



GENOME RESEARCH

Nanopore sequencing identifies high-frequency somatic structural variations in laryngeal squamous cell carcinoma genomes

Xuyan Liu, Lin Xia, Yixin Qiao, et al.

Genome Res. published online April 8, 2026

Access the most recent version at doi:[10.1101/gr.281046.125](https://doi.org/10.1101/gr.281046.125)

P<P Published online April 8, 2026 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Open Access Freely available online through the *Genome Research* Open Access option.

Creative Commons License This manuscript is Open Access. This article, published in *Genome Research*, is available under a Creative Commons License (Attribution 4.0 International license), as described at <http://creativecommons.org/licenses/by/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

CRISPR and RNAi Genetic Screening.
Your new superpower.

LEARN MORE

CELLECTA

The advertisement features a woman in a red and white superhero costume with a red mask. To her right is a green molecular structure icon. The text "LEARN MORE" is enclosed in a white box, and the Cellecta logo is at the bottom right.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Nanopore sequencing identifies high-frequency somatic structural**
2 **variations in laryngeal squamous cell carcinoma genomes**

3 Xuyan Liu^{1#}, Lin Xia^{1#}, Yixin Qiao^{2#}, Yang Li¹, Yan Huang¹, Bingyan Yue², Xi Liang², Xin
4 Yang², Honghui Zhang², Jiaxun Zhang¹, Xiao Chen¹, Dan Xie^{1*}, Jifeng Liu^{2,3*}

5 1. Laboratory of Omics Technology and Bioinformatics, Frontiers Science Center for
6 Disease-related Molecular Network, State Key Laboratory of Biotherapy, West China
7 Hospital, Sichuan University, Chengdu, Sichuan, 610041, China

8 2. Department of Otolaryngology-Head & Neck Surgery, West China Hospital, Sichuan
9 University, Chengdu, Sichuan, 610041, China

10 3. Deep Underground Space Medical Center, West China Hospital, Sichuan University,
11 Chengdu, Chengdu, Sichuan, 610041, China

12

13 # These authors contributed equally

14 * Corresponding author

15 E-mail: danxie@scu.edu.cn (DX), liujifeng777@wchscu.cn (JFL)

16

17 **Running title**

18 Somatic variations in LSCC by long-read sequencing

19

20 Abstract

21 Laryngeal squamous cell carcinoma (LSCC) is an aggressive cancer with poor quality of life.
22 Understanding the somatic mutations in its genome is crucial for elucidating its occurrence
23 and progression. Although somatic structural variations (SVs) have been documented in LSCC,
24 conventional short-read sequencing lacks the sensitivity to effectively detect high-frequency
25 SVs shared across multiple samples—variants that play a crucial role in tumorigenesis. Here,
26 we present SomaGauss-SV, a somatic SVs detection workflow leveraging nanopore long-read
27 sequencing data. Benchmarking against five paired tumor cell line datasets shows
28 SomaGauss-SV consistently achieves a balanced high precision and recall. SomaGauss-SV
29 applied to 15 paired LSCC tumor-blood samples uncovers a comprehensive SVs landscape
30 and reveals a significant positive correlation between somatic deletion burden and smoking
31 intensity. Furthermore, a high-frequency somatic simple repeat expansion is identified in
32 28/39 (71.79%) of LSCC patients, which upregulates the expression of genes *TP53BP2* and
33 *FBXO28* through spatial proximity. These findings underscore the potential of long-read
34 sequencing and SomaGauss-SV for uncovering recurrent somatic SVs in LSCC, providing
35 valuable resources for biomarker discovery.

36

37 Introduction

38 Laryngeal squamous cell carcinoma (LSCC) is the second most common subtype of head and
39 neck squamous cell carcinoma (HNSCC), characterized by marked gender differences and
40 regional variability in both incidence and mortality rates. The incidence rate in males is
41 approximately eight times that of females(Han et al. 2024) LSCC pathogenesis is closely
42 related to modifiable risk factors, such as tobacco exposure and alcohol consumption(Bray et
43 al. 2024; Sung et al. 2021; Chen et al. 2024; Johnson et al. 2020). However, the molecular
44 mechanisms by which these exogenous factors induce genomic instability remain unclear.
45 Early symptoms of LSCC are often overlooked, resulting in more than 50% of patients being
46 diagnosed at advanced stages, particularly those with supraglottic tumors(Steuer et al. 2017;
47 Nocini et al. 2020). This delay in diagnosis typically requires aggressive treatments that
48 severely impair patients' respiratory and vocal functions, thereby significantly diminishing
49 their quality of life(R et al. 2019). Despite advances in treatment regimens, combining surgery
50 with radiotherapy and chemotherapy, the five-year overall survival rate has shown limit

51 improvement(Cavaliere et al. 2021; Wang et al. 2020), highlighting the urgent need for
52 molecular biomarkers specific to LSCC. Compared to other HNSCC subtypes, LSCC has the
53 unique biological characteristics and epidemiological features(Cavaliere et al. 2021; Pan et al.
54 2025), making specialized research essential to address its unique diagnostic challenges and
55 the critical need for organ function preservation.

56
57 In recent years, pan-cancer genomic studies have highlighted the pivotal role of structural
58 variants in the occurrence and progression of cancer(Cosenza et al. 2022). Structural variants
59 (SVs), which encompass large genomic rearrangements exceeding 50 bp, include deletions
60 (DEL), insertions (INS), duplications (DUP), inversions (INV), and translocations (TRA) (van
61 Belzen et al. 2021). These variants drive oncogenesis through mechanisms such as gene
62 disruption and fusion, copy number alterations, and three-dimensional genome
63 reorganization(Beroukhim et al. 2016). Large-scale cancer cohort studies have shown that
64 somatic SVs are the most prevalent class of driver mutations in cancer, vastly outnumbering
65 single nucleotide variants (SNVs) and small insertions/deletions (indels)(Cosenza et al. 2022).
66 Given their abundance and diverse roles in cancer genomes, a comprehensive understanding
67 of somatic SV's scope and mechanisms is essential for uncovering key cancer mutations in
68 patient tumors.

69
70 However, despite the widespread clinical application of somatic SV detection methods based
71 on second-generation short-read sequencing, these approaches have inherent limitations in
72 sensitivity for smaller structural variants (≤ 10 kb), and in resolving repetitive sequence
73 variants (ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020; Choo et al.
74 2023). The advent of nanopore long-read sequencing technology offers a promising solution to
75 overcome this bottleneck. With read length exceeding 10 kb, nanopore sequencing can span
76 SVs breakpoint regions in their entirety, enabling comprehensive mapping of somatic SV
77 landscapes. Recently, several somatic SVs algorithms, such as nanomonsv(Shiraishi et al.
78 2023), Severus(Keskus et al. 2025), and SVision-pro(Sahlin et al. 2023), have been developed
79 to harness this potential. However, existing tools still face challenge in correcting alignment
80 biases in repetitive regions, distinguishing mixed signals of multi-scale INs and DELs, and
81 their practical application in clinical samples of sufficient scale remains limited(Sahlin et al.
82 2023). Therefore, developing detection workflows suitable for different research scenarios is a

83 key step for advancing the clinical interpretation of SVs.

84
85 Here, we present SomaGauss-SV, a nanopore-based somatic SV caller developed to address
86 this limitation. We rigorously benchmark its performance in paired tumour–normal cell lines
87 using complementary, orthogonal validation strategies to establish robustness before clinical
88 deployment. We then apply this framework to delineate the somatic SV landscape of primary
89 LSCC and interrogate potential biological and etiological associations. Ultimately, this study
90 aims to deliver a high-performance workflow for somatic SV discovery, elucidate recurrent
91 SV-linked mechanisms in LSCC, and generate a comprehensive, clinically annotated SV
92 resource to enable future mechanistic and translational investigations.

93
94 **Results**
95
96 **SomaGauss-SV demonstrates excellent performance in detecting somatic structural**
97 **variations**

98
99 Current long-read sequencing algorithms for detecting SVs often exhibit reduced accuracy in
l00 repetitive genomic regions. This limitation arises from inherent biases in traditional seed-chain
l01 alignment methods(Sahlin et al. 2023) (Supplemental Fig. S1). The lack of systematic
l02 alignment error correction in existing SVs detection tools further constrains their utility for
l03 somatic SVs calling. To address these challenges, we present SomaGauss-SV, a novel
l04 workflow (Fig. 1A) that integrates a mixed Gaussian distribution model with alignment error
l05 correction to improve the accuracy of somatic SVs detection. Its core innovation lies in
l06 applying tailored Gaussian modeling for different SV types coupled with stringent
l07 population-based filtering to distinguish somatic events (see Methods for full details).

l08
l09 To systematically assess the performance of SomaGauss-SV in somatic SVs detection,
l10 we utilized a public resource that provides five tumor-normal cell line pairs
l11 (HCC1395/HCC1395BL, HCC1937/HCC1937BL, HCC1954/HCC1954BL, H1437/H1437BL,
l12 H2009/H2009BL,Supplementary Table 1) with matched sequencing data and a curated SV
l13 ensemble as the ground truth (Keskus et al. 2025). We employed two independent evaluation
l14 strategies, Truvari (English et al. 2022) and minda (Keskus et al. 2025), to compare
l15 SomaGauss-SV against established somatic SV callers: nanomonsv (Shiraishi et al. 2023),

l16 SVision-pro (Wang et al. 2025), and Severus (Keskus et al. 2025). For a fair comparison in this
l17 controlled setting, we omitted the panel-of-normal (PON) filtration step.

l18
l19 SomaGauss-SV consistently achieved the highest F1 scores across all five cell line pairs using
l20 the Truvari evaluation strategy, achieving a mean F1-score of 62.8% ($\pm 19.2\%$ SD),
l21 significantly outperforming Severus (51.5% $\pm 18.0\%$), SVision-pro (42.0% $\pm 16.3\%$),
l22 and nanomonsv (28.9% $\pm 11.1\%$) (Fig. 1B). Notably, SomaGauss-SV achieved peak
l23 performance (84.9% F1-score) in the H2009 cell line, which is characterized by
l24 insertion-dominant SVs profiles (mean comparator performance: 57.3%,
l25 Supplemental Fig. S2A), highlighting its particular strength in resolving this
l26 challenging variant class. In contrast, the HCC1395 cell line, with its high SVs burden,
l27 showed narrower performance differentials between SomaGauss-SV (74.9%) and
l28 Severus (72.0%), both substantially outperforming the other tools (mean: 42.5%). The
l29 genome-doubled HCC1954 cell line exhibited universally lower F1-scores across all
l30 platforms (mean: 28.2% vs. 50.8% in non-polyploid lines), suggesting that
l31 polyploidization introduces systematic challenges in SVs detection. A critical analysis
l32 revealed the precision-recall trade-offs of each tool (Fig. 1C): nanomonsv exhibited
l33 higher precision (59.1% vs. competitors' mean precision 50.7%), while Severus and
l34 SVision-pro prioritized recall (64.3% vs. competitors' mean recall 40.6%).
l35 SomaGauss-SV uniquely balanced both metrics (mean precision: 64.0%, mean recall:
l36 61.9%). Considering practical application, we also evaluated the computational space
l37 and time complexity. SomaGauss-SV exhibited the smallest Maximum Resident Set
l38 Size, and its Elapsed Real Time was significantly less than that of Severus, despite
l39 their similar F1-scores (Supplemental Fig. S2B).

l40
l41 The inter-cell line performance variation strongly correlated with the composition of
l42 somatic SVs types. To explore this relationship further, we conducted stratified
l43 analyses by SVs types using Truvari with type-specific ensemble settings
l44 (Supplemental Fig. S2C). SomaGauss-SV consistently outperformed other tools across
l45 all SVs types. Notably, in INS detection, SomaGauss-SV achieved an F1-score of
l46 59.3%, outperforming comparator tools by 13.9–93.1 % (mean F1-score of
l47 comparators: 25.5%). For DUP, SomaGauss-SV exhibited performance comparable to

l48 Severus (67.7% vs. 69.1%). While all tools showed relatively higher efficacy in DEL
l49 detection (mean F1-scores: 45.4–70.2%), both breakend (BND) and inversion (INV)
l50 detection proved challenging across platforms (mean F1 < 5.7%). This limitation was
l51 most pronounced in the BND-rich HCC1954 dataset, where BNDs constituted 64.2% of SVs
l52 (Supplemental Fig. S2C). To ensure a robust and unbiased evaluation, we conducted a
l53 complementary assessment using *minda* (see Supplementary Table 2 for full details),
l54 which matches variants based on breakpoint coordinates and length, independent of
l55 SV type annotation. Re-evaluation with *minda* yielded two key insights. First, under
l56 this more lenient and type-agnostic matching scheme, F1-scores for all tools increased,
l57 with the most substantial improvements (15.56%–48.47%) observed in the BND-rich
l58 HCC1954 dataset (Supplemental Fig. S2D), highlighting the inherent difficulty in
l59 computationally distinguishing BND. Second, and more importantly,
l60 SomaGauss-SV’s advantage in precision became even more pronounced. It achieved
l61 the highest average precision (79.20%) across four cell lines (Supplemental Fig. S2E),
l62 demonstrating its ability to generate a highly reliable, low false-positive call set—a
l63 critical feature for downstream biological discovery. In light of the inherent difficulty
l64 in computationally detecting complex SVs like INV and BND, our final analysis
l65 pipeline incorporates a mandatory manual curation step for these variant types in
l66 tumor–normal pairs to ensure accuracy (see Methods). Beyond type-specific
l67 performance, we analyzed the concordance between different SV detection tools using
l68 the HCC1395/HCC1395BL cell line pair (Fig. 1D). SomaGauss-SV demonstrated
l69 superior consensus validity (92.2% overlap) and the highest true positive rate (23.94%)
l70 among its unique calls (Fig. 1E). Although consensus calling improved true positive
l71 rates, it incurred substantial variant loss (52.3%) and systematic class biases (e.g., the
l72 proportion of DEL detection increased by 29.0%; Supplemental Fig. S2F),
l73 underscoring its limitations for comprehensive SV profiling. Collectively, these
l74 findings position SomaGauss-SV as the optimal solution for our LSCC study,
l75 balancing precision and recall while preserving both high-confidence consensus
l76 variants and biologically relevant unique calls - critical advantages for downstream
l77 clinical interpretation.

