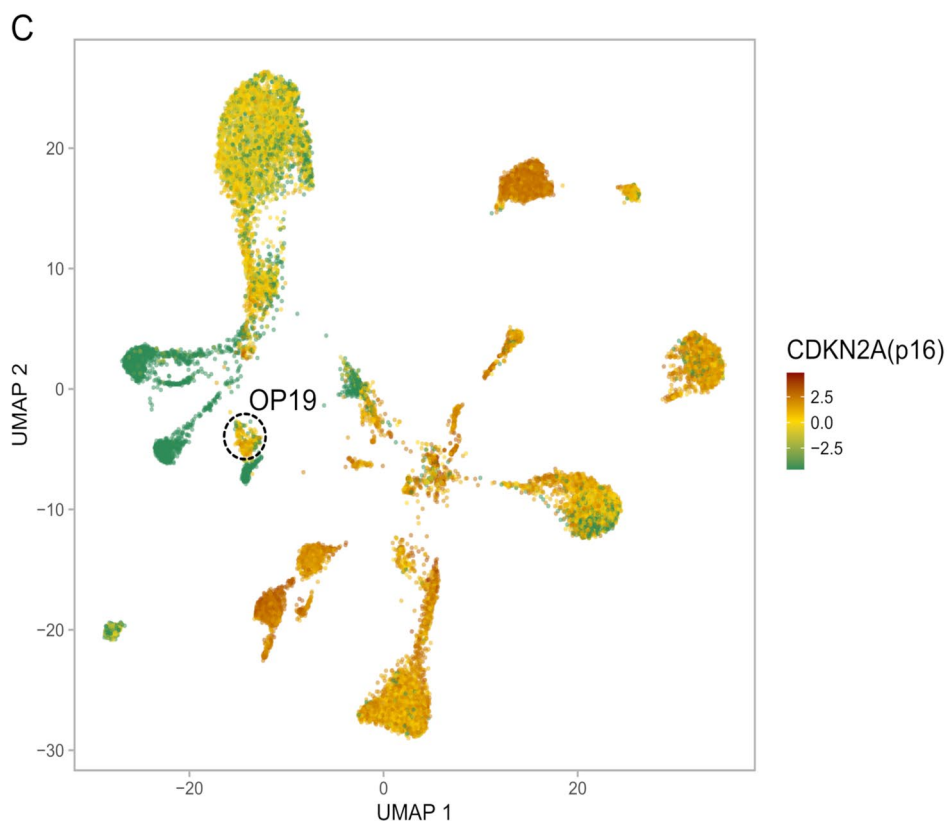
expression of the OP9 mesenchymal signature (left panel) and the TCGA HNSCC mesenchymal signature (right panel) in four subsets of cells; 300 cells were randomly sampled from each subset to ensure equal-sized groups. **(e)** Dot plot showing variability in HPV gene expression between subclones in one patient, OP4. The size of each dot represents the fraction of cells with at least one read for that gene in each subclone, while the color represents the fraction of HPV reads in one subclone that reflect the corresponding gene. For the latter metric, HPVtotal is set to 1.



Extended Data Fig. 3 | See next page for caption.

Extended Data Fig. 3 | CNA-based detection of invasive malignant cells, related to Fig. 2. (a) Histologic section of the lateral margin from OP34, stained by H&E. A piece of mucosa was taken beyond this histologically clear (pathologically negative) margin for scRNA-seq (labeled 'margin'). Staining was repeated three independent times with similar results. Scale bar = 1000 μ m. (b) CNA signal and correlation scatter plot of OP34. Cells are colored by their expression of the HPV-positive tumor score. Epithelial cells from the margin sample are circled. (c) CNA plot of OP34. Cells were randomly sampled from all subclones in equal numbers to ensure equal-sized groups. Column at the right shows the origin of cells from the tumor core and margin samples. (d) Heatmap of differentially expressed genes in the three epithelial cell subsets of lung adenocarcinoma sample TH179 – normal epithelial cells, invasive malignant

