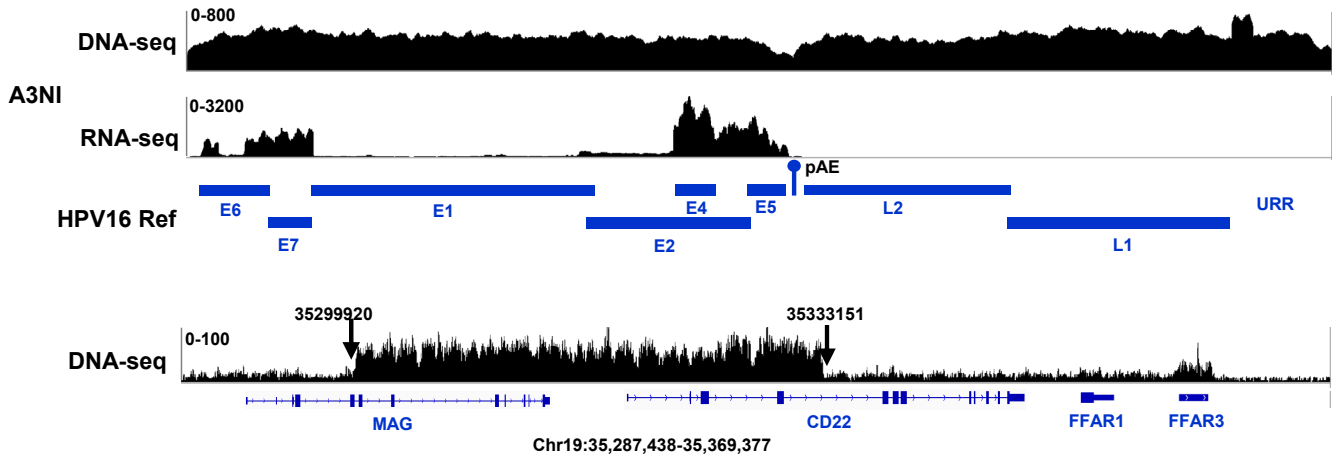
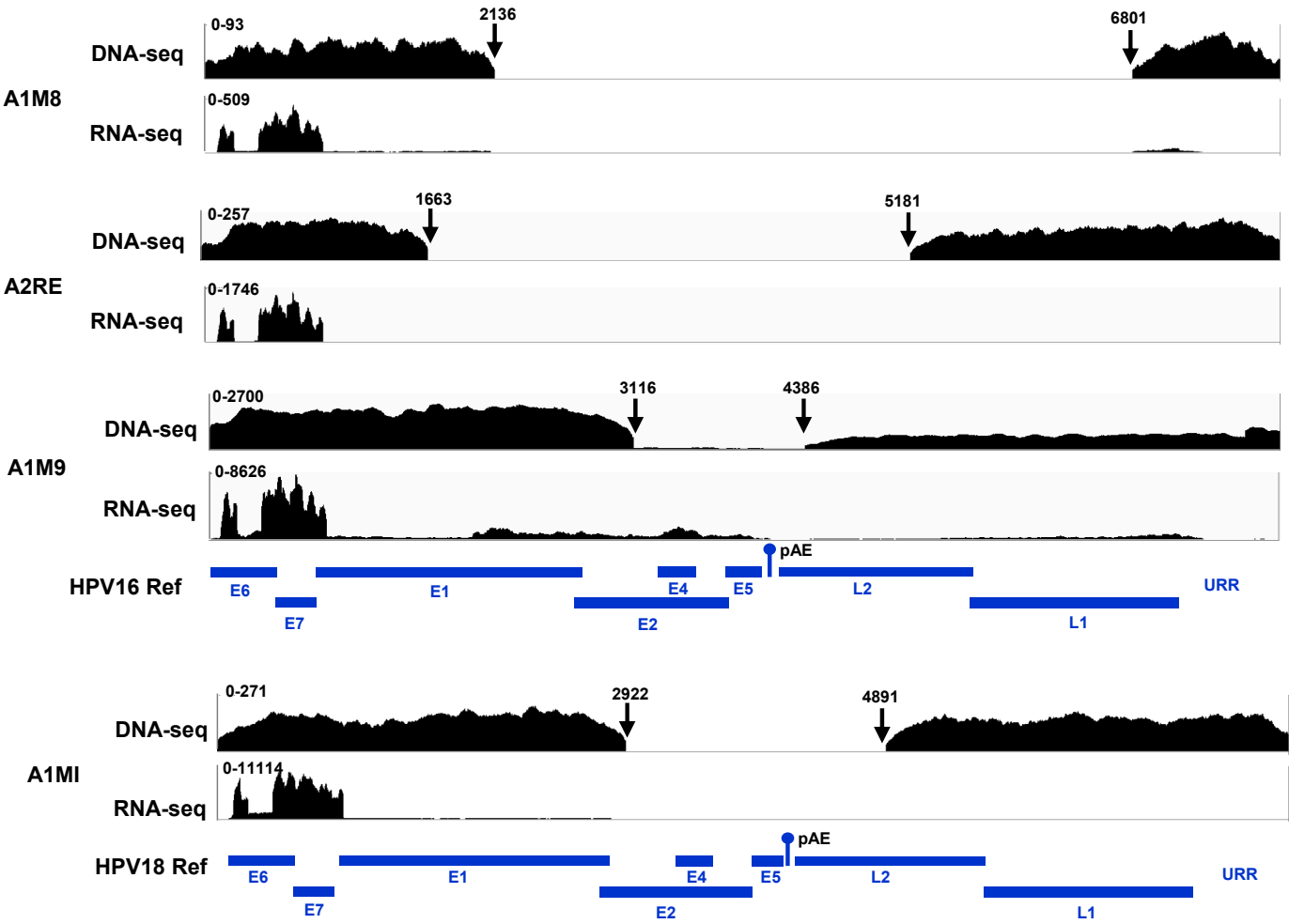
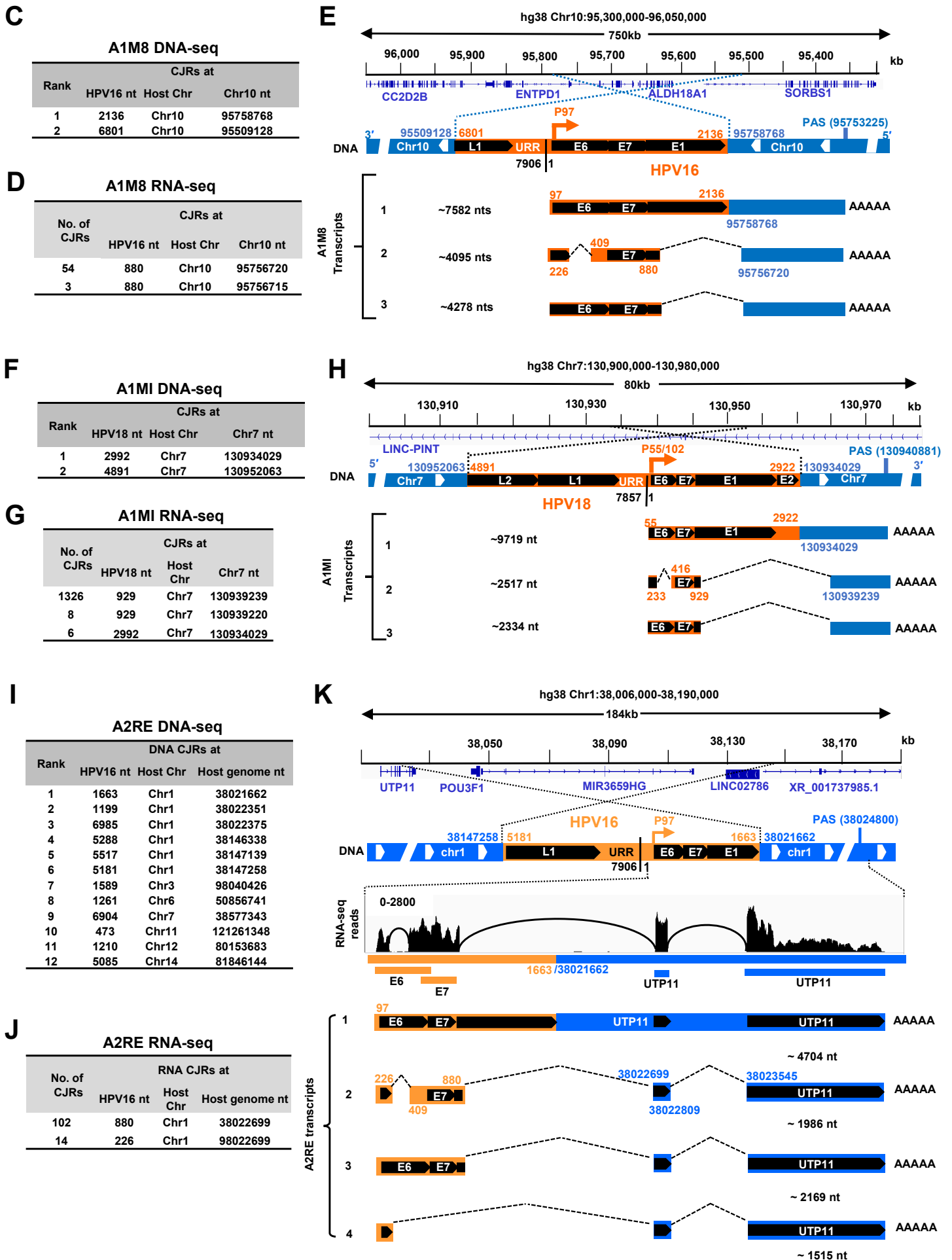


A



B



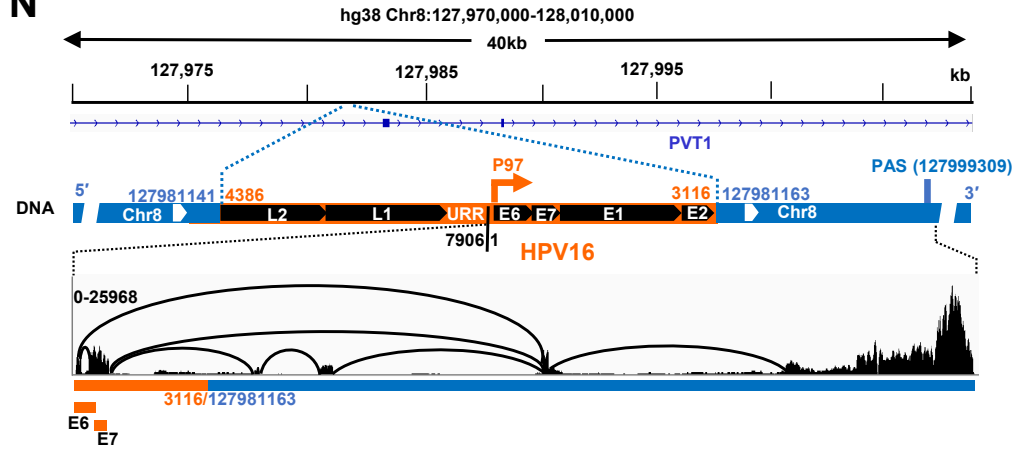


L

A1M9 DNA-seq

Rank	CJRs at		
	HPV16 nt	Host Chr	Chr nt
1	4386	Chr8	127981141
2	3116	Chr8	127981163
3	1749	Chr8	127983497
4	7547	Chr8	128000542
5	1785	Chr8	128001577
6	1684	Chr8	128005189
7	2068	Chr4	119703111
8	7614	Chr12	29998504
9	5495	Chr17	39718111
10	2448	Chr20	17390600

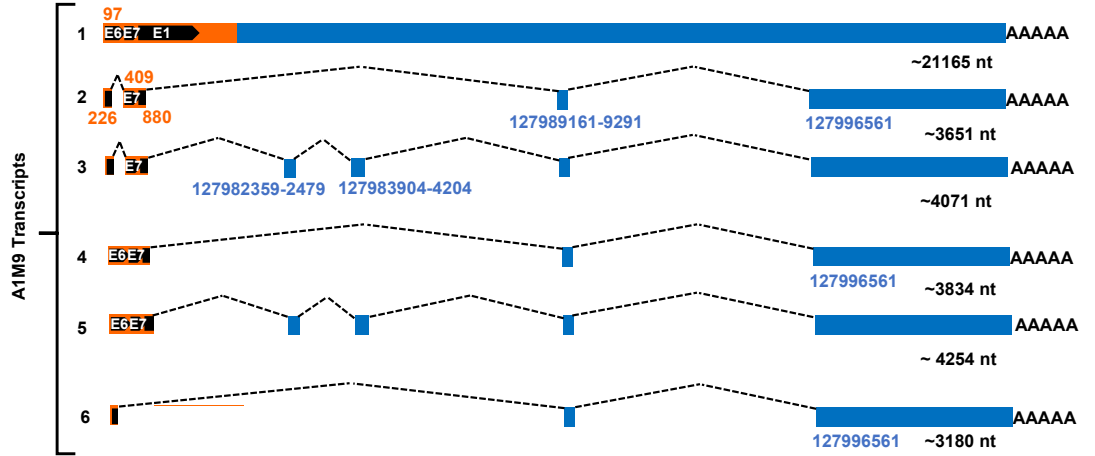
N



M

A1M9 RNA-seq

No. of CJRs	CJRs at		
	HPV16 nt	Host Chr	Chr8 nt
717	880	Chr8	127989161
262	226	Chr8	127989161
127	880	Chr8	127982359
20	3116	Chr8	127981163
2	4386	Chr8	127981141

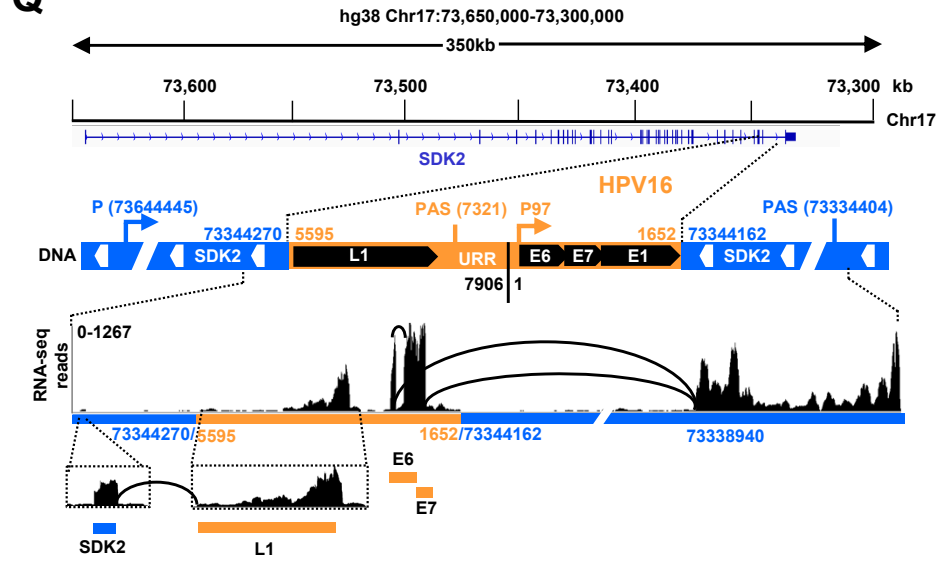


O

T1074 targeted DNA-seq

Rank	Host Chr	No. of mapped CJR locations
1	Chr3	22
2	Chr5	22
3	Chr2	16
4	Chr4	16
5	ChrX	13
6	Chr1	11
7	Chr7	11
8	Chr13	11
9	Chr18	11
10	Chr6	9
11	Chr8	9
12	Chr14	8
13	Chr10	7
14	Chr11	6
15	Chr15	6
16	Chr16	6
17	Chr9	5
18	Chr12	5
19	Chr17	5
20	Chr19	5
21	Chr20	3
22	Chr21	3
23	Chr22	3

