



GENOME RESEARCH

Multisite long-read sequencing reveals the early contribution of somatic structural variations to HBV-related hepatocellular carcinoma tumorigenesis

Tianfu Zeng, Haotian Liao, Lin Xia, et al.

Genome Res. published online March 4, 2025

Access the most recent version at doi:[10.1101/gr.279617.124](https://doi.org/10.1101/gr.279617.124)

P<P Published online March 4, 2025 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-NonCommercial 4.0 International license), as described at <http://creativecommons.org/licenses/by-nc/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](#).



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

Published by Cold Spring Harbor Laboratory Press

1 **Multisite Long-Read Sequencing Reveals the Early Contribution of Somatic Structural**
2 **Variations to HBV-related Hepatocellular Carcinoma Tumorigenesis**

3

4 Tianfu Zeng^{1#}, Haotian Liao^{2#}, Lin Xia^{1#}, Siyao You¹, Yanqun Huang¹, Jiaxun Zhang¹, Yahui
5 Liu¹, Xuyan Liu¹, Dan Xie^{1*}

6

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

10 2. Division of Liver Surgery, Department of General Surgery and Laboratory of Liver Surgery,
11 and State Key Laboratory of Biotherapy and Collaborative Innovation Center of Biotherapy, West
12 China Hospital, Sichuan University, Chengdu, Sichuan, 610041, China.

13 # These authors equally contributed to this work.

14

15 * Corresponding author

16 Dan Xie

17 Frontiers Science Center for Disease-related Molecular Network, West China Hospital, Sichuan
18 University, Chengdu, Sichuan, China.

19 Email: danxie@scu.edu.cn

20

21 **Running title**

22 Somatic Variations in HCC by Long-Read Sequencing

23

24

25 Abstract

26 Somatic structural variations (SVs) represent a critical category of genomic mutations in
27 hepatocellular carcinoma (HCC). However, the accurate identification of somatic SVs using
28 short-read high-throughput sequencing (HTS) is challenging. Here, we applied long-read
29 nanopore sequencing and multisite sampling in a cohort of 42 samples from five patients. We
30 discovered a prominent presence of somatic SVs in adjacent nontumor tissues, which significantly
31 differed from somatic single nucleotide variants (SNVs) and copy number variations (CNVs). The
32 types of SVs were markedly different between adjacent nontumor and tumor tissues, with somatic
33 insertions (INSS) and deletions (DELS) serving as early genomic alterations associated with HCC.
34 Notably, hepatitis B virus (HBV) DNA integration frequently resulted in the generation of
35 somatic SVs, particularly inducing inter-chromosomal translocations. While HBV DNA
36 integration into the liver genome occurs randomly, multisite shared HBV-induced SVs are
37 implicated as early driving events in the pathogenesis of HCC. Long-read RNA sequencing
38 revealed that some HBV-induced SVs impact cancer-associated genes, with translocations being
39 capable of inducing the formation of fusion genes. These findings enhance our understanding of
40 somatic SVs in HCC and their role in early tumorigenesis.

41 Introduction

42 As the most prevalent form of primary liver cancer, hepatocellular carcinoma (HCC), which
43 mostly forms after years of chronic liver disease, against a background of severe liver scarring
44 and typically cirrhosis, has been ranked as one of the leading causes of cancer-related death
45 worldwide (Villanueva 2019; Muller et al. 2020). More than half of HCC cases globally occur in
46 China, where chronic hepatitis B virus (HBV) infection is the major etiological factor, accounting
47 for over 60% of cases (Llovet et al. 2021; Bray et al. 2024; Chen et al. 2024a). Somatic DNA
48 alterations are known as pivotal drivers of HCC tumorigenesis and progression. Mutations in the

49 TERT promoter are the most prevalent (~60%) genetic alterations in HCC(Schulze et al. 2016),
50 and the TERT promoter is a recurrent insertion site of HBV DNA. Other recurrently mutated
51 genes primarily harbor somatic single nucleotide variants (SNVs) within their coding regions,
52 including TP53 (~30%), CTNNB1 (~30%), and ARID1A (~10%), which affect the cell cycle,
53 WNT signaling, or chromatin remodeling (Khemlina et al. 2017). Nevertheless, these most
54 prevalent mutations remain undruggable at present (Zucman-Rossi et al. 2015), and there are no
55 mutations available in clinical practice to predict therapeutic response, highlighting the
56 incomplete understanding of the mutational landscape of HCC.

57 Structural variants (SVs) are large genomic alterations (>50 bp), distinct from small variants
58 like SNVs and short insertions and deletions (indels), as they often arise from different
59 mechanisms (Abyzov et al. 2015). SVs are generally defined as insertions (INS), deletions (DEL),
60 duplications (DUP), inversions (INV), and translocations (TRA), and are known to play an
61 important role in cancer pathogenesis(George et al. 2015; Waddell et al. 2015). Over the past
62 decade, utilizing short-read high-throughput sequencing (HTS) technologies, studies have
63 revealed that somatic SVs can drive malignant phenotypes by altering the expression or function
64 of oncogenes (The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020;
65 Cosenza et al. 2022). Additionally, certain somatic SVs within tumors have been used to make
66 therapeutic decisions, such as the structural rearrangements involving the ALK gene in non-small
67 cell lung cancer and the translocation leading to BCR-ABL1 fusion in chronic myeloid leukemia
68 (Soda et al. 2007; Hochhaus et al. 2017a; Hochhaus et al. 2017b). However, limited by the short
69 read lengths of HTS, the current understanding of somatic SVs in HCC is still incomplete.

70 Long-read sequencing (LRS) technology, represented by Oxford Nanopore Technologies
71 (ONT) and Pacific Bioscience (PacBio), can produce continuous reads longer than 10,000 base
72 pairs. These reads can completely span across genomic repetitive segments and complex genomic
73 variant regions. Thus, long-read sequencing is progressively being applied to studies on genomic

74 SVs in tumors. For instance, utilizing nanopore sequencing, a complex SV named “cancerous
75 local copy-number lesions” (CLCLs) has been identified within lung cancer, and the precise
76 junctions of these complex SVs are difficult to identify using short-read sequencing (Sakamoto et
77 al. 2020). While previous studies have revealed distinct mechanisms in the generation of germline
78 and somatic SVs in liver cancer by analyzing complete breakpoint sequences derived from long
79 reads, the genome-wide landscape of somatic SV characteristics within HCC has yet to be fully
80 elucidated (Fujimoto et al. 2021). A recent study revealed notable occurrences of somatic SVs
81 within cirrhotic tissues (Brunner et al. 2019). Given that cirrhosis from any etiology is the
82 strongest risk factor for HCC, this further underscores the importance of deciphering somatic SVs
83 in elucidating the pathogenesis of HCC.

84 **Results**

85 **Detection of somatic SVs using long-read sequencing.**

86 We performed multisite sampling of tumor tissues and adjacent nontumor tissues from five
87 patients diagnosed with HBV-positive HCC. The nontumor tissues were categorized as proximal
88 nontumor (located within proximal 1 cm of the tumor) and distal nontumor (>1 cm distal to the
89 tumor), with both exhibiting fibrosis. Blood samples collected before surgical resection were used
90 as matched normal controls. In total, twenty-seven tumor tissue samples, 10 matched nontumor
91 liver samples (comprising 6 proximal and 4 distal), and 5 matched blood samples were obtained
92 from the 5 patients (Fig. 1A). Detailed clinical information of included patients is presented in
93 Supplemental Table S1.

94 We generated long-read whole-genome sequencing data from all samples using the
95 PromethION (Oxford Nanopore Technologies) platform. To obtain high quality sequencing data,
96 we excluded reads with lengths less than 1 kilobase (kb) and mean base quality lower than 7. The
97 mean N50 of long reads was 22.41 kb (19.22 kb ~ 31.02 kb) after quality control (Supplemental

98 Fig. S1A, Supplemental Table S2). The average sequencing depth reached 36.72X (30.50X ~
99 44.22X) for tumor and nontumor tissues, 22.19X (19.84X ~ 25.34X) for blood samples
100 (Supplemental Fig. S1B). Considering the relatively high single-base error rate of long-read
101 sequencing, short-read whole-genome sequencing using the Illumina NGS platform was
102 additionally conducted on all samples to facilitate the detection of somatic SNVs.

103 Utilizing long-read sequencing, we developed a rigorous in-house bioinformatics pipeline to
104 identify somatic SVs and HBV DNA integration events, facilitating the characterization of SV
105 types associated with HBV integration (Supplemental Fig. S1C). Clean long reads were aligned
106 against a custom reference genome consisting of the human genome and 21 HBV genomes using
107 minimap2. Subsequently, the alignment results from paired samples were merged. The widely
108 used variant caller Sniffles2 was used to cluster the long reads supporting the same SV. In
109 instances where supporting reads for SVs span both the HBV and human genomes in alignments,
110 SV types were re-inferred based on the alignment patterns of human sequences flanking the
111 breakpoints. After conducting quality control and manual inspection, we identified an average of
112 253 somatic SVs across all samples (range from 122 to 405, a median of 252, Fig. 1B). This count
113 notably exceeds previous reports based on short-read sequencing, where only approximately 20%
114 of HCC cases displayed somatic SV counts greater than 200 (Fujimoto et al. 2016; Li et al.
115 2020b). These results underscore the superior sensitivity of long-read sequencing in the detection
116 of SVs (Aganezov et al. 2020; Xu et al. 2023).

117 The number of somatic SVs rapidly decreases as the SV length increases, consistent with
118 previous studies (Beyter et al. 2021). We identified two distinct peaks, located at approximately
119 170 bp and 320 bp (Fig. 1C). The predominant types of SVs at these peaks were INS and DEL,
120 accounting for 98.47% and 99.17%, respectively. To further characterize these sequences, we
121 performed sequence annotation using RepeatMasker. The SVs around 170 bp were predominantly
122 composed of satellite elements (Supplemental Fig. S2A), which are involved in the maintenance

123 of chromosomal integrity and have been found to be overexpressed in various cancer types (Ting
124 et al. 2011; Ho et al. 2017). In contrast, the sequences near 320 bp are primarily composed of
125 SINEs (Supplemental Fig. S2B), which are associated with SV mechanisms associated with the
126 mobilization of active transposable elements (Kolomietz et al. 2002; Robberecht et al. 2013).
127 Notably, somatic INSs and DELs were significantly shorter in length compared to somatic
128 duplications (DUPs) and inversions (INVs), with the majority of INVs and DUPs exceeding
129 10,000 bp (Wilcoxon rank-sum test, p value $< 2.2 \times 10^{-16}$, Fig. 1D).

130 The mean variant allele frequency (VAF) of somatic SVs across all samples was 0.14.
131 Notably, the VAF of somatic SVs in the adjacent nontumor samples from each patient was
132 significantly lower than that in the tumor samples (Supplemental Fig. S3A), a phenomenon also
133 observed in somatic SNVs in this project (Supplemental Fig. S3B) and prior studies (Huang et al.
134 2017; Strandgaard et al. 2020). This indicates that the fibrotic liver tissue adjacent to the tumor
135 does not reflect a normal genomic state and has undergone accumulation of low-frequency
136 somatic mutations.