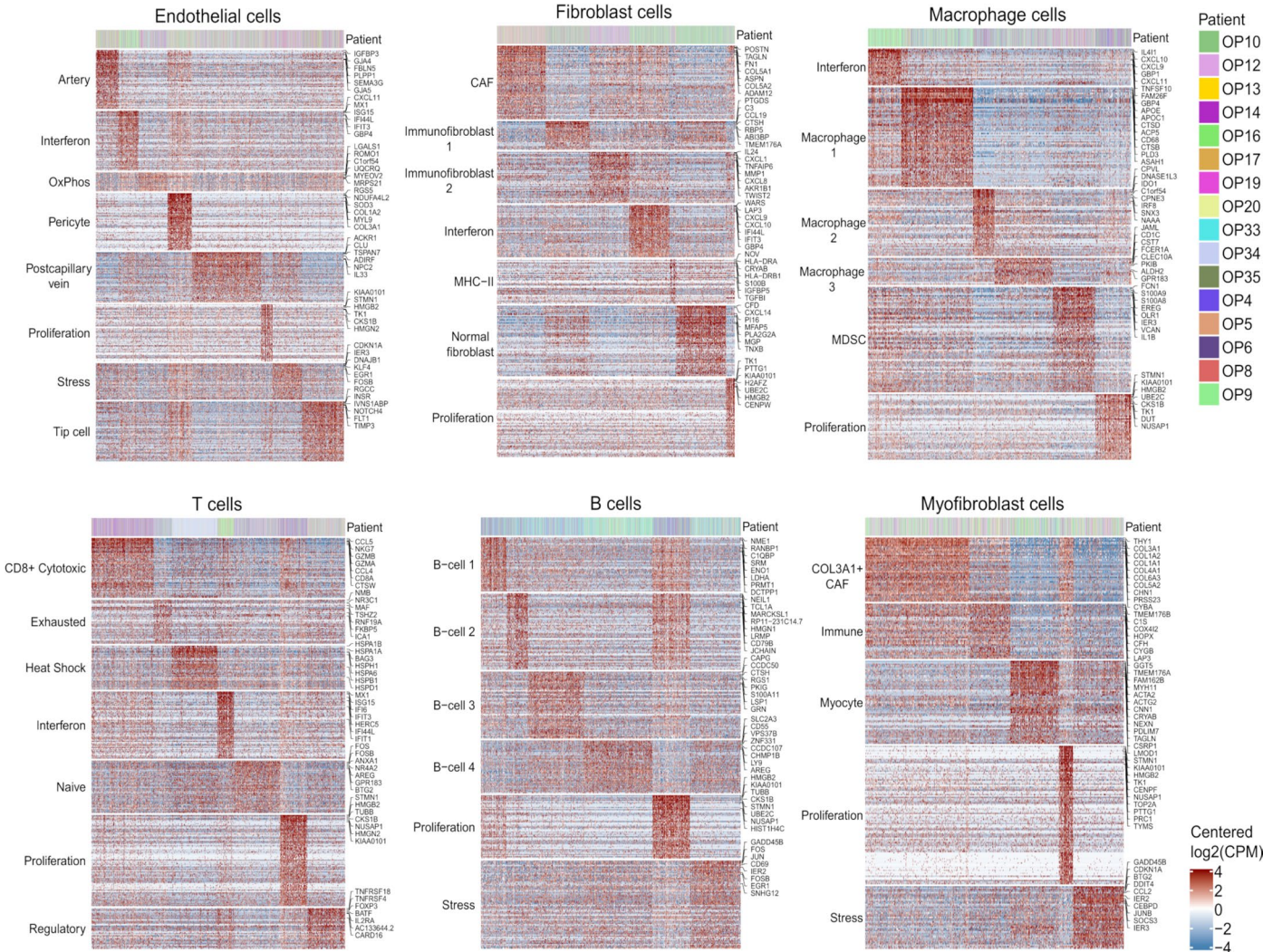
cells and malignant cells from the tumor core. Rows are genes, columns are cells. Cells were randomly sampled from the normal and core subsets to ensure equal-sized groups. (e) CNA plot of lung adenocarcinoma sample TH179. Column at the right shows the origin of cells from the tumor core and margin samples. (f) HPV expression in normal epithelial cells. Violin plots showing values for CNA signal and CNA correlation for the 51 HPV-positive and 779 HPV-negative negative nonmalignant epithelial cells from HPV-positive patients, as well as for 830 randomly sampled cancer cells from the same patients, one cancer cell per patient sampled per nonmalignant epithelial cell. (g) Volcano plot of differentially expressed genes between nonmalignant epithelial cells (defined by lack of CNAs) with or without HPV expression. P-value derived from two-sided t-test adjusted for multiple comparisons.