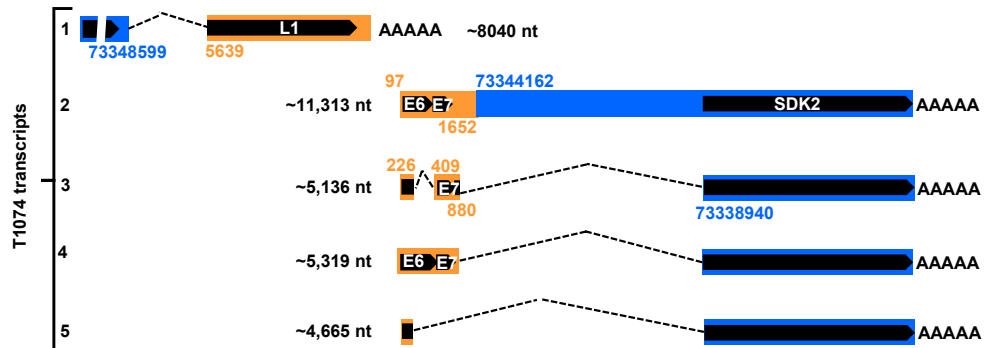
Q



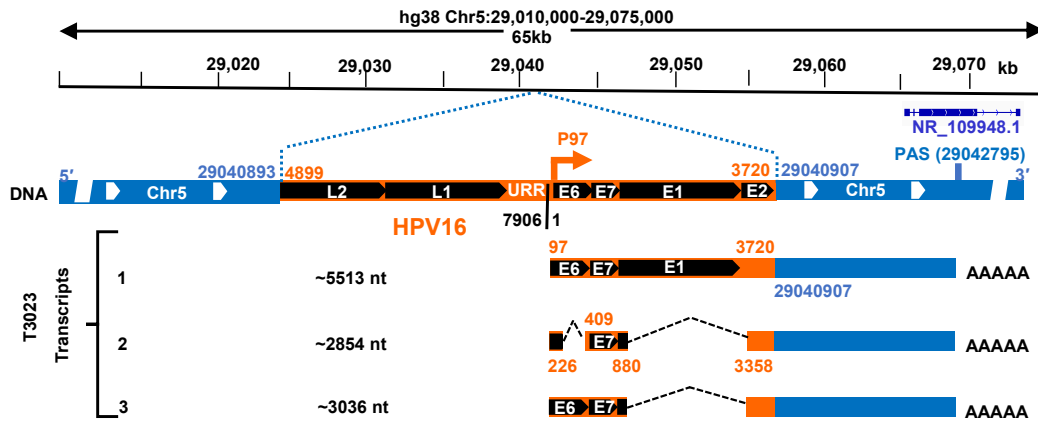
P

T1074 RNA-seq

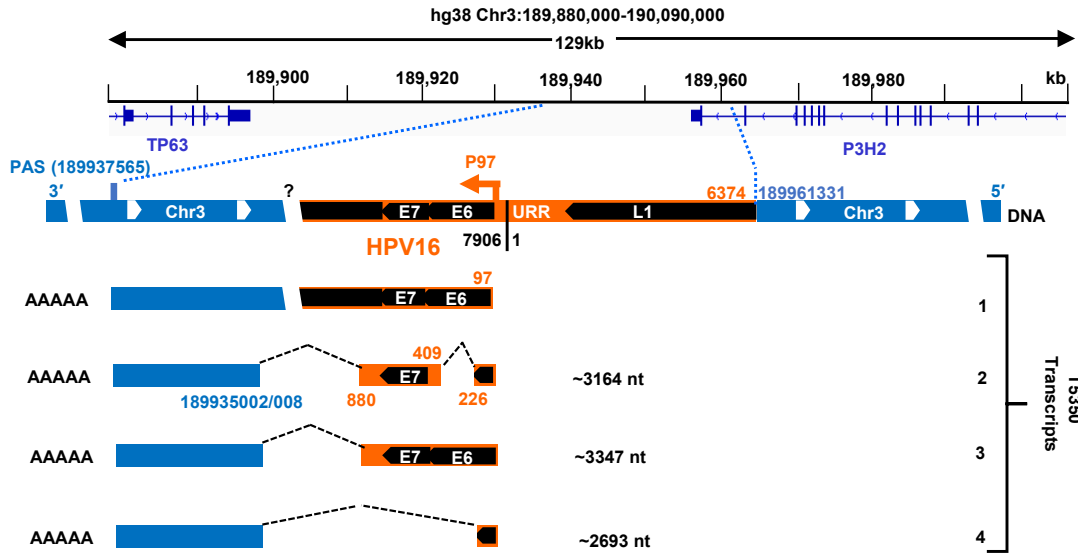
No. of CJRs	RNA CJRs at		
	HPV16 nt	Host Chr	Host genome nt
86	880	Chr17	73338940
22	226	Chr17	73338940
9	5639	Chr17	73348599



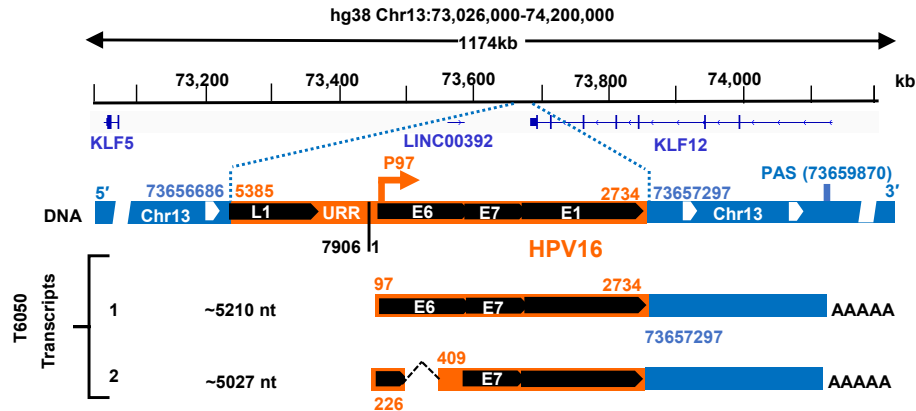
R



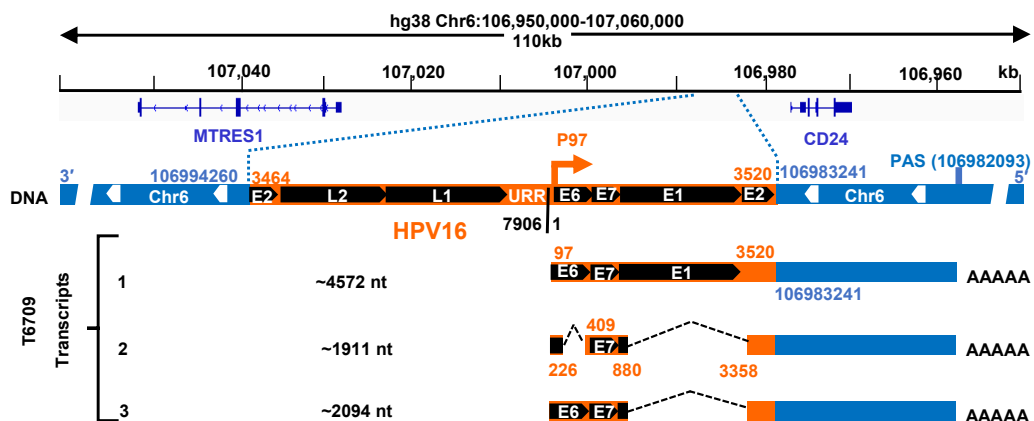
S



T



U



V

T6570 targeted DNA-seq

Order	DNA CJRs at		
	HPV18 nt	Host chr	Host chr nt
1	3631	Chr13	73350892
2	3455	Chr13	73658495
3	6264	Chr19	30548314
4	6864	Chr2	227416717

W

T6570 RNA-seq

No. of CJRs	RNA CJRs at		
	HPV18 nt	Host Chr	Host chr nt
67	3631	Chr13	73350892
50	929	Chr13	73648631
2	3284	Chr13	73655946
2	929	Chr13	73655946

X

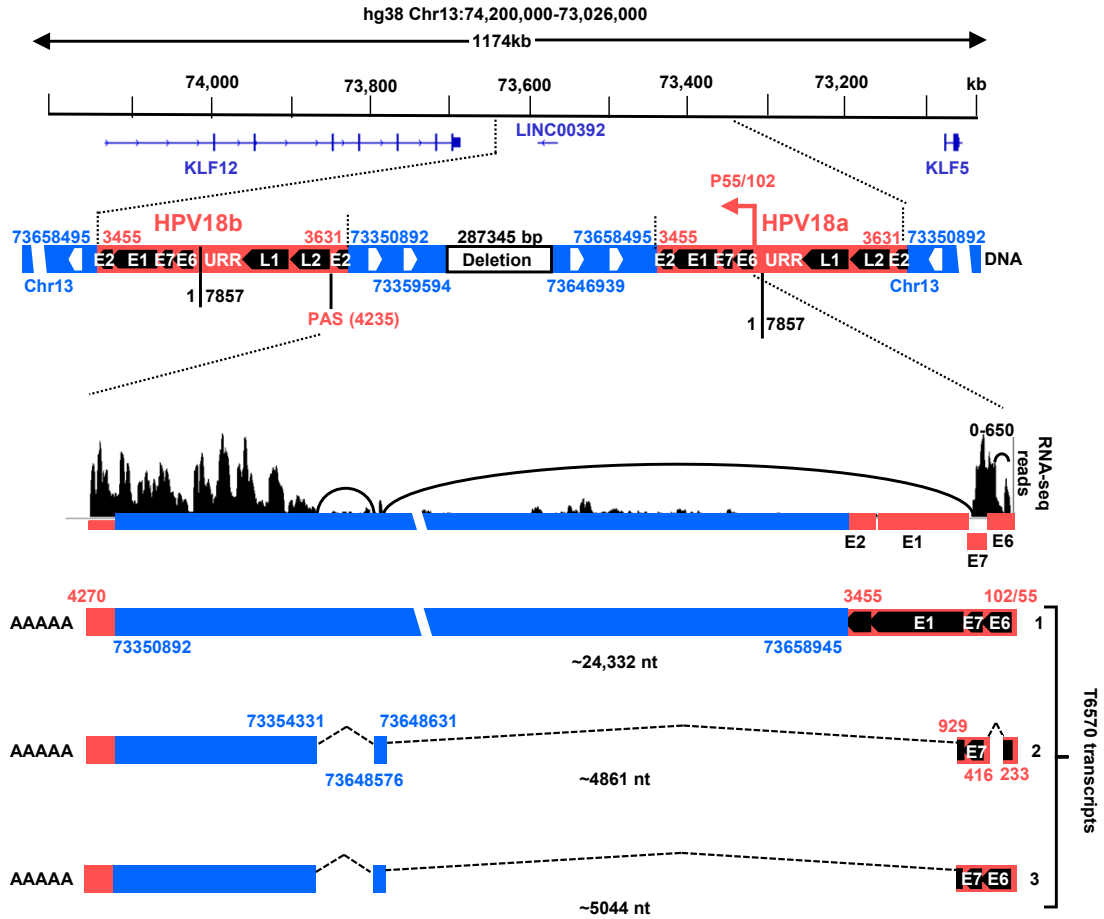
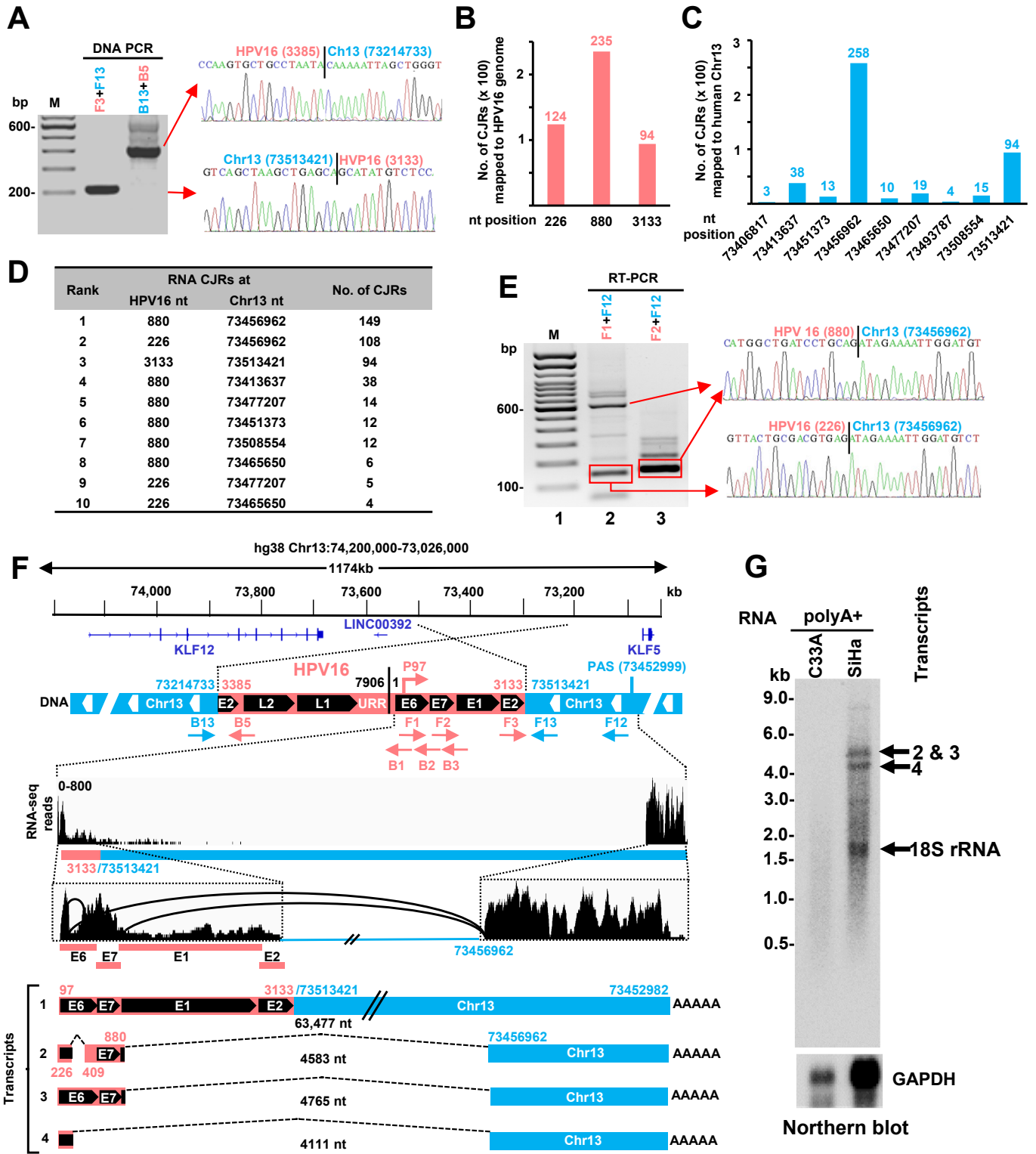