137 **Somatic SVs were prominent in adjacent nontumor tissues.**

138 A recent study revealed that cirrhotic livers exhibit a higher mutation burden compared to
139 normal livers, with the most significant divergence observed in SVs (Brunner et al. 2019). To
140 explore the process of somatic mutation accumulation in HCC, we compared the burden of
141 mutation between adjacent nontumor tissues and tumor tissues. The genome was segmented into
142 500-bp bins, and the number of bins harboring each mutation type was determined. We observed
143 a significantly lower count of somatic SNVs in the adjacent nontumor liver samples compared to
144 the tumor samples (Wilcoxon rank-sum test, p value $= 5.7 \times 10^{-9}$, Fig. 2A,). Besides, there were
145 also fewer somatic CNV events in the adjacent nontumor liver samples, with 50% of these
146 samples exhibiting no CNVs. However, somatic SVs were prominent in adjacent nontumor
147 tissues, with no significant difference observed when compared to tumor tissues (Fig. 2A). These

148 results indicate that somatic SVs follow a distinct development process compared to somatic
149 CNVs and somatic SNVs during the progression of HCC.

150 To compare the mutation profiles across different sampling sites, hierarchical clustering of
151 all samples was performed on mutation matrices of 500 bp for both somatic SNVs and SVs across
152 all samples. The clustering analysis of somatic SNVs revealed that tumor samples grouped
153 according to their patient origin, whereas nontumor samples formed a distinct outgroup, separate
154 from the tumor cohort (Fig. 2B). This separation likely arises from the significantly lower number
155 of somatic SNVs in adjacent nontumor tissues compared to tumor tissues (an average of 851
156 versus 24,234), leading to their clustering based on their mutual dissimilarity to the tumors. In
157 contrast, clustering analysis based on somatic SVs revealed a distinct pattern, with somatic SVs
158 clustering according to patients, and both tumor and nontumor samples clustered together within
159 individual (Fig. 2B). Due to the absence of somatic CNVs in half of the adjacent nontumor
160 samples, we directly compared the mutational profiles of somatic CNVs across different sampling
161 sites within the same patient (Supplemental Fig. S4A-E). Notably, there was a high degree of
162 consistency in CNV changes among various tumor sampling sites. For instance, all 5 tumor
163 sampling sites from HCC8 exhibited amplification of the entire chromosome 2 and 5, despite the
164 absence of CNV alterations in the adjacent nontumor samples (Supplemental Fig. S4A). These
165 findings underscore the similarity of nontumor tissues to tumors in terms of somatic SVs, while
166 highlighting a marked distinction in terms of somatic CNVs and SNVs.

167 To further examine the similarities in mutations among samples in detail, we performed
168 pairwise comparisons across different sampling sites (Supplemental Fig. S5A-E). The proportion
169 of shared mutations between adjacent nontumor and tumor tissues exhibited significant
170 differences among the categories of somatic SVs, SNVs, and CNVs. Specifically, in both
171 nontumor and tumor samples, the proportion of shared somatic SVs between adjacent nontumor
172 and tumor tissues was significantly higher than that of somatic SNVs (p value = 5.3×10^{-8} and

173 2.2×10^{-16} , *t*-test, Supplemental Fig. S6A). Notably, approximately 12.98% of somatic SVs
174 detected in the tumors were also present in the adjacent nontumor tissues, whereas somatic SNVs
175 accounted for only 0.45%. There were no shared somatic CNVs between adjacent nontumor and
176 tumor tissues. This observation aligns with the current understanding of tumor evolution,
177 indicating that the acquisition of CNVs occurs late in clonal expansion (Li et al. 2020a),
178 representing a consequence of the final malignant transformation of genomic stability.

179 To further understand tumor evolution, we constructed phylogenetic trees for each patient
180 based on the mutation matrix of 500bp windows, and branch evolution patterns correlated with
181 tumor heterogeneity were observed in all patients (Supplemental Fig. S7A). Notable branches
182 shared between tumor and adjacent nontumor tissues were observed on the phylogenetic tree
183 based on somatic SVs, whereas no such pattern was evident in the case of somatic SNVs
184 (Supplemental Fig. S7B). These results suggest that somatic SVs might serve as the primary
185 mutation type in the context of chronic liver disease during the early stages of HCC progression.

186 **HCC was characterized by somatic INS and DEL at an early stage.**

187 The adjacent nontumor tissue in HCC is not entirely normal (Carlessi et al. 2023; Chen et al.
188 2024b). A transcriptomic study examining HCC alongside seven other tumor types has
189 demonstrated that adjacent nontumor tissues exhibit a unique intermediate state between health
190 and tumor (Aran et al. 2017). To investigate the role of shared SVs between adjacent nontumor
191 and tumor tissues in HCC pathogenesis, we categorized somatic SVs into three distinct categories
192 within individual: shared (present in both adjacent nontumor and tumor tissues), T-specific
193 (exclusive to tumor tissues), and N-specific (exclusive to adjacent nontumor tissues). For each
194 patient, we found that the shared somatic SV categories exclusively comprised DELs and INSs,
195 whereas translocations (TRAs) were unique to the T-specific category (Fig. 3A). Although DUPs
196 and INVs were present in both T-specific and N-specific categories, the relative proportions
197 varied significantly among patients. For instance, HCC9 exhibited a high proportion of DUPs,

198 whereas HCC13 had a higher proportion of INVs. Notably, while HCC9 lacked DUPs in the
199 adjacent nontumor tissue, there was a high consistency among DUPs across different tumor sites,
200 with approximately 68.46% being mutually shared among samples (Supplemental Fig. S8A). In
201 contrast, although INVs were present in the adjacent nontumor tissue of HCC13, no shared INVs
202 were identified between the tumor and adjacent tissues (Supplemental Fig. S8B). The
203 accumulation of INVs in the tumor of HCC13 was localized to Chromosome 9, suggesting that
204 this chromosome may have experienced a catastrophic event, such as chromothripsis
205 (Supplemental Fig. S8C). Moreover, the percentage of DUPs, INVs, and TRAs within the T-
206 specific category exhibited a significantly higher prevalence than those in the N-specific category
207 (chi-square test, p value < 0.01). These specific types of SVs are frequently associated with
208 chromosomal rearrangements in tumor evolution. These results indicate a significant difference in
209 the types of SVs between adjacent nontumor tissues and tumors, suggesting that the tumor
210 genome has undergone more complex genomic alterations.

211 We further investigated the differences in somatic INs and DELs between the shared and
212 N/T-specific categories. Notably, the lengths of somatic INs and DELs shared between adjacent
213 nontumor and tumor samples were significantly greater than those exclusive to either category
214 (Fig. 3B, Wilcoxon rank-sum test, $< 2.2 \times 10^{-16}$). Additionally, the shared somatic INs and
215 DELs exhibited a higher prevalence of complex repetitive sequences and satellite elements, while
216 non-repetitive and simple repeat sequences were less represented (Supplemental Fig. S9A,B).

217 Next, we analyzed the landscape of shared somatic SVs between tumor and nontumor
218 samples within individual. Notably, we identified a subset of recurrent somatic SVs, present in
219 two or more patients, affecting a total of 42 genes (Fig. 3C). These include tumor-related genes
220 such as *PGA5* and *DPP6* (Shen et al. 2020; Choy et al. 2021; Munkhjargal et al. 2023), as well as
221 newly identified candidate genes like *EVAIC*, which is present in four out of five patients. *EVAIC*
222 is a membrane protein-encoding gene that plays a role in inflammatory and immunobiological

223 processes (Hu and Qu 2021). The SV observed on the *EVAIC* gene manifests as a deletion event
224 spanning approximately 160 base pairs in length. Moreover, this deletion is consistently located in
225 close proximity across different samples, with breakpoint variations confined within 10 bp. Short-
226 read sequencing in the samples failed to detect this deletion, primarily due to the tandem-repeated
227 *Alu* sequences present within the corresponding region of the reference genome (Supplemental
228 Fig. S10A,B).

229 To assess the functional impact of these recurrent somatic SVs, we conducted long-read
230 RNA sequencing on samples HCC10_T7 and HCC10_N1 (Supplemental Table S3). Transcript
231 loss was observed for *EVAIC*, *GSTM1*, and *GSM2* in samples harboring DELs (Supplemental Fig.
232 11A,B), despite the DEL in *EVAIC* localized exclusively to an intronic region. Additionally,
233 qPCR validation revealed that the expression level of the *EVAIC* gene in samples with DEL was
234 significantly lower than in samples without DEL (Supplemental Fig. S11C). Furthermore, based
235 on gene expression and clinical data obtained from the Cancer Genome Atlas (TCGA) project
236 (The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium 2020), survival analysis
237 showed that the low expression of *EVAIC* was associated with poor survival of patients
238 (Supplemental Fig. S11D, $P = 0.0055$). These findings indicate that somatic SVs shared between
239 tumors and adjacent nontumor tissues have the potential to induce abnormal expression of genes
240 associated with cancer.

241

242 **HBV integrations induced the formation of SVs.**

243 Viral DNA integration into the host genome can lead to the malignant transformation of
244 tumors, resulting in an increased burden of mutations near the integration breakpoints (Zapatka et
245 al. 2020). HBV infection is frequently linked to proliferative HCC, which is characterized by
246 chromosomal instability (Villanueva 2019). To further explore the connection between somatic

247 SVs and HBV integrations, we developed a novel pipeline to comprehensively identify HBV
248 integration events using long-read sequencing. We primarily identified two categories of long
249 reads for detecting HBV integration events. First, chimeric reads that mapped simultaneously to
250 the human genome and the HBV genome. Secondly, due to disparities in minimap2 alignment,
251 certain HBV sequences were either mapped as insertions or remained unmapped due to clipping
252 (Supplemental Fig. S12A). Basing on these two kinds of reads, we identified an average of 10
253 HBV integration breakpoints in the tumor samples (range from 0 to 37, median: 8) and 5 in the
254 adjacent nontumor samples (range from 0 to 15, median: 4), which was significantly more than
255 previous short-read datasets (Sung et al. 2012; Alvarez et al. 2021). Furthermore, 61.38% of HBV
256 DNA integrated into the human genome exceeded 500 bp in length. (Supplemental Fig. S12B).
257 In short-read sequencing, aligning these insert sequences larger than the DNA library fragment
258 size is challenging, and complete insert sequences cannot be detected (Craven et al. 2022; Rajaby
259 et al. 2023).

