

## Supplementary Information

### **Whole Genome Sequencing of Liver Cancers Identifies Etiological Influences on Mutation Patterns and Recurrent Mutations in Chromatin Regulators**

Akihiro Fujimoto\*, Yasushi Totoki\*, Tetsuo Abe, Keith A Boroevich, Fumie Hosoda, Ha Hai Nguyen, Masayuki Aoki, Naoya Hoshono, Michiaki Kubo, Fuyuki Miya, Yasuhito Arai, Hiroyuki Takahashi, Takuya Shirakihara, Masao Nagasaki, Tetsuo Shibuya, Kaoru Nakano, Kumiko Watanabe-Makino, Hiroko Tanaka, Hiromi Nakamura, Jun Kusuda, Hidenori Ojima, Kazuaki Shimada, Takuji Okusaka, Masaki Ueno, Yoshinobu Shigekawa, Yoshiiku Kawakami, Koji Arihiro, Hideki Ohdan, Kunihito Gotoh, Osamu Ishikawa, Shunichi Ariizumi, Masakazu Yamamoto, Terumasa Yamada, Kazuaki Chayama, Tomoo Kosuge, Hiroki Yamaue, Naoyuki Kamatani, Satoru Miyano, Hitoshi Nakagama, Yusuke Nakamura, Tatsuhiko Tsunoda, Tatsuhiro Shibata, and Hidewaki Nakagawa

\*These authors contributed equally to this work.

Correspondence should be addressed to Hidewaki Nakagawa (hidewaki@ims.u-tokyo.ac.jp) or Tatsuhiro Shibata (tashibat@ncc.go.jp).

## Supplementary Note

### Somatic point mutation and short indel calls

Read sequences were mapped by BWA<sup>1</sup> to the human reference genome (GRCh37). Possible PCR duplicated reads were removed by SAMtools<sup>2</sup> and an in-house program. After filtering by pair mapping distance, mapping uniqueness and orientation between paired reads, the mapping result files were converted into pileup format by SAMtools.

We used three kinds of read filters: set1; both read pairs were uniquely mapped with consistent orientation and pair distance (within average  $\pm 3$  s.d.), set2; at least one read pair was uniquely mapped with consistent orientation and pair distance and, set3; all uniquely mapped paired reads and set2, as described elsewhere<sup>3</sup>. Mutation calls were done using all three sets of filtered reads, and mutations identified in the all three sets were considered as candidates.

To identify point mutations, we used the following criteria; (1) non-reference calls with a frequency  $\geq 0.15$  after removing bases calls with base quality  $< 10$ , and mapping quality  $< 20$ , (2) supported by at least two base calls including one base call with base quality  $\geq 30$ , (3) a SAMtools consensus quality  $\geq 20$  and root mean square mapping quality  $\geq 40$ , (4) if three or more single nucleotide variants (SNVs) were found within any 10bp windows, or distance from nearest indel was less than 5bp, we discarded all SNVs. (5) if candidate non-coding SNVs were in a tandem repeat region suggested by tandem repeat finder<sup>4</sup>, we discarded the SNVs. (6) if candidate SNVs were in RepeatMasker repeat regions (<http://www.repeatmasker.org>) within 1Mb from the

centromeric or telomeric gaps, we discarded the SNVs. (7) if a base with consensus quality lower than 20 occurs within 3bp on either side of the target SNV, we discarded the SNVs. After SNV calling in the tumor samples, candidate SNVs were filtered based on the lymphocyte sequence of the same patient; (1) candidate SNV alleles with a frequency  $\geq 0.03$  after removing reads with base quality  $< 15$ , and mapping quality  $< 20$ , (2) depth of coverage in lymphocyte  $\leq 5$ , (3) depth of coverage in lymphocyte  $\leq 10$  and candidate SNV allele was represented in the dbSNP database v131 (<http://www.ncbi.nlm.nih.gov/projects/SNP/>).

Short indels were identified based on gaps in a read's alignment by BWA. We defined indels using the following criteria; (1) if indels were supported by a frequency  $\geq 0.1$  and  $\geq 4$  reads after removing reads with mapping quality  $< 20$ , and root mean square mapping quality  $\geq 40$ , (2) if candidate non-coding indels were in repeat regions suggested by tandem repeat finder or RepeatMasker, we discarded the indels. After indel calling in tumor samples, the candidates were filtered based on the lymphocyte sequence of the same patient using the following criteria; (1) depth of coverage in lymphocyte  $\leq 7$ , (2) for coding and non-coding region, if any indels were identified within 5bp or 10bp region in lymphocyte, respectively, the candidate indel was discarded. Gene information was based on UCSC<sup>5</sup> annotation.

The proportion of mutant alleles was estimated by the number of mutant and reference base calls. The mutation caller showed a false positive rate  $\leq 0.05$  in point mutations and  $\leq 0.1$  in short indels.

### **Estimation of false positive and false negative rates in somatic mutation calling**

To examine accuracy of our analysis pipeline, we performed Sanger sequence validation for randomly selected point mutation and indel candidates. The false positive rate was 0/124 in coding point mutations, 2/44 (4.5%; 95% CI 0.6-15.5%) in non-coding point mutations, 1/16 (6.3%; 95% CI 0.16-30.2%) in coding indels and 3/32 (9.3%; 95% CI 0.20-25.0%) in non-coding indels. Note that false positive rates in repeat regions may be higher than these values because mutations in repeat regions were difficult to examine by PCR-Sanger sequencing method.

False negative rate (FNR) is very important, but difficult to estimate in the cancer genome sequence, because it is difficult to obtain sufficient number of already identified somatic mutations. Therefore, we used SNP array genotype for FNR estimation. We genotyped one cancer (HB6 tumor) and lymphocyte sample from another patient (HB5 lymphocyte) with Illumina OmniExpress array, and then ran our analysis pipeline by using HB6 tumor and HB5 lymphocyte data. We compared the results from our analysis pipeline with these from genotyping arrays. If our analysis pipeline did not call HB6 tumor-specific SNP alleles, i.e. heterozygous in HB6 tumor and homozygous in HB5 lymphocyte, we considered that SNP to be a false negative in our pipeline. By the comparison of HB6 tumor data and HB5 lymphocyte data, we identified 61,077 SNPs that were homozygous with reference allele on GRCh37 in HB5 lymphocyte and heterozygous in HB6 tumor by the genotyping array. Of these, 56,074 were identified by WGS data analysis, therefore FNR was estimated to be 8.2% (95% C.I. 8.0-8.4%). Four-fifth of false negatives (6.5%) were caused by insufficient depth of

coverage in HB5 lymphocyte.

### **Identification of common mutation in MCTs**

We compared somatic mutation sites and mutation patterns of the two pairs of MC tumors (HC3 and HC7). In protein-coding regions, no common point somatic mutation was identified. Three genes, *ATM*, *FSIP2*, and *LRFN5*, were mutated in both the HC3-1 and HC3-2 tumors, but the mutations occurred at different positions. In non-coding regions, WGS identified 30 and 37 common somatic point mutations and indels in the HC3 and HC7 MCT HCC genomes, respectively. However, most of these occurred in repetitive regions and all candidates that could be sequenced by Sanger sequencing method (n=20) were found to be germline variants. This high false positive rate is due to the use of a common lymphocyte sample as a control. False negatives of germline variants from the lymphocyte sample could cause common false positive somatic mutations in both tumors. There was also no common structural alteration found in these MCT HCC genomes. These findings strongly suggest that these synchronous MCT tumors developed through an accumulation of a completely different set of genetic alterations.

### **Permutation test on the distance between MCTs**

To evaluate the distance between the paired MCTs, we performed permutation tests based on the PCA (**Fig. 2c and 2d**). We randomly selected two sets of tumors from the 25 tumors and calculated the three-dimensional distance between tumors in a set. Then

the average distance of two sets was calculated. We repeated this process 100,000 times to obtain the null distribution. The average distance of the MC pair was tested under this null distribution.

### **Association between mutations in chromatin regulators and clinical factors**

We examined the association of chromatin regulator mutations with several clinical factors. Tests of significance for tumor size, Edmondson grade, and liver fibrosis were performed by the Pearson's correlation test. Tests for portal vein invasion and hepatic vein invasion were carried out by the Fisher's exact test. We found that the stage of liver fibrosis in the background liver was associated with mutations and CNAs of the chromatin regulator genes ( $P$ -value = 0.026) or ARID family genes ( $P$ -value = 0.0086) (**Supplementary Table 11**). We also found that the presence of invasion to hepatic vein in tumors was associated with mutations and CNAs of *ARID1A* ( $P$ -value = 0.042) (**Supplementary Table 11**). Although we did not observe any significant association after the Bonferroni correction, mutations of the chromatin regulators in HCC may contribute to poor prognosis of HCC patients.

### **Copy number alternation (CNA)**

To investigate copy number alterations, we analyzed the ratio of the depth of coverage across 5kb windows of the tumor sample to that of the matched lymphocytes. In total, we detected 294 deleted regions ( $\log_2R$  ratio  $\leq -1$ ) and 20 amplified regions ( $\log_2R$  ratio  $\geq 2$ ), of which 39 deletions and one amplification occurred recurrently. We also

observed 397 low-level losses ( $\log_2R$  ratio  $\leq -0.6$ ) and one low-level gain ( $\log_2R$  ratio  $\geq 1$ ) with a frequency greater than 15% ( $n \geq 5$ ), most of which have been reported previously<sup>6,7</sup> (**Supplementary Table 3**). Among the significantly mutated genes (**Table 1, Supplementary Table 7**), *TP53* harbored both a missense point mutation and low-level loss in five tumors, while HC6 had a 1bp-deletion in the coding region and low-level loss of *ARID1A*. Additionally, we identified several deletions or low-level losses of the chromatin regulators; *ARID1A* in 5 tumors, *MLL* in 4 tumors, *ARID1B* in one tumor and *MLL3* in one tumor (**Fig. 4a**).

### Genomic rearrangement

Using inconsistently-mapped reads, we detected an average of 20.8 genomic rearrangements per tumor with a wide range among samples (0-59) (**Supplementary Table 2**). Among the 561 PCR-validated rearrangements (**Supplementary Table 2, Supplementary Fig. 11**), breakpoints of 351 rearrangements were in intragenic regions and 30 rearrangements involved possible gene fusion events. Thirteen genes (*ACVR2A*, *AUTS2*, *CADPS2*, *EPAS1*, *FHIT*, *MBD5*, *NR3C1*, *PLAG1*, *PLEKHG3*, *RAD51L1*, *RBFOX3*, *STK39* and *XRCC4*) had rearrangement breakpoints in multiple HCCs. Seventy genes had both protein-altering point mutations and a rearrangement. Fifteen genes found to have rearrangement breakpoints also contained recurrent protein-altering point mutations (**Supplementary Tables 4 and 7**), including *CSMD1*, the chromatin regulator *ARID2* and two significantly mutated genes, *ALB* and *ATM*. We did not observe any recurrent genomic rearrangements that generated an

in-frame fusion gene in this cohort. There is no apparent difference of the observed genomic structure changes between HBV- and HCV-related HCCs (**Supplementary Fig. 2**), and there is no common rearrangement in the two sets of the MCTs (HC3 pair and HC7 pair) (**Fig. 2a**).

### **Recurrent somatic mutations in non-coding regions**

We tested the number of mutations in the 5'UTR, 3'UTR, promoter regions (1kb region from transcription start sites) and non-coding RNA regions while considering the regional differences of the mutation rate and region lengths. As a result, 74 regions were identified with  $\geq 3$  point mutations or indels. Of these, 63 were considered significant after adjustment for multiple testing ( $q$ -value  $\leq 0.05$ ), which included regulatory regions of coding genes as well as multiple small and long intergenic non-coding RNAs (**Fig 2a, Supplementary Table 12**). Among them, *PHF1*, which is involved in DNA damage response and epigenetic modification and reported to be rearranged in rare uterine sarcomas<sup>8</sup>, was significantly mutated in the promoter region as well as the protein-coding regions. Additionally, *MED1*, which is reported to play an important role in liver regeneration and liver carcinogenesis in animal models<sup>9</sup>, was also significantly mutated in the 3'UTR as well as its protein-coding region.

### **References**

1. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).
2. Li, H. et al. The Sequence Alignment/Map format and SAMtools.



- Bioinformatics* **25**, 2078-9 (2009).
3. Totoki, Y. et al. High-resolution characterization of a hepatocellular carcinoma genome. *Nat Genet* **43**, 464-9 (2011).
  4. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res* **27**, 573-80 (1999).
  5. Fujita, P.A. et al. The UCSC Genome Browser database: update 2011. *Nucleic Acids Res* **39**, D876-82.
  6. Laurent-Puig, P. & Zucman-Rossi, J. Genetics of hepatocellular tumors. *Oncogene* **25**, 3778-86 (2006).
  7. Tsai, W.L. & Chung, R.T. Viral hepatocarcinogenesis. *Oncogene* **29**, 2309-24 (2010).
  8. Micci, F., Panagopoulos, I., Bjerkehagen, B. & Heim, S. Consistent rearrangement of chromosomal band 6p21 with generation of fusion genes JAZF1/PHF1 and EPC1/PHF1 in endometrial stromal sarcoma. *Cancer Res* **66**, 107-12 (2006).
  9. Matsumoto, K. et al. Critical role for transcription coactivator peroxisome proliferator-activated receptor (PPAR)-binding protein/TRAP220 in liver regeneration and PPARalpha ligand-induced liver tumor development. *J Biol Chem* **282**, 17053-60 (2007).