Figure S1. Identification of the mapped chromosomal HPV16 or HPV18 DNA integration sites expressing a bicistronic E6E7 RNA in cervical cancer tissues. (A, B) Coverage and distribution of HPV16- or HPV18-specific reads along with the viral genome in the indicated cervical cancer tissues. HPV16 DNA integration in the A3NI induced the genomic amplification of a Chr19 region covering MAG and CD22 genes. Arrows indicate the integration-mediated genome amplification or deletion. (C-H) E6E7 expression in the tissues A1M8 (HPV16, C-E) and A1MI (HPV18, F-H) bearing only a single integrated HPV DNA copy. The CJRs identified by both DNA-seq and RNA-seq were mapped to the HPV16 DNA integrated in the region between ENTPD1 and SORBS1 genes on Chr10 in the A1M8 (E) and the HPV18 DNA in the intron of LINC-PINT on Chr7 in the A1MI (H). The truncated viral DNA integrated in a reversed orientation expressed viral E6E7 RNA by using a viral early promoter P97 (HPV16) or P55/102 (HPV18) and a host polyadenylation signal (PAS) at Chr10 nt 95,753,225 in the A1M8 (E) and Chr7 nt 130,940,881 in the A1MI (H) to produce three virus-host chimeric transcripts from alternative RNA splicing in each tissue. (I-U) E6E7 expression selectively from a single integrated HPV16 DNA copy in cervical cancer tissues bearing integrated HPV16 DNAs on multiple chromosomes (Figure 1, Supplemental Tables S1-3). A bicistronic E6E7 RNA in each tissue was transcribed from a viral early promoter P97 of the integrated HPV16 DNA, alternatively spliced to produce multiple RNA isoforms, and polyadenylated by using a host PAS accessible downstream of the integrated viral DNA. A number below or on the right or left end of the isoform RNA indicates the estimated length of the RNA. URR, upstream regulatory region; nt 7906/1, tail-to-head junction. RNA-seq reads-coverage and Sashimi plots for splice junctions are visualized by IGV. In the A2RE tissue (I-K), DNA CJRs were mapped to multiple host chromosomes (I) but RNA CJRs only to a single HPV16 DNA integration site on Chr1 (J). The truncated HPV16 DNA was integrated in a reversed orientation in a region between UTP11 and AL139158.3 genes on Chr1 (K). The E6E7 RNA in the tissue expressed from the integrated HPV16 DNA was polyadenylated by using a host PAS at Chr1 nt 38,024,800 to produce four chimeric RNA transcripts. In the A1M9 tissue (L-N), DNA CJRs were mapped to multiple chromosomes (L) but the RNA CJRs only to a single chromosomal HPV DNA integration site on Chr8 (M). The integrated HPV16 DNA in the intron region of PVT1 on Chr8 expressed the E6E7 RNA by using a host PAS at nt 127,999,309 to produce six virus-host chimeric transcripts (N). In the T1074 tissue (O-Q), DNA CJRs were mapped to multiple host chromosomes (O) but RNA CJRs only to the integrated HPV16 DNA on Chr17 (P). The truncated HPV16 DNA integrated in the last intron of the SDK2 gene on Chr17 (Q) expressed the E6E7 RNA by using a host PAS at Chr17 nt 73334404 to produce five chimeric virus-SDK2 RNA isoforms. In the tissues T3023 (R), T5350 (S), T6050 (T), and T6709 (U) with multiple chromosomal HPV16 DNA integration sites, E6E7 RNA in the T3023 (R) was expressed only from the integrated HPV16 DNA in the intron region of XR_001742622 gene on Chr5 and polyadenylated by using a host PAS at Chr5 nt 29,042,795 to produce three virus-host chimeric transcripts; in the tissue 5350 (S), expressed only from the integrated HPV16 DNA in an intergenic region between TP63 and P3H2 on Chr3 and terminated at a host PAS at Chr3 nt 189,937,565 to produce four virus-host chimeric transcripts; in the tissue T6050 (T), expressed from the integrated HPV16 DNA in an intergenic region between KLF5 and KLF12 on Chr13 and terminated at a host PAS at Chr3 nt 73,659,870 to produce two virus-host chimeric transcripts; and in the tissue T6709 (U), expressed from the integrated HPV16 DNA in an intergenic region between MTRES1 and CD24 on Chr6 and terminated at a host PAS at Chr6 nt 106,982,093 to produce three virus-host chimeric transcripts. (V-X) E6E7 expression selectively from a single integrated HPV18 DNA in the tissue T6570. DNA CJRs were mapped to multiple host chromosomes (V) but RNA CJRs only to a single HPV18 DNA integration site on Chr13 (W). (X) Two tandem copies of the truncated HPV18 DNA (HPV18a and HPV18b) with 176 bp deletion in the E2 region were integrated in a rearranged intergenic region between KLF12 and KLF5 on Chr13 but separated by a ~20-kb host DNA remnant from a large (~287-kb) deletion. Viral E6E7 RNA transcribed from the viral early promoter (P55/ P102) in the HPV18a and terminated at the viral early PAS at nt 4235 in the HPV18b to produce three virus-host chimeric transcripts.



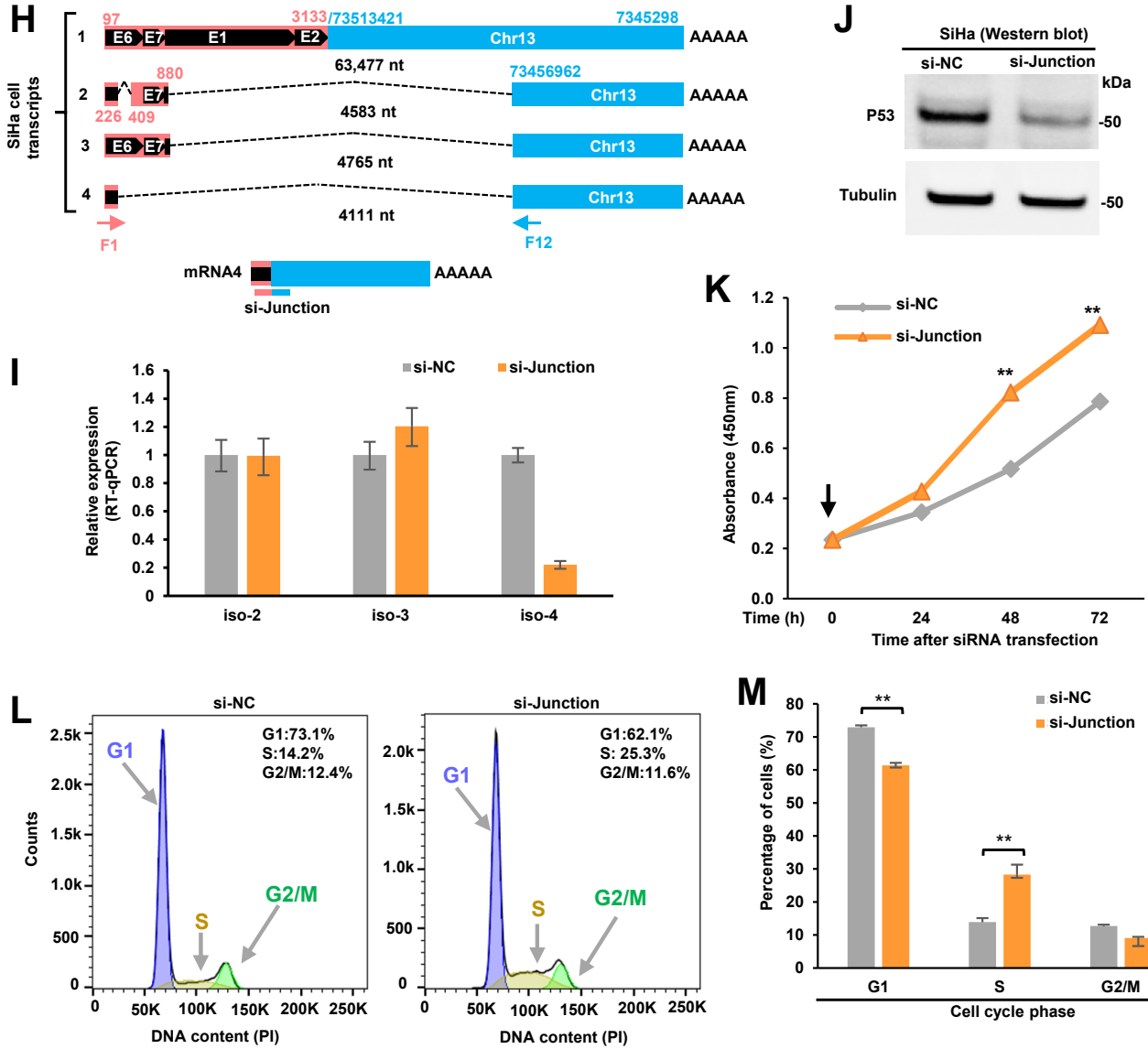


Figure S2. Bicistronic E6E7 RNA expression from a single integrated viral DNA copy in HPV16-positive SiHa cells and potential function of the virus-host chimeric transcript-4. (A) Two integration junctions of the HPV16 DNA integration site were verified by PCR (see two primer sets F3+F13 and B13+B5 in F) and by Sanger-sequencing, with the sequenced integration junction on the right. (B-C) Virus-host RNA CJRs were mapped to the HPV16 genome (B) and to host Chr13 (C). (D) Distribution of the top 10 RNA CJRs mapped to HPV16 and host Chr13. (E) Validation of two RNA splice junctions by RT-PCR using two primer sets F1+F12 and F2+F12 in F and Sanger-sequencing, with the sequenced splice junctions on the right. (F) The illustration shows the truncated HPV16 DNA integrated in a rearranged region (dashed lines) on Chr13 between nt 73,214,733 and nt 73,513,421 (top), with the primers used for validation of the integrated viral DNA. Total RNA-seq reads-coverage along with the chimeric host-virus genome region and Sashimi plots for splice junctions are visualized by IGV (middle). The integrated viral DNA expresses the E6E7 RNA by using the viral early promoter P97 and a host PAS at Chr13 nt 73,452,999 to produce four virus-host chimeric RNA isoforms (bottom). The number in nt below each isoform RNA indicates the estimated length of the RNA without a pA tail. (G) Identification of the major viral RNA transcripts in SiHa cells by Northern blot analysis using polyA⁺ RNA selected from 100 µg of total cell RNA. A mixture of three viral probes (B1, B2 and B3) located in the HPV16 E6 and E7 regions (F) were used for Northern blot analysis. RNA from HPV-negative C33A cells served as a negative control. GAPDH RNA served as a sample loading control. (H-M) Reduced expression of a viral-host fusion transcript (transcript-4) promotes the proliferation of SiHa cells by increasing cell cycle G1 to S phase transition. (H) Major viral-host fusion transcripts in SiHa cells and detection strategies by RT-PCR with the indicated primers (arrows). A synthetic siRNA was designed to specifically target transcript-4 (si-Junction). (I-K), Reduced expression of the transcript-4 level in SiHa cells led to a reduction of p53 protein and enhancement of cell proliferation. SiHa cells were harvested at 48 h after siRNA transfection. The knockdown efficiency of the transcript-4 was confirmed by RT-qPCR using a virus-host specific junction TaqMan probe (I). The effects of the knockdown of transcript-4 on p53 expression was measured by Western blot with tubulin serving as a sample loading control (J). (K) SiHa cell proliferation was examined at the indicated times (h) after transfection of 40 nM of the indicated siRNAs (arrow) by counting cell number (mean ± SE from three independent experimental repeats). ** indicates p<0.01 by unpaired, two-tailed Student's *t* test. (L and M) Reduction of transcript-4 expression by the junction-specific siRNA promotes cell entry into S phase. Flow cytometry analysis of the fraction of cells in different phases of the cell cycle was performed after indicated siRNA knockdown. SiHa cells 24 h after plating were transfected with 40 nM of the indicated siRNAs and harvested for flow cytometry 48 h after siRNA transfection. Data are from a representative of three separate experiments, each in triplicate (L). Percentage of cells in the indicated cell cycle after siRNA knockdown of the transcript-4 expression are shown as bar graphs (M). Data are mean ± SD; n=3. **, p<0.05 by unpaired, two-tailed Student's *t* test.

