

Resolving complex structures at oncovirus integration loci with conjugate graph

Wenlong Jia^{1,+}, Chang Xu^{1,+}, and Shuai Cheng Li^{1,*}

¹Department of Computer Science, City University of Hong Kong, Kowloon Tong, Hong Kong

⁺These authors contributed equally to this work

^{*}Correspondence should be addressed to S.C.L. (shuaicli@cityu.edu.hk).

ABSTRACT

Oncovirus integrations cause complex structural variations (SVs) on host genomes. We propose a conjugate graph model to reconstruct the rearranged local genomic map (LGM) at integrated loci. Simulation tests prove the reliability and credibility of the algorithm. Applications of the algorithm to whole-genome sequencing data of Human papillomavirus (HPV) and hepatitis B virus (HBV)-infected cancer samples gained biological insights on oncovirus integrations. We observed five affection patterns of oncovirus integrations from the HPV and HBV-integrated cancer samples, including the exon loss, promoter gain, hyper-amplification of tumor gene, the viral cis-regulation inserted at the single intron and at the intergenic region. We found that the focal duplicates and host SVs are frequent in the HPV-integrated LGMs, while the focal deletions and complex virus SVs are prevalent in HBV-integrated LGMs. Furthermore, with the results yields from our method, we found the enhanced microhomology-mediated end joining (MMEJ) might lead to both HPV and HBV integrations, and conjectured that the HPV integrations might mainly occur during the DNA replication process.

1 INTRODUCTION

2 Oncovirus integrations (VITs) could induce cancer genome instability¹⁻⁵, emerging as copy number
3 variations (CNVs) and complex structural variations (SVs). Adjacent genes to the integrated sites are
4 frequently dysregulated, such as *TERT* in hepatocellular carcinoma^{5,6} and *MYC* in cervical carcinoma^{7,8}.
5 However, the explicit structure of the rearranged local genomic map (LGM) remains elusive. Some studies
6 have reported the local maps at the integration site of human papillomavirus (HPV). Adey *et al.* resolved
7 the haplotype of the HeLa cell line and HPV18-integrated genomic locus by consolidating results from a
8 range of sequencing technologies, including the whole-genome shotgun sequencing, the fosmid mate-pair
9 library with large insert size, and PacBio single-molecule sequencing data¹. Akagi *et al.* constructed
10 manually preliminary LGMs flanking integration sites in HPV-positive samples², including head and neck
11 squamous cell carcinomas (HNSCC) and cervical cancers. However, oncovirus integrations demand the
12 resolved local maps systematically, which would lead to more insights. Accurate LGMs analysis lays
13 the foundations to understand how oncovirus integrations emerge and compromise the host genome's
14 functioning. A computational model and algorithm are necessary to obtain the local genomic maps.
15 Although a growing number of studies examine the oncovirus integration effects on the expression
16 regulations and genome stability in the flanking host genomic regions⁶⁻¹³, our understanding of how
17 inserted viral DNA influences the local genome remains limited. The LGMs can provide the haplotyping
18 information to promote the investigation, such as the pairwise integrations, host and viral rearrangements,
19 inserted viral regulatory elements, and the copy number changes. Further, the LGM linear structure can
20 facilitate the viral integration mechanism research with junction breakpoints phased together, as well as

21 the alternative splicing and virus-host fusions. Several factors hinder the LGMs accurate construction,
22 including the sequencing fluctuations, short genomic segments, repeated elements, and insufficient junction
23 reads. Here, we introduce a graph-based algorithm to overcome these difficulties and construct LGMs at
24 the oncovirus integration sites from cancer whole genome sequencing (WGS) data. Using multiple real and
25 *in silico* datasets, we show that the algorithm is able to construct credible LGMs of oncovirus integrations
26 and complex rearrangements on cancer genome. The LGM results further facilitate our investigation on
27 the regulatory effects and formation mechanism accompanied by virus integrations.

28 RESULTS

29 Conjugate graph to construct local genomic map

30 Based on the breakpoints of VITs and SVs, both the host local genome and viral genome are divided into
31 segments (Figure 1a-b), as vertices in the graph model (Figure 1d). Copy numbers (CNs) of segments
32 are as the weights of vertices. Junctions of SVs and VITs are modeled as directed edges connecting
33 two vertices, and their split-reads counts are converted to CN as the weights of the edges (Figure 1c-d).
34 Integer linear programming is applied to correct weights of vertices and edges to reduce segment depth
35 fluctuations, as well as compliment insufficient split-reads (see Methods). As genome DNA has double
36 strands, the vertices or edges are conjugated, i.e, a vertex or an edge and their reverse complementary
37 counterparts and share the respective weights. We then use an algorithm modified from the *Eulerian*
38 *Circuits* traversal rule to find available circular paths (called unit-cycles, Figure 1e), and then merge them
39 to construct the local genomic map consisting of host and virus segments (Figure 1f).

40 In our algorithm, the first upstream and last downstream host segments are called the *source* and *sink*
41 vertices of the conjugate graph (Figure 1d-e) respectively. The algorithm first finds the main path, which
42 starts from the *source* vertex and extends along the edges with non-zero weight until it reaches the *sink*
43 vertex (Figure 1e). If a vertex or an edge is included in the path, its weight on the graph decreases by one.
44 Then, starting at a vertex with non-zero weight left, our algorithm forms a unit-cycle by randomly walking
45 along edges with non-zero weight until the vertex itself has reached again. More unit-cycles can be found
46 iteratively until the all the weights are drained. Finally, all unit-cycles are merged with the main path to
47 get putative LGMs. Due to the random walking and merging of the unit-cycles, there could be multiple
48 LGM candidates generated from the conjugate graph, and our algorithm will report all possible LGMs.

49 Evaluation on *in silico* data

50 To evaluate the reliability of the conjugate graph algorithm, we simulated a set of oncovirus integration
51 LGMs at genomic loci of well-known hotspot genes, including HBV integrations (*TERT*, *CCNE1*, and
52 *KMT2B*)¹¹ and HPV integrations (*LRP1B*, *FHIT*, and *MYC*)⁷. Combining the different SV types, copy
53 numbers, and viral segments, we generated eight LGM simulated modes in each gene locus, such as
54 substitutive integrations, inversed insertions, and complex tandem duplications (Supplementary Figure
55 1, see Methods). In total, we prepared simulated sequencing data of 48 LGM cases for evaluation
56 (Supplementary Table 1). The conjugate graph algorithm successfully recovered all simulated LGMs
57 at HBV and HPV integrated sites. The integer linear programming obtained the copy numbers of all
58 segments and junctions (SVs and VITs) with zero offset from the expected values. Moreover, for complex
59 structures with high CN segments, the algorithm reported other possible candidate genomic maps that
60 also satisfy the observed conjugate graph. The performance on *in silico* data supports the reliability of the
61 conjugate graph to resolve diverse local genomic maps.

62 On WGS data of HPV-integrated cervical cancers

63 We first applied the developed method to reconstruct the LGMs at HPV integration loci in four cervical
64 carcinoma samples (T4931, T6050, HeLa, and SiHa)⁷. All samples have arm-level copy number changes
65 of the HPV-integrated chromosomes (Supplementary Figure 2). The conjugate graph algorithm constructed
66 LGMs of all validated HPV integrations with most CN offsets less than one (73.3%, Supplementary Table
67 2 and 3).

68 First, the LGM on unbalanced polyploidy was successfully constructed in sample T4931. The HPV16
69 integrations locate at gene *GLI2* locus on chr2q (heterozygous pentaploid, minor CN = 2, Supplementary
70 Figure 2a). As the integration might locate on the paternal or maternal chromosomes, our method algorithm
71 provides two LGM candidates (Figure 2 and Supplementary Figure 3). One LGM has less CN offsets and
72 it includes two unit-cycles containing three virus integrations (VITs), virus circularization (V-CIRC), one
73 viral inversion, and one host tandem duplication. The HPV segments insert to the intron region and carry
74 viral long control region (LCR).

75 Second, we solved the LGM with asymmetric copy of the *source* and *sink* vertices. The chr8q of HeLa
76 cell line shows different copy numbers at the HPV18 integrated sites: the upstream and downstream are
77 tetraploid and hexaploid respectively, both of which are heterozygous (minor CN = 1, Supplementary
78 Figure 2b). Previous study has reported that the HPV18 integrates on the major allele of the tetraploid¹.
79 The constructed LGM contains three unit-cycles consisting of four VITs and the V-CIRC (Supplementary
80 Figure 4). The unit-cycles are evenly distributed in the three copies of the major allele. The LGM structure
81 is similar to the reported¹, but with different copy numbers and break-points of the unit-cycles.

82 Third, two LGMs appeared as diploid with loss of heterozygosity (LOH). Both of the T6050 sample
83 and SiHa cell line have HPV16 integrations at chr13q (LOH, Supplementary Figure 2c-d). The LGM
84 of T6050 contains one unit-cycle consisting of two VITs, the V-CIRC, and one host tandem duplication
85 at the gene *KLF12* locus (Supplementary Figure 5). Interestingly, this unit-cycle has different copy
86 times (two and three) in the remained diploid alleles, which implies the possible existence of circular
87 extrachromosomal DNA (ecDNA). The SiHa cell line has two unit-cycles in the LGM at the HPV16
88 integration sites (downstream of the gene *KLF12*, Supplementary Figure 6). The unit-cycles contain two
89 VITs, the V-CIRC, one viral deletion, and one host deletion. The structure of the unit-cycles is the same as
90 in the previous reports².

91 On WGS data of HPV-integrated HNSCC cancers

92 We then applied the conjugate graph algorithm on two HPV16-positive HNSCC cell lines (UM-SCC-47
93 and UPCI-SCC090)². The resolved HPV16-integrated LGMs are more complex than the four cervical
94 cancers. Different from the published preliminary results², our method reported all detected genomic
95 segments and breakpoints.

96 The first LGM contains tandem duplication of partial gene body of tumor suppressor, and might lead
97 to the truncation of the coding-frame. The UM-SCC-47 cell line has HPV16 integrations at the tumor
98 suppressor *TP63* locus on chr3q^{14,15}. The broad region (3q26.31-q29) shows heterozygous triploid (minor
99 CN = 1), and the focal region (in q28) at the HPV16 integrations has prominent high copy number gains
100 (Supplementary Figure 7a). In total, six HPV16 integrations and four endogenous rearrangements are
101 considered to construct two LGMs (Supplementary Table 4) corresponding to the minor (Supplementary
102 Figure 8) and major alleles (Supplementary Figure 9). Our algorithm successfully obtained five unit-cycles.
103 Interestingly, all segments of the NO.1 unit-cycle are from the host genome and form a tandem duplication
104 structural variation with an inner deletion (the 'H3' segment). The NO.3 and NO.4 unit-cycles both have
105 high copy numbers (5 and 49 respectively) and mainly contribute to the focal gain on the host genome.
106 The minor allele and major allele LGMs share the same unit-cycles and differ only in copy numbers.

107 Another LGM has focal hyper-amplification inducing the copy number gain of oncogene. The UPCI-
108 SCC090 cell line has HPV16 integrations in chr6p21.2. The chr6 is heterozygous triploid (minor CN = 1),
109 and the HPV-integrated focal region shows a higher copy number ([Supplementary Figure 7b](#)). The focal
110 host region contains seven HPV16 integrations and six endogenous rearrangements. These breakpoints
111 divide the host and virus genome into a total of 28 segments ([Supplementary Table 4](#)), significantly more
112 than the other samples mentioned above. Two LGMs are constructed for virus integrations on the minor
113 ([Supplementary Figure 10](#)) and major allele ([Supplementary Figure 11](#)), respectively. The minor allele
114 LGM is formed by seven unit-cycles. The NO.1 unit-cycle crosses the host and viral genomes five times
115 with inversions and integrations. The NO.6 and NO.7 unit-cycles mainly contribute to the high copy
116 number of this focal region with 19 and 14 copies, respectively. The major allele LGM has a similar
117 global structure but one more unit-cycle. This LGM involves multiple genes, such as *P116*, *FGD2*, *PIM1*,
118 and *SNORD112*. The proto-oncogene *PIM1* is completely located in the ‘H9’ segment with a high copy
119 number (CN = 58). In addition, after the LGM construction, the remaining HPV16 genome copies formed
120 two types of free viral genome by our conjugate graph algorithm: the complete (CN = 655) and with
121 deletion (CN = 22). Moreover, the remaining host segments might exist as the ecDNAs.

122 On WGS data of HBV-integrated hepatocellular carcinomas

123 Next, we applied the conjugate graph algorithm on the HBV-infected hepatocellular carcinomas (HCC).
124 Six HCC samples were selected from the previous study that reported the genome-wide detection of HBV
125 integrations¹¹. Four samples have arm-level copy number changes of the HBV-integrated chromosomes,
126 but none shows the focal gain at the integration sites ([Supplementary Figure 12](#)). The constructed LGMs
127 further reveal that most of the inserted HBV DNA segments replace the host segments and consequently
128 lead to the deletions on the human genome ([Supplementary Table 5 and 6](#), [Figure 3](#), [Supplementary](#)
129 [Figure 13-17](#)), different from the tandem duplications induced by the HPV integrations. Additionally, the
130 remaining HBV segment copies imply the free viral genomes in five samples.

131 There are two HBV-integrated LGMs located on the heterozygous diploid (62T and 182T, [Supplemen-](#)
132 [tary Figure 12a-b](#)). The sample 62T has three HBV integrations at the *MLL4* oncogene on chr19q13.12
133 ([Figure 3](#)). The LGM contains two unit-cycles. The NO.1 unit-cycle duplicates the ‘H2’ host segment
134 covering the whole gene body of *ZBTB32*, the promoter and first two exons of *MLL4*, where the regulatory
135 element marks are abundant¹⁶. The structure is predicted to have no truncation effect on gene *MLL4*
136 as the coding frame is completely preserved (‘H2’ to HBV to ‘H3’). Furthermore, the inserted HBV
137 segments (‘V2’ to ‘V5’) carry the enhancer II and core promoter^{17,18}, which potentially promote the
138 *MLL4* expression¹¹. The HBV DNA segments inserted at gene *NBAS* on chr2p24.3 of sample 182T, and
139 resulted in the deletion of the host segment (‘H2’) containing 14 exons of *NBAS* ([Supplementary Figure](#)
140 [13](#)). The resolved LGM has a simple structure that the ‘H2’ host segment is replaced by the majority of
141 the HBV genome.

142 We found that the sample 13T and 260T both have copy number gains of broad regions at the
143 HBV integration loci. The 13T has HBV integrations at the *TERT* gene locus on chr5p15.33, which is
144 heterozygous triploid extending to p15.2 ([Supplementary Figure 12c](#)). As the second host segment (‘H2’,
145 CN = 2, intron region) has one copy loss, we determined that the HBV integrated on the minor allele and
146 replaced the ‘H2’ segment ([Supplementary Figure 14](#)). The HBV genome is divided into 16 segments
147 by breakpoints of two integrations and seven viral rearrangements, which are all considered in the LGM
148 construction. The inserted HBV DNA segments formed a complex structure consisting of four unit-cycles.
149 The unit-cycle NO.2 and NO.3 carry the HBV enhancers and core promoters, which might explain the
150 elevated gene expression of *TERT* in this sample¹¹. In sample 260T, the HBV integrated at the *ANK3*
151 gene on chr10 (heterozygous tetraploid, minor CN = 2, [Supplementary Figure 12d](#)). The resolved LGM

152 contains two unit-cycles including two HBV integrations and three HBV SVs (Supplementary Figure 15).
153 These two unit-cycles are completely composed of inserted HBV DNA segments and substitute the ‘H2’
154 host segment in two allele copies. Notably, our algorithm reports the unit-cycle NO.1 only in one allele,
155 and it might exist in the form of ecDNA.

156 The remaining two HCC samples displayed arm-level LOH of HBV-integrated chromosomes (101T
157 and 261T). The sample 101T has two HBV integrations at the *GLRA2* gene locus in the chrXp22.2 (CN =
158 1). The remained heterozygous state of the chrXq supports the chrXp is an LOH (Supplementary Figure
159 12e). The LGM is simple and the only unit-cycle is the inserted ‘V2’ HBV segment that replaces the ‘H2’
160 host segment (partial intron of *GLRA2*, Supplementary Figure 16). In sample 261T, there are two HBV
161 integrations located in the intron of *BBS2* gene on chr16 (CN = 1, Supplementary Figure 12f). Similar to
162 the sample 13T, the HBV genome is divided into 17 segments by two integrations and seven HBV SVs.
163 Our algorithm reports the LGM consisting of three unit-cycles with all segments from the viral genome
164 (Supplementary Figure 17). These unit-cycles have complex structures composed of deletion, tandem
165 duplication, and inversions. Again, one host deletion (‘H2’) is caused by the inserted HBV DNA.

166 Diverse affection patterns of oncovirus integrations

167 According to the resolved LGMs, we summarized five affection patterns of the oncovirus integrations
168 in our samples. Firstly, the integrations locate in a single intron, including the HPV16 in T4931 (gene
169 *GLI2*), and the HBV in 13T (*TERT*), 101T (*GLRA2*), 260T (*ANK3*), and 261T (*BBC2*). The inserted
170 viral DNA segments carry the epigenomic regulator elements, such as the long control region (LCR)
171 of HPV, and the enhancers and core promoter of HBV, which might contribute to the local expression
172 dysregulations. The abnormal alternative splicings are also reported in researches of both HPV and HBV
173 integrations^{2,13}. Secondly, the LGM covers multiple exons in sample UM-SCC-47 (HPV16 at *TP63*)
174 and 182T (HBV at *NBAS*), both likely lead to the truncation of the coding-frame. Thirdly, in sample
175 UPCI-SCC090, the LGM contains focal hyper-amplifications harboring the complete genes (*PIMI* and
176 *SNORD112*), which commonly relate to elevated expression and might induce tumorigenesis¹⁹. Fourthly,
177 in sample 62T, tandem duplication doubles the promoter and maintains the complete coding-frame of
178 *MLL4*, which is also possible subject to the inserted HBV enhancer and core promoter. Fifthly, both of
179 T6050 sample and SiHa cell line have HPV integrations at the intergenic regions between gene *KLF5*
180 and *KLF12*, whose expressions are probably affected by the amplified HPV LCR in the LGM. The HeLa
181 cell line has been reported that the inserted HPV segments have long-range chromatin interactions with
182 downstream *MYC* gene, which has high RNA expression phased in the HPV-integrated haplotype¹. This
183 implicates the cis-regulation of HPV integration is possibly mediated by the epithelium specific viral
184 enhancer remained in the LGMs.

185 Enhanced Microhomology-mediated mechanism of HPV and HBV integration

186 The resolved LGMs enable us to analyze the details of breakpoints and investigate the oncovirus integration
187 mechanisms. Here, we propose that the slipped microhomology and the small junctional insertion
188 (see Method), which are the defining characteristic of the microhomology-mediated end joining^{20,21}.
189 A previous study on cervical cancer has reported the enrichment of aligned microhomology at HPV
190 integration sites⁷. In total, 54.05% (20/37) of virus integrations have aligned microhomology at the
191 junction site. Meanwhile, the percentage became to 83.78% (31/37) when the slipped microhomology
192 and the junctional insertion are considered (Figure 4, Supplementary Table 8, Supplementary Figure
193 19-26), and it is more prevalent in HPV integrations (91.67%, 22/24) comparing with HBV (69.23%,
194 9/13). Moreover, microhomology and insertions are found at the junction sites of the majority (84.84%,
195 28/33) of the endogenous SVs: 100% (9/9) for host SVs in HPV LGMs, 83.33% (5/6) for HPV SVs, and

196 77.78% (14/18) for HBV SVs. Furthermore, the complex SVs in sample GBM0152 also shows similar
197 phenomenon that 60% (12/20) of the SV junctions carry the microhomology although the insertion data
198 is not provided (Supplementary Table 7 and 8, Supplementary Figure 18, 27 and 28). Microhomology
199 enrichment supports the suggested the replication-based mechanism, which generates these complex
200 rearrangements in this sample²².

201 We further propose that the enhanced microhomology-mediated end joining (MMEJ) may be the
202 underlying mechanism of both HPV and HBV integrations together with the endogenous host and viral
203 rearrangements in the local region. Each unit-cycle in the resolved LGMs has at least one boundary
204 harboring the microhomology or junctional insertion (Figure 5 for HPV16, Supplementary Figure 29 for
205 HPV18, and Figure 6 for HBV), suggesting that the formation of each unit-cycle is possibly initialized via
206 the replication-based template-switching event. The well-known alternative mechanisms of MMEJ include
207 the fork stalling and template switching (FoSTeS)²³ and the microhomology-mediated break-induced
208 replication (MMBIR)²⁴.

209 We hypothesize that the HPV integrations might mainly occur during the DNA replication process.
210 All of the HPV-integrated LGMs harbor the focal duplications of human genomic segments (Figure 5
211 and Supplementary Figure 29), while most (83%, 5/6) of the HBV-integrated LGMs only contain focal
212 deletions caused by the substitution of the inserted viral segments (Figure 6). Additionally, most (83%,
213 5/6) of the HPV-integrated LGMs have host SVs, but none of the HBV-integrated LGMs has. These
214 differences further enlightens us to hypothesize that the HPV integrations might mainly occur during the
215 DNA replication process (i.e., the S phase during the mitosis) and form complex structure via the rolling
216 circle replication^{2,25}, and the HBV integrations likely arise under the DNA repairing circumstance and are
217 prone to produce simple structure.

218 DISCUSSION

219 The detection of the oncovirus integrations and structural variations is widely developed, and the effects
220 and formation mechanisms of integrations are extensively studied. However, the investigation in the
221 local genomic region with the linear complex structure remains limited. Our conjugate graph algorithm
222 constructs the local genomic maps of HPV and HBV integrations. It provides details to research the
223 regulatory influences on local genes and the underlying formation mechanisms, combining with the copy
224 changes of rearranged host and viral genomic segments.