Nature Genetics

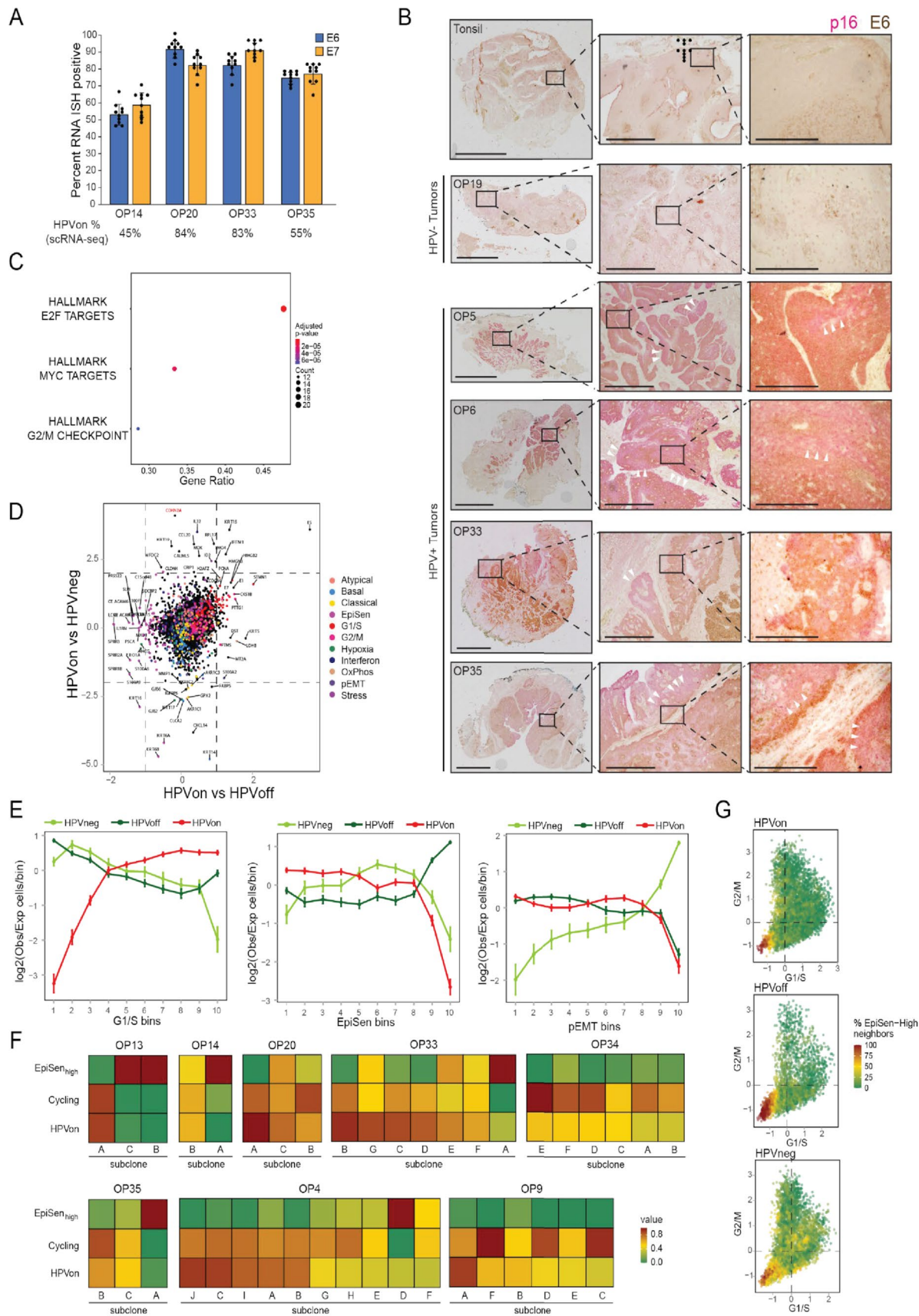
Extended Data Fig. 4 | Diversity of malignant cells across tumors, related to Fig. 3. (a) Heatmap showing relative expression of differentially expressed genes (rows) across all tumor samples (columns). Selected genes include the top 50 preferentially expressed genes from each tumor. **(b)** Hierarchical clustering of 'pseudobulk' tumor profiles (defined by averaging all malignant cells per sample). Shown are Pearson correlations, ordered by the clustering of samples.

Bottom panels show additional tumor characteristics with the same tumor ordering as in the heatmap, including (from top to bottom): the percentage of cells with detected HPV reads, the clinical HPV status (defined by p16 staining), three TCGA subtype scores, and scores for all meta-programs defined in Fig. 3c, d. **(c)** UMAP of all malignant cells, colored by mRNA expression of CDKN2A (encoding for p16). OP19 is circled.