260 Within the 335 identified HBV integration breakpoints, 85.67% (287/335) were derived from
261 HBV integration events fully captured by long reads, characterized by HBV integration sequences
262 flanked by human sequences on both sides within the same read. These HBV integration events
263 provided an opportunity for further investigation into the specific patterns of HBV integration into
264 the HCC genome. Based on the alignment positions, orientations and order of human sequences
265 flanking the HBV sequences, we inferred the integration patterns of HBV (Supplemental Fig.
266 S12C-G). We observed that simple insertion of HBV accounted for only 20.20% of events, while
267 the majority of HBV integrations resulting in chromosomal rearrangements within the HCC
268 genome. Among these genomic alterations, translocations were the predominant type, constituting
269 44.80%, followed by inversions (22.20%), duplications (8.62%), deletions (4.12%), and complex
270 SVs (1.78%). Notably, the complex SVs identified were all duplication-inversion events,
271 originating from five distinct tumor sampling sites in HCC9. These rearrangements were

272 characterized by HBV integration resulting in a fold-back inversion, with human sequences on
273 both sides of the integration site being inverted and duplicated (Fig. 4A). Moreover, a 76.23 Mb
274 copy number amplification (Chr7:124,879,000-48,644,500) was observed 779 bp downstream of
275 the HBV integration site, which is consistent with the breakage-fusion-bridge (BFB) inversion
276 model.

277 Next, we further explored the relationship between HBV integration breakpoints and CNV
278 boundaries. We found significant differences in the proximity of various HBV-induced SV types
279 to CNV boundaries. INs were located further from CNV boundaries, with all IN breakpoints
280 situated more than 100 kb away from the CNV borders. In contrast, except for five cases of
281 complex (duplication-inversion) events, TRAs and DUPs were found closer to the CNV
282 boundaries, with 6.4% of TRA breakpoints located within 1 kb of the CNV boundary
283 (Supplemental Fig. S13A). For instance, within the HBV integration event observed in
284 HCC13_T3, approximately 2914 base pairs (bp) of HBV sequence integrated into the HCC
285 genome, giving rise to a translocation event that linked Chr8:31,957,777 and Chr22:50,545,277.
286 Notably, a CNV is located 214 bp downstream of the TRA breakpoint on Chromosome 8. This
287 CNV corresponds to a 5.36 Mb gain on Chromosome 8p12 (Chr8: 31,957,991–37,314,348), and
288 this CNV region harbors several tumor-associated genes, such as *NRG1*, *UNC5D*, and *KCNU1*
289 (Repana et al. 2019). (Fig. 4B). These findings indicate that HBV-induced SVs may play a role in
290 CNV formation and directly affect critical oncogenes involved in tumorigenesis.

291 In addition, in comparison to a random genomic background, we found a significant
292 enrichment of somatic SVs in proximity to HBV integration sites (two-sided binomial test, p
293 value < 0.01 , Supplemental Fig. S13B). The above findings suggested that the HBV integration
294 was closely related to the formation of SVs within the HCC genome.

295

296 Multisite shared HBV-induced SVs emerge as distinctive signatures of HCC.

297 Based on the HBV-induced SVs identified above, we explored their breakpoint location
298 characteristics within the HCC genome. Among all samples, it was observed that 56.2% of
299 breakpoints were shared across multiple samples (Supplemental Fig. S15A). It is noteworthy that
300 these breakpoints were shared among tumor samples, with no shared breakpoints detected in the
301 adjacent nontumor tissues (Supplemental Fig. S15B). Specifically, we observed a pronounced
302 aggregation of HBV-induced SV breakpoints within tumors sampled from diverse regions of the
303 same patient, while HBV-induced SV breakpoints in adjacent nontumor tissues exhibited a
304 stochastic distribution (Fig. 5A). Furthermore, annotation of the HBV genome revealed that the
305 HBV-induced SVs shared across the five samples consistently harbored similar HBV integration
306 sequences, with an average sequence similarity of 96.78% (Fig. 5B). These integration sequences
307 were derived from proximal loci within the HBV genome, with an average breakpoint position
308 difference of 35 bp. (Supplemental Fig. S14A). For instance, all five tumor sampling sites in
309 HCC9 revealed a consistent translocation linked Chr2:74,273,672 and Chr6:70,300,819, with the
310 integrated HBV sequences derived from the HBV X gene (HBV:1544-1481). These findings
311 suggest that HBV DNA integration into the liver occurs randomly, whereas integration at specific
312 loci may lead to hepatocellular carcinogenesis.

313 Given the PCR-free nature of nanopore sequencing, each read is likely to represent a distinct
314 single-cell origin. The probability of the same HBV DNA integration event occurring
315 independently in multiple cells is minimal. Therefore, the shared HBV-induced SVs observed
316 across different tumor sampling sites may indicate a common clonal origin for these tumor
317 regions. To elucidate the clonal characteristics, we examined the mutational timing of the shared
318 HBV-induced SVs within individual. Variant allele frequency (VAF) analysis revealed a
319 significantly higher VAF of HBV-induced somatic SVs that shared within tumors compared to
320 those appearing independently (Supplemental Fig. S15C, Wilcoxon rank-sum test, p value = $1.8 \times$

321 10^{-24}). To further assess the cancer cell fraction of these shared HBV-induced SVs, we modified
322 the existing algorithm SVclone (Cmero et al. 2020) to suit long-read sequencing data. By
323 integrating established timing algorithms MutationTimeR (Gerstung et al. 2020), these shared
324 HBV-induced SVs were identified as early clonal events, preceding copy number gains (Fig. 5C).

325 **HBV-induced SVs influence genes associated with tumorigenesis.**

326 Next, we studied the potentially effects of HBV-induced SVs on genes, including those
327 genes directly intersecting with breakpoints or located closest to the breakpoints. Pathway
328 enrichment analysis of the relevant genes revealed their significant association with key HCC
329 pathways, such as Proteoglycans in cancer, PI3K/AKT signaling pathway, and Viral
330 carcinogenesis (Supplemental Fig. S16A, p value < 0.01). In addition, we found that shared HBV-
331 induced SVs within individual impact a total of 24 genes, including genes associated with
332 tumorigenesis, such as *SLC4A5*, *ABC13*, *NRG1*, and *SORCS1* (Fig. 6A). In addition, compared to
333 breakpoints occurring independently, breakpoints of shared HBV-induced SVs were significantly
334 concentrated in gene bodies and their flanking 5 kilobases (kb) regions (Supplemental Fig. S16B,
335 chi-square test, p value = 2.8×10^{-5}). This finding suggests that shared HBV-induced SVs were
336 more likely to have impact on gene regulation and expression. Notably, among these shared
337 HBV-induced SVs, TRAs account for 71.43%, highlighting the evolutionary advantage of this
338 type of SV.

339 We next investigated the functional consequences of HBV-induced TRAs on genes located
340 near the integration breakpoints. Among these genes, *NRG1*, *KLHDC7B*, and *KLHDC7B-DT*,
341 which have been implicated in various cancers (Maillet et al. 2018; Li et al. 2021; Yahiro et al.
342 2021), were directly affected by the translocations. In HCC13, a shared HBV-induced TRA was
343 detected across five sampling sites, with breakpoints located within the intronic region between
344 the first and second exons of *NRG1*, as well as upstream of *KLHDC7B* (620 bp) and *KLHDC7B-*
345 *DT* (2359 bp). While several previous studies have suggested their relevance to cancer (Martin-

346 Pardillos and Cajal 2019; Rosas et al. 2021), mutations within these genes have not been
347 identified in the context of HCC. Thus, we further investigated the impact of the HBV-induced
348 TRA on the expression of genes situated on either side. Long-read RNA sequencing was
349 performed on both the HCC13_T1 tumor sample and its corresponding adjacent nontumor tissue.
350 To obtain high-quality sequences, we conducted local assembly of the TRA using reads
351 supporting this integration events. By aligning RNA reads spanning the integration breakpoints to
352 the assembled translocation sequence, we identified a fusion gene between the *KLHDC7B*-DT
353 gene and the inserted HBV sequence in the tumor sample. However, no evidence of this fusion
354 gene was detected in the adjacent nontumor sample (Fig. 6B). Subsequent PCR validation
355 confirmed the presence of this fusion gene in the tumor sample and its absence in the adjacent
356 nontumor sample (Fig. 6C). Additionally, quantitative PCR (qPCR) revealed that the expression
357 levels of *KLHDC7B* in samples with HBV integration were significantly reduced, while
358 *KLHDC7B-DT* expression was markedly increased (Fig. 6D,E). These findings suggest that HBV
359 integration may disrupt the promoter of *KLHDC7B*, leading to decreased expression of this gene.
360 Furthermore, qPCR demonstrated that HBV integration led to reduced expression of *NRG1* in the
361 tumor (Supplemental Fig. S16C), potentially attributed to the integration disrupting the structural
362 integrity of the gene. Overall, these findings indicate that HBV-induced TRA can directly impact
363 genes adjacent to integration sites, including the formation of fusion genes and disruption of gene
364 expression.

365

366 Discussion

367 Somatic SVs play a pivotal role in HCC, yet their accurate identification using short-read
368 high-throughput sequencing poses challenges. To overcome limitations related to intratumor
369 heterogeneity (ITH) and refine the estimation of the mutational landscape, we employed multisite
370 sampling in tandem with long-read nanopore sequencing for the comprehensive genome-wide

371 analysis of somatic SVs. We depicted the characteristics of somatic SVs in both adjacent
372 nontumor and tumor tissues, revealing somatic SVs as early genomic alterations that may
373 contribute to the complexity of HCC development. Furthermore, we elucidated that HBV
374 integration serves as a frequent driver of somatic SVs, particularly in the context of inter-
375 chromosomal translocations. Our findings indicated that shared HBV-induced SVs represent early
376 clonal events in the progression of HCC, with the potential to directly impact the transcription of
377 genes relevant to tumorigenesis.

378 Building on these findings, it is important to note that while our study provides compelling
379 insights into the role of somatic SVs in HCC, some of the results require additional validation to
380 establish their broader significance. Specifically, we observed that the proportion of somatic
381 DUPs, INVs, and TRAs was significantly higher in patients with a history of smoking compared
382 to non-smokers. However, this observation warrants further validation through the inclusion of
383 additional patient samples to improve statistical robustness. Utilizing multiple tumor sampling
384 sites, we identified that HCC13 exhibits the highest number of somatic INVs in tumor sites T1,
385 T2, and T5, which markedly differ from those at T3 and T4. The morphological characteristics of
386 HCC13 indicate that T1, T2, and T5 are spatially distinct from T3 and T4, suggesting that the
387 differing INV patterns may reflect their clonal origins. This highlights the potential of somatic
388 SVs as biomarkers for pathological subtyping, offering insights into the clonal evolution of
389 tumors and their underlying mechanisms. Further investigation into these patterns could enhance
390 our understanding of tumor heterogeneity and its implications for personalized therapeutic
391 strategies.