l78

l79

!80 **Pathological information and sequencing quality data display of ONT WGS samples**

!81
!82 To interrogate the Somatic SVs landscape of LSCC, we included tumor samples,
!83 comprising with matched blood samples from 15 treatment-naive with primary,
!84 HPV-negative LSCC (Supplementary Table 3, Fig. 2A). These tumor samples
!85 represented a broad spectrum of clinical stages and anatomical subtypes (Fig. 2B and
!86 2C). All patients were male, with a mean age of 63 years (median: 60). Smoking
!87 exposure was assessed using the smoking index, calculated as the number of cigarettes
!88 per day multiplied by years of smoking(Zhou et al. 2023). 14/15 patients had a history
!89 of smoking, with an average smoking index of 648 (median: 610) (Supplementary
!90 Table 3, Fig. 2D). These cohort characteristics align with previous studies,
!91 demonstrating that LSCC predominantly affects middle-aged and elderly males and is
!92 strongly associated with smoking(Johnson et al. 2020; Lechien et al. 2022), reflecting the
!93 typical demographic profile of LSCC patients.

!94
!95 Whole-genome sequencing of paired tumor and blood samples was performed using
!96 the PromethION platform (R9: 9 samples, R10: 6 samples; Supplementary Table 3),
!97 generating 2.10 Tb of data. After quality control (sequencing quality > 7, read length >
!98 1 kb), tumor samples yielded an average N50 of 28 kb (range: 20 kb to 36 kb) and an
!99 average of 94 billion bases per sample, equating to a mean sequencing depth of 30×
!00 (range: 28× to 33×, Fig. 2E). Matched blood samples generated an average N50 of 28
!01 kb (range: 19 kb to 40 kb) and an average of 49 billion bases per sample, resulting in a
!02 mean sequencing depth of 15× (range: 12× to 18×, Fig. 2E). Alignment of the clean
!03 long-read sequencing data to the hg38 reference genome using minimap2 yielded
!04 average alignment rate of 99.95% for tumor samples and 99.90% for blood samples
!05 (Supplementary Table 4, Fig. 2F).

!06 !07 **Discovery of the positive relation between somatic DEL burden and tobacco exposure**

!08
!09 Using SomaGauss-SV, 8128 somatic SVs were identified across the 15 paired samples,
!10 with an average of 541 somatic SVs per tumor samples. Somatic insertions (INS) were
!11 the most common variant, comprising 54.7% of the total, followed by deletions (DEL,
!12 19.7%), rearrangements (BND, 18.0%), inversions (INV, 4.0%), and duplications

!13 (DUP, 3.6%) (Supplemental Fig. S3A). The distribution of SV types in individual
!14 samples was consistent with this overall pattern (Fig. 3A). Notably, samples with
!15 higher smoking indices displayed a significantly greater somatic SV burden (Pearson
!16 correlation coefficient, $p < 0.05$, $r = 0.54$; Supplemental Fig. S3B). DEL contributed
!17 most prominently to this trend, exhibiting the strongest correlation with smoking
!18 index among all SV types (Pearson correlation coefficient, $p < 0.05$, $r = 0.67$; Fig. 3B).
!19 To dissect this relationship more rigorously, we performed negative binomial
!20 regression, treating the burden of specific SV types as the dependent variable while
!21 controlling for TNM stage, age, and sequencing depth. This multivariate analysis
!22 identified smoking index as a highly significant, independent positive predictor for
!23 DEL burden specifically ($p = 0.000596$; Supplemental Fig. S3C). To strengthen the
!24 generalizability of this finding, we validated the association in an independent head
!25 and neck squamous cell carcinoma cohort from the ICGC ($n = 44$, 19 of whom had no
!26 smoking history). After controlling for TNM stage, age, and sex via negative binomial
!27 regression, smoking index remained significantly associated with DEL burden ($p =$
!28 0.0498 ; Supplemental Fig. S3D). This finding similar with previous study indicating
!29 that smoking increases the somatic mutation burden in bronchial epithelial cells and
!30 LSCC (Yoshida et al. 2020; Degawa et al. 1994), suggested that tobacco exposure may
!31 promote somatic DEL by introducing DNA double-strand breaks (Yoshida et al. 2020;
!32 Degawa et al. 1994; Aw et al. 2021).

!33
!34 To investigate the functional impact of somatic deletions in smoking-related LSCC,
!35 we focused on the 42 somatic DEL regions (after merging adjacent windows, as
!36 described in Methods: Sample Enrichment Analysis) that were shared by three or
!37 more samples from smokers. Using GREAT (McLean et al. 2010), these loci were
!38 mapped to 49 putative distal regulatory target genes. Pathway enrichment analysis
!39 revealed that six of the top ten enriched pathways ($q < 0.05$) were immune-related,
!40 including Dectin-2 family (Graham and Brown 2009), Translocation of ZAP-70 to
!41 Immunological Synapse (Blanchard et al. 2002), and Complement
!42 Cascade (Afshar-Kharghan 2017) (Fig. 3C). GO analysis further supported this
!43 immune-focused pattern, with nine of the top ten biological processes ($q < 0.05$)
!44 related to lymphocyte activation and cytokine responses (Supplemental Fig. S3E).

!45 These findings supported the role of smoking in tumorigenesis through
!46 immune-inflammatory responses(Aw et al. 2021), suggested that somatic DEL
!47 frequently carried by high-smoking LSCC patients may regulate the
!48 smoking-associated immune microenvironment.

!49
!50 Among the immune-related somatic DEL regions, a particular a somatic DEL, shared
!51 by four smoking-related LSCC samples, may affect the expression of *HLA-DRA* and
!52 *HLA-DRB5*, which are key immune modulators involved in tumor immune evasion.
!53 This high-frequency DEL is located near a CTCF binding site, identified in the
!54 MCF-7 and A673 cell lines, 265 bp apart (Fig. 3D). The CTCF site is predicted to
!55 regulate *HLA-DRA* and *HLA-DRB5* (The ENCODE Project Consortium 2012). We
!56 propose that the somatic DEL disrupts CTCF binding, leading to altered chromatin
!57 architecture and dysregulated expression of *HLA-DRA* and *HLA-DRB5*. Additionally,
!58 we identified a somatic DEL in the *MUC3A* gene, which encodes mucin and has been
!59 associated with smoking-induced changes in chronic obstructive pulmonary disease
!60 (COPD)(Merikallio et al. 2023). *MUC3A* expression is reported to be higher in smokers
!61 than in non-smokers due to tobacco exposure(Gendler and Spicer 1995). The somatic
!62 DEL, located at exon 2 of *MUC3A* gene (DEL length: 1120 bp to 1123 bp), contains a
!63 simple tandem repeat sequence (motif: ACC; Fig. 3E). PCR analysis of cDNA from
!64 clinical samples demonstrated that this genomic deletion results in a deletion within
!65 the *MUC3A* transcript, resulting in a shorter band (~346-349 bp) compared to the
!66 expected full-length product (1469 bp) (Supplemental Fig. S3F). This exon encodes a
!67 peptide rich in threonine (Thr) and serine (Ser), which are essential for glycosylation
!68 and for mucin 3A's role in protecting the epithelial barrier from environmental
!69 damage (Bhatia et al. 2019; Kufe 2009). Therefore, the somatic DEL is likely to cause
!70 truncation of this critical peptide domain, potentially impairing mucin 3A function.

!71

!72 **Charateristics of somatic SVs**

!73

!74 We characterized somatic SVs in LSCC by analyzing SV size, sequence context, and
!75 breakpoint distribution. Our results highlight distinct length preferences across SV
!76 types (Fig. 4A). Somatic DEL and INS were characterized by shorter lengths and

177 narrower distributions. DEL (mean: 1493 bp, range: 50 bp to 30 kb), had a peak at
178 size of 70 bp, while INS (mean 1107 bp, range: 50 bp to 24 kb), exhibiting a tri-modal
179 distribution with peaks at 70 bp, 170 bp, and 300 bp. In contrast, DUP and INV
180 displayed larger, more dispersed length distributions (DUP mean: 39 kb, INV mean:
181 5.3 Mb). These patterns are similar with prior study on somatic SVs in HCC(Zeng et al.
182 2025).

183
184 To investigate the sequence determinants of somatic SVs, we analyzed the sequence
185 feature of SVs at characteristic size peaks. Given the relatively uniform length
186 distributions of DUP and INV (Fig. 4A and Supplemental Fig. S4A), we focused on
187 DEL and INS (Fig. 4B and 4C). At the 70 bp peak, both somatic DEL and INS were
188 significant enriched in simple repeat sequences. Notably, AAAG motifs predominated
189 in somatic DEL (30%, Supplemental Fig. S4B), whereas somatic INS exhibited a
190 strong AT-richness (58.6% AT, Supplemental Fig. S4C). The 170-bp INS peak
191 showed a notable association with centromeric regions. Specifically, 45.4% of these
192 170-bp INS harbored Satellite/CEN (monomers: ~170 bp (Aldrup-MacDonald et al.
193 2016), Pearson's Chi-squared test, p -value = 1.1×10^{-6}). Among these
194 satellite/CEN-containing INS, 55.6% were directly located within centromeres, a
195 significant enrichment compared to the genome-wide background of 12.3% (Pearson's
196 Chi-squared test, p -value < 0.001; Supplemental Fig. S4D). Even those not residing
197 inside centromeres clustered in close proximity, with 62.7% found within 1 Mb of
198 centromeric regions (Wilcoxon test, p -value < 0.001, Supplemental Fig. S4D). This
199 pericentromeric enrichment pattern was further confirmed by re-alignment of
200 Chromosome 20 reads to the complete T2T-CHM13 assembly (Wilcoxon test, p -value
201 < 0.001, Supplemental Fig. S4E). At the 300-bp peak, SINE/*Alu* elements, primarily
202 the *AluY* subclass, were predominant in INS, comprising 43.9% of the total (Pearson's
203 Chi-squared test, p -value = 3×10^{-10}) (Supplemental Fig. S4F). This observation may
204 be attributed to the *AluY* subfamily, the youngest and most active member of the *Alu*
205 family(Hormozdiari et al. 2011; Witherspoon et al. 2010), whose high activity has been
206 identified in both cancer and healthy populations(Wu et al. 2021; Zeng et al. 2025; Ma
207 and Pl 2002).

208

309 Somatic SVs breakpoints were non-uniformly distributed across the genome, with
310 significant enriched in repetitive elements (TE enrichment score: 0.85; repeat
311 enrichment score: 10.20; non-repetitive regions enrichment score: 0.82).
312 Hyper-repetitive regions (including satellite, simple repeats, low complexity regions
313 and SVA elements) showed substantial enrichment (mean enrichment score > 6.56,
314 Fisher's exact test, p -value < 0.05, Fig. 4D). Notably, somatic DEL and INS were
315 more frequently observed in centromeric satellite sequences (enrichment score > 6.4,
316 Fisher's exact test, p -value < 0.05), whereas larger SVs, such as INV and DUP, were
317 underrepresented in these regions (enrichment score < 0.54). This pattern may be due
318 to satellite sequences' susceptibility to replication fork slippage or mispairing,
319 combined with the compact chromatin structure of centromeric regions that limits
320 DNA accessibility and favors conserved repair mechanisms, reducing the occurrence
321 of large-scale SVs (Kashi and King 2006). Sequence annotation revealed that over 60%
322 of INS inserted in satellite/center, satellite, and simple repeat regions harbored
323 inserted sequences identical to flanking repeats (Supplemental Fig. S4G). A
324 representative case in simple repeats demonstrated a 12-fold amplification of the
325 reference TTCT motif in tumor samples (Supplemental Fig. S4H), indicating that
326 somatic INS preferentially propagate existing repetitive architectures through
327 replication- or repair-associated templated insertion.

