

1 **Genetic signatures for lineage/sublineage classification of**  
2 **HPV16, 18, 52 and 58 variants**

3 **Running Title: Genetic signatures for four HPV types**

4

5 **Zihua Ou,<sup>1,2</sup> Zigui Chen,<sup>3</sup> Yanping Zhao,<sup>1,2</sup> Haorong Lu,<sup>1,4,5</sup> Wei Liu,<sup>1,2,6</sup>**

6 **Wangsheng Li,<sup>1,4,5</sup> Chunyu Geng,<sup>1,7</sup> Guohai Hu,<sup>1,4,5</sup> Xiaman Wang,<sup>1,8</sup> Peidi**

7 **Ren,<sup>1,2</sup> Na Liu,<sup>1,8</sup> Shida Zhu,<sup>1,8,9</sup> Ling Lu,<sup>1,2</sup> Junhua Li<sup>1,2,6\*</sup>**

8

9 <sup>1</sup> BGI-Shenzhen, Shenzhen 518083, China.

10 <sup>2</sup> Shenzhen Key Laboratory of Unknown Pathogen Identification, BGI-Shenzhen,  
11 Shenzhen 518083, China.

12 <sup>3</sup> Department of Microbiology, Faculty of Medicine, The Chinese University of Hong  
13 Kong, Hong Kong SAR, China.

14 <sup>4</sup> China National GeneBank, BGI-Shenzhen, Shenzhen 518120, China.

15 <sup>5</sup> Shenzhen Key Laboratory of Environmental Microbial Genomics and Application,  
16 BGI-Shenzhen, Shenzhen 518083, China.

17 <sup>6</sup> School of Biology and Biological Engineering, South China University of Technology,  
18 Guangzhou, China.

19 <sup>7</sup> MGI, BGI-Shenzhen, Shenzhen, 518083, China.

20 <sup>8</sup> BGI Genomics, BGI-Shenzhen, Shenzhen, 518083, China.

21 <sup>9</sup> Shenzhen Engineering Laboratory for Innovative Molecular Diagnostics, BGI-  
22 Shenzhen, Shenzhen, 518120, China.

23 \*Corresponding author.

24 Word counts of the abstract: 147

25 Word counts of the text: 3463

26 **Abstract**

27 Increasing evidences indicate that high-risk HPV variants are heterogeneous in  
28 carcinogenicity and ethnic dispersion. In this work, we identified genetic signatures  
29 for convenient determination of lineage/sublineage of HPV16, 18, 52 and 58  
30 variants. Using publicly available genomes, we found that E2 of HPV16, L2 of  
31 HPV18, L1 and LCR of HPV52, and L2, LCR and E1 of HPV58 contain the proper  
32 genetic signature for lineage/sublineage classification. Sets of hierarchical signature  
33 nucleotide positions (SNPs) were further confirmed for high accuracy (>98%) by  
34 classifying HPV genomes obtained from Chinese females, which included 117  
35 HPV16 variants, 48 HPV18 variants 117 HPV52 variants and 89 HPV58 variants.  
36 The circulation of HPV variants posing higher cancer risk in Eastern China, such as  
37 HPV16 A4 and HPV58 A3, calls for continuous surveillance in this region. The  
38 marker genes and signature nucleotide positions may facilitate cost-effective  
39 diagnostic detections of HPV variants in clinical settings.

40

41

42 **Keywords:** Human papillomavirus; lineage; sublineage; classification; signature;  
43 genome; sequencing; Chinese female; cervical cancer; detection.

44

45

46

47

## 48 **Background**

49 Human papillomaviruses (HPVs) are a heterogeneous group of double-stranded  
50 DNA viruses mainly infecting epithelial surfaces of human beings. Currently, more  
51 than 200 types of human papillomaviruses have been identified, with the majority  
52 clustering into a limited set of phylogenetic genera (e.g., Alpha-, Beta-, Gamma-, Mu-  
53 and Nu-PV) [1]. The genital high-risk carcinogenic HPV types (e.g., HPV types 16  
54 and 18) causing cervical cancer are part of a monophyletic clade within the genus  
55 Alpha-PV [2,3]. Cervical cancer is the 4<sup>th</sup> most common cancer