392 Recent studies have shown that adjacent nontumor tissues are not healthy and may harbor
393 molecular features associated with tumors (Carlessi et al. 2023; Zhu et al. 2023; Chen et al.
394 2024b). A recent study has found that the complexity of HCC arises during the progression to
395 chronic liver disease and subsequent malignant transformation (Brunner et al. 2019). However, to

our knowledge, the association of genomic mutation features between adjacent nontumor and tumor tissues remains to be further elucidated. In this study, we observed a comparable prevalence of somatic SVs in adjacent nontumor samples and tumor samples, contrasting with the distinct patterns observed for somatic SNVs and CNVs. This could be attributed to the instability of the host genome caused by HBV DNA integration (Zhao et al. 2016). The number of somatic SNVs and somatic CNVs in adjacent nontumor tissues was significantly lower than that observed in tumor tissues, suggesting that these alterations may share similar developmental processes in HCC (Huang et al. 2017). Notably, the differences in somatic CNVs observed between HCC and adjacent nontumor tissues are consistent with findings from a study on urothelial cell carcinoma (Li et al. 2020a), suggesting that this may represent a general feature of tumor evolution.

Furthermore, a considerable proportion of somatic SVs were found to be shared between the adjacent nontumor and tumor tissues. INs and DELs have emerged as shared SV types between tumor and adjacent nontumor tissues, potentially exerting a direct influence on genes integral to normal hepatic function, including *TCPI10L* and *EVAIC*. A possible explanation is that these somatic SVs may arise from pre-malignant intermediate state cells or tumor cells within the adjacent nontumor tissues (Aran et al. 2017; Carlessi et al. 2023). Additionally, a significant proportion of the insertions and deletions consisted of transposons and other repetitive sequences, indicating that these elements may play an important role in tumorigenesis. The challenge of accurately mapping such repetitive sequences using short-read sequencing likely accounts for the increased detection of somatic SVs in our study. Given that these repetitive sequences are often refractory to accurate mapping using short-read sequencing, this could explain the higher number of somatic SVs detected in our study. However, the functional implications of these shared somatic SVs between adjacent nontumor and tumor tissues require further investigation, particularly from the perspectives of transcriptomics and epigenetics.

420 The integration of the HBV was considered a major factor driving the development of HCC.
421 Studies had shown that HBV integration could lead to an increased burden of mutations in the
422 genomic region within tens to hundreds of kilobases around the integration sites (Zapatka et al.
423 2020). However, the specific characteristics of HBV integration patterns were not yet clearly
424 understood. Taking advantage of long-read sequencing, we could detect the complete sequence of
425 HBV integration. Our findings indicated that HBV integration demonstrated a propensity towards
426 generating somatic SVs linked to chromosomal rearrangements, particularly inter-chromosomal
427 translocations, while simple insertions made up only about 20.20% of the occurrences. Notably,
428 shared HBV-induced SVs emerged as distinctive signatures of HCC, exhibiting a clustering
429 tendency proximal to various cancer-associated genes. However, we did not detect HBV
430 integration in the promoter region of TERT, despite previous studies reporting its presence in
431 approximately 30% of HCC cases (Chen et al. 2024a). This discrepancy may be stem from
432 limitations in our sample size or sequencing depth. In comparison to adjacent nontumor tissues,
433 HBV DNA integration in HCC leads to a higher incidence of gene fusion events (Zhuo et al.
434 2021). Using long-read RNA sequencing, we confirmed the formation of fusion transcripts
435 between HBV and human genes. Previous studies have suggested that HBV DNA methylation
436 could influence the expression patterns of viral genes (Fernandez et al. 2009; Mirabello et al.
437 2012). However, we did not observe significant differences between adjacent nontumor and
438 tumor tissues in the methylation status of integrated HBV DNA sequences. This highlights the
439 need for more advanced techniques, such as Cas9-targeted nanopore sequencing (Goldsmith et al.
440 2021), to enhance the detection sensitivity of HBV DNA and facilitate a deeper exploration of its
441 potential epigenetic regulation.

442 In summary, our comprehensive analysis of somatic SVs in HCC, including those induced
443 by HBV integration, shed light on the early genomic alterations and clonal events contributing to
444 HCC development. The intricate relationship between viral integration and somatic SVs

445 emphasizes the importance of considering these factors in understanding the molecular landscape
446 of HCC. This study extends our understanding of somatic SVs in HCC, providing valuable
447 insights into their characteristics and significance in early tumorigenesis.

448 **Methods**

449 **Long-read whole-genome sequencing**

450 Tumor, nontumor, and blood samples from each patient were frozen in liquid nitrogen for
451 DNA isolation. Genomic DNA from each sample was extracted with the MagAttract HMW DNA
452 Kit (QIAGEN). DNA libraries were constructed using the SQK-LSK109 library preparation kit
453 (ONT, UK) following the manufacturer's instructions. Then, the prepared libraries were
454 sequenced on the PromethION sequencer (ONT, UK). Raw sequencing data were basecalled
455 using Guppy 3.2.8 with the default parameters during sequencing.

456 **Short-read whole-genome sequencing**

457 Genomic DNA was isolated from fresh-frozen tumor, nontumor, and blood tissues using the
458 DNeasy Blood & Tissue Kit (QIAGEN) according to the manufacturer's protocol. The purity and
459 integrity of DNA were assessed via agarose gel electrophoresis. The DNA concentration was
460 measured using Qubit 2.0 (Invitrogen). The DNA was fragmented to approximately 350 bp using
461 the Covaris ultrasonicator. The library was constructed using the established Illumina paired-end
462 protocols. Following library quality control, sequencing was performed using the Illumina
463 NovaSeq 6000 platform (Illumina).

464 **Identification of somatic SVs**

465 First, clean reads were obtained by excluding those with an average base quality < 7 and a
466 length < 1000 bp. Then, minimap2 (v.2.17) (Li 2018) was used to align the clean reads to a
467 custom reference genome (composed of hg38 and 21 HBV genomes (Yan et al. 2015; Zheng et al.

2021)), with the following parameters: `minimap2 -N 10 -p 0.3 -ax map-ont --MD`. To calculate somatic SVs, blood samples were used as normal controls, and the BAM files from paired tumor or adjacent nontumor samples were merged with the BAM file from the blood sample after tagging. The merged files were then processed with Sniffles (v.2.0.6) (Smolka et al. 2024) to identify SV breakpoints and discern supporting reads, with the following parameters: `Sniffles -i bam_merge -v vcf --tandem-repeats -t 36 --minsupport 1 --mapq 20 --min-alignment-length 500 -minsvlen 50 --output-mnames`. All SVs with supporting reads ≥ 3 and a length ≥ 50 bp were identified. If all supporting reads originated from tumor or adjacent nontumor samples, and none from the normal control, the SV was categorized as a candidate somatic SV.

Next, we rigorously filtered the candidate somatic SVs. ONT sequencing data from the liver tissues of five healthy individuals (Pascarella et al. 2022), along with data from 63 healthy individuals from the HPRC cohort (Liao et al. 2023), were obtained (Supplemental Table S4). Together with the five blood samples from this study, SVs were identified using Sniffles with the following parameters: `Sniffles -i bam_merge -v vcf --tandem-repeats -t 36 --minsupport 3 --mapq 20 --min-alignment-length 500 --minsvlen 50`. Additionally, SVs derived from a cohort of 405 individuals from the Chinese population (Wu et al. 2021), also sequenced via nanopore technology, as well as SVs from the HGSC database (Ebert et al. 2021), were incorporated. These SVs were used to construct a Panel of Normals (PON), which served to filter out false positives in somatic SV calls. Such false positives could arise due to insufficient sequencing depth in paired normal samples, biases associated with ONT technology, or tissue-specific SVs in normal liver tissue. Any candidate somatic SVs present in the PON were excluded. We also removed SVs occurring in unreliable alignment regions, including those where over 50% of reads in the SV region had mapping quality < 20 , and SVs with variant allele frequency below 0.01. Furthermore, using `genomview` (Abeel et al. 2012) to visualize SVs, we manually inspected all candidate somatic SVs. Following the methodology employed in a previous study (Fujimoto et al.

493 2021), candidate somatic SVs displaying matching breakpoints in corresponding normal samples
494 were discarded from further analysis.

495 **VCF merging of SVs and identification of shared somatic SVs**

496 The Jasmine (v1.1.5) software was used to merge the VCF files of somatic SVs across
497 different samples, incorporating PON filtering and identification of shared SVs. This tool is
498 specifically designed for SV detection from long-read sequencing, considering factors such as the
499 type, length, breakpoint positions, and sequence context of the SVs. Jasmine compares and
500 merges SV calls across samples by representing variants as points in a multidimensional space
501 and constructing a proximity graph based on their breakpoints and lengths. We employed Jasmine
502 with the following parameters: `jasmine file_list out_file --ignore_strand`. The default parameters
503 in Jasmine merge SVs of the same type if they are within 100 bp of each other and have length
504 differences of less than 50%. SVs that could be merged are considered as shared SVs.

505 **Detection of HBV integrations and HBV-induced SVs**

506 Although the alignment results of long reads against the customized reference genome
507 already contain partial information concerning HBV integration, notably the presence of chimeric
508 reads aligning to both HBV and human genomes, our examination revealed that certain HBV
509 DNA sequences present on some long reads failed to align to the 21 HBV reference genomes.
510 These HBV sequences either received low alignment scores in minimap2, resulting in being
511 labeled as INS, or were too short and directly clipped without alignment. Based on these
512 characteristics of reads, we developed a workflow for detecting HBV integration using long-read
513 sequencing, mainly consisting of the following steps: (1) extraction of chimeric reads, (2)
514 extraction of INS sequences and clipped sequences from all long reads, (3) re-alignment of INS
515 and clipped sequences to the custom reference genome using BLAST, and (4) determination of
516 HBV DNA integration breakpoints.

517 Benefiting from the advantages of long-read sequencing, long reads can simultaneously
518 contain HBV DNA sequences and their adjacent human sequences. We utilized these spanning
519 reads to further deduce the structural consequences of HBV DNA integration on the human
520 genome. Based on the alignment direction, position, and order of human sequences flanking HBV
521 DNA, five types of HBV-induced structural variations were identified. Insertion (INS) denotes
522 direct insertion of HBV DNA into the human genome without altering the flanking human
523 sequences. Deletion (DEL) corresponds to the loss of human DNA segments caused by HBV
524 DNA integration. Duplication (DUP) refers to the presence of duplicated human DNA segments
525 flanking the HBV DNA integration breakpoint. Inversion (INV) occurs when the human DNA
526 adjacent to one side of the HBV DNA integration breakpoint is reversed. Translocation (TRA)
527 involves human sequences flanking HBV DNA originating from different chromosomes.

528 **Identification of somatic SNVs and somatic CNVs**

529 The clean reads were aligned to the human reference genome (hg38) using Burrows-Wheeler
530 Aligner (BWA, v.0.7.17)(Li 2013), with the following parameters: BWA-MEM -M -Y -t 38 -R.
531 SAMtools (v.1.9) (Li et al. 2009) was employed for sorting and indexing the resulting BAM files.
532 PCR duplicate reads were marked using the MarkDuplicates module of the Genome Analysis
533 Toolkit (GATK, v.4.1.2.0) (McKenna et al. 2010). Local realignment and base quality
534 recalibration were conducted using the GATK BaseRecalibrator and ApplyBQSR modules.