328
329 To elucidate the regulatory consequences of somatic SVs, we performed an
330 enrichment analysis of SV breakpoints across functional genomic regions. Somatic
331 INS, DUP, DEL, and INV exhibited significant enrichment within proximal
332 transcriptional regulatory domains (Fig. 4E), particularly in promoter-proximal (2.5kb)
333 and termination regulatory regions (downstream 2.5 kb), suggested that these SVs
334 may alter gene expression through *cis*-regulatory interference (Cosenza et al. 2022).
335 Specifically, INS and BND breakpoints were broadly localized to *cis*-regulatory
336 elements, with the strongest enrichment observed at CTCF-binding sites (enrichment
337 score > 3.19, Fisher's exact test, p -value < 0.05, Supplemental Fig. S4I), suggested
338 potential disruption of topologically associating domain boundary. In contrast, DUP
339 and DEL breakpoints preferentially disrupted enhancers and promoters near
340 transcription/termination sites, suggesting they preferred modify the regulatory

relationships between elements and genes (enrichment score >1.40 , Fisher's exact test, p -value < 0.05 , Supplemental Fig. S4I).

343

344 **Identify extensive high-frequency somatic SVs regions in LSCC**

345

346 To detect recurrent somatic SVs across genome, we performed a non-overlapping
347 sliding window approach for sample carrier rate quantification. We identified 226
348 high-frequency SV regions (sample carrier rate $> 26\%$) exhibiting non-random
349 chromosomal distribution (Fig. 5, Circos plot (Krzywinski et al. 2009)), and significant
350 co-localization with known tumor-related genes (Tate et al. 2019) (Chi-squared test,
351 p -value < 0.05 , Supplementary Table S5). To assess robustness, we performed
352 orthogonal validation using an independent SV caller (Severus) on the same dataset.
353 Notably, two INS hotspots identified by SomaGauss-SV were concordantly confirmed
354 by Severus, underscoring their biological relevance beyond algorithmic artifacts
355 (Supplemental Fig. S5).

356

357 We further characterized two high-frequency somatic INS regions,
358 Chr7:121,603,263–121,603,506 (shared by 9 LSCC samples, Supplemental Fig. S6A)
359 and Chr15:29,945,811–29,950,000 (shared by 7 LSCC samples, Supplemental Fig.
360 S6B), which overlapped with recurrent repeat expansions (rREs) previously identified
361 in the TCGA pan-cancer cohort (Erwin et al. 2023). The former INS aligned with a rRE
362 specifically detected in lung squamous cell carcinoma, while the latter corresponds to
363 a rRE characteristic of prostate cancers. The inserted sequences of these somatic INS
364 matched the motifs described in previous study, suggested that these somatic INS
365 have the potential to serve as pan-cancer biomarker across multiple malignancies.
366 Moreover, we identified several high-frequency DEL that may influence nearby genes
367 by disrupting regulatory elements. For instance, the region
368 Chr11:117,260,000–117,262,041 harbored somatic DEL (235-353 bp) in six tumors
369 (Supplemental Fig. S6C). This region is located within an intron of *RNF214* and
370 overlapped with an enhancer element predicted by ENCODE (Je et al. 2020). This
371 enhancer is predicted to regulate *RNF214* expression in PC-3 prostate cancer cells (Je
372 et al. 2020). *RNF214* is an E3 ubiquitin ligase. Its overexpression promotes tumor cell

373 proliferation, migration, and invasion(Lin et al. 2024). Additionally, in
374 Chr12:54,206,733–54,210,000, somatic DEL (163-178 bp) in five samples overlapped
375 with an enhancer predicted to regulate *SMUG1* expression in Panc1 pancreatic cancer
376 cells(Je et al. 2020) (Supplemental Fig. S6D). *SMUG1* is involved in base excision
377 repair and genomic stability maintenance(Iakovlev et al. 2019). High *SMUG1*
378 expression correlates with poorer survival in the TCGA HNSCC cohort (Log-rank test,
379 p -value = 0.003) (Supplemental Fig. S6E).

380
381 We also identified eight somatic INV shared by more than four samples
382 (Supplemental Fig. S6F). The most frequent event occurred in region
383 Chr7:77,175,000–102,599,637 (9/15 samples). This interval encompasses 197
384 protein-coding genes, including 11 annotated proto-oncogenes. Pathway enrichment
385 analysis revealed significant enrichment for the Cytochrome P450 metabolic pathway
386 ($q < 0.05$), which includes genes within this region such as *CYP3A5* and *CYP3A7*.
387 Notably, CYP family enzymes can be induced by polycyclic aromatic hydrocarbons in
388 tobacco smoke, enhancing the metabolic activation of procarcinogens and thereby
389 promoting tumorigenesis(Zevin and Benowitz 1999).Analysis of INV breakpoints
390 showed clustering at similar genomic positions in 7 of the 9 positive samples,
391 resulting in INV fragments of 18–30 Mb. These clustered into two distinct groups:
392 LSCC1, 3, 6, 14 and LSCC5, 7, 11 (Supplemental Fig. S7). Of particular note, the
393 breakpoints in the first group (LSCC1, 3, 6, 14) localized within an intron
394 of *CCDC146* and exon 7 of *RASA4* (Supplemental Fig. S7). *RASA4* encodes a
395 calcium-dependent GTPase-activating protein that inactivates RAS and suppresses the
396 MAPK pathway, functioning as a tumor suppressor. Its loss of function is implicated
397 in malignant progression across various cancers(Chen et al. 2021; Prior et al. 2020).
398 Since large inversions can disrupt gene architecture (Xu et al. 2023), this recurrent
399 event may contribute to LSCC pathogenesis by altering *RASA4* transcription or
400 function.

401

402 **Analysis and validation of a somatic SRE in 71.8% LSCC patients**

403

404 A somatic INS, observed in 66% (10/15) of LSCC tumors, was the most frequent

105 recurrent SV in our study (Fig. 6A). It was located at Chr1: 224,011,691-224,016,553,
106 a simple repeat region annotated by RepeatMasker(Tarailo-Graovac and Chen 2009). The
107 region initially consisted of a short tandem repeat with the (AATGG) repeat motif in
108 normal samples, while the somatic alteration in 9/10 LSCC samples involved an
109 expansion of this pattern, specifically a simple repeat expansion (SRE) (Fig. 6B).The
110 LSCC11 sample differed, harboring Human Satellite II as its repeat motif (HSATII , a
111 satellite DNA derived from the (CATTC)_n repeat (Altemose et al. 2022)). The somatic
112 SRE exhibited low supporting read ratios in six samples (10%-30%), while four
113 exhibited stronger support (>50%) (Supplemental Fig. S8A). SRE lengths varied
114 among tumors (median: 581 bp, Fig. 6C and S8A). LSCC15 had a 3,610 bp INS, while
115 other samples showed peak SRE lengths around 450 bp (mean 517 bp, SD 382 bp).
116 Additionally, two patients harbored multiple INS of different lengths. For example,
117 the LSCC1 sample harbored 390 bp, 650 bp, and 1,030 bp INS (Supplemental Fig.
118 S8A and Fig. 6A), indicating more tumor heterogeneity.

119
120 To validate the presence of somatic SRE identified through LRS WGS, target PCR
121 amplification of the SRE region was performed on DNA from available 14 paired
122 tumor and blood samples (including 9 tumors with INS identified by LRS WGS). Gel
123 electrophoresis indicated that 6/9 tumors displayed longer fragments compared to
124 blood samples (Fig. 6D and 6E). The remaining 3 samples presented challenges due to
125 low SRE supporting ratio (20% and 30% for LSCC10 and LSCC11) and a long
126 insertion length in LSCC15 (3,610 bp)(Supplemental Fig. S8A), complicating the gel
127 electrophoresis validation. Despite these challenges, our results confirm the
128 authenticity and reliability of these high-frequency somatic SRE.

129
130 To assess the prevalence of the high-frequency somatic SRE in LSCC tumors, we
131 further collected 24 paired LSCC samples for PCR validation and ONT Sequence
132 (Supplemental Fig. S8B-S8G, S9-S19). Gel electrophoresis and LRS data of the target
133 PCR products revealed that 18/24 LSCC samples carried the somatic SRE with an
134 (AATGG) repeat motif, with insertion lengths ranging from 130 bp to 3,693 bp
135 (Supplemental Fig. S9-S19, Supplementary Table S6). Furthermore, this SRE was
136 detected in the peritumoral tissues of 15/21(71.42%) LSCC tumors (Supplemental Fig.

137 S8G, Supplementary Table S6). Together, the high prevalence of this somatic SRE in
138 peritumoral tissue suggested it is an early genomic event, while its significantly
139 elevated allele frequency in tumor compare with peritumoral tissue (t -test, p -value =
140 0.0046; Supplemental Fig. S8H) indicates positive selection of the SRE-harboring
141 clone during tumor evolution.

142
143 To further explore the functional role of the high-frequency somatic SRE, we analysis
144 the chromatin environment surrounding the SRE. In the A673, A549, and SK-N-MC
145 cell lines, the SRE region overlapped with DNase-H3K4me3 signals, which are known
146 to mark transcriptional start sites(Je et al. 2020)(Fig. 6C). While the target genes
147 regulated by this DNase-H3K4me3 mark have not been extensively characterized,
148 Hi-C data from the hypopharyngeal cancer FaDu cell line suggested that the region
149 resides within the same topological domain (TAD) as the nearby coding genes
150 *TP53BP2* and *FBXO28* (Supplemental Fig. S20A). Furthermore, data from the
151 Four-Dimensional (4D) Nucleosome Consortium(Zhu et al. 2022) supported the
152 presence of three-dimensional chromatin interactions between this DNase-H3K4me3
153 region and the genes *TP53BP2* and *FBXO28* (Fig. 6C). Critically, analysis of matched
154 TCGA LSCC samples confirmed that, among genes near the SRE, the expression of
155 *FBXO28* and *TP53BP2* exhibited the most significant tumor-specific upregulation
156 (t -test, p -value < 0.001; Supplemental Fig. S20B). These convergent lines of evidence
157 raised the possibility that this somatic SRE may alter the regulation of these two genes.
158 Consistent with this hypothesis, qPCR analysis of our clinical LSCC samples revealed
159 that tumors harboring the SRE exhibited significantly higher expression levels of both
160 *TP53BP2* and *FBXO28* compared to SRE-negative tumors (t -test, p -value = 0.00039
161 and p -value = 8.1×10^{-6} , Fig. 6F), directly linking the presence of the SRE to
162 upregulated expression of its putative target genes. *TP53BP2* encodes the
163 apoptosis-stimulating protein, which interacts with P53 family members to modulate
164 cell apoptosis and growth(Huo et al. 2023). *FBXO28*, which encodes a ubiquitin ligase,
165 has been reported to be aberrantly overexpressed in human epithelial cancer cell
166 lines(Song et al. 2024).

167

168 We next examine the regulatory effect of the high-frequency somatic SRE on gene
169 expression using cell line models. With gel electrophoresis, we identified longer
170 fragments in two LSCC cell lines (SNU46 and SNU899) compared to a normal human
171 bronchial epithelial cell line (BEAS-2B) (Supplemental Fig. S8B). Nanopore
172 sequencing of the PCR products further confirmed that both LSCC cell lines harbored
173 an expected INS with AATGG repeat motif. Specifically, SNU46 and SNU899 carried
174 homozygous INS with lengths of 165 bp and 240 bp, respectively (Supplemental Fig.
175 S20C and S20D). In contrast, no significant sequence length variation was observed in
176 BEAS-2B cells compared to the reference genome. We then selected the SNU899 cell
177 line, which has a homozygous SRE with longer insertion, and the BEAS-2B cell line
178 for quantitative real-time PCR (QPCR) to assess whether the SRE affects the
179 expression levels of *TP53BP2* and *FBXO28*. Consistent with findings in primary
180 tumors, We found that *TP53BP2* and *FBXO28* expression levels were significantly
181 higher in the SNU899 cell line than in the BEAS-2B cell line (Fig. 6G, *t*-test,
182 *TP53BP2*: *p*-value < 0.01; *FBXO28*: *p*-value < 0.001). To rule out tissue-specific
183 baseline differences as the cause, we analyzed normal tissue data from TCGA (LSCC
184 vs. LUSC) and cell line data from DepMap. This confirmed that expression of these
185 genes in laryngeal contexts was not inherently higher than in pulmonary contexts,
186 supporting the SRE's specific role in upregulation (Supplemental Fig. S20E, F). To
187 directly test whether the SRE alters local chromatin architecture, we performed
188 3C-PCR with an anchor primer (Anchor-S) placed downstream of the simple repeat
189 region where the SRE occurs. In SNU899 (SRE-positive), the Anchor-S locus
190 exhibited strong, specific spatial proximity to the F1 regulatory site upstream of
191 *FBXO28*. In BEAS-2B (SRE-negative), the same Anchor-S locus showed a markedly
192 weaker interaction with *FBXO28*-F1. Instead, its predominant interactions were with
193 sites near *TP53BP2* (T3/T8) (Supplemental Fig. S20G). This reciprocal pattern
194 demonstrates that the presence of the SRE is associated with a distinct
195 three-dimensional chromatin conformation, specifically promoting spatial proximity
196 to the regulatory region of *FBXO28*. Together, these findings suggest that the somatic
197 SRE may upregulate *TP53BP2* and *FBXO28* expression in LSCC by altering the local
198 three-dimensional chromatin architecture.