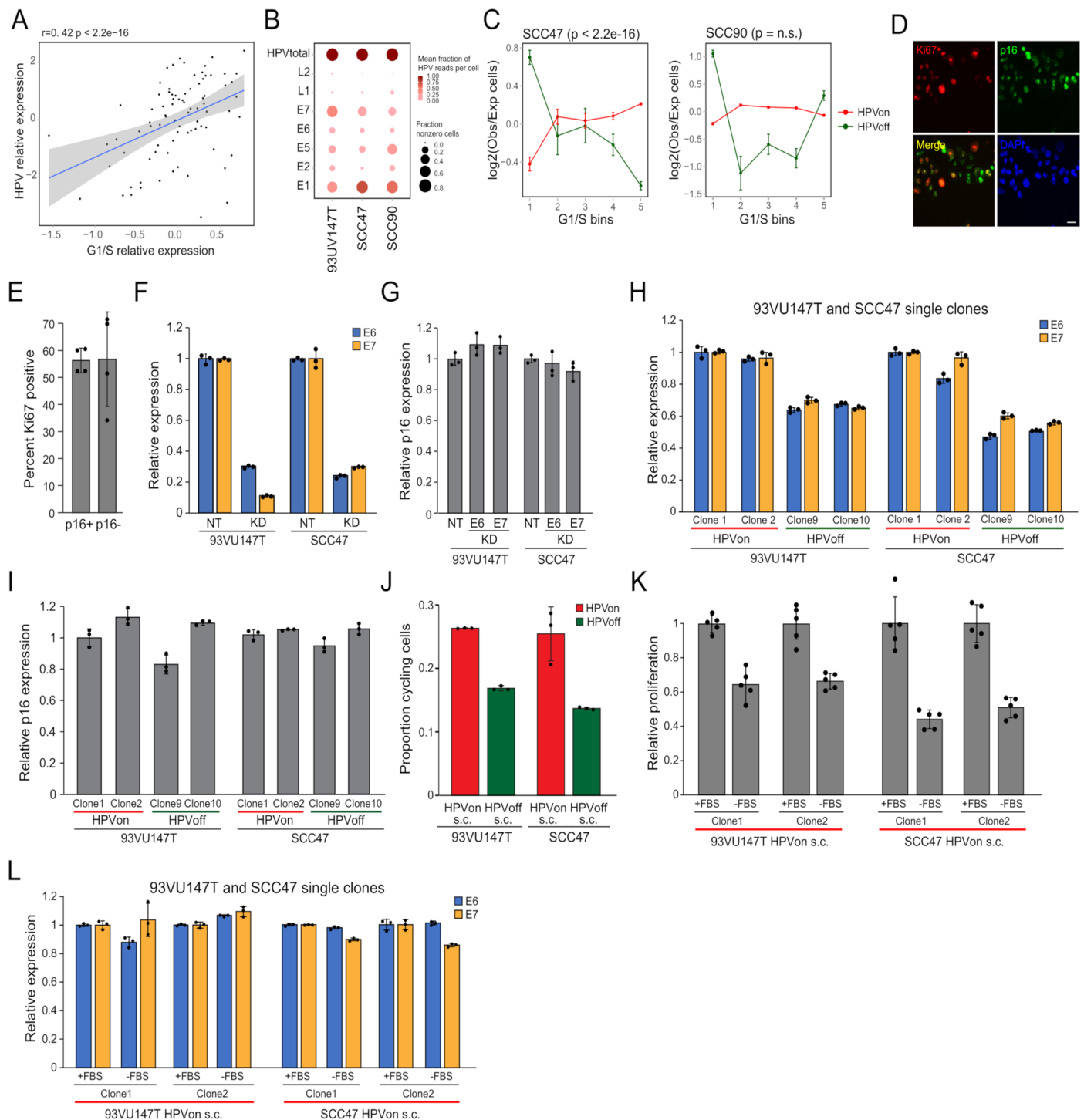
Extended Data Fig. 5 | Heterogeneity among common cell types in the OPSC microenvironment, related to Fig. 3. For each of the common cell types in the OPSC microenvironment (endothelial cells, fibroblasts, macrophages, T cells, B cells, and myofibroblasts), the corresponding panel shows meta-programs, as

defined using the same approach as performed for malignant cells and shown in Fig. 3d. Shown are the relative expression levels of meta-program genes (rows) in all cells of the corresponding cell types (columns). Top panels indicate the patient of origin for all cells.