535 Somatic SNVs were called from paired tumor and matched blood samples using GATK
536 Mutect2, followed by variant filtering using FilterMutectCalls. Next, annotation of the filtered
537 somatic SNVs was performed via ANNOVAR (Wang et al. 2010). To ensure high-quality
538 somatic SNVs, the following quality control criteria were applied: (1) the variants supported by at
539 least 7% of total reads in tumor sample and less than 2% of total reads in normal sample were
540 retained. (2) the variants with sequencing depths less than 8 in tumor sample or less than 6 in
541 normal sample were excluded. (3) Only variants with more than 3 mutation reads in the tumor

542 were considered. Additionally, potential germline variants were filtered out by excluding variants
543 present in the 1000 Genomes Project database with a minor allele frequency ≥ 0.001 .

544 Somatic CNVs were detected in paired tumor and normal samples using FACETS (v0.6.2)
545 (Shen and Seshan 2016), alongside evaluation of tumor purity and ploidy. The analysis was
546 performed with the following parameters: `cnv_facets -t tumor_bam -n normal_bam -vcf snps --`
547 `snp-nprocs 30 --depth 15 4000 --cval 25 400 --nbhd-snp 500`. Copy numbers characterized by a
548 total copy number of 2 and a lesser copy number of 1 were subsequently filtered out, as this is
549 consistent with the normal diploid copy number state. To focus on large-scale CNV alterations,
550 any CNV smaller than 100 kb in length was removed from the analysis.

551 **Phylogenetic Tree Construction and Clustering Analysis**

552 The genome was divided into 500-bp bins, and binary presence/absence matrices were
553 constructed based on SV breakpoints and SNVs. A distance-based Neighbor-Joining (NJ) method
554 was then applied to construct the phylogenetic trees. The resulting phylogenetic trees were
555 subsequently visualized using the `ggtree` (Yu et al. 2018) package.

556 For clustering analysis, Jaccard indices were computed pairwise between samples. The
557 resulting matrix of Jaccard indices was subjected to hierarchical clustering, using the `pheatmap`
558 (v1.0.12) package in R (R Core Team 2023) with the “Pearson’s correlation” parameter.

559 **Timing of somatic SVs**

560 SVclone (v1.1.2) (Cmero et al. 2020) is a computational tool designed to infer the cancer cell
561 fraction (CCF) of SV breakpoints based on short read sequencing data. To adapt SVclone for SVs
562 detected from nanopore sequencing, several modifications were implemented: (1) SV supporting
563 reads (including split and spanning reads), along with normal reads crossing the SV breakpoints,
564 were utilized to calculate the variant allele fraction (VAF) of SV breakpoints. (2) As each long
565 read in nanopore sequencing represents a single DNA molecule and does not overestimate the

566 number of normal reads in cases of DNA gain, the step of adjusting the number of normal reads
567 was omitted. By integrating somatic SNVs, somatic CNVs, tumor purity, and ploidy, SVclone
568 performed CCF estimation and clustering using a co-clustering mode. Next, clonality of somatic
569 SVs was determined based on the somatic CNVs and SVclone clusters using the timing analysis
570 algorithm MutationTimeR (Gerstung et al. 2020). This step assigned somatic SVs to four different
571 timing categories , include early clonal, clonal not specified, late clonal and subclonal.

572 **Long-read RNA sequencing (RNA-seq)**

573 RNA from tumor and nontumor tissues was extracted using TRIzol Reagent (Invitrogen)
574 according to the standard instructions. The RNA concentration was determined using Qubit RNA
575 HS Assay Kit (Thermo Fisher Scientific2), and the purity and integrity were assessed by
576 NanoDrop One, Qubit, and agarose gel electrophoresis. The cDNA library was prepared
577 according to the SQK-DCS109 and EXP-NBD104 protocol (ONT). Then, the cDNA library was
578 sequenced on a PromethION machine (ONT). Basecalling from raw data was performed in
579 batches using Guppy software 3.2.8.

580 **PCR validation of human-HBV fusion expression**

581 Primers were designed to target the *KLHDC7B-DT* gene and HBV DNA sequences
582 identified in the long-read RNA sequencing (Supplemental Table S5). A total of 35 ng SMART
583 cDNA was mixed with 1 μ L of forward primer, 1 μ L of reverse primer, 12.5 μ L of 2 \times Taq
584 Master Mix (Vazyme), and nuclease-free water to a total volume of 25 μ L. The thermal cycling
585 protocol was as follows: 3 min at 95 $^{\circ}$ C for initial denaturation, 35 cycles of 15 s at 95 $^{\circ}$ C, 15 s at
586 58 $^{\circ}$ C, and 4 min at 72 $^{\circ}$ C, followed by 5 min at 72 $^{\circ}$ C. Next, DNA was electrophoresed on a 1%
587 agarose gel and sent to Youkang Bio for first-generation sequencing.

588 **qPCR validation of gene expression**

589 cDNA was prepared from 1 μ g of RNA using the PrimeScript RT reagent Kit with gDNA
590 Eraser (TAKARA) following the manufacturer's protocol. The qPCR primers were designed
591 using Primer3Plus. Each pair of primers was designed to target exons to avoid amplification of
592 genomic DNA. Oligos for qPCR are provided in Supplemental Table S5. For each sample, three
593 replicates were performed using iTaq universal SYBR Green supermix (Bio-rad) following the
594 manufacturer's instructions. The thermal cycling protocol on the CFX96™ Real-Time System
595 (Bio-rad) was as follows: 3 min at 95 °C for initial denaturation, 41 cycles of (10 s at 95 °C, 30 s
596 at 60 °C, and 15 s at 72 °C), followed by a melt curve from 65 °C to 95 °C in 0.5 °C increments.
597 β -actin was employed as the internal control, and the Cq value was used for quantitative analysis.
598 The experiment was repeated three times.

599 **Data Access**

599 The raw sequence data reported in this paper have been deposited in the Genome Sequence
599 Archive in National Genomics Data Center(CNCB-NGDC Members Partners 2024), China
599 National Center for Bioinformatics/ Beijing Institute of Genomics, Chinese Academy of Sciences
599 (GSA-Human: HRA007177) that are publicly accessible at <https://ngdc.cncb.ac.cn/gsa-human>.
599 All code used is open source and available at https://github.com/tianfuzeng/ONT_SV, as well as
599 in the Supplemental Code.

599 **Competing interests**

599 The authors declare no potential conflicts of interest.

599 **Acknowledgments**

599 We thank Ranlei Wei from the Laboratory of Omics Technology and Bioinformatics, West
599 China Hospital of Sichuan University, Kefei Yuan from the West China Hospital of Sichuan
599 University for their technical assistance. This work was supported by grants from the National
599 Natural Science Foundation of China (82173383 to D. Xie, 32200508 to L. Xia, 82202260 to H.

513 Liao), the Sichuan Province Science and Technology Program (2022NSFSC1553 to L. Xia), the
514 1-3-5 project for disciplines of excellence, West China Hospital, Sichuan University (ZYYC23024
515 to D. Xie), and China Postdoctoral Science Foundation (2022TQ0221 to H. Liao).

516

517 **Author contributions**

518 T.Z., H.L. and L.X. were responsible for the experimental design, execution, data analysis.
519 T.Z. wrote the paper, with contributions from all other authors. Y.H. and J.Z. were responsible for
520 RNA sequencing and validation. S.Y., Y.L. and X.L. were responsible for data downloading,
521 organization, and processing of raw data. D.X. was responsible for supervision of research, data
522 interpretation and manuscript preparation.

523 **References**

- 524 Abeel T, Van Parys T, Saeys Y, Galagan J, Van de Peer Y. 2012. GenomeView: a next-
525 generation genome browser. *Nucleic Acids Res* **40**: e12.
- 526 Abyzov A, Li S, Kim DR, Mohiyuddin M, Stutz AM, Parrish NF, Mu XJ, Clark W, Chen K,
527 Hurles M et al. 2015. Analysis of deletion breakpoints from 1,092 humans reveals details
528 of mutation mechanisms. *Nat Commun* **6**: 7256.
- 529 Aganezov S, Goodwin S, Sherman RM, Sedlazeck FJ, Arun G, Bhatia S, Lee I, Kirsche M,
530 Wappel R, Kramer M et al. 2020. Comprehensive analysis of structural variants in breast
531 cancer genomes using single-molecule sequencing. *Genome Res* **30**: 1258-1273.
- 532 Alvarez EG, Demeulemeester J, Otero P, Jolly C, Garcia-Souto D, Pequeno-Valtierra A, Zamora
533 J, Tojo M, Temes J, Baez-Ortega A et al. 2021. Aberrant integration of Hepatitis B virus
534 DNA promotes major restructuring of human hepatocellular carcinoma genome
535 architecture. *Nat Commun* **12**: 6910.
- 536 Aran D, Camarda R, Odegaard J, Paik H, Oskotsky B, Krings G, Goga A, Sirota M, Butte AJ.
537 2017. Comprehensive analysis of normal adjacent to tumor transcriptomes. *Nat Commun*
538 **8**: 1077.
- 539 Beyter D, Ingimundardottir H, Oddsson A, Eggertsson HP, Bjornsson E, Jonsson H, Atlason BA,
540 Kristmundsdottir S, Mehringer S, Hardarson MT et al. 2021. Long-read sequencing of
541 3,622 Icelanders provides insight into the role of structural variants in human diseases and
542 other traits. *Nat Genet* **53**: 779-786.
- 543 Bray F, Laversanne M, Sung H, Ferlay J, Siegel RL, Soerjomataram I, Jemal A. 2024. Global
544 cancer statistics 2022: GLOBOCAN estimates of incidence and mortality worldwide for
545 36 cancers in 185 countries. *CA Cancer J Clin* doi:10.3322/caac.21834.
- 546 Brunner SF, Roberts ND, Wylie LA, Moore L, Aitken SJ, Davies SE, Sanders MA, Ellis P, Alder
547 C, Hooks Y et al. 2019. Somatic mutations and clonal dynamics in healthy and cirrhotic
548 human liver. *Nature* **574**: 538-542.
- 549 Carlessi R, Denisenko E, Boslem E, Kohn-Gaone J, Main N, Abu Bakar NDB, Shirolkar GD,
550 Jones M, Beasley AB, Poppe D et al. 2023. Single-nucleus RNA sequencing of pre-