## Supplementary Figures Contents

**Supplementary Figure 1:** Average depth of coverage with uniquely mapped reads in each WGS sample.

**Supplementary Figure 2:** Difference in mutation counts between HBV- and HCV-related tumors.

**Supplementary Figure 3:** Somatic substitution pattern of each tumor.

**Supplementary Figure 4:** The number of somatic substitutions on the transcribed strand and the untranscribed strand.

**Supplementary Figure 5:** The repair on the transcribed strand.

**Supplementary Figure 6:** PCA of somatic substitution patterns.

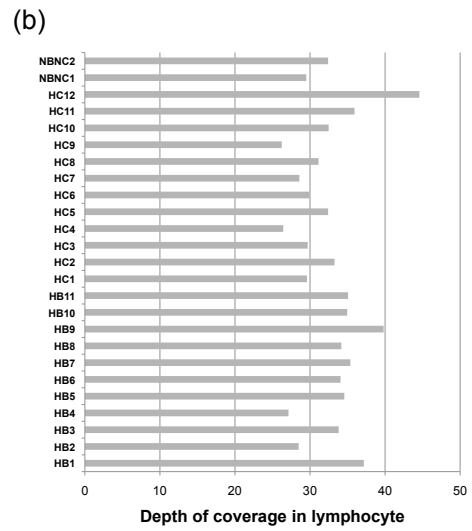
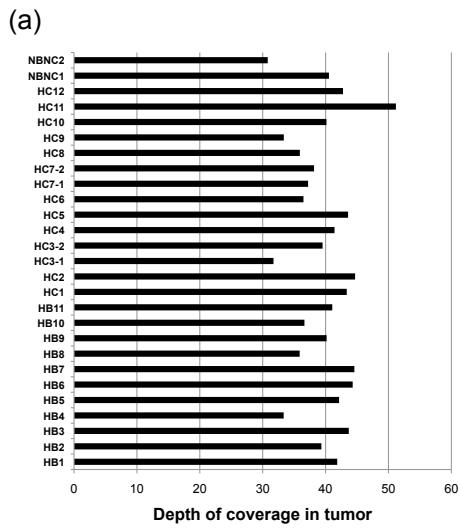
**Supplementary Figure 7:** Reduction of the target gene expression by siRNA.

**Supplementary Figure 8:** Down-regulation of *ERRF1* and *ARID1B* genes did not affect proliferation of JHH5 cells which lack expression of these two genes.

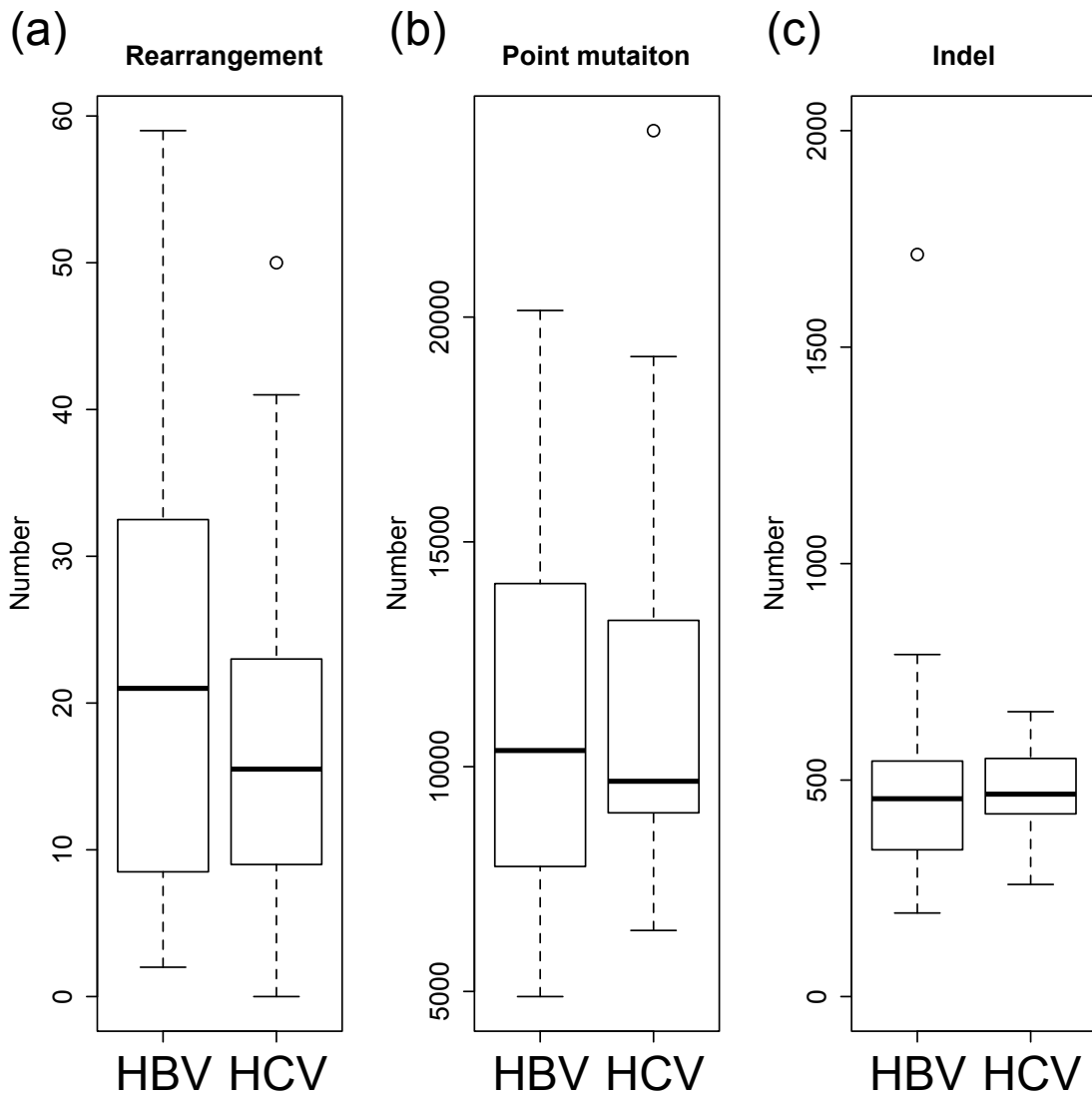
**Supplementary Figure 9:** Integration sites on HBV genome.

**Supplementary Figure 10:** HBV genome integration sites into the *TERT* region.

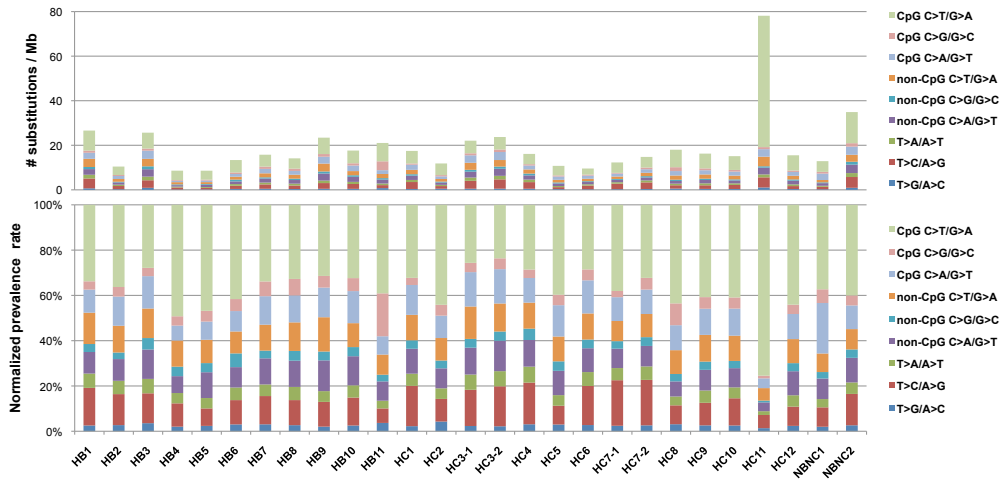
**Supplementary Figure 11:** Graphical representation of 27 HCC genomes.



**Supplementary Figure 1:** Average depth of coverage with uniquely mapped reads in each sample. We calculated depth of coverage of each sample by using uniquely mapped reads after removing PCR-duplication.



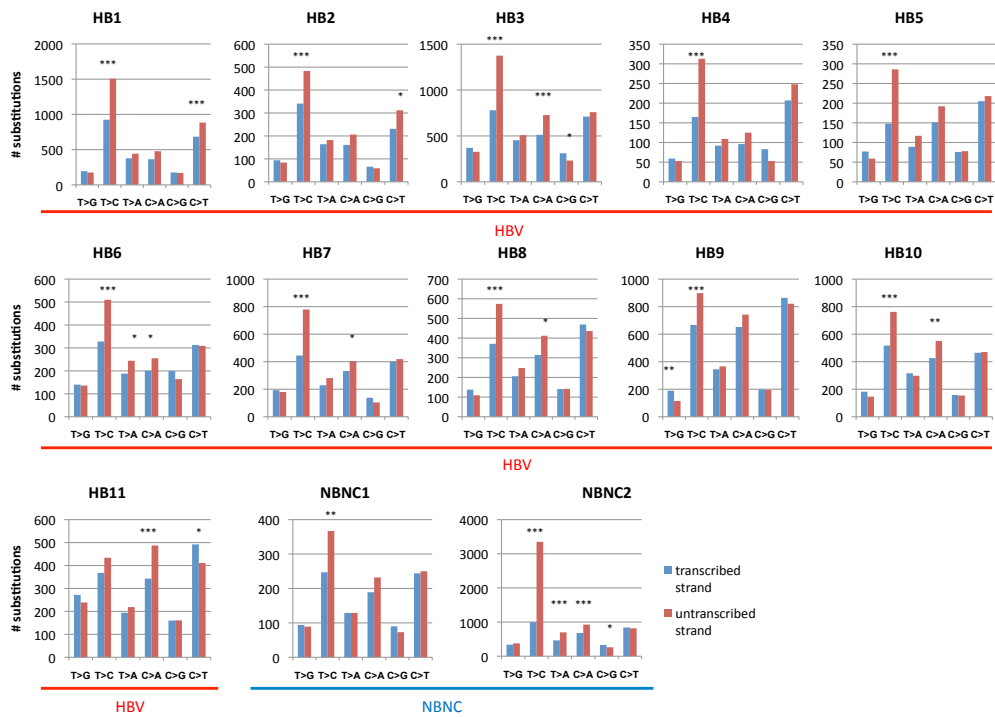
**Supplementary Figure 2:** Difference in mutation counts between HBV- and HCV-related tumors. (a) Number of rearrangements in HBV- and HCV-related tumors. (b) Number of point mutations in HBV- and HCV-related tumors. (c) Number of indels in HBV- and HCV-related tumors. No significant difference was observed in these comparisons.



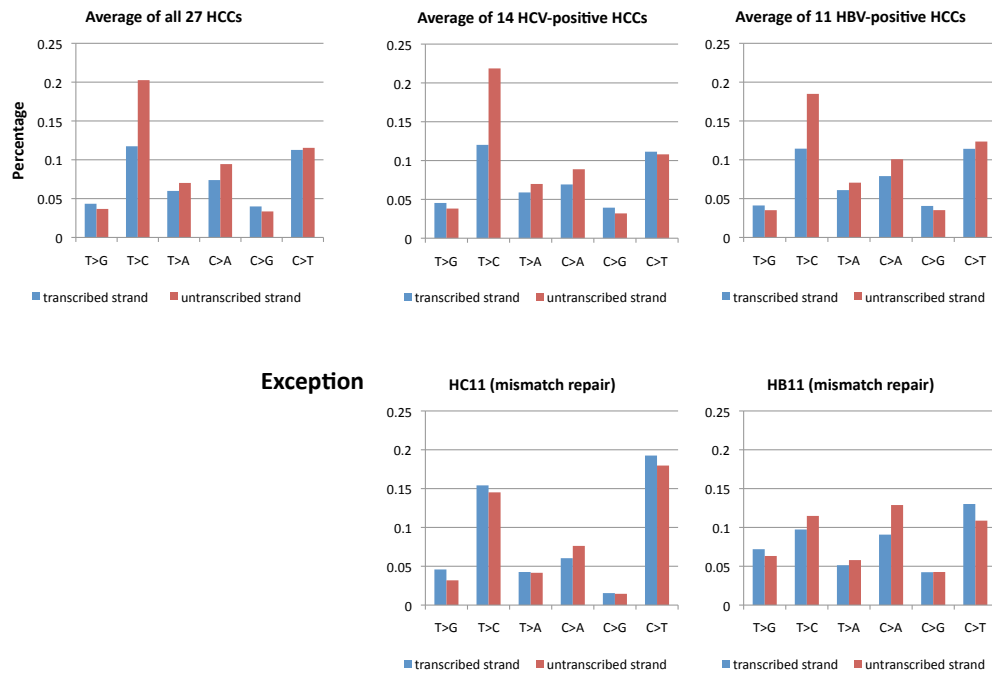
**Supplementary Figure 3:** Somatic substitution pattern of each tumor. (Top) Number of somatic substitutions in each of the nine classes per Mb in 27 HCC genomes. (Bottom) Frequency of the nine classes of somatic substitutions in 27 HCC genomes.