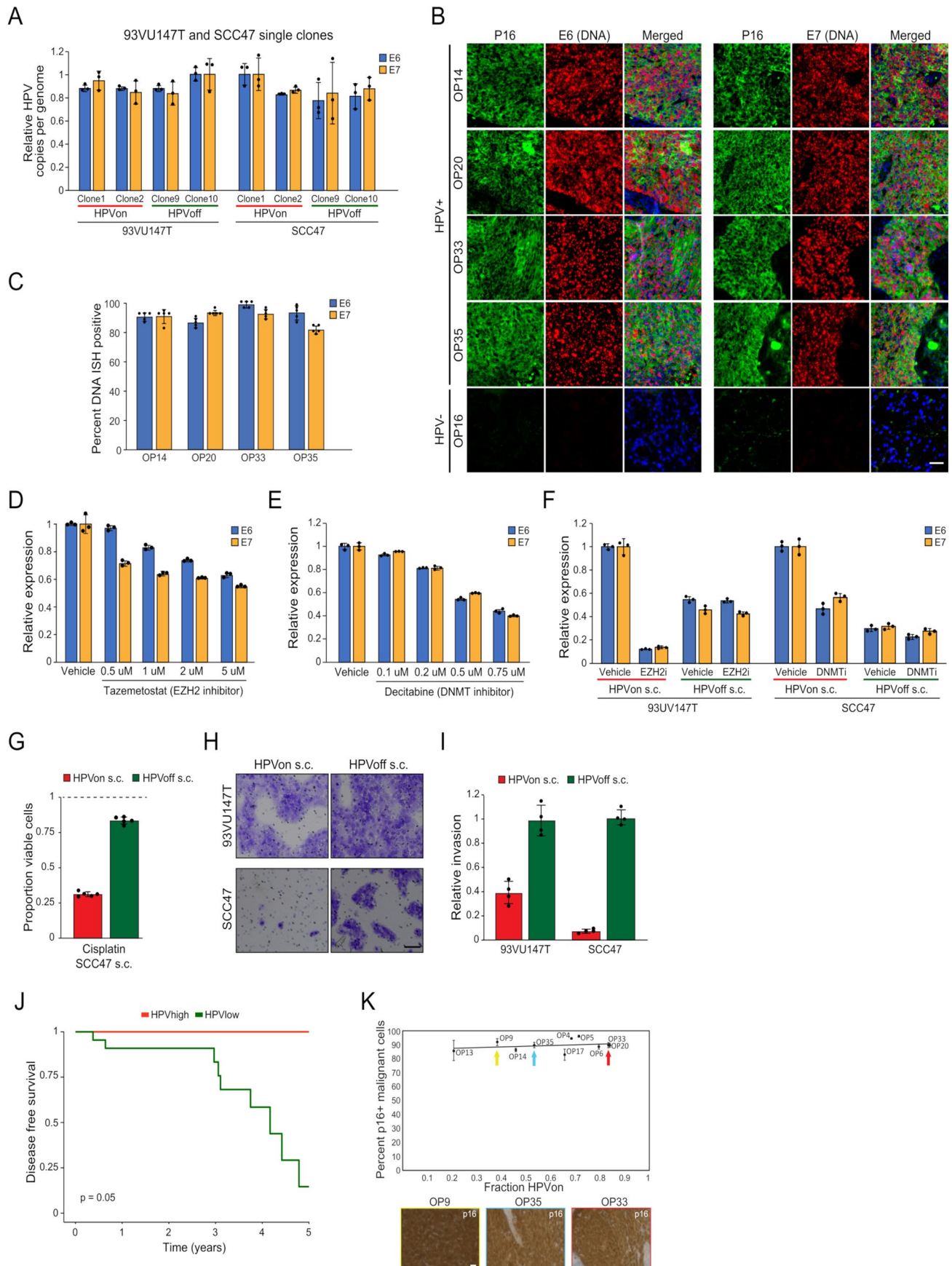
Extended Data Fig. 6 | Characteristics of *HPVoff* cells, related to Fig. 4. (a) Percentage of cells positive for *E6* or *E7* in RNA ISH analyses ($n = 4$ tumors, shown are mean and standard error across nine regions per tumor). Percentage of *HPVon* cells by scRNA-seq (bottom) correlates with RNA ISH values ($P < 0.01$, ANOVA). (b) IHC of representative HPV-positive (OP5, OP6, OP33, and OP35) and HPV-negative (OP19) tumors and normal tonsil stained for malignant-cell specific marker p16 (pink) and viral E6 protein (brown). Similar results were obtained in three independent experiments. White arrowheads denote p16 positivity without E6 expression. Scale bars: Low magnification = 10 mm (tonsil, OP5, OP6), 5 mm (OP19), 7.5 mm (OP35); intermediate magnification = 1000 μm ; highest magnification = 250 μm . (c) Enriched MSigDB Hallmark gene-sets among genes significantly overexpressed in *HPVon* versus *HPVoff* cells. X-axis: fraction of significantly upregulated genes in the gene set. (d) Differential expression of all analyzed genes between HPV-related classes of malignant cells. X-axis: difference between *HPVon* and *HPVneg* cells; Y-axis: difference between *HPVon*