- 551 malignant liver reveals disease-associated hepatocyte state with HCC prognostic potential.
552 *Cell Genom* **3**: 100301.
- 553 Chen L, Zhang C, Xue R, Liu M, Bai J, Bao J, Wang Y, Jiang N, Li Z, Wang W et al. 2024a.
554 Deep whole-genome analysis of 494 hepatocellular carcinomas. *Nature* **627**: 586-593.
- 555 Chen Z, Shi Q, Zhao Y, Xu M, Liu Y, Li X, Liu L, Sun M, Wu X, Shao Z et al. 2024b. Long-read
556 transcriptome landscapes of primary and metastatic liver cancers at transcript resolution.
557 *Biomark Res* **12**: 4.
- 558 Choy TK, Wang CY, Phan NN, Khoa Ta HD, Anuraga G, Liu YH, Wu YF, Lee KH, Chuang JY,
559 Kao TJ. 2021. Identification of Dipeptidyl Peptidase (DPP) Family Genes in Clinical
560 Breast Cancer Patients via an Integrated Bioinformatics Approach. *Diagnostics (Basel)*
561 **11**.
- 562 Cmero M, Yuan K, Ong CS, Schroder J, Evolution P, Heterogeneity Working G, Corcoran NM,
563 Papenfuss T, Hovens CM, Markowitz F et al. 2020. Inferring structural variant cancer cell
564 fraction. *Nat Commun* **11**: 730.
- 565 CNCB-NGDC Members Partners. 2024. Database Resources of the National Genomics Data
566 Center, China National Center for Bioinformation in 2024. *Nucleic Acids Res* **52**: D18-
567 D32.
- 568 Cosenza MR, Rodriguez-Martin B, Korbel JO. 2022. Structural Variation in Cancer: Role,
569 Prevalence, and Mechanisms. *Annu Rev Genomics Hum Genet* **23**: 123-152.
- 570 Craven KE, Fischer CG, Jiang L, Pallavajjala A, Lin MT, Eshleman JR. 2022. Optimizing
571 Insertion and Deletion Detection Using Next-Generation Sequencing in the Clinical
572 Laboratory. *J Mol Diagn* **24**: 1217-1231.
- 573 Ebert P, Audano PA, Zhu Q, Rodriguez-Martin B, Porubsky D, Bonder MJ, Sulovari A, Ebler J,
574 Zhou W, Serra Mari R et al. 2021. Haplotype-resolved diverse human genomes and
575 integrated analysis of structural variation. *Science* **372**.
- 576 Fernandez AF, Rosales C, Lopez-Nieva P, Grana O, Ballestar E, Ropero S, Espada J, Melo SA,
577 Lujambio A, Fraga MF et al. 2009. The dynamic DNA methylomes of double-stranded
578 DNA viruses associated with human cancer. *Genome Res* **19**: 438-451.
- 579 Fujimoto A, Furuta M, Totoki Y, Tsunoda T, Kato M, Shiraishi Y, Tanaka H, Taniguchi H,
580 Kawakami Y, Ueno M et al. 2016. Whole-genome mutational landscape and
581 characterization of noncoding and structural mutations in liver cancer. *Nat Genet* **48**: 500-
582 509.
- 583 Fujimoto A, Wong JH, Yoshii Y, Akiyama S, Tanaka A, Yagi H, Shigemizu D, Nakagawa H,
584 Mizokami M, Shimada M. 2021. Whole-genome sequencing with long reads reveals
585 complex structure and origin of structural variation in human genetic variations and
586 somatic mutations in cancer. *Genome Med* **13**: 65.
- 587 George J, Lim JS, Jang SJ, Cun Y, Ozretic L, Kong G, Leenders F, Lu X, Fernandez-Cuesta L,
588 Bosco G et al. 2015. Comprehensive genomic profiles of small cell lung cancer. *Nature*
589 **524**: 47-53.
- 590 Gerstung M, Jolly C, Leshchiner I, Dentre SC, Gonzalez S, Rosebrock D, Mitchell TJ, Rubanova
591 Y, Anur P, Yu K et al. 2020. The evolutionary history of 2,658 cancers. *Nature* **578**: 122-
592 128.
- 593 Goldsmith C, Cohen D, Dubois A, Martinez MG, Petitjean K, Corlu A, Testoni B, Hernandez-
594 Vargas H, Chemin I. 2021. Cas9-targeted nanopore sequencing reveals epigenetic
595 heterogeneity after de novo assembly of native full-length hepatitis B virus genomes.
596 *Microb Genom* **7**.
- 597 Ho XD, Nguyen HG, Trinh LH, Reimann E, Prans E, Koks G, Maasalu K, Le VQ, Nguyen VH,
598 Le NTN et al. 2017. Analysis of the Expression of Repetitive DNA Elements in
599 Osteosarcoma. *Front Genet* **8**: 193.

- 700 Hochhaus A, Larson RA, Guilhot F, Radich JP, Branford S, Hughes TP, Baccarani M, Deininger
701 MW, Cervantes F, Fujihara S et al. 2017a. Long-Term Outcomes of Imatinib Treatment
702 for Chronic Myeloid Leukemia. *N Engl J Med* **376**: 917-927.
- 703 Hochhaus A, Masszi T, Giles FJ, Radich JP, Ross DM, Gomez Casares MT, Hellmann A,
704 Stentoft J, Conneally E, Garcia-Gutierrez V et al. 2017b. Treatment-free remission
705 following frontline nilotinib in patients with chronic myeloid leukemia in chronic phase:
706 results from the ENESTfreedom study. *Leukemia* **31**: 1525-1531.
- 707 Hu Z, Qu S. 2021. EVA1C Is a Potential Prognostic Biomarker and Correlated With Immune
708 Infiltration Levels in WHO Grade II/III Glioma. *Front Immunol* **12**: 683572.
- 709 Huang A, Zhao X, Yang XR, Li FQ, Zhou XL, Wu K, Zhang X, Sun QM, Cao Y, Zhu HM et al.
710 2017. Circumventing intratumoral heterogeneity to identify potential therapeutic targets in
711 hepatocellular carcinoma. *J Hepatol* **67**: 293-301.
- 712 Khemlina G, Ikeda S, Kurzrock R. 2017. The biology of Hepatocellular carcinoma: implications
713 for genomic and immune therapies. *Mol Cancer* **16**: 149.
- 714 Kolomietz E, Meyn MS, Pandita A, Squire JA. 2002. The role of Alu repeat clusters as mediators
715 of recurrent chromosomal aberrations in tumors. *Genes Chromosomes Cancer* **35**: 97-112.
- 716 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
717 *bioRxiv*: Preprint at <https://arxiv.org/abs/1303.3997v1302>.
- 718 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094-
719 3100.
- 720 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
721 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and
722 SAMtools. *Bioinformatics* **25**: 2078-2079.
- 723 Li MX, Wang HY, Yuan CH, Ma ZL, Jiang B, Li L, Zhang L, Xiu DR. 2021. KLHDC7B-DT
724 aggravates pancreatic ductal adenocarcinoma development via inducing cross-talk
725 between cancer cells and macrophages. *Clin Sci (Lond)* **135**: 629-649.
- 726 Li R, Du Y, Chen Z, Xu D, Lin T, Jin S, Wang G, Liu Z, Lu M, Chen X et al. 2020a.
727 Macroscopic somatic clonal expansion in morphologically normal human urothelium.
728 *Science* **370**: 82-89.
- 729 Li Y, Roberts ND, Wala JA, Shapira O, Schumacher SE, Kumar K, Khurana E, Waszak S, Korbel
730 JO, Haber JE et al. 2020b. Patterns of somatic structural variation in human cancer
731 genomes. *Nature* **578**: 112-121.
- 732 Liao WW, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, Lu S, Lucas JK, Monlong J, Abel HJ et
733 al. 2023. A draft human pangenome reference. *Nature* **617**: 312-324.
- 734 Llovet JM, Kelley RK, Villanueva A, Singal AG, Pikarsky E, Roayaie S, Lencioni R, Koike K,
735 Zucman-Rossi J, Finn RS. 2021. Hepatocellular carcinoma. *Nat Rev Dis Primers* **7**: 6.
- 736 Maillet V, Boussetta N, Leclerc J, Fauveau V, Foretz M, Viollet B, Couty JP, Celton-Morizur S,
737 Perret C, Desdouets C. 2018. LKB1 as a Gatekeeper of Hepatocyte Proliferation and
738 Genomic Integrity during Liver Regeneration. *Cell Rep* **22**: 1994-2005.
- 739 Martin-Pardillos A, Cajal SRY. 2019. Characterization of Kelch domain-containing protein 7B in
740 breast tumours and breast cancer cell lines. *Oncol Lett* **18**: 2853-2860.
- 741 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
742 Altshuler D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce
743 framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297-
744 1303.
- 745 Mirabello L, Sun C, Ghosh A, Rodriguez AC, Schiffman M, Wentzensen N, Hildesheim A,
746 Herrero R, Wacholder S, Lorincz A et al. 2012. Methylation of human papillomavirus
747 type 16 genome and risk of cervical precancer in a Costa Rican population. *J Natl Cancer*
748 *Inst* **104**: 556-565.