**Supplementary Figure 4:** The number of somatic substitutions on the transcribed strand and the untranscribed strand (**a**) in 14 HCV-related HCCs. \*\*\*;  $P$ -value <  $10^{-6}$ , \*\*;  $P$ -value < 0.0001, \*;  $P$ -value < 0.01.

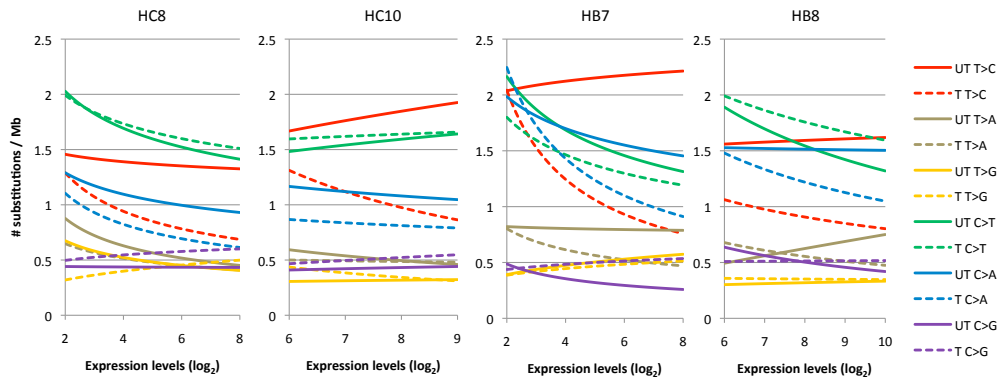


**Supplementary Figure 4:** The number of somatic substitutions on the transcribed strand and the untranscribed strand (**b**) in 11 HBV-related HCCs, and 2 HCCs without HBV and HCV infections (NBNC). \*\*\*,  $P$ -value <  $10^{-6}$ , \*\*,  $P$ -value < 0.0001, \*,  $P$ -value < 0.01.

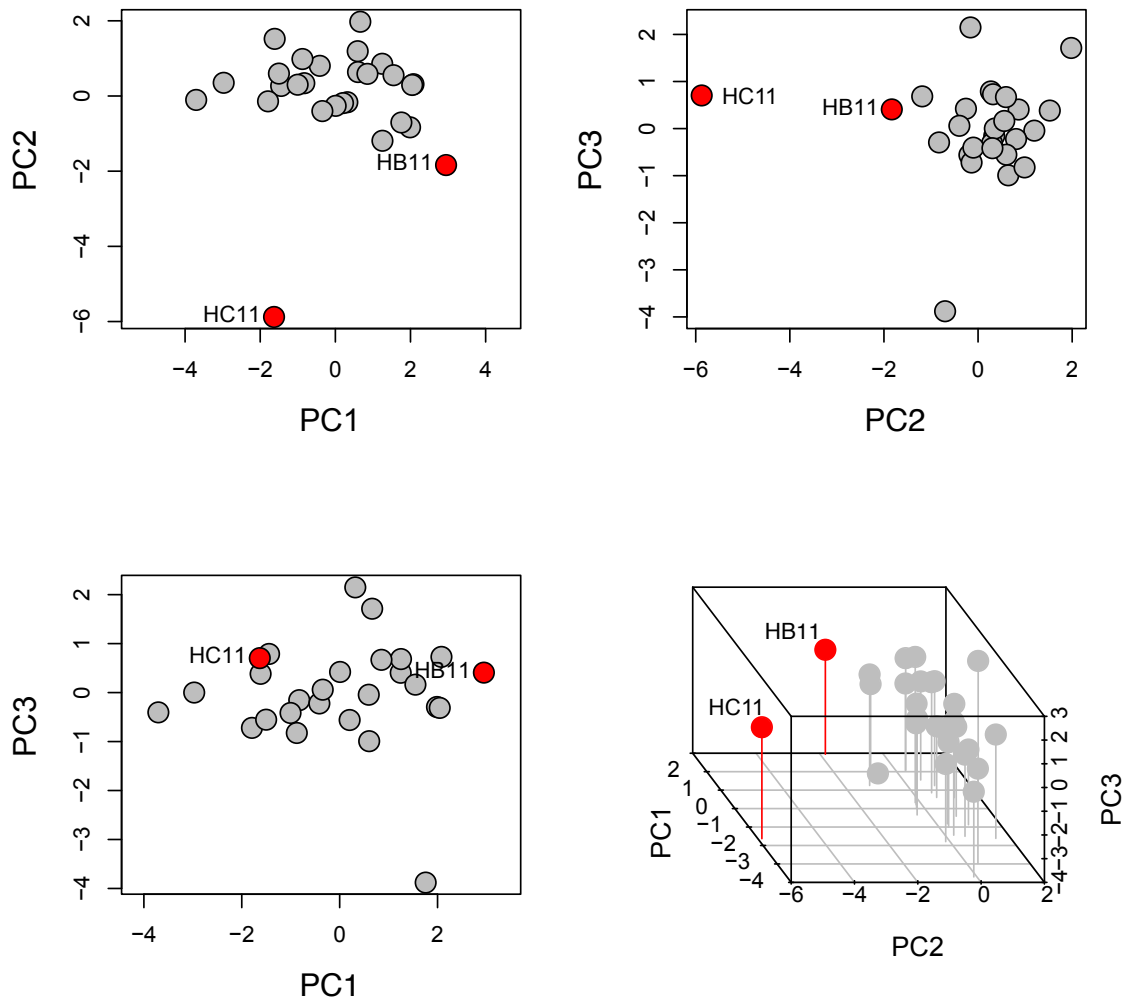


**Supplementary Figure 4:** The number of somatic substitutions on the transcribed strand and the untranscribed strand. **(c)** The average percentage of somatic substitutions in all 27 HCCs, 14 HCV-related HCCs and 11 HBV-related HCCs (top). Two exceptions (bottom). TCR did not occur in the mismatch-repair deficient tumor (HC11) with somatic nonsense mutation of *MLH1*. Another case (HB11) had a familial disposition to cancer and exhibited a distinct mutation signature (increased indels, suggesting mismatch-repair deficiency, less dominant T>C/A>G transitions, and a decreased effect of TCR at T>C transition).

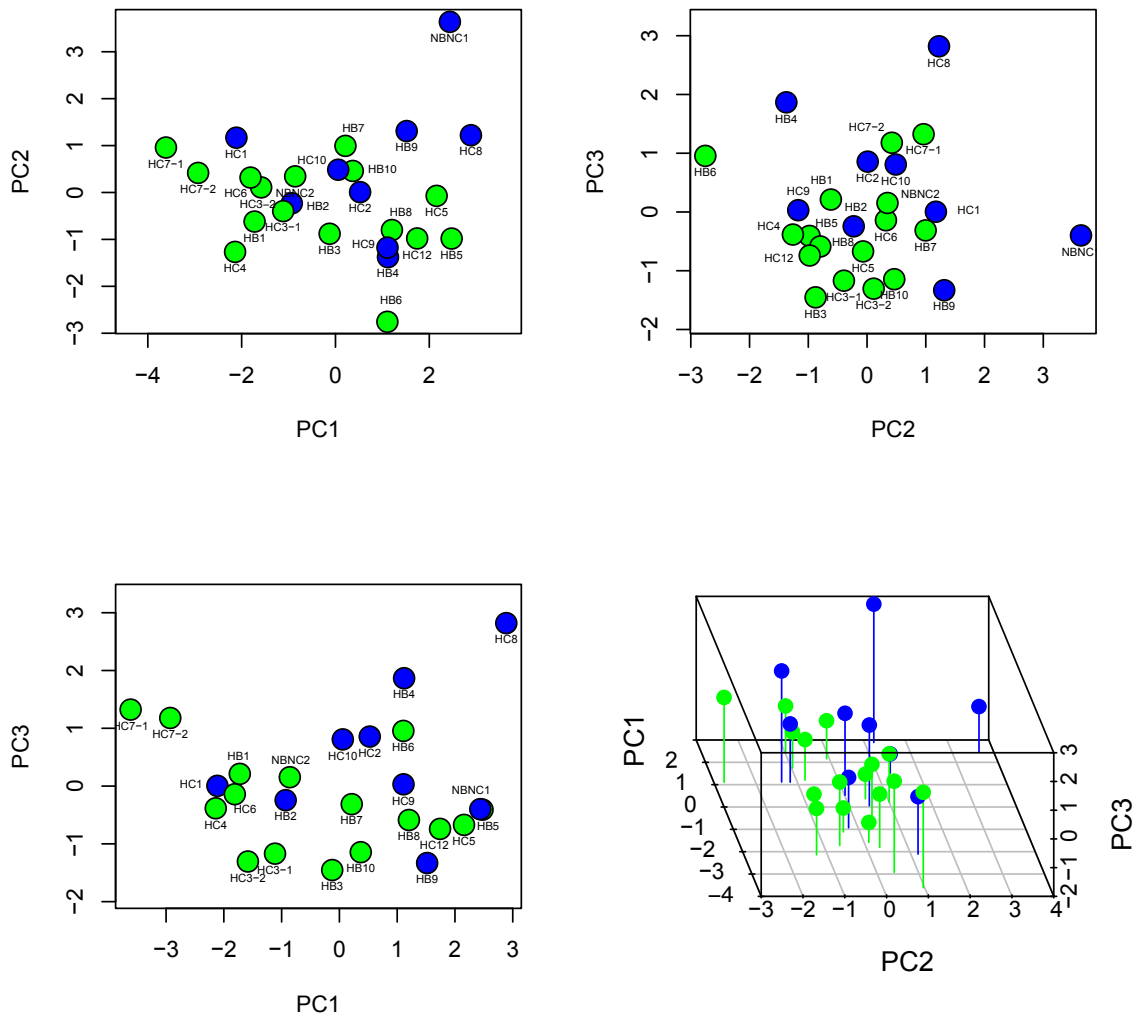




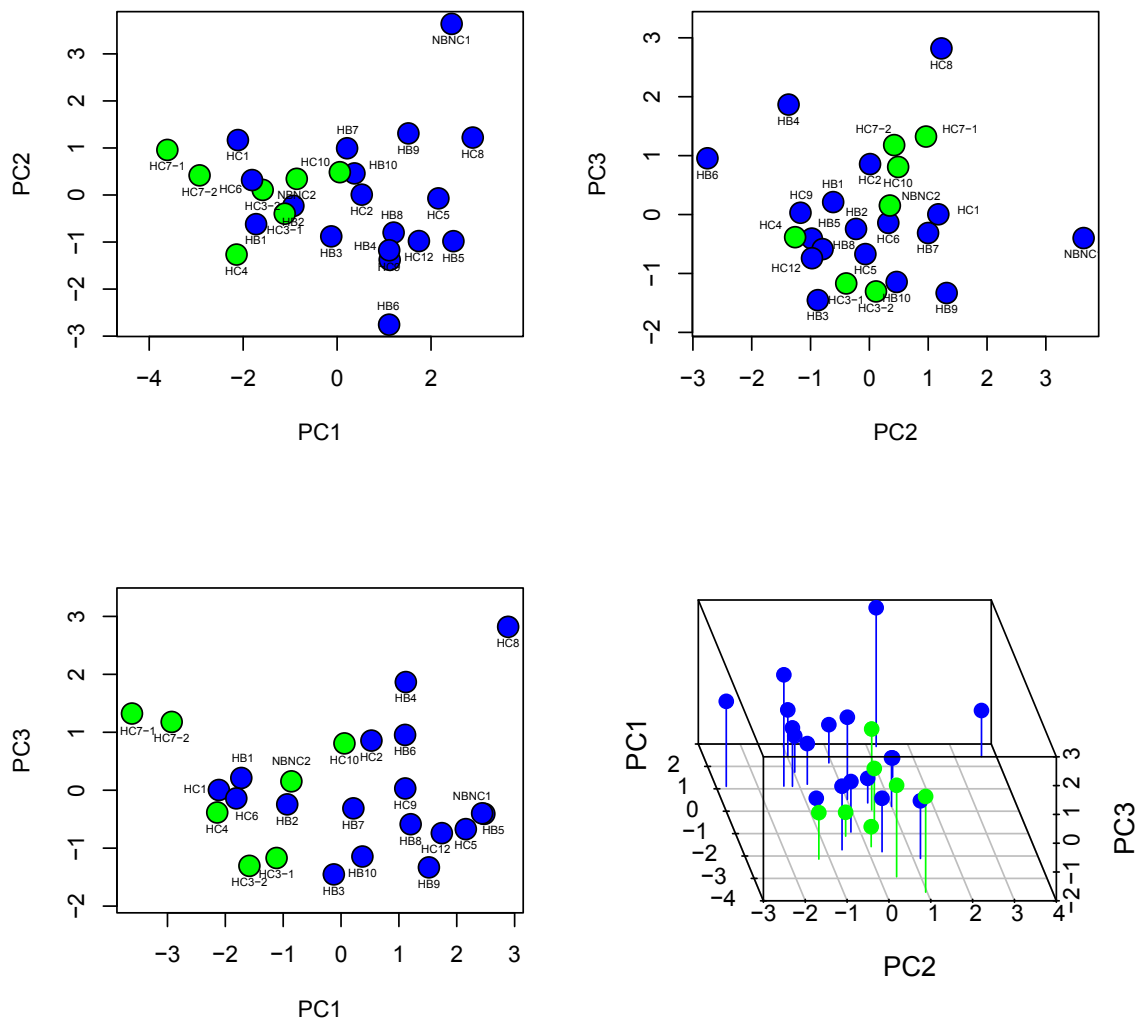
**Supplementary Figure 5:** Repair on the transcribed strand. Fitted curves showing the effect of gene expression and strand bias on substitution prevalence. We used Agilent microarray expression data (Whole Human Genome (8×60K) Oligonucleotide Microarray) in this TCR analysis and Expression level (log<sub>2</sub>) indicate Agilent microarray intensity level units with log<sub>2</sub> scale. UT, untranscribed strands; T, transcribed strands.