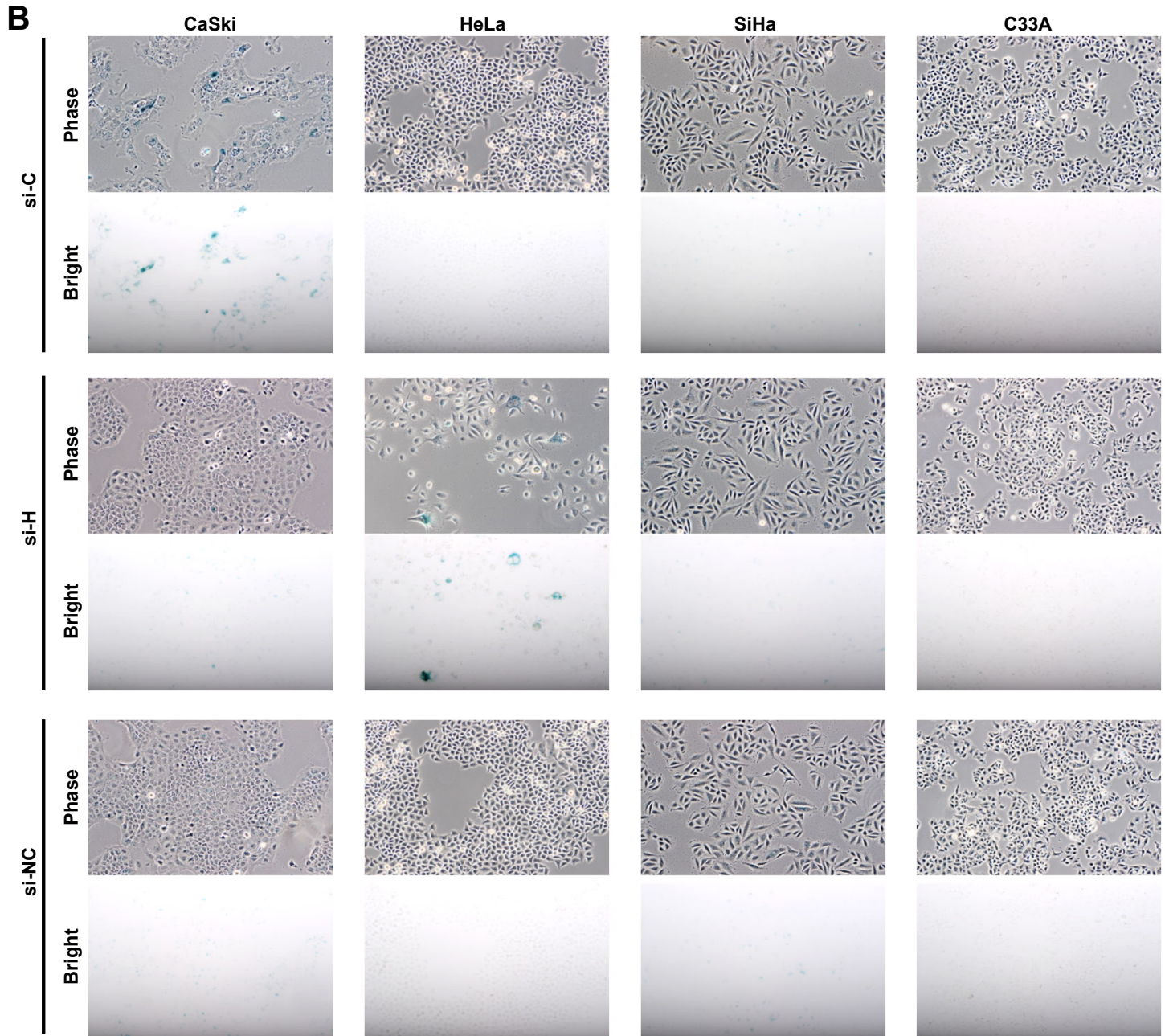
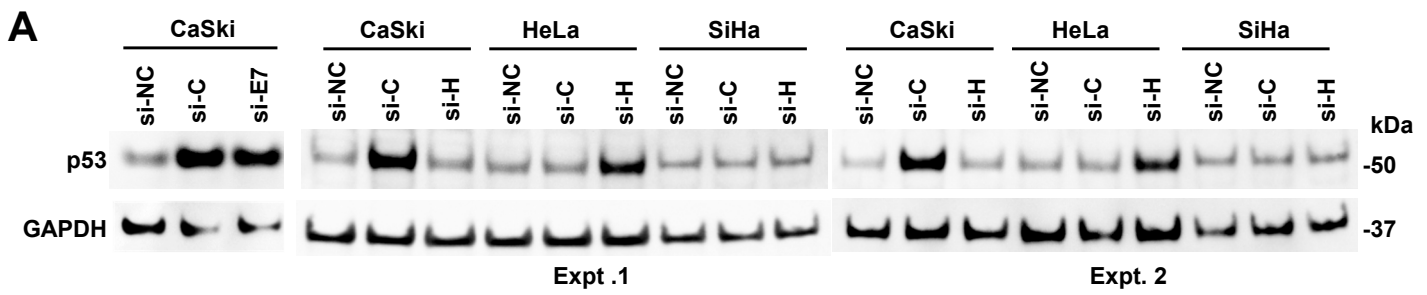
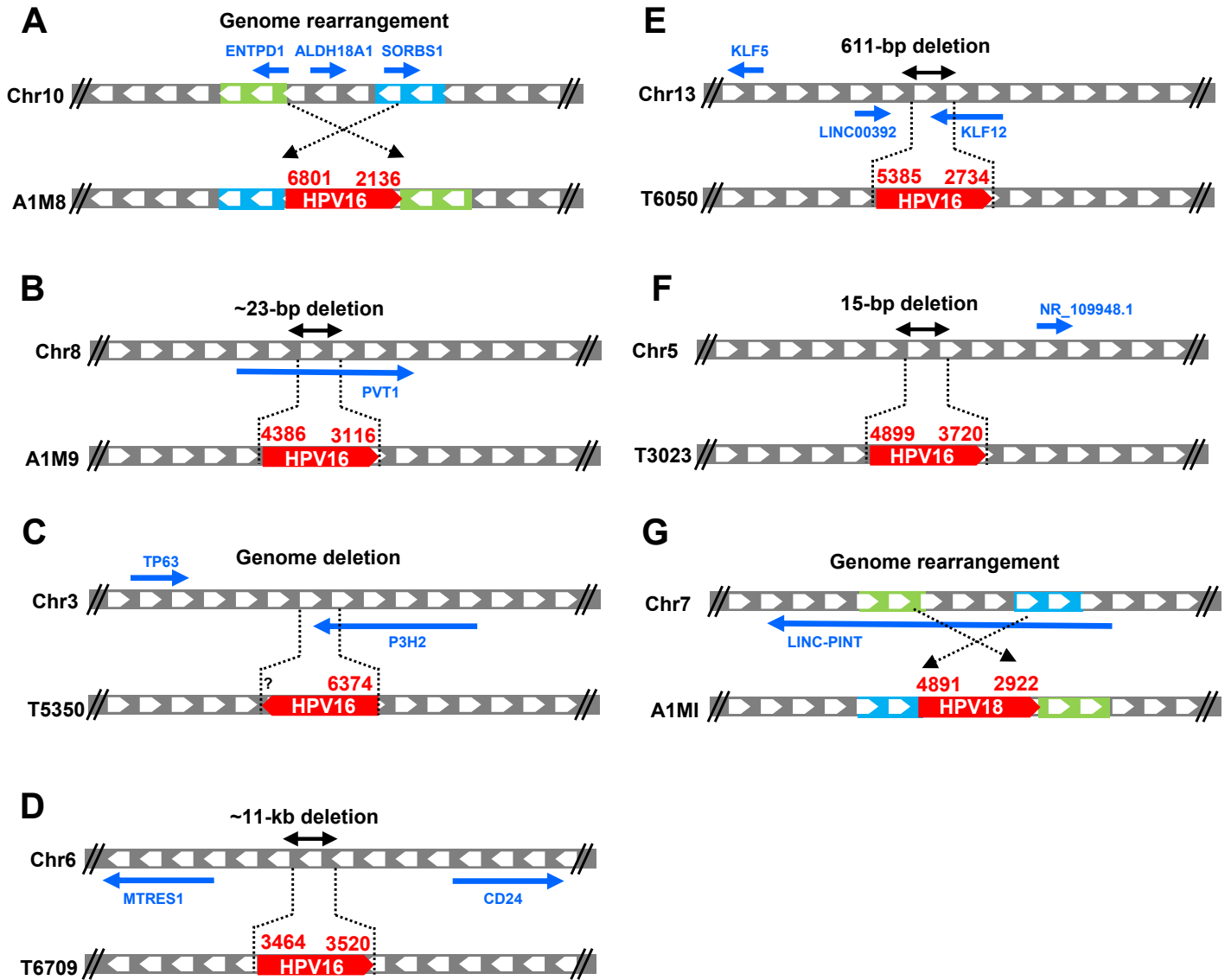
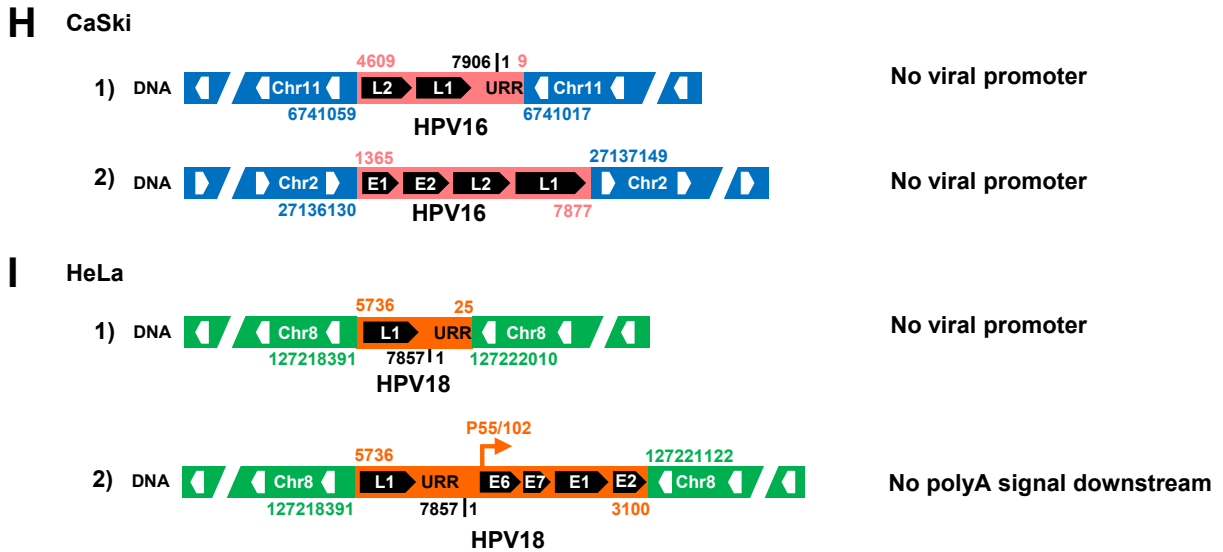


Figure S3. A siRNA targeting the host portion of the virus-host chimeric E6E7 RNA stabilizes p53 and promotes cell senescence in CaSki and HeLa cells, but not in SiHa or C33A cells. (A) Immunoblot of p53 protein expression in CaSki, HeLa and SiHa cells upon indicated siRNAs: si-NC, non-targeting siRNA control; si-C, CaSki-specific; si-E7, HPV16 E7-specific siRNA 198 (6); si-H, HeLa-specific. **(B)** Representative images of β -galactosidase (β -gal) staining. Cells were transfected twice with 40 nM of the indicated siRNAs before β -gal staining. The cells were plated on day 1, transfected with the siRNAs on day 2 and passed on day 4, and transfected again with the siRNAs on day 5. The β -gal staining was performed on day 8. The stained cells were examined on day 9 by both phase field and bright field microscopy.

Viral E6E7 expressible integration sites

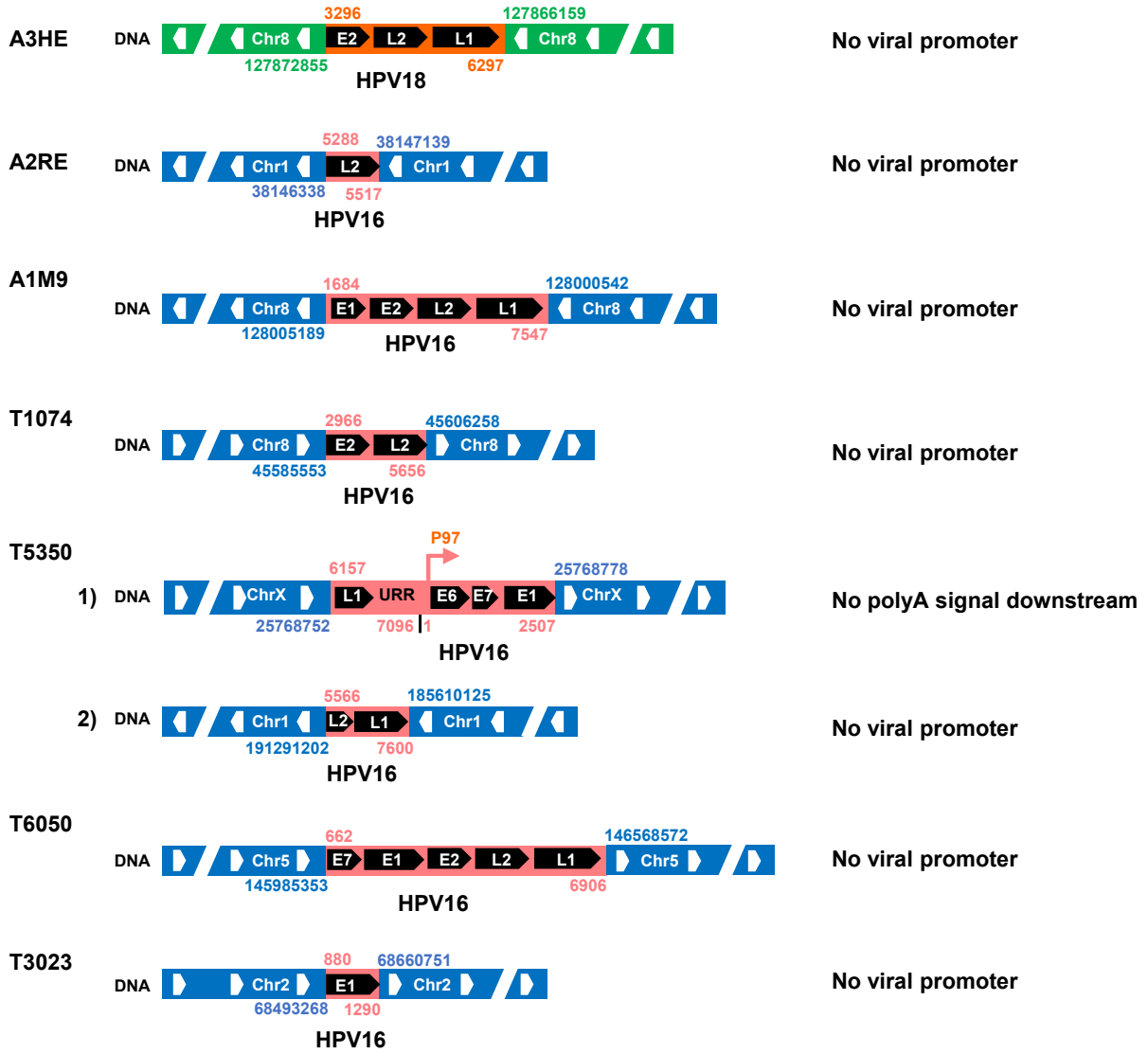


Viral E6E7 non-expressible integration sites



Viral E6E7 non-expressible integration sites

J Cervical cancer tissues



K

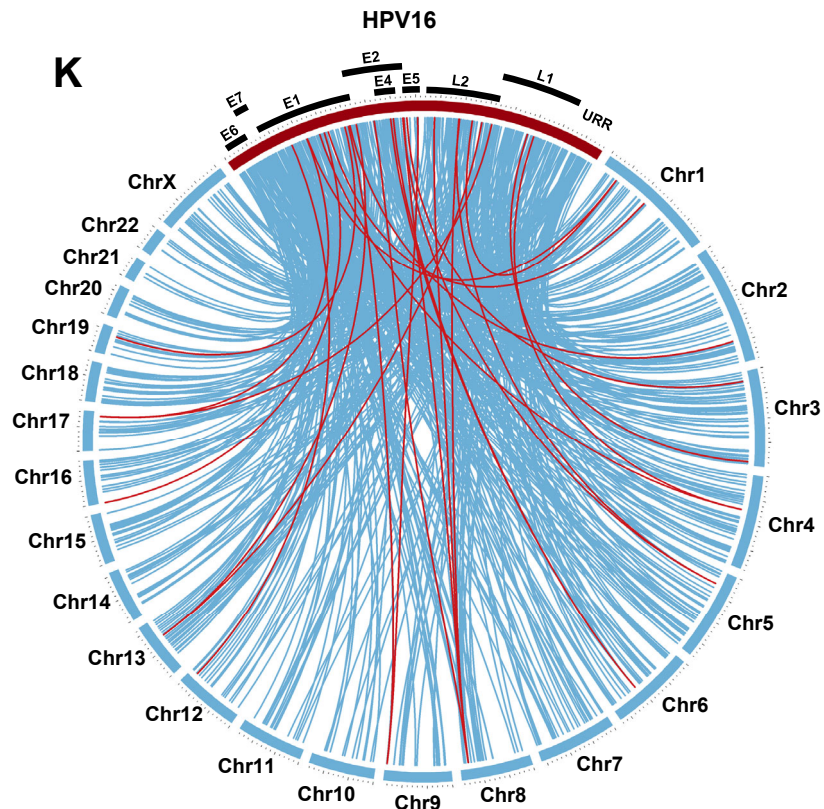


Figure S4. Simplified structures of viral E6E7 expressible or non-expressible HPV16 and HPV18 integration sites in cervical cancer tissues and derived cell lines. (A-G) E6E7 expressible HPV16 (**A-F**) and HPV18 (**G**) DNA integration sites. (**H-J**) Examples of E6E7 non-expressible HPV16 and HPV18 integration sites in CaSki (**H**), HeLa (**I**) and cervical cancer tissues (**J**). No viral promoter means no detectable viral early promoter P97 from the integrated HPV16 or P55/102 from HPV18 DNA. No PAS for the integrated viral DNA means an optimal poly A signal (UUAUUU or its equivalents) not available 10 kb downstream of the integrated viral DNA by the online Poly(A) Finder tool (<http://dnafsmineer.bic.nus.edu.sg/PolyA.html>). An optimal poly(A) signal within the searched region must be determined by a matrix score ≥ 0.5 and the presence of additional UGUA motifs upstream and G/U- or U-rich downstream elements (DSE) of the PAS. (**K**) Genome-wide distributions of HPV16 integration events in 17 HPV16-positive cervical cancer tissues. A total of 525 expressible and non-expressible viral DNA integration junctions were mapped from the HPV16 genome (red) to individual human chromosomes (blue). Relative viral ORF positions are shown above the linearized HPV16 genome in a scale of 100,000x. URR, upstream regulatory region. The size of human chromosomes (Chr1-22 + ChrX) is in proportion to their length. Each “tick” in an internal circle indicates the genome positions of a detected virus-host DNA junction (red, expressible; blue non-expressible) from the virus to host genome. The circular map was created using Circos package (doi:10.1101/gr.092759.109).