and *HPVoff* cells, averaged across all HPV-positive patients. Genes are colored by their assignment to meta-program (right legend). CDKN2A (p16, highlighted in red) was not significantly different between *HPVon* versus *HPVoff* cells, but was the most overexpressed gene in *HPVon* cells compared to *HPVneg* cells. (e) For three meta-programs (panels), cells were divided into 10 bins of equal size, ranked by average expression from low (*left*) to high (*right*). Y-axis: mean ratio of cells belonging to an HPV subset versus the expected number assuming random distribution across bins. Error bars reflect SEM based on 100 re-sampling runs ($n = 5$ patients for *HPVneg*, $n = 11$ patients for *HPVon* and *HPVoff*). P-values are based on chi-square test. (f) Fractions of cycling cells, EpiSen-high cells and *HPVon* cells across genetic subclones. Subclones with a high fraction of *HPVon* cells tend to also have higher proliferation ($p < 0.05$ for correlations in OP13, OP33 and OP35). (g) G1/S (X-axis) and G2/M (Y-axis) scores of all malignant cells, colored by the percentage of cycling cells among their neighbors (20 closest cells in this plot).



Extended Data Fig. 7 | Regulation and function of HPVoff cells, related to Fig. 5. (a) HPV expression and G1/S gene expression across cervical squamous cell carcinoma TCGA samples. Shown are residuals after regression (Supplementary Table 3). (b) Variability in HPV expression between cell lines. Dot size and color represent fraction of cells with at least one read and fraction of HPV reads that reflect the corresponding gene, respectively. (c) Cells were divided into 5 bins by average G1/S expression from low (left) to high (right). Y-axis: mean ratio of cells in an HPV subset versus expected number assuming random distribution. Error bars are SEM by 100 resampling runs. P-value based on chi-square test. (d) Immunocytochemistry of 93VU147T cells probed with Ki67 (red), p16 (green), and DAPI (blue). Scale bar = 100 μm . (e) Percentage of Ki67 positive cells among p16 positive and negative cells. 50 cells were counted across four fields ($n = 4$). (f) Relative expression of E6 and E7 in non-target, control (NT) compared to E6 or E7 CRISPRi knockdown (KD) 93VU147T (left) or SCC47 (right) lines

($n = 3$; $P < 0.0001$, t-test). (g) Relative expression of p16 in same lines as in (f). Data are presented as mean \pm SEM. There was no change in p16 upon E6 or E7 knockdown ($n = 3$). (h) Relative expression of E6 and E7 among HPVon and HPVoff single clones derived from 93VU147T (left) and SCC47 (right) after three weeks of culture and numerous passages. HPVon and HPVoff clones maintained relatively high and low expression states ($n = 3$; $P < 0.005$, t-test). (i) Relative expression of p16 in same clones as in (h). (j) Proportion of cycling cells in HPVon and HPVoff single clones in 93VU147T (left) and SCC47 (right) by flow cytometry ($n = 3$; $P < 0.05$, t-test). (k) Relative proliferation of HPVon single clones from 93VU147T (left) and SCC47 (right) cultured under normal growth conditions (+FBS) or serum starvation (-FBS) for 48 hours. Proliferation was reduced with serum starvation ($n = 5$; $P < 0.001$, t-test). (l) Relative expression of E6 and E7 in HPVon single clones in 93VU147T (left) and SCC47 (right) under normal growth conditions (+FBS) or serum starvation (-FBS) for 48 hours ($n = 3$).



Extended Data Fig. 8 | See next page for caption.

Extended Data Fig. 8 | Functional impact of HPVoff cells and p16 expression, related to Fig. 5. (a) HPV copies per genome of *E6* and *E7* (normalized to albumin) for *HPVon* and *HPVoff* single clones from 93VU147T (left) and SCC47 (right). (b) DNA ISH (DNAScope) of representative HPV-positive (OP14, OP20, OP33, and OP35) and HPV-negative (OP16) tumors for viral *E6* (left) and *E7* (right) DNA (red) with immunofluorescence co-staining for regions of tumor as marked by p16 protein (green) and nuclei by DAPI (blue). HPV-positive tumors display p16 positive malignant cells with homogenous *E6* and *E7* DNA signal. HPV-negative tumors do not have signal for p16 protein or *E6* or *E7* DNA. Scale bar = 1000 μ m. (c) Percentage of cells positive for *E6* or *E7* DNA among p16 positive malignant cells in DNA ISH analyses ($n = 4$ tumors, five areas per tumor). Nearly all p16 positive malignant cells demonstrated *E6* or *E7* DNA signal. (d) Relative expression of *E6* and *E7* in 93VU147T cells treated with vehicle or tazemetostat ($n = 3$). All doses did not significantly affect cell viability. (e) Relative expression of *E6* and *E7* in SCC47 cells treated with vehicle or escalating concentrations of decitabine ($n = 3$). All doses did not significantly affect cell viability. (f) Relative expression of *E6* and *E7* in *HPVon* and *HPVoff* single clones from 93VU147T (left)