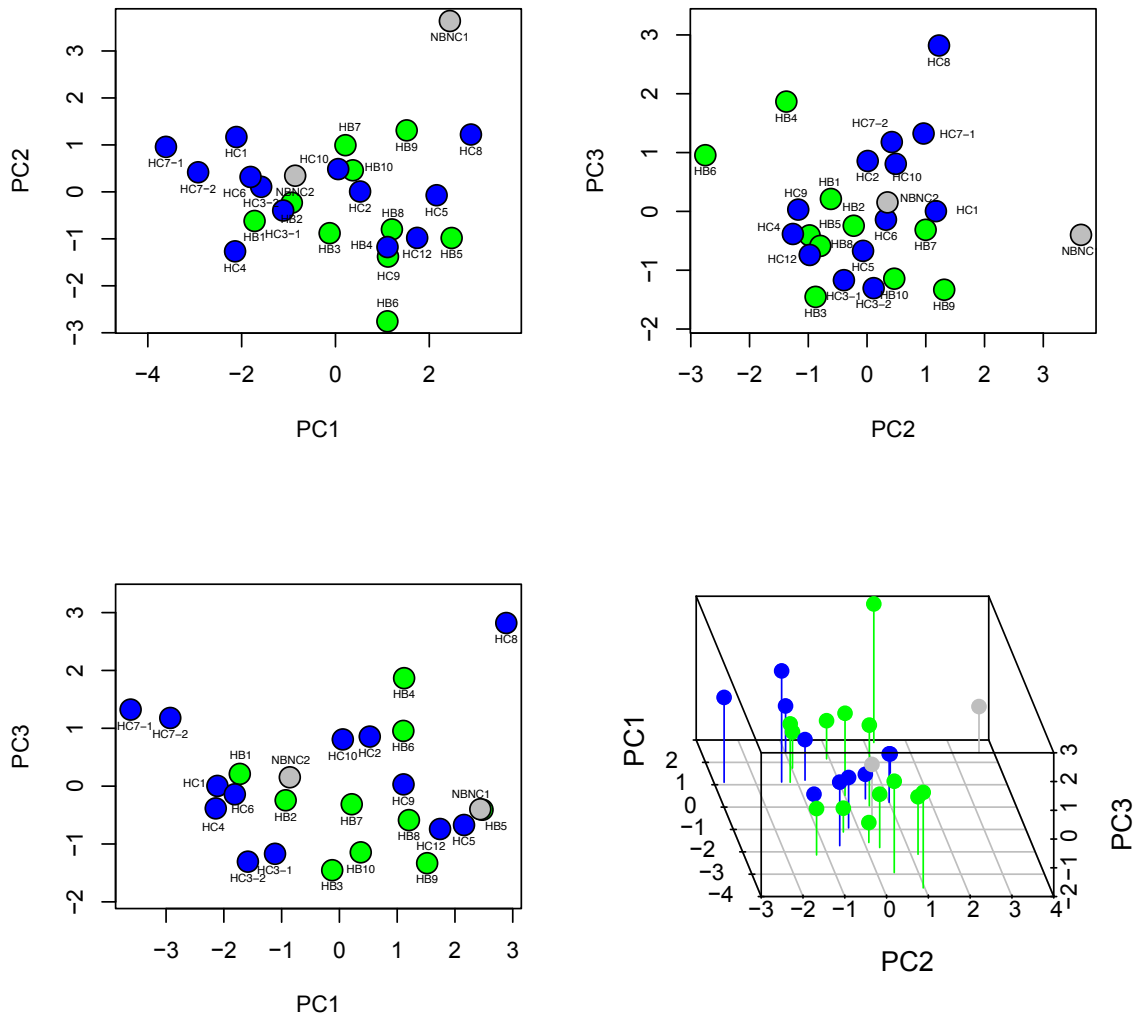
**Supplementary Figure 6: PCA of somatic substitution patterns. (a)** HB11 and HC11 are outliers, probably because they apparently show mismatch-repair deficient phenotype due to a somatic nonsense mutation (p.E234\*) of *MLH1* for HC11 and unknown factors for HB11



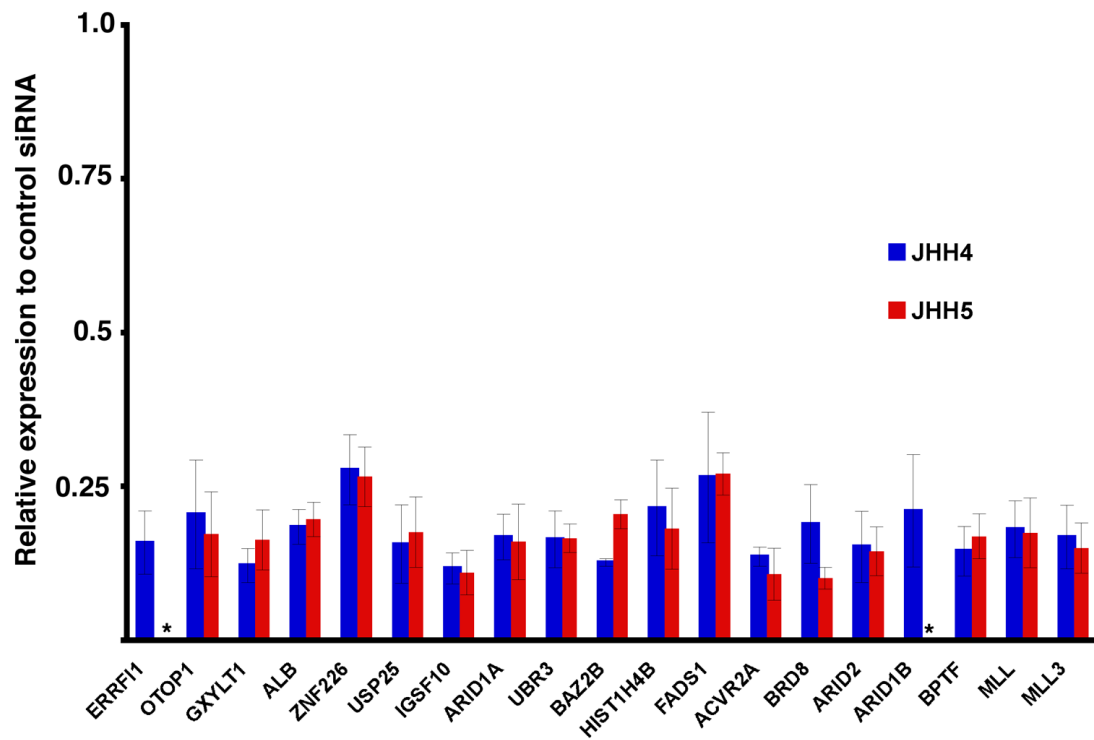
**Supplementary Figure 6:** PCA of somatic substitution patterns. **(b)** Comparison between tumors from habitual alcohol drinking (green) and non-drinking (blue) HCC patients. Habitual alcohol drinking was significantly correlated with somatic substitution patterns in HCCs ( $P$ -value = 0.028).



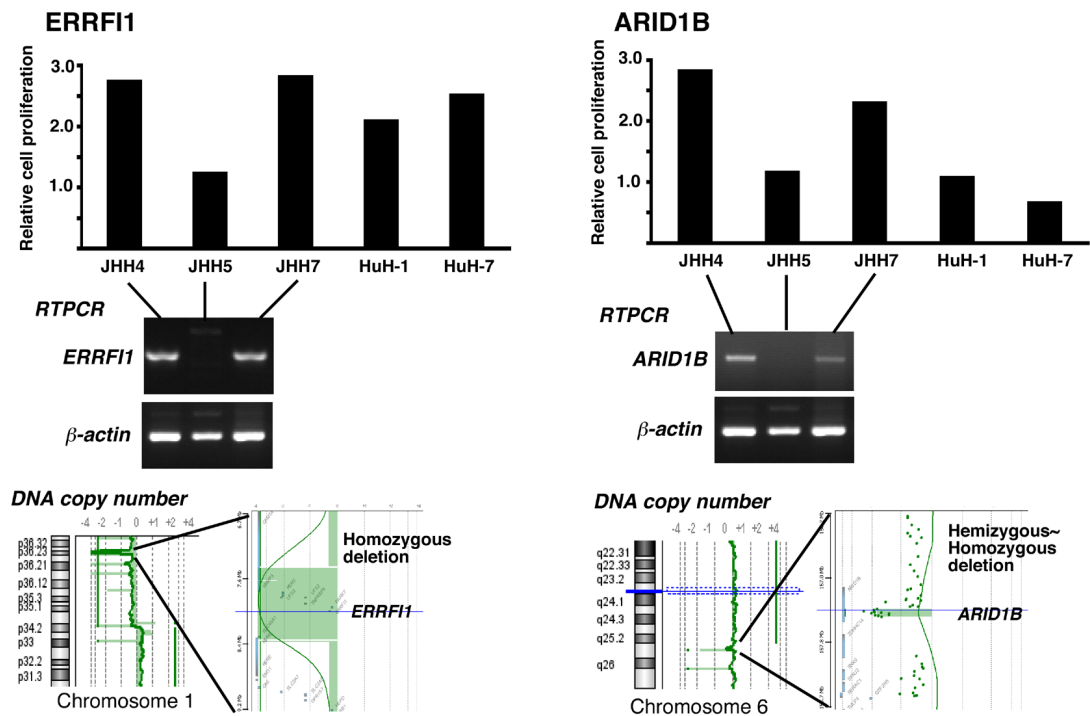
**Supplementary Figure 6:** PCA of somatic substitution patterns. **(c)** Comparison between synchronous/metachronous multiple liver nodules (green) and the other (blue) HCCs. The occurrence of synchronous or metachronous multiple HCCs, most of which are likely to be MC, was significantly correlated with somatic substitution patterns ( $P$ -value = 0.016)



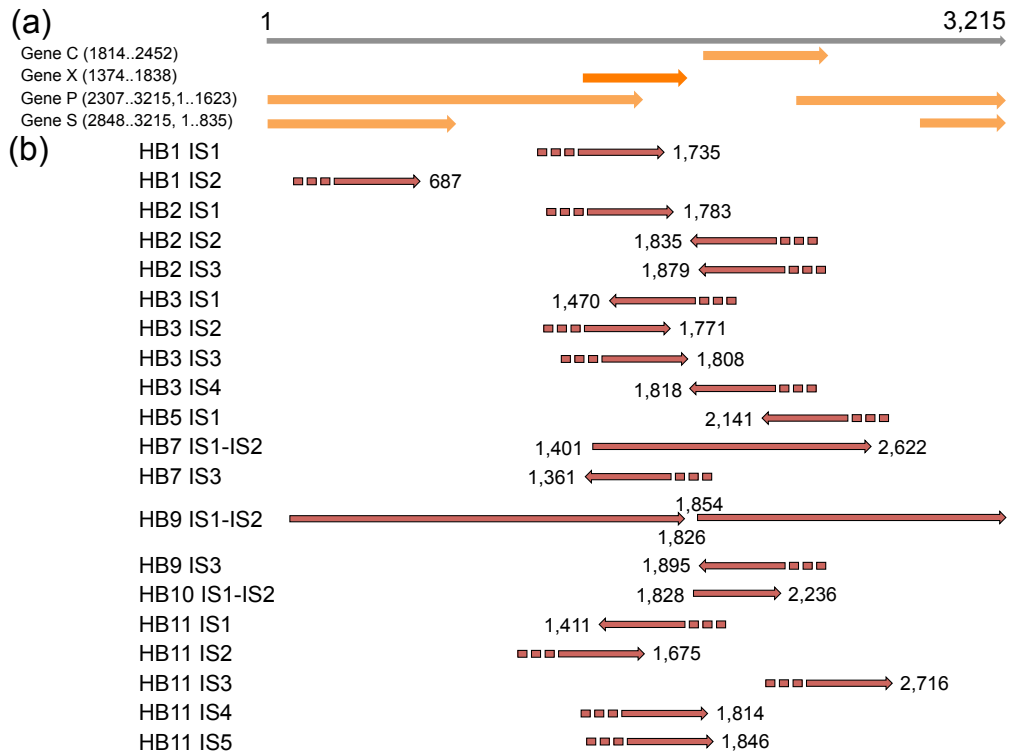
**Supplementary Figure 6:** PCA of somatic substitution patterns. **(d)** Comparison between HBV-related (green), HCV-related (blue), and NBNC (grey) HCCs. Different substitution patterns between the HBV-, HCV-related and NBNC HCCs were observed ( $P$ -value = 0.091, canonical correlation analysis, and  $P$ -value = 0.020, correlation analysis on the second component: PC2).