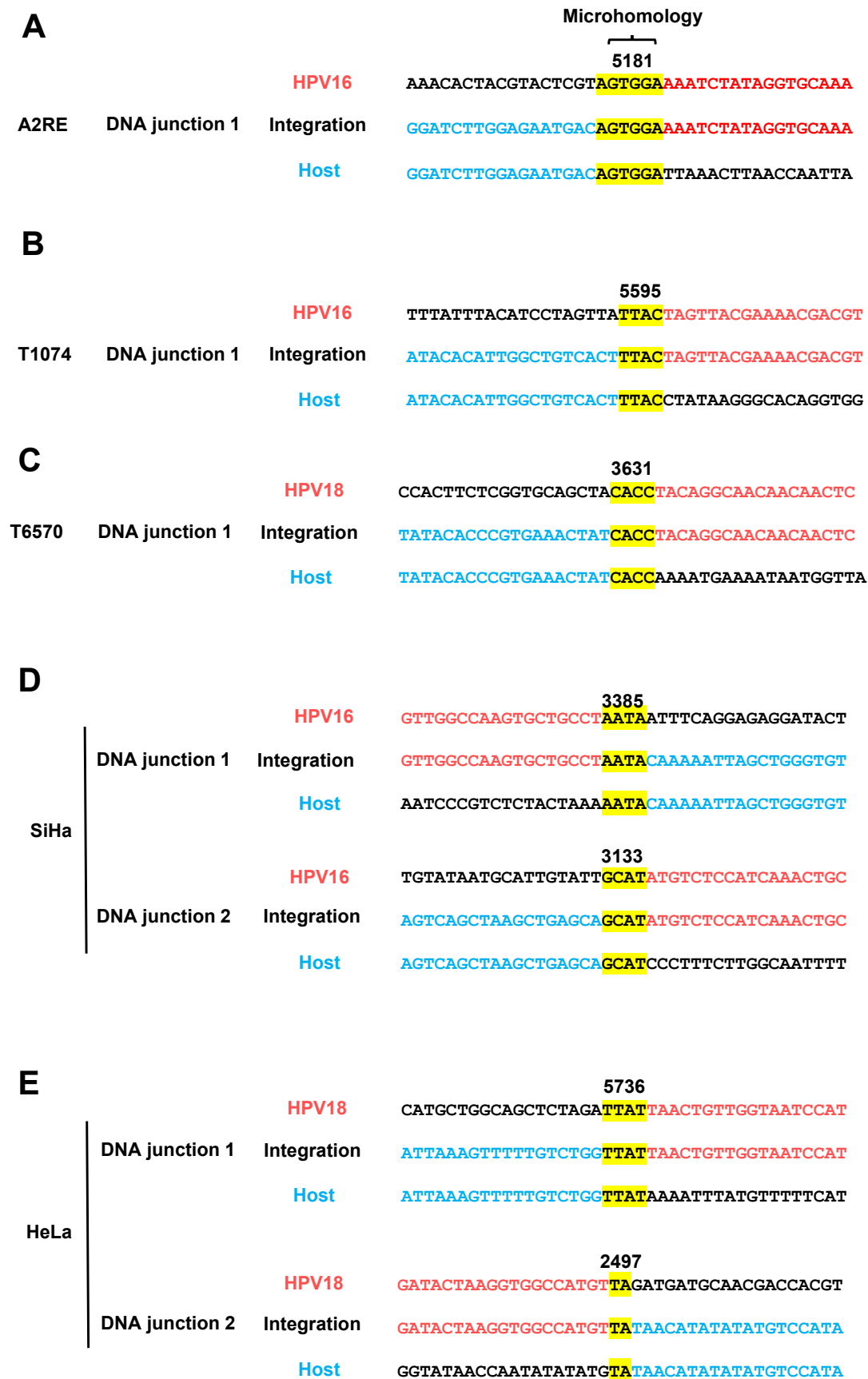
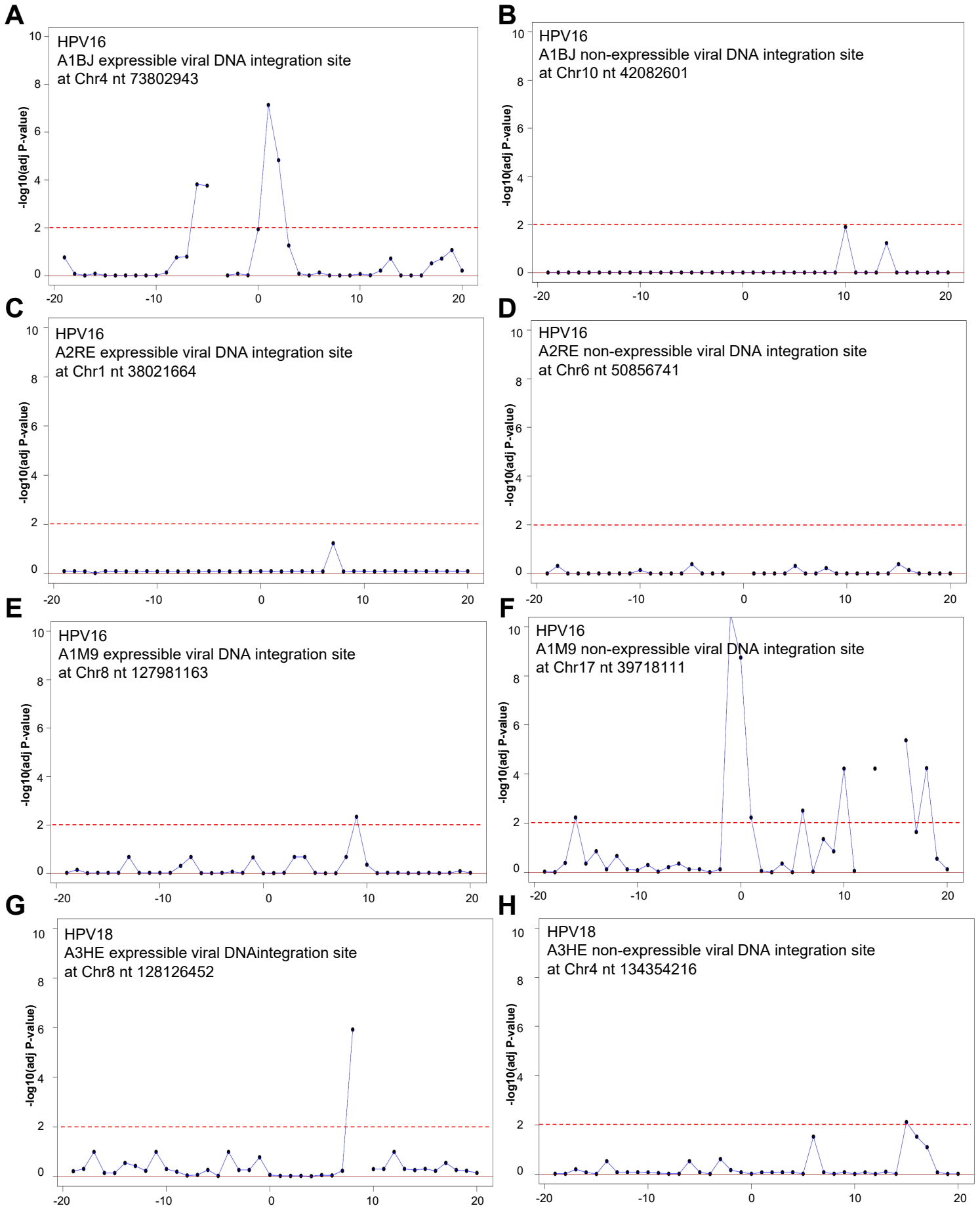
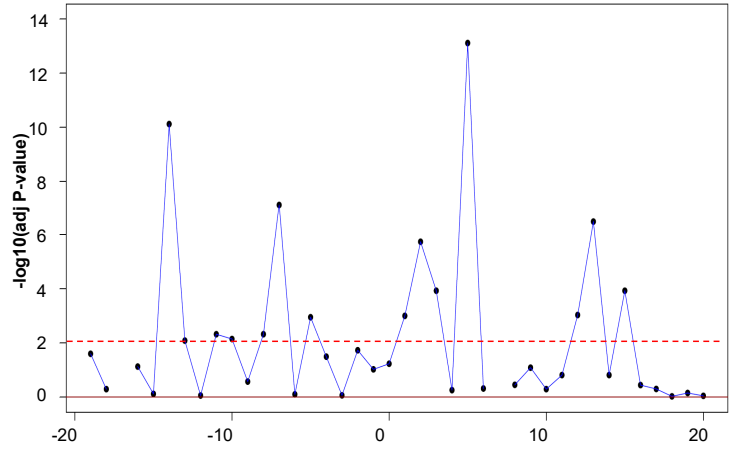
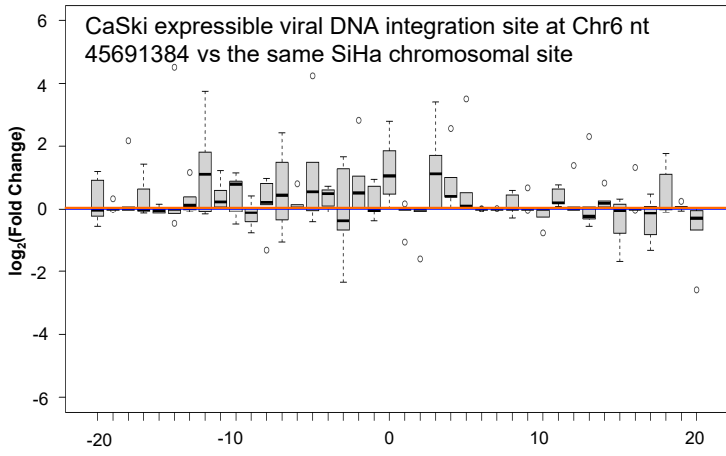


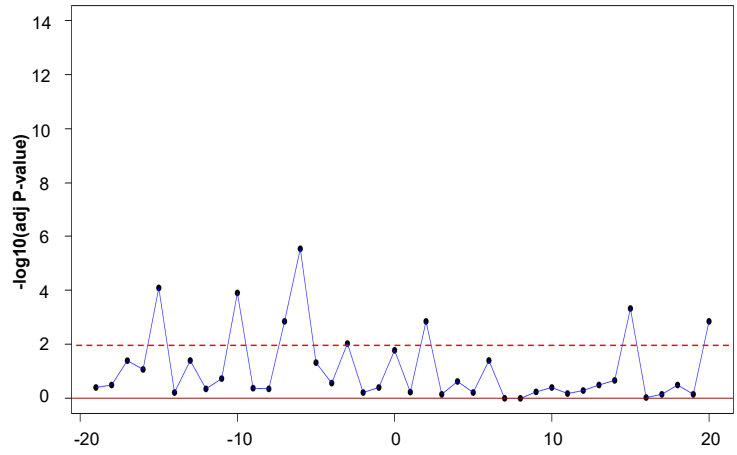
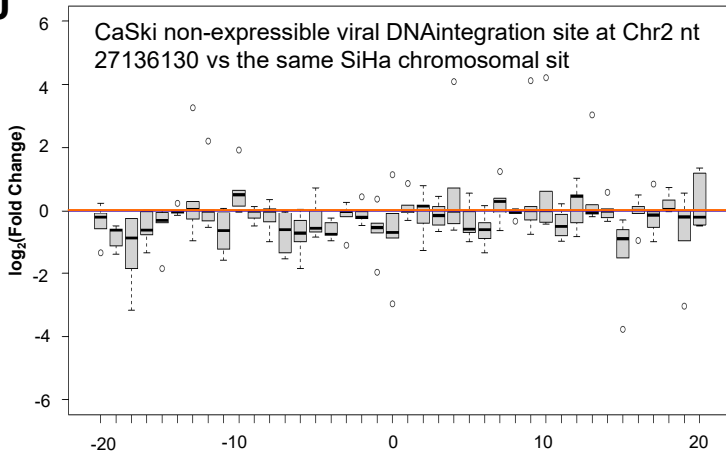
Figure S5. Regions of the microhomology sequence (MHs) were identified at virus-host integration junctions in the cervical cancer tissues A2RE (A), T1074 (B) and T6570 (C) and in SiHa (D) and HeLa (E) cells. Yellow color shows the MHs at the integration junction between HPV16/18 and human (hg38) reference genomes. Number above the yellow color indicates the last nt position in the virus genome at the integration junction.