199 Discussion

200
201 Somatic structural variations (SVs) are critical hallmarks of tumorigenesis (Cosenza et al.
202 2022). However, technical limitations of SRS in sensitively detecting somatic SVs,
203 compounded by the tissue heterogeneity characteristic of HNSCC, have left the landscape of
204 high-frequency somatic SVs in LSCC underexplored. This gap in knowledge has hindered the
205 identification of potential LSCC-specific molecular biomarkers for effective cancer
206 management. To address these challenges, we developed SomaGauss-SV, a LRS-optimized
207 computational framework for somatic SVs detection. Benchmarking against established tools
208 revealed that SomaGauss-SV offered a superior balance of sensitivity and specificity, with an
209 F1-score improvement of 11.3-33.9%. SomaGauss-SV excels in two key areas: (1) it reduces
210 alignment errors in repetitive genomic regions through a multi-step filtering approach, and (2)
211 it accurately distinguishes true somatic INS, DUP and DEL in the presence of size-variable
212 INS/DUP/DEL using a Gaussian mixture model. Despite these advances, significant
213 technological challenges persist. Current LRS platforms still struggle with resolving complex
214 breakpoint architectures, particularly those involving nested rearrangements and ultra-long
215 SVs (>100 kb)(Keskus et al. 2025), limitations shared across existing detection tools,
216 including SomaGauss-SV. Furthermore, the current benchmarking landscape relies
217 heavily on consensus-based “gold standard” sets derived from multiple sequencing
218 platforms and computational tools rather than orthogonal experimental validation.
219 This approach may introduce methodological bias, especially for complex SVs, which
220 remain difficult and costly to validate experimentally, resulting in limited and
221 potentially incomplete ground-truth datasets (e.g., the COLO829 benchmark provides
222 only 68 gold-standard SVs (Espejo Valle-Inclan et al. 2022)). We also observed that the
223 choice of benchmarking strategy (minda vs Truvari) itself significantly impacts
224 performance evaluation. There is currently no consensus in the field on whether SV
225 type should be incorporated as a matching criterion during benchmarking. While
226 classifying complex SVs can be ambiguous, applying a consistent typing standard is
227 necessary for the robust evaluation of simple somatic variants (e.g., INS, DEL, DUP),
228 as type-agnostic matching may misclassify genuine events. However, recent innovations
229 in high-throughput ultra-long-read sequencing technologies (Wang et al. 2021) and

graph-based or assembly-based long-read alignment algorithms (Li and Durbin 2024) offer promise for overcoming these obstacles.

We systematically analyzed the LRS WGS data of 15 paired LSCC samples using SomaGauss-SV, revealing the distribution and sequence features of somatic SVs. We identified a dose-dependent positive correlation between smoking intensity and somatic DEL burden in LSCC. Previous studies indicated tobacco exposure is the primary driver of HNC, with tobacco users exhibiting higher frequencies of somatic SNVs and indels (Alexandrov et al. 2016; Torrens et al. 2025). Notably, tobacco-related mutational signatures show higher burdens and frequencies in laryngeal cancers compared to other HNC subsites (Alexandrov et al. 2016; Torrens et al. 2025). However, the link between tobacco and somatic SVs remains less clearly defined (Yoshida et al. 2020; Alexandrov et al. 2016; Torrens et al. 2025; Jethwa and Khariwala 2017). Our study fills this gap by demonstrating the role of chemical carcinogens in chromosomal instability. Specifically, we suggested that the accumulation of somatic DEL may result from the synergistic effects of multiple carcinogens in tobacco smoke. DNA adducts from polycyclic aromatic hydrocarbons and nitrosamines can block replication forks and cause DNA double-strand breaks (Yoshida et al. 2020). These breaks, when repaired poorly, may lead to chromosomal deletions. Also, free radicals and reactive oxygen species from tobacco can cause more DNA damage through oxidative stress, promoting deletions during faulty repair (Yoshida et al. 2020). Tobacco components might also disrupt cell division or cycle regulation (Aw et al. 2021), increasing chromosomal instability. Although high expression of *CYP1A1*, the key enzyme involved in the metabolism of benzopyrene and nitrosamines, has been proposed as a potential cause of the elevated tobacco-related mutation signatures in laryngeal tumors (Degawa et al. 1994; Hecht and Hatsukami 2022; Yamazaki et al. 1992), the specific tobacco carcinogens responsible for the observed DEL burden remain to be identified in vitro. Clinically, DEL burden could serve as a potential biomarker for risk stratification and therapeutic target screening in LSCC. Further investigation into high-frequency deletions associated with smoking may provide evidence for preventive interventions targeting high-dimensional genomic regions. For example, somatic DEL might disrupt CTCF binding in 3D chromatin architecture, influencing the expression of immune genes like *HLA-DRB5* and *HLA-DRA*. However, the limited sample size in this study restricts the generalizability of the observed association between somatic SVs burden and smoking.

Other potential confounders, such as alcohol consumption and dietary factors, may also influence the results. Future studies with larger cohorts and the incorporation of multivariable models or *in vivo/vitro* experiments are needed to confirm these findings.

In addition, 226 high-frequency SV loci were identified, exhibiting selective advantages during tumor clonal evolution. Cross-cancer comparison analysis revealed that high-frequency SVs at Chr7:121,603,263-121,603,506 and Chr15:29,945,811-29,950,000 overlap with hotspot somatic repeat expansions reported in pan-cancer studies (Erwin et al. 2023), suggesting their involvement in conserved carcinogenic pathways across various cancer types. Significant CNV regions in laryngeal cancer from the TCGA database also overlap with the high-frequency large DELs, DUPs and INSSs, further validating the unique strength of LRS technology in resolving multi-scale genomic rearrangements. A key issue to explore further is how to define the boundary between CNVs identified through sequencing depth and DEL/INS/DUP detected by long reads. Furthermore, we investigated co-occurring SV patterns within individual reads, which may represent complex structural rearrangements. For instance, in sample LSCC1, a single read simultaneously contained an insertion of a simple repeat (AT), a deletion within a LINE/L1 element, and an insertion of a SINE/*AluY* sequence (Supplemental Fig. S21). In LSCC3, we identified an event structurally resembling an “inverted duplication” as previously reported (Schloissnig et al. 2025), involving both INV and DUP on the same allele (Supplemental Fig. S22). These examples illustrate the capacity of long-read sequencing to resolve complex, multi-variant rearrangements.

Notably, analysis of both WGS and an independent validation cohort identified a somatic high-frequency SRE in 71.79% (28/39) of LSCC patients. Among the 21 patients with available peritumoral tissue, the SRE was detected in both tumor and matched peritumoral samples in 15 individuals. The allele frequency of the SRE was significantly higher in tumor tissue than in peritumoral tissue, suggested that this alteration may arise early during tumorigenesis and undergo subsequent positive selection. We also observed heterogeneity in SRE lengths within the same tumor, indicative of intratumoral diversity and potentially distinct selective advantages associated with different expansion sizes. This dynamic evolution was particularly

593 evident in a patient with lymph node metastasis, where the SRE was detected in the
594 primary tumor (400 bp), peritumoral tissue (125 bp), and metastatic lymph node (200
595 bp). This suggested that the SRE evolves dynamically throughout tumor progression
596 and metastasis, conferring varying selective advantages in different
597 microenvironments. However, the mechanisms linking SRE length variation to
598 phenotypic diversity in tumor cells require further investigation.

599
600 The recurrent somatic SRE identified in this study holds significant potential as a biomarker
601 for LSCC. Integrating 3C-PCR and qPCR results from our cell line models, we propose a
602 hypothesis in which the SRE may regulate the expression of neighboring genes by altering
603 local three-dimensional chromatin architecture. Specifically, the SRE-mediated sequence
604 elongation could increase spatial proximity between a distal regulatory element and the
605 promoters of *TP53BP2* and *FBXO28*, thereby enhancing their transcription. Although external
606 data suggested that baseline expression of these genes in laryngeal contexts is not inherently
607 higher than in pulmonary tissues, the inherent tissue-of-origin heterogeneity between the lung
608 epithelial cell line BEAS-2B and the LSCC-derived SNU899 line warrants caution in
609 interpreting cross-tissue comparisons. Nevertheless, the functional relevance of these genes in
610 cancer is well-supported: elevated *FBXO28* expression is linked to poor prognosis in ovarian
611 cancer, where it activates the TGF- β 1/Smad2/3 pathway to promote proliferation and invasion
612 (Song et al. 2024). *TP53BP2*, a member of the ASPP family, regulates apoptosis, proliferation,
613 and autophagy, and its upregulation correlates with poor prognosis in colorectal cancer and
614 other malignancies (Huo et al. 2023). To definitively test our mechanistic hypothesis, future
615 studies employing CRISPR-Cas9 to delete the SRE in relevant models would be essential to
616 directly assess its impact on *TP53BP2* and *FBXO28* expression and downstream phenotypes,
617 thereby excluding the influence of tissue-of-origin heterogeneity.

618
619 In conclusion, the bioinformatics pipeline developed in this study achieves an optimal balance
620 between sensitivity and specificity, offering a robust tool for future research on tumor somatic
621 SVs. This study represents the first comprehensive genome-wide survey of somatic SVs in
622 LSCC using LRS, with the generated whole-genome sequencing data and identified SVs
623 offering a valuable resource for clinical research in the field.

624

325 **Method**

326 **Sample and patient information collection**

327 This study involved 39 patients with LSCC who underwent surgical resection at West
328 China Hospital, Sichuan University. Tumor samples and matched blood samples were
329 collected from these patients during surgery and venipuncture. Throughout the tissue
330 collection and utilization process, this study strictly adhered to the Declaration of
331 Helsinki and received formal approval from the Ethics Review Committee of West
332 China Hospital, Sichuan University. All patients participating in this study have
333 signed written informed consent forms. In addition, key clinical information,
334 including clinical staging, smoking history, age, and gender, was thoroughly
335 documented for each patient.

336

337 **Cell culture**

338 Human cell lines, obtained from Procell (Wuhan Pricella Biotechnology Co., Ltd.) were
339 utilized in this study. The cells were grown in DMEM medium (Gibco, C11995500BT) or
340 RPMI 1640 medium (Gibco, C11875500BT), supplemented with 10% fetal bovine serum
341 (FBS) and 100 U/mL Penicillin-Streptomycin (P/S). All cell lines were maintained at 37°C
342 and 5% CO₂. Periodic mycoplasma testing was conducted using the Mycoplasma Detection
343 Kit (Vazyme, D101). Cells were passaged at 70–90% confluency with 0.25% Trypsin-EDTA
344 (Gibco, 25200-072).

345

346 **Gene expression, chromatin conformation, and genomic validation assays**

347 For gene expression analysis, total RNA was extracted, reverse transcribed, and quantified by
348 qRT-PCR using gene-specific primers. To investigate chromatin interactions between the SRE
349 region and target genes (*FBXO28* and *TP53BP2*), Chromatin Conformation Capture (3C)
350 libraries were generated from cross-linked SNU899 cells via EcoRI digestion and
351 intra-molecular ligation, and specific interactions were quantified by 3C-qPCR. Additionally,
352 genomic DNA was used to validate the presence of the SRE in clinical samples and cell lines
353 by PCR amplification followed by Sanger sequencing. The primer sequences are listed in the
354 oligonucleotide table (Supplementary Table S7). Detailed protocols, including all primer
355 sequences and reaction conditions, are provided in the Supplemental Methods.

356

357 **Nanopore sequencing library preparation**

358 For whole-genome sequencing, high-molecular-weight genomic DNA was extracted and used
359 for library preparation with Oxford Nanopore ligation kits (SQK-LSK109 or SQK-LSK114),
360 followed by sequencing on PromethION flow cells. For targeted sequencing of the SRE locus,
361 PCR amplicons were prepared into libraries using either the Oxford Nanopore or Qitan Tech
362 nanopore sequencing platforms and sequenced on GridION X5 or QNome-3841 instruments,
363 respectively. Detailed protocols for library preparation, including specific reagents, reaction
364 conditions, and purification steps, are provided in the Supplemental Methods.

365

366 **Preprocessing of long-read sequencing data**

367 For data preprocessing, raw long-read sequencing data were quality-filtered and assessed
368 using NanoFilt (version 2.8.0)(Lee et al. 2021) and NanoPlot (version 1.39.0)(De Coster and
369 Rademakers 2023) , respectively. Filtered reads were aligned to the hg38 reference genome
370 using minimap2 (version 2.26)(Li 2018). Resulting SAM files were converted, sorted, and
371 read-tagged using SAMtools (version 1.9)(Danecek et al. 2021) to generate final BAM files.
372 For the analysis of SRE-targeted PCR products , a similar preprocessing pipeline was applied,
373 with specific length filtering to ensure complete spanning of the repeat region. The processed
374 alignments were visualized using the Integrative Genomics Viewer (IGV) (Robinson et al.
375 2011). Detailed parameters and steps are provided in the Supplemental Methods.