- 749 Muller M, Bird TG, Nault JC. 2020. The landscape of gene mutations in cirrhosis and
750 hepatocellular carcinoma. *J Hepatol* **72**: 990-1002.
- 751 Munkhjargal B, Kondo K, Soejima S, Tegshee B, Takai C, Kawakita N, Toba H, Takizawa H.
752 2023. Aberrant methylation of dipeptidyl peptidase-like 6 as a potential prognostic
753 biomarker for lung adenocarcinoma. *Oncol Lett* **25**: 206.
- 754 Pascarella G, Hon CC, Hashimoto K, Busch A, Luginbuhl J, Parr C, Hin Yip W, Abe K, Kratz A,
755 Bonetti A et al. 2022. Recombination of repeat elements generates somatic complexity in
756 human genomes. *Cell* **185**: 3025-3040 e3026.
- 757 R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. R Foundation for
758 Statistical Computing, Vienna. <https://www.R-project.org/>.
- 759 Rajaby R, Liu DX, Au CH, Cheung YT, Lau AYT, Yang QY, Sung WK. 2023. INSurVeyor:
760 improving insertion calling from short read sequencing data. *Nat Commun* **14**: 3243.
- 761 Repana D, Nulsen J, Dressler L, Bortolomeazzi M, Venkata SK, Tourna A, Yakovleva A,
762 Palmieri T, Ciccarelli FD. 2019. The Network of Cancer Genes (NCG): a comprehensive
763 catalogue of known and candidate cancer genes from cancer sequencing screens. *Genome*
764 *Biol* **20**: 1.
- 765 Robberecht C, Voet T, Zamani Esteki M, Nowakowska BA, Vermeesch JR. 2013. Nonallelic
766 homologous recombination between retrotransposable elements is a driver of de novo
767 unbalanced translocations. *Genome Res* **23**: 411-418.
- 768 Rosas D, Raez LE, Russo A, Rolfo C. 2021. Neuregulin 1 Gene (NRG1). A Potentially New
769 Targetable Alteration for the Treatment of Lung Cancer. *Cancers (Basel)* **13**.
- 770 Sakamoto Y, Xu L, Seki M, Yokoyama TT, Kasahara M, Kashima Y, Ohashi A, Shimada Y,
771 Motoi N, Tsuchihara K et al. 2020. Long-read sequencing for non-small-cell lung cancer
772 genomes. *Genome Res* **30**: 1243-1257.
- 773 Schulze K, Nault JC, Villanueva A. 2016. Genetic profiling of hepatocellular carcinoma using
774 next-generation sequencing. *J Hepatol* **65**: 1031-1042.
- 775 Shen R, Seshan VE. 2016. FACETS: allele-specific copy number and clonal heterogeneity
776 analysis tool for high-throughput DNA sequencing. *Nucleic Acids Res* **44**: e131.
- 777 Shen S, Li H, Liu J, Sun L, Yuan Y. 2020. The panoramic picture of pepsinogen gene family with
778 pan-cancer. *Cancer Med* **9**: 9064-9080.
- 779 Smolka M, Paulin LF, Grochowski CM, Horner DW, Mahmoud M, Behera S, Kalef-Ezra E,
780 Gandhi M, Hong K, Pehlivan D et al. 2024. Detection of mosaic and population-level
781 structural variants with Sniffles2. *Nat Biotechnol* doi:10.1038/s41587-023-02024-y.
- 782 Soda M, Choi YL, Enomoto M, Takada S, Yamashita Y, Ishikawa S, Fujiwara S, Watanabe H,
783 Kurashina K, Hatanaka H et al. 2007. Identification of the transforming EML4-ALK
784 fusion gene in non-small-cell lung cancer. *Nature* **448**: 561-566.
- 785 Strandgaard T, Nordentoft I, Lamy P, Christensen E, Thomsen MBH, Jensen JB, Dyrskjöt L.
786 2020. Mutational Analysis of Field Cancerization in Bladder Cancer. *Bladder Cancer* **6**:
787 253-264.
- 788 Sung WK, Zheng H, Li S, Chen R, Liu X, Li Y, Lee NP, Lee WH, Ariyaratne PN, Tennakoon C
789 et al. 2012. Genome-wide survey of recurrent HBV integration in hepatocellular
790 carcinoma. *Nat Genet* **44**: 765-769.
- 791 The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Consortium. 2020. Pan-cancer
792 analysis of whole genomes. *Nature* **578**: 82-93.
- 793 Ting DT, Lipson D, Paul S, Brannigan BW, Akhavanfard S, Coffman EJ, Contino G, Deshpande
794 V, Iafrate AJ, Letovsky S et al. 2011. Aberrant overexpression of satellite repeats in
795 pancreatic and other epithelial cancers. *Science* **331**: 593-596.
- 796 Villanueva A. 2019. Hepatocellular Carcinoma. *N Engl J Med* **380**: 1450-1462.

- 797 Waddell N, Pajic M, Patch AM, Chang DK, Kassahn KS, Bailey P, Johns AL, Miller D, Nones K,
798 Quek K et al. 2015. Whole genomes redefine the mutational landscape of pancreatic
799 cancer. *Nature* **518**: 495-501.
- 300 Wang K, Li M, Hakonarson H. 2010. ANNOVAR: functional annotation of genetic variants from
301 high-throughput sequencing data. *Nucleic Acids Res* **38**: e164.
- 302 Wu Z, Jiang Z, Li T, Xie C, Zhao L, Yang J, Ouyang S, Liu Y, Li T, Xie Z. 2021. Structural
303 variants in the Chinese population and their impact on phenotypes, diseases and
304 population adaptation. *Nat Commun* **12**: 6501.
- 305 Xu L, Wang X, Lu X, Liang F, Liu Z, Zhang H, Li X, Tian S, Wang L, Wang Z. 2023. Long-read
306 sequencing identifies novel structural variations in colorectal cancer. *PLoS Genet* **19**:
307 e1010514.
- 308 Yahiro K, Ogura K, Tsutsuki H, Iyoda S, Ohnishi M, Moss J. 2021. A novel endoplasmic stress
309 mediator, Kelch domain containing 7B (KLHDC7B), increased Harakiri (HRK) in the
310 SubAB-induced apoptosis signaling pathway. *Cell Death Discov* **7**: 360.
- 311 Yan H, Yang Y, Zhang L, Tang G, Wang Y, Xue G, Zhou W, Sun S. 2015. Characterization of
312 the genotype and integration patterns of hepatitis B virus in early- and late-onset
313 hepatocellular carcinoma. *Hepatology* **61**: 1821-1831.
- 314 Yu G, Lam TT, Zhu H, Guan Y. 2018. Two Methods for Mapping and Visualizing Associated
315 Data on Phylogeny Using Ggtree. *Mol Biol Evol* **35**: 3041-3043.
- 316 Zapatka M, Borozan I, Brewer DS, Iskar M, Grundhoff A, Alawi M, Desai N, Sultmann H, Moch
317 H, Pathogens P et al. 2020. The landscape of viral associations in human cancers. *Nat*
318 *Genet* **52**: 320-330.
- 319 Zhao LH, Liu X, Yan HX, Li WY, Zeng X, Yang Y, Zhao J, Liu SP, Zhuang XH, Lin C et al.
320 2016. Genomic and oncogenic preference of HBV integration in hepatocellular carcinoma.
321 *Nat Commun* **7**: 12992.
- 322 Zheng B, Liu XL, Fan R, Bai J, Wen H, Du LT, Jiang GQ, Wang CY, Fan XT, Ye YN et al.
323 2021. The Landscape of Cell-Free HBV Integrations and Mutations in Cirrhosis and
324 Hepatocellular Carcinoma Patients. *Clin Cancer Res* **27**: 3772-3783.
- 325 Zhu H, Lin Y, Lu D, Wang S, Liu Y, Dong L, Meng Q, Gao J, Wang Y, Song N et al. 2023.
326 Proteomics of adjacent-to-tumor samples uncovers clinically relevant biological events in
327 hepatocellular carcinoma. *Natl Sci Rev* **10**: nwad167.
- 328 Zhuo Z, Rong W, Li H, Li Y, Luo X, Liu Y, Tang X, Zhang L, Su F, Cui H et al. 2021. Long-read
329 sequencing reveals the structural complexity of genomic integration of HBV DNA in
330 hepatocellular carcinoma. *NPJ Genom Med* **6**: 84.
- 331 Zucman-Rossi J, Villanueva A, Nault J-C, Llovet JM. 2015. Genetic Landscape and Biomarkers
332 of Hepatocellular Carcinoma. *Gastroenterology* **149**: 1226-1239.e1224.
- 333

334 **Figure 1. Identification of somatic SVs through multi-region long-read sequencing.**

335 A. Schematic diagram showing the sampling locations for five patients. Blue dots represent
336 adjacent nontumor tissues. Red dots represent tumor tissues.

337 B. Bar plots depicting the number of somatic SVs across individual samples, including insertions
338 (INS), deletions (DEL), duplications (DUP), inversions (INV), and translocations (TRA).

339 C. Stacked plots showing the distribution of somatic SV lengths categorized into two ranges:

340 (50bp, 1kb) and (1kb, 10kb), stratified based on SV types.

341 D. The proportion of different SV size ranges within each somatic SV type.

342

343 **Figure 2. Somatic SVs were prominent in adjacent nontumors.**

344 A. Bar plots showing the number of somatic mutations (classified into SVs, SNVs, and CNVs) in
345 adjacent nontumor and tumor tissues from all patients.

346 B. Heatmap depicting the hierarchical clustering of 37 tumor and adjacent nontumor samples
347 based on Jaccard indexes calculated from somatic SVs.

348 C. Heatmap depicting the hierarchical clustering of 37 tumor and adjacent nontumor samples
349 based on Jaccard indexes calculated from somatic SNVs. The white dashed box outlines a cluster
350 consisting of adjacent nontumor samples.

351

352 **Figure 3. The characteristics of shared somatic SVs between adjacent nontumor and tumor**
353 **tissues.**

354 A. The frequency of each somatic SV type in different categories was shown: shared (in both
355 adjacent nontumor and tumor tissues), T-specific (only in tumor tissues), and N-specific (only in
356 adjacent nontumor tissues).

357 B. The length distribution of shared somatic SVs between adjacent nontumor and tumor tissues,
358 as well as the SVs unique to either adjacent nontumor or tumor tissues.

359 C. An overview of genes overlapped by recurrently shared somatic SVs in adjacent nontumor and
360 tumor tissues across patients. The left bar chart shows the mutation percentage across patients.
361 The right bar chart displays the mutation percentage across samples.

362

363 **Figure 4. Inferring SV types based on HBV DNA integration patterns.**

364 A. In sample HCC9_T3, an HBV DNA insertion induces a fold-back inversion in the host DNA.

365 The top panel illustrates the CNV increase occurring approximately 76.23 Mb downstream of the
366 integration site. Sequencing data demonstrate that long reads encompass the sequences flanking
367 the HBV integration, revealing identical breakpoints and reverse duplication of the adjacent
368 sequences.

369 B. An example of HBV DNA integration resulting in a translocation (HCC13_T2). The IGV plot
370 below illustrating the alignment of long reads to Chr22, HBV, and Chr8. The plot above depicting
371 the concordance of its breakpoint on Chr8 (red dashed line) with a breakpoint of approximately
372 5.36 Mbp CNV (gray shade).

373

374 **Figure 5. Integration breakpoints of HBV DNA in human and HBV genome.**

375 A. Integration breakpoints of HBV DNA within the human genome. The same color scheme
376 represents the same patient. Circular points represent tumor samples, while triangular points
377 represent adjacent nontumor samples.

378 B. A schematic diagram of shared HBV-induced SVs at all tumor sampling sites in the patients.
379 The horizontal line represents the relative position of the integrated HBV sequence within the
380 HBV genome.

381 C. In HCC13, the clonal status of HBV integration sites was in the context of somatic SVs (top
382 panel), with green indicating early clonal events. The red asterisk symbol represents the shared
383 HBV integration sites. The middle panel represents the CNV status, while the lower panel
384 displays the inferred mutation timeline.

385

386 **Figure 6. Genes disrupted by HBV-induced SVs**

387 A. Manhattan plot showing the genes affected by HBV-induced SVs. Tumor samples are
388 represented by orange points, while adjacent nontumor samples are represented by blue points.
389 Tumor genes annotated by COSMIC or NCG6 are marked in red.