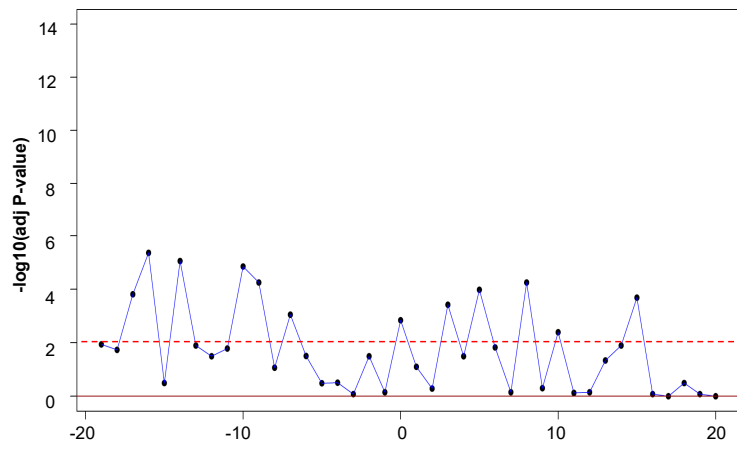
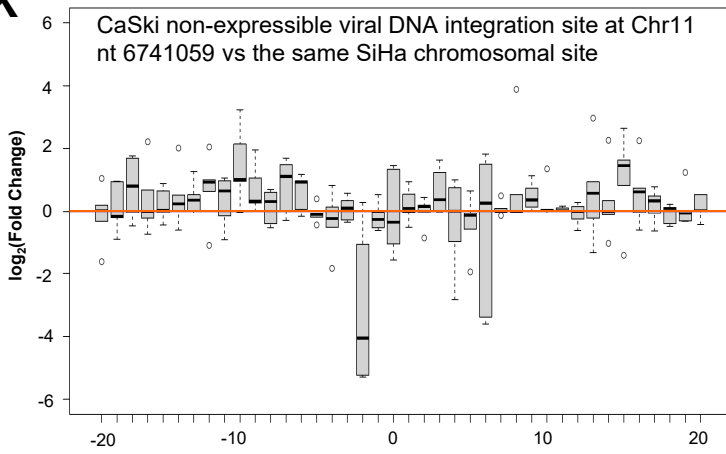
I



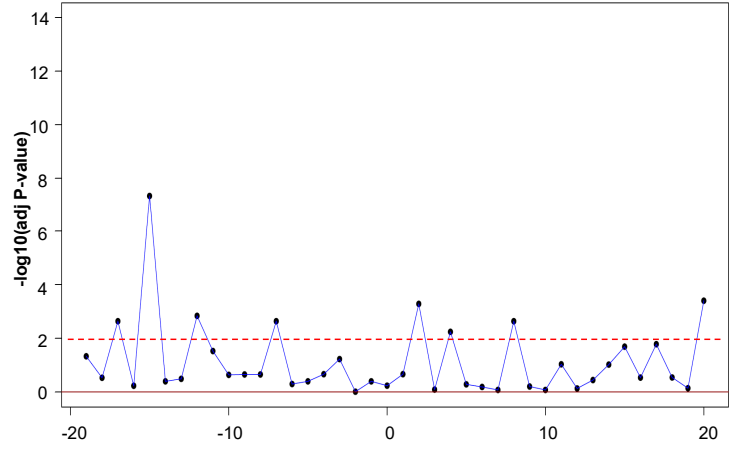
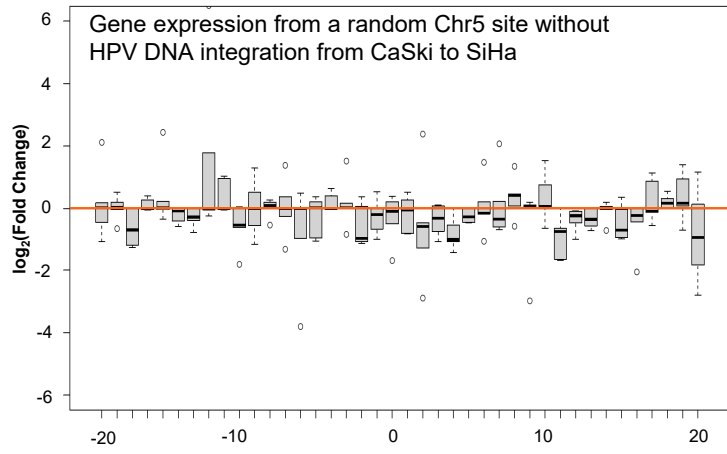
J



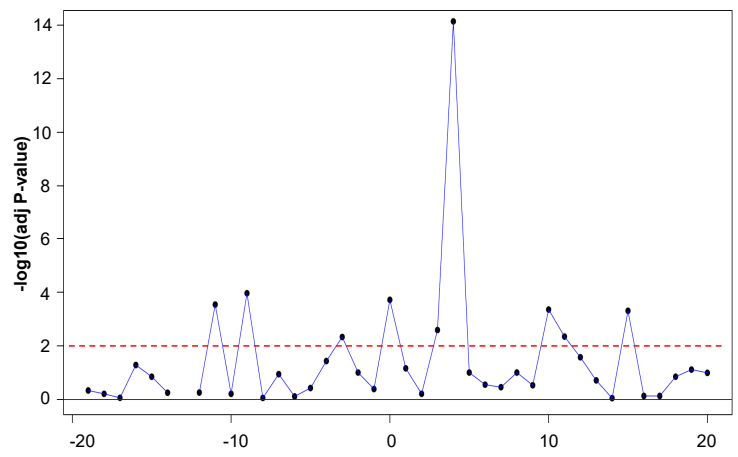
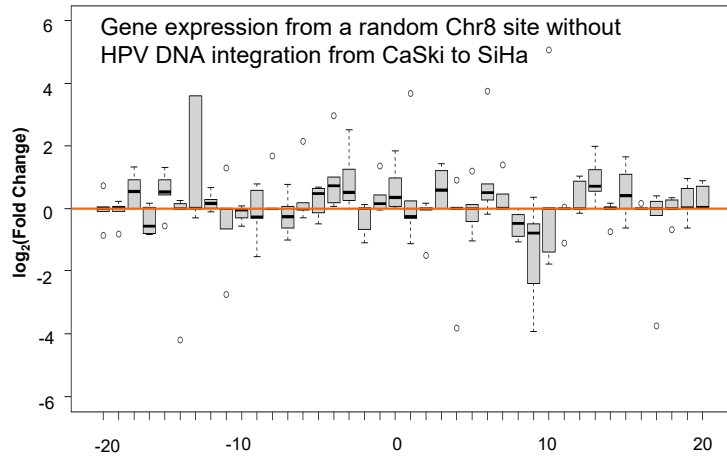
K



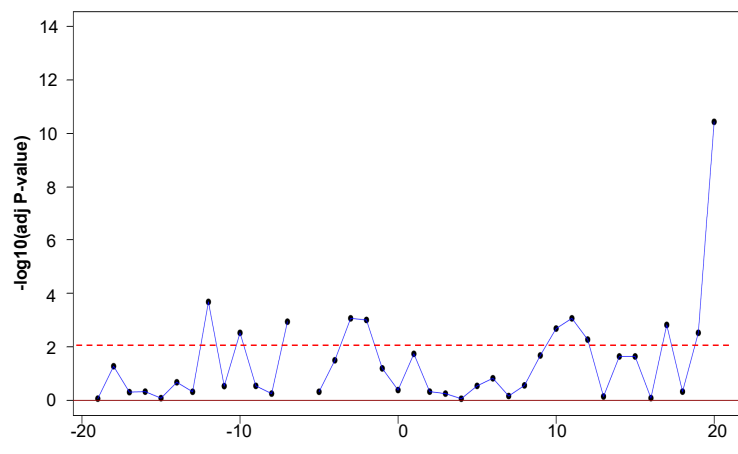
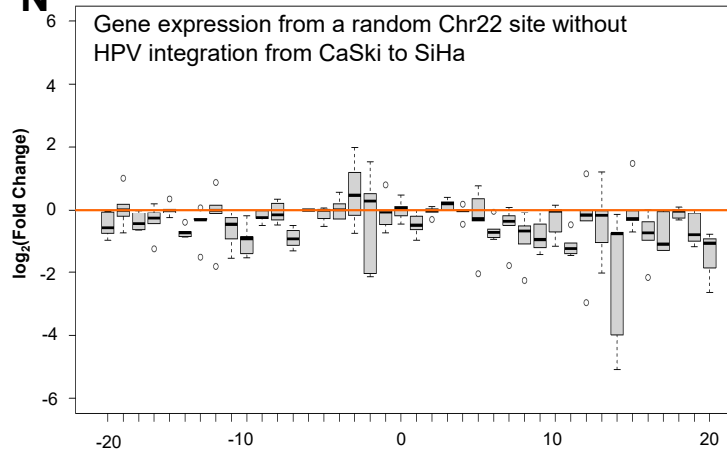
L

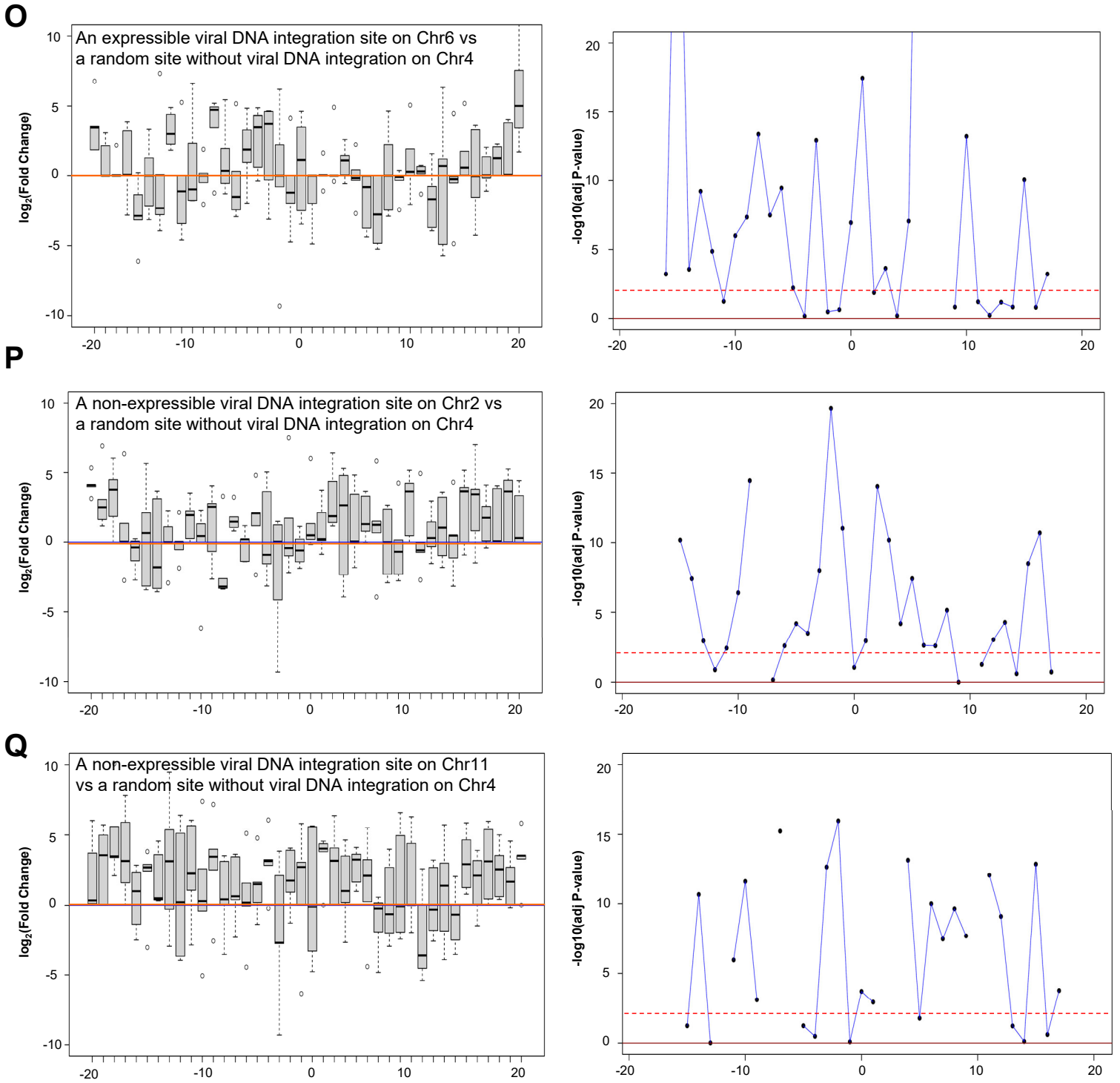


M

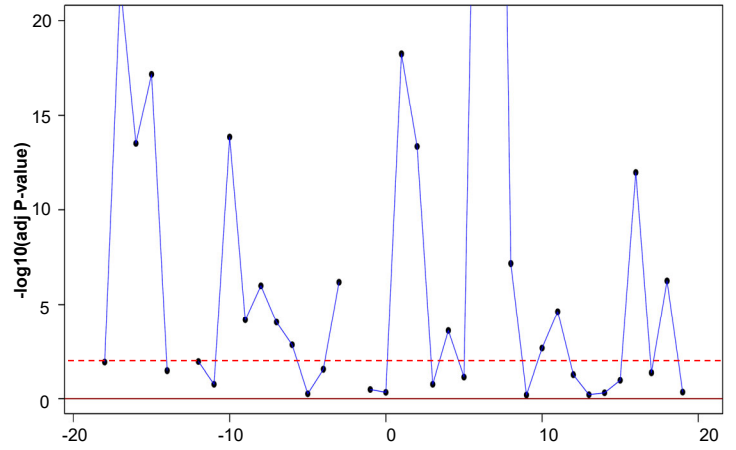
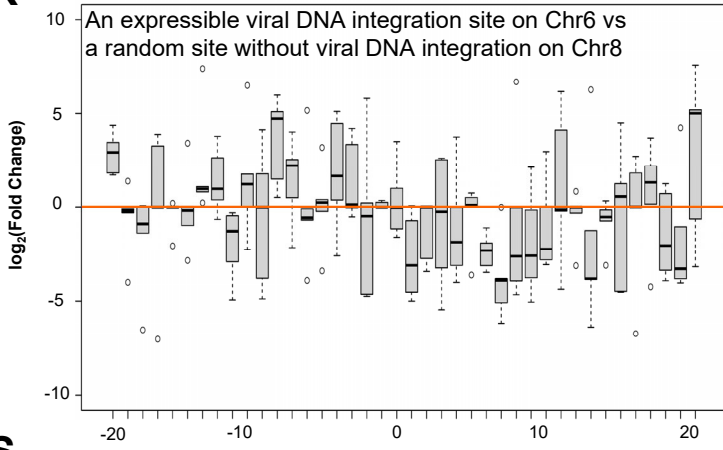


N

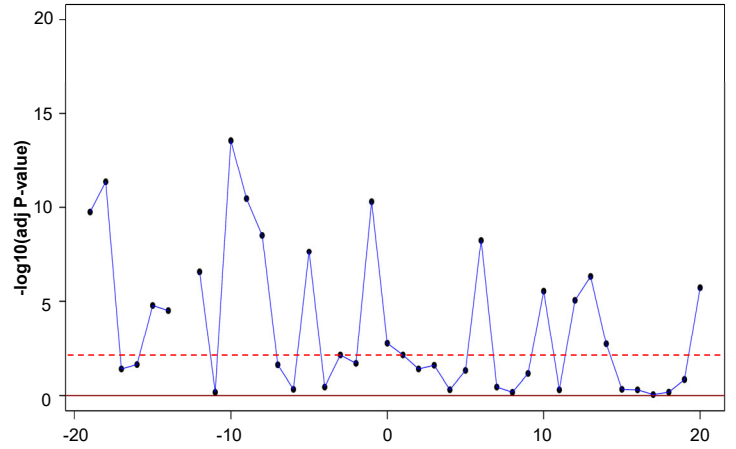
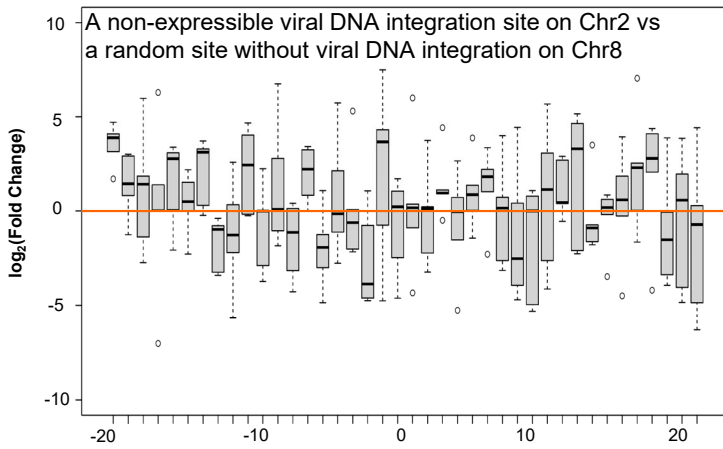