376

377 **Acquisition and preprocessing of nanopore sequencing data and gold standard somatic 378 structural variations data from cell lines**

379 According to the reference(Keskus et al. 2025), we downloaded the nanopore
380 sequencing data for five paired cell lines: HCC1395/HCC1395BL, H1437/H1437BL,
381 H2009/H2009BL, HCC1937/HCC1937BL, and HCC1954/HCC1954BL, as well as the
382 corresponding gold standard sets of somatic structural variations for these cell lines.
383 To ensure the accuracy of the analysis, we filtered the gold standard data, retaining
384 only those somatic structural variations that could be detected by the nanopore
385 sequencing platform, totaling 4058, including DEL 1493, INS 906, DUP 385, INV 48,
386 and BND 1226. Subsequently, following the standard workflow for nanopore
387 sequencing data preprocessing, we aligned and processed these data, ultimately
388 obtaining BAM files with "blood" and "tumor" tags, respectively. The LRS WGS data
389 and benchmark SV sets of five cell line pairs were downloaded from NCBI SRA BioProject

390 PRJNA1086849.

391

392 **SomaGauss-SV workflow**

393 The SomaGauss-SV workflow for somatic SV detection involved four key stages.

394 First, merged tumor-normal BAM files were analyzed with Sniffles2(version

395 2.0.6)(Smolka et al. 2024) to call candidate variants. Second, a type-specific screening

396 strategy was applied: candidate DEL were validated using a Gaussian mixture model

397 on locally extracted reads, while candidate INS and DUP were re-analyzed using

398 Straglr(Chiu et al. 2021), a length-based Gaussian modeling tool optimized for

399 repetitive regions. Third, for clinical samples, variants were filtered against a Panel of

700 Normals derived from Chinese population SV data (Wu et al. 2021). using Jasmine

701 (version 1.1.5)(Kirsche et al. 2023) to remove common germline polymorphisms.

702 Finally, complex variants (INV and BND) were manually validated through

703 visualization using the Integrative Genomics Viewer (IGV) (Robinson et al. 2011).

704 Detailed parameters and scripts are provided in the Supplemental Methods.

705

706 To validate our foundational tool selection and demonstrate the incremental value of
707 each step in the SomaGauss-SV pipeline, we conducted systematic benchmarking.

708 First, Sniffles2 was empirically selected over cuteSV(Jiang et al. 2020) as the initial

709 variant caller due to its superior performance (mean F1 score: 68.7% vs. 53.4%,

710 Supplemental Fig. S23 A and B). Subsequently, a stepwise comparison of three

711 workflows: (1) Original Sniffles: The merged VCFs obtained by running Sniffles2

712 alone, (2) Filter Reads: The SV set after further filtering out variants with any

713 supporting reads in the matched normal (blood) sample, (3) SomaGauss (Full

714 Pipeline): The final output after applying our Gaussian Mixture Model (GMM)-based

715 filtering step to the "Filter Reads" set. It significantly enhanced precision (from 14.5%

716 to 71.51%) and F1 scores (overall +39.63%, with a more pronounced improvement of

717 +47.86% in challenging VNTR regions) by effectively removing subtle false positives,

718 particularly those arising from alignment errors in repetitive sequences. This confirms

719 that each component of SomaGauss-SV adds substantial value, culminating in a

720 highly accurate and robust somatic SV detection workflow (Supplemental Fig. S23 C-F).

721

722 **Benchmarking and performance evaluation of SV detection tools**

723 For benchmarking comparisons, somatic structural variants were independently
724 detected using three established calls: nanomonsv(Shiraishi et al. 2023), Severus
725 (Keskus et al. 2025), and SVision-pro (Wang et al. 2025), following each call's
726 recommended parameters and workflow. The performance of all callers was then
727 evaluated against the gold-standard variant set using two evaluation tools:
728 Truvari(English et al. 2022) and minda (Keskus et al. 2025). Detailed command-line
729 parameters are provided in the Supplemental Methods.

730

731 **Evaluating software runtime efficiency**

732 To evaluate the runtime efficiency of each software tool, we used the command
733 `"/usr/bin/time -v"` followed by the respective software command line to measure both
734 execution time and maximum memory usage. Given that most tools involve multiple
735 distinct steps in their workflows—for example, Severus includes separate phases for
736 SNP calling and phasing of the normal BAM, haplotagging of normal and tumor
737 BAMs, and somatic SV calling, while SomaGauss-SV comprises steps such as
738 merging BAM files, calling SVs from the merged BAM, and running Straglr. We
739 adopted the following practical approach to summarize resource consumption: Total
740 Elapsed Real Time: We summed the Elapsed Real Time from each step of a given tool
741 to represent its overall runtime. Maximum resident set size: We recorded the
742 maximum Resident Set Size observed across all steps as the tool's peak memory
743 footprint. This method provides a consolidated and realistic view of computational
744 resource requirements for each somatic SV detection pipeline under evaluation.

745

746 **Variant sequence annotation and genomic context analysis**

747 For sequence annotation and refinement, the consensus sequences of somatic INS were first
748 corrected using Iris (version 1.0.4)(Kirsche et al. 2023) and then annotated for repetitive
749 elements using RepeatMasker (version 4.1.2)(Tarailo-Graovac and Chen 2009). To
750 characterize the genomic context of somatic DEL, their breakpoint regions were intersected
751 with a RepeatMasker-annotated genome using BEDTools (version 2.30.0)(Quinlan and Hall
752 2010) to identify overlaps with repetitive sequences. Detailed parameters for each tool are
753 provided in the Supplemental Methods.

754

755 Genomic enrichment and recurrence analysis of somatic SV

756 To assess genomic enrichment of somatic SVs, we performed two complementary analyses.
757 First, SV breakpoints were intersected with functional genomic elements (genes and candidate
758 *cis*-regulatory elements from UCSC (Cm et al. 2020)) using BEDTools. Enrichment fold and
759 statistical significance were calculated using Fisher's exact test based on a contingency table of
760 variant counts within and outside these elements relative to genomic background. Second, to
761 address alignment ambiguity in repetitive regions (Supplemental Fig. S1, (Sahlin et al. 2023;
762 Ahsan et al. 2023)), we constructed a non-uniform window BED file by combining annotated
763 simple repeat regions with 500 bp/5 kb windows in non-repetitive regions. Somatic SV loci
764 were then intersected with these windows using BEDTools to calculate sample carrier rates,
765 and adjacent regions with shared variants were merged prior to downstream analysis. Detailed
766 procedures are provided in the Supplemental Methods.

767

768 Data access

769 The BAM files generated in this study have been submitted to the Genome Sequence Archive
770 (Genomics, Proteomics & Bioinformatics 2025) in the National Genomics Data Center
771 (Nucleic Acids Res 2025), China National Center for Bioinformation / Beijing Institute of
772 Genomics, Chinese Academy of Sciences (GSA-Human; <https://ngdc.cncb.ac.cn/gsa-human/>)
773 under accession number HRA017212.

774

775 Code availability

776 SomaGauss-SV is available at <https://github.com/XuYan000131/SomaGauss-SV> and as
777 Supplemental Code.

778

779 Acknowledgements

780 We thank Ranlei Wei for managing the computational resources required for this study. This
781 work was supported by grants from the National Natural Science Foundation of China
782 (82371883 and 12441517 to JFL, 82173383 to DX, 32200508 to LX), the Noncommunicable
783 Chronic Diseases-National Science and Technology Major Project (2026ZD0553400,
784 2026ZD0553404 to DX), the National Key Science and Technology Special Project for Deep
785 Earth Research (2024ZD1000600 and 2024ZD1000606 to JFL), the National Key Research
786 and Development Program of China (2024YFF1207002,4 to JFL and YXQ) , the Foundation

787 of Sichuan Provincial Science and Technology Program (2025ZNSFSC0728 to JFL), the 1·3·5
788 project for disciplines of excellence, West China Hospital, Sichuan University (ZYYC23024
789 to DX), the Project Supported by Scientific Reserch Fund of SiChuan Provincial Education
790 Department (2024NSFC1187 to YL).

791

792 **Authors' contributions**

793 LX and XYL conceived and developed the methodology, performed bioinformatic analyses.
794 JFL and YXQ collected the clinical samples and clinical data. BYY, XL, XY, XC and HHZ
795 collected the clinical samples. YH and YL designed and performed experiments. LX, XYL,
796 YXQ, YH wrote the manuscript. DX and JFL responsible for supervision of research, data
797 interpretation, and manuscript preparation. All authors read and approved the final manuscript.

798

799 **Competing interests**

300 Dan Xie is the co-founder of Qitan Technology. Other authors declare no competing interests.

301

302 **References**

- 303 Afshar-Kharghan V. 2017. The role of the complement system in cancer. *J Clin Invest* 127:
304 780–789.
- 305 Ahsan MU, Liu Q, Perdomo JE, Fang L, Wang K. 2023. A survey of algorithms for the
306 detection of genomic structural variants from long-read sequencing data. *Nat Methods*
307 20: 1143–1158.
- 308 Aldrup-MacDonald ME, Kuo ME, Sullivan LL, Chew K, Sullivan BA. 2016. Genomic
309 variation within alpha satellite DNA influences centromere location on human
310 chromosomes with metastable epialleles. *Genome Res* 26: 1301–1311.
- 311 Alexandrov LB, Ju YS, Haase K, Van Loo P, Martincorena I, Nik-Zainal S, Totoki Y,
312 Fujimoto A, Nakagawa H, Shibata T, et al. 2016. Mutational signatures associated with
313 tobacco smoking in human cancer. *Science* 354: 618–622.
- 314 Altemose N, Glennis A, Bzikadze AV, Sidhwani P, Langley SA, Caldas GV, Hoyt SJ, Uralsky
315 L, Ryabov FD, Shew CJ, et al. 2022. Complete genomic and epigenetic maps of
316 human centromeres. *Science* 376: eabl4178.
- 317 Aw C, S T, A B. 2021. Relationships among smoking, oxidative stress, inflammation,
318 macromolecular damage, and cancer. *Mutation research Reviews in mutation research*
319 787. <https://pubmed.ncbi.nlm.nih.gov/34083039/> (Accessed March 25, 2025).
- 320 Beroukhim R, Zhang X, Meyerson M. 2016. Copy number alterations unmasked as enhancer
321 hijackers. *Nat Genet* 49: 5–6.
- 322 Bhatia R, Gautam SK, Cannon A, Thompson C, Hall BR, Aithal A, Banerjee K, Jain M,
323 Solheim JC, Kumar S, et al. 2019. Cancer-associated mucins: role in immune
324 modulation and metastasis. *Cancer Metastasis Rev* 38: 223–236.
- 325 Blanchard N, Di Bartolo V, Hivroz C. 2002. In the immune synapse, ZAP-70 controls T cell
326 polarization and recruitment of signaling proteins but not formation of the synaptic
327 pattern. *Immunity* 17: 389–399.