225 The local genomic map sometimes contains segments shorter than the insert-size of the paired-end
226 sequencing data or harboring repeated DNA sequence, leading to the underestimated copy number
227 calculation from the average depth of the relevant segment. Similarly, the insufficient split-reads count
228 also reduces the copy number of the SV and VIT junctions. These inaccurate copy numbers will hinder the
229 accuracy of the LGM results. The integer linear programming in our algorithm overcomes such difficulties
230 by assigning smaller weights for relevant vertex and edge in the graph model, and successfully construct
231 the optimal LGM with minimum offsets.

232 The unit-cycles in the conjugate graph model are the basic structural units representing the connections
233 of segments and junctions. Unit-cycles are merged in ways to form different LGM solutions. Due to
234 the short reads in NGS data, the factual orders of merging cannot be determined technically. Here, we
235 select the solution with the simplest structure to report, following Ockham's razor. The candidate LGMs
236 could be further filtered based on the phased structure obtained from the advanced DNA sequencing,
237 such as 10x long-range linked-reads and single-molecule reads. Furthermore, when the virus-integrated
238 chromosome is heterozygous unbalanced polyploidy, it is difficult to determine whether the paternal or
239 maternal allele carries the viral insertions based on short reads. The mutations phasing between the local

240 genomic segments and the flanking regions of the constructed LGMs will help solve the problem. Multiple
241 sequencing technology can provide this additional information, including the mate-pair library with large
242 insert-size¹, 10x linked-reads²⁶, and Hi-C²⁷.

243 When the virus-integrated allele has more than one copy, our algorithm will try to distribute all
244 relevant unit-cycles evenly among the multiple copies to minimize the differences. However, the actual
245 count of the unit-cycles might change within different allele copies (e.g., [Supplementary Figure 5d and](#)
246 [6d](#)). According to the circular structure and common hyperamplification, the unit-cycles can be the
247 circular extrachromosomal DNA, which needs more investigations in the future study, for example, the
248 inter-chromosomal interactions from Hi-C data²⁸. Recent studies have proved that the ecDNA promotes
249 oncogene expressions and drives tumor evolution^{29,30}. The speculated oncovirus integrations related
250 ecDNAs probably have similar carcinogenic functions.

251 From the constructed LGMs, we can get the structural investigation on the formation mechanism
252 of the oncovirus integrations and endogenous SVs in the local genomic region. To our knowledge, this
253 is the first time to combine the aligned/slipped microhomology and junctional insertion together in the
254 mechanism research of viral integrations. The enrichment of microhomology and insertion indicates that
255 both the HPV and HBV integrations may be generated via the microhomology-mediated end joining
256 (MMEJ), promoted during the HPV16 infection³¹. Furthermore, the focal duplicates and host SVs are
257 frequent in the HPV-integrated LGMs, while the HBV-integrated LGMs are abundant in focal deletions
258 and complex virus SVs. This apparent difference suggests that the formation circumstances of HPV and
259 HBV integration may be different. To answer this question, it needs more samples and experiments in a
260 future study.

261 In summary, we present a conjugate graph algorithm to construct the local genomic map of virus
262 integrations mediated host and viral rearrangements from the whole genome sequencing data. The
263 algorithm can normalize imperfect sequencing depths and elucidate the linear DNA structure at virus
264 integration sites. The simulation tests and the applications on WGS data of cancer samples prove the
265 reliability and credibility of our model. The results shed light on the genomic dysregulation and formation
266 mechanism of the oncovirus integrations.

267 **DATA AVAILABILITY**

268 LGM analysis data of all samples are available in the [Supplementary Data 1](#) zip file.

269 **CODE AVAILABILITY**

270 The conjugate graph algorithm code and LGM construction pipeline are available in github repository
271 (<https://github.com/deepomicslab/FuseSV>).

272 **AUTHOR CONTRIBUTIONS**

273 S.L. proposed the conjugate graph model and designed algorithm to obtain the LGM paths. W.J. and
274 C.X. implemented and optimized the algorithms. W.J. and C.X. performed the WGS data analysis. C.X.
275 performed the simulation work. W.J. developed the FuseSV. W.J. summarised and compared the viral
276 integration features and mechanisms. W.J., C.X. and S.L. wrote and revised the manuscript. All authors
277 reviewed the article and approved the final manuscript.

278 **FUNDING**

279 This work was supported by the General Research Fund (GRF) Projects 9042348 (CityU 11257316).

280 CONFLICT OF INTEREST

281 Conflict of interest statement. None declared.

282 METHODS

283 Conjugate graph

284 We first consider the problem of reconstructing the disrupted genome with perfect NGS data and breakpoint
 285 detection, giving rise to no ambiguity. Based on the breakpoints, we derive two sets of reference regions
 286 respectively from within (1) the human/host genome (denoted $H = \{h_1, h_2, \dots, h_m\}$), and (2) the virus
 287 genome (denoted $V = \{v_1, v_2, \dots, v_n\}$). Each region $s \in H \cup V$ can be either on the positive or the negative
 288 DNA chain, denoted s^+ and s^- respectively. Associated with each region s is a number $c(s)$ which
 289 indicates the number of times the region is referenced in the NGS data, that is, the copy number of each
 290 region.

291 We assume that the NGS data consists of a set J of junctions, each of which connects two regions
 292 in $H \cup V$. Each $\mu \in J$ is an ordered pair $\langle u_i^{s_i}, u_j^{s_j} \rangle$, where $s_i, s_j \in \{+, -\}$ and $u_i, u_j \in H \cup V$. As DNA
 293 possesses a double helix structure, a junction segment is equivalent to its reverse complementary segment.
 294 Therefore, two junction segments $\langle u_i^{s_i}, u_j^{s_j} \rangle$ and $\langle u_j^{-s_j}, u_i^{-s_i} \rangle$ are equivalent, where $-s_i$ denotes the reverse
 295 complementary chain of s_i . Note that junction segments are directed, that is, $\langle u_i^{s_i}, u_j^{s_j} \rangle$ is different from
 296 $\langle u_j^{s_j}, u_i^{s_i} \rangle$. Meanwhile, each $\mu \in J$ may also occur multiple times in NGS data, denoted by $c(\mu)$, reflecting
 297 the number of sequencing reads crossing a junction breakpoint, i.e., the split-reads.

298 We construct a junction graph $G = (U, E)$, a directed graph with vertices $U = \{u^+, u^- | u \in H \cup V\}$
 299 and edges $E = \{\mu^+, \mu^- | \mu \in J\}$. We assume that G is connected, and that $\langle u_i^{s_i}, u_j^{s_j} \rangle \in E$ if and only if
 300 $\langle u_j^{-s_j}, u_i^{-s_i} \rangle \in E$. We refer to $\langle u_i^{s_i}, u_j^{s_j} \rangle$ and $\langle u_j^{-s_j}, u_i^{-s_i} \rangle$ as conjugate edges, and the vertices u^+ and u^- as
 301 conjugate vertices; we also call this type of graph as *conjugate graph*.

302 Local genomic map.

303 Local genomic map (LGM) is defined as a path over the conjugate graph, where each $u \in H \cup V$ is visited
 304 exactly $c(u)$ times, and each $\mu \in J$ is traversed exactly $c(\mu)$ times. Hence an LGM of length ℓ gives us a
 305 sequence of visited vertices $(u_1^{s_1}, u_2^{s_2}, \dots, u_\ell^{s_\ell})$ where $u_i \in H \cup V$ and $s_i \in \{+, -\}$. Without loss of generality,
 306 we assume that both the starting and ending point of an LGM are from the positive chain of the host, that
 307 is, $u_1^{s_1} = h_1^+$, $u_\ell^{s_\ell} = h_m^+$. Note that we assume the path is along the positive strand. If the local genomic map
 308 contains a vertex of a negative strand, it implies an inversion event.

309 Existence of LGM in the conjugate graph.

310 We show that it is possible to know if a conjugate graph contains an LGM. Define the in-copy of a vertex
 311 u^s as the total copy number of edges with u^s as target vertex; that is, $in(u^s) = \sum c(\langle u_i^{s_i}, u^s \rangle)$. Similarly, we
 312 can define $out(u^s) = \sum c(\langle u^s, u_i^{s_i} \rangle)$.

313 Clearly, if G contains an LGM, then it satisfies the following properties:

314 (1) $c(u) = in(u^-) + out(u^+) = out(u^-) + in(u^+), \forall u \in H \cup V - \{h_1, h_m\};$

315 (2) $c(h_1) = in(h_1^-) + out(h_1^+) = out(h_1^-) + in(h_1^+) + 1$, and

316 $c(h_m) = in(h_m^-) + out(h_m^+) + 1 = out(h_m^-) + in(h_m^+),$

317 (3) a path exists from h_1^+ to u^+ or u^- , and a path exists from u^+ or u^- to h_m^+ .

318 Note that, to form a path, we add a putative edge links *sink* vertex back to *source* vertex (property (2)). We
319 refer to the properties (1) and (2) as degree balance, and properties (3) as reachability.

320 An LGM of a conjugate graph shares similarities with an Eulerian path of an Eulerian Graph. The
321 differences lie in that the vertices and edges are conjugated in LGMs. We adopt an algorithm derived from
322 the Hierholzer's algorithm to find LGMs.

323 Similar to Hierholzer's algorithm, the algorithm for finding LGMs starts from an arbitrary vertex u_i^s ,
324 and recursively follow a trail of out-going edges from the vertex until returning to the vertex u_i^s itself,
325 forming a circuit C . Due to the conjugation property, after we traverse an edge $\langle u_i^{s_i}, u_{i+1}^{s_{i+1}} \rangle$, the next edge
326 with unused copy numbers we can choose from is either in the form of $\langle u_{i+1}^{s_{i+1}}, * \rangle$ or $\langle *, u_{i+1}^{-s_{i+1}} \rangle$, where $*$
327 represents a wildcard. Similar to the Eulerian path argument, as long as we have a vertex in the graph with
328 unused copy numbers, the degree balance property ensures that it is always possible to find a circle that
329 starts from and ends in the vertex. Therefore, given a vertex u_j^s in the circle C that still has unused copy
330 numbers, we can find a new circle C' which can be merged with C to form a longer circular path. The
331 above process is performed iteratively until we drain the copy numbers in the conjugate graph.

332 **Theorem 1** *A conjugate graph contains a local genomic map if and only if vertex u_i^+ or u_i^- of any region*
333 *$u \in H \cup V$ is in a path from h_1^+ to h_m^+ or a path from h_m^- to h_1^- , and furthermore, u satisfies the degree*
334 *balance condition.*

335 **Correcting the copy numbers.**

336 In the case of imperfect NGS data, the estimated copy numbers may contain fluctuations. The conjugate
337 graph constructed would then lead to erroneous LGMs or no LGMs at all. To prevent such cases, we
338 perform corrections to the graph under the principle that a corrected graph would contain an LGM, that is,
339 the resultant conjugate graph should fulfill the two properties of reachability and degree balance.

340 It is possible to have missing edges or vertices due to reasons such as low sequencing depth and failure
341 in detecting breakpoints which could break the property of reachability. To solve this in a parsimonious
342 manner, we only consider the insertion of additional edges that are likely to exist in a normal genome, such
343 as $\langle h_1^+, h_2^+ \rangle, \langle h_2^+, h_3^+ \rangle, \dots, \langle h_{m-1}^+, h_m^+ \rangle, \langle v_1^+, v_2^+ \rangle, \langle v_2^+, v_3^+ \rangle, \dots, \langle v_{n-1}^+, v_n^+ \rangle$, and $\langle v_n^+, v_1^+ \rangle$ if the virus in question
344 has a circular DNA structure.

345 We propose an integer linear programming approach to address the degree balance property. Assign
346 each region u and each junction μ to a target copy number $t(u)$ or $t(\mu)$ (Equations 1d and 1e), respectively,
347 to satisfy the degree balance property (Equation 1b). The objective is to minimize the disagreement
348 between the observed copy number and the target copy number (Equation 1a).

minimize

$$\sum_u \varepsilon_u + \sum_\mu \varepsilon_\mu \tag{1a}$$

subject to

$$in(u) = out(u) = t(u), \quad \forall u \in H \cup V \quad (1b)$$

$$t(h_1) = t(h_m), \quad (1c)$$

$$-\varepsilon_u \leq c(u) - t(u) \leq \varepsilon_u, \quad \forall u \in H \cup V \quad (1d)$$

$$-\varepsilon_\mu \leq c(\mu) - t(\mu) \leq \varepsilon_\mu, \quad \forall \mu \in J \quad (1e)$$

$$\varepsilon_u, \varepsilon_\mu \in \mathbb{R}^+ \quad (1f)$$

$$t(u), t(\mu) \in \mathbb{Z}^+ \quad (1g)$$

349 Additional domain knowledge can be incorporated to further refine the linear programming process.
350 For instance, the copy number of a virus DNA segment can be easily overestimated due to the existence of
351 free virus in cells and the copy number of DNA segments too short can be underestimated due to alignment
352 difficulties. In those circumstances, we may assign lower penalty to the change of their copy numbers
353 during linear programming.

354 Finding LGM.

355 Our basic algorithm identifies an LGM path by iteratively selecting edges with non-zero weight at random
356 to find circuits (called unit-cycles) and merge them. When determining the order of unit-cycles to be
357 merged, the basic algorithm has no preference for any one possible LGM over the others. Principles of
358 parsimony can be adopted; for instance, edges which create no breakpoints are preferred over edges that
359 can not be found in the reference.

360 Simulation evaluation.

361 For repeated testing, we referred to the literature and chose one subtype (B2) of HBV and two subtypes
362 (HPV16 and HPV18) of HPV to simulate virus integrations. Virus reference genome was randomly
363 divided into segments. Integration sites in human reference genome (GRCh37) were randomly chosen
364 inside reported integration hotspot regions (Near gene *TERT*, *CCNE1*, *KMT2B* for HBV, *LRP1B*, *FHIT* for
365 HPV16 and at the upstream of *MYC* for HPV18). Lengths of segments are random values within empirical
366 ranges based on our research (2000-35000bp for each human segment, 1000-4000bp and 500-2000bp
367 for each HPV and HBV segment respectively). Human and virus segments were then concatenated with
368 regard to the different integration modes ([Supplementary Figure 1](#)) and saved as FASTA format files.

369 Simulated sequencing was performed using SAMtools wgsim utility at 30X average depth and 150bp
370 read length³². 15X coverage is also introduced to complex integration cases to test algorithm's sensitivity
371 in lower coverage situations. Since wgsim simulates all bases with the same quality score, we have also
372 replaced the quality value of each sequence with qualities from actual 150bp sequencing data to resemble
373 real sequencing results. Simulated sequencing reads are aligned back to human reference genome. Depth
374 of each segment is calculated using SAMTools depth command on uniquely mappable regions. Sequence
375 mappability is determined by performing a 150bp read length simulated sequencing and aligning back to
376 human genome. Regions covered by multiple-mapped reads are excluded in depth calculation.

377 All simulated samples are considered to have homogeneous integrations. Integration breakpoints
378 are called by FuseSV. Structural variations on human and virus genome are called by Meerkat²². Local
379 genomic map of all simulated cases are recovered by the algorithm. Complex cases can have multiple
380 possible LGMs.

381 **Evaluation of integer linear programming.**

382 Test were performed to evaluate the robustness of the linear programming approach. We chose six simulated
383 samples with high segment counts and tandem duplications, and then randomly assign fluctuations to
384 segments' depth. The fluctuation values are within $[-c_{avg}/2, c_{avg}/2]$ where c_{avg} is the average haploid
385 sequencing depth. The randomization is done 1000 times for each chosen samples and we performed
386 linear programming for each iteration. 98% of all 6000 cases are corrected to our expected copy number
387 profile. The other 2% cases had cycles with one more or one less copy.

388 **WGS data analysis.**

389 Tumor samples from cell lines have no corresponding normal sample. From our data repository, we selected
390 a normal sample with no CNV at known integration regions. WGS data was aligned to human reference
391 genome (GRCh37) using BWA³³. HPV and HBV integrations are detected by FuseSV. Meerkat²² and
392 seeksv³⁴ are used to call structural variations on human and viral genome. CNV profile, tumor purity and
393 ploidy are generated by Patchwork³⁵. Segment depths were normalized against GC bias by deepTools³⁶.

394 **Copy number calculation of segments in LGM.**

395 In tumor and control samples, the average depth of each segment was calculated from GC-bias corrected
396 alignments using SAMtools depth utility with minimum base and mapping quality set as 5 and 10,
397 respectively. Call the purity of tumor-cell in tumor sample as $Purity_t$, the average ploidy of pure tumor-cell
398 in tumor sample as $Ploidy_t$, both $Purity_t$ and $Ploidy_t$ are from Patchwork result³⁵. And the average ploidy
399 of pure normal cell as $Purity_n$ (assumed as 2). The average ratio of the DNA from normal cells ($Ratio_n$) in
400 each genomic region was calculated as below:

$$Ratio_n = 1 - \frac{Purity_t \times Ploidy_t}{Purity_t \times Ploidy_t + (1 - Ploidy_t) \times Purity_n} \quad (2a)$$

For each segment, the germline copy ratio of control normal sample ($CopyRatio_n$) can be obtained from WGS data (generally is one). The depth of this segment from pure tumor-cell in tumor sample ($Depth_t$) could be obtained from the observed segment average depth ($Depth_o$) and the whole genome average depth ($Depth_g$) of this segment via formula

$$Depth_t = Depth_o - Depth_g \times Ratio_n \times CopyRatio_n. \quad (3a)$$

401 The major and minor copy numbers of LGM region are determined from Patchwork CNV result³⁵.
402 The bilateral segments (*source* and *sink*) of LGM are used to obtain the haploidy depth. The copy numbers
403 of other segments and junctions of VITs and SVs in LGM are calculated according to the segment average
404 depth and split-reads count.

405 **Construction of virus genome in each sample.**

406 We extract the unmapped and soft-clipped reads from alignment results of WGS and realign these reads to
407 HPV and HBV references from the NCBI nucleotide database. The viral subtype with the most uniquely
408 aligned reads and coverage of at least 20% was selected as the major one. Mutations (SNV and InDel)
409 are detected in an iterative process to amend the viral genome sequence until no more mutations can be
410 identified. This procedure is also adopted in other study⁴. The constructed viral genome is applied in
411 the subsequent viral integration and LGM analyses of the corresponding cancer sample. Two reads types
412 (span-reads and split-reads) are sought to support the viral integrations.

413 **Microhomology at the junction sites.**

414 Microhomology bases (MHs) are detected at the junction sites of viral integrations and structural variations.
415 Two types of MH bases are considered: the aligned and slipped MHs. The slipped MH bases are required
416 to have $\geq 3bp$ size and at least one base overlap. Criteria to determine the MH existence: i) the MH
417 bases ($\geq 2bp$) must cover or locate next to the junction site, or ii) locate in flanking region ($\leq 5bp$)
418 of the junction site and the MH size must be $\geq 4bp$ according our previous finding⁷, or iii) merge
419 1bp-gaped neighbour MH bases and subject to the second criterion.

420 **References**

- 421 1. Adey, A. *et al.* The haplotype-resolved genome and epigenome of the aneuploid hela cancer cell line.
422 *Nature* **500**, 207–211 (2013).
- 423 2. Akagi, K. *et al.* Genome-wide analysis of hpv integration in human cancers reveals recurrent, focal
424 genomic instability. *Genome research* **24**, 185–199 (2014).
- 425 3. Satou, Y. *et al.* The retrovirus htlv-1 inserts an ectopic ctf-binding site into the human genome.
426 *Proc. Natl. Acad. Sci.* **113**, 3054–3059 (2016).
- 427 4. Chen, X.-P. *et al.* Viral integration drives multifocal hcc during the occult hbv infection.
428 *J. Exp. & Clin. Cancer Res.* **38**, 261 (2019).
- 429 5. Kan, Z. *et al.* Whole-genome sequencing identifies recurrent mutations in hepatocellular carcinoma.
430 *Genome Res.* **23**, 1422–1433 (2013).
- 431 6. Zhao, L.-H. *et al.* Genomic and oncogenic preference of hbv integration in hepatocellular carcinoma.
432 *Nat. Commun.* **7**, 12992 (2016).
- 433 7. Hu, Z. *et al.* Genome-wide profiling of hpv integration in cervical cancer identifies clustered genomic
434 hot spots and a potential microhomology-mediated integration mechanism. *Nat. Genet.* **47**, 158–163
435 (2015).
- 436 8. Network, C. G. A. R. *et al.* Integrated genomic and molecular characterization of cervical cancer.
437 *Nature* **543**, 378 (2017).
- 438 9. Murakami, Y. *et al.* Large scaled analysis of hepatitis b virus (hbv) dna integration in hbv related
439 hepatocellular carcinomas. *Gut* **54**, 1162–1168 (2005).
- 440 10. Zhao, X. *et al.* Dr. vis: a database of human disease-related viral integration sites. *Nucleic Acids Res.*
441 *gkr1142* (2011).
- 442 11. Sung, W.-K. *et al.* Genome-wide survey of recurrent hbv integration in hepatocellular carcinoma.
443 *Nat. Genet.* **44**, 765–769 (2012).
- 444 12. Ojesina, A. I. *et al.* Landscape of genomic alterations in cervical carcinomas. *Nature* **506**, 371–375
445 (2014).
- 446 13. Nault, J.-C. *et al.* Recurrent aav2-related insertional mutagenesis in human hepatocellular carcinomas.
447 *Nat. Genet.* **47**, 1187 (2015).
- 448 14. Flores, E. R. *et al.* Tumor predisposition in mice mutant for p63 and p73: evidence for broader tumor
449 suppressor functions for the p53 family. *Cancer Cell* **7**, 363–373 (2005).
- 450 15. Melino, G. p63 is a suppressor of tumorigenesis and metastasis interacting with mutant p53.
451 *Cell Death & Differ.* **18**, 1487–1499 (2011).