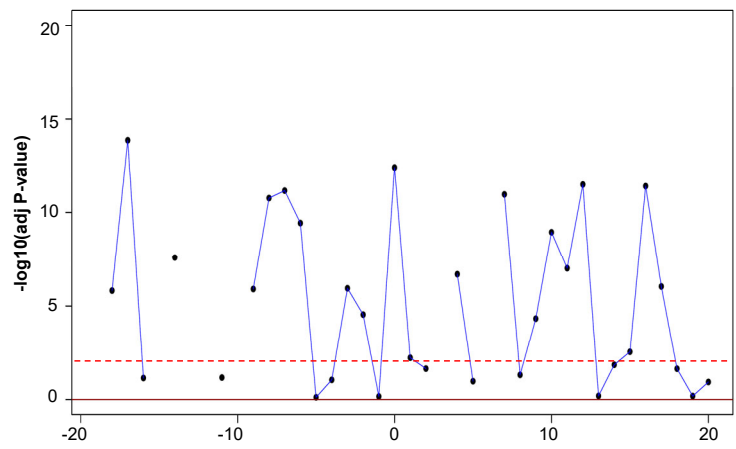
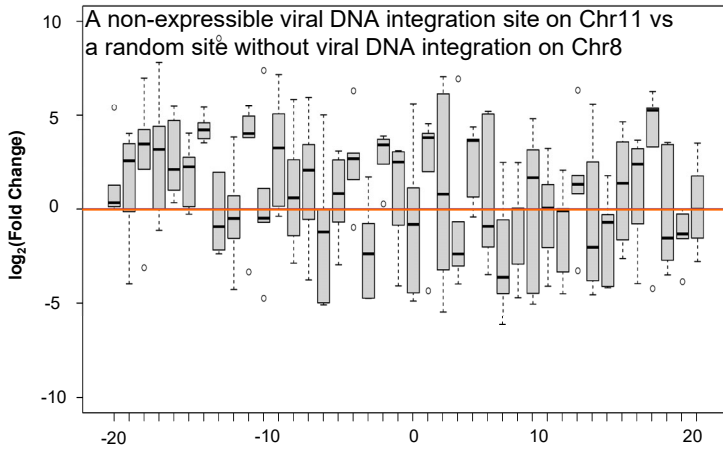
R



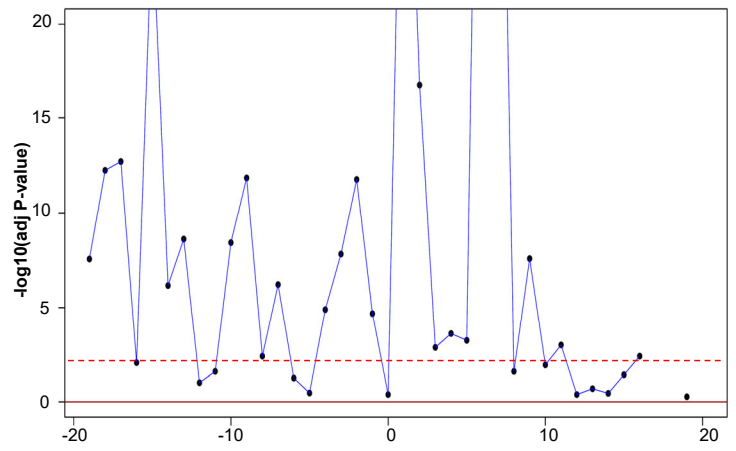
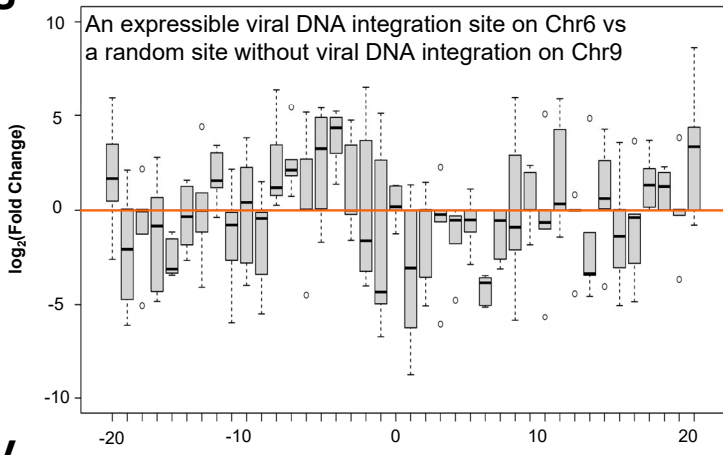
S



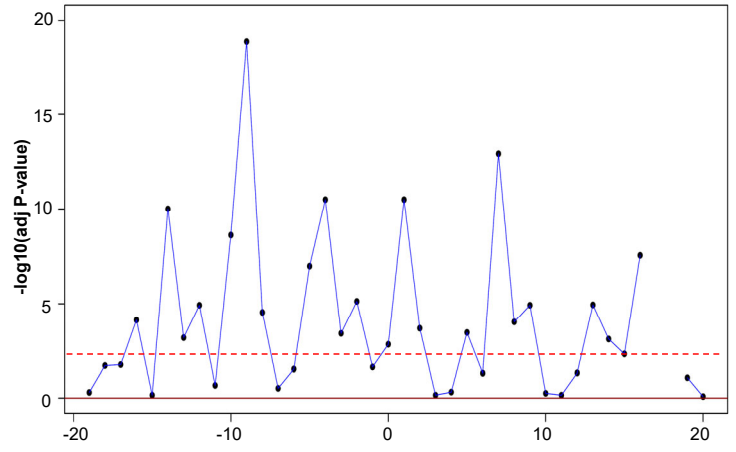
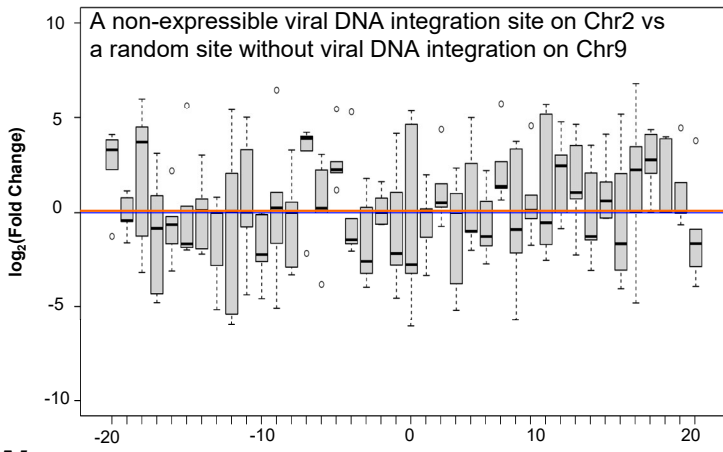
T



U



V



W

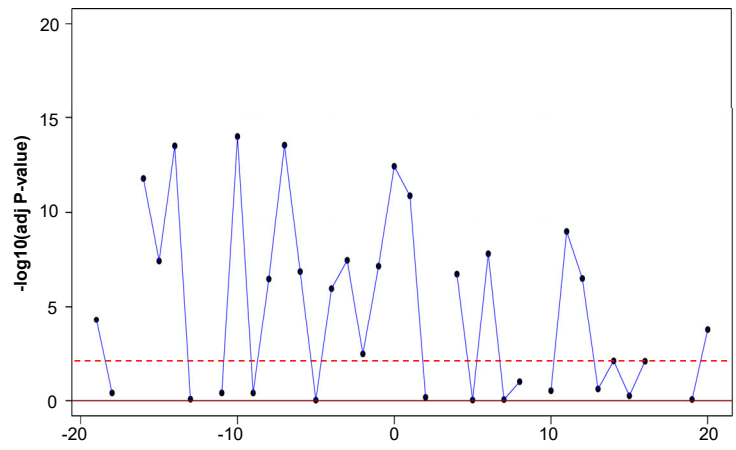
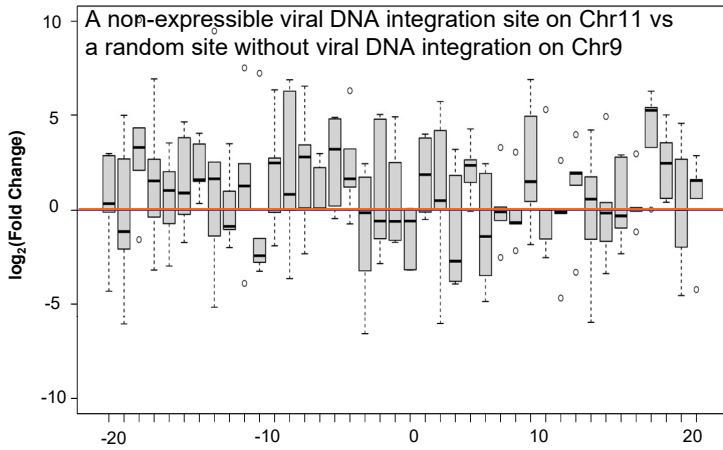


Figure S6. HPV DNA integration may or may not impact the expression of the nearby host genes in cervical cancer tissues. (A-H) Significance of the variance by F-test with Bonferroni correction within each bin shown in Figure 7 for individual expressible vs non-expressible HPV DNAs integrated in the mapped tissue chromosomal regions. Each point represents a five-gene bin, corresponding to those in Figure 7. The missing point in each bin indicates the P-value close to zero. The points above the dashed red line indicate $P < 0.01$. **(I-N)** Comparison of the host gene expression nearby an expressible or non-expressible HPV integration site in CaSki cells over the same chromosomal position in SiHa cells **(I-K)** or from the same randomly selected chromosomal site without HPV integration from CaSki to SiHa cells **(L-N)**. See Figure 7 and above for other description details. The genome range for each panel is shown as below: Chr6 nt 37,482,938-75,250,323 **(I)**, Chr2 nt 885,290-43,767,987 **(J)**, Chr11 nt 3,671,083-18,452,063 **(K)**, Chr5 nt 39,105,252-93,594,611 **(L)**, Chr8 nt 22,138,020- 66,178,464 **(M)**, Chr22 nt 18,149,899-35,394,207 **(N)**. **(O-W)** Comparison of host gene expression nearby the expressible or non-expressible HPV DNA integration site to a randomly selected chromosomal site without HPV DNA integration in CaSki cells. See Figure 7 and above for other description details. The genome ranges for comparison in each panel are shown below: Chr6 nt 37,482,938-75,250,323 vs Chr4 nt 9,224,896-70,058,845 **(O)**, Chr8 nt 10,482,878-42,796,392 **(R)** or Chr9 nt 93,058,688-123,268,576 **(U)**; Chr2 nt 885,290-43,767,987 vs Chr4 nt 9,224,896-70,058,845 **(P)**, Chr8 nt 10,482,878-42,796,392 **(S)** or Chr9 nt 93,058,688-123,268,576 **(V)**; Chr11 nt 3,671,083-18,452,063 vs Chr4 nt 9,224,896-70,058,845 **(Q)**, Chr8 nt 10,482,878-42,796,392 **(T)** or Chr9 nt 93,058,688-123,268,576 **(W)**. Panels on the right of **I, J, K, L, M, N, O, P, Q, R, S, T, U, V,** and **W** show the correspondent statistics of the variance by F-test with Bonferroni correction within each bin (five genes per bin). The missing point in each bin indicates the P-value close to zero. The points above the red dashed line indicate the $P < 0.01$.

A

Number	Genome nt position	Reference	CaSki
1	131	A	G
2	350	T	G
3	1252	G	A
4	1418	C	A
5	1522	T	A
6	2457	C	T
7	2609	G	T
8	2938	A	G
9	3384	T	C
10	3410	C	T
11	3684	C	A

B

Number	Genome nt position	Reference	SiHa
1	350	T	G
2	442	A	C
3	645	A	C
4	1194	A	G
5	1842	A	G
6	3068	G	A

C

Number	Genome nt position	Reference	HeLa
1	104	T	C
2	485	T	C
3	549	C	A
4	751	C	T
5	806	G	A
6	1012	A	T
7	1194	C	A
8	1353	T	A
9	1807	T	C
10	1843	T	G
11	2269	C	T
12	5875	C	A
13	6401	A	G
14	7258	T	A
15	7486	C	T
16	7529	C	A
17	7567	A	C
18	7592	T	C
19	7670	A	T

D

HPV16 E6 protein	
HPV16 Ref.	MFQDPQERPR <u>K</u> LPLQLCTELQTTIHDIILECVYCKQQLLRREVDFAFRDLCIVYRDGNPYAVCDKCLKFYISKISEYRHYCY <u>S</u> LYGTTLEQQYNKPLCDLLIRCIN CQKPLCP <u>E</u> EKQRHLDKKQRFHNIRGRWTGRCM <u>S</u> CCRSSRTRRETQL*
CaSki	MFQDPQERPR <u>G</u> KLPLQLCTELQTTIHDIILECVYCKQQLLRREVDFAFRDLCIVYRDGNPYAVCDKCLKFYISKISEYRHYCY <u>S</u> VYGTTLEQQYNKPLCDLLIRCIN CQKPLCP <u>E</u> EKQRHLDKKQRFHNIRGRWTGRCM <u>S</u> CCRSSRTRRETQL*
SiHa	MFQDPQERPR <u>K</u> LPLQLCTELQTTIHDIILECVYCKQQLLRREVDFAFRDLCIVYRDGNPYAVCDKCLKFYISKISEYRHYCY <u>S</u> VYGTTLEQQYNKPLCDLLIRCIN CQKPLCP <u>D</u> EKQRHLDKKQRFHNIRGRWTGRCM <u>S</u> CCRSSRTRRETQL*

E

HPV16 E7 protein	
HPV16 Ref.	MHGDTPTLHEYMLDLQPETTDLYCYEQ <u>L</u> NDSS EE EEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRLCVQSTHVDIRTLLEDLLMGTLGIVCPICSQKP*
SiHa	MHGDTPTLHEYMLDLQPETTDLYCYEQ <u>F</u> NDSS EE EEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRLCVQSTHVDIRTLLEDLLMGTLGIVCPICSQKP*

F

HPV18 E7 protein	
HPV18 Ref.	MHGPKATLQDIVLHLEPQNEIPVDLLCHEQLSDSEEENDEIDGVNHQHLPA RRAEPQRHTMLCMCKCEAR <u>I</u> ELVV ESS ADDLRAFQQLFLNTLSFVCPWC ASQQ*
HeLa	MHGPKATLQDIVLHLEPQNEIPVDLLCHEQLSDSEEENDEIDGVNHQHLPA RRAEPQRHTMLCMCKCEAR <u>K</u> LV ESS ADDLRAFQQLFLNTLSFVCPWC ASQQ*

G

Number	HPV16 Genome nt position	Reference	A1BJ	A1M9	A1M8
1	109	T	C	C	T
2	178	T	T	T	G
3	350	T	G	T	T
4	647	A	A	A	G
5	840	C	C	C	T
6	846	T	T	T	C