- 328 Bray F, Laversanne M, Sung H, Ferlay J, Siegel RL, Soerjomataram I, Jemal A. 2024. Global
329 cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide
330 for 36 cancers in 185 countries. *CA Cancer J Clin* 74: 229–263.
- 331 Cavaliere M, Bisogno A, Scarpa A, D’Urso A, Marra P, Colacurcio V, De Luca P, Ralli M,
332 Cassandro E, Cassandro C. 2021. Biomarkers of laryngeal squamous cell carcinoma: a
333 review. *Ann Diagn Pathol* 54: 151787.
- 334 Chen J, Huang J, Huang Q, Li J, Chen E, Xu W. 2021. RASA4 inhibits the HIF α signaling
335 pathway to suppress proliferation of cervical cancer cells. *Bioengineered* 12:
336 10723–10733.
- 337 Chen JP, Diekmann C, Wu H, Chen C, Della Chiara G, Berrino E, Georgiadis KL, Bouwman
338 BAM, Viridi M, Harbers L, et al. 2024. scCircle-seq unveils the diversity and
339 complexity of extrachromosomal circular DNAs in single cells. *Nat Commun* 15:
340 1768.
- 341 Chiu R, Rajan-Babu I-S, Friedman JM, Birol I. 2021. Straglr: discovering and genotyping
342 tandem repeat expansions using whole genome long-read sequences. *Genome Biology*
343 22: 224.
- 344 Choo Z-N, Behr JM, Deshpande A, Hadi K, Yao X, Tian H, Takai K, Zakusilo G, Rosiene J,
345 Da Cruz Paula A, et al. 2023. Most large structural variants in cancer genomes can be
346 detected without long reads. *Nat Genet* 55: 2139–2148.
- 347 Cm L, Gp B, J C, H C, M D, Jn G, As H, Bt L, Lr N, Cc P, et al. 2020. UCSC Genome
348 Browser enters 20th year. *Nucleic acids research* 48.
349 <https://pubmed.ncbi.nlm.nih.gov/31691824/> (Accessed March 30, 2025).
- 350 Cosenza MR, Rodriguez-Martin B, Korbel JO. 2022. Structural Variation in Cancer: Role,
351 Prevalence, and Mechanisms. *Annu Rev Genomics Hum Genet* 23: 123–152.
- 352 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
353 McCarthy SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools.
354 *Gigascience* 10: giab008.
- 355 De Coster W, Rademakers R. 2023. NanoPack2: population-scale evaluation of long-read
356 sequencing data. *Bioinformatics* 39: btad311.
- 357 Degawa M, Stern SJ, Martin MV, Guengerich FP, Fu PP, Ilett KF, Kaderlik RK, Kadlubar FF.
358 1994. Metabolic activation and carcinogen-DNA adduct detection in human larynx.
359 *Cancer Res* 54: 4915–4919.
- 360 ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the
361 human genome. *Nature* 489: 57–74.
- 362 English AC, Menon VK, Gibbs RA, Metcalf GA, Sedlazeck FJ. 2022. Truvari: refined
363 structural variant comparison preserves allelic diversity. *Genome Biol* 23: 271.
- 364 Erwin GS, Gürsoy G, Al-Abri R, Suriyaprakash A, Dolzhenko E, Zhu K, Hoerner CR, White
365 SM, Ramirez L, Vadlakonda A, et al. 2023. Recurrent repeat expansions in human
366 cancer genomes. *Nature* 613: 96–102.
- 367 Espejo Valle-Inclan J, Besselink NJM, de Bruijn E, Cameron DL, Ebler J, Kutzera J, van
368 Lieshout S, Marschall T, Nelen M, Priestley P, et al. 2022. A multi-platform reference
369 for somatic structural variation detection. *Cell Genom* 2: 100139.
- 370 Gendler SJ, Spicer AP. 1995. Epithelial mucin genes. *Annu Rev Physiol* 57: 607–634.
- 371 Graham LM, Brown GD. 2009. The Dectin-2 family of C-type lectins in immunity and
372 homeostasis. *Cytokine* 48: 148–155.
- 373 Han B, Zheng R, Zeng H, Wang S, Sun K, Chen R, Li L, Wei W, He J. 2024. Cancer
374 incidence and mortality in China, 2022. *Journal of the National Cancer Center* 4:
375 47–53.

- 376 Hecht SS, Hatsukami DK. 2022. Smokeless tobacco and cigarette smoking: chemical
377 mechanisms and cancer prevention. *Nat Rev Cancer* 22: 143–155.
- 378 Hormozdiari F, Alkan C, Ventura M, Hajirasouliha I, Malig M, Hach F, Yorukoglu D, Dao P,
379 Bakhshi M, Sahinalp SC, et al. 2011. Alu repeat discovery and characterization within
380 human genomes. *Genome Res* 21: 840–849.
- 381 Huo Y, Cao K, Kou B, Chai M, Dou S, Chen D, Shi Y, Liu X. 2023. TP53BP2: Roles in
382 suppressing tumorigenesis and therapeutic opportunities. *Genes Dis* 10: 1982–1993.
- 383 Iakovlev DA, Alekseeva IV, Vorobjev YN, Kuznetsov NA, Fedorova OS. 2019. The Role of
384 Active-Site Residues Phe98, His239, and Arg243 in DNA Binding and in the Catalysis
385 of Human Uracil-DNA Glycosylase SMUG1. *Molecules* 24: 3133.
- 386 ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. 2020. Pan-cancer analysis
387 of whole genomes. *Nature* 578: 82–93.
- 388 Je M, Mj P, He P, Cb E, N S, J A, T K, Ca D, A D, R K, et al. 2020. Expanded encyclopaedias
389 of DNA elements in the human and mouse genomes. *Nature* 583.
390 <https://pubmed.ncbi.nlm.nih.gov/32728249/> (Accessed April 10, 2025).
- 391 Jethwa AR, Khariwala SS. 2017. Tobacco-related carcinogenesis in head and neck cancer.
392 *Cancer Metastasis Rev* 36: 411–423.
- 393 Johnson DE, Burtneß B, Leemans CR, Lui VWY, Bauman JE, Grandis JR. 2020. Head and
394 neck squamous cell carcinoma. *Nat Rev Dis Primers* 6: 92.
- 395 Kashi Y, King DG. 2006. Simple sequence repeats as advantageous mutators in evolution.
396 *Trends Genet* 22: 253–259.
- 397 Keskus AG, Bryant A, Ahmad T, Yoo B, Aganezov S, Goretsky A, Donmez A, Lansdon LA,
398 Rodriguez I, Park J, et al. 2025. Severus detects somatic structural variation and
399 complex rearrangements in cancer genomes using long-read sequencing. *Nat*
400 *Biotechnol*. <https://www.nature.com/articles/s41587-025-02618-8> (Accessed April 22,
401 2025).
- 402 Kirsche M, Prabhu G, Sherman R, Ni B, Battle A, Aganezov S, Schatz MC. 2023. Jasmine
403 and Iris: population-scale structural variant comparison and analysis. *Nat Methods* 20:
404 408–417.
- 405 Krzywinski M, Schein J, Birol İ, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA.
406 2009. Circos: An information aesthetic for comparative genomics. *Genome Res* 19:
407 1639–1645.
- 408 Kufe DW. 2009. Mucins in cancer: function, prognosis and therapy. *Nat Rev Cancer* 9:
409 874–885.
- 410 Lechien JR, Sadoughi B, Hans S. 2022. Laryngeal cancers in paediatric and young adult
411 patients: epidemiology, biology and treatment. *Curr Opin Otolaryngol Head Neck Surg*
412 30: 145–153.
- 413 Lee S, Nguyen LT, Hayes BJ, Ross EM. 2021. Prowler: a novel trimming algorithm for
414 Oxford Nanopore sequence data. *Bioinformatics* 37: 3936–3937.
- 415 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34:
416 3094–3100.
- 417 Li H, Durbin R. 2024. Genome assembly in the telomere-to-telomere era. *Nat Rev Genet* 25:
418 658–670.
- 419 Lin M, Zheng X, Yan J, Huang F, Chen Y, Ding R, Wan J, Zhang L, Wang C, Pan J, et al.
420 2024. The RNF214-TEAD-YAP signaling axis promotes hepatocellular carcinoma
421 progression via TEAD ubiquitylation. *Nat Commun* 15: 4995.
- 422 Ma B, Pl D. 2002. Alu repeats and human genomic diversity. *Nature reviews Genetics* 3.
423 <https://pubmed.ncbi.nlm.nih.gov/11988762/> (Accessed March 22, 2025).

- 324 McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G.
325 2010. GREAT improves functional interpretation of cis-regulatory regions. *Nat*
326 *Biotechnol* 28: 495–501.
- 327 Merikallio H, Pincikova T, Kotortsi I, Karimi R, Li C-X, Forsslund H, Mikko M, Nyrén S,
328 Lappi-Blanco E, Wheelock ÅM, et al. 2023. Mucins 3A and 3B Are Expressed in the
329 Epithelium of Human Large Airway. *Int J Mol Sci* 24: 13546.
- 330 Nocini R, Molteni G, Mattiuzzi C, Lippi G. 2020. Updates on larynx cancer epidemiology.
331 *Chin J Cancer Res* 32: 18–25.
- 332 Pan X, Che Q, Liu D, Xie Y, Li B, Zhang S, Li T, Li G, Li X, Zheng Q, et al. 2025.
333 Development and validation of a novel endoplasmic reticulum stress-related lncRNA
334 signature in laryngeal squamous cell carcinoma. *Sci Rep* 15: 12497.
- 335 Prior IA, Hood FE, Hartley JL. 2020. The Frequency of Ras Mutations in Cancer. *Cancer Res*
336 80: 2969–2974.
- 337 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
338 features. *Bioinformatics* 26: 841–842.
- 339 R O, M R, C T. 2019. The Treatment of Laryngeal Cancer. *Oral and maxillofacial surgery*
340 *clinics of North America* 31. <https://pubmed.ncbi.nlm.nih.gov/30449522/> (Accessed
341 March 24, 2025).
- 342 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP.
343 2011. Integrative genomics viewer. *Nat Biotechnol* 29: 24–26.
- 344 Sahlin K, Baudeau T, Cazaux B, Marchet C. 2023. A survey of mapping algorithms in the
345 long-reads era. *Genome Biology* 24: 133.
- 346 Schloissnig S, Pani S, Ebler J, Hain C, Tsapalou V, Söylev A, Hüther P, Ashraf H, Prodanov
347 T, Asparuhova M, et al. 2025. Structural variation in 1,019 diverse humans based on
348 long-read sequencing. *Nature*. <https://www.nature.com/articles/s41586-025-09290-7>
349 (Accessed August 4, 2025).
- 350 Shiraishi Y, Koya J, Chiba K, Okada A, Arai Y, Saito Y, Shibata T, Kataoka K. 2023. Precise
351 characterization of somatic complex structural variations from tumor/control paired
352 long-read sequencing data with nanomonsv. *Nucleic Acids Res* 51: e74.
- 353 Smolka M, Paulin LF, Grochowski CM, Horner DW, Mahmoud M, Behera S, Kalef-Ezra E,
354 Gandhi M, Hong K, Pehlivan D, et al. 2024. Detection of mosaic and population-level
355 structural variants with Sniffles2. *Nat Biotechnol* 42: 1571–1580.
- 356 Song G, Sun Z, Chu M, Zhang Z, Chen J, Wang Z, Zhu X. 2024. FBXO28 promotes cell
357 proliferation, migration and invasion via upregulation of the TGF-beta1/SMAD2/3
358 signaling pathway in ovarian cancer. *BMC Cancer* 24: 122.
- 359 Steuer CE, El-Deiry M, Parks JR, Higgins KA, Saba NF. 2017. An update on larynx cancer.
360 *CA Cancer J Clin* 67: 31–50.
- 361 Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, Bray F. 2021. Global
362 Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide
363 for 36 Cancers in 185 Countries. *CA Cancer J Clin* 71: 209–249.
- 364 Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive elements in
365 genomic sequences. *Curr Protoc Bioinformatics Chapter 4*: 4.10.1-4.10.14.
- 366 Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, Boutselakis H, Cole CG,
367 Creatore C, Dawson E, et al. 2019. COSMIC: the Catalogue Of Somatic Mutations In
368 Cancer. *Nucleic Acids Res* 47: D941–D947.
- 369 Torrens L, Moody S, de Carvalho AC, Kazachkova M, Abedi-Ardekani B, Cheema S, Senkin
370 S, Cattiaux T, Cortez Cardoso Penha R, Atkins JR, et al. 2025. The complexity of
371 tobacco smoke-induced mutagenesis in head and neck cancer. *Nat Genet* 57: 884–896.

- 372 van Belzen IAEM, Schönhuth A, Kemmeren P, Hehir-Kwa JY. 2021. Structural variant
373 detection in cancer genomes: computational challenges and perspectives for precision
374 oncology. *NPJ Precis Oncol* 5: 15.
- 375 Wang N, Yan H, Wu D, Zhao Z, Chen X, Long Q, Zhang C, Wang X, Deng W, Liu X. 2020.
376 PRMT5/Wnt4 axis promotes lymph-node metastasis and proliferation of laryngeal
377 carcinoma. *Cell Death Dis* 11: 864.
- 378 Wang S, Lin J, Jia P, Xu T, Li X, Liu Y, Xu D, Bush SJ, Meng D, Ye K. 2025. De novo and
379 somatic structural variant discovery with SVision-pro. *Nat Biotechnol* 43: 181–185.
- 380 Wang Y, Zhao Y, Bollas A, Wang Y, Au KF. 2021. Nanopore sequencing technology,
381 bioinformatics and applications. *Nat Biotechnol* 39: 1348–1365.
- 382 Witherspoon DJ, Xing J, Zhang Y, Watkins WS, Batzer MA, Jorde LB. 2010. Mobile element
383 scanning (ME-Scan) by targeted high-throughput sequencing. *BMC Genomics* 11:
384 410.
- 385 Wu Z, Jiang Z, Li T, Xie C, Zhao L, Yang J, Ouyang S, Liu Y, Li T, Xie Z. 2021. Structural
386 variants in the Chinese population and their impact on phenotypes, diseases and
387 population adaptation. *Nat Commun* 12: 6501.
- 388 Xu L, Wang X, Lu X, Liang F, Liu Z, Zhang H, Li X, Tian S, Wang L, Wang Z. 2023.
389 Long-read sequencing identifies novel structural variations in colorectal cancer. *PLoS*
390 *Genet* 19: e1010514.
- 391 Yamazaki H, Inui Y, Yun CH, Guengerich FP, Shimada T. 1992. Cytochrome P450 2E1 and
392 2A6 enzymes as major catalysts for metabolic activation of N-nitrosodialkylamines
393 and tobacco-related nitrosamines in human liver microsomes. *Carcinogenesis* 13:
394 1789–1794.
- 395 Yoshida K, Gowers KHC, Lee-Six H, Chandrasekharan DP, Coorens T, Maughan EF, Beal K,
396 Menzies A, Millar FR, Anderson E, et al. 2020. Tobacco smoking and somatic
397 mutations in human bronchial epithelium. *Nature* 578: 266–272.
- 398 Zeng T, Liao H, Xia L, You S, Huang Y, Zhang J, Liu Y, Liu X, Xie D. 2025. Multisite
399 long-read sequencing reveals the early contributions of somatic structural variations to
400 HBV-related hepatocellular carcinoma tumorigenesis. *Genome Res*.
- 401 Zevin S, Benowitz NL. 1999. Drug interactions with tobacco smoking. An update. *Clin*
402 *Pharmacokinet* 36: 425–438.
- 403 Zhou Y, Hu Y, Yan X, Zheng Y, Liu S, Yao H. 2023. Smoking index and COPD duration as
404 potential risk factors for development of osteoporosis in patients with non-small cell
405 lung cancer – A retrospective case control study evaluated by CT Hounsfield unit.
406 *Heliyon* 9: e20885.
- 407 Zhu X, Zhang Y, Wang Y, Tian D, Belmont AS, Swedlow JR, Ma J. 2022. Nucleome
408 Browser: an integrative and multimodal data navigation platform for 4D Nucleome.
409 *Nat Methods* 19: 911–913.