- 452 **16.** Consortium, E. P. *et al.* An integrated encyclopedia of dna elements in the human genome. *Nature*
453 **489**, 57–74 (2012).
- 454 **17.** Kramvis, A. & Kew, M. The core promoter of hepatitis b virus. *J. Viral Hepat.* **6**, 415–427 (1999).
- 455 **18.** Seeger, C. & Mason, W. S. Hepatitis b virus biology. *Microbiol. Mol. Biol. Rev.* **64**, 51–68 (2000).
- 456 **19.** Narlik-Grassow, M., Blanco-Aparicio, C. & Carnero, A. The pim family of serine/threonine kinases
457 in cancer. *Medicinal Res. Rev.* **34**, 136–159 (2014).
- 458 **20.** Bennardo, N., Cheng, A., Huang, N. & Stark, J. M. Alternative-nhej is a mechanistically distinct
459 pathway of mammalian chromosome break repair. *PLoS Genet.* **4**, e1000110 (2008).
- 460 **21.** McVey, M. & Lee, S. E. Mmej repair of double-strand breaks (director’s cut): deleted sequences and
461 alternative endings. *Trends Genet.* **24**, 529–538 (2008).
- 462 **22.** Yang, L. *et al.* Diverse mechanisms of somatic structural variations in human cancer genomes. *Cell*
463 **153**, 919–929 (2013).
- 464 **23.** Lee, J. A., Carvalho, C. M. & Lupski, J. R. A dna replication mechanism for generating nonrecurrent
465 rearrangements associated with genomic disorders. *Cell* **131**, 1235–1247 (2007).
- 466 **24.** Hastings, P., Ira, G. & Lupski, J. R. A microhomology-mediated break-induced replication model for
467 the origin of human copy number variation. *PLoS Genet.* **5**, e1000327 (2009).
- 468 **25.** Kusumoto-Matsuo, R., Kanda, T. & Kukimoto, I. Rolling circle replication of human papillomavirus
469 type 16 dna in epithelial cell extracts. *Genes to Cells* **16**, 23–33 (2011).
- 470 **26.** Zheng, G. X. *et al.* Haplotyping germline and cancer genomes with high-throughput linked-read
471 sequencing. *Nat. Biotechnol.* (2016).
- 472 **27.** Selvaraj, S., Dixon, J. R., Bansal, V. & Ren, B. Whole-genome haplotype reconstruction using
473 proximity-ligation and shotgun sequencing. *Nat. Biotechnol.* **31**, 1111 (2013).
- 474 **28.** Harewood, L. *et al.* Hi-c as a tool for precise detection and characterisation of chromosomal
475 rearrangements and copy number variation in human tumours. *Genome Biol.* **18**, 125 (2017).
- 476 **29.** Turner, K. M. *et al.* Extrachromosomal oncogene amplification drives tumour evolution and genetic
477 heterogeneity. *Nature* **543**, 122–125 (2017).
- 478 **30.** Wu, S. *et al.* Circular ecdna promotes accessible chromatin and high oncogene expression. *Nature*
479 **575**, 699–703 (2019).
- 480 **31.** Leeman, J. E. *et al.* Human papillomavirus 16 promotes microhomology-mediated end-joining.
481 *Proc. Natl. Acad. Sci.* **116**, 21573–21579 (2019).
- 482 **32.** Li, H. *et al.* The sequence alignment/map format and samtools. *Bioinformatics* **25**, 2078–2079 (2009).
- 483 **33.** Li, H. & Durbin, R. Fast and accurate short read alignment with burrows–wheeler transform.
484 *Bioinformatics* **25**, 1754–1760 (2009).
- 485 **34.** Liang, Y. *et al.* Seeksv: an accurate tool for somatic structural variation and virus integration detection.
486 *Bioinformatics* **33**, 184–191 (2017).
- 487 **35.** Mayrhofer, M., DiLorenzo, S. & Isaksson, A. Patchwork: allele-specific copy number analysis of
488 whole-genome sequenced tumor tissue. *Genome Biol.* **14**, R24 (2013).
- 489 **36.** Benjamini, Y. & Speed, T. P. Summarizing and correcting the gc content bias in high-throughput
490 sequencing. *Nucleic Acids Res.* **40**, e72–e72 (2012).

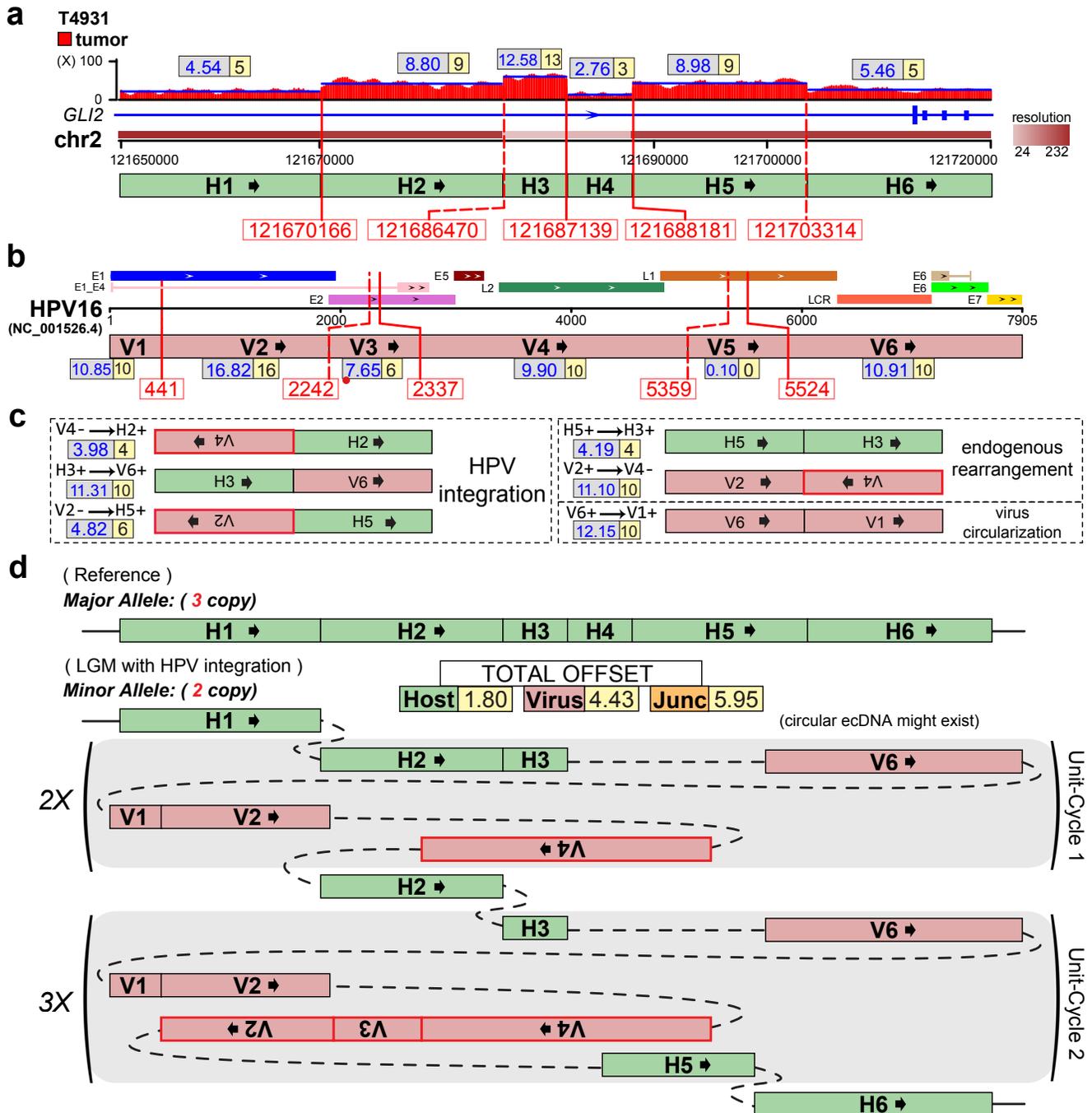


Figure 2. Presentative LGM at HPV integration sites (gene *GLI2*) on chr2 (minor allele) of the T4931 sample⁷. **(a)** Human genomic region flanking HPV integrations are divided into six segments (H1~H6, in different resolutions) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. **(b)** Segmentation (V1~V6) of HPV16 genome by VITs and SVs. The segment (V3) less than 100bp is marked with red dot. **(c)** Variant segment junctions utilized in Conjugate graph. **(d)** Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (major allele) and that harbours HPV16 integrations (minor allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the unit-cycles might exist as circular extrachromosomal DNA (ecDNA).

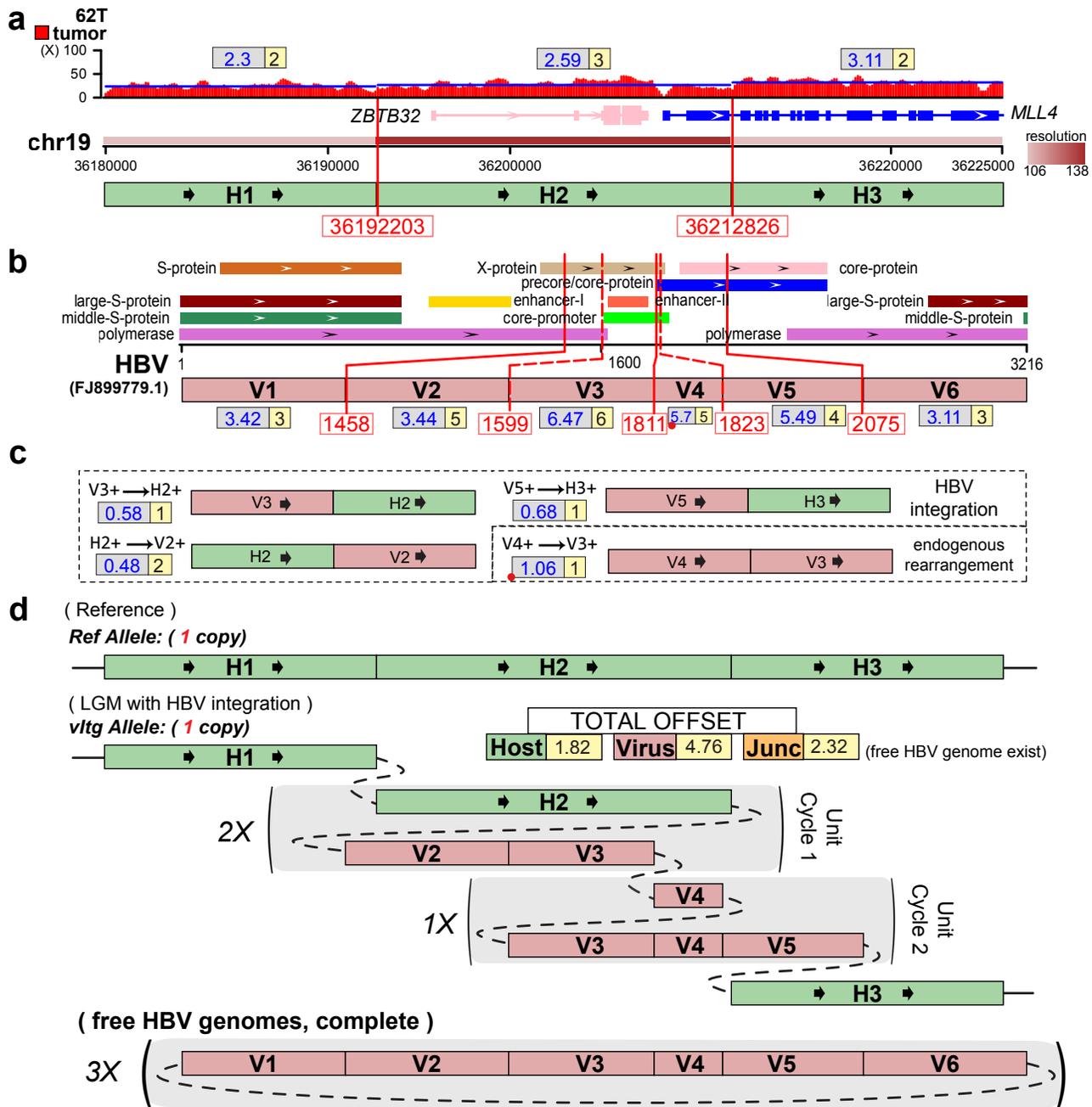


Figure 3. Presentative LGM at HBV integration sites (gene *MLL4*) on chr19 (minor allele) of the 62T HCC sample¹¹. **(a)** Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. **(b)** Segmentation (V1~V6) of HBV genome by VITs and SVs. The segment less than 100bp is marked with red dot. **(c)** Variant segment junctions utilized in Conjugate graph. **(d)** Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele and that harbours HBV integrations. Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the free HBV genomes might exist.

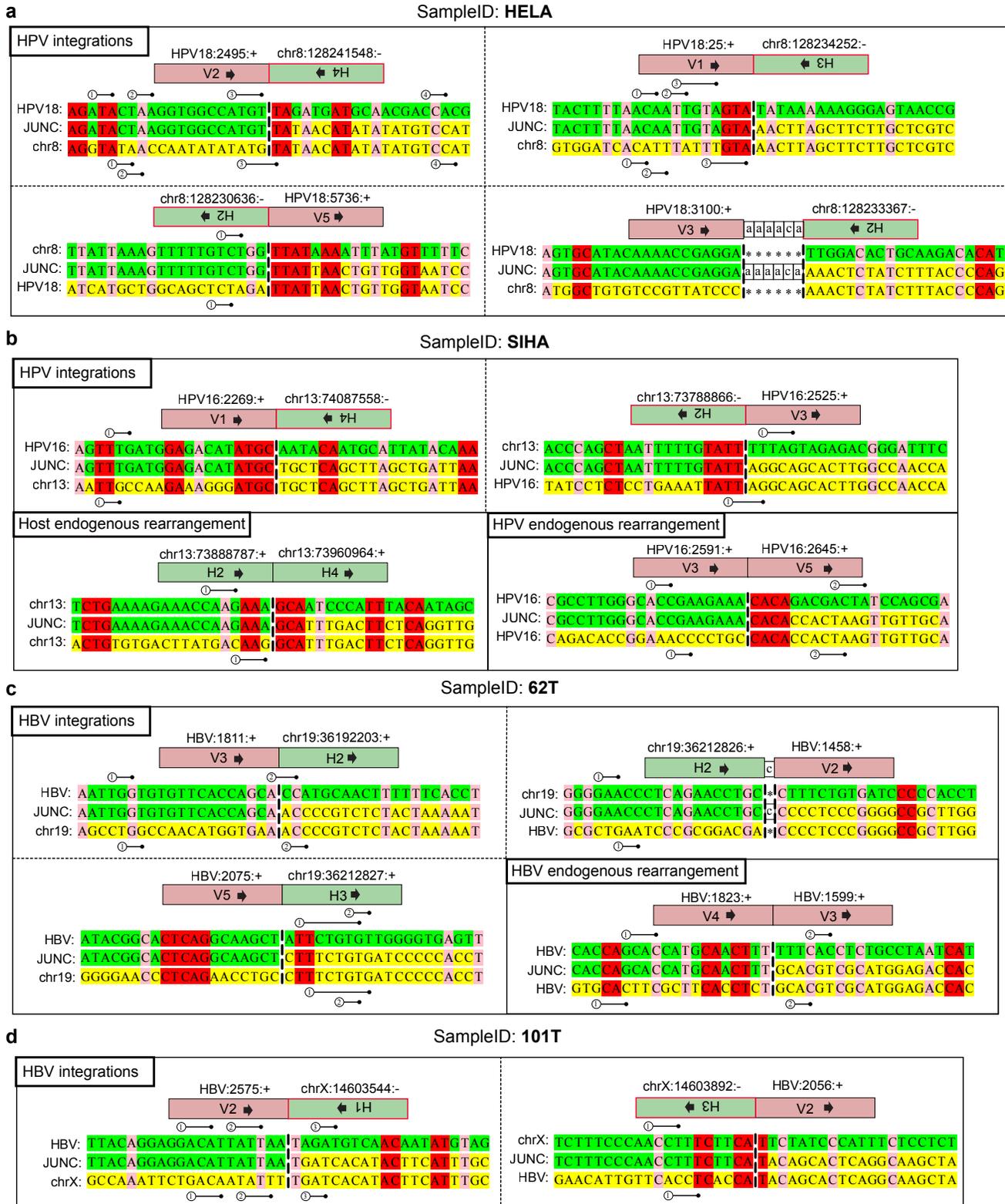


Figure 4. Alignment of the sequence around the integration site between the host and virus genome, and the endogenous rearrangements, in the samples (a) HeLa, (b) SiHa, (c) 62T, (d) 101T. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved LGMs.

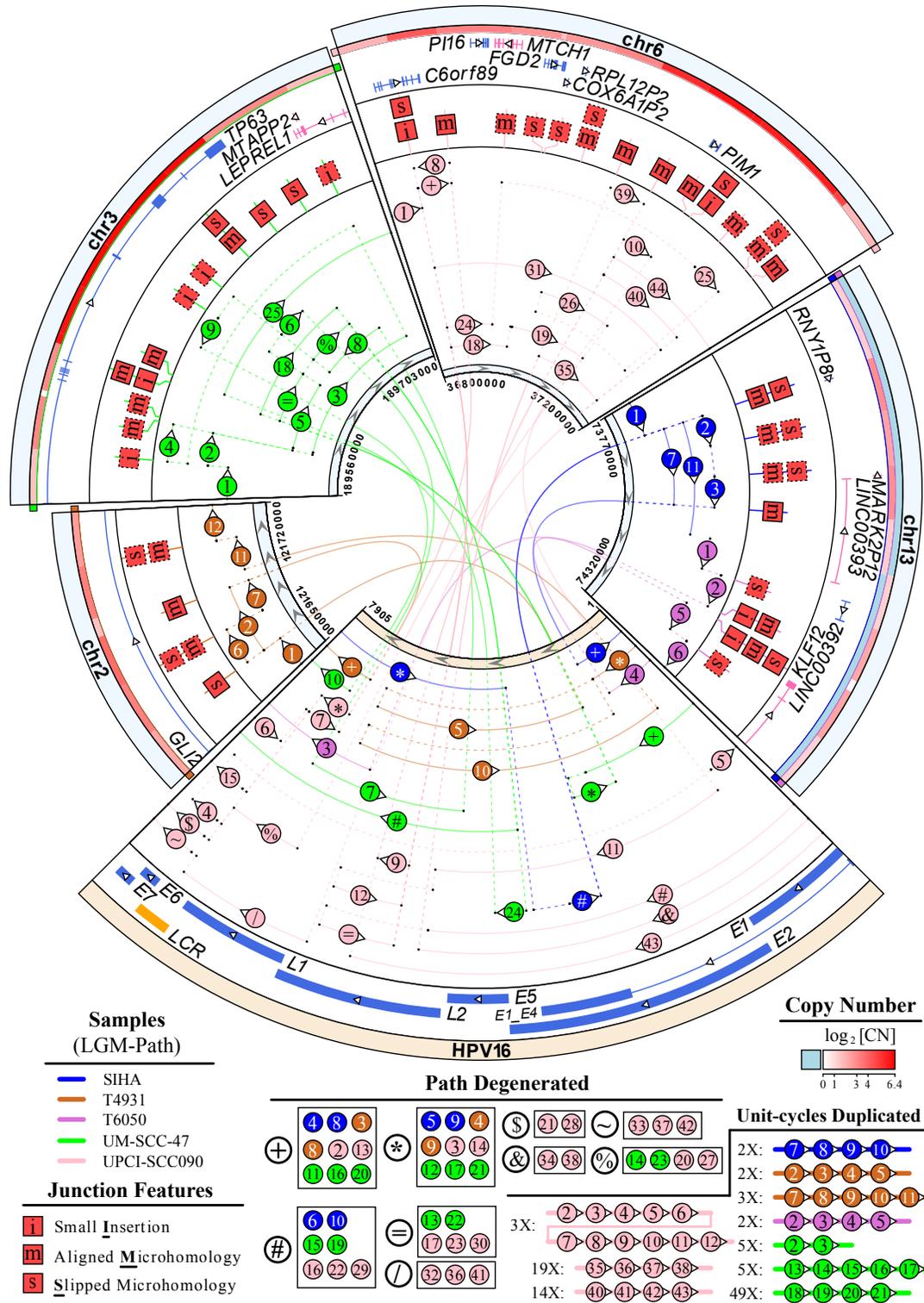


Figure 5. Features of HPV16 integrated LGMs in five samples. Human genomic segments related to HPV16 integrated LGMs are shown as sectors with their relevant LGM path in sample specific colours. Segments of LGMs are denoted by circled numbers in sequence, where some are degenerated by symbols for simplification. The numbered arcs might represent multiple sequential segments (Supplementary Table 9). Repeat times of unit-cycles in LGMs are stated in figure legend. Features of HPV16 integration (solid box) and host rearrangement (dotted) sites are depicted as single-letter icons. DNA copy number (CN) is displayed in gradient red colour (light-blue for regions outside of the LGM), with bilateral labels in relevant sample colour. The HPV16 genome reference is NC_001526.4 from the NCBI Nucleotide database.