and SCC47 (right) treated with tazemetostat, decitabine, or vehicle. *HPVon* clones show reduction in *E6* and *E7* expression upon tazemetostat or decitabine treatment compared to *HPVoff* clones ($n = 3$; $P < 0.00001$, t-test). (g) Proportion of viable cells after treatment of SCC47 *HPVon* and *HPVoff* single cell clones with cisplatin, relative to cells treated with vehicle (dashed line). *HPVon* clones were more susceptible to cisplatin compared to *HPVoff* clones ($n = 5$; $P < 0.00001$, t-test). (h) Invasion of *HPVon* and *HPVoff* single clones from 93VU147T (top) and SCC47 (bottom). Scale bar = 100 μ m. (i) Relative invasion of *HPVon* and *HPVoff* single clones from 93VU147T (left) and SCC47 (right) cells. *HPVoff* cells were more invasive than *HPVon* ($n = 4$; $P < 0.05$, t-test). (j) Improved disease-free survival in HPVhigh compared to HPVlow samples, among TCGA p16+ oropharyngeal samples ($n = 28$; $P = 0.05$). (k) Top: percentage of p16 positive malignant cells (by IHC) and proportion of *HPVon* cells (by scRNA-seq). Bottom: p16 staining from tumors with low (OP9), intermediate (OP35) and high (OP20) proportions of *HPVon* cells (bottom). No correlation between *HPVon* proportion and percentage of p16 positive cells ($n = 10$ tumors). Scale bar = 100 μ m.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- ☐ ☒ The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- ☐ ☒ A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☐ ☒ The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- ☐ ☒ A description of all covariates tested
- ☐ ☒ A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- ☐ ☒ A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- ☐ ☒ For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- ☒ ☐ For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- ☒ ☐ For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- ☒ ☐ Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

NovaSeq S4 (200 cycles: 50-10-16-150)
Sequencing Analysis Viewer v2.4.7
Cellranger v4.0

Data analysis

Data analysis was performed in R v 4.1.0. Packages used, what they were used for, and versions, follow below:
Package Description (specific use in this project) Version
caTools Calculating moving average 1.18.2
circlize Creating colour palettes 0.4.13
class Classifying cells by kNN 7.3-19
clusterProfiler Enrichment analysis 4.0.2
ComplexHeatmap Plotting heatmaps 2.8.0
dplyr Data handling 1.0.7
FNN Create kNN graph 1.1.3
ggplot2 Plotting 3.3.5
ggrepel Separating text labels in plots 0.9.1
gtools Random permutation 3.9.2
igraph Louvain clustering 1.2.6
Matrix Working with sparse matrices 1.3-4
Matrix.utils Working with sparse matrices 0.9.8
msigdb Enrichment analysis 7.4.1
NMF Performing NMF 0.23.0
parallel Parallelising computation 4.1.0
reshape2 Data handling 1.4.4
scDbFinder Doublet detection 1.7.1

scds Doublet detection 1.8.0
 scan Doublet detection 1.22.1
 SingleCellExperiment Data handling 1.14.1
 stringdist Similarity between character strings 0.9.6.3
 uwot Creating UMAP 0.1.10

Other Software used:

Vartrix Calling SNVs in single-cell data 1.1.14
 bwa-mem Genome alignment 0.7.15
 GATK Genome Analysis Toolkit 4.1.7.0
 Mutect2 Somatic variant calling 4.1.7.0
 HaplotypeCaller Germline variant calling 4.1.7.0
 FlowJo Analysing flow cytometry data 9.0