**Supplementary Figure 7:** Reduction of the target gene expression by siRNA. Relative expression of the target genes compared to the control siRNA treated cells was measured by quantitative RT-PCR. Asterisk shows the genes whose expression was not detected because of homozygous gene deletion (see **Supplementary Figure 8**). Data indicates mean +/- SD.

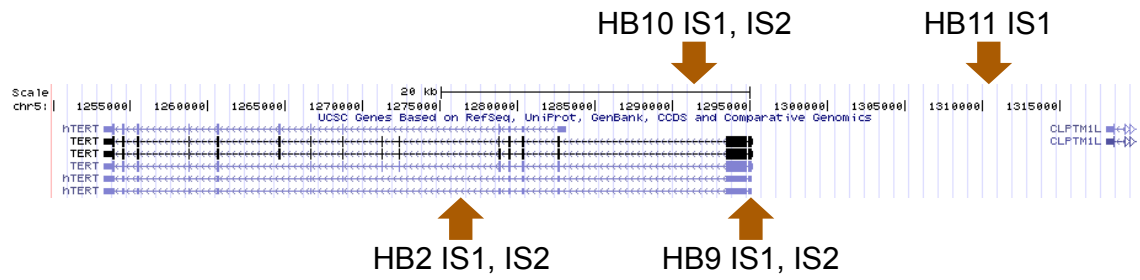


**Supplementary Figure 8:** Down-regulation of *ERRF1* and *ARID1B* genes did not affect proliferation of JHH5 cells which lack expression of these two genes. (Top) Relative cell proliferation of five HCC cell lines treated with siRNAs against *ERRF1* and *ARID1B* genes. Note that down-regulation of these genes did not affect proliferation of JHH5 cells. (Middle) RT-PCR analysis of *ERRF1*, *ARID1B* and  $\beta$ -Actin (control) genes in three HCC cell lines. *ERRF1* and *ARID1B* mRNAs were not detected in JHH5 cells. (Bottom) Array CGH analysis of *ERRF1* and *ARID1B* loci in JHH5 cells. A homozygous deletion including the whole *ERRF1* gene was detected on *1p36*. A hemizygous deletion with a minute homozygous deletion within the *ARID1B* gene was detected on *6q25*. A high-density oligonucleotide array containing 244,000 independent probes (Human Genomic Microarray 44B, Agilent Technologies) was used for array comparative genomic hybridization employing the manufacturer's protocol.

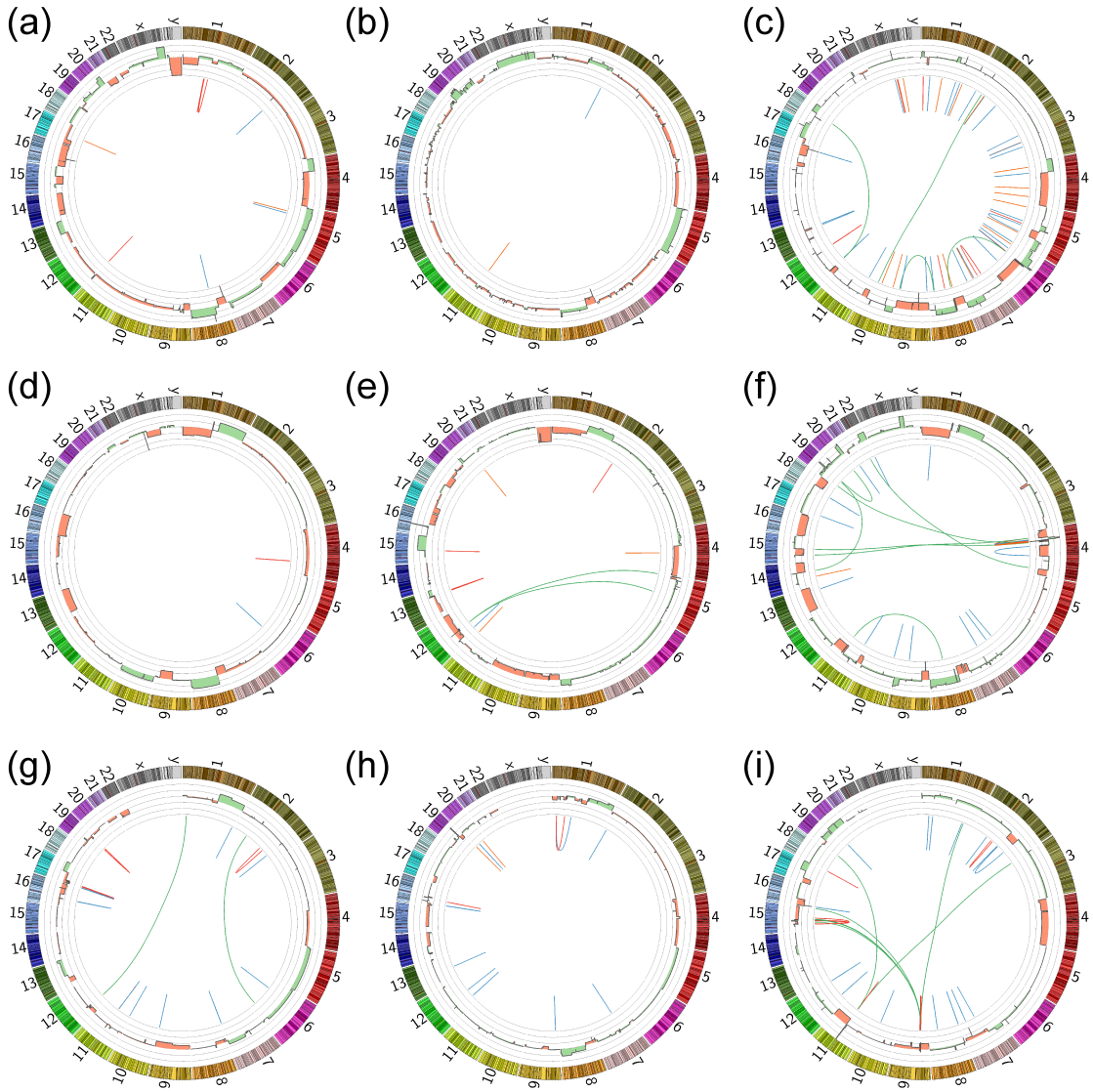


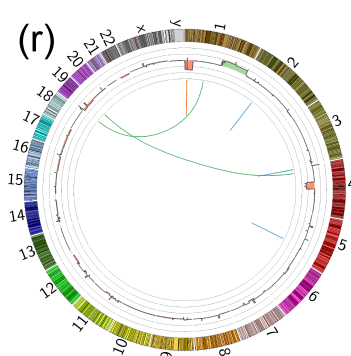
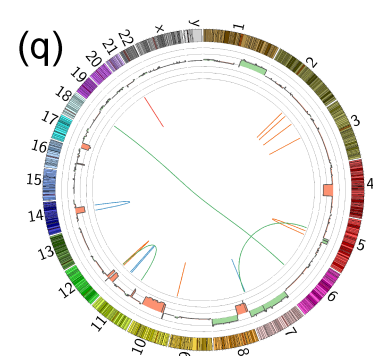
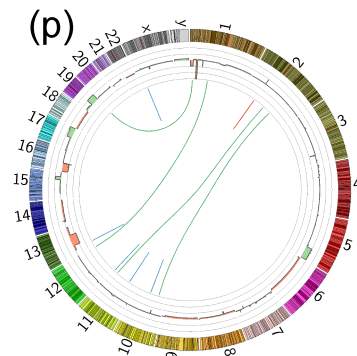
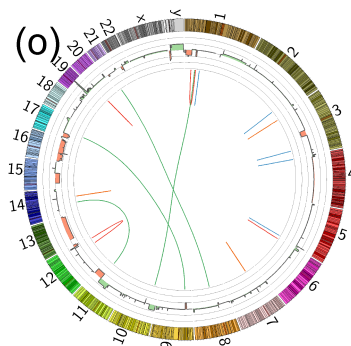
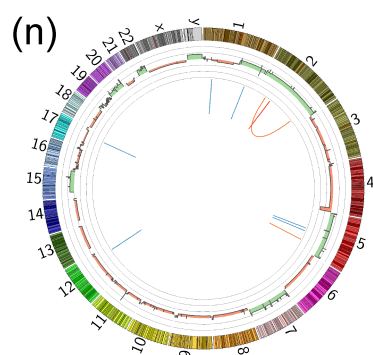
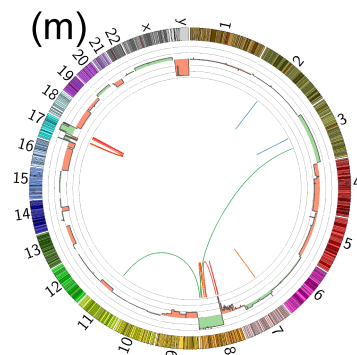
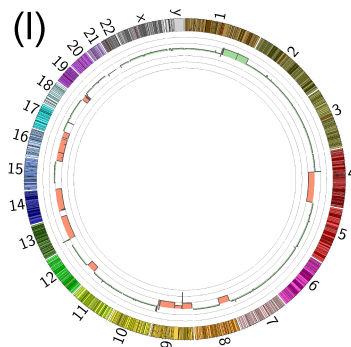
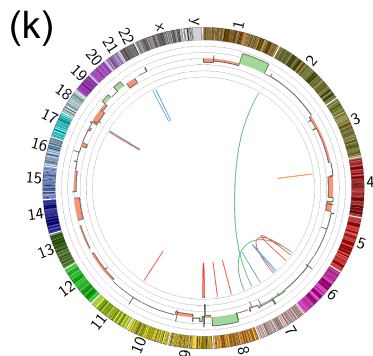
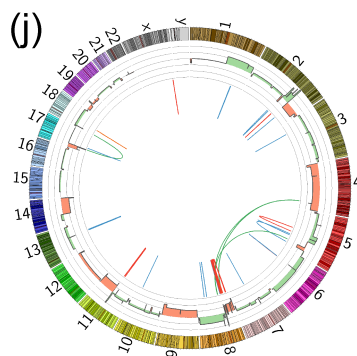
**Supplementary Figure 9:** Integration sites on HBV genome. **(a)** Structure of HBV genome. Arrows indicate gene structure. **(b)** Integration site on HBV genome. Arrows represents direction and location of integrated HBV genome. Integrated sequences were completely identified on HB9 IS1-IS2 and HB10 IS1-IS2. Integration site in host genome suggested that HB2 IS1 and IS2, and HB2 IS3 and IS4 are both ends of single integration events (**Supplementary Table 13**).