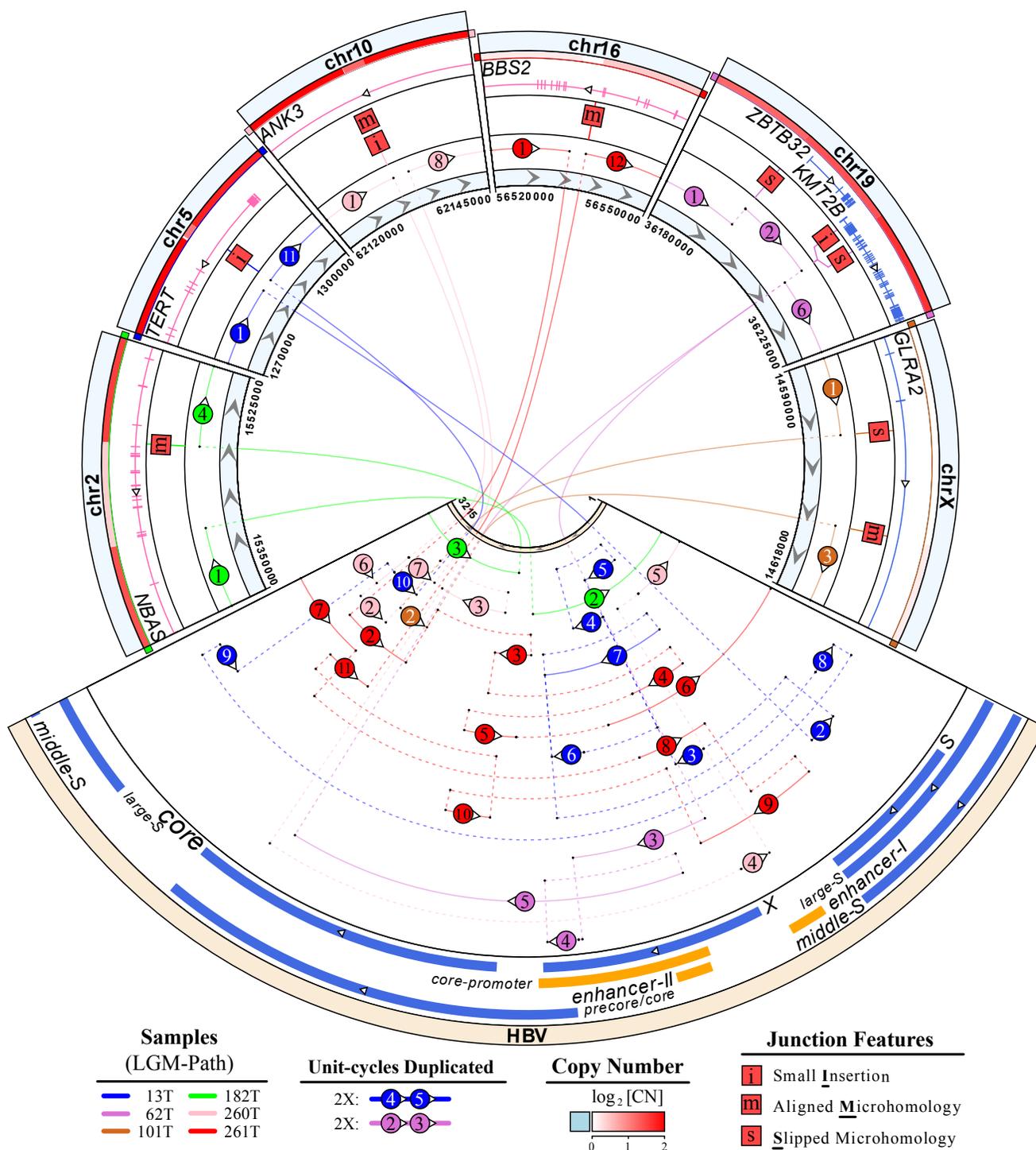
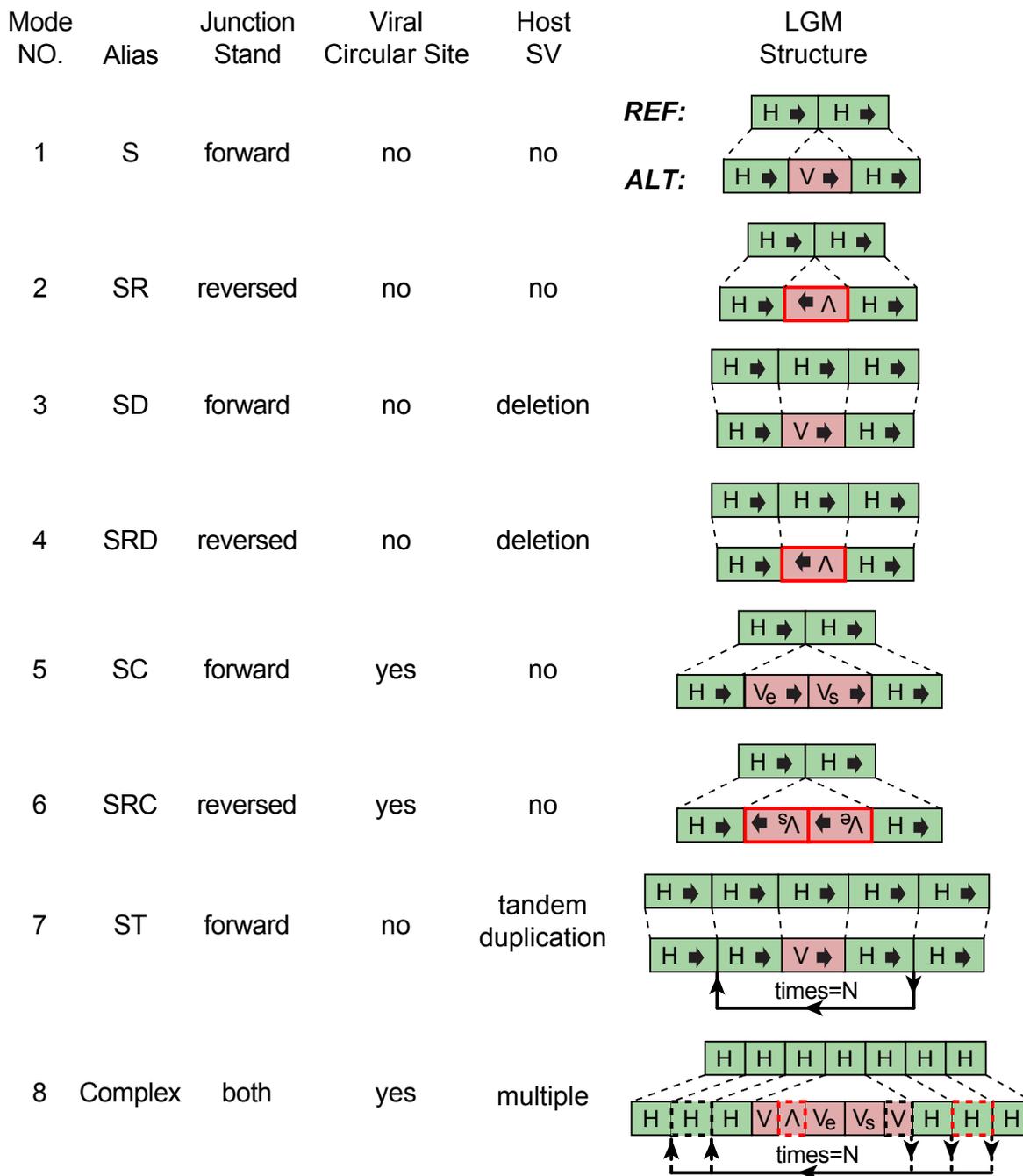
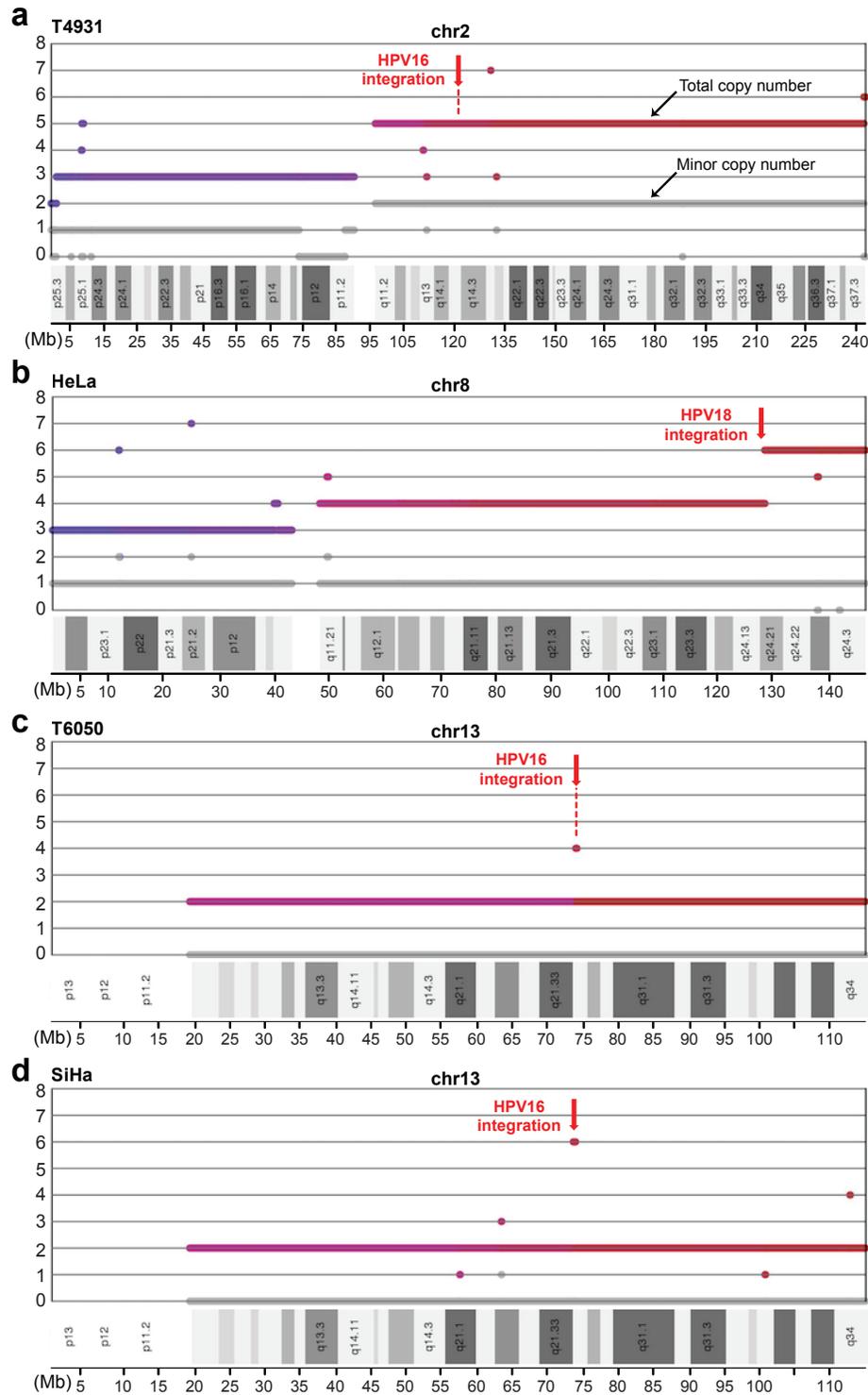


Figure 6. Features of HBV integrated LGMs in six HCC samples. Human genomic segments related to HBV integrated LGMs are shown as sectors with their relevant LGM path in sample specific colours. Segments of LGMs are denoted by circled numbers in sequence, where some are degenerated by symbols for simplification. The numbered arcs might represent multiple sequential segments (Supplementary Table 9). Repeat times of unit-cycles in LGMs are stated in figure legend. Features of HBV integrated sites are depicted as single-letter icons. DNA copy number (CN) is displayed in gradient red colour (light-blue for regions outside of the LGM), with bilateral labels in relevant sample colour. The HBV genome reference is FJ899779.1 from the NCBI Nucleotide database.

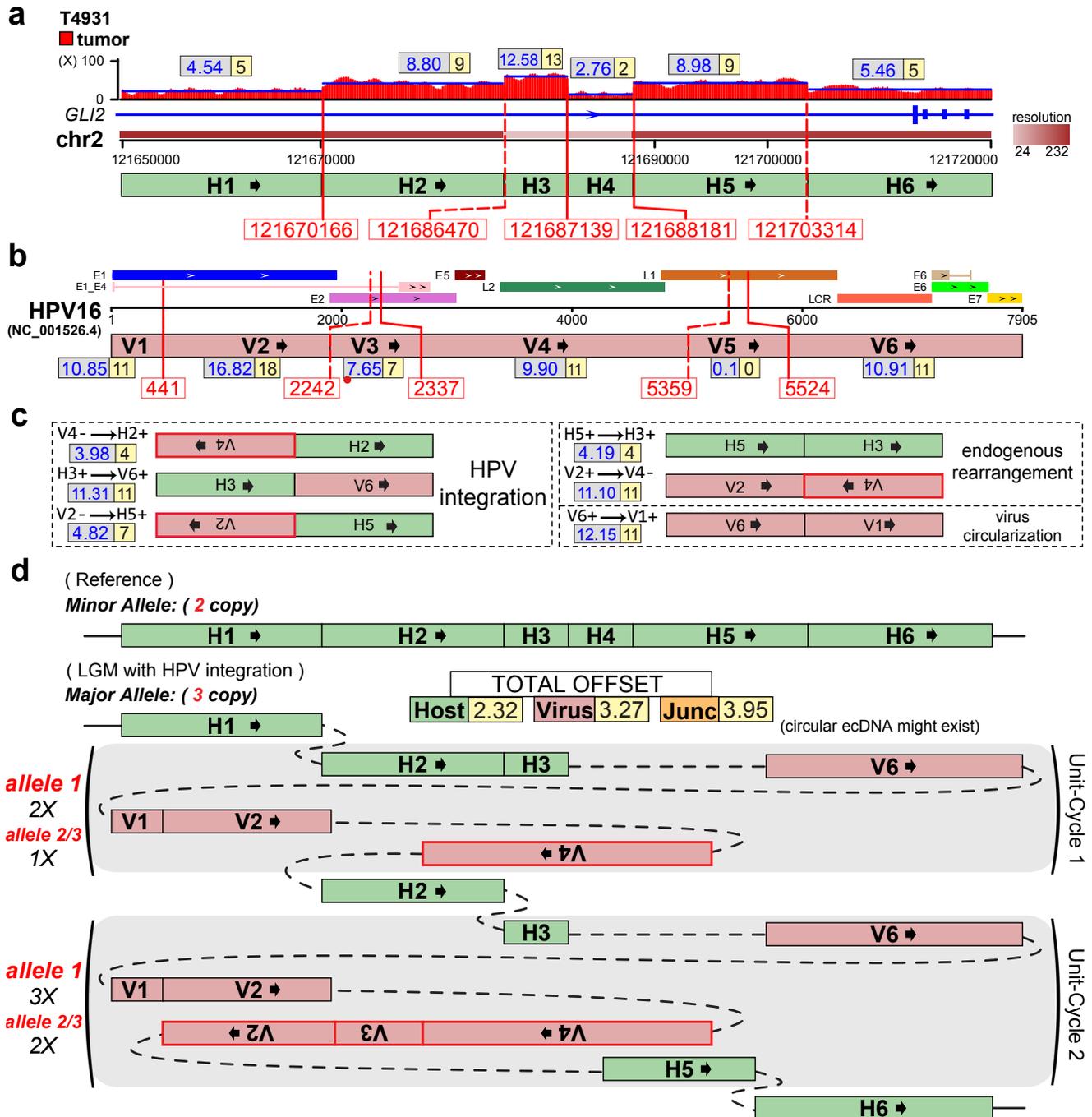
Supplementary Figures



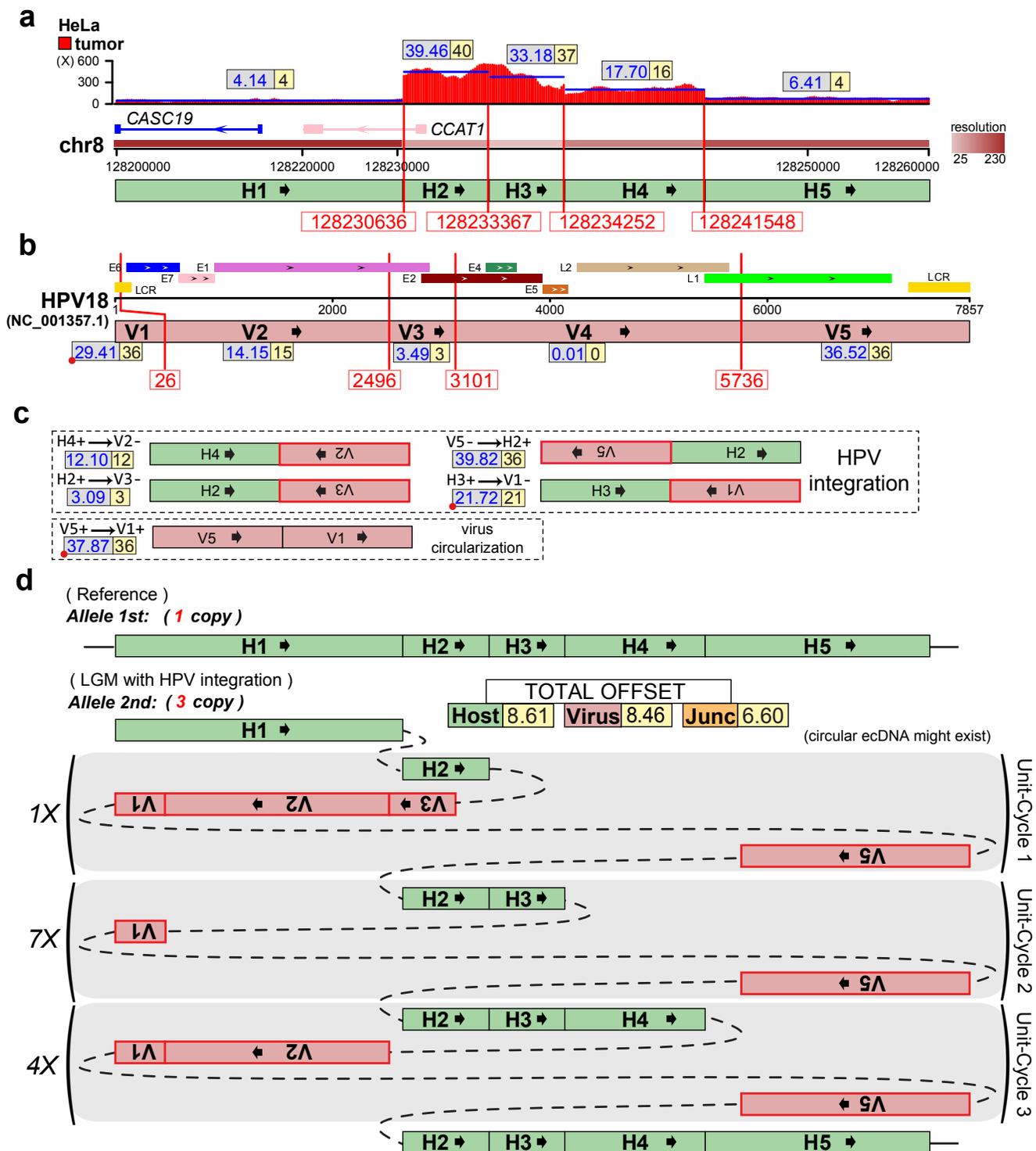
Supplementary Figure 1. The eight LGM modes in simulation work are denoted with features and represented by basic structures. The alteration sequences with host ('H') and virus ('V') segments are aligned with the reference. Segments in red frame are reverse complementary counterpart. The first and last viral segments are denoted by 'V_s' and 'V_e', respectively. Dashed frame and directed lines means segments and paths randomly selected.



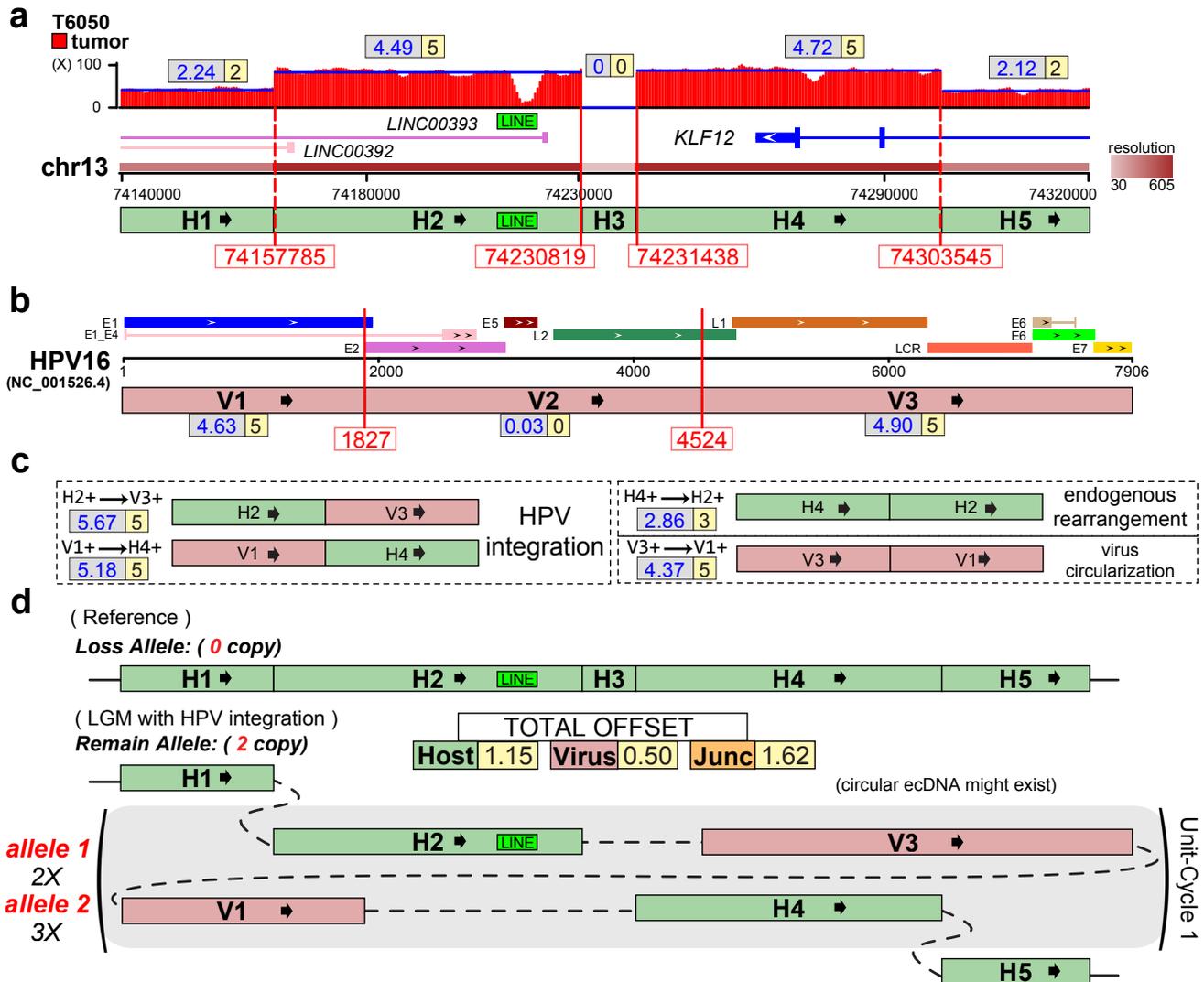
Supplementary Figure 2. Copy number distribution along the HPV-integrated chromosome in four cervical cancer samples⁷. Figures are from patchwork results. HPV integration sites are denoted by red arrows. Colored and grey bold lines indicate total and minor copy number respectively.