410
411
412
413 **Fig. 1. Overview of somatic SV detection softwares evaluation.** (A).Schematic of the
414 SomaGauss-SV workflow, which consists of three main steps: (1) detection of candidate
415 somatic SVs, screening of somatic SVs through a merge method, (2)filtering of insertions
416 (INS) and deletions (DEL) using a Gaussian Mixture Model (GMM), and (3) filtering of SVs
417 using data from the Chinese normal population. (B).Bar chart shwoing the F1 scores for
418 different software. The height of the bars represents the average F1 score for each software,
419 and different shapes of points indicate the F1 scores for different cell lines. (C).Contour plot of

120 software evaluation. The X-axis represents precision, the Y-axis represents recall, the contours
 121 in the plot represent F1 scores, and software and cell lines are distinguished by the color and
 122 shape of the points, respectively. **(D)**. Venn diagram of software detection results in the
 123 HCC1395 cell line. Numbers indicate the number of SVs in the intersections between software,
 124 and red lines indicate the intersections involving SomaGauss-SV with other software.
 125 **(E)**. UPSet plot of true positive (TP) proportions for different software in the HCC1395 cell
 126 line. Yellow indicates TP, gray indicates false positives (FP), red boxes indicate the proportion
 127 of TP unique to each software, and black circles below indicate intersections between
 128 software.

129
 130
 131
 132
 133

134 **Fig. 2. Overview of pathological information and sequencing data.** **(A)**. Display of
 135 sampling sites for LSCC. There are 3 cases of supraglottic type, 11 cases of glottic type, and 1
 136 case of subglottic type. **(B)**. Pie chart of the proportion of sampling sites for LSCC. **(C)**. Pie
 137 chart of the proportion of TNM staging for LSCC patients. **(D)**. Pathological scatter plot for
 138 LSCC patients. The X-axis represents the smoking index of patients (calculated as: number of
 139 cigarettes per day \times smoking years), the Y-axis represents patient age, and the shape of the
 140 points indicates TNM staging. **(E)**. Reads length distribution plot of nanopore sequencing. The
 141 upper part shows the distribution of average read lengths (bp) for tumor (pink) and blood
 142 (blue), and the lower part shows the N50 distribution for tumor and blood. **(F)**. Box plot of
 143 read Mapping rates. Blue represents blood, and pink represents tumor.

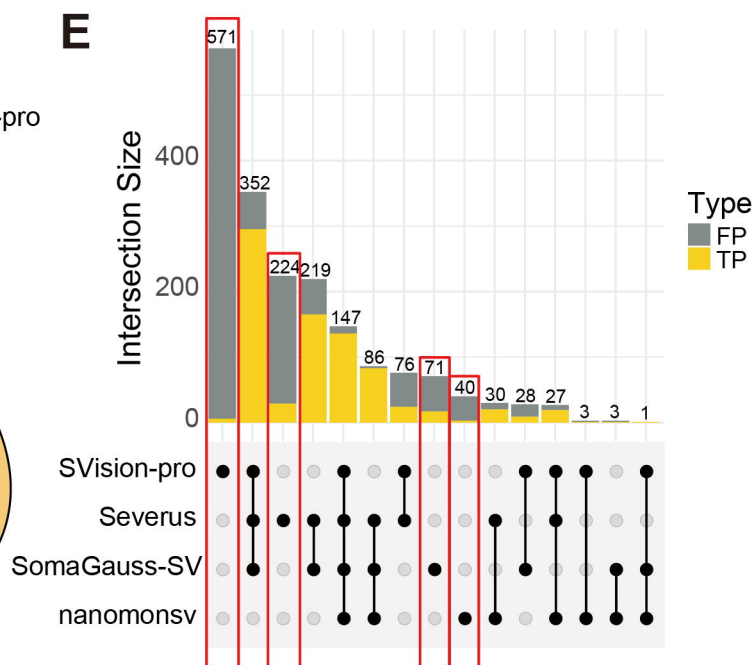
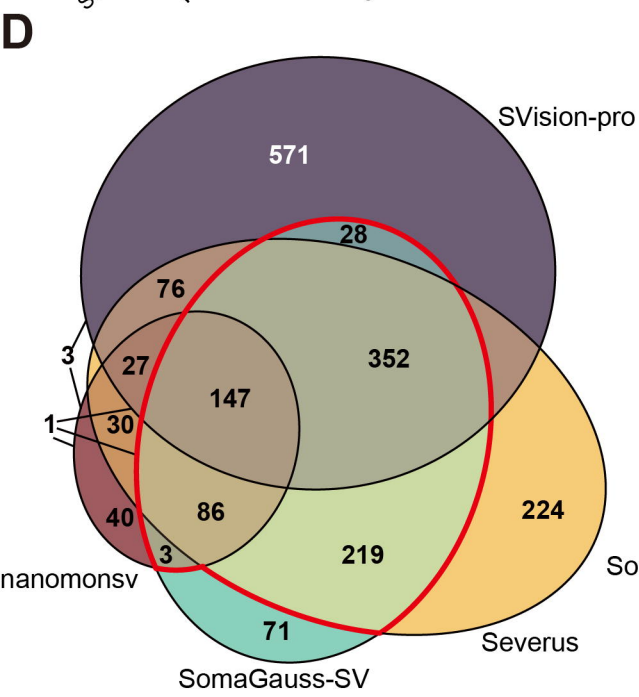
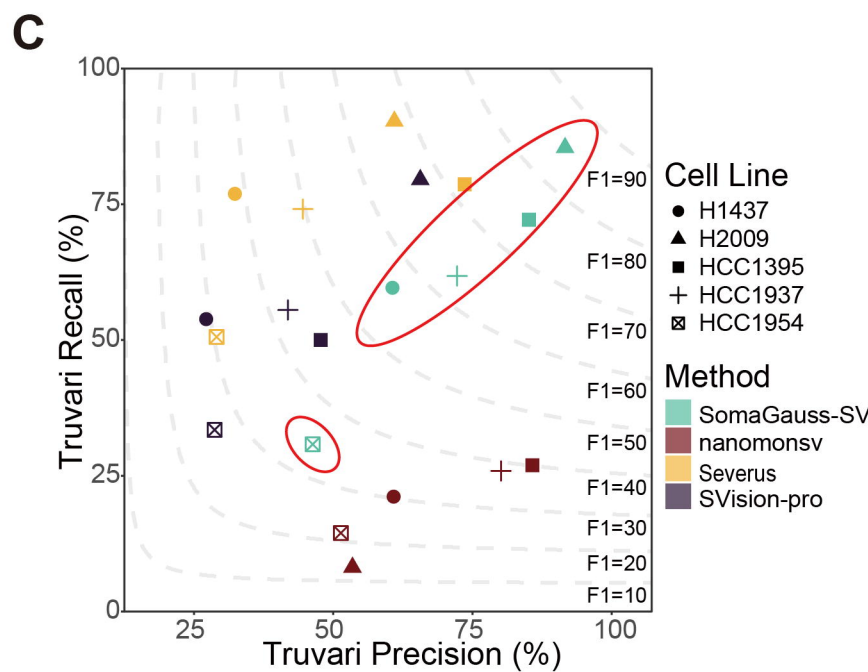
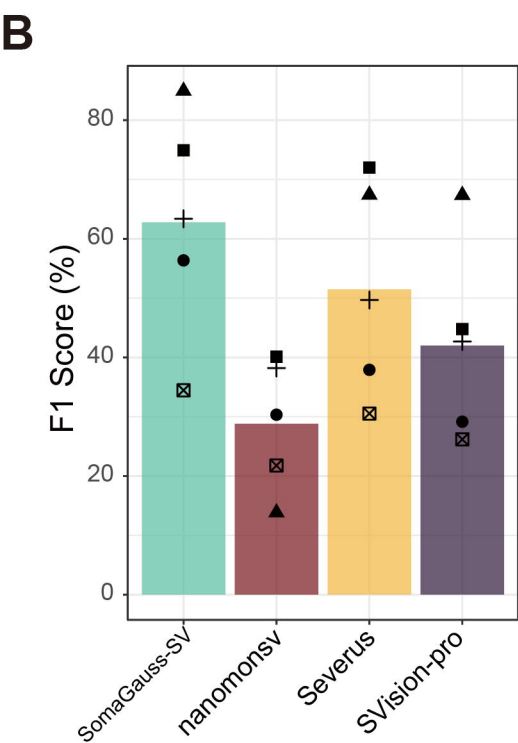
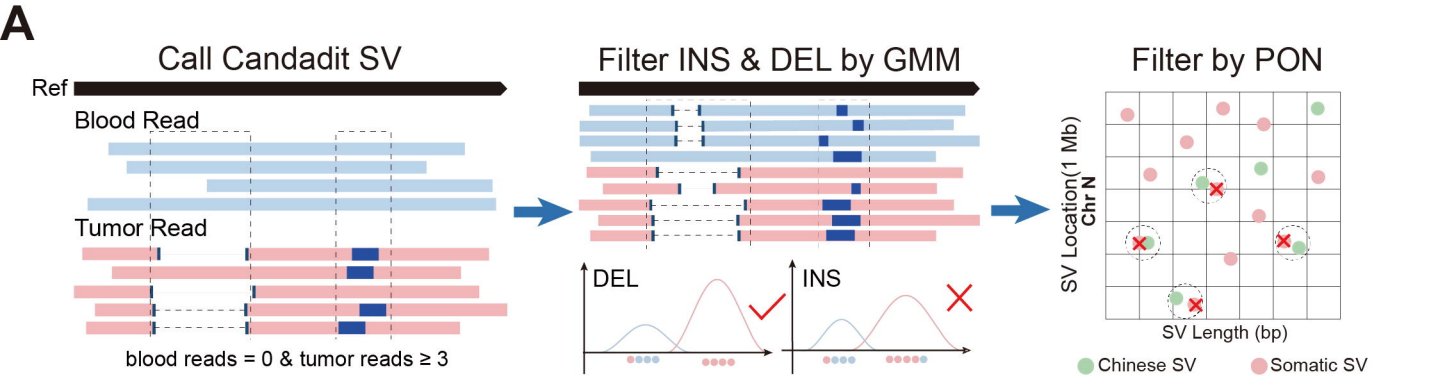
144
 145 **Fig. 3. Relationship between somatic DEL and tobacco exposure.** **(A)**. Stacked bar chart of
 146 somatic structural variation counts and smoking index distribution. The upper stacked bar
 147 chart shows the number of different SVs for each sample, and the lower bubble chart shows
 148 the smoking index for corresponding samples, where higher indices are represented by redder
 149 colors. **(B)**. Correlation plot of somatic DEL counts and smoking index. The X-axis represents
 150 the smoking index, and the Y-axis represents the number of DEL, with degrees of freedom,
 151 p -value, and Pearson correlation coefficient marked on the plot. **(C)**. Reactome pathway
 152 enrichment plot for genes affected by common DEL in multiple samples from smokers
 153 (q -value < 0.05). The X-axis represents $-\log_{10}(q\text{-value})$, and the Y-axis represents the names of
 154 enriched pathways, with immune-related pathways highlighted in red. The size of the circles
 155 indicates the number of enriched genes. **(D)**. Browser plot of somatic DEL
 156 (Chr6:32,477,000-32,477,500) illustrating the following epigenomic signals: sample DEL
 157 location, ENCODE-annotated gene transcript, and Candidate *cis*-regulatory Elements (cCREs)
 158 predicted by ENCODE. **(E)**. Browser plot of somatic DEL (Chr7:100,953,500-100,954,000)
 159 illustrating the following epigenomic signals: sample DEL location, ENCODE-annotated gene
 160 transcript, and Simple Tandem Repeats identified by TRF.