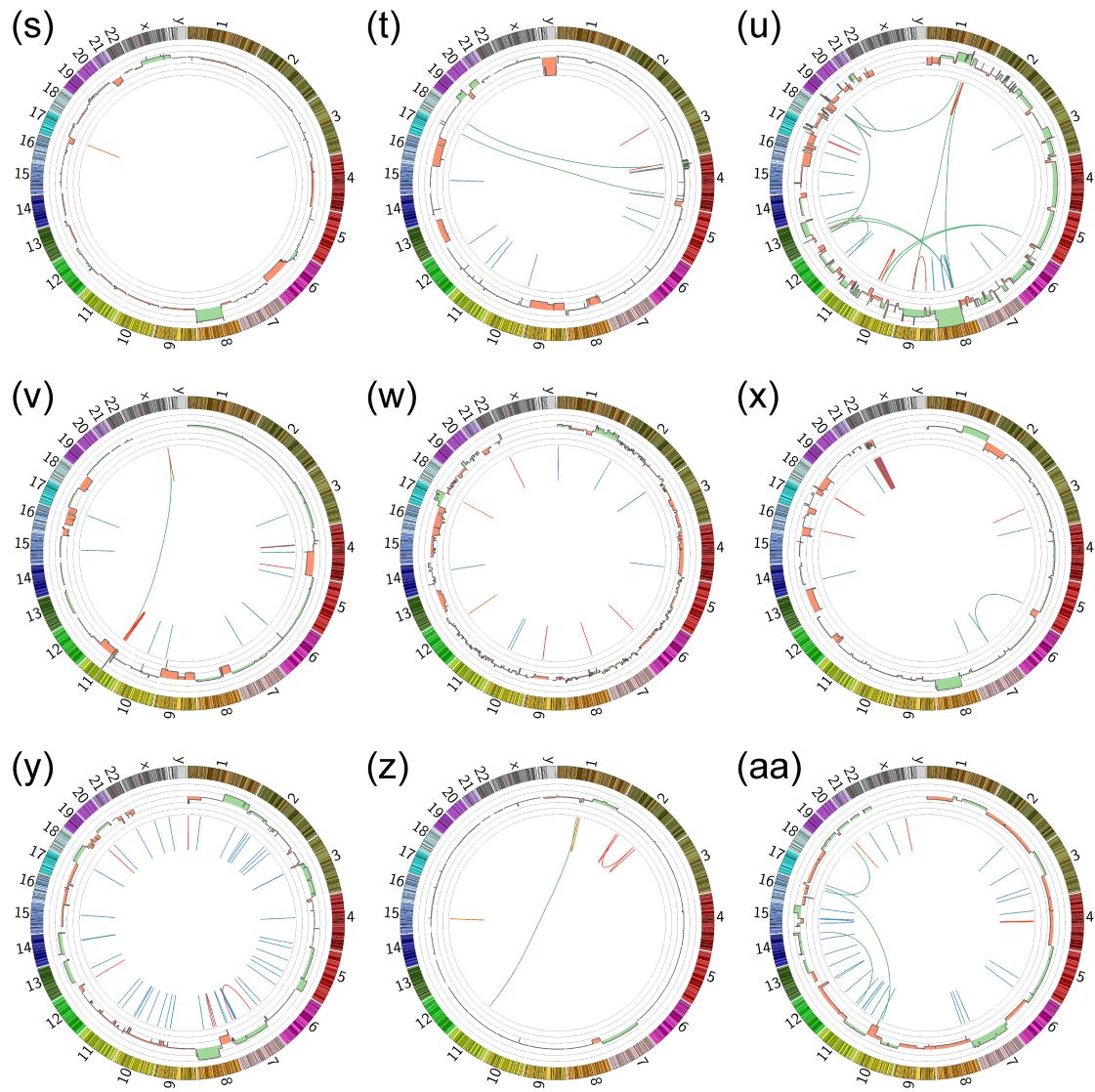




**Supplementary Figure 10:** HBV genome integration sites into the *TERT* locus. HBV genome integration was observed within or in the upstream region of the *TERT* gene in four HBV-related HCCs.







**Supplementary Figure 11:** Graphical representation of 27 HCC genomes. (a) HB1, (b) HB2, (c) HB3, (d) HB4, (e) HB5, (f) HB6, (g) HB7, (h) HB8, (i) HB9, (j) HB10, (k) HB11, (l) HC1, (m) HC2, (n) HC3-1, (o) HC3-2, (p) HC4, (q) HC5, (r) HC6, (s) HC7-1, (t) HC7-2, (u) HC8, (v) HC9, (w) HC10, (x) HC11, (y) HC12, (z) NBNC1, (aa) NBNC2. Each circle plot represents validated rearrangements (inner arcs), and copy number alternations (inner rings). In rearrangements, green, blue, orange and red lines show translocation, deletion, inversion and tandem duplication. Copy number gain and loss region were shown in green and red.

**Supplementary Table 1: Clinical and pathological features of 27 HCCs from 25 patients analyzed by whole genome sequencing**

ID	Age	Gender	Hepatitis	Alcohol drinking	Tumor size (mm)	Pretreatment	Edmondson grade <sup>d</sup>	Portal vein invasion
HB1	56	M	HBV	+	33	-	II	-
HB2	46	M	HBV	-	24	-	II	-
HB3	74	M	HBV	+	60	-	III	+
HB4	61	M	HBV	-	30	-	II-III	-
HB5	38	M	HBV	+	18	-	II-III	-
HB6	57	M	HBV	+	30	-	III	+
HB7	58	M	HBV	+	22	-	II	+
HB8	68	F	HBV	+	45	TACE <sup>b</sup>	II-III	+
HB9	60	M	HBV	-	80	-	II-III	+
HB10	56	M	HBV	+	63	-	III	+
HB11	41	F	HBV	+	30	-	II-III	-
HC1	62	M	HCV	-	87	-	II-III	+
HC2	71	M	HCV	-	45	-	III	-
HC3-1 <sup>a</sup>	69	M	HCV	+	20	-	II	-
HC3-2 <sup>a</sup>	69	M	HCV	+	20	-	II	-
HC4	61	M	HCV	+	18	-	II	-
HC5	58	M	HCV	+	24	-	II	-
HC6	61	M	HCV	+	35	TACE	II	-
HC7-1 <sup>a</sup>	64	M	HCV	+	25	-	II	-
HC7-2 <sup>a</sup>	64	M	HCV	+	25	-	II	-
HC8	73	F	HCV	-	21	-	II	-
HC9	66	M	HCV	-	40	-	II-III	-
HC10	62	F	HCV	-	20	-	III	-
HC11	71	F	HCV	-	35	-	III	+
HC12	57	M	HCV	+	43	-	I	-
NBNC1	81	F	NBNC	-	70	TACE+PEIT <sup>c</sup>	II	-
NBNC2	73	M	NBNC	+	60	TACE	II	-

**Supplementary Table 1 (continued)**

ID	Hepatic vein invasion	Liver fibrosis <sup>e</sup>	Platelet count (x10 <sup>4</sup> )/mm <sup>2</sup>	Serum AFP level (ng/mL)	Other diseases	Family history of cancer	Multiple liver nodules
HB1	-	0	19.0	106.8	-	mother:HCC	-
HB2	-	4	11.3	15.6	pancreatic tumor	-	-
HB3	+	1	16.6	99.9	esophageal cancer	-	-
HB4	+	2	28.3	382.4	-	-	-
HB5	-	3	15.0	6.9	-	mother: breast cancer	-
HB6	-	4	8.9	2302.0	-	mother & brother: HCC	-
HB7	-	3	21.1	4.3	-	-	-
HB8	-	4	6.1	nd	-	-	-
HB9	-	2	14.4	1630.3	-	-	-
HB10	-	3	11.9	48.6	colon cancer	-	-
HB11	-	4	23.1	1678	-	mother:thyroid cancer, father: cancer of unknown origin	synchronous
HC1	+	4	12.4	1457.4	-	-	-
HC2	-	3	11.3	27.7	-	-	-
HC3-1 <sup>a</sup>	-	4	14.4	7.0	-	-	synchronous
HC3-2 <sup>a</sup>	-	4	14.4	7.0	-	-	synchronous
HC4	-	2	12.9	<5	HCC	father: stomach cancer	metachronous
HC5	-	3	15.3	<5	-	-	-
HC6	+	4	7.9	<5	-	mother: HCC, father: lung cancer	-
HC7-1 <sup>a</sup>	-	4	22.6	3209.0	diabetes	-	synchronous
HC7-2 <sup>a</sup>	-	4	22.6	3209.0	diabetes	-	synchronous
HC8	-	4	15.1	20	-	-	-
HC9	-	3	13.2	nd	-	-	-
HC10	-	4	11.3	26.8	lung cancer	-	metachronous
HC11	-	2	14.2	nd	-	-	-
HC12	-	2	22.7	5.2	diabetes	-	-
NBNC1	-	4	16.8	6.9	-	-	-
NBNC2	-	0	14.3	96.2	diabetes, colon cancer	-	metachronous

a; Multicentric occurrence or intrahepatic metastasis was determined by considering their location, size, and histopathological features.

b; Transcatheter arterial chemoembolization. c; Percutaneous ethanol injection.

d; Edmondson tumor grading. e; Fibrosis in non-cancerous liver tissue is determined according to the New Inuyama Classification.

**Supplementary Table 2: Summary of somatic mutations in 27 HCC genomes**

Sample	Hepatitis	Point mutation					Somatic indel		
		Total	Coding	Synonymous	Nonsynonymous	Splice site	Total	Coding indel	Splice site
HB1	HBV	20,147	133	38	95	5	616	5	0
HB2	HBV	7,083	49	12	37	2	472	6	0
HB3	HBV	19,586	120	30	90	4	790	4	0
HB4	HBV	4,974	26	11	15	1	315	0	0
HB5	HBV	4,886	35	9	26	0	193	2	0
HB6	HBV	8,477	62	14	48	3	470	5	0
HB7	HBV	10,619	84	26	58	3	435	6	0
HB8	HBV	9,583	81	22	59	1	252	1	0
HB9	HBV	16,186	136	38	98	2	363	2	0
HB10	HBV	11,959	103	26	77	2	457	3	0
HB11	HBV	10,359	91	18	73	1	1,714	10	1
HC1	HCV	13,055	69	15	54	4	502	2	0
HC2	HCV	7,120	45	13	32	1	466	5	0
HC3-1	HCV	17,294	115	30	85	0	461	6	0
HC3-2	HCV	19,122	112	26	86	2	550	9	0
HC4	HCV	13,254	68	12	56	3	469	4	1
HC5	HCV	6,357	36	7	29	0	336	2	0
HC6	HCV	7,248	55	13	42	0	259	2	0
HC7-1	HCV	9,006	49	11	38	1	391	6	0
HC7-2	HCV	11,370	63	14	49	0	485	2	1
HC8	HCV	9,426	57	7	50	2	561	2	1
HC9	HCV	9,925	91	25	66	2	658	4	0
HC10	HCV	9,293	66	17	49	1	426	2	0
HC11	HCV	24,147	376	109	267	2	8,950	58	0
HC12	HCV	8,973	74	18	56	3	422	3	0
NBNC1	nBnC	6,464	25	5	20	0	430	2	0
NBNC2	nBnC	22,964	242	62	180	3	496	8	0
Total		318,877	2,463	628	1,835	48	21,939	161	4
Max		24,147	376	109	267	5	8,950	58	1
Min		4,886	25	5	15	0	193	0	0
Average		11,810	91	23	68	2	813	6	0

**Supplementary Table 2 (continued)**

Sample	Hepatitis	Somatic rearrangement				
		Total	Deletion	Inversion	Tandem duplicaiton	Translocation
HB1	HBV	8	3	3	2	0
HB2	HBV	2	1	0	1	0
HB3	HBV	59	27	11	16	5
HB4	HBV	2	1	1	0	0
HB5	HBV	9	1	3	3	2
HB6	HBV	25	14	1	2	8
HB7	HBV	21	12	6	1	2
HB8	HBV	15	10	2	3	0
HB9	HBV	40	18	12	0	10
HB10	HBV	48	20	20	2	6
HB11	HBV	22	7	11	2	2
HC1	HCV	0	0	0	0	0
HC2	HCV	16	2	6	6	2
HC3-1	HCV	9	6	1	2	0
HC3-2	HCV	19	5	7	3	4
HC4	HCV	10	5	1	0	4
HC5	HCV	15	4	1	7	3
HC6	HCV	7	3	0	1	3
HC7-1	HCV	2	1	0	1	0
HC7-2	HCV	22	11	4	5	2
HC8	HCV	41	19	12	0	10
HC9	HCV	23	11	10	1	1
HC10	HCV	15	7	6	2	0
HC11	HCV	40	15	23	1	1
HC12	HCV	50	43	7	0	0
NBNC1	nBnC	6	0	3	2	1
NBNC2	nBnC	35	27	3	2	3
Total		561	273	154	65	69
Max		59	43	23	16	10
Min		0	0	0	0	0
Average		21	10	6	2	3