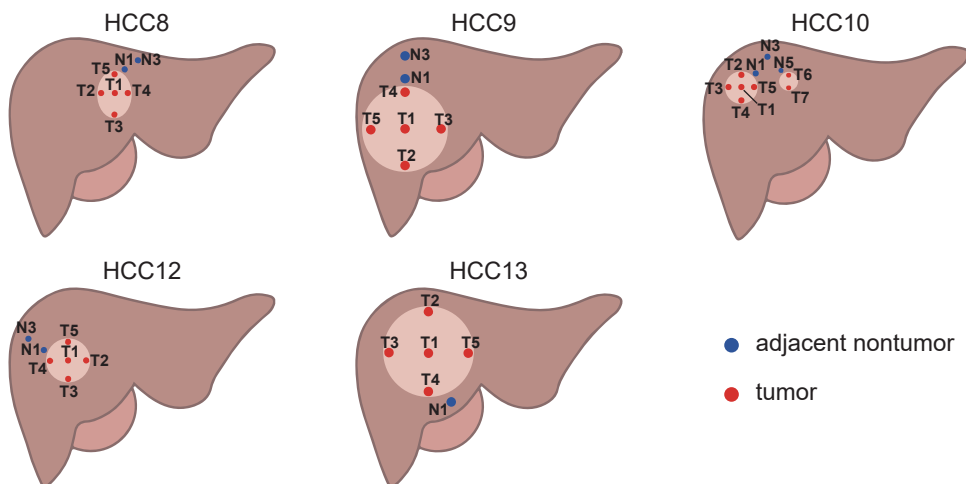
390 B. IGV plot depicting the fusion transcript formed between the *KLHDC7B-DT* gene and the HBV
391 DNA sequence. The top track represents the assembled TRA sequence, while the bottom track
392 illustrates gene annotations for the TRA sequence. The middle tracks include long-read RNA
393 sequencing from both tumor and adjacent nontumor samples, and Sanger sequencing reads for the
394 fusion transcript PCR product.

395 C. The PCR products of fusion transcript were detected by agarose gel electrophoresis.

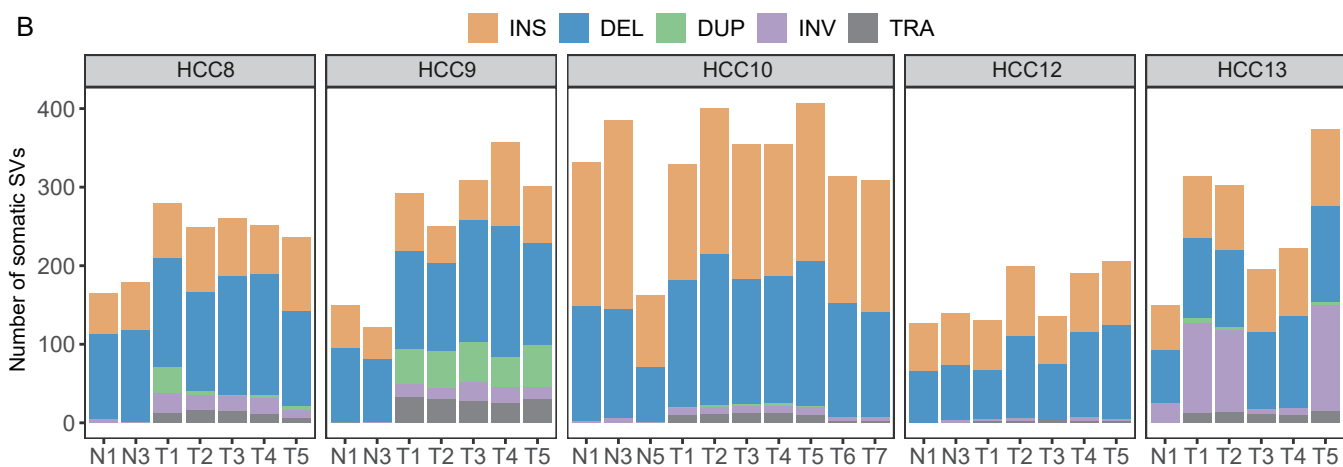
396 D. The relative expression levels of *KLHDC7B* determined using RT-PCR in samples with and
397 without the HBV-induced TRA.

398 E. The relative expression levels of *KLHDC7B-DT* determined using RT-PCR in samples with
399 and without the HBV-induced TRA.

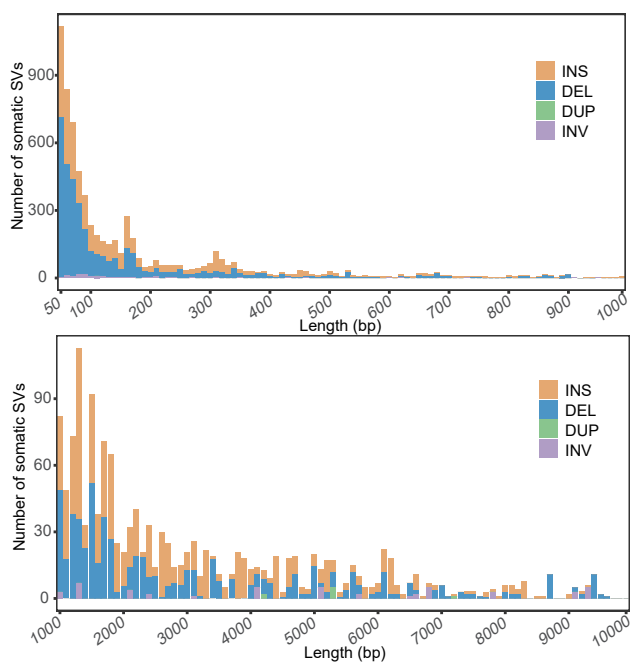
A



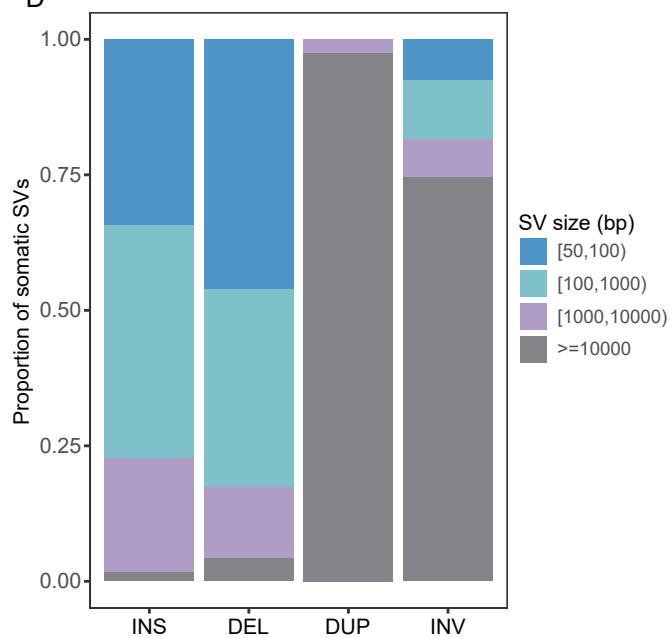
B



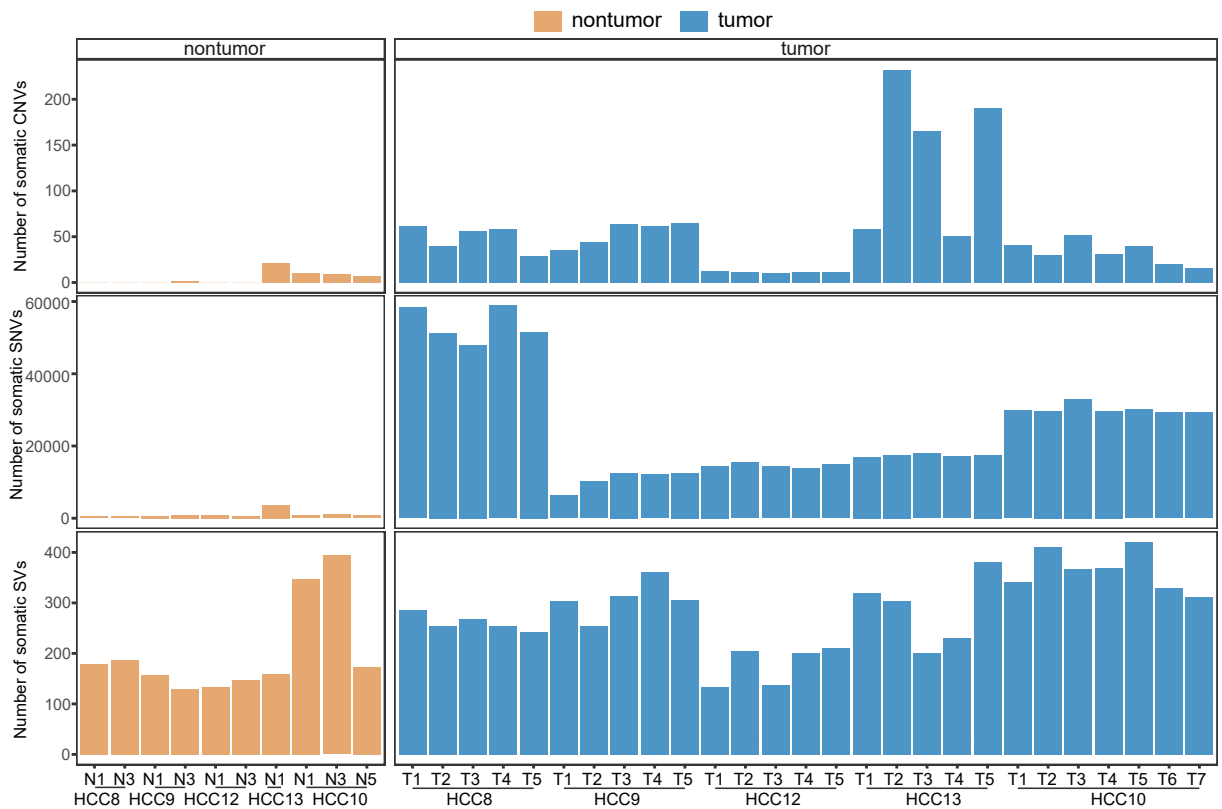
C



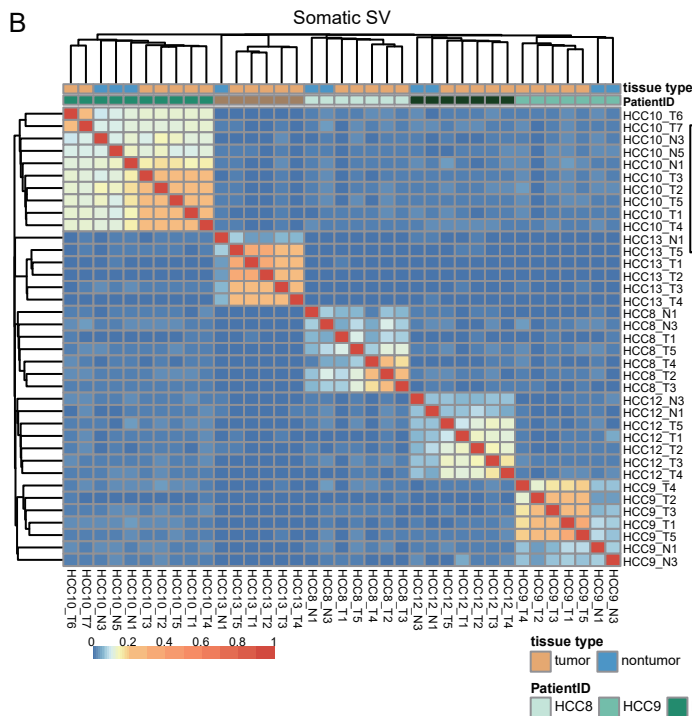
D



A



B



C