H

Number	HPV18 Genome nt position	Reference	A3HE	A1MI
1	104	T	G	G
2	485	T	C	C
3	549	C	A	A
4	751	C	T	T

I

Number	HPV16 Genome nt position	Reference	T1074	T6050	T6709	T5350	T3023
1	178	T	G	G	G	G	T
2	256	C	C	C	C	C	T
3	350	T	T	T	T	T	G
4	647	A	G	G	G	G	A
5	843	T	C	C	T	T	T
4	846	T	C	C	C	C	T

J

HPV16 E6 protein	
HPV16 Ref.	MFQDPQERPRKLPQLCTELQTTIH <u>D</u> IILECVYCKQQLLRREVDFAFRDLCIVY RDGNPYAVCDKCLKFYISKISEYRHYCYS <u>L</u> YGTTLQYQNKPLCDLLIRCINCQ KPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*
A1M8	MFQDPQERPRKLPQLCTELQTTIH <u>E</u> IILECVYCKQQLLRREVDFAFRDLCIVY RDGNPYAVCDKCLKFYISKISEYRHYCYS <u>L</u> YGTTLQYQNKPLCDLLIRCINCQ KPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*
A1BJ	MFQDPQERPRKLPQLCTELQTTIH <u>I</u> IILECVYCKQQLLRREVDFAFRDLCIVY RDGNPYAVCDKCLKFYISKISEYRHYCYS <u>V</u> YGTTLQYQNKPLCDLLIRCINCQ KPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*

K

HPV16 E7 protein	
HPV16 Ref.	MHGDTPTLHEYMLDLQPETTDLYCYEQL <u>N</u> DSSEEEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRRCVQSTHVDIRTLLEDLLMGTGIVCPICSQKP*
A1M8	MHGDTPTLHEYMLDLQPETTDLYCYEQL <u>S</u> DSSEEEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRRCVQSTHVDIRTLLEDLLMGTGIVCPICSQKP*

L

HPV16 E6 protein	
HPV16 Ref.	MFQDPQERPRKLPQLCTELQTTIH <u>D</u> IILECVYCKQQLLRREVDFAFRDLCIV YRDGNPYAVCDKCLKFYISKISEYRHYCYS <u>L</u> YGTTLQYQNKPLCDLLIRCIN CQKPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*
T1074 T6050 T6709 T5350	MFQDPQERPRKLPQLCTELQTTIH <u>E</u> IILECVYCKQQLLRREVDFAFRDLCIVY RDGNPYAVCDKCLKFYISKISEYRHYCYS <u>L</u> YGTTLQYQNKPLCDLLIRCINCQ KPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*
T3023	MFQDPQERPRKLPQLCTELQTTIH <u>I</u> IILECVYCKQQLLRREVDFAFRDLCIVY RDGNPYAVCDKCLKFYISKISEYRHYCYS <u>V</u> YGTTLQYQNKPLCDLLIRCINCQ KPLCPEEKQRHLDKKQRFHNIRGRWTGRCMSSCRSSRTRRETQL*

M

HPV16 E7 protein	
HPV16 Ref.	MHGDTPTLHEYMLDLQPETTDLYCYEQL <u>N</u> DSSEEEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRRCVQSTHVDIRTLLEDLLMGTGIVCPICSQKP*
T1074 T6050 T6709 T5350	MHGDTPTLHEYMLDLQPETTDLYCYEQL <u>S</u> DSSEEEDEIDGPAGQAEPDR AHYNIVTFCKCDSTLRRCVQSTHVDIRTLLEDLLMGTGIVCPICSQKP*

N

SiHa virus-host fusion protein 1 Transcript-3: HPV16:97-880/Chr13:73456962-2979

Coding region: HPV16 E1:865-880/Chr13:73456962-73456814 (165 nts / 54 aa)

Protein: MADPADRKLDVCLANWIKGQDPSVCCIQETHLTCKDTHRIKIKGWRKIYQANGK*

O

SiHa virus-host fusion protein 2 Transcript-4: HPV16:97-226/Chr13:73456962-2979

Coding region: HPV16 E6:104-226/Chr13:73456962-73456814 (159 nts / 52 aa)

Protein: MFQDPQERPRKLPQLCTELQTTIHDIILECVYCKQQLLRREIENWMSAWQIG*

P

HeLa virus-host fusion protein 1 Transcript-2: Chr8:127219268-18810/HPV18:24-233/416-929/Chr8:127229130-28559

Coding region: Chr8:127218909-127218810/HPV18:24-118 (195 nts / 64 aa)

Protein: MNFGLLCNPSGAFTDNIDFEVALTPALPLTRLVYKKGSNRKRSGPKTVYIKDVRNTPQYHGAL*

Q

HeLa truncated HPV18 E1 protein Transcript-3: HPV18:102-233/416-2497/Chr8:127229301-8559

Coding region: HPV18:914-2497/Chr8:127229301-127229299 (1587 nts / 528 aa)

Protein: MADPEGTDGEGTGCNGWFYVQAIVDKKTGDVISDDDEDENATDTGSDMVDIFDTQGTFCQAELETAQALFHAQEVH
 NDAQVLHVLKRKFAGGSKENSPLGERLEVDTLSPRLQEISLNSGQKKAKRRLFTISDSGYGCSEVEATQNQVTTNGEHGG
 NVCSGGSTEIDNGGTEGNSSVDGTSDNSNIENVPQCTIAQLKDLLKVNKQGAMLAVFKDITYGLSFTDLVRNFKSDKTT
 CTDWVTAIFGVNPTIAEGFKLIQPFILYAHIQCLDCKWGVLLALLRYKCGKSRLTVAKGLSTLLHVPETCMLIQPPKLRSSVA
 ALYWYRTGISNISEVMGDTPEWIRQLTIHQHIDDSNFDLSEMVQWAFDNLDESMAFEYALLADSNSNAAFLKSNCQAK
 YLKDCATMCKHYRRAQKRQMNMSQWIRFRCSKIDEGGDWRPQVFLRYQQIEFITFLGALKSFLKGTPKKNCLVFCGPANTG
 KSYFGMSFIHQGAVISFVNSTSHFWLEPLTDTKVAML*

R

A1BJ virus-host fusion protein 1 Transcript-4: HPV16:97-226/Chr4:73777731-73774860

Coding region: HPV16 E6:104-226/Chr4:73777731-73777699 (150 nts / 49 aa)

Protein: MFQDPQERPRKLPQLCTELQTTIHDIILECVYCKQQLLRREGQDDNEPS*

S

A2RE virus-host fusion protein 1 Transcript-4: HPV16:97-226/Chr1:38022699-38023808

Coding region: HPV16 E6:104-226/Chr1:38022699-38022809/38023545-38023625 (318 nts / 105 aa)

Protein:

MFQDPQERPRKLPQLCTELQTTIHDIILECVYCKQQLLRREERIAKERQKQYNCLTQRIEREKLFVIAQKIQTRKDLMDKTQKV
 KVKKETVNSPAIYKFSRRKR*

T

T1074 host-virus fusion protein 1 Transcript-1: [Chr17:73644445-73348599/HPV16: 5639-7321](#)Coding region: [Chr17:73644088-73348599/HPV16 L1: 5639-7156](#) (7683 nts / 2560 aa)

Protein:

MWGLLIWTLALHQIRAARAQDDVSPYFKTEPVRTQVHLEGNRLVLTCAEGSWPLEFKWLHNNRELTKFSLEYRYMITSLD
 RTHAGFYRCIVRNRMGALLQRQTEVQVAYMGSFEEGEGKHQSVSHGEAAVIRAPRIASFPQPQVTFWRDGRKIPSSRIAITLE
 NTLVILSTVAPDAGRYVQAVNDKNGDNKTSQPITLTVENVGGPADPIAPTIIPPKNTSVVAGTSEVTLECVANARPLIKLHIIWK
 KDGVLLSGGISDHNRLTIPNPTGSDAGYEECEAVLRSSSVPSVVRGAYLSVLEPPQFVKEPERHITAEEMKVVDPICQAKGVP
 PPSITWYKDAAVVEVEKLRFRQRNDGGLQISGLVPDDTGMFQCFAAAGEVQTSTYLAVTSIAPNITRGPLDSTVIDGMSV
 LACETSGAPRAITWQKGERILASGSVQLPRFTPLESGSLLISPTHISDAGTYTCLATNSRGVDEASADLVVWARTRITKPPQD
 QSVIKGTQASMVCGVTHDPRVTIRYIWEKDGATLGTESHPRIRLDRNGSLHISQTSWSDIGTYTCRVISAGGNDRSRSHLRVRQ
 LPHAPEHPVATLSTVERRAINLTWTKPFDGNSPLIRYILEMSENNAPWTVLLASVDPKATSVTVKGLVPARSYQFRLCAVNDVG
 KGQFSKDERVSLPEEPPTAPPQNVIASGRTNQSIMIQWQPPPEHQNGILKGYIIRYCLAGLPVGYQFKNITDADVNNLLLEDL
 IIWTNYEIEVAAYNSAGLGVYSSKVTEWTLQGVPTVPPGNVHAEATNSTTIRFTWNAPSPQFINGINQGYKLIawePEQEEV
 MVTARPNFQDSIHVGFVSGLKKFTEYFTSVLCFTTPGDGPRSTPQLVRTHEDVPGPVGHLSFSEILDTSLKVSWQEPGEKNGIL
 TGYRISWEEYNRTNTRVTHYLPNVTLLEYRVTGLTALTTYIEVAAMTSKGGQVQVASTISSGVPELPGPPTNLGISNIGRPSVT
 LQFRPGYDGKTSISRWLVEAQVGVVGEWLLIHLQSLNEPDARSMEVPDLNPFCTCYFRMRQVNIIVGTSPSPQSRKIQTLQ
 APPDMAPANVSLRTASETSLWLRWMLPEMEYNGNPESVGYKIKYRSRSDGHGKTLSHVVQDRVERDYTIEDLEEWTEYRVQ
 VQAFNAIGSGPWSQTVVGRTRRESVPSSGPTNVSAALATSSSMLVRWSEVPEADRNLVLYGYKVMYKEKSDTQPRFWLVEG
 NSSRSAQLTGLGKYVLYEVQVLAFTRIGDGSHPILERTLDDVPGPPMGILFPEVRTTSVRLIWQPPAAPNGIILAYQITHRLN
 TTTANTATVEVLAPSARQYATGLKPESVYLFRIQAQTRKGWGEAAEALVVTTEKDRPQPPSRPMVQQEDVRARSVLLSWE
 PGSDGLSPVRYTIQTRELPSGRWALHSASVSHNASSFIVDRPKPFTSYKFRVKATNDIGDSEFSEESESLTLQAAPDEAPTL
 SVTPHTTTSVLRWQPPAEDKINGILLGFRIRYRELLYELRGLFTLRGINNPGATWAELTSMYSMRNLSRPSLTQYELDNLNKH
 RRYEIRMSVYNVAVGEGSPSPQEVFVGEAVPTAARNVVHGHATATQLDVTWEPPLDSQNGDIQGYKIYFWEAQRGNLTER
 VKTLFLAENSVKLNLTGYTAYMVSVAAFNAAGDGPRSTPTQGTQQAAPSAPSSVKFSELTTTSVNVSWEAPQFPNGILEG
 YRLVYEPSPVDGVSKIVTVDVKGNLPLWLVKDLAEGVYRFRIRAKTFTYGPPIEAVTTGPGEGAPGPPGVPIIVRYSSAIA
 IHWSSGDPGKGPITRYVIEARPSDEGLWDILIKDIPKEVSSYTFMSDILKPGVSYDFRVIANDYGFGTSPSSPSQSVPAQKANPF
 YEEWWFLVVIALVGLIFILLVFLIIRGQSKKYAKKTDGNSAKSGALGHSEMMSLDESSFALELNNRRLSVKNSFCRKNGL
 YTRSPRRPSPGSLHYSDVDVTKYNDLIPAESSLTKPSEISDSQMSLWLPSEATVYLPVPVSKVSTDEYVARTNIYYHAGT
 SLLAVGHYPYFIKPPNNKILVPKVSGLQYRVFRIHLDPNKFQFPDTSFYNPDTQRLVWACVGVVGRGQPLGVGISGHPL
 LNKLDDENASAYAANAGVDNRECISMDYKQTQLCLIGCKPPIGEHWGKGSPTNAVAVNPGDCPPELINTVIQDGMVDTGF
 GAMDFTTLQANKSEVPLDICTSICKYPDYIKMVSEPYGDSLFFYLRRQMFVRHLFNRAGTVGENVPDDLVIKGSSTANLASS
 NYFPTSPGSMVTSDAQIFNKPYWLQRAQGHNNGICWGNQLFVTVDTRTNSMLCAAISTSETTYKNTNFKEYLRHGEEYDL
 QFIFQLCKITLTADVMTYIHSMNSTILEDWNFGLQPPPGTLEDTYRFVTSQAIACQKHTPPAPKEDPLKKYTFWEVNLKEKFS
 ADLDQFPLGRKFLQAGLKAKPKFTLGKRKATPTTSSTSTAKRKRKL*

U

T1074 virus-host fusion protein 2 Transcript-5: [HPV16:97-226/Chr17:73338940-73334403](#)Coding region: [HPV16 E6:104-226/Chr13:73338940-73338587](#) (477 nts / 158 aa)

Protein:

MFQDPQERPRKLPQLCTELQTTIHEIILECVYCKQQLLRREIGSDSEYEVDNSHQKAHSFVNHYISDPTYNSWRRQKQGISRA
 QAYSYTESDSGEPDHTTVTNSTSTQQGSLFRPKASRTPQNPSPSSQSTLYRPPSSLAPGSRAPAGFSSFV*

V

T5350 virus-host fusion protein 1 Transcript-4: [HPV16:97-226/Chr3:1899355002-189937565](#)Coding region: [HPV16 E6:104-226/Chr3:1899355002-1899355193](#) (318 nts / 105 aa)

Protein:

MFQDPQERPRKLPQLCTELQTTIHEIILECVYCKQQLLRREIQKKCNPRRKGQMSQSPVSSWVCMCPGSLGTSNLNLYALDM
 WLGLLRHKVGELEDERLQTLRSMAL*

Figure S7. Point mutations that cause missense mutations in E6 and E7 proteins from an integrated, expressible viral DNA and predicted virus-host fusion proteins from the transcribed virus-host chimeric transcripts in cervical cancer cell lines and tissues. (A-C) Point mutations in an integrated, expressible viral DNA identified in CaSki (**A**), SiHa (**B**), and HeLa (**C**) cells and ten cervical cancer tissues (**G-I**) identified by RNA-seq. Most mutations are silent. The missense mutations that cause amino acid changes in E6 and/or E7 proteins are indicated in red for CaSki (**D**), SiHa (**E**) and HeLa (**F**) cells, and eight cervical cancer tissues (**J-M**). The normal amino acid residues in the reference proteins are underlined. (**N-V**) Amino acid sequences of the predicted viral-host fusion proteins expressed from the chimeric virus-host fusion RNA transcripts detected in SiHa (**N, O**) and HeLa (**P, Q**) cells and in the cervical cancer tissues A1BJ (**R**), A2RE (**S**), T1074 (**T, U**) and T5350 (**V**). In HeLa cells, transcript-3 (**Q**) retains a truncated E1 ORF (528 aa) using a host stop codon at the integration junction site. A normal full size HPV18 E1 has 658 aa residues.