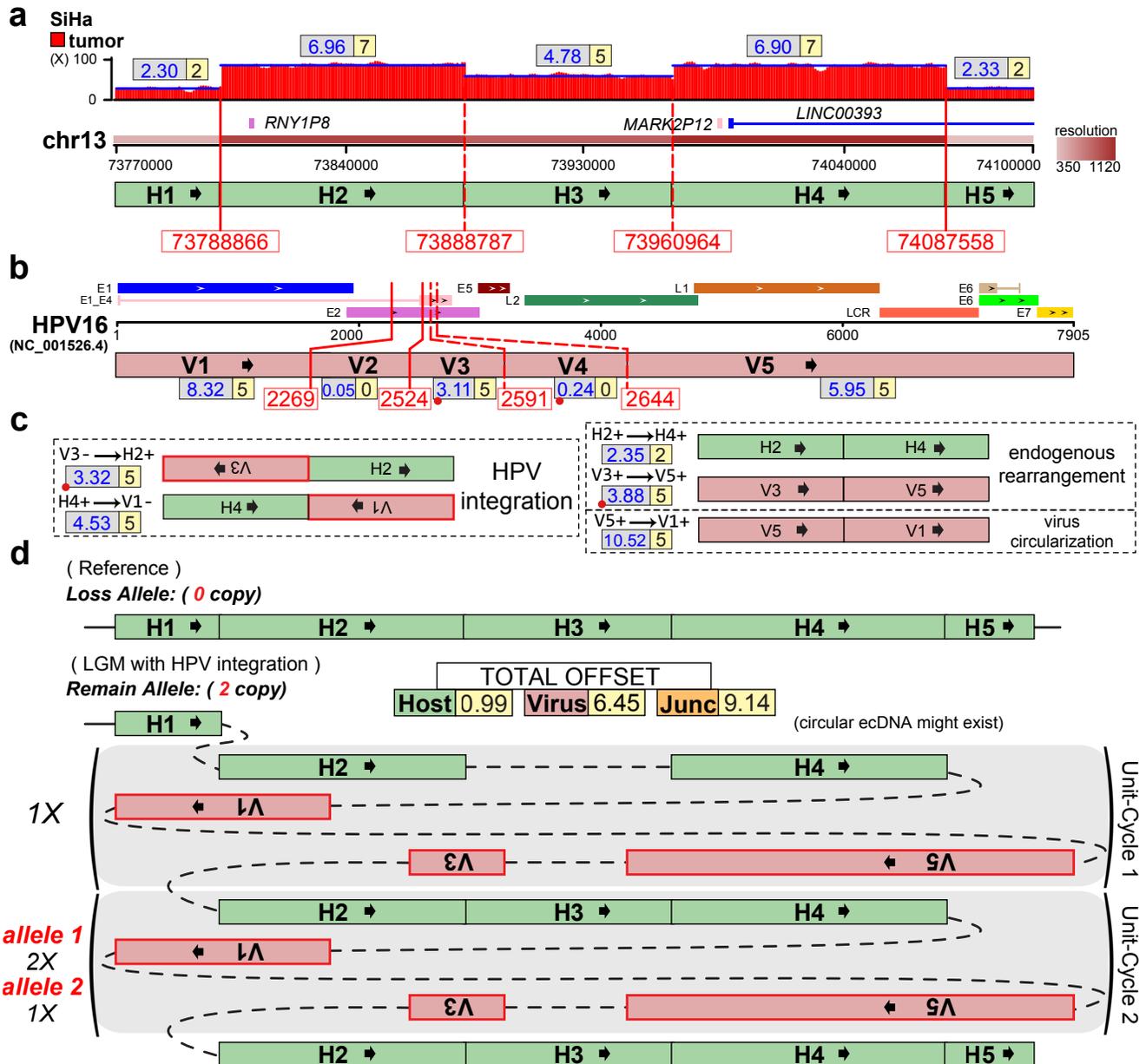
Supplementary Figure 3. Presentative LGM at HPV integration sites (gene *GLI2*) on chr2 (major allele) of the T4931 sample⁷. (a) Human genomic region flanking HPV integrations are divided into six segments (H1~H6) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V6) of HPV16 genome by VITs and SVs. The segment (V3) less than 100bp is marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the 'Simplest LGM' are indicated as string of coloured segments with copy times, including reference allele (minor allele) and that harbours HPV16 integrations (major allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shadowed areas) in LGM are denoted with repeat time. Note that the HPV-integrated alleles might have different copies of unit-cycles, which might exist as circular ecDNAs.



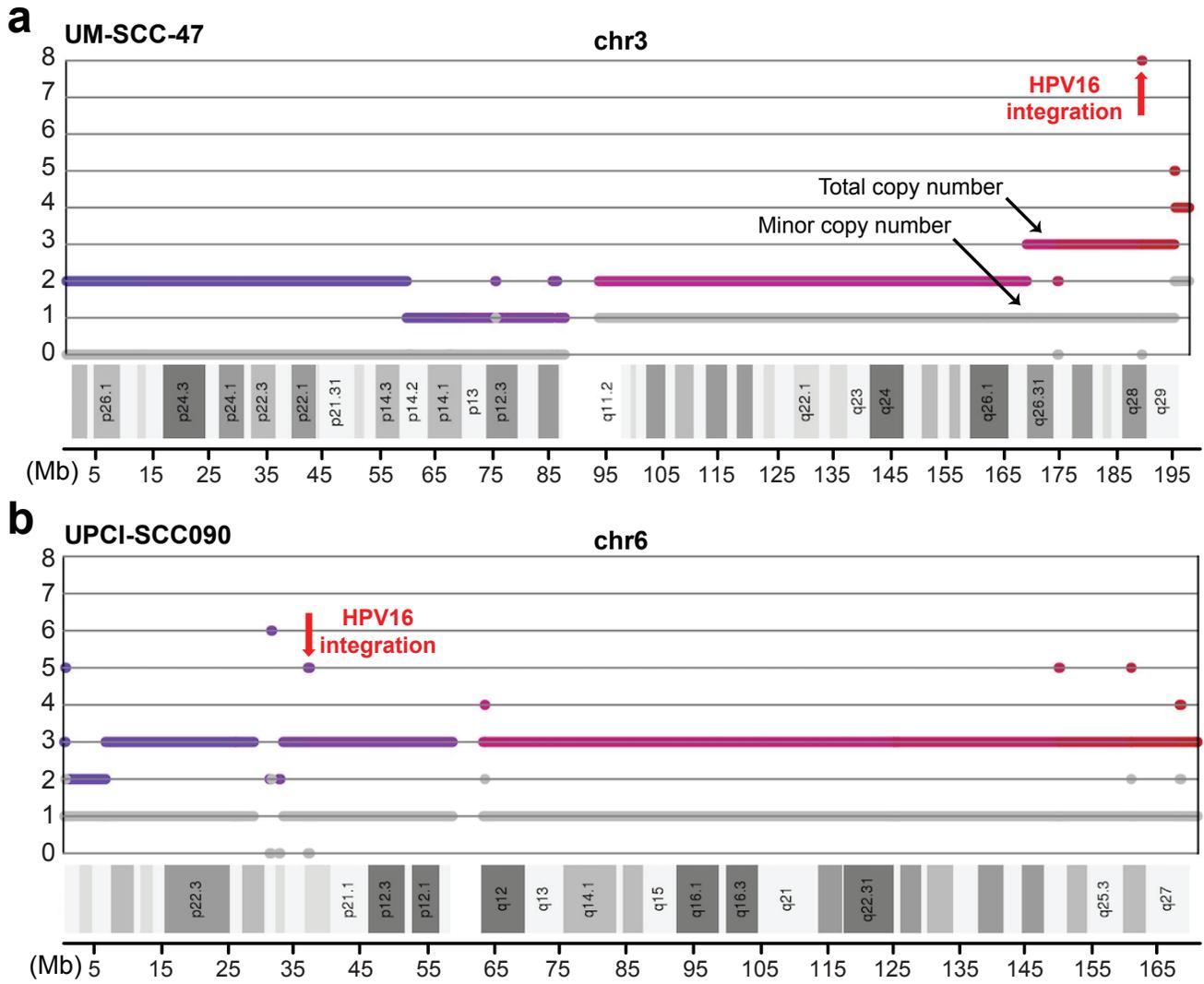
Supplementary Figure 4. Presentative LGM at HPV integration sites (upstream of gene *MYC*) on chr8 (major allele) of the HeLa cell line⁷. (a) Human genomic region flanking HPV integrations are divided into five segments (H1~H5) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V5) of HPV18 genome by VITs. The segment (V1) less than 100bp is marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the 'Simplest LGM' are indicated as string of coloured segments with copy times, including reference allele (minor allele) and that harbours HPV18 integrations (major allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. The unit-cycles might exist in the form of circular ecDNA.



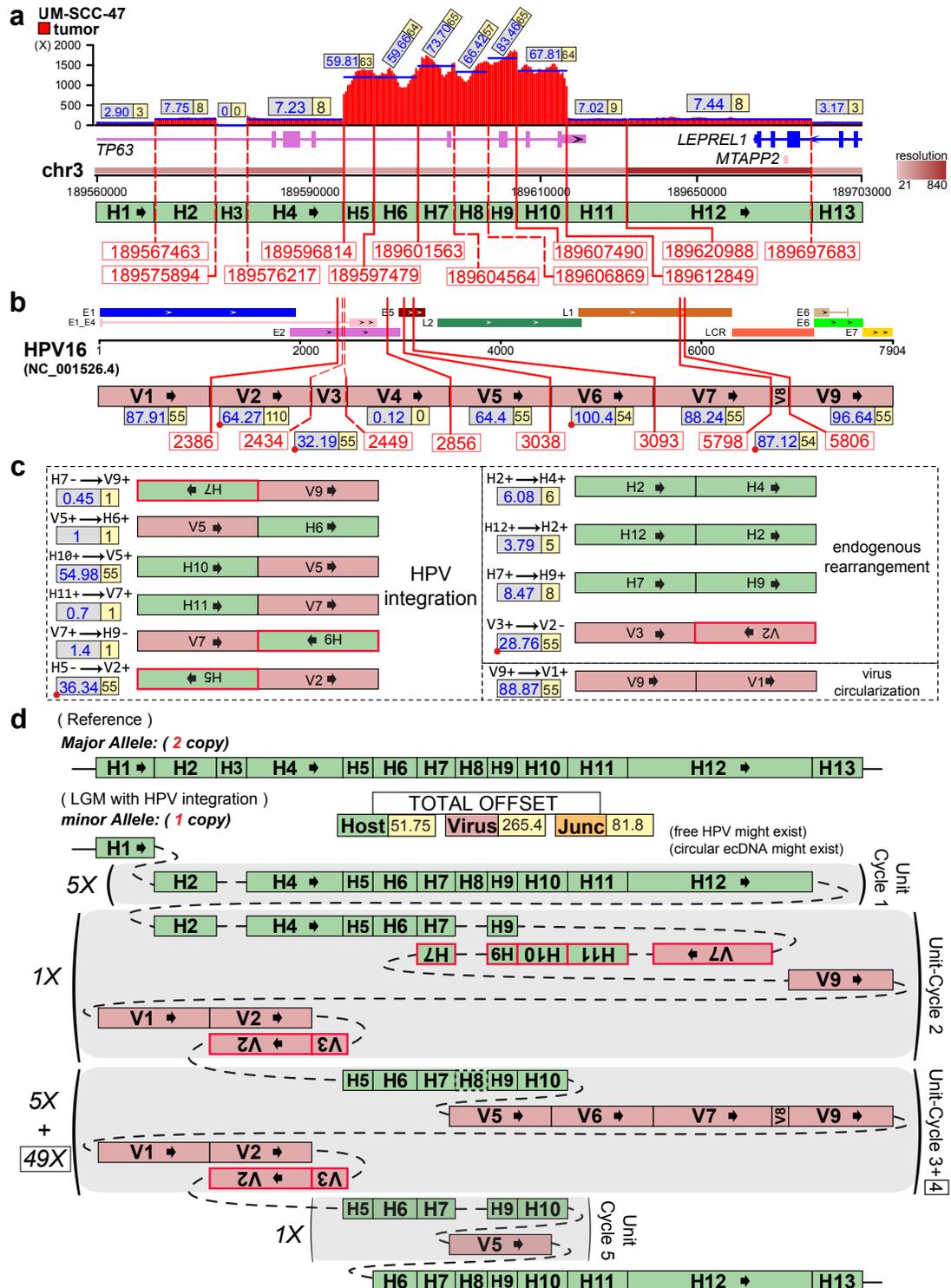
Supplementary Figure 5. Presentative LGM at HPV integration sites (gene *KLF12*) on chr13 (LOH) of the T6050 sample⁷. (a) Human genomic region flanking HPV integrations are divided into five segments (H1~H5) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V3) of HPV16 genome by VITs and SVs. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (loss allele) and that harbours HPV16 integrations (remain allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the HPV-integrated alleles might have different copies of unit-cycles, which might exist as circular ecDNAs.



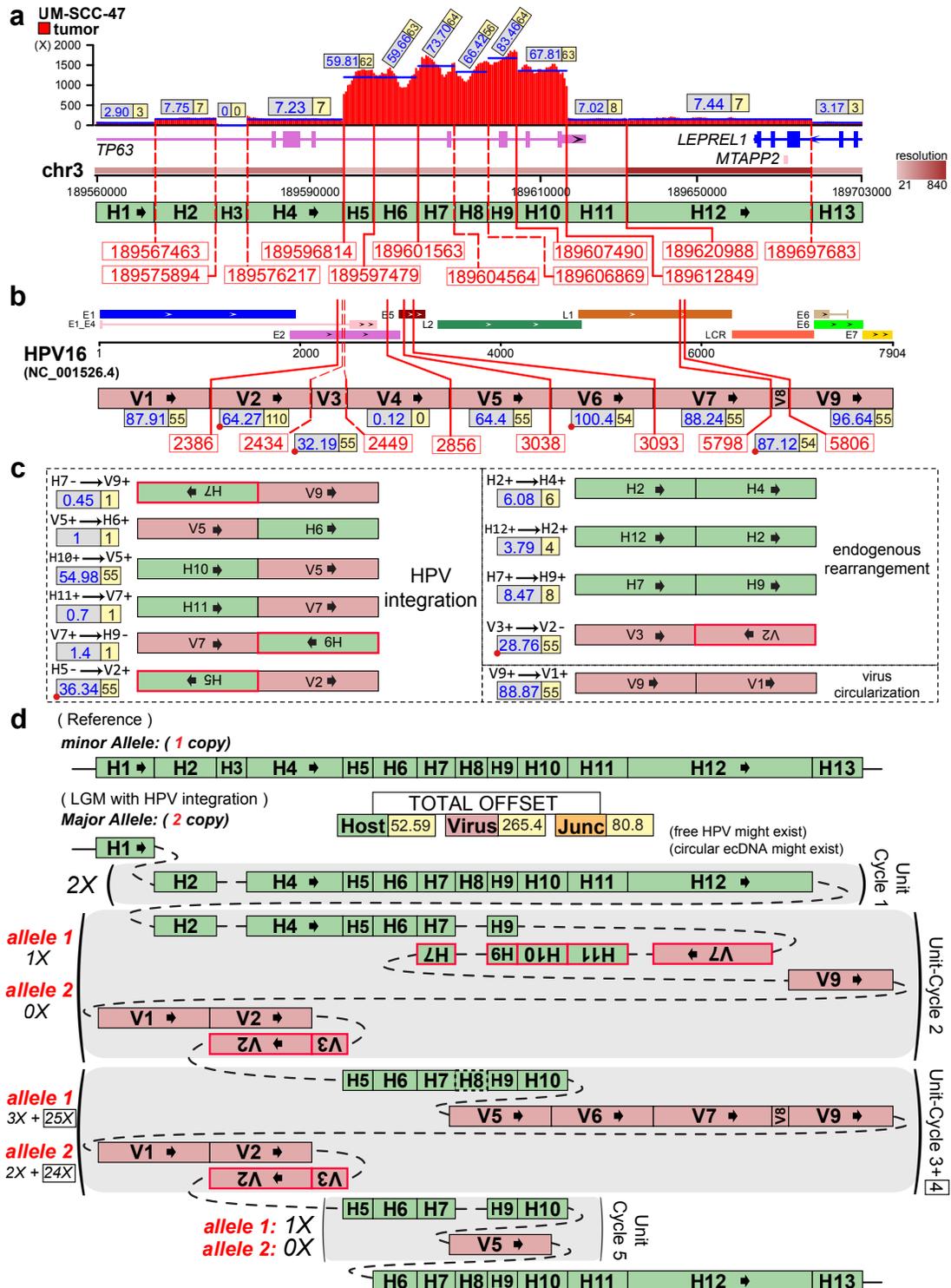
Supplementary Figure 6. Presentative LGM at HPV integration sites (upstream of gene *KLF12*) on chr13 (LOH) of the SiHa cell line⁷. (a) Human genomic region flanking HPV integrations are divided into five segments (H1~H5) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V5) of HPV16 genome by VITs and SVs. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (loss allele) and that harbours HPV16 integrations (remain allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the HPV-integrated alleles might have different copies of unit-cycles, which might exist as circular ecDNAs.



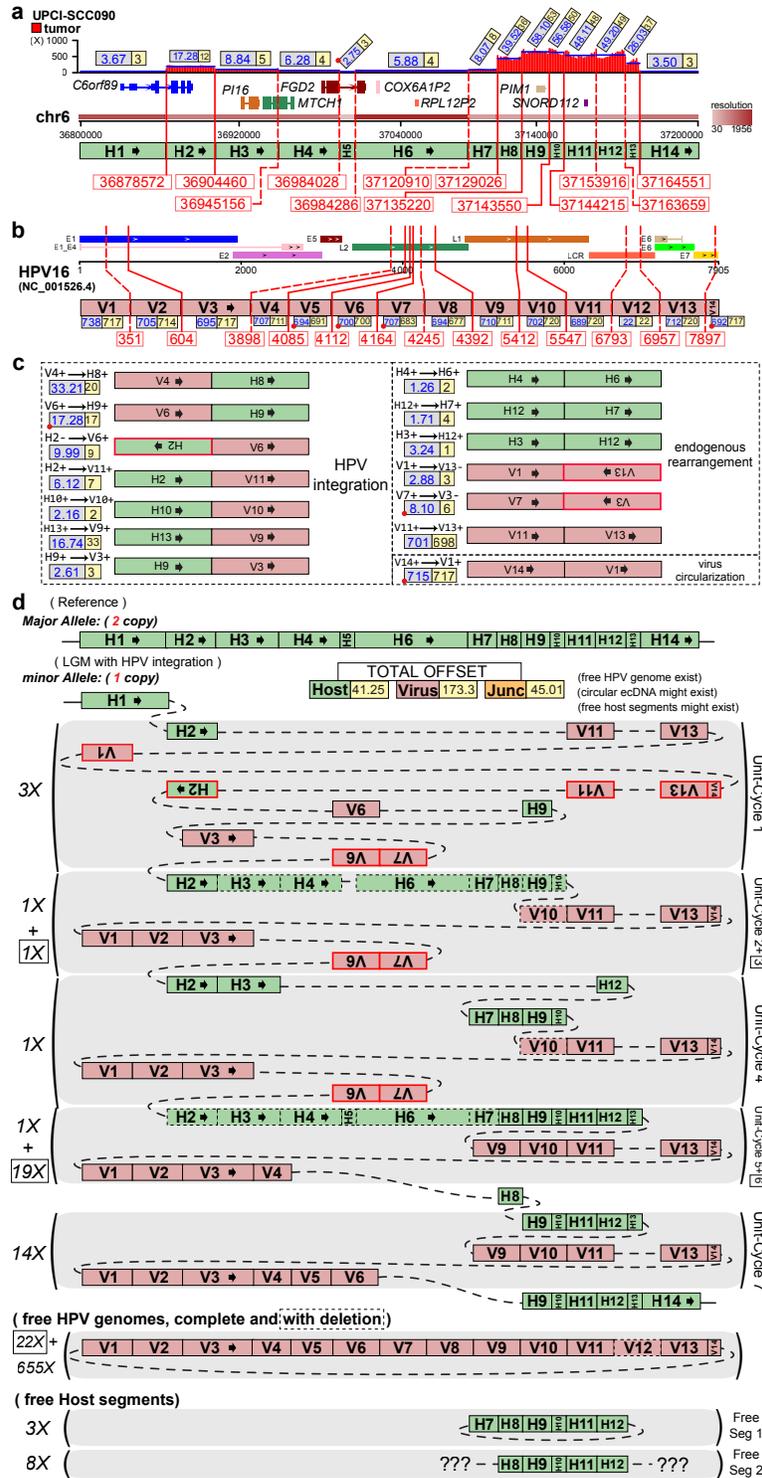
Supplementary Figure 7. Copy number distribution along the HPV-integrated chromosome in two HNSCC cell lines². Figures are from patchwork results. HPV integration sites are denoted by red arrows. Colored and grey bold lines indicate total and minor copy number respectively.



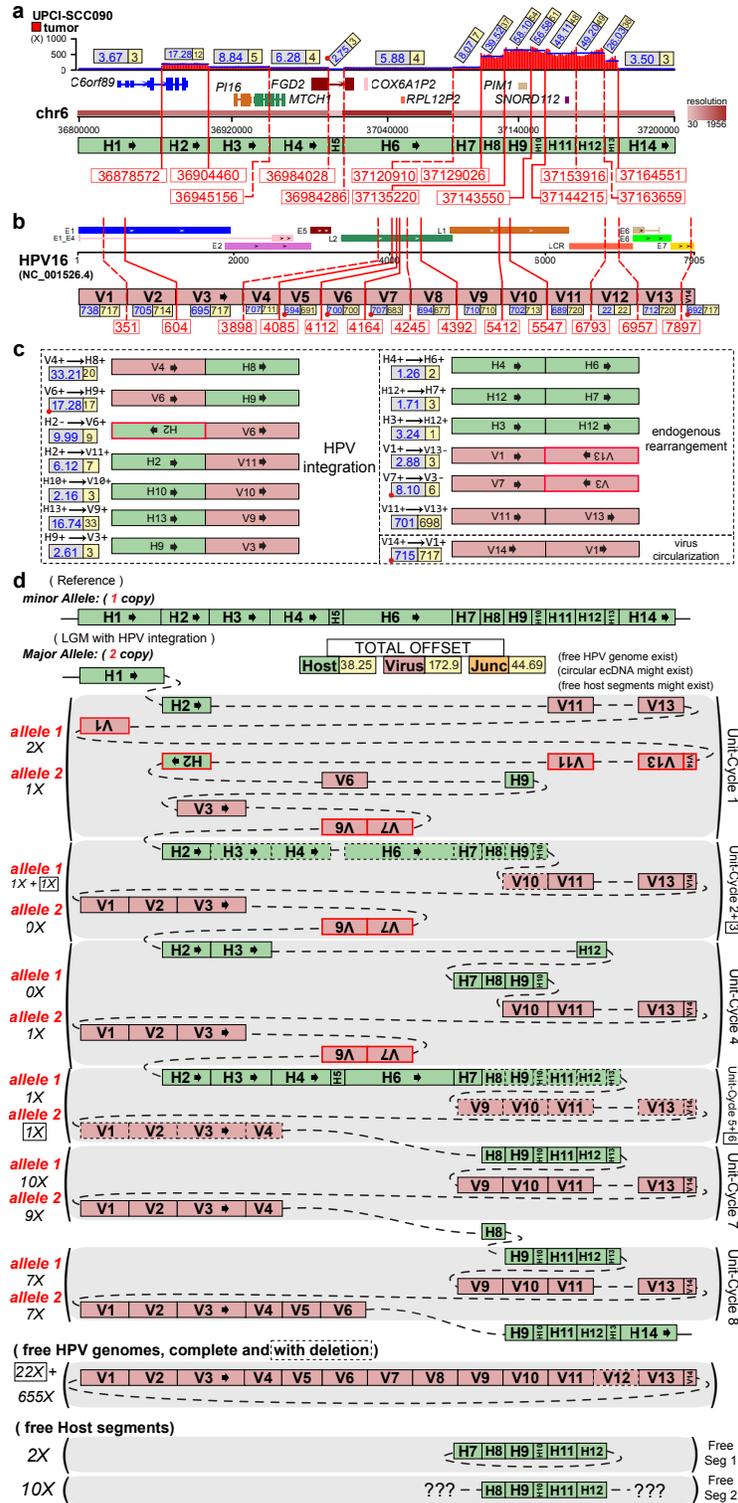
Supplementary Figure 8. Presentative LGM at HPV integration sites (gene *TP63*) on chr3 (minor allele) of the UM-SCC-47 cell line². **(a)** Human genomic region flanking HPV integrations are divided into 13 segments (H1~H13) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. **(b)** Segmentation (V1~V9) of HPV16 genome by VITs and SVs. The segments less than 100bp are marked with red dot. **(c)** Variant segment junctions utilized in Conjugate graph. **(d)** Resolved alleles of the 'Simplest LGM' are indicated as string of coloured segments with copy times, including reference allele (major allele) and that harbours HPV16 integrations (minor allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. The frame enclosed copy number is corresponding to unit-cycle (the NO.4) only containing the segments in solid frame. The unit-cycles might exist in the form of circular ecDNA.



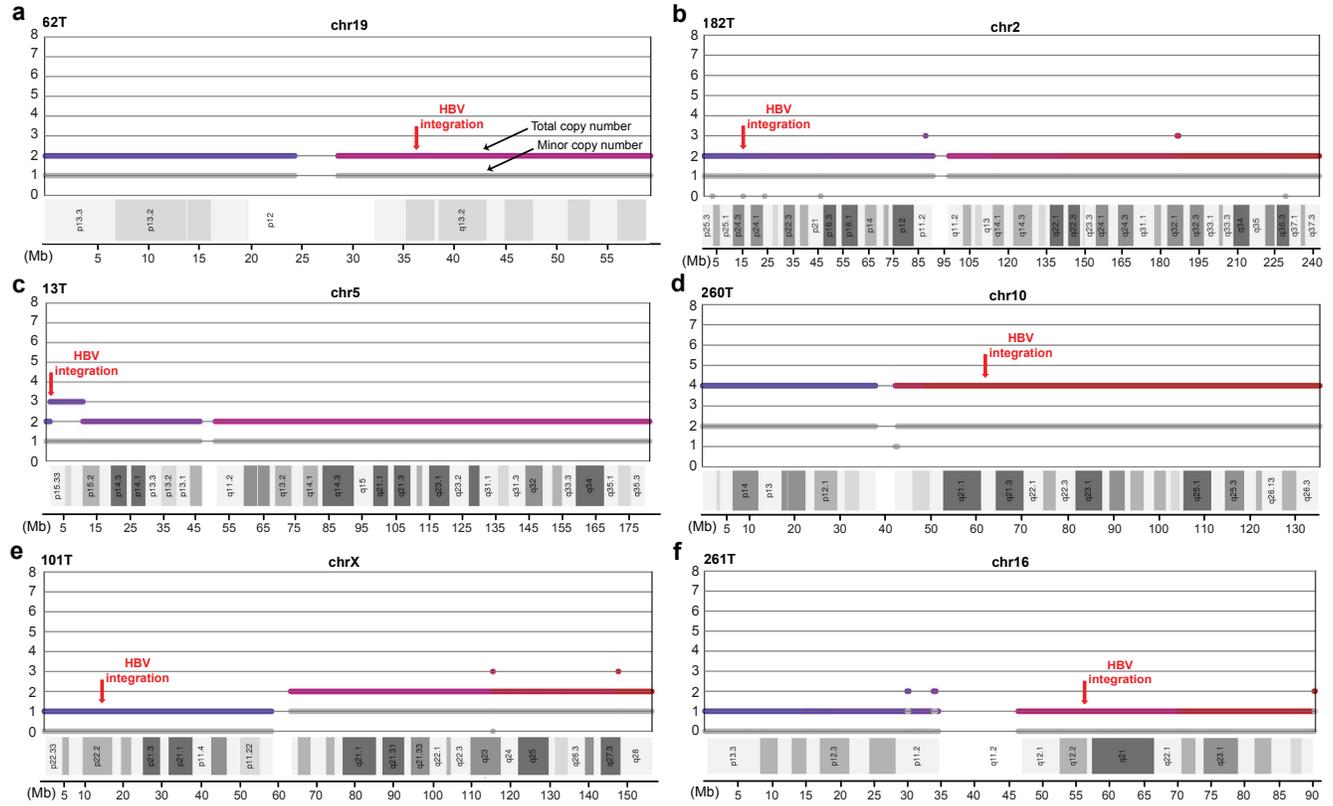
Supplementary Figure 9. Presentative LGM at HPV integration sites (gene *TP63*) on chr3 (major allele) of the UM-SCC-47 cell line². (a) Human genomic region flanking HPV integrations are divided into 13 segments (H1~H13) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V9) of HPV16 genome by VITs and SVs. The segments less than 100bp are marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (minor allele) and that harbours HPV16 integrations (major allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. The frame enclosed copy number is corresponding to unit-cycle (the NO.4) only containing the segments in solid frame. Note that the HPV-integrated alleles might have different copies of unit-cycles, which might exist as circular ecDNAs.



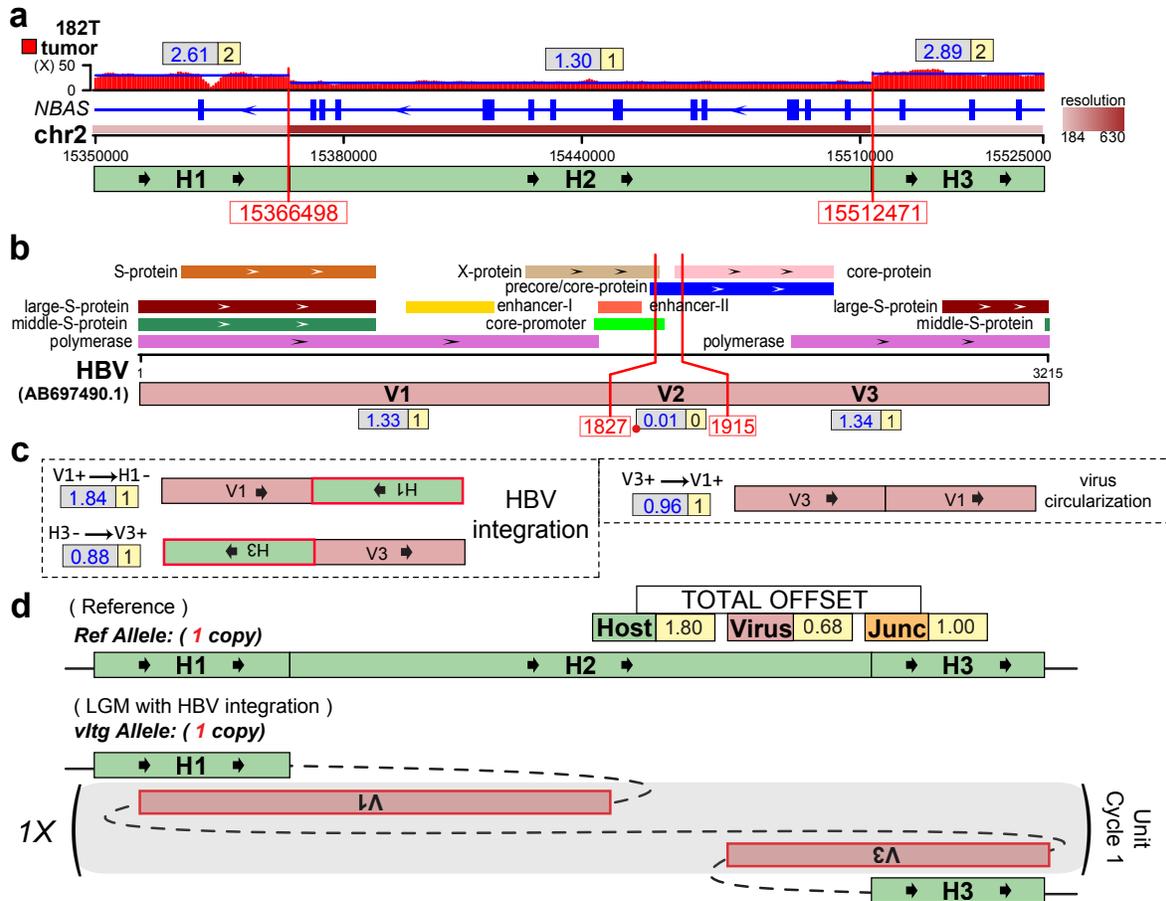
Supplementary Figure 10. Presentative LGM at HPV integration sites (gene *PIM1*) on chr6 (minor allele) on chr6 of the UPCI-SCC090 cell line². (a) Human genomic region flanking HPV integrations are divided into 14 segments (H1~H14) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V14) of HPV16 genome by VITs and SVs. The segments less than 100bp are marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (major allele) and that harbours HPV16 integrations (minor allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that both free HPV genomes and host segments exist with circular and linear structures. The frame enclosed copy number is corresponding to the unit-cycle only containing the segments in solid frame.



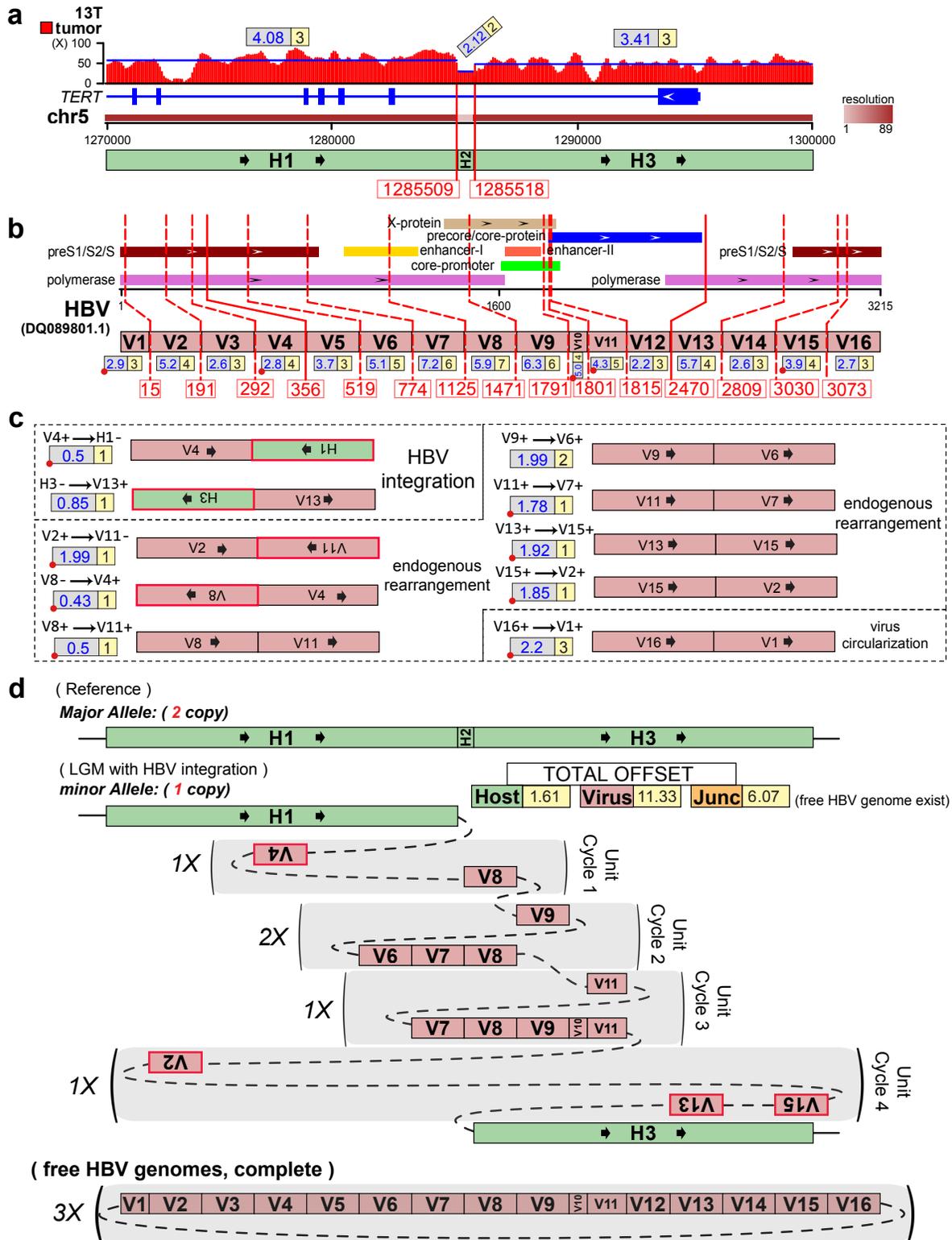
Supplementary Figure 11. Presentative LGM at HPV integration sites (gene *PIM1*) on chr6 (major allele) of the UPCI-SCC090 cell line². **(a)** Human genomic region flanking HPV integrations are divided into 14 segments (H1~H14) by VITs (red solid-line) and SVs (red dashed-line) denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. **(b)** Segmentation (V1~V14) of HPV16 genome by VITs and SVs. The segments less than 100bp are marked with red dot. **(c)** Variant segment junctions utilized in Conjugate graph. **(d)** Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (minor allele) and that harbours HPV16 integrations (major allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that both free HPV genomes and host segments exist with circular and linear structures. The frame enclosed copy number is corresponding to the unit-cycle only containing the segments in solid frame.



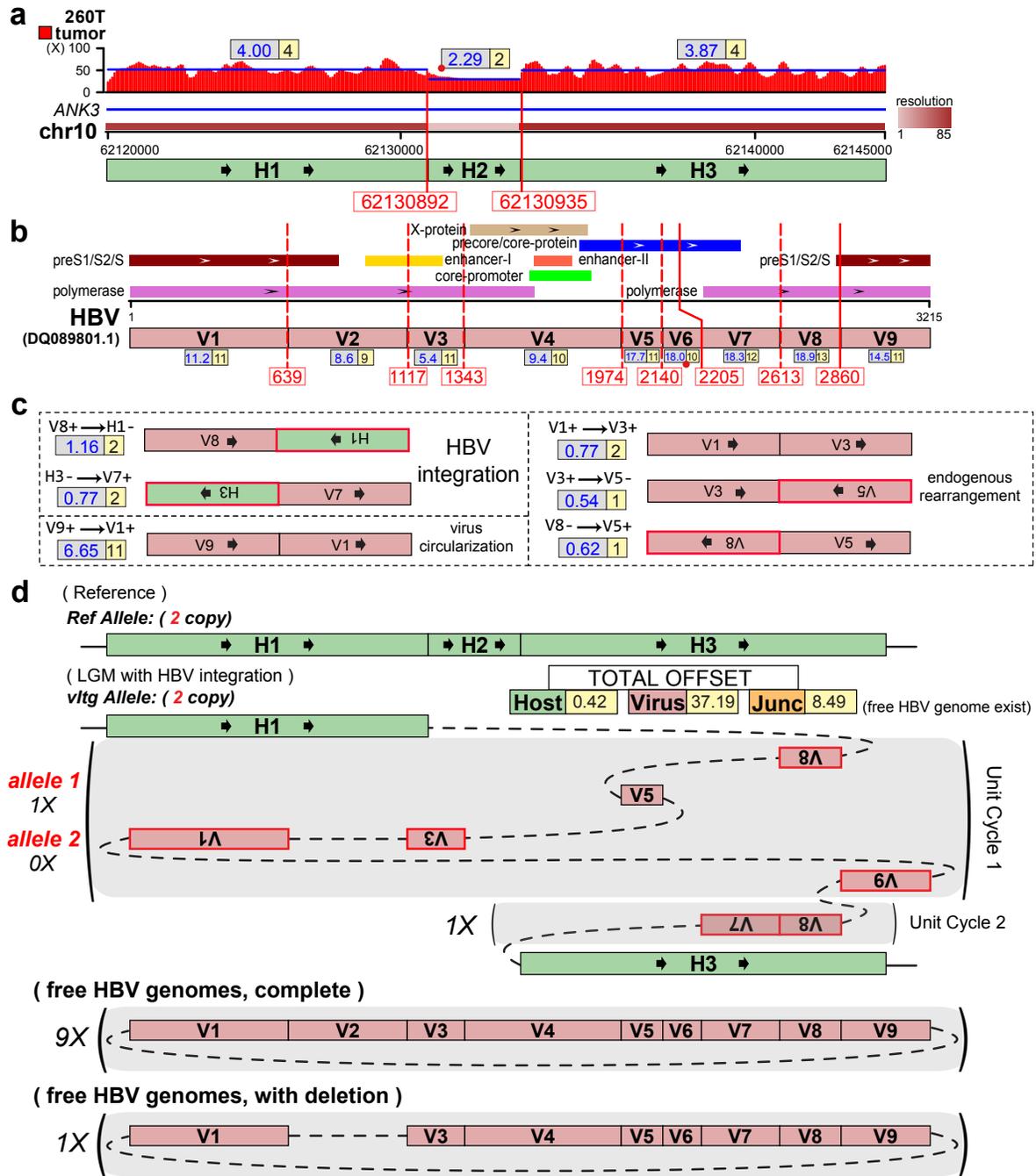
Supplementary Figure 12. Copy number distribution along the HBV-integrated chromosome in six HCC samples¹¹. Figures are from patchwork results. HBV integration sites are denoted by red arrows. Colored and grey bold lines indicate total and minor copy number respectively.



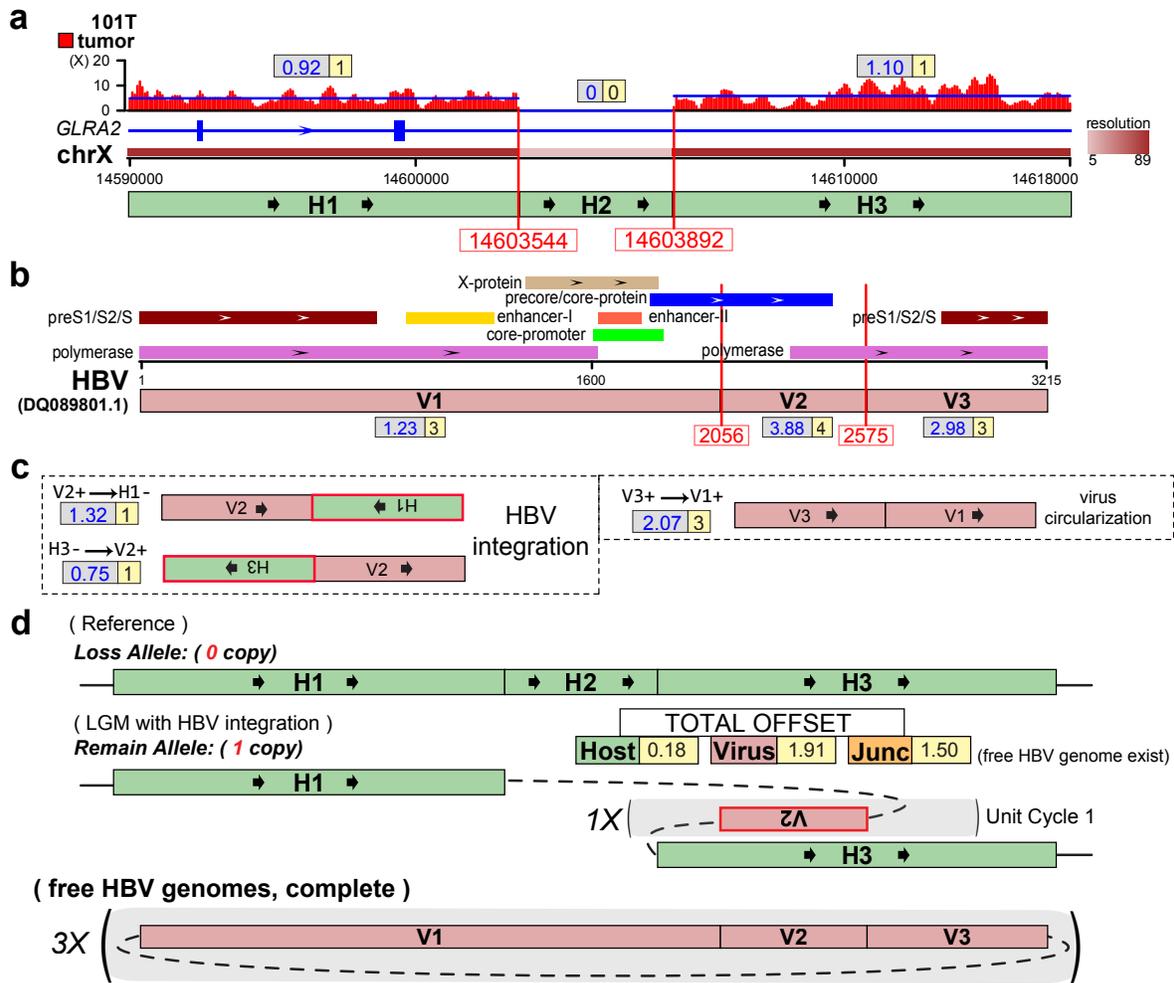
Supplementary Figure 13. Presentative LGM at HBV integration sites (gene *NBAS*) on chr2 of the 182T HCC sample¹¹. **(a)** Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. **(b)** Segmentation (V1~V3) of HBV genome by VITs and SVs. The segment less than 100bp is marked with red dot. **(c)** Variant segment junctions utilized in Conjugate graph. **(d)** Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele and that harbours HBV integrations. Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time.



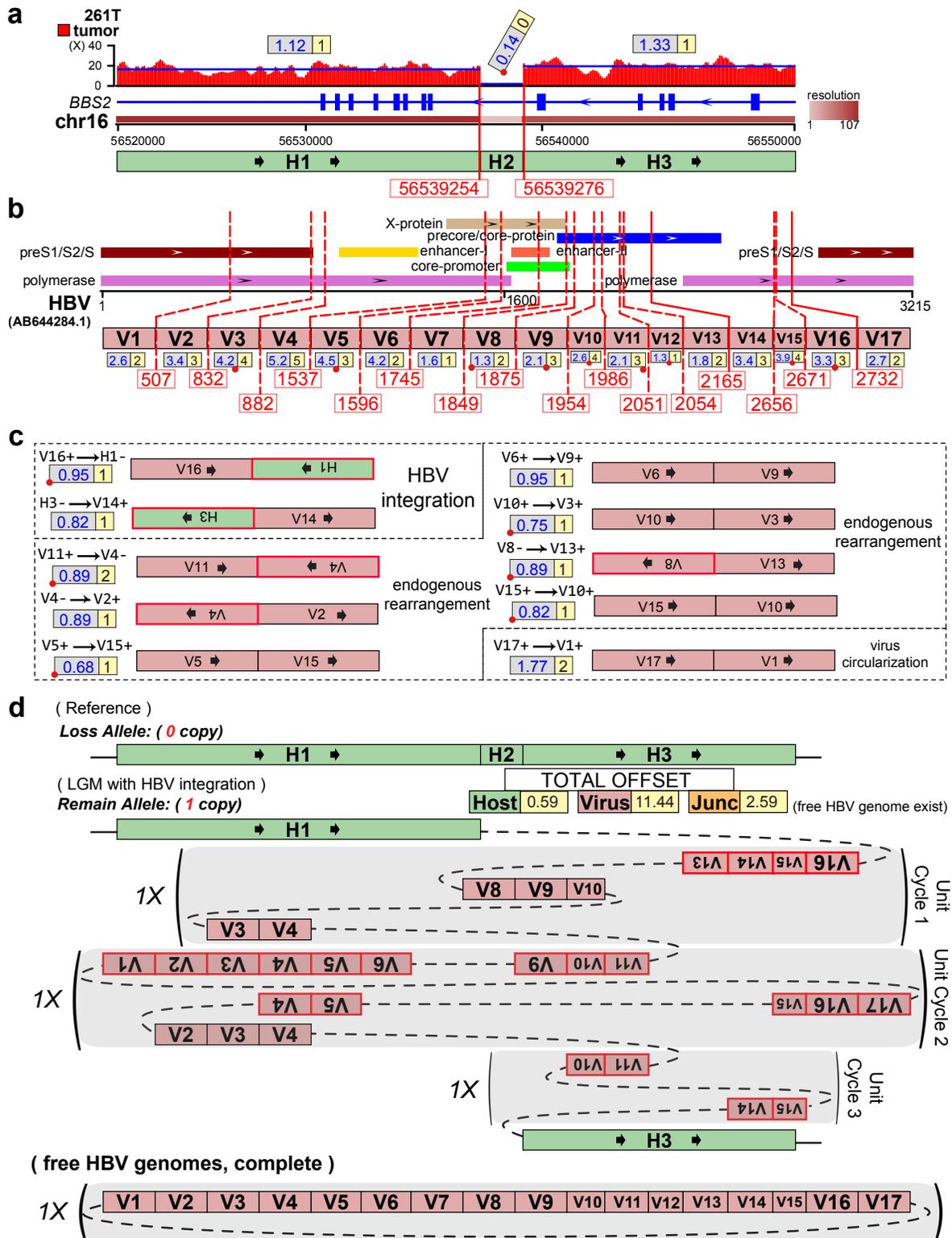
Supplementary Figure 14. Presentative LGM at HBV integration sites (gene *TERT*) on chr5 (minor allele) of the 13T HCC sample¹¹. (a) Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V16) of HBV genome by VITs and SVs. The segment less than 100bp is marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (major allele) and that harbours HBV integrations (minor allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shadowed areas) in LGM are denoted with repeat time. Note that the free HBV genomes might exist.



Supplementary Figure 15. Presentative LGM at HBV integration sites (gene *ANK3*) on chr10 of the 260T HCC sample¹¹. (a) Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V9) of HBV genome by VITs and SVs. The segment less than 100bp is marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele and that harbours HBV integrations. Sum of absolute offset of segments and junctions are shown. Unit-cycles (shadowed areas) in LGM are denoted with repeat time. Note that the HBV-integrated alleles might have different copies of unit-cycles, which might exist as circular ecDNA. The HBV free genome (both complete and with deletion) might exist.

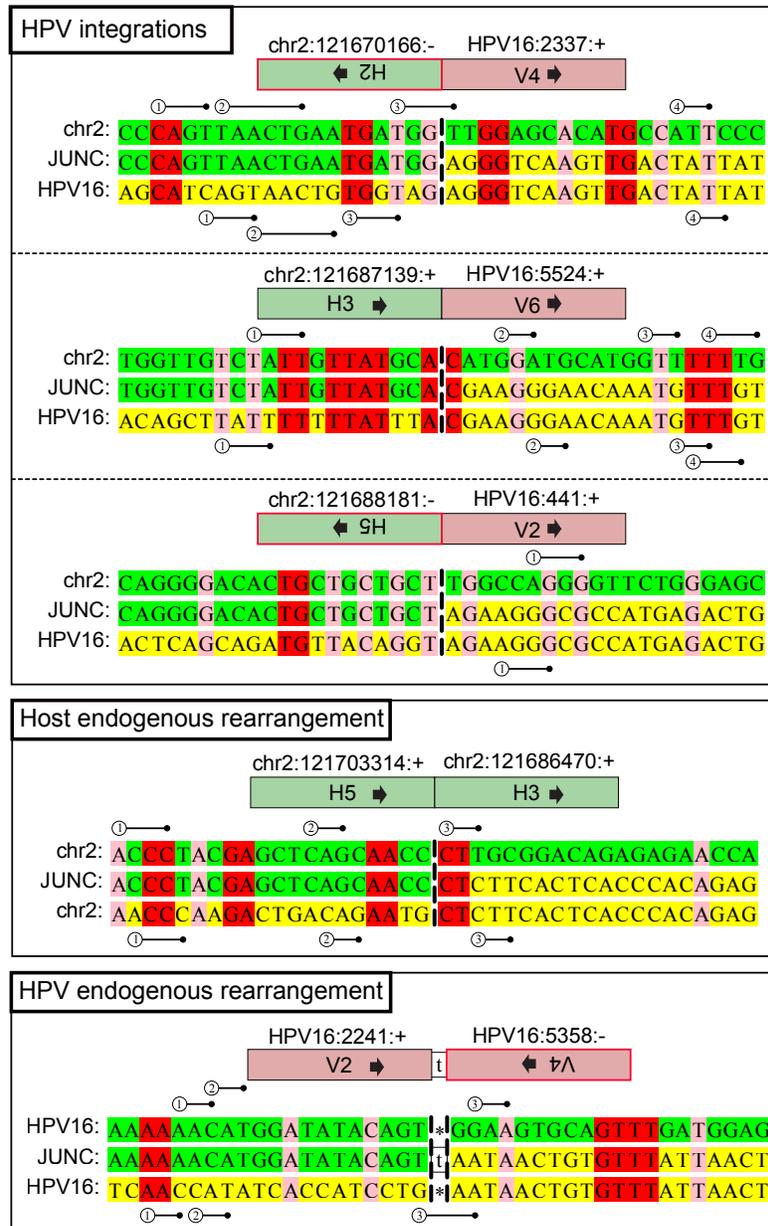


Supplementary Figure 16. Presentative LGM at HBV integration sites (gene *GLRA2*) on chrX of the 101T HCC sample¹¹. (a) Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V3) of HBV genome by VITs. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (loss allele) and that harbours HBV integrations (remain allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the free HBV genomes might exist.



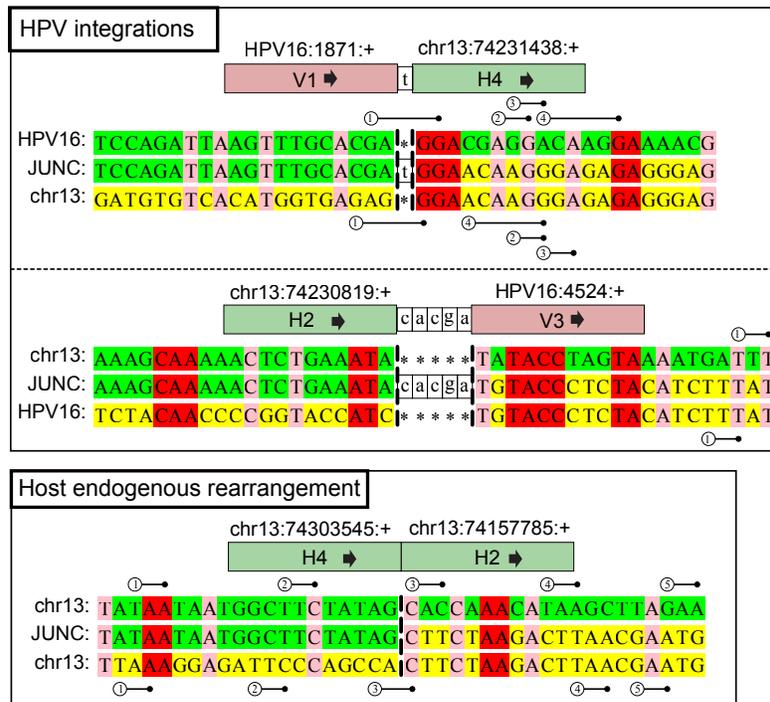
Supplementary Figure 17. Presentative LGM at HBV integration sites (gene *BBS2*) on chr16 of the 261T HCC sample¹¹. (a) Human genomic region flanking HBV integrations are divided into three segments (H1~H3) by VITs denoted with breakpoints. Depth spectrum is displayed with original (grey frame) and ILP-adjusted (yellow frame) copy numbers of segments. (b) Segmentation (V1~V17) of HBV genome by VITs and SVs. The segment less than 100bp is marked with red dot. (c) Variant segment junctions utilized in Conjugate graph. (d) Resolved alleles of the ‘Simplest LGM’ are indicated as string of coloured segments with copy times, including reference allele (loss allele) and that harbours HBV integrations (remain allele). Sum of absolute offset of segments and junctions are shown. Unit-cycles (shaded areas) in LGM are denoted with repeat time. Note that the free HBV genomes might exist.

SampleID: T4931



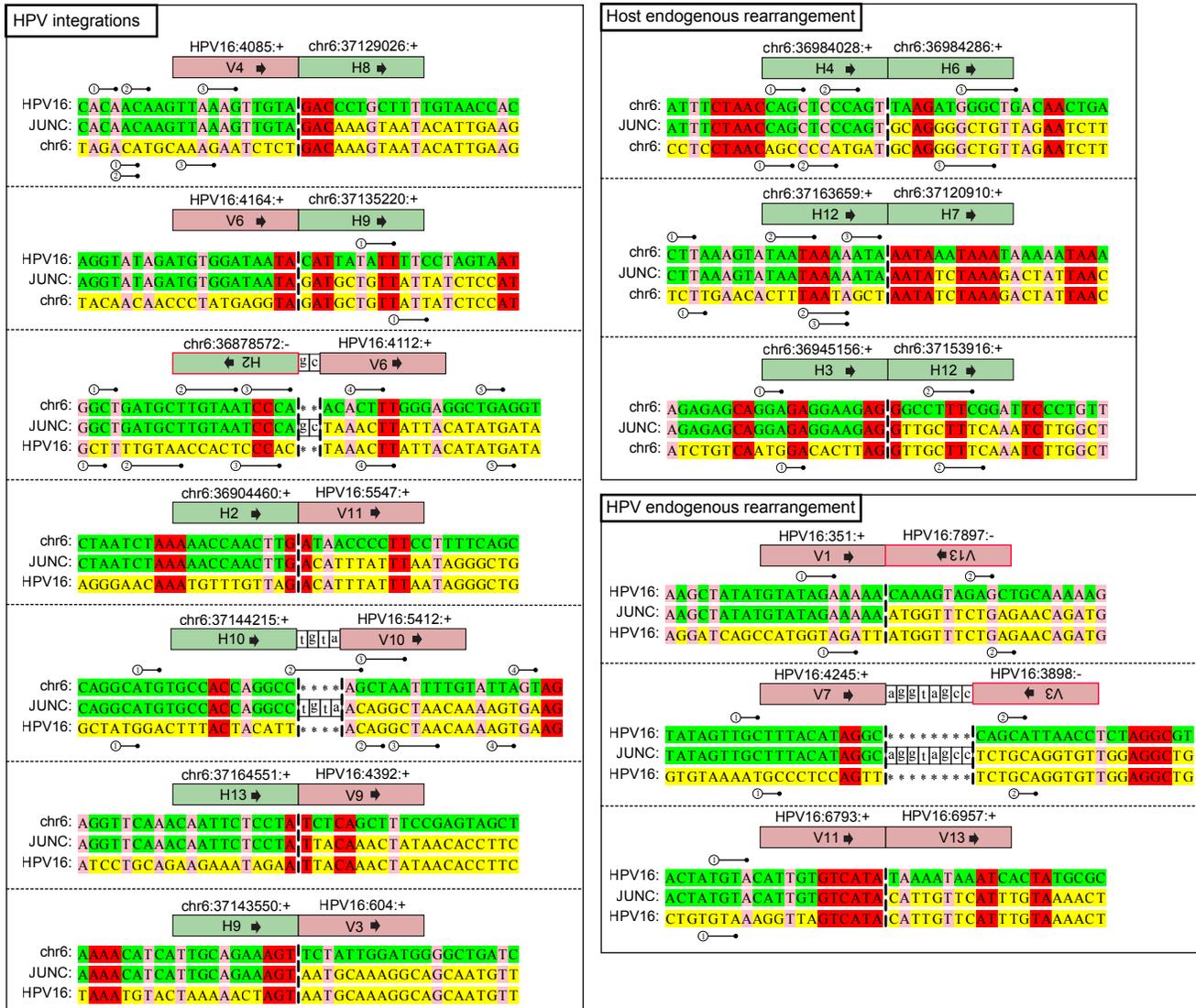
Supplementary Figure 19. Alignment of the sequence around the integration site between the human genome and the HPV16 genome, and the endogenous rearrangements on human genome and the HPV16 genome respectively, in the T4931 sample⁷. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologies. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Figure 2 and Supplementary Figure 3, Supplementary Table 3).

SampleID: **T6050**



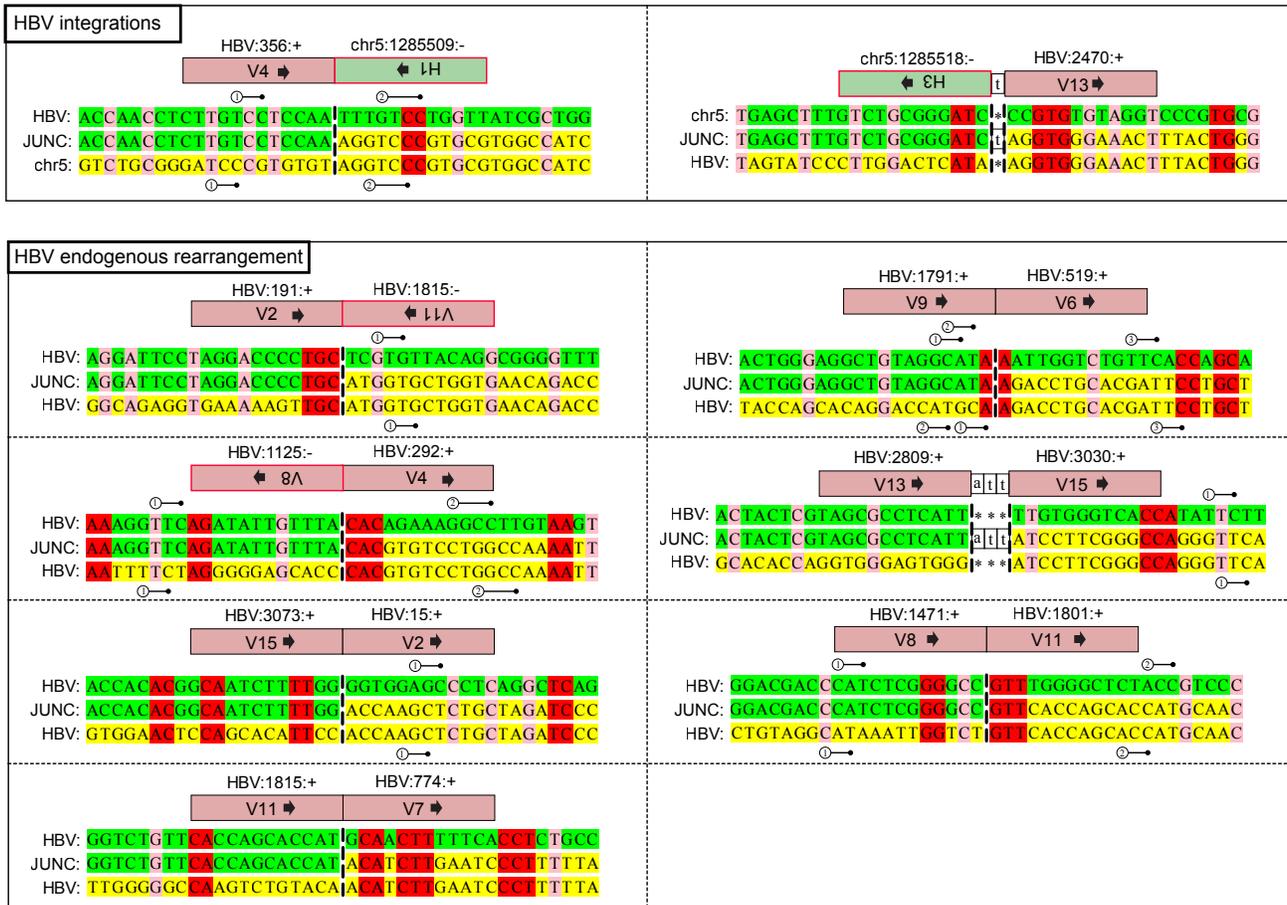
Supplementary Figure 20. Alignment of the sequence around the integration site between the human genome and the HPV16 genome, and the endogenous rearrangements on human genome, in the T6050 sample⁷. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microhomologous bases); numbered stick, slipped microhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome map (Supplementary Figure 5, Supplementary Table 3).

SampleID: UPCI-SCC090



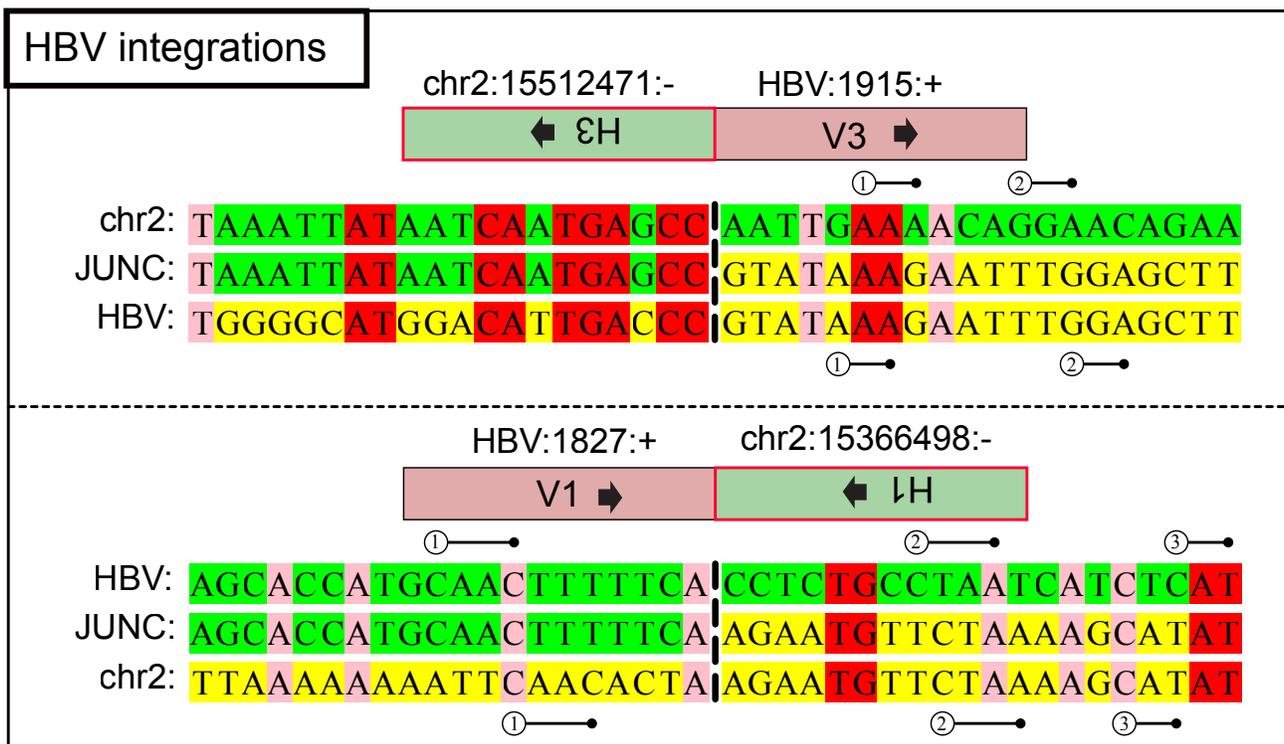
Supplementary Figure 22. Alignment of the sequence around the integration site between the human genome and the HPV16 genome, and the endogenous rearrangements on human genome and the HPV16 genome respectively, in the UPCI-SCC090 cell line². The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 10 and 11, Supplementary Table 4).

SampleID: 13T



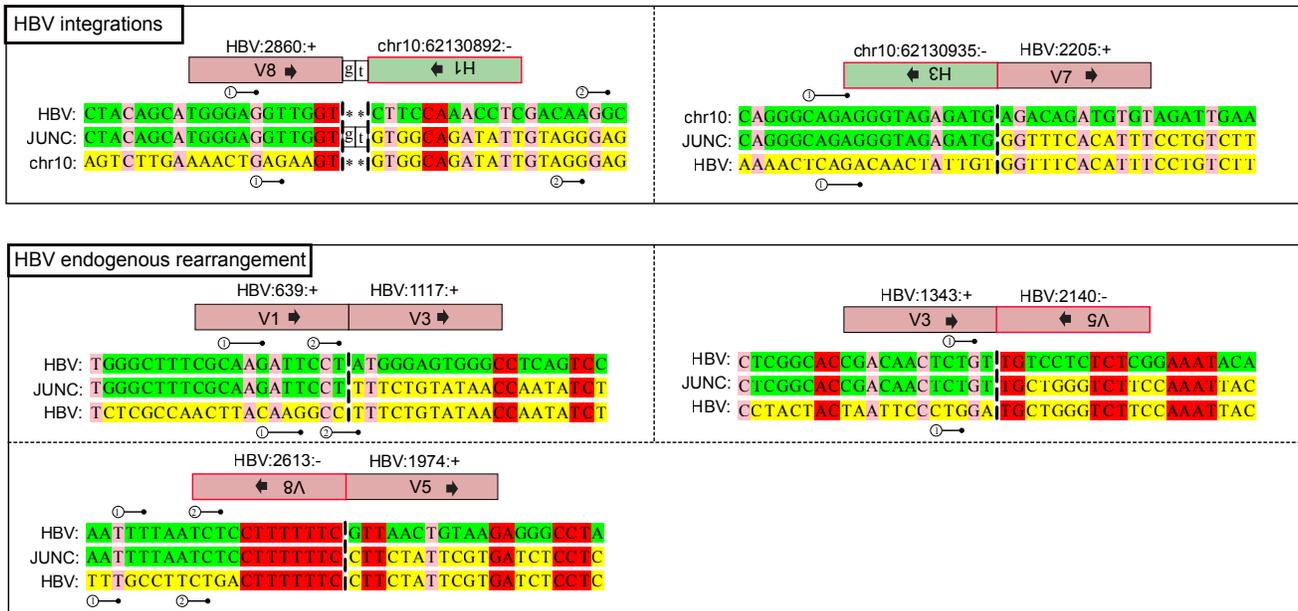
Supplementary Figure 23. Alignment of the sequence around the integration site between the human genome and the HBV genome, and the endogenous rearrangements on HBV genome, in the 13T HCC sample¹¹. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 14, Supplementary Table 6).

SampleID: 182T



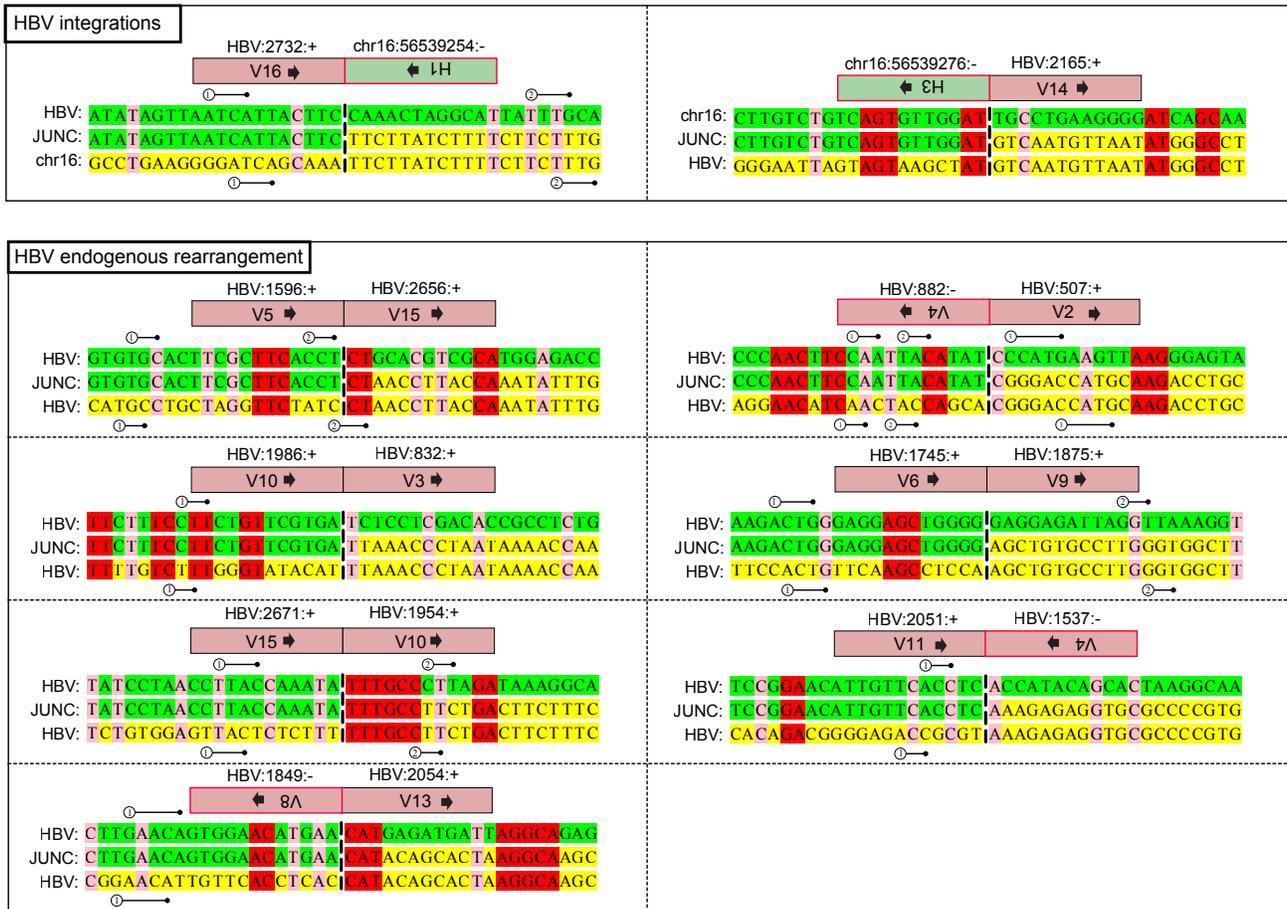
Supplementary Figure 24. Alignment of the sequence around the integration site between the human genome and the HBV genome in the 182T HCC sample¹¹. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 13, Supplementary Table 6).

SampleID: **260T**



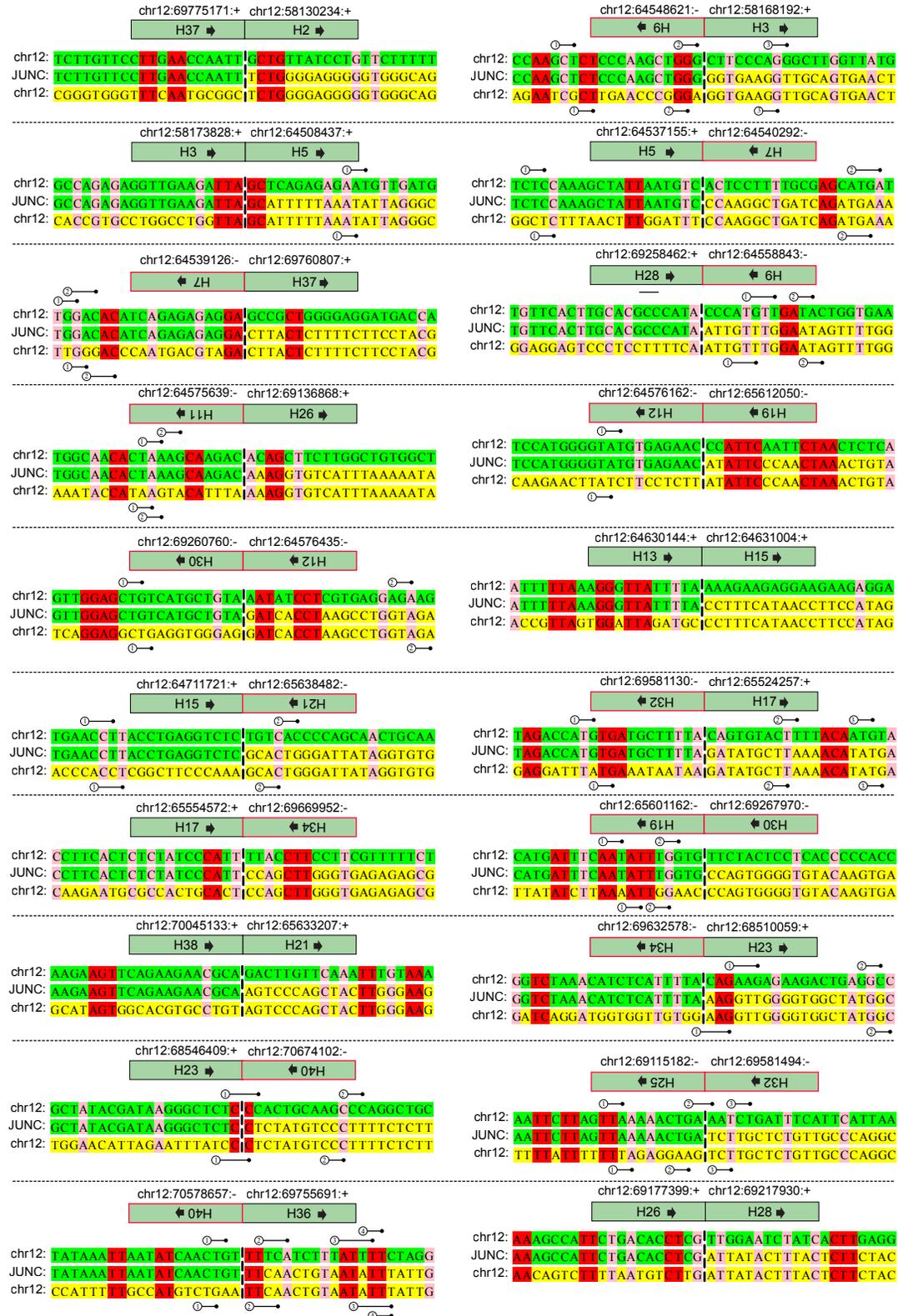
Supplementary Figure 25. Alignment of the sequence around the integration site between the human genome and the HBV genome, and the endogenous rearrangements on HBV genome, in the 260T HCC sample¹¹. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 15, Supplementary Table 6).

SampleID: 261T

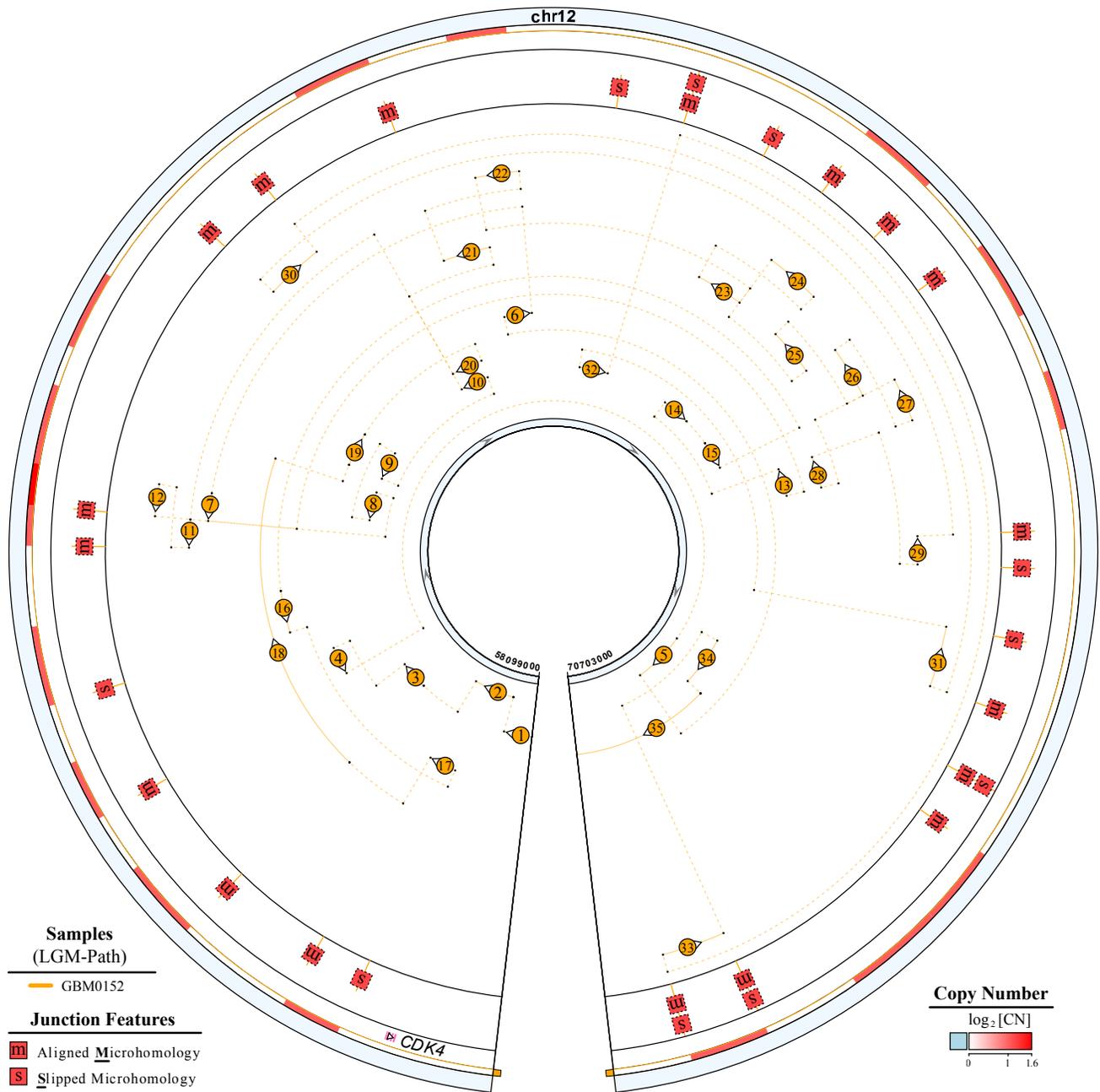


Supplementary Figure 26. Alignment of the sequence around the integration site between the human genome and the HBV genome, and the endogenous rearrangements on HBV genome, in the 261T HCC sample¹¹. The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned micorhomologous bases); numbered stick, slipped micorhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 17, Supplementary Table 6).

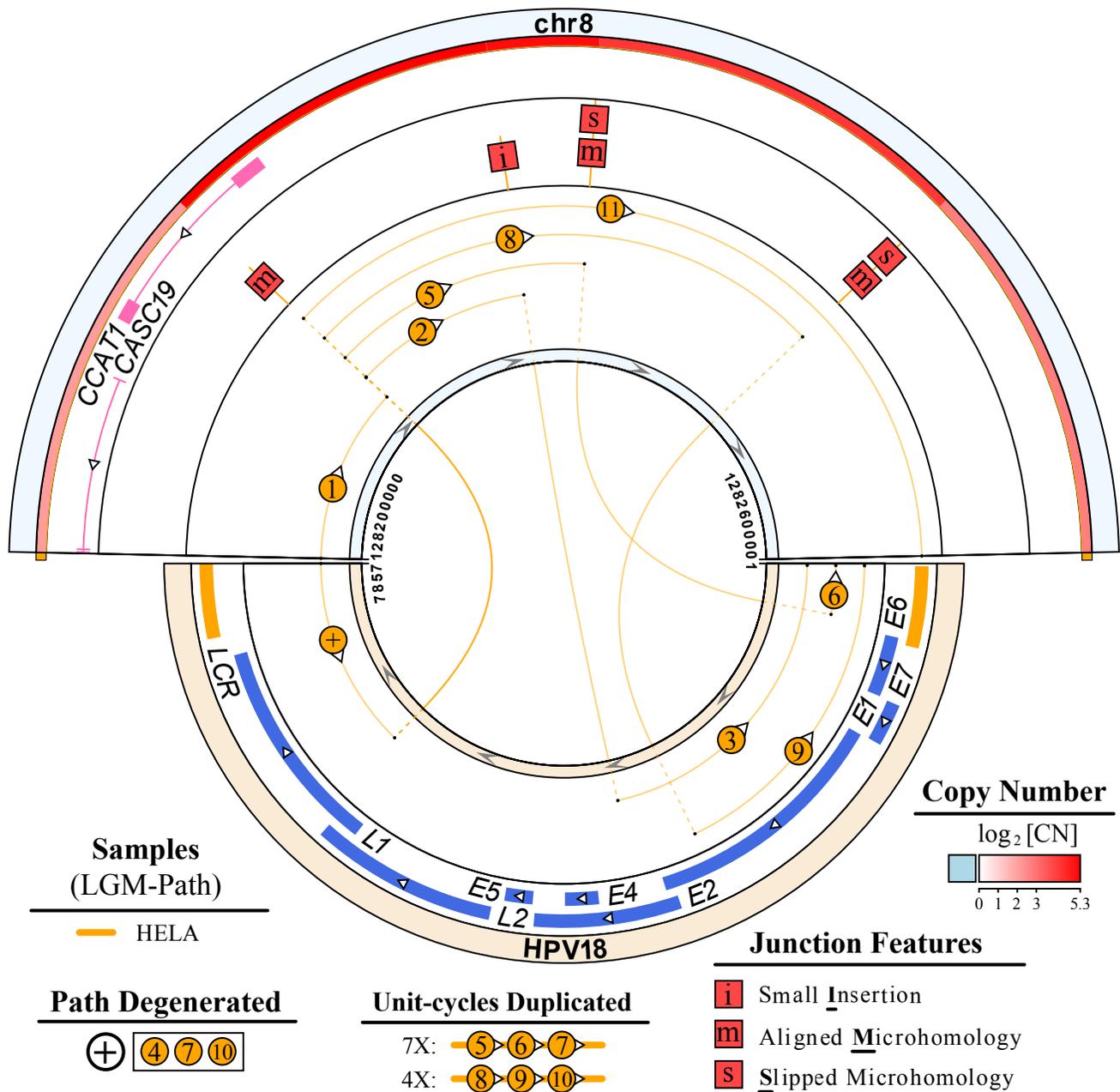
SampleID: **GBM0152**



Supplementary Figure 27. Alignment of the sequence around the endogenous complex rearrangements on human genome in the GBM0152 cancer sample²². The junction boundaries are shown as vertical dashed lines. All viral sequences are from the reference strand. Green, upstream partner; yellow, downstream partner; red, nucleotides that vertically align to both reference sequences (aligned microrhomologous bases); numbered stick, slipped microrhomologous bases. The junction segment IDs are corresponding to the segments in the resolved local genome maps (Supplementary Figure 18, Supplementary Table 7).



Supplementary Figure 28. Features of the complex SV LGM in GBM0152 sample. Human genomic segments related to complex SVs LGMs are shown as sectors with their relevant LGM path. Segments of LGMs are denoted by circled numbers in sequence. The numbered arcs might represent multiple sequential segments (Supplementary Table 9). Features of rearrangement sites are depicted as single-letter icons. DNA copy number (CN) is displayed in gradient red colour.



Supplementary Figure 29. Features of HPV18 integrated LGM in HeLa cell line. Human genomic segments are shown as sectors with the LGM path. Segments of LGMs are denoted by circled numbers in sequence, where some are degenerated by symbols for simplification. The numbered arcs might represent multiple sequential segments (Supplementary Table 9). Repeat times of unit-cycles in LGMs are stated in figure legend. Features of HPV18 integration sites are depicted as single-letter icons. DNA copy number (CN) is displayed in gradient red colour. The HPV18 genome reference is NC_001357.1 from the NCBI Nucleotide database.