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

All single cell RNA-seq data produced by this study is available through the Gene Expression Omnibus with GEO accession GSE182227. TCGA bulk RNAseq and clinical data for head and neck and cervical cancer is available through the Broad Genome Data Analysis Center Firehouse (<https://gdac.broadinstitute.org/>). Single-cell datasets reanalyzed to compare proliferation rates are available through the Gene Expression Omnibus with accession numbers GSE150430 (nasopharyngeal carcinoma), GSE131907 (lung carcinoma), GSE132465, GSE132257, GSE144735 (CRC), GSE125449 (HCC), through CNCB with accession GSA: CRA001160 and through EMBL-EBI ArrayExpress with accession numbers E-MTAB-8107, E-MTAB-6149 and E-MTAB-6653 (breast, lung and ovarian cancer). Cell line data used for validation analysis is available through GEO with accession number GSE157220. The NSCLC dataset used to validate finding malignant cells in normal samples is deposited as an NCBI BioProject with accession number PRJNA591860. Raw data for figures 1b-d, 2b-c, 2f, 3b, 3e, 4a, 4c-g, 5a, 5c and extended data figures 1d-e, 2b-c, 2e-f, 3b, 6c-f, 7b-c and 8j are made available as source data.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	16 patients were chosen to provide a reasonable representation of both HPV+ (n=11) and HPV- (n=5) tumors. These sample sizes are similar to other single cell sequencing analyses of head and neck squamous cell carcinoma previously performed in the field. Because the vast majority of our analyses focused on HPV+ oropharyngeal tumors, our sample size emphasized these tumors.
Data exclusions	Cells that did not pass QC, were classified as doublets or had an unresolved cell type were excluded from further analysis, as described in the Methods section
Replication	This was an exploratory study, where replication is not applicable. The main findings are validated through external datasets and cell line experiments as described in the manuscript.
Randomization	Impact of treatment on patients was not studied, and all samples were from treatment-naïve patients. Thus, randomization is not relevant.
Blinding	Patient identity was not relevant to the study, and tissue samples were used in a de-identified fashion. Blinding was not relevant to the study.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input type="checkbox"/>	<input checked="" type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Antibodies

Antibodies used	Ki-67 1:500 dilution (D2H0 rabbit mAb, Cell Signaling), E6 1:100 dilution (mouse anti-virus, clone C1P5, Invitrogen), E7 1:100 dilution (mouse anti-virus, clone TVG701Y, Invitrogen), HPV TYPE 16/18 E6 Mouse Monoclonal antibody (1:50 dilution, Thermofisher, cat# MA1-46057), P16-INK4A polyclonal antibody (1:75 dilution, Thermofisher, cat# 10883-1-AP)
Validation	Antibodies were validated for immunohistochemistry using a no primary control to confirm that signal was specific to the primary antibody, with HPV- tumors (which lack E6, E7, and p16) used as a negative control

Eukaryotic cell lines

Policy information about [cell lines](#)

Cell line source(s)	SCC47, 93VU147T, SCC90 and 293T cell lines were obtained from Dr. James Rocco (co-author)
Authentication	All cell lines were validated by STR analysis and confirmed to be accurate prior to being utilized.
Mycoplasma contamination	All cell lines were tested for mycoplasma.
Commonly misidentified lines (See ICLAC register)	No commonly misidentified cell lines were used in the study.

Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics	All population characteristics are found in extended data tables 1-2
Recruitment	Patients were preoperatively asked to provide tissue samples for the study. No compensation was involved. Since our study uses each patient as its own internal control, comparing cells with or without HPV expression, any potential recruiting biases are highly unlikely to impact the results.
Ethics oversight	Washington University School of Medicine Institutional Review Board

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Flow Cytometry

Plots

Confirm that:

- ☒ The axis labels state the marker and fluorochrome used (e.g. CD4-FITC).
- ☒ The axis scales are clearly visible. Include numbers along axes only for bottom left plot of group (a 'group' is an analysis of identical markers).
- ☒ All plots are contour plots with outliers or pseudocolor plots.
- ☒ A numerical value for number of cells or percentage (with statistics) is provided.

Methodology

Sample preparation	As described in methods, for flow cytometric cell cycle analyses, cells were fixed with 70% ice cold ethanol and stained with propidium iodide (30 µg/ml of PI (Sigma) with 200 µg/ml of RNase (Sigma) in 0.1% of Triton-X-100 (Sigma) in PBS) for one hour at room temperature. Cell cycle analysis was completed with at least 10,000 cells using CytoFLEX Flow Cytometer and data were analyzed using FlowJo v9.0 software.
Instrument	CytoFlex Flow Cytometer

Software

FlowJo v9.0

Cell population abundance

All cells were analyzed for their cell cycle phase after exclusion of dead cells and debris.

Gating strategy

Dead cells and debris were used to gate cells based on FSC and SCC using standard approaches with the help of the Siteman Flow Cytometry Core.

☒ Tick this box to confirm that a figure exemplifying the gating strategy is provided in the Supplementary Information.