<i>DCC</i>	18	49,867,158	51,057,023	4,456	0	2	0	0	2	0.1734	0.2088
<i>TRANK1</i>	3	36,869,766	36,986,300	8,866	0	2	0	0	2	0.1734	0.2088
<i>BPTF</i>	17	65,821,841	65,978,404	8,828	1	1	0	0	2	0.1734	0.2088
<i>MAP1B</i>	5	71,403,359	71,501,066	7,431	0	2	0	0	2	0.1948	0.2320
<i>ITPR1</i>	3	4,558,176	4,887,909	8,464	0	2	0	0	2	0.1948	0.2320
<i>NCKAP5</i>	2	133,430,862	134,275,097	5,798	0	2	0	0	2	0.2166	0.2564
<i>AHNAK</i>	11	62,284,216	62,303,570	17,681	0	3	0	0	3	0.2192	0.2577
<i>PREX2</i>	8	68,864,630	69,143,613	4,977	0	2	0	0	2	0.2202	0.2577
<i>ODZ3</i>	4	183,245,174	183,721,504	8,204	0	2	0	0	2	0.2239	0.2605
<i>PRUNE2</i>	9	79,229,486	79,520,879	9,339	1	1	0	0	2	0.2495	0.2888
<i>C12orf51</i>	12	112,600,191	112,757,291	13,033	1	1	0	0	2	0.2532	0.2914
<i>DCHS1</i>	11	6,643,010	6,662,844	9,973	0	2	0	0	2	0.2569	0.2940
<i>HERC1</i>	15	63,901,280	64,067,822	14,890	0	1	1	0	2	0.2606	0.2949
<i>PCDH7</i>	4	30,723,045	31,144,471	3,752	0	2	0	0	2	0.2606	0.2949
<i>RELN</i>	7	103,113,259	103,629,803	10,639	0	2	0	0	2	0.2716	0.3041
<i>MLL</i>	11	118,307,228	118,392,887	12,050	1	1	0	0	2	0.2716	0.3041
<i>TRRAP</i>	7	98,478,774	98,609,978	11,769	0	2	0	0	2	0.2936	0.3270
<i>RNF213</i>	17	78,237,481	78,367,298	16,039	0	2	0	0	2	0.2973	0.3293
<i>FAT4</i>	4	126,237,567	126,412,923	15,010	1	2	0	0	3	0.2990	0.3294
<i>ABCA13</i>	7	48,211,081	48,685,108	15,421	0	3	0	0	3	0.3071	0.3365
<i>LAMA2</i>	6	129,204,391	129,837,492	9,625	0	2	0	0	2	0.3338	0.3638
<i>SPTA1</i>	1	158,581,054	158,656,307	7,464	0	1	1	0	2	0.3374	0.3658
<i>DNAH6</i>	2	84,744,951	85,046,532	12,777	0	2	0	0	2	0.3589	0.3871
<i>LRP2</i>	2	169,985,173	170,218,909	14,280	1	1	0	0	2	0.4012	0.4305
<i>ANK3</i>	10	61,802,449	62,149,296	13,302	0	2	0	0	2	0.4422	0.4707
<i>LRP1B</i>	2	140,990,755	142,888,298	14,160	0	3	0	0	3	0.4433	0.4707
<i>CSMD1</i>	8	2,796,107	4,851,938	10,971	0	2	0	0	2	0.4488	0.4742
<i>EYS</i>	6	64,430,492	66,205,303	9,591	1	1	0	0	2	0.4718	0.4959
<i>FBN2</i>	5	127,595,147	127,873,296	8,995	0	2	0	0	2	0.4783	0.5001
<i>MLL3</i>	7	151,833,917	152,132,871	14,968	0	2	0	0	2	0.5160	0.5369
<i>USH2A</i>	1	215,799,123	216,595,678	15,889	0	2	0	0	2	0.5460	0.5652
<i>DNAH7</i>	2	196,602,645	196,933,435	12,331	0	2	0	0	2	0.5747	0.5920
<i>DYNC2H1</i>	11	102,980,304	103,349,981	13,301	0	2	0	0	2	0.5858	0.6004
<i>CSMD3</i>	8	113,237,000	114,449,083	11,404	0	2	0	1	3	0.5967	0.6085
<i>MUC16</i>	19	8,959,608	9,091,814	43,856	1	2	0	0	3	0.6641	0.6739
<i>PCLO</i>	7	82,387,891	82,791,908	15,525	0	3	0	0	3	0.6773	0.6839
<i>GPR98</i>	5	89,854,713	90,459,717	19,277	0	2	0	0	2	0.7307	0.7342
<i>TTN</i>	2	179,391,739	179,669,369	101,512	0	4	1	0	5	0.9935	0.9935



**Supplementary Table 9: Gene set enrichment (GSE) analysis for genes with nonsense, coding indel and splice-site mutations**

Category	Term	Count	%	<i>P</i> -value	List Total	Pop Hits	Pop Total	Fold Enrichment	<i>q</i> -value
SP_PIR_KEYWORDS	phosphoprotein	162	55.1	6.92E-10	292	7263	19235	1.469292662	0.00000023
INTERPRO	IPR013032:EGF-like region, conserved site	19	6.5	8.02E-07	260	293	16659	4.154909425	0.00049
INTERPRO	IPR000742:EGF-like, type 3	15	5.1	2.14E-06	260	194	16659	4.954103886	0.0006
INTERPRO	IPR006210:EGF-like	15	5.1	3.25E-06	260	201	16659	4.781572905	0.0007
SP_PIR_KEYWORDS	egf-like domain	15	5.1	1.15E-05	292	230	19235	4.296083979	0.0019
UP_SEQ_FEATURE	domain:EGF-like 1	12	4.1	2.25E-06	292	120	19113	6.545547945	0.0039
SMART	SM00181:EGF	15	5.1	3.00E-05	175	201	9079	3.871641791	0.0043
SP_PIR_KEYWORDS	polymorphism	208	70.7	4.75E-05	292	11550	19235	1.18628951	0.0052
SP_PIR_KEYWORDS	calcium	28	9.5	9.42E-05	292	803	19235	2.2969515	0.0078
UP_SEQ_FEATURE	sequence variant	218	74.1	1.06E-05	292	11992	19113	1.189901144	0.0092
SP_PIR_KEYWORDS	<b>bromodomain</b>	6	2.0	2.89E-04	292	39	19235	10.13435195	0.019
SP_PIR_KEYWORDS	<b>chromatin regulator</b>	12	4.1	4.17E-04	292	213	19235	3.711171136	0.023
INTERPRO	IPR006209:EGF	10	3.4	1.66E-04	260	127	16659	5.045124167	0.025
SP_PIR_KEYWORDS	disease mutation	42	14.3	5.45E-04	292	1591	19235	1.738955426	0.026
INTERPRO	<b>IPR001487:Bromodomain</b>	6	2.0	3.69E-04	260	40	16659	9.610961538	0.044
SP_PIR_KEYWORDS	tumor suppressor	9	3.1	0.001157126	292	137	19235	4.327442256	0.047

q-value was obtained by Benjamini and Hochberg's FDR method

Supplementary Table 10: Mutation list of the chromatin regulator genes

Gene	Sample	Chr	Genomic position	Genotype	RNA_id	Exon	aa position	Type	Reference aa	Mutant aa
<i>ARID1B</i>	HB9	6	157,528,540	CG	NM_020732	20	2,089	nonsynonymous	L	V
<i>ARID1B</i>	HC5	6	157,454,194	AG	NM_020732	8	802	nonsynonymous	M	V
<i>ARID1B</i>	validation study	6	157,099,508	CG	NM_020732	1	149	nonsynonymous	Q	E
<i>ARID1B</i>	validation study	6	157,100,285	CT	NM_020732	1	408	nonsynonymous	Q	Stop
<i>ARID1B</i>	validation study	6	157,099,981	-GGC	NM_020732	1	306	coding indel	-	-
<i>ARID1B</i>	validation study	6	157,100,046	-AGC	NM_020732	1	328	coding indel	-	-
<i>ARID1B</i>	validation study	6	157,100,137	-G	NM_020732	1	358	coding indel	-	-
<i>ARID1B</i>	validation study	6	157,100,222	+C	NM_020732	1	387	coding indel	-	-
<i>ARID1B</i>	validation study	6	157,100,309	-C	NM_020732	1	416	coding indel	-	-
<i>ARID1B</i>	validation study	6	157,505,455	-T	NM_020732	13	1,146	coding indel	-	-
<i>ARID2</i>	HC9	12	46,244,449	+T	NM_152641	15	848	coding indel	-	-
<i>ARID2</i>	HB3	12	46,243,834	AC	NM_152641	15	643	nonsynonymous	S	Stop
<i>ARID2</i>	validation study	12	46,123,699	AT	NM_152641	1	27	nonsynonymous	H	L
<i>ARID2</i>	validation study	12	46,231,200	CG	NM_152641	9	374	nonsynonymous	G	R
<i>ARID2</i>	validation study	12	46,244,897	GT	NM_152641	15	997	nonsynonymous	Q	H
<i>ARID2</i>	validation study	12	46,285,641	GT	NM_152641	17	1,667	nonsynonymous	C	W
<i>ARID2</i>	validation study	12	46,123,699	-C	NM_152641	1	27	coding indel	-	-
<i>ARID2</i>	validation study	12	46,285,661	-T	NM_152641	17	1,674	coding indel	-	-
<i>ARID2</i>	validation study	12	46,298,726	-A	NM_152641	21	1,791	coding indel	-	-
<i>MLL</i>	HC11	11	118,344,186	-C	NM_005933	3	771	coding indel	-	-
<i>MLL</i>	HC7-2	11	118,342,410	CG	NM_005933	3	179	nonsynonymous	P	R
<i>MLL</i>	validation study	11	118,307,445	AG	NM_005933	1	73	nonsynonymous	G	E
<i>MLL</i>	validation study	11	118,352,504	AG	NM_005933	7	1,237	nonsynonymous	V	M
<i>MLL3</i>	HC3-1	7	151,853,128	CT	NM_170606	46	3,943	nonsynonymous	G	R
<i>MLL3</i>	HC4	7	151,859,936	CT	NM_170606	43	3,576	nonsynonymous	I	V
<i>MLL3</i>	validation study	7	151,921,595	CT	NM_170606	19	1,028	nonsynonymous	C	Y
<i>MLL3</i>	validation study	7	151,932,906	CG	NM_170606	16	922	nonsynonymous	P	R
<i>MLL3</i>	validation study	7	151,932,946	CT	NM_170606	16	909	nonsynonymous	R	G
<i>MLL3</i>	validation study	7	151,962,189	AGT	NM_170606	8	373	nonsynonymous	I	T, K
<i>MLL3</i>	validation study	7	151,879,611	-T	NM_170606	36	1,778	coding indel	-	-
<i>BRD8</i>	HB1	5	137,496,591	AC	NM_139199	18	806	nonsynonymous	E	Stop
<i>BRD8</i>	HB6	5	137,496,741	CT	NM_139199	18	756	nonsynonymous	T	A
<i>BRD8</i>	validation study	5	137,497,513	CT	NM_139199	17	741	nonsynonymous	A	T
<i>BPTF</i>	HB10	17	65,907,062	AG	NM_004459	13	1,147	nonsynonymous	D	G
<i>BPTF</i>	HC3-1	17	65,871,775	-C	NM_004459	5	656	coding indel	-	-
<i>BPTF</i>	validation study	17	65,907,062	AG	NM_004459	13	1,147	nonsynonymous	D	G
<i>BPTF</i>	validation study	17	65,919,058	CT	NM_004459	16	2,013	nonsynonymous	V	A
<i>HIST1H4B</i>	HB11	6	26,027,251	+A	NM_003544	1	77	coding indel	-	-
<i>HIST1H4B</i>	HC4	6	26,027,170	AT	NM_003544	1	104	nonsynonymous	Stop	L
<i>HIST1H4B</i>	validation study	6	26,027,170	AT	NM_003544	1	104	nonsynonymous	Stop	K
<i>BRE</i>	HC11	2	28,248,273	-A	NM_004899	4	161	coding indel	-	-
<i>BRE</i>	HC4	2	28,464,244	CG	NM_004899	8	279	nonsynonymous	L	V

**Supplementary Table 11: Association between mutations in chromatin regulators and clinical factors**

data	Tumor size	Edmondson grade	Portal vein invasion
Point mutation and indel in the chromatin regulator genes	0.67	0.75	1.00
Point mutation, indel and CNA in the chromatin regulator genes	0.78	0.62	1.00
Point mutation and indel in the <i>ARID</i> genes ( <i>ARID1A</i> , <i>ARID2</i> , and <i>ARID1B</i> )	0.56	0.94	0.36
Point mutation, indel and CNA in the <i>ARID</i> genes ( <i>ARID1A</i> , <i>ARID2</i> , and <i>ARID1B</i> )	0.80	0.83	0.38
Point mutation and indel in <i>ARID1A</i> gene	0.44	0.20	0.59
Point mutation, indel and CNA in <i>ARID1A</i> gene	0.39	0.67	0.36

**Supplementary Table 11 (continued)**

data	Hepatic vein invasion	Liver fibrosis
Point mutation and indel in the chromatin regulator genes	1.00	0.061
Point mutation, indel and CNA in the chromatin regulator genes	0.61	0.026
Point mutation and indel in the <i>ARID</i> genes ( <i>ARID1A</i> , <i>ARID2</i> , and <i>ARID1B</i> )	0.27	0.16
Point mutation, indel and CNA in the <i>ARID</i> genes ( <i>ARID1A</i> , <i>ARID2</i> , and <i>ARID1B</i> )	0.27	0.0086
Point mutation and indel in <i>ARID1A</i> gene	0.39	0.63
Point mutation, indel and CNA in <i>ARID1A</i> gene	0.042	0.16

*P-values* for Tumor size, Edmondson grade, and Liver fibrosis were calculated by the Pearson's correlation test.

*P-values* for Portal vein invasion and Hepatic vein invasion were obtained by the Fisher's exact test.