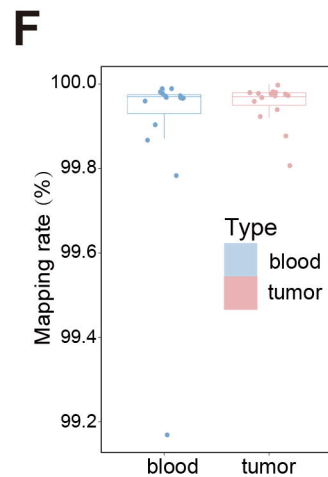
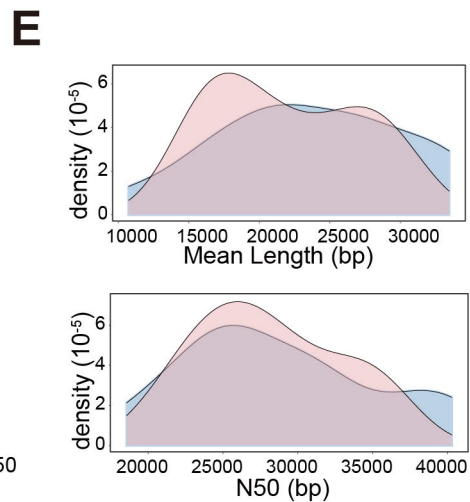
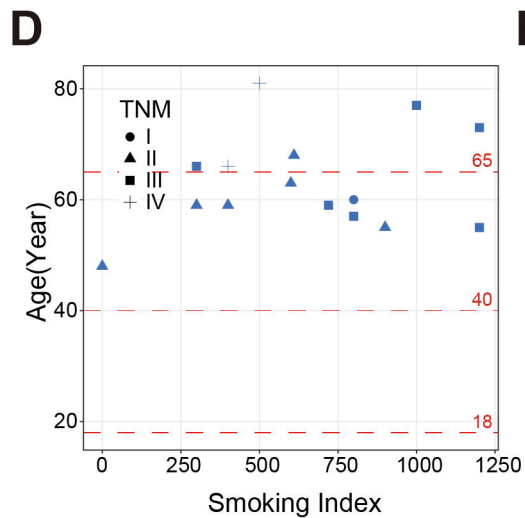
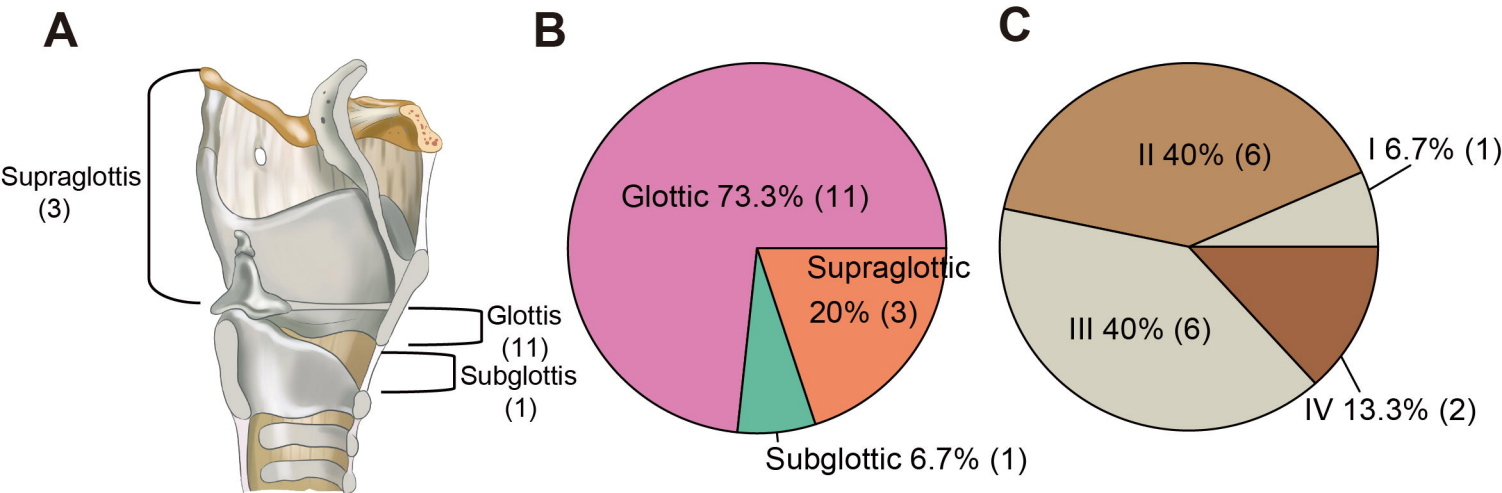
161
 162 **Fig. 4. Overview of SV length and site enrichment.** **(A)**. Length distribution plot for different
 163 SV types. From top to bottom are DEL, INS, DUP, and INV, with red dashed lines marking the
 164 peaks for each SV type. **(B)**. Stacked bar chart of sequence annotation proportions at DEL
 165 peaks. **(C)**. Stacked bar chart of sequence annotation proportions at INS peaks. **(D)**. Heatmap of
 166 somatic SVs site enrichment in repeat regions. Each row represents a different SV type, and

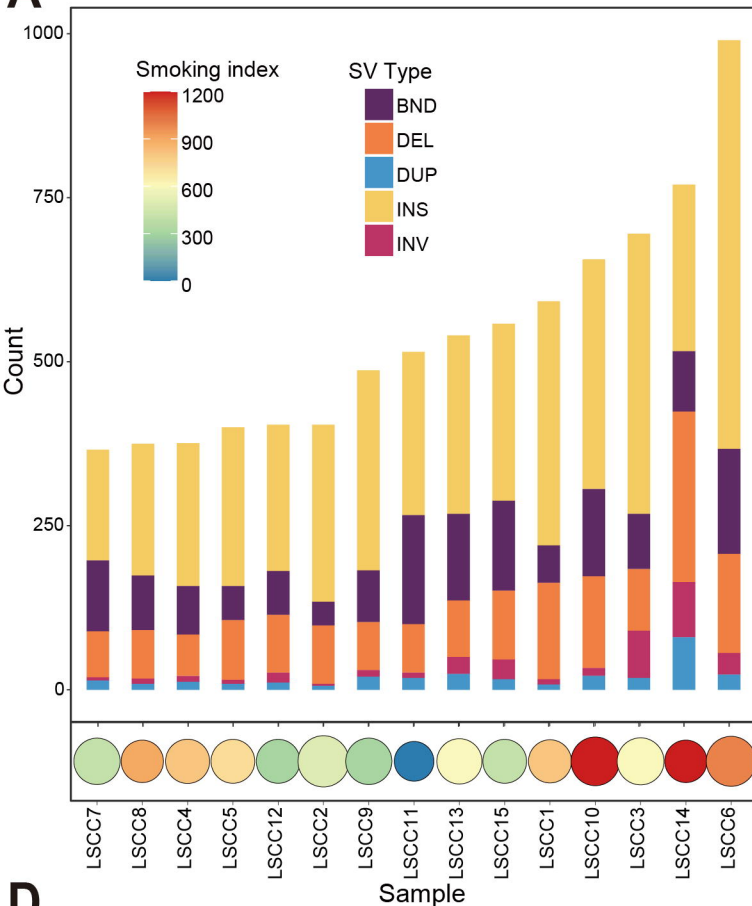
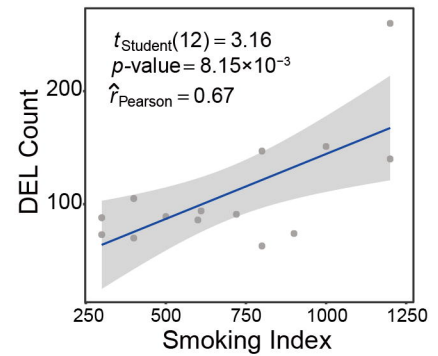
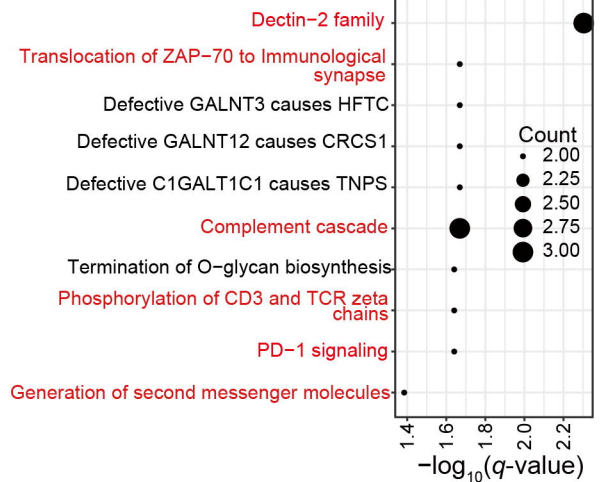
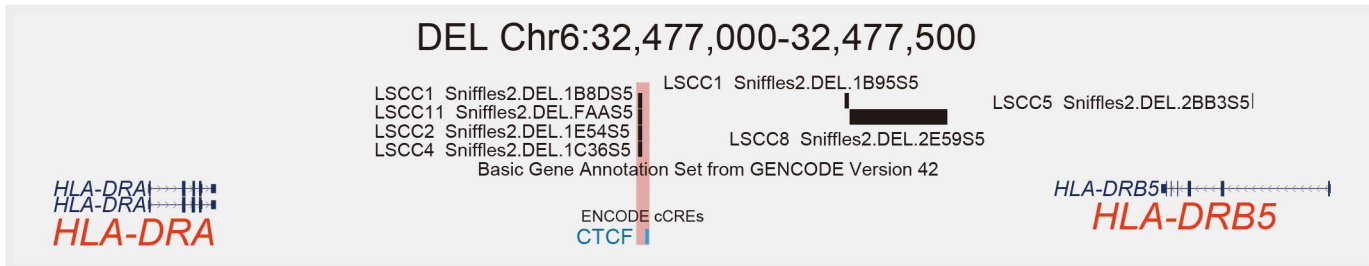
167 each column represents a different repeat region type, with color intensity indicating \log_2 (Fold
 168 Enrichment). Repeat regions are categorized into TE (transposable elements region), Repeat
 169 (repeat region), and de novo (non repetitive regions). **(E)**.Bar chart of somatic SVs site
 170 enrichment for gene elements. The X-axis represents \log_2 (Fold Enrichment), the Y-axis
 171 represents gene element names, and facets represent different SV types (Fisher's exact test,
 172 * p -value < 0.05, ** p -value < 0.01, *** p -value < 0.001).
 173
 174

175 **Fig. 5. Circos plot of genome-wide sample enrichment.** The tracks in the figure represent
 176 the enrichment of different SVs types in samples, from the innermost to the outermost being
 177 INV, BND, DUP, INS, and DEL. The plot highlights regions with the top 2 sample counts for
 178 INS and DEL that may affect genes, the top 3 regions for INV sample counts, and regions
 179 intersecting with those reported in previous studies.
 180
 181

182 **Fig. 6. Overview of somatic SRE shared by 66% of samples.** **(A)**.Oncoplot plot of INS
 183 sample enrichment regions. Rows represent samples, and columns represent regions. **(B)**.IGV
 184 plot of somatic SRE (LSCC 13, Chr1:224,009,042-224,018,679). Light blue represents reads
 185 from blood samples, pink represents reads from tumor samples, and dark blue indicates the
 186 Simple repeat region with the AATGG motif. **(C)**.Browser plot of somatic SRE
 187 (Chr1:224,009,042-224,018,679) illustrating the following epigenomic signals: a stacked bar
 188 chart of SRE positions and inserted sequences, with different colors representing different
 189 bases; the plot concludes with the AATGG motif of the repeat, the ENCODE-annotated gene
 190 transcript, and Candidate *cis*-regulatory Elements (cCREs) predicted by ENCODE. Also
 191 included are H3K4me1, H3K4me3, and H3K27ac marks predicted by ENCODE. In the upper
 192 left corner, the predicted 4D nucleosome structure around the LSCC-enriched SRE, *TP53BP2*,
 193 and *FBXO28* is displayed. **(D)**.Horizontal gel electrophoresis images 1 showing DNA
 194 fragments from blood and tumor tissues of LSCC samples sequenced using nanopore
 195 technology. Each sample's blood and tumor are labeled separately, and green dots and red lines
 196 mark the longer bands in the tumor relative to the blood. **(E)**.Horizontal gel electrophoresis
 197 images 2 showing DNA fragments from blood and tumor tissues of LSCC samples sequenced
 198 using nanopore technology. Each sample's blood and tumor are labeled separately, and green
 199 dots and red lines mark the longer bands in the tumor relative to the blood. **(F)**.Relative
 200 expression levels of *TP53BP2* and *FBXO28* determined by QCPCR in SRE and No SRE
 201 Clinical Samples (*t*-test). **(G)**.Relative expression levels of *TP53BP2* and *FBXO28* determined
 202 by QCPCR in BEAS 2B and SNU899 cell lines (*t*-test, * p -value < 0.05, ** p -value < 0.01,
 203 *** p -value < 0.001).
 204
 205





A**B****C****D****E**