

1           **The chromatin insulator CTCF regulates HPV18 transcript splicing and**  
2           **differentiation-dependent late gene expression**

3  
4 Jack Ferguson<sup>a^</sup>, Karen Campos Leon<sup>a^</sup>, leisha Pentland<sup>a</sup>, Joanne Stockton<sup>b</sup>,  
5 Thomas Günther<sup>c</sup>, Andrew Beggs<sup>a,b</sup>, Sally Roberts<sup>a</sup>, Boris Noyvert<sup>a,d,†</sup>, Joanna L.  
6 Parish<sup>a,†,#</sup>

7  
8 <sup>a</sup>Institute of Cancer and Genomic Sciences, College of Medical and Dental Sciences,  
9 University of Birmingham, B15 2TT, UK

10 <sup>b</sup>Genomics Birmingham, University of Birmingham, B15 2TT, UK

11 <sup>c</sup>Heinrich-Pette Institute, Leibniz Institute for Experimental Virology, Hamburg,  
12 Germany

13 <sup>d</sup>CRUK Birmingham Centre and Centre for Computational Biology, University of  
14 Birmingham, B15 2TT, UK

15  
16 <sup>^</sup>Equal Contribution

17 <sup>†</sup>Joint Senior Authors

18  
19 <sup>#</sup>Corresponding Author:

20 Joanna L. Parish

21 Email: [j.l.p parish@bham.ac.uk](mailto:j.l.p parish@bham.ac.uk)

22  
23 Running Title: Dynamic regulation of HPV18 transcription by CTCF

24 Abstract word count: 246

25 Main text word count: 5529

26

27

28

29

30

31

32

33

34

35 **ABSTRACT**

36 The ubiquitous host protein, CCCTC-binding factor (CTCF), is an essential regulator  
37 of cellular transcription and functions to maintain epigenetic boundaries, stabilise  
38 chromatin loops and regulate splicing of alternative exons. We have previously  
39 demonstrated that CTCF binds to the E2 open reading frame (ORF) of human  
40 papillomavirus (HPV) 18 and functions to repress viral oncogene expression in  
41 undifferentiated keratinocytes by co-ordinating an epigenetically repressed chromatin  
42 loop within HPV episomes. Cellular differentiation, which is necessary for HPV life  
43 cycle completion disrupts CTCF-dependent chromatin looping of HPV18 episomes  
44 inducing enhanced activity of the HPV18 early promoter P<sub>105</sub> and increased viral  
45 oncogene expression.

46 To further characterise CTCF function in HPV transcription control we utilised direct,  
47 long-read Nanopore RNA-sequencing which provides information on the structure  
48 and abundance of full-length transcripts. Nanopore analysis of primary human  
49 keratinocytes containing HPV18 episomes before and after synchronous  
50 differentiation allowed quantification of viral transcript species in these cultures,  
51 including the identification of low abundance novel transcripts. Comparison of  
52 transcripts produced in wild type HPV18 genome-containing cells to those identified  
53 in CTCF-binding deficient genome-containing cells (HPV18-ΔCTCF) identifies CTCF  
54 as a key regulator of differentiation-dependent late promoter activation, required for  
55 efficient E1<sup>^</sup>E4 and L1 protein expression. Furthermore, our data show that CTCF  
56 binding at the E2 ORF of HPV18 promotes usage of the downstream weak splice  
57 donor (SD) sites SD3165 and SD3284, to the dominant E4 splice acceptor site at  
58 nucleotide 3434. These findings demonstrate importance of CTCF-dependent  
59 transcription regulation at multiple stages of the HPV life cycle.

60

61 **IMPORTANCE**

62 Oncogenic human papillomavirus (HPV) infection is the cause of a subset of  
63 epithelial cancers of the uterine cervix, other anogenital areas and the oropharynx.  
64 HPV infection is established in the basal cells of epithelia where a restricted  
65 programme of viral gene expression is required for replication and maintenance of  
66 the viral episome. Completion of the HPV life cycle is dependent on the maturation  
67 (differentiation) of infected cells which induces enhanced viral gene expression and  
68 induction of capsid production. We previously reported that the host cell

69 transcriptional regulator, CTCF, is hijacked by HPV to control viral gene expression.  
70 In this study, we use long-read mRNA sequencing to quantitatively map the variety  
71 and abundance of HPV transcripts produced in early and late stages of the HPV life  
72 cycle and to dissect the function of CTCF in controlling HPV gene expression and  
73 transcript processing.

74

## 75 **INTRODUCTION**

76 Human papillomaviruses (HPVs) are a family of small, double-stranded DNA viruses  
77 that infect cutaneous and mucosal epithelia. Most HPV types cause benign epithelial  
78 hyperproliferation, which is usually resolved by host immune activation. However,  
79 persistent infection with a subset of HPV types (e.g., HPV16 and 18) are the cause  
80 of epithelial tumours including cervical and other anogenital cancers, and carcinoma  
81 of the oropharyngeal tract (1).

82 The viral episome is maintained and replicated in the cell nucleus as an  
83 extrachromosomal, chromatinised episome which allows the epigenetic regulation of  
84 viral transcription in an equivalent manner to host genes (2). The regulation of HPV  
85 gene expression in differentiating epithelia is tightly regulated and is a key strategy in  
86 the maintenance of persistent infection. Several distinct transcriptional start sites  
87 (TSSs) have been identified including the major early and late promoters, the E8  
88 promoter ( $P_{E8}$ ) and less well-defined TSSs around nucleotide 520 ( $P_{520}$ ) and 3000  
89 ( $P_{3000}$ ). The relative activity of these promoters is dependent on the differentiation  
90 status of the host keratinocyte (3-5). Establishment of HPV infection occurs in the  
91 undifferentiated basal keratinocytes of epithelia where viral genome copy number  
92 and transcription are maintained at low levels, presumably to prevent host immune  
93 activation. We and others have shown that the viral episome is maintained in an  
94 epigenetically repressed state in undifferentiated keratinocytes, characterised by low  
95 abundance of trimethylation of lysine 4 (H3K4Me3) and enrichment of trimethylation  
96 of lysine 27 (H3K27Me3) on histone H3, which attenuates viral gene expression (5,  
97 6). The host cell chromatin-organising and transcriptional insulation factor, CCCTC-  
98 binding factor (CTCF) is important in the maintenance of the epigenetic repression of  
99 the HPV genome through the stabilisation of a chromatin loop. CTCF binds to a  
100 conserved site in the E2 open reading frame (ORF) of HPV18 approximately 3,000  
101 base pairs downstream of the viral transcriptional enhancer situated in the long  
102 control region (LCR) (7). Although the major CTCF binding site and the viral

103 enhancer are physically separated, we demonstrated that abrogation of CTCF  
104 binding resulted in inappropriate epigenetic activation of the HPV18 enhancer and  
105 early promoter (termed P<sub>105</sub> in HPV18) and increased expression of the viral  
106 oncoproteins E6 and E7 (E6/E7) (6, 7). CTCF physically associates with the  
107 transcriptional repressor Ying Yang 1 (YY1) (8) and we subsequently showed that  
108 CTCF-dependent epigenetic repression of the HPV18 episome was through  
109 interaction with YY1 bound at the viral LCR, such that CTCF and YY1 co-operate to  
110 stabilise an epigenetically repressed chromatin loop within the early gene region (6).  
111 While the association of CTCF with the HPV18 episome is not significantly altered by  
112 keratinocyte differentiation, YY1 protein expression and binding to the HPV18  
113 genome is dramatically reduced in differentiated keratinocytes leading to loss of  
114 CTCF-YY1 dependent chromatin loop stabilisation (6). This differentiation-dependent  
115 topological change in the HPV episome is coincident with epigenetic activation of the  
116 P<sub>105</sub> promoter and increased expression of the HPV E6/E7 oncoproteins.  
117 Activation of the major late promoter (termed P<sub>811</sub> in HPV18) in part occurs through  
118 epigenetic derepression of the HPV episome upon keratinocyte differentiation (5, 6,  
119 9) and reviewed in (10). This restricts expression of the viral capsid proteins L1 and  
120 L2 to the upper compartment of infected epithelia, limiting their potential for host  
121 immune activation (4, 11, 12). The late promoter also regulates expression of viral  
122 intermediate genes including E1, E2, E1<sup>^</sup>E4 and E5, which are important for viral  
123 genome amplification in the upper layers of the infected epithelia (13, 14). The  
124 mechanisms underlying the differentiation-dependent epigenetic activation of late  
125 promoter activity are not clear, but it has been shown that the viral enhancer in the  
126 LCR is required for late promoter activation (15) and that differentiation-dependent  
127 enhancement of transcription elongation may play a key role in late promoter  
128 activation (16).  
129 Further enhancing the complexity of HPV gene expression regulation, the  
130 polycistronic HPV mRNA is subject to extensive post-transcriptional splicing, which  
131 gives rise to an array of transcripts that each encode a distinct subset of full length,  
132 and/or fusion proteins. While studies have mapped the HPV18 transcriptome (17,  
133 18), the quantification of HPV promoter activity and the abundance of each mature  
134 transcript has not been reported. Cellular splicing factors are utilised and  
135 manipulated by the virus to co-ordinate differentiation dependent viral transcript  
136 splicing, including the serine-arginine rich (SR) proteins and heterogeneous

137 ribonucleoproteins (hnRNPs) (19, 20). In addition to its functions in chromatin  
138 looping and epigenetic isolation, CTCF can play an important role in regulating  
139 alternative gene splicing, most likely through multiple mechanisms. In the host cell  
140 *CD45* locus, CTCF binding within exon 5 promotes inclusion of upstream exons by  
141 creating a “roadblock” to pause RNA polymerase II progression, allowing more  
142 efficient recognition of weak exons by the splicing machinery (21). It has also been  
143 shown that DNA methylation-dependent binding of CTCF within normally weak  
144 exons promotes inclusion during co-transcriptional splicing (22). To support these  
145 findings, a significant enrichment of CTCF binding sites in close proximity to  
146 alternatively spliced exons has been reported (23). However, CTCF binding at  
147 distant sites can also influence alternative exon usage through the stabilisation of  
148 intragenic chromatin loops (24). Our early analysis of CTCF-dependent control of  
149 HPV18 transcript splicing indicated an important role for this factor in maintaining the  
150 complexity of slicing events (7) but the global effect of CTCF on HPV18 transcript  
151 processing was not analysed.

152 Next generation sequencing (NGS) has revolutionised virology research by providing  
153 nucleotide resolution data on existing and emerging pathogens, prevalence, and  
154 evolution. However, conventional Illumina-based RNA sequencing (RNA-Seq)  
155 methods are limited in that information on the structure of full-length transcripts,  
156 including alternative splicing is sacrificed to preserve accuracy and read depth (25).  
157 Direct, long-read Nanopore sequencing overcomes this limitation by providing  
158 quantitative data on the abundance of individual mRNA isoforms (26).

159 In this study, we use Nanopore sequencing to quantify the spectrum of HPV18  
160 transcripts in HPV18 episome-containing primary human keratinocytes and to map  
161 differentiation-induced changes in promoter usage, splicing and transcript  
162 abundance. Furthermore, we characterise the global effect of CTCF binding to the  
163 HPV18 genome on transcript splicing and early and late promoter activity.

164

## 165 **METHODS**

### 166 *Ethical approval*

167 The collection of neonatal foreskin tissue for the isolation of primary human foreskin  
168 keratinocytes (HFKs) for investigation of HPV biology was approved by Southampton  
169 and South West Hampshire Research Ethics Committee A (REC reference number  
170 06/Q1702/45). Written consent was obtained from the parent/guardian. The study

171 was approved by the University of Birmingham Ethical Review Committee (ERN 16-  
172 0540).

173

#### 174 *Cell culture, methylcellulose differentiation and organotypic raft culture*

175 Normal primary HFKs from neonatal foreskin epithelia were transfected with  
176 recircularised HPV18 wild type (WT) or  $\Delta$ CTCF genomes and maintained on  
177 irradiated J2-3T3 fibroblasts in complete E medium (27) as previously described (7).  
178 For methylcellulose-induced keratinocyte differentiation,  $3 \times 10^6$  HPV18-WT or -  
179  $\Delta$ CTCF genome containing keratinocytes were suspended in E-media supplemented  
180 with 10 % FBS and 1.5 % methylcellulose and incubated at 37 °C, 5 % CO<sub>2</sub> for 48  
181 hrs. Cells were then harvested by centrifugation at 250 x g followed by washing with  
182 ice-cold PBS. Cells were then either suspended in medium containing 1 %  
183 formaldehyde to cross-link for chromatin immunoprecipitation (ChIP) as described  
184 below, or RNA and protein was extracted from cell pellets as previously described  
185 (7).

186

187 Organotypic raft cultures were prepared as previously described (7). Rafts were  
188 cultured for 14 days in E medium without epidermal growth factor to allow cellular  
189 stratification. Raft cultures were fixed in 3.7 % formaldehyde and paraffin embedded  
190 and sectioned by ProPath Ltd (Hereford, United Kingdom).

191

#### 192 *Antibodies*

193 Anti-CTCF (61311) and anti-H4Ac (39925) antibodies was purchased from Active  
194 Motif, HPV18 L1 (ab69) antibody purchased from Abcam and anti- $\beta$ -actin (AC-74)  
195 purchased from Sigma. E1<sup>A</sup>E4 antisera were produced as previously described (28).  
196 HRP-conjugated anti-mouse and anti-rabbit antibodies (Jackson Laboratories) were  
197 used for Western blotting and Alexa-488 and -594 conjugated anti-rabbit/mouse  
198 antibodies (Invitrogen) were used for immunofluorescence staining.

199

#### 200 *Chromatin immunoprecipitation-qPCR (ChIP-qPCR)*

201 ChIP-qPCR assays were performed using the ChIP-IT Express Kit (Active Motif) as  
202 per the manufacturer's protocol. Briefly, cells were fixed in 1 % formaldehyde for 5  
203 mins at room temperature with gentle rocking, quenched in 0.25 M glycine and  
204 washed with ice-cold PBS. Nuclei were released using a Dounce homogeniser.

205 Chromatin shearing was carried out by sonication at 25 % amplitude for 30 secs  
 206 on/30 secs off for a total time of 15 mins using a Sonics Vibracell sonicator fitted with  
 207 a microprobe. ChIP efficiency was assessed by qPCR using SensiMix SYBR master  
 208 mix using a Stratagene Mx3005P (Agilent Technologies, Santa Clara, CA, USA).  
 209 Primer sequences for ChIP experiments are shown in Table 1. Cycle threshold ( $C_T$ )  
 210 values were used to calculate fold enrichment compared to a negative control FLAG  
 211 antibody with the following formula:

212

213 
$$\text{Fold binding over IgG} = (2^{\Delta C_T \text{ Target}})/(2^{\Delta C_T \text{ IgG}})$$

214

215 Where  $\Delta C_T \text{ target} = \text{Input } C_T - \text{Target } C_T$  and  $\Delta C_T \text{ IgG} = \text{Input } C_T - \text{IgG } C_T$ . Each  
 216 independent experiment was performed in technical triplicate and data shown are  
 217 the mean and standard deviation of three independent repetitions.

218

219

Primer pair (amplicon mid-point)	Fw (5' – 3')	Rev (5' – 3')	Ta (°C)
4440	GGGGTCGTACAGGGTACATT	GATGTTATATCAAACCCAGACG TG	56
5381	TCTGCCTCTTCCTATAGTAAT GTAACG	GGAATAAAATAATATAATGGCC ACAAA	56
5655	CCTCCTTCTGTGGCAAGAGT	GGTCAGGTAAGTGCACCCTAA	56
6659	AGTCTCCTGTACCTGGGCAA	AACACCAAAGTTCCAATCCTCT	58
7301	GTGTGTTATGTGGTTGCGCC	GGATGCTGTAAGGTGTGCAG	58
7746	ACTTTCATGTCCAACATTCTGT CT	ATGTGCTGCCCAACCTATTT	56
115	TGTGCACGGAAGTGAACACT	CAGCATGCGGTATACTGTCTC	58
751	CGAACCAACAACGTCACACAAT	ACGGACACACAAAGGACAGG	58
1500	GCAATGTATGTAGTGGCGGC	TACTGCTGTTGTTGCCCT	58
2819	TGCAGACACCGAAGGAAACC	CATTTTCCCAACGTATTAGTTGC C	58
2926	GGCAACTAATACGTTGGGAAA A	TGTCTTGCAGTGTCCAATCC	56
3165	AGGTGGCCAAACAGTACAAGT	GCCGTTTTGTCCCATGTTCC	58
3381	TGGGAAGTACATTTTGGGAAT AA	TCCACAGTGTCCAGGTCGT	56
3971	TATGTGTGCTGCCATGTCCC	CTGTGGCAGGGGACGTTATT	56

220

221 **Table 1: Primer sequences used for ChIP-qPCR experiments.** Ta, annealing  
 222 temperature



223

#### 224 *ChIP-Seq*

225 ChIP and respective input samples were used for generation of ChIP-Seq libraries  
226 as described (29). Briefly, 2-10 ng DNA was used in conjunction with the NEXTflex  
227 Illumina ChIP-Seq library prep kit (Cat# 5143-02) as per the manufacturer's protocol.  
228 Samples were sequenced on a HiSeq 2500 system (Illumina) using single read  
229 (1x50) flow cells. Sequencing data was aligned to the HPV18 genome (accession  
230 number: AY262282.1) using Bowtie (30) with standard settings and the -m1 option  
231 set to exclude multi mapping reads (31).

232

#### 233 *RNA sequencing and data analysis*

234 For RNA-Seq, libraries were prepared using Tru-Seq Stranded mRNA Library Prep  
235 kit for NeoPrep (Illumina, San Diego, CA, USA) using 100ng total RNA input  
236 according to manufacturer's instructions. Libraries were pooled and run as 75-cycle-  
237 pair end reads on a NextSeq 550 (Illumina) using a high-output flow cell. Sequencing  
238 reads were aligned to human (GRCh37) and HPV18 (AY262282.1) genomes with  
239 STAR aligner (v2.5.2b) (32). The computations were performed on the CaStLeS  
240 infrastructure (33) at the University of Birmingham. Sashimi plots were generated in  
241 Integrative Genomics Viewer (IGV), Broad Institute  
242 (<http://software.broadinstitute.org/software/igv/>).

243

#### 244 *Nanopore direct RNA sequencing and data analysis*

245  $8 \times 10^7$  cells from undifferentiated or methylcellulose differentiated keratinocytes  
246 containing HPV18 (WT or  $\Delta$ CTCF) samples for RNA extraction using the RNeasy  
247 Plus Mini Kit (Qiagen) according to the manufacturer's instructions and DNaseI  
248 treated (Promega). 500 ng of polyA+ RNA was used in conjunction with the direct  
249 RNA sequencing kit (Oxford Nanopore technologies, Oxford, UK [SQK-RNA002]). All  
250 protocol steps are as described in (34). The reads were aligned to the human  
251 (GRCh37) and HPV18 (AY262282.1) genomes using minimap2 (35) with options "-  
252 ax splice -uf -k14" for nanopore direct RNA mapping. The splicing coordinates were  
253 extracted from the bam files using custom scripts.

254

255

#### 256 *Cell lysis and western blotting*



257 Cells were lysed with urea lysis buffer (ULB; 8 M urea, 100 mM Tris-HCl, pH 7.4, 14  
258 mM  $\beta$ -mercaptoethanol, protease inhibitors) and protein concentration determined.  
259 Protein extracts from organotypic raft cultures were harvested using ULB and  
260 homogenised using a Dounce homogeniser contained within a category II biological  
261 safety cabinet. Lysates were incubated on ice for 20 mins before centrifugation at  
262 16,000 x g for 20 mins at 4 °C. Supernatant was transferred to a fresh tube and  
263 protein concentration assessed by Bradford Assay. For Western blotting, equal  
264 quantities of protein lysates were separated by SDS-PAGE and western blotting was  
265 carried out by conventional methods. Chemiluminescent detection was carried out  
266 using a Fusion FX Pro and densitometry performed with Fusion FX software.

267

### 268 *Immunofluorescence*

269 Immunofluorescence was carried out on paraffin embedded organotypic raft culture  
270 sections using the agitated low temperature epitope retrieval (ALTER) method as  
271 previously described (36). Briefly, slides were sequentially immersed in HistoClear  
272 (Scientific Laboratory Supplies) and 100 % IMS and incubated at 65 °C in 1 mM  
273 EDTA (pH 8.0), 0.1 % Tween 20 overnight with agitation. Slides were then blocked in  
274 PBS containing 20 % heat-inactivated normal goat serum and 0.1 % BSA (Merck).  
275 Primary antibodies were diluted in block solution and incubated overnight at 4 °C  
276 followed by 3x PBS washes. Fluorophore-conjugated secondary antibodies were  
277 diluted in block buffer and added to slides which were incubated at 37 °C for 1 hour.  
278 Slides were subsequently washed 4x 10 mins in PBS with Hoechst 33342 solution  
279 (10  $\mu$ g/ml) added to the final PBS wash. Slides were mounted in Fluoroshield  
280 (Sigma-Aldrich) and visualised using a Nikon inverted Epifluorescent microscope  
281 fitted with a 40x oil objective. Images were captured using a Leica DC200 camera  
282 and software.

283

### 284 **Results**

285 We have previously characterised a CTCF binding site within the E2 open reading  
286 frame (ORF) of HPV18 which is strongly bound by CTCF in a primary HFK model of  
287 the HPV18 life cycle (6, 7). Although the E2-CTCF binding site was the most CTCF  
288 enriched region of the HPV18 genome in our ChIP-qPCR analysis, there did appear  
289 to be other regions of the viral genome that were bound at a lower level by CTCF. In  
290 addition, CTCF binding sites have been predicted in the late gene region of HPV18

291 and other high-risk HPV types and binding has been demonstrated in HPV31  
292 episomes (7, 37). To analyse CTCF binding to the HPV18 genome with greater  
293 sensitivity, we opted to map CTCF binding peaks using ChIP-sequencing (ChIP-  
294 Seq). Anti-CTCF immunoprecipitated chromatin harvested from HFKs harbouring  
295 HPV18 episomes was subject to Illumina next generation sequencing. Reads were  
296 aligned to the HPV18 genome revealing robust enrichment of CTCF in the E2 ORF  
297 with maximal binding between nucleotides 2960-3020, corresponding to the  
298 previously identified E2-CTCF binding site (**Fig.1A**). No other distinct CTCF peaks  
299 were observed in the HPV18 genome. In addition, ChIP-Seq analysis of CTCF  
300 enrichment in mutant HPV18 genomes in which the E2-CTCF binding site was  
301 mutated to prevent CTCF binding (HPV18- $\Delta$ CTCF), revealed a complete loss of  
302 CTCF binding to the E2-ORF with no evidence of enhanced binding at secondary  
303 sites (**Fig.1A**), confirming our previous ChIP-qPCR analysis of this mutant virus.  
304 Abrogation of CTCF binding at the HPV18 E2 ORF resulted in increased  
305 transcriptional activity of the HPV18 early promoter ( $P_{102}$ ) and a concomitant  
306 increase in E6/E7 protein expression (6, 7). These studies also revealed alterations  
307 in the splicing of early transcripts, indicated by a significant reduction in the  
308 abundance of transcripts spliced at 233^3434 upon amplification by semi-  
309 quantitative RT-PCR (7). To confirm these findings and to further characterise  
310 CTCF-dependent regulation of HPV18 transcript splicing, we utilised high-depth  
311 Illumina RNA-Seq data in HPV18-wild type and - $\Delta$ CTCF transfected primary HFKs to  
312 quantify individual splicing events (**Fig.1B**). While there were a similar number of  
313 splicing events at 233^3434 in the wild type and mutant HPV18 genome-containing  
314 cells (403 and 407 events, respectively), splicing at 233^416 was increased in  
315 HPV18- $\Delta$ CTCF genome containing cells in comparison to wild type (28,918 events  
316 compared to 16,557 events respectively, Fisher's test p-value <0.00001), which  
317 could account for the observed relative reduction in amplification of transcripts  
318 spliced at 233^3434 by qRT-PCR (7). Interestingly, we also noted a reduction in  
319 splicing at 3284^3434, previously proposed to encode a truncated form of the E2  
320 protein, E2C (18), and a complete loss of splicing at 3165^3434 in HPV18- $\Delta$ CTCF  
321 genome containing cells compared to wild type HPV18. Found at relatively low  
322 abundance, splicing at 3165^3434 has been previously described and predicted to  
323 encode a novel E2^E4 fusion protein termed E2^E4L (38). Similarly, splicing at  
324 2853^3434 has been proposed to encode a shorter form of E2^E4 fusion protein,

325 E2<sup>E4S</sup> (38), however, this splice was not detected in our Illumina RNA-Seq data.  
326 Splicing from the 929 SD site to the weak acceptor site at 3440 was observed in  
327 HPV18- $\Delta$ CTCF genome containing cells (391 events) but not in wild type HPV18.  
328 These findings suggest that CTCF may play a role in controlling acceptor site usage  
329 downstream of the E2-CTCF binding site.  
330 While individual splicing events can be quantified using conventional short-read RNA  
331 sequencing methods, the evaluation of the structure of individual transcripts and the  
332 multiple splicing events that may occur within a single transcript is not possible. To  
333 fully characterise and, for the first time, quantify the relative abundance of individual  
334 HPV18 transcripts in primary HFKs, purified and polyA<sup>+</sup> enriched RNA was analysed  
335 by direct long-read MinION sequencing. Cells were either grown in monolayer  
336 culture on feeder cells (undifferentiated) or embedded in semi-solid methylcellulose  
337 containing medium for 48 hours, to induce synchronous differentiation, and the  
338 spectrum of transcripts quantified as reads per million (RPM) for each sample. To  
339 confirm that morphological keratinocyte differentiation was induced by suspension in  
340 methylcellulose, differentiation-dependent changes to cellular markers of  
341 keratinocyte differentiation were analysed. A notable increase in involucrin (IVL)  
342 expression was observed (**Fig.2A**; Fisher's test p-value < 0.00001). In addition, an  
343 alteration in transcript splicing of the keratinocyte-specific extracellular matrix protein,  
344 Ecm1, upon keratinocyte differentiation has been reported (39). Undifferentiated  
345 keratinocytes express full length Ecm1 transcript 2 but expression of a shorter,  
346 alternatively spliced transcript (transcript 3) is induced upon keratinocyte  
347 differentiation. Analysis of Ecm1 transcripts in our MinION sequencing data  
348 demonstrated the appearance of Ecm1 transcript 3 which lacks exon 7 in  
349 methylcellulose differentiated keratinocytes only (**Fig.2B**).  
350 Virus host fusion transcripts were identified at very low abundance (<2% of total HPV  
351 reads), indicative of low-level viral integration. Nonetheless, these fusion transcripts  
352 were removed from our data set prior to analysis to include only those transcripts  
353 derived from HPV episomes. Data were then normalised to the total number of reads  
354 in each sample to calculate reads per million (RPM) of each viral transcript species.  
355 In agreement with previous reports (17, 18), five clear groupings of transcriptional  
356 start sites were identified in undifferentiated HPV18 genome containing cells, which  
357 originated between nucleotides 1-350 (P<sub>102</sub>), 351-700 (P<sub>520</sub>), 701-900 (P<sub>811</sub>), 1000-  
358 1400 (P<sub>1193</sub>) and 2800-4000 (P<sub>3000</sub>) (**Fig.2C**), which were used to define transcript

359 species in subsequent quantifications. Keratinocyte differentiation resulted in a  
360 significant change in TSS usage characterised by activation of the P<sub>811</sub> major late  
361 promoter (**Fig.2C**). In undifferentiated HPV18 wild type genome-containing cells, the  
362 most abundant transcript was initiated at the P<sub>102</sub> promoter and spliced at 233^416-  
363 929^3434 (transcript 3; **Fig.3**). This transcript has the potential to encode E6\*1, E7,  
364 E1^E4 and E5. Several novel transcripts were identified at low abundance in the  
365 undifferentiated wild type HPV18 cells including transcripts 13 and 25, which encode  
366 E6\*1 and E7 respectively along with E2 and E5. Interestingly, splicing at both  
367 3165^3434 and 3284^3434 was observed in undifferentiated cells (transcripts 11 and  
368 12; **Fig.3**), however, these transcripts originated from the P<sub>3000</sub> promoter and  
369 therefore lack the E2 start codon at nt2816 and more likely encode E5 in the basal  
370 keratinocytes rather than E2^E4 fusion proteins as previously suggested (38). A  
371 single, previously described transcript, spliced at 929^2779-3165^3434 was  
372 identified in differentiated WT HPV18 genome containing cells that originated at the  
373 P<sub>105</sub> promoter, which contains the E2 start codon and could therefore encode  
374 E2^E4L (transcript 6; **Fig.3**)(38). We also identified very low abundance of  
375 transcripts spliced at 929^2779-3284^3434, which is in frame with the E4 ORF and  
376 we predict would encode a novel E2^E4 fusion protein, denoted E2^E4XL but the  
377 presence of this theoretical protein in HPV18 infected cells is yet to be confirmed  
378 (transcript 5; **Fig.3**).

379 Comparison of viral transcripts in WT and HPV18-ΔCTCF genome-containing cells  
380 revealed a significant increase in abundance of the major early transcript originating  
381 from the P<sub>105</sub> promoter and spliced at 233^416-929^3434, which encodes E6\*1, E7,  
382 E1^E4 and E5 (transcript 3; **Fig.3**, Fisher's test p-value < 0.00001). A more modest  
383 increase in the second most abundant transcript in undifferentiated cells, originating  
384 from the P<sub>105</sub> promoter and spliced at 929^3434 was also observed, which has the  
385 potential to encode full length E6 as well as E7, E1^E4 and E5 (transcript 4; **Fig.3**,  
386 Fisher's test, non-significant). The increased abundance of these major early viral  
387 transcripts corroborates the previously observed increase in E6 and E7 protein  
388 expression when CTCF binding site is ablated (6, 7).

389 Notably, splicing at both 3165^3434 and 3284^3434 (transcripts 11 and 12; **Fig.3**)  
390 was significantly reduced in HPV18-ΔCTCF genome containing cells compared to  
391 WT (Fisher's test p-value < 0.00001 and 0.01, respectively) corroborating our finding  
392 in Illumina RNA-Seq datasets that CTCF may function to enhance the activity of

393 downstream weak SD sites in the HPV18 genome. Transcripts spliced at 929<sup>^</sup>3440  
394 (transcripts 13, 14 and 15) were also detected at low abundance.

395 Transcripts that originate from the P<sub>811</sub> late promoter were abundantly expressed in  
396 undifferentiated cells; transcripts originating from this promoter and spliced at  
397 929<sup>^</sup>3434 to encode E1<sup>^</sup>E4 and E5 proteins (transcript 9; **Fig.3**) were the second  
398 most abundant transcript in undifferentiated cells. As expected, the abundance of  
399 this transcript was dramatically increased around 50-fold (Fisher's test p-value <  
400 0.00001) upon differentiation of the WT HPV18 cells in methylcellulose. However,  
401 while differentiation of HPV18- $\Delta$ CTCF genome-containing cells similarly resulted in  
402 an increase in abundance of this major E1<sup>^</sup>E4 encoding transcript, the overall  
403 abundance of this transcript was reduced by around 50 % compared to WT. It is also  
404 interesting to note that transcripts encoding the L1/L2 capsid proteins (transcripts 31-  
405 34; **Fig.3**) were induced upon cellular differentiation in WT genome-containing cells,  
406 albeit at a low level, but these transcripts were all lower in abundance in HPV18-  
407  $\Delta$ CTCF cells. These data suggest that recruitment of CTCF to the HPV18 genome at  
408 the E2-ORF may be important for differentiation-dependent activation of the viral late  
409 promoter.

410 The major transcriptional promoters in the HPV18 genome have been previously  
411 mapped using 5' RACE (17). Although transcript sequencing by Nanopore does not  
412 provide nucleotide resolution accuracy in mapping transcription start sites (40), the  
413 clustering of the 5' end of viral transcripts was clearly enriched at the previously  
414 annotated transcriptional start sites (**Fig.2C**). Therefore, to characterise the  
415 differential activity of the major viral promoters in HPV18 WT and - $\Delta$ CTCF cells, the  
416 5' end of each viral read in our Nanopore datasets was mapped and quantified. The  
417 5' end of most transcripts (>95 %) mapped near four previously described  
418 promoters; P<sub>102</sub>, P<sub>520</sub>, P<sub>811</sub> and P<sub>3000</sub> (**Fig.4**). Interestingly, the 5' end of transcripts  
419 that originated from both the P<sub>102</sub> and P<sub>811</sub> promoters clustered as a sharp peak at  
420 the previously annotated transcriptional start site whereas the 5' end of transcripts  
421 originating from either the P<sub>520</sub> or P<sub>3000</sub> promoter regions were more broadly  
422 distributed (**Fig.4A-D**). As expected, the P<sub>102</sub> promoter was the most active promoter  
423 in HPV18 WT genome-containing undifferentiated cells with very few transcripts  
424 originating from the P<sub>811</sub> late promoter. Differentiation of these cells resulted in a  
425 dramatic increase in transcripts originating from the P<sub>811</sub> promoter (Fisher's test p-  
426 value < 0.00001), coincident with a slight increase in P<sub>102</sub> activity (Fisher's test p-

427 value  $< 0.00001$ ) (**Fig.4A and C**). Transcripts originating from the  $P_{102}$  promoter  
428 were ~30 % more abundant in HPV18- $\Delta$ CTCF genome containing cells than WT,  
429 which was further activated upon cellular differentiation confirming enhanced activity  
430 of the early promoter in the absence of CTCF recruitment. Interestingly, the activity  
431 of the  $P_{811}$  late promoter was notably lower in differentiated HPV18- $\Delta$ CTCF genome  
432 containing cells compared to WT (Fisher's test p-value  $< 0.00001$ ), providing  
433 evidence that the activity of the late promoter in differentiated cells is attenuated  
434 when CTCF recruitment is abrogated. The  $P_{520}$  promoter had reduced activity  
435 compared to the  $P_{102}$  and  $P_{811}$  promoters, but interestingly this promoter was more  
436 active in differentiated HPV18 WT genome-containing cells than undifferentiated  
437 (Fisher's test p-value  $1.8E-4$ ). In HPV18- $\Delta$ CTCF cells, the  $P_{520}$  promoter was more  
438 active than WT in undifferentiated cells (Fisher's test p-value 0.016) but was not  
439 further activated by cellular differentiation (**Fig.4B**). Very few transcripts originated  
440 from  $P_{3000}$  in undifferentiated cells, however this promoter was strongly activated  
441 following cellular differentiation in HPV18 WT genome containing cells. As was  
442 observed at  $P_{811}$ , differentiation-dependent activation of  $P_{3000}$  was reduced in  
443 HPV18- $\Delta$ CTCF genome containing cells compared to WT. The  $P_{E8}$  promoter ( $P_{1193}$ )  
444 was only weakly active with less than 3 % of transcripts originating at this promoter  
445 in undifferentiated cells. Furthermore, the activity of  $P_{E8}$  was not affected by  
446 keratinocyte differentiation or mutation of the E2-CTCF binding site (data not shown).  
447 Analysis of TSS usage in the bulk population of viral transcripts revealed that while  
448 there was a greater proportion of transcripts which initiated from the  $P_{102}$  early  
449 promoter in HPV18- $\Delta$ CTCF episomes than WT (indicated by tighter density  
450 grouping), this did not reach significance ( $p$  0.16) (**Fig.5A**). In contrast, highly  
451 significant differences were observed between TSS usage in HPV18- $\Delta$ CTCF  
452 episomes compared to WT following keratinocyte differentiation ( $p < 1E-16$ ). While in  
453 WT HPV18 cells, the TSS usage density was highly enriched at the  $P_{811}$  promoter,  
454 transcripts in  $\Delta$ CTCF-HPV18 genome-containing cells were less abundant at the  $P_{811}$   
455 promoter, and the  $P_{102}$  promoter was proportionately more active than in WT-HPV18  
456 episomes (**Fig.5B**). These analyses demonstrate that differentiation-dependent  
457 stimulation of  $P_{811}$  major late promoter activity is facilitated by recruitment of CTCF to  
458 the E2 ORF.



459 We previously demonstrated that in undifferentiated cells, HPV18- $\Delta$ CTCF episomes  
460 had increased trimethylation of lysine 4 in histone 3 (H3K4Me3) at the P<sub>102</sub> early  
461 promoter compared to WT, correlating with increased promoter activity. However,  
462 while differentiation of HPV18 WT genome-containing cells resulted in a significant  
463 enrichment of H3K4Me3 at the P<sub>811</sub> late promoter, no such enrichment was observed  
464 in HPV18- $\Delta$ CTCF episomes (6). Enhanced acetylation of histones is also indicative  
465 of enhanced activation of transcription by facilitating increased chromatin  
466 accessibility and the recruitment of transcriptional activators (41). We therefore  
467 assessed the changes in histone 4 acetylation (H4Ac) in HPV18 episomes induced  
468 by keratinocyte differentiation. H4Ac abundance in the viral genome in  
469 undifferentiated cells was detectable at low levels, consistent with restricted virus  
470 transcription (**Fig.6A**). Differentiation of the cells in methylcellulose resulted in a  
471 dramatic increase in H4Ac abundance throughout the WT-HPV18 genome, with an  
472 over 10-fold enrichment at the P<sub>811</sub> late promoter, consistent with increased  
473 production of late transcripts (**Fig.6A**). However, HPV18- $\Delta$ CTCF episomes were  
474 devoid of H4Ac with only a small, insignificant increase in H4Ac abundance at the  
475 P<sub>811</sub> following differentiation (**Fig.6B**). Together, these findings suggest that CTCF  
476 recruitment to the E2-ORF is necessary for appropriate epigenetic programming of  
477 the viral chromatin and differentiation-dependent transcriptional activation of P<sub>811</sub>.  
478 To determine whether the reduced differentiation-dependent activation of P<sub>811</sub> in  
479 HPV18- $\Delta$ CTCF genomes resulted in reduced late protein expression, we analysed  
480 E1<sup>E4</sup> protein in methylcellulose differentiated cultures. Western blotting of lysates  
481 harvested from HPV-18 WT and - $\Delta$ CTCF genome containing cells before and after  
482 differentiation revealed an induction of involucrin protein expression. However, there  
483 was a significant attenuation of E1<sup>E4</sup> protein expression when CTCF binding to the  
484 viral genome was abrogated (**Fig.7A and B**). Since L1 protein is not robustly  
485 expressed in methylcellulose differentiated keratinocytes, we analysed L1 protein  
486 expression by immunostaining organotypic raft culture sections derived from two  
487 independent donors of HPV18-WT and - $\Delta$ CTCF genome containing cells. L1 positive  
488 cells were visible in the upper layers of HPV18-WT genome containing rafts but were  
489 barely detectable in HPV18- $\Delta$ CTCF rafts and this difference was significant (**Fig.7C**  
490 **and D**). While the total number of E1<sup>E4</sup> positive cells in the upper layers of HPV18-  
491  $\Delta$ CTCF rafts was not altered, the intensity of E1<sup>E4</sup> staining was notably reduced  
492 (**Fig.7C**). Western blot analysis of protein lysates harvested from three independent



493 raft cultures confirmed a significant reduction in E1<sup>E4</sup> protein abundance in HPV18-  
494  $\Delta$ CTCF genome containing raft cultures in comparison to WT (**Fig.7E**).

495

## 496 **DISCUSSION**

497 The differentiation-dependent regulation of papillomavirus transcription is  
498 fundamental to the productivity and persistence of infection. Previous studies have  
499 shown that the viral early ( $P_{105}$ ) promoter is active in basal keratinocytes and  
500 becomes further activated as the cells enter terminal differentiation (5, 6). In contrast,  
501 the viral late promoter ( $P_{811}$ ) is repressed in undifferentiated basal cells and strongly  
502 activated upon induction of cellular differentiation (4, 5, 9, 16, 42). In this study, we  
503 have utilised direct, long-read RNA sequencing to quantitatively analyse HPV18  
504 promoter activity and to dissect the role of CTCF in regulating viral transcription at  
505 key stages of the virus life cycle. Our findings confirm the differentiation-dependent  
506 model of HPV transcription control; transcripts that originate from the  $P_{105}$  promoter  
507 are dominant in undifferentiated cells and further increased in abundance upon  
508 cellular differentiation. The abundance of transcripts originating from the  $P_{811}$  late  
509 promoter is low in undifferentiated cells but is dramatically upregulated when cells  
510 are differentiated. Transcription originating from the  $P_{520}$  and  $P_{3000}$  promoter regions  
511 is also activated by cellular differentiation but overall, these promoters are far less  
512 active than either the  $P_{105}$  or  $P_{811}$  promoters. The  $P_{E8}$  promoter is equally weak in  
513 both undifferentiated and differentiated cells with only two transcript species that  
514 originate from this TSS. The most dominant transcript identified from the  $P_{E8}$   
515 promoter was spliced at 1357<sup>3434</sup> and encodes E8<sup>E2</sup> and E5. The second  
516 transcript, spliced at 1357<sup>3465</sup> to encode E5 only, was slightly increased in  
517 expression in differentiated cell cultures.

518 Transcripts that encode fusion products between the E2 and E4 ORFs (E2<sup>E4</sup>) have  
519 been previously described (38). These transcripts were reported to originate  
520 upstream of the E2 start code at position 2816 in HPV18 and therefore encode a  
521 protein fusion between the N-terminus of E2 and the C-terminus of E4. E2<sup>E4S</sup>  
522 encoding transcripts, spliced at 2853<sup>3434</sup>, were not identified in any of our  
523 Nanopore or RNA-Seq datasets. We did however detect transcripts spliced at  
524 3165<sup>3434</sup>, which have been previously described to encode a fusion protein termed  
525 E2<sup>E4L</sup> (38). However, this transcript was detected at very low abundance (~1 RPM)  
526 and only in differentiated keratinocytes. Interestingly, we also identified a third

527 transcript that may encode a previously uncharacterised E2<sup>E4</sup> fusion protein, which  
528 we termed E2<sup>E4XL</sup>. This transcript originated from the P<sub>105</sub> promoter and was  
529 spliced at 929<sup>2779</sup> and 3284<sup>3434</sup>. Like E2<sup>E4L</sup>, this transcript retains the E2 start  
530 codon but potentially encodes amino acids 1-156 of E2 fused to amino acid 6-88 of  
531 E4, but it remains to be determined if this transcript encodes a bone fide E2<sup>E4</sup>  
532 fusion protein. Interestingly, most of the transcripts that originated from the P<sub>3000</sub>  
533 promoter were also spliced 3165<sup>3434</sup> or 3284<sup>3434</sup>. These transcripts were in  
534 higher abundance than those originating from the P<sub>105</sub> promoter in both  
535 undifferentiated and differentiated cells, but since they lack the E2 start codon, they  
536 are likely to encode E5 protein only. Supporting this hypothesis, splicing of  
537 transcripts originating from the P<sub>3000</sub> promoter 3165<sup>3434</sup> and 3284<sup>3434</sup>  
538 respectively removes several intronic ATG codons (7 and 11, respectively),  
539 potentially facilitating enhanced translation of E5.

540 Comparison of the HPV18 transcript map between WT and  $\Delta$ CTCF genome-  
541 containing cells revealed several important phenotypes. Firstly, abrogation of CTCF  
542 binding resulted in enhanced production of transcripts originating from the P<sub>105</sub>  
543 promoter, in agreement with our previous findings (6, 7). The increased P<sub>105</sub> activity  
544 resulted in an increase in transcripts spliced at 233<sup>416-929</sup><sup>3434</sup> (encoding E6<sup>I</sup>,  
545 E7, E1<sup>E4</sup> and E5) and 929<sup>3434</sup>, (encoding E6, E7, E1<sup>E4</sup> and E5) while there  
546 was a small decrease in transcripts spliced solely at 233<sup>416</sup> (encoding E6<sup>I</sup>, E1, E7  
547 and E2) and 233<sup>3434</sup> (the only known transcript to encode E6<sup>II</sup>), confirming our  
548 previous observation that abrogation of CTCF binding to the HPV18 genome  
549 reduces the abundance of transcripts spliced at 233<sup>3434</sup> (7). In addition, a marked  
550 decrease in transcripts spliced at 3165<sup>3434</sup> and 3284<sup>3434</sup> was observed in  
551  $\Delta$ CTCF-HPV18 genome containing cells in comparison to WT, confirming our initial  
552 analysis of HPV18 transcript splicing by conventional RNA-Seq. These data indicate  
553 that CTCF plays a key role in splice donor choice when splicing to the dominant  
554 splice acceptor site at nucleotide 3434 in the HPV18 genome.

555 A functional role for CTCF in influencing cellular co-transcriptional alternative splicing  
556 has been previously demonstrated. CTCF binding within or downstream of weak  
557 exons can promote exon inclusion by creating a roadblock to pause RNA  
558 polymerase II progression, allowing greater splicing efficiency (22, 23, 43).  
559 Interestingly, CTCF-mediated chromatin loop stabilisation between gene promoter  
560 and exon regions also plays a key role in regulating alternative splicing events.

561 Exons downstream of a CTCF stabilised promoter-exon loop are more likely to be  
562 included in the nascent mRNA, providing a functional link between three-dimensional  
563 chromatin organisation and splicing regulation (24). Notably, in this study we show  
564 that this mechanism of splicing regulation is recapitulated in the HPV18 genome.  
565 CTCF binding within the HPV18 E2 ORF, upstream of weak splice donor sites at 3165  
566 and 3284 results in YY1-dependent stabilisation of a distinct chromatin loop with the  
567 upstream viral promoter (6) and correlates with increased splicing at both 3165^3434  
568 and 3284^3434 to produce E5 encoding transcripts.

569 As expected, cellular differentiation strongly induced P<sub>811</sub> promoter activation in WT  
570 HPV18 episomes. However, HPV18- $\Delta$ CTCF genome containing cells displayed a  
571 notable reduction in the abundance of transcripts originating from this promoter  
572 following differentiation. Differentiation dependent activation of the P<sub>3000</sub> promoter  
573 was also attenuated in virus unable to bind CTCF. In contrast the P<sub>520</sub> promoter in  
574  $\Delta$ CTCF-HPV18 was active in both undifferentiated and differentiated cells, albeit at a  
575 low level. Activity of P<sub>520</sub> was induced by cellular differentiation in WT HPV18 cells  
576 but was not further activated in  $\Delta$ CTCF-HPV18 cells. In agreement with the observed  
577 differentiation induced activation of the P<sub>811</sub> and P<sub>3000</sub> promoters in WT HPV18  
578 episomes, we demonstrated a marked increase in H4Ac enrichment, particularly in  
579 around the P<sub>811</sub> and P<sub>3000</sub> promoters. Interestingly, H4Ac enrichment following  
580 differentiation was not recapitulated in  $\Delta$ CTCF-HPV18 episomes, indicating that  
581 CTCF binding to the E2-ORF is important for enhanced transcriptional activation in  
582 the late stages of the virus life cycle. Importantly, attenuation of differentiation-  
583 dependent late promoter activation in  $\Delta$ CTCF-HPV18 resulted in significantly  
584 reduced E1^E4 protein expression following methylcellulose differentiation and an  
585 almost complete loss of L1 protein expression in stratified epithelia. These results  
586 demonstrate for the first time that CTCF has essential functions in differentiation-  
587 dependent transcriptional dynamics in the late stages of the HPV life cycle.

588

## 589 **ACKNOWLEDGEMENTS**

590 This work was funded by grants from the Medical Research Council awarded to JLP  
591 and SR (MR/R022011/1, MR/T015985/1 and MR/N023498/1). IP was supported by a  
592 Cancer Research UK non-clinical PhD studentship awarded to JLP and SR. BN is  
593 funded through the Cancer Research UK Birmingham Centre award  
594 C17422/A25154. The funders had no role in study design, data collection and

595 interpretation, or the decision to submit the work for publication. We thank Dr.  
596 Joseph Spitzer and his patients for the collection and donation of foreskin tissue.

597

598

### 599 **Figure Legends**

600 **Figure 1: Abrogation of CTCF recruitment to the HPV18 E2 ORF alters early**  
601 **transcript splicing.** (A) Enrichment of CTCF in the HPV18 genome was assessed  
602 by ChIP-Seq in either WT (blue) or  $\Delta$ CTCF-HPV18 (red) genome-containing  
603 keratinocytes. Next generation sequencing data were IGV. The position of HPV18  
604 ORFs and LCR are indicated below the alignment profiles. (B) Exon-exon junctions  
605 in Illumina-based RNA-Seq data sets of either WT (blue) or  $\Delta$ CTCF-HPV18 (red)  
606 genome-containing keratinocytes were identified and quantified in IGV and  
607 represented in Sashimi plots. The co-ordinates of splice donor and acceptor sites  
608 and annotated ORFs are indicated. The number of reads at each exon-exon junction  
609 is indicated. \*denotes splicing event identified in WT HPV18 but reduced or lost in  
610  $\Delta$ CTCF-HPV18 genome containing cells.

611

612 **Figure 2: Analysis of differentiation-dependent host cell gene expression and**  
613 **HPV transcriptional start site usage.** HPV18-HFK were synchronously  
614 differentiated in methylcellulose for 48 hrs (green). The host and viral transcriptomes  
615 were analysed by long read RNA-Seq and compared to undifferentiated HPV18-HFK  
616 (blue). (A) Enhanced involucrin (IVL) expression following keratinocyte differentiation  
617 and (B) enhanced ECM1 expression combined with differentiation-induced exon 7  
618 skipping. (C) Clustered HPV18 TSS usage in undifferentiated and differentiated  
619 keratinocytes showing differentiation-dependent alteration of the major early ( $P_{102}$ )  
620 and major late ( $P_{811}$ ) promoter usage. \*\*\*\* $p < 0.0001$  (Fisher's test).

621

622 **Figure 3: Quantitative analysis of the HPV18 transcriptome in undifferentiated**  
623 **and differentiated keratinocytes and alterations induced by abrogation of**  
624 **CTCF binding.** Alignment of direct RNA sequencing data to the HPV18 genome  
625 facilitated the characterisation of all HPV-specific transcripts. Transcripts were only  
626 included in the data set if they were represented by two or more reads. The relative  
627 abundance of each transcript type was calculated in reads per million (RPM) of the  
628 total reads in each sample. Relative abundance (RPM) of each transcript is shown

629 for WT (blue) and  $\Delta$ CTCF-HPV18 (red) genome-containing cells in undifferentiated  
630 keratinocytes (left) and for WT (green) and  $\Delta$ CTCF-HPV18 (purple) in differentiated  
631 keratinocytes (right). The identified splice donor (blue) and acceptor (green) sites are  
632 indicated above the transcript map and HPV18 ORFs encoded by each transcript are  
633 shown. \*denotes transcripts that have previously been identified (17, 18).

634

635 **Figure 4: Quantitative analysis of transcription start site usage in**  
636 **undifferentiated and differentiated keratinocytes and CTCF-dependent**  
637 **regulation of promoter activity.** The 5' end of each HPV18 transcript was identified  
638 in Nanopore RNA sequencing data sets and relative abundance calculated as reads  
639 per million (RPM). Total counts at each nucleotide position were binned into 10 (A,  
640 B, C) or 100 (D) nucleotide regions in the data shown. Transcripts originating around  
641 the P<sub>105</sub> (A), P<sub>520</sub> (B), P<sub>811</sub> (C) and P<sub>3000</sub> (D) promoters were identified in wild type  
642 and  $\Delta$ CTCF-HPV18 cells in undifferentiated (blue and red, respectively) and  
643 methylcellulose differentiated (green and purple, respectively) cultures. Relevant  
644 HPV18 genome features are shown alongside each panel.

645

646 **Figure 5: CTCF regulates efficient differentiation-dependent HPV18 late**  
647 **promoter activation.**

648 The TSS of each viral transcript was identified and the distribution shown in violin  
649 plots in (A) undifferentiated and (B) differentiated keratinocytes containing WT  
650 HPV18 (blue and green, respectively) and  $\Delta$ CTCF-HPV18 (red and purple,  
651 respectively) episomes. Data distribution are shown by the kernel shape and median  
652 indicated with a vertical solid line. The wider sections of the violin plot indicate a high  
653 probability of TSS usage within that region of the HPV18 genome. The shape of the  
654 distribution indicates the concentration of data points in a particular region; the  
655 steeper the side of each bubble indicates a greater concentration of data points. ns,  
656 not significant; \*\*\*\*p<0.0001 (Fisher's test).

657

658 **Figure 6: Keratinocyte differentiation induces increased H4Ac abundance at**  
659 **the HPV18 late promoter in wild type but not  $\Delta$ CTCF-HPV18 genome-containing**  
660 **cells.** HPV18 WT and HPV18- $\Delta$ CTCF genome-containing primary keratinocytes  
661 grown in monolayer (undifferentiated; blue and green, respectively) or differentiated  
662 in methylcellulose-containing media for 48 hrs (green and purple, respectively).

663 Enrichment of H4Ac was assessed by ChIP-qPCR. Each bar in the chart represents  
664 the mid-point for primer pairs used to amplify immunoprecipitated chromatin. Fold  
665 binding over IgG control was calculated. The data shown are the mean and standard  
666 deviation of three independent replicates. Annotation of the HPV18 LCR, promoters  
667 and ORFs is provided below.

668

669 **Figure 7: Abrogation of CTCF binding to the HPV18 genome causes a**  
670 **significant reduction in differentiation-dependent late protein abundance.** (A)

671 HPV18 genome containing keratinocytes (WT or  $\Delta$ CTCF) grown in monolayer  
672 (undifferentiated, 0h) or differentiated in methylcellulose (48h) and E1<sup>E4</sup>, involucrin  
673 (IVL) and GAPDH protein expression analysed by Western blotting. Molecular weight  
674 markers are indicated on the left (kDa). (B) Relative E1<sup>E4</sup> protein expression in  
675 comparison to GAPDH was quantified in three independent experiments by  
676 densitometry. Data are the mean +/- standard deviation. \* denotes  $p < 0.05$ . (C)  
677 E1<sup>E4</sup> (red) and L1 (green) protein abundance was analysed by indirect  
678 immunofluorescence in epithelia derived from wild type and  $\Delta$ CTCF-HPV18 genome-  
679 containing keratinocytes grown in organotypic raft culture. Cellular nuclei are shown  
680 in blue, and the basal layer indicated with white arrows. Scale bar indicates 10  $\mu$ m.

681 (D) The total number of L1 positive cells per section of three independent raft  
682 cultures grown from two independent keratinocyte donors was counted. Data show  
683 the mean +/- standard deviation. \*\*\*  $p < 0.001$ , \*\*\*\*  $p < 0.0001$ . (E) E1<sup>E4</sup> protein  
684 expression in organotypic raft cultures was assessed by Western blotting lysates  
685 harvested from three independent raft cultures alongside GAPDH loading control.  
686 Molecular weight markers are indicated on the left.

687

## 688 REFERENCES

689

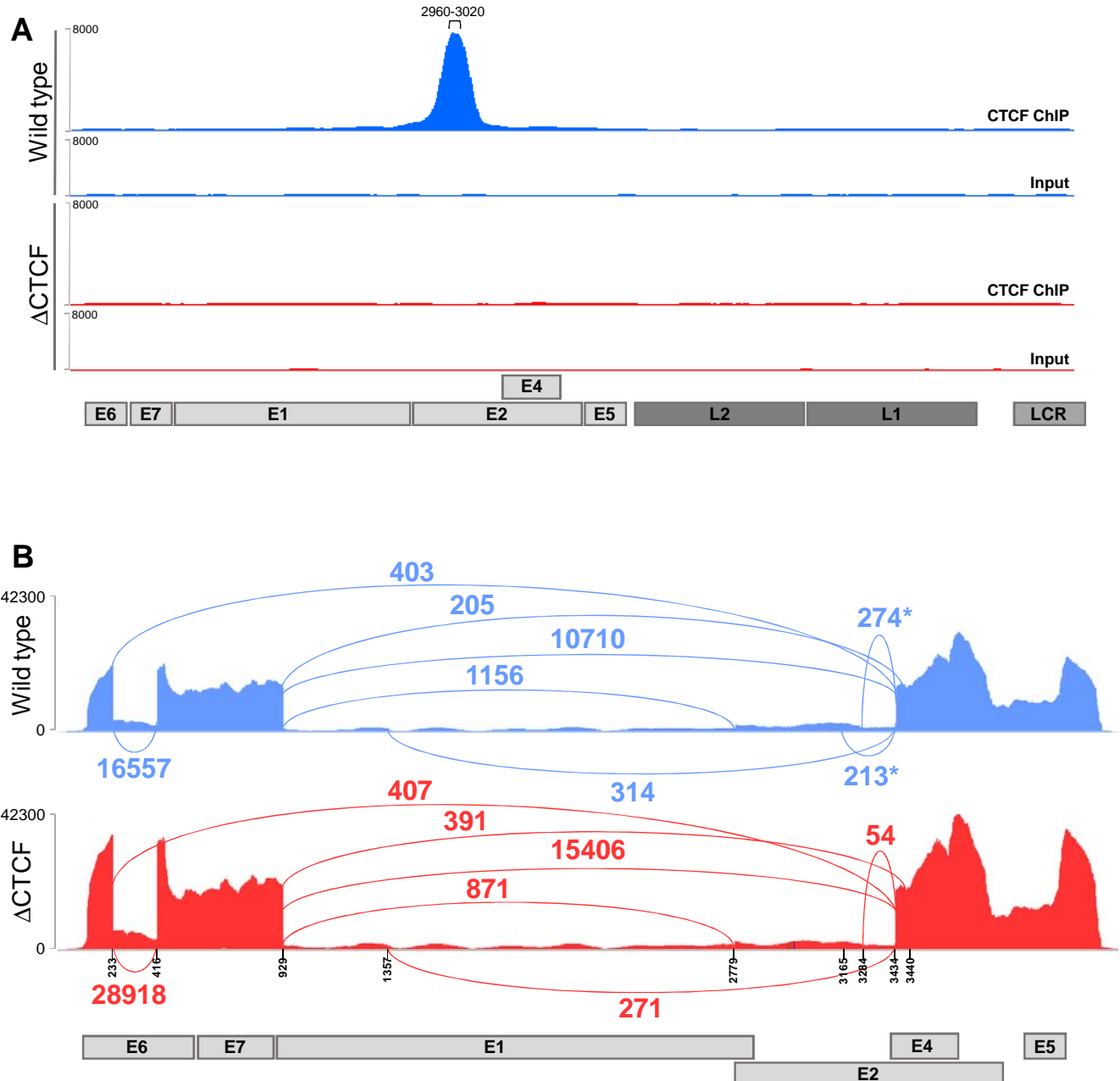
- 690 1. Tommasino M. 2014. The human papillomavirus family and its role in carcinogenesis.  
691 *Semin Cancer Biol* 26:13-21.
- 692 2. Stunkel W, Bernard HU. 1999. The chromatin structure of the long control region of  
693 human papillomavirus type 16 represses viral oncoprotein expression. *J Virol*  
694 73:1918-30.
- 695 3. Hummel M, Lim HB, Laimins LA. 1995. Human papillomavirus type 31b late gene  
696 expression is regulated through protein kinase C-mediated changes in RNA  
697 processing. *J Virol* 69:3381-8.

- 698 4. Ruesch MN, Stubenrauch F, Laimins LA. 1998. Activation of papillomavirus late gene  
699 transcription and genome amplification upon differentiation in semisolid medium is  
700 coincident with expression of involucrin and transglutaminase but not keratin-10. *J*  
701 *Virology* 72:5016-24.
- 702 5. Wooldridge TR, Laimins LA. 2008. Regulation of human papillomavirus type 31 gene  
703 expression during the differentiation-dependent life cycle through histone  
704 modifications and transcription factor binding. *Virology* 374:371-80.
- 705 6. Pentland I, Campos-Leon K, Cotic M, Davies KJ, Wood CD, Groves IJ, Burley M,  
706 Coleman N, Stockton JD, Noyvert B, Beggs AD, West MJ, Roberts S, Parish JL. 2018.  
707 Disruption of CTCF-YY1-dependent looping of the human papillomavirus genome  
708 activates differentiation-induced viral oncogene transcription. *PLoS Biol*  
709 16:e2005752.
- 710 7. Paris C, Pentland I, Groves I, Roberts DC, Powis SJ, Coleman N, Roberts S, Parish JL.  
711 2015. CCCTC-binding factor recruitment to the early region of the human  
712 papillomavirus 18 genome regulates viral oncogene expression. *J Virology* 89:4770-85.
- 713 8. Beagan JA, Duong MT, Titus KR, Zhou L, Cao Z, Ma J, Lachanski CV, Gillis DR, Phillips-  
714 Cremins JE. 2017. YY1 and CTCF orchestrate a 3D chromatin looping switch during  
715 early neural lineage commitment. *Genome Res* doi:10.1101/gr.215160.116.
- 716 9. del Mar Pena LM, Laimins LA. 2001. Differentiation-dependent chromatin  
717 rearrangement coincides with activation of human papillomavirus type 31 late gene  
718 expression. *J Virology* 75:10005-13.
- 719 10. Burley M, Roberts S, Parish JL. 2020. Epigenetic regulation of human papillomavirus  
720 transcription in the productive virus life cycle. *Semin Immunopathol* 42:159-171.
- 721 11. Grassmann K, Rapp B, Maschek H, Petry KU, Iftner T. 1996. Identification of a  
722 differentiation-inducible promoter in the E7 open reading frame of human  
723 papillomavirus type 16 (HPV-16) in raft cultures of a new cell line containing high  
724 copy numbers of episomal HPV-16 DNA. *J Virology* 70:2339-49.
- 725 12. Hummel M, Hudson JB, Laimins LA. 1992. Differentiation-induced and constitutive  
726 transcription of human papillomavirus type 31b in cell lines containing viral  
727 episomes. *J Virology* 66:6070-80.
- 728 13. Wilson R, Fehrmann F, Laimins LA. 2005. Role of the E1--E4 protein in the  
729 differentiation-dependent life cycle of human papillomavirus type 31. *J Virology*  
730 79:6732-40.
- 731 14. Peh WL, Brandsma JL, Christensen ND, Cladel NM, Wu X, Doorbar J. 2004. The viral  
732 E4 protein is required for the completion of the cottontail rabbit papillomavirus  
733 productive cycle in vivo. *J Virology* 78:2142-51.
- 734 15. Bodily JM, Meyers C. 2005. Genetic analysis of the human papillomavirus type 31  
735 differentiation-dependent late promoter. *J Virology* 79:3309-21.
- 736 16. Songcock WK, Scott ML, Bodily JM. 2017. Regulation of the human papillomavirus  
737 type 16 late promoter by transcriptional elongation. *Virology* 507:179-191.
- 738 17. Wang X, Meyers C, Wang HK, Chow LT, Zheng ZM. 2011. Construction of a full  
739 transcription map of human papillomavirus type 18 during productive viral infection.  
740 *J Virology* 85:8080-92.
- 741 18. Toots M, Mannik A, Kivi G, Ustav M, Jr., Ustav E, Ustav M. 2014. The transcription  
742 map of human papillomavirus type 18 during genome replication in U2OS cells. *PLoS*  
743 *One* 9:e116151.



- 744 19. Mole S, McFarlane M, Chuen-Im T, Milligan SG, Millan D, Graham SV. 2009. RNA  
745 splicing factors regulated by HPV16 during cervical tumour progression. *J Pathol*  
746 219:383-91.
- 747 20. McFarlane M, MacDonald AI, Stevenson A, Graham SV. 2015. Human Papillomavirus  
748 16 Oncoprotein Expression Is Controlled by the Cellular Splicing Factor SRSF2 (SC35).  
749 *J Virol* 89:5276-87.
- 750 21. Shukla S, Kavak E, Gregory M, Imashimizu M, Shutinoski B, Kashlev M, Oberdoerffer  
751 P, Sandberg R, Oberdoerffer S. 2011. CTCF-promoted RNA polymerase II pausing  
752 links DNA methylation to splicing. *Nature* 479:74-79.
- 753 22. Lopez Soto EJ, Lipscombe D. 2020. Cell-specific exon methylation and CTCF binding in  
754 neurons regulate calcium ion channel splicing and function. *Elife* 9.
- 755 23. Agirre E, Bellora N, Allo M, Pages A, Bertucci P, Kornblihtt AR, Eyras E. 2015. A  
756 chromatin code for alternative splicing involving a putative association between  
757 CTCF and HP1alpha proteins. *BMC Biol* 13:31.
- 758 24. Ruiz-Velasco M, Kumar M, Lai MC, Bhat P, Solis-Pinson AB, Reyes A, Kleinsorg S, Noh  
759 KM, Gibson TJ, Zaugg JB. 2017. CTCF-Mediated Chromatin Loops between Promoter  
760 and Gene Body Regulate Alternative Splicing across Individuals. *Cell Syst* 5:628-637  
761 e6.
- 762 25. Goodwin S, McPherson JD, McCombie WR. 2016. Coming of age: ten years of next-  
763 generation sequencing technologies. *Nature Reviews Genetics* 17:333-351.
- 764 26. van Dijk EL, Jaszczyszyn Y, Naquin D, Thermes C. 2018. The Third Revolution in  
765 Sequencing Technology. *Trends in Genetics* 34:666-681.
- 766 27. Wilson R, Laimins LA. 2005. Differentiation of HPV-containing cells using organotypic  
767 "raft" culture or methylcellulose. *Methods Mol Med* 119:157-69.
- 768 28. Roberts S, Hillman ML, Knight GL, Gallimore PH. 2003. The ND10 Component  
769 Promyelocytic Leukemia Protein Relocates to Human Papillomavirus Type 1 E4  
770 Intranuclear Inclusion Bodies in Cultured Keratinocytes and in Warts. *Journal of*  
771 *Virology* 77:673-684.
- 772 29. Günther T, Fröhlich J, Herrde C, Ohno S, Burkhardt L, Adler H, Grundhoff A. 2019. A  
773 comparative epigenome analysis of gammaherpesviruses suggests cis-acting  
774 sequence features as critical mediators of rapid polycomb recruitment. *PLOS*  
775 *Pathogens* 15:e1007838.
- 776 30. Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient  
777 alignment of short DNA sequences to the human genome. *Genome Biol* 10:R25.
- 778 31. Günther T, Grundhoff A. 2010. The epigenetic landscape of latent Kaposi sarcoma-  
779 associated herpesvirus genomes. *PLoS Pathog* 6:e1000935.
- 780 32. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,  
781 Gingeras TR. 2012. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29:15-  
782 21.
- 783 33. Thompson S, Thompson S, Cazier J. 2019. CaStLeS (Compute and Storage for the Life  
784 Sciences): a collection of compute and storage resources for supporting research at  
785 the University of Birmingham. Zenodo.
- 786 34. Schwenzer H, Abdel Mouti M, Neubert P, Morris J, Stockton J, Bonham S,  
787 Fellermeier M, Chettle J, Fischer R, Beggs AD, Blagden SP. 2021. LARP1 isoform  
788 expression in human cancer cell lines. *RNA Biol* 18:237-247.
- 789 35. Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*  
790 34:3094-3100.

- 791 36. Reynolds G, Deshmukh NS, Mangham D. 2000. Agitated low temperature epitope  
792 retrieval (ALTER): Effective antigen retrieval for immunohistochemistry with  
793 excellent morphological preservation. *The Journal of Pathology* 190:51A-51A.
- 794 37. Mehta K, Gunasekharan V, Satsuka A, Laimins LA. 2015. Human papillomaviruses  
795 activate and recruit SMC1 cohesin proteins for the differentiation-dependent life  
796 cycle through association with CTCF insulators. *PLoS Pathog* 11:e1004763.
- 797 38. Tan CL, Gunaratne J, Lai D, Carthage L, Wang Q, Xue YZ, Quek LS, Doorbar J,  
798 Bachelier F, Thierry F, Bellanger S. 2012. HPV-18 E2<sup>E4</sup> chimera: 2 new spliced  
799 transcripts and proteins induced by keratinocyte differentiation. *Virology* 429:47-56.
- 800 39. Smits P, Poumay Y, Karperien M, Tylzanowski P, Wauters J, Huylebroeck D, Ponc M,  
801 Merregaert J. 2000. Differentiation-dependent alternative splicing and expression of  
802 the extracellular matrix protein 1 gene in human keratinocytes. *J Invest Dermatol*  
803 114:718-24.
- 804 40. Donovan-Banfield Ia, Turnell AS, Hiscox JA, Leppard KN, Matthews DA. 2020. Deep  
805 splicing plasticity of the human adenovirus type 5 transcriptome drives virus  
806 evolution. *Communications Biology* 3:124.
- 807 41. LeRoy G, Rickards B, Flint SJ. 2008. The double bromodomain proteins Brd2 and Brd3  
808 couple histone acetylation to transcription. *Mol Cell* 30:51-60.
- 809 42. Spink KM, Laimins LA. 2005. Induction of the human papillomavirus type 31 late  
810 promoter requires differentiation but not DNA amplification. *J Virol* 79:4918-26.
- 811 43. Shukla S, Kavak E, Gregory M, Imashimizu M, Shutinoski B, Kashlev M, Oberdoerffer  
812 P, Sandberg R, Oberdoerffer S. 2011. CTCF-promoted RNA polymerase II pausing  
813 links DNA methylation to splicing. *Nature* 479:74-9.
- 814



**Figure 1**

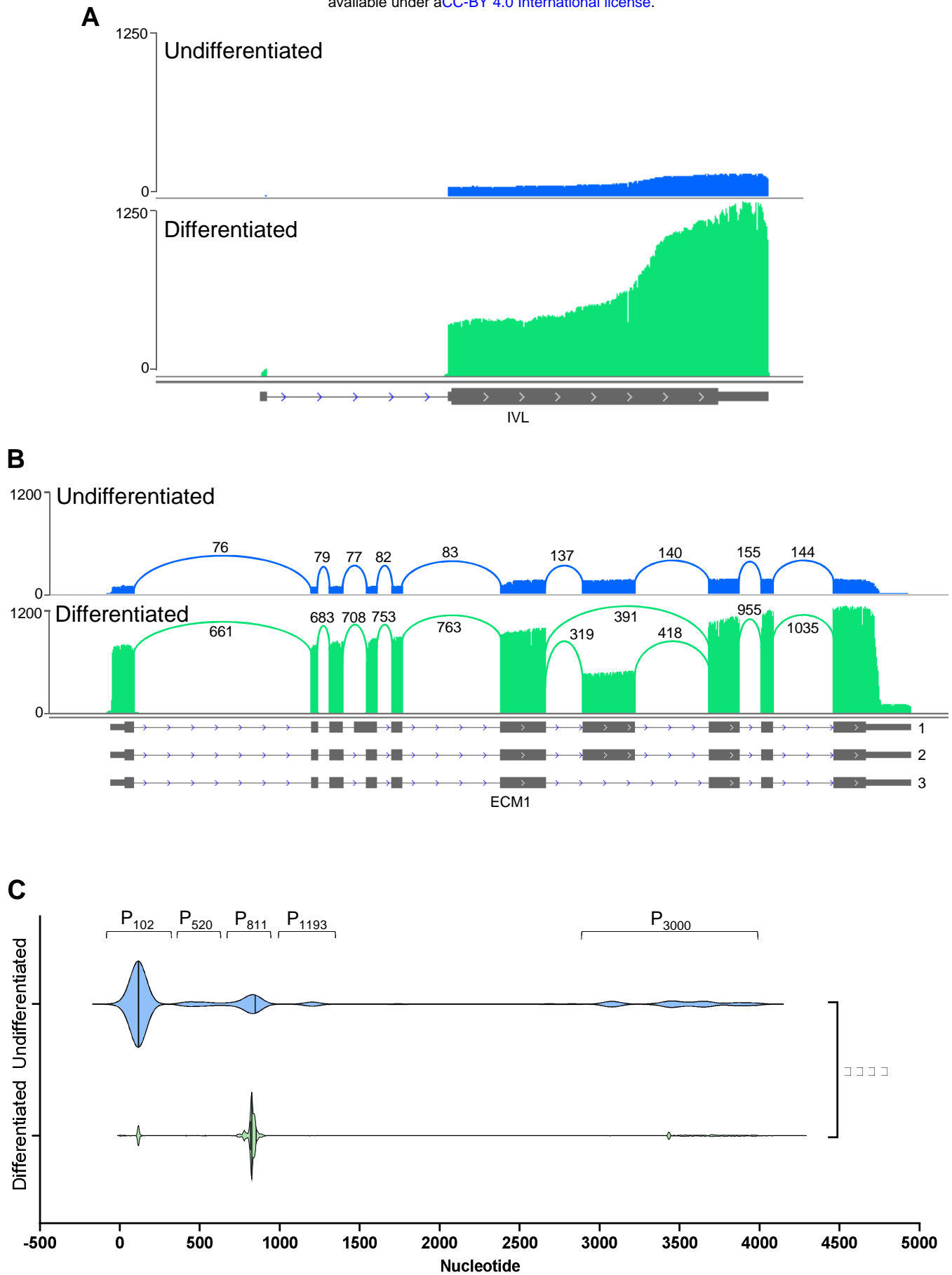


Figure 2

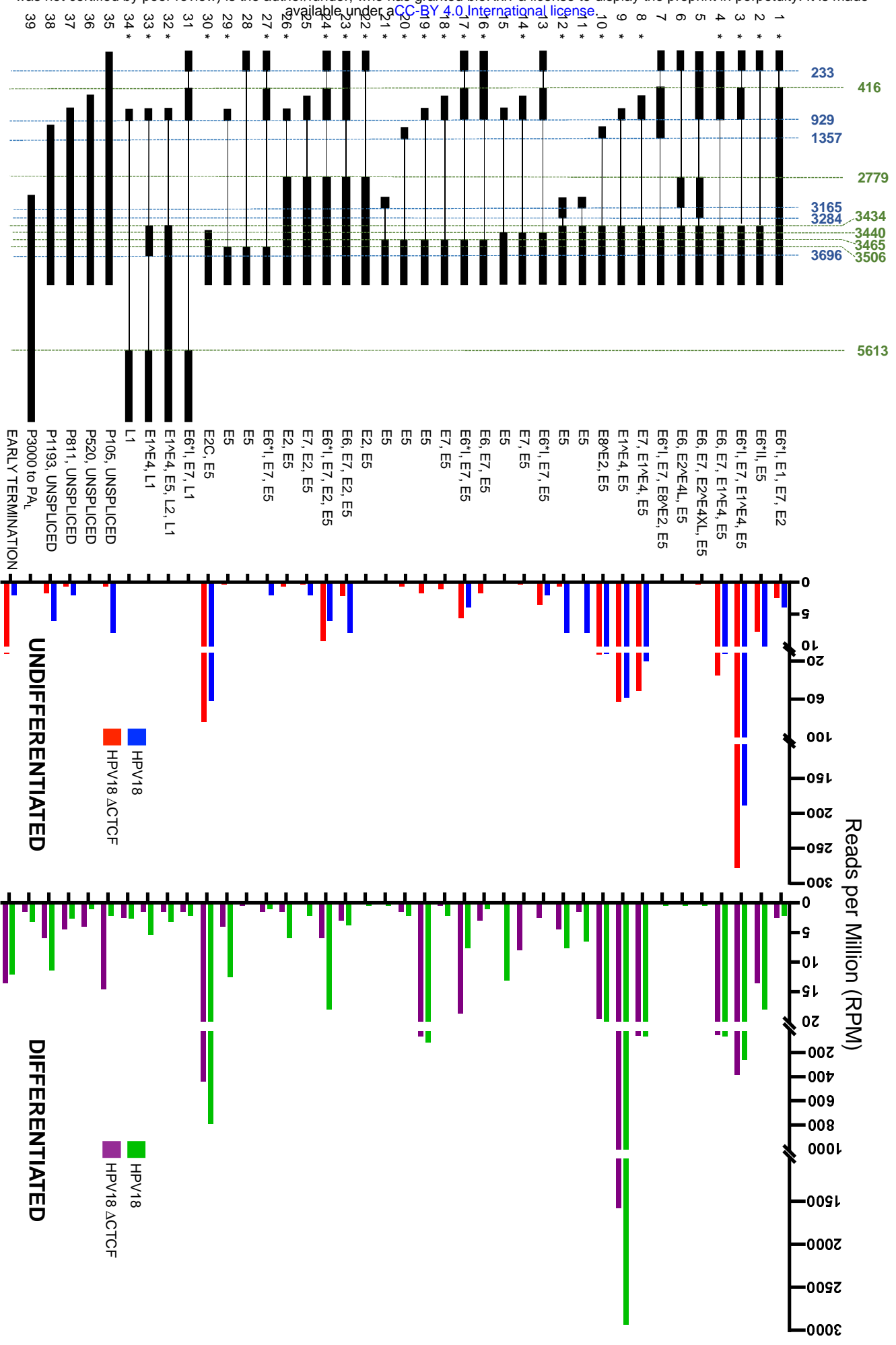
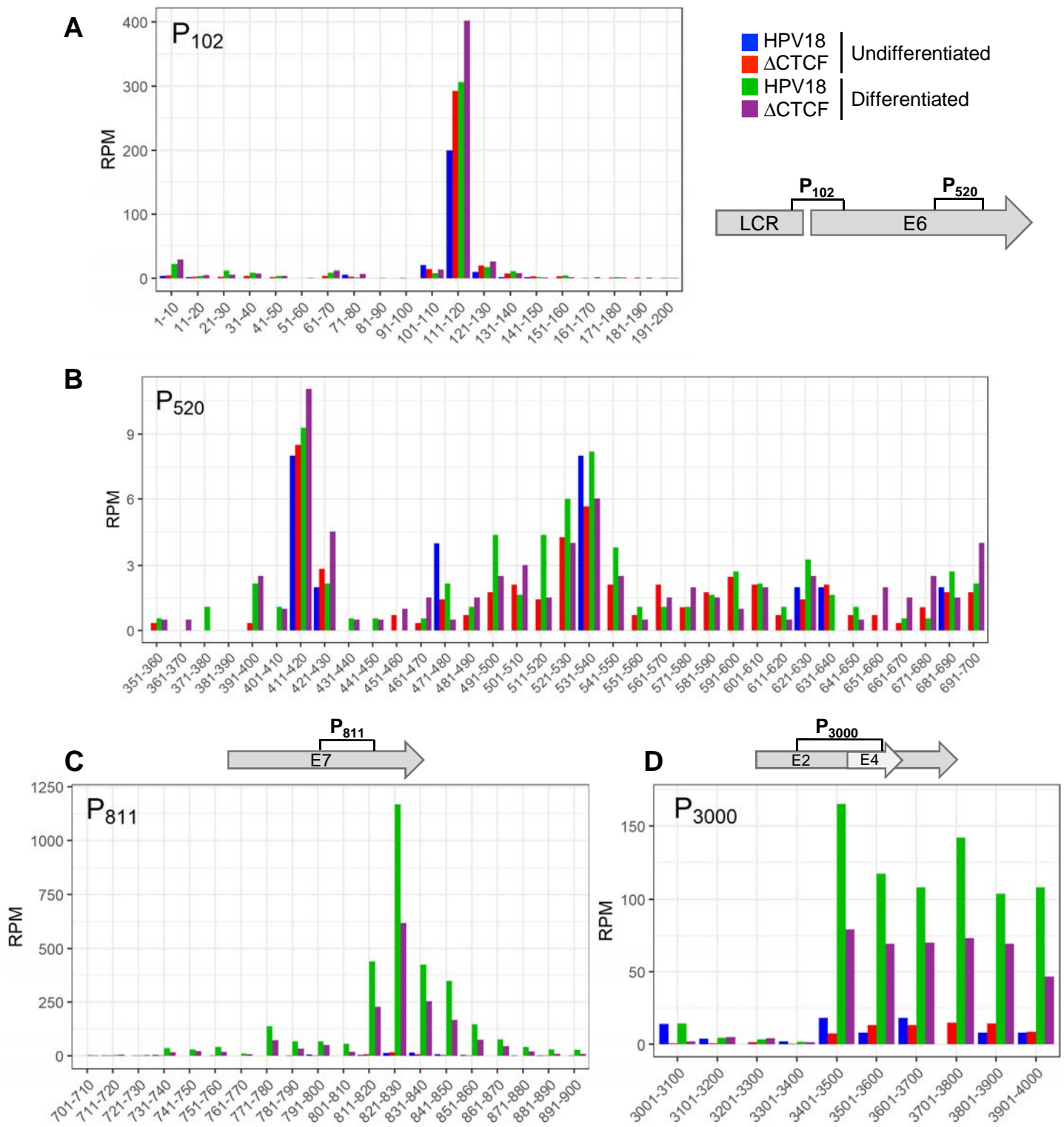
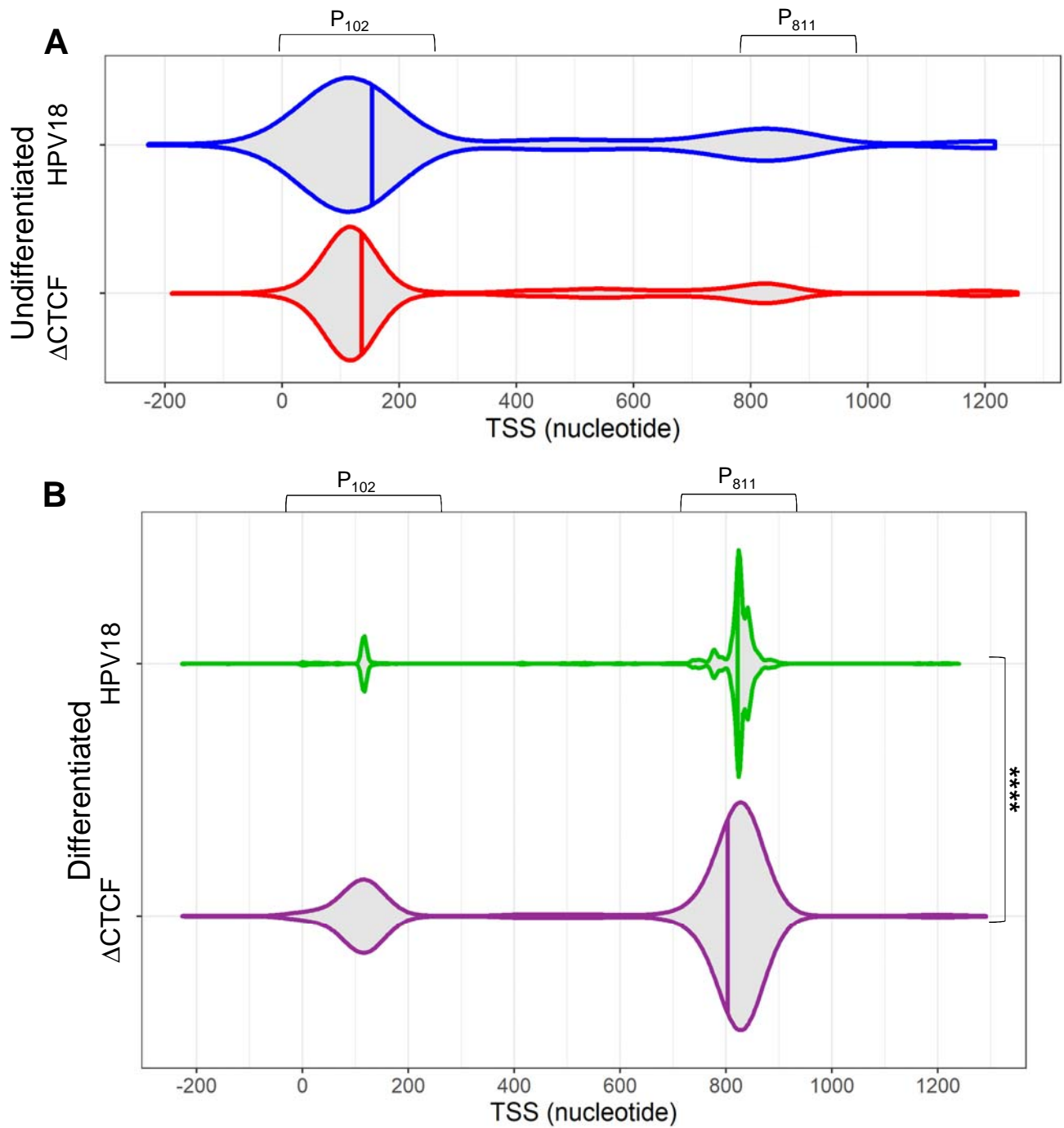


Figure 3

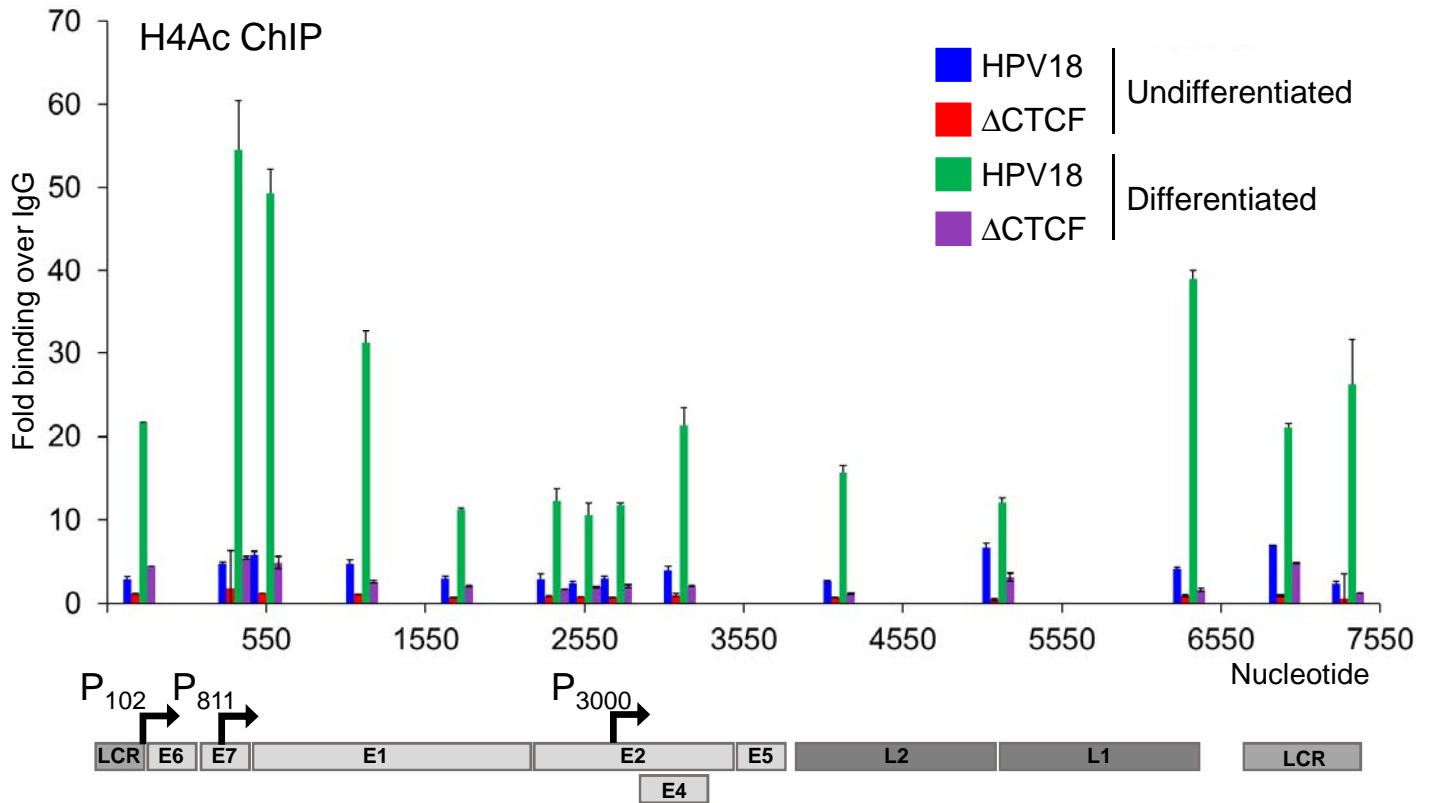


**Figure 4**



**Figure 5**





**Figure 6**

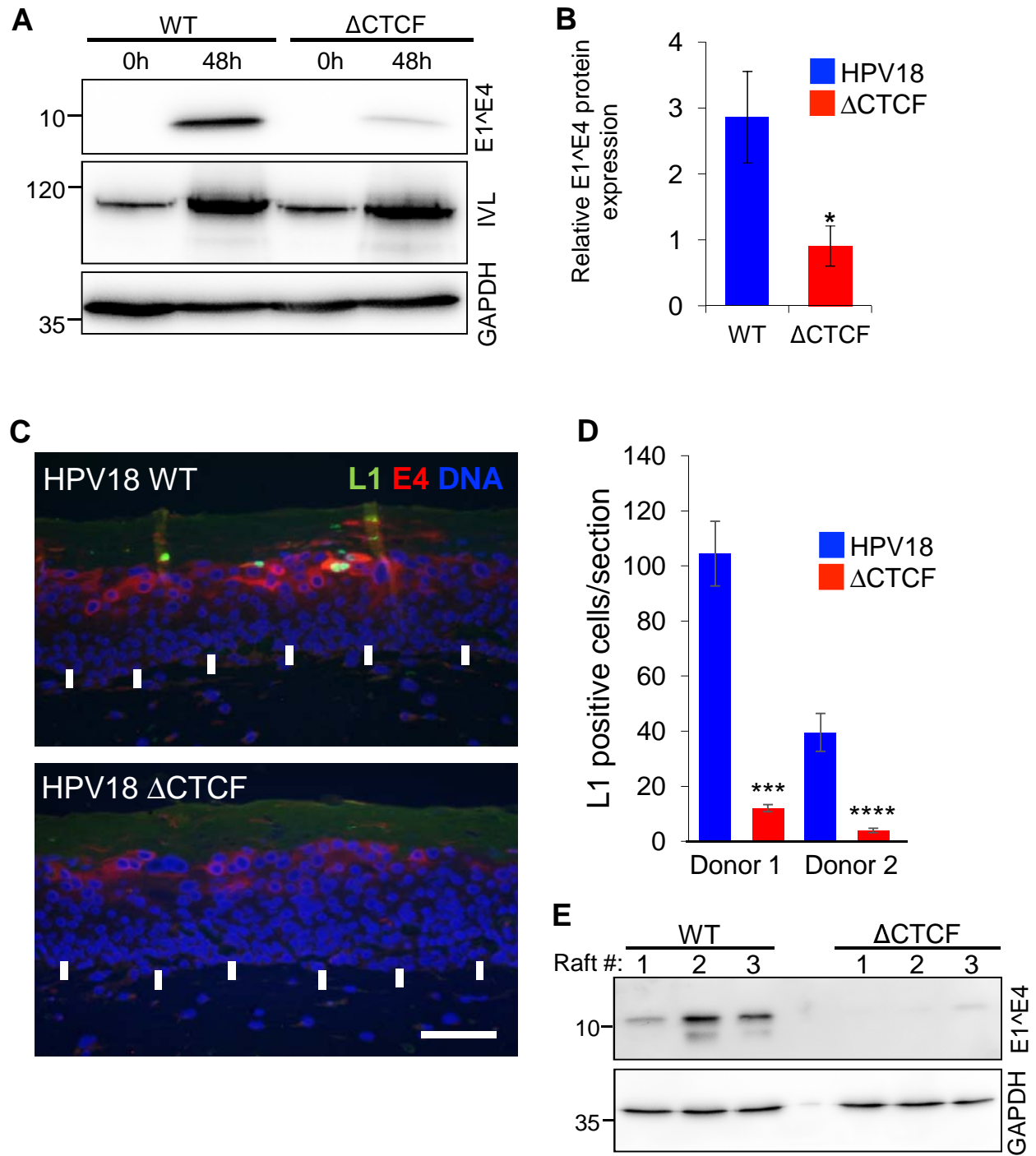


Figure 7