**Supplementary Table 12: Recurrent somatic mutations of non-coding regions (n ≥ 3)**

Symbol	Type	Chr	Start	Last	SNV	indel	Length (bp)	q-value
<i>ANKRD30BL</i>	ncRNA	2	132,905,164	133,015,542	8	6	110,378	0
<i>AF146191.4</i>	lincRNA	4	190,742,266	190,838,477	12	0	96,211	8.10463E-15
<i>U2</i>	snRNA	10	103,124,602	103,124,792	4	0	190	1.59648E-07
<i>RP5-875O13.1</i>	lincRNA	1	16,860,381	16,862,144	3	3	1,763	1.69297E-07
<i>AD000090.1</i>	miRNA	19	36,066,627	36,066,721	3	0	94	2.41516E-06
<i>AC020926.1</i>	lincRNA	16	20,685,853	20,693,496	6	0	7,643	7.84436E-06
<i>HMHA1</i>	UTR5	19	1,067,174	1,086,627	3	0	19,453	1.19856E-05
<i>RNU2-2</i>	snRNA	11	62,609,091	62,609,281	2	1	190	1.19856E-05
<i>LOC100130581</i>	PRO	17	41,447,213	41,466,266	2	2	19,053	3.12021E-05
<i>CTBP2</i>	UTR3	10	126,676,418	126,849,103	2	2	172,685	5.73097E-05
<i>NEAT1</i>	lincRNA	11	65,190,269	65,213,011	9	0	22,742	9.66691E-05
<i>ANKRD36BP2</i>	lincRNA	2	89,065,385	89,106,126	7	0	40,741	9.91213E-05
<i>PRR15</i>	UTR3	7	29,603,427	29,606,911	3	0	3,484	0.000193272
<i>DKFZp434L192</i>	PRO	7	56,563,916	56,564,977	4	0	1,061	0.000357403
<b>PHF1</b>	PRO	6	33,378,773	33,384,230	3	0	5,457	0.000391151
<i>HLA-DRB5</i>	PRO	6	32,485,154	32,498,006	4	0	12,852	0.000664692
<b>MEDI</b>	UTR3	17	37,560,538	37,607,527	4	0	46,989	0.000664692
<i>AC079305.1</i>	lincRNA	2	178,107,103	178,111,509	2	2	4,406	0.000724295
<i>PRSS3</i>	PRO	9	33,795,559	33,799,229	3	0	3,670	0.000745744
<i>CTD-2384A14.1</i>	lincRNA	14	29,400,776	29,437,598	3	0	36,822	0.000745744
<i>LOC100130581</i>	ncRNA	17	41,453,296	41,466,266	3	1	12,970	0.000813248
<i>CTB-113P19.1</i>	lincRNA	5	151,056,506	151,067,471	3	1	10,965	0.001123863
<i>NIP7</i>	UTR3	16	69,373,415	69,377,013	3	0	3,598	0.001307338
<i>SMG1</i>	UTR3	16	18,816,175	18,937,726	3	0	121,551	0.001768832
<i>RIBC1</i>	UTR5	X	53,449,839	53,456,776	3	0	6,937	0.001768832
<i>IL22RA2</i>	UTR3	6	137,464,957	137,494,785	3	0	29,828	0.002785282
<i>URB1</i>	UTR3	21	33,683,330	33,765,312	3	1	81,982	0.003785807
<i>SRSF3</i>	UTR3	6	36,562,090	36,572,244	1	2	10,154	0.003791246
<i>LOC150622</i>	PRO	2	6,072,819	6,120,350	3	0	47,531	0.003791246
<i>TRY6</i>	PRO	7	142,478,757	142,482,399	3	0	3,642	0.006317702
<i>LOC121952</i>	ncRNA	13	103,532,449	103,548,383	3	0	15,934	0.006317702
<i>HCG22</i>	lincRNA	6	31,021,984	31,027,653	2	2	5,669	0.006388021
<i>PIK3R3</i>	UTR3	1	46,505,812	46,598,380	3	0	92,568	0.007245836
<i>KIAA1430</i>	UTR3	4	186,080,819	186,125,182	3	0	44,363	0.007507514
<i>PTAR1</i>	UTR3	9	72,324,438	72,374,876	5	0	50,438	0.007507514
<i>LOC220980</i>	lincRNA	10	45,306,472	45,455,137	3	0	148,665	0.007994232
<i>SCD</i>	UTR3	10	102,106,772	102,124,588	3	0	17,816	0.009002124
<i>SMAD6</i>	ncRNA	15	66,994,674	67,074,337	3	0	79,663	0.009002124
<i>SRSF3</i>	ncRNA	6	36,562,090	36,572,244	1	2	10,154	0.009522649
<i>PRPF40A</i>	UTR3	2	153,508,107	153,573,975	4	0	65,868	0.009808589
<i>MAG13</i>	UTR3	1	113,933,475	114,228,545	2	1	295,070	0.009808589
<i>RP11-622O11.2</i>	lincRNA	8	126,953,373	126,963,433	3	0	10,060	0.011986707
<i>FLJ39653</i>	lincRNA	4	16,228,286	16,259,810	3	0	31,524	0.012564096
<i>NFIA</i>	UTR3	1	61,547,980	61,928,460	4	0	380,480	0.017052264
<i>EDNRB</i>	UTR3	13	78,469,616	78,492,966	3	0	23,350	0.017648691
<i>AP000459.4</i>	lincRNA	21	24,733,426	24,757,182	3	0	23,756	0.018307238
<i>HCG2P7</i>	ncRNA	6	29,866,808	29,870,431	3	0	3,623	0.01897249
<i>AC011477.2</i>	lincRNA	19	19,945,686	20,008,579	3	0	62,893	0.019644199
<i>CSGALNACT1</i>	ncRNA	8	19,261,672	19,540,261	3	0	278,589	0.021434704
<i>ATXN7L3B</i>	UTR3	12	74,931,551	74,935,232	3	0	3,681	0.024476649
<i>PTGER3</i>	UTR3	1	71,471,537	71,513,491	4	0	41,954	0.028065538
<i>AFF4</i>	UTR3	5	132,211,071	132,299,354	3	0	88,283	0.029118099
<i>RORA</i>	UTR3	15	60,780,483	60,884,707	3	1	104,224	0.029118099
<i>MYEF2</i>	UTR3	15	48,431,629	48,470,558	3	0	38,929	0.029118099
<i>FAM126B</i>	UTR3	2	201,838,441	201,936,392	3	0	97,951	0.029384429
<i>RLIM</i>	UTR3	X	73,802,811	73,834,461	3	0	31,650	0.032953863
<i>ADCY2</i>	UTR3	5	7,396,343	7,830,194	3	0	433,851	0.032953863
<i>KLF6</i>	ncRNA	10	3,818,188	3,827,473	3	0	9,285	0.032953863
<i>ZC3H6</i>	UTR3	2	113,033,178	113,097,640	2	1	64,462	0.03371726

<i>LOC400940</i>	ncRNA	2	6,122,110	6,128,364	3	0	6,254	0.045347117
<i>AP000476.1</i>	lincRNA	21	25,801,054	25,862,621	3	0	61,567	0.045347117
<i>PAM</i>	ncRNA	5	102,201,527	102,366,808	3	0	165,281	0.046140706

---

UTR3; 3'UTR, UTR5; 5'UTR, PRO; promoter

**Supplementary Table 13: HBV integration sites (IS) identified in HBV-related HCC genomes**

Sample	Candidate name	Chr1	Pos1	Chr2	Pos2	Number of support read-pair	Affected gene
HB1	IS1	HBV	1,735	5	88,432,442	14	
	IS2	HBV	687	5	88,754,909	18	
HB2	IS1	HBV	1,783	5	1,275,381	3	<i>TERT</i>
	IS2	HBV	1,835	5	1,275,390	8	<i>TERT</i>
	IS3	HBV	1,879	5	173,125,603	9	
HB3	IS1	HBV	1,470	6	68,664,024	11	
	IS2	HBV	1,771	8	39,200,450	14	<i>ADAM5P</i>
	IS3	HBV	1,808	4	79,048,012	19	<i>FRAS1</i>
	IS4	HBV	1,818	4	79,048,008	15	<i>FRAS1</i>
HB4	-	-	-	-	-	-	
HB5	IS1	HBV	2,141	3	24,499,331	11	<i>THRB</i>
HB6	-	-	-	-	-	-	
HB7	IS1	HBV	1,401	7	52,802,419	11	
	IS2	HBV	2,622	7	52,802,387	5	
	IS3	HBV	1,361	3	121,439,129	3	<i>HNFI1A</i>
HB8	-	-	-	-	-	-	
HB9	IS1	HBV	1,854	5	1,295,172	26	<i>TERT</i>
	IS2	HBV	1,726	5	1,295,200	13	<i>TERT</i>
	IS3	HBV	1,895	7	70,267,023	4	
HB10	IS1	HBV	1,828	5	1,293,404	37	<i>TERT</i>
	IS2	HBV	2,236	5	1,293,416	25	<i>TERT</i>
HB11	IS1	HBV	1,411	5	1,310,429	74	<i>TERT-CLPTMIL<sup>a</sup></i>
	IS2	HBV	1,675	17	15,412,626	34	<i>FAM18B2</i>
	IS3	HBV	2,716	20	30,502,092	18	<i>TLL9</i>
	IS4	HBV	1,814	9	24,396,870	14	
	IS5	HBV	1,846	4	151,954,261	4	

Integration site was not identified in HB4, HB6, and HB8.

a; Intergenic region between *TERT* and *CLPTMIL*



**Supplementary Table 14: Association between etiological factors and mutation pattern**

Data set	Distance between MC tumors <sup>j</sup>	HBV vs. HCV vs. NBNC	Liver fibrosis	Alcohol drinking	Edmondson grade <sup>k</sup>	Age
Point mutation (6 category) <sup>a</sup>	0.0073	0.3812	0.9328	0.0051	0.5866	0.7804
Point mutation <sup>b</sup>	0.0005	0.0909	0.5716	0.0283	0.5324	0.2423
Point mutation+indel <sup>c</sup>	0.0008	0.0685	0.5252	0.0408	0.4137	0.2375
Point mutation+indel+Rearrangement <sup>d</sup>	0.0001	0.0682	0.3835	0.0128	0.4149	0.2716
Point mutation+indel length <sup>e</sup>	0.0003	0.1147	0.4427	0.0049	0.6310	0.4401
Point mutation+indel length+Rearrangement <sup>f</sup>	0.0002	0.1891	0.5016	0.0055	0.8440	0.6076
Point mutation+indel length+ns_rate <sup>g</sup>	0.0006	0.0216	0.2497	0.0256	0.5051	0.2610
Point mutation+indel+Rearrangement+ns_rate <sup>h</sup>	0.0003	0.0146	0.2151	0.0892	0.3218	0.2058
Point mutation+total_num <sup>i</sup>	0.0006	0.0877	0.1529	0.0398	0.5388	0.2494

**Supplementary Table 14 (continued)**

Data set	Tumor size	Platelet count (x10 <sup>4</sup> )/mm <sup>2</sup>	Portal vein invasion	Hepatic vein invasion	Serum AFP level <sup>l</sup>	Multiple liver nodules
Point mutation (6 category) <sup>a</sup>	0.0411	0.3601	0.5337	0.8128	0.0681	0.0118
Point mutation <sup>b</sup>	0.0784	0.5838	0.2971	0.8416	0.1420	0.0161
Point mutation+indel <sup>c</sup>	0.1066	0.7748	0.4683	0.8277	0.2155	0.0168
Point mutation+indel+Rearrangement <sup>d</sup>	0.0939	0.7317	0.3926	0.8394	0.2707	0.0135
Point mutation+indel length <sup>e</sup>	0.1029	0.6230	0.3990	0.7856	0.2601	0.0173
Point mutation+indel length+Rearrangement <sup>f</sup>	0.0810	0.6289	0.4068	0.7854	0.2805	0.0204
Point mutation+indel length+ns_rate <sup>g</sup>	0.2538	0.4773	0.4026	0.6012	0.1704	0.0186
Point mutation+indel+Rearrangement+ns_rate <sup>h</sup>	0.1576	0.8231	0.4465	0.7129	0.2054	0.0108
Point mutation+total_num <sup>i</sup>	0.0475	0.7064	0.1603	0.8449	0.4671	0.0119

a; number of substitutions without considering CpG. b; number of substitutions with CpG. c; number of substitutions and indels. d; number of substitutions, indels and rearrangements.

e; number of substitutions and length distribution of indels (number of 1bp, 2bp, 3bp, 5bp and >=5bp indels). f; number of substitutions, length distribution of indels and rearrangements.

g; number of substitution, length distribution of indels and the ratio of number of nonsynonymous to that of coding mutations.

h; number of substitution, indel, rearrangement and the ratio of number of nonsynonymous to that of coding mutations.

i; number of substitutions with CpG and total number of point mutations. j; Distance between MC tumors are tested by permutation test.

k; Edmondson tumor grading. l; Samples were classified into two groups with AFP level (>10ng/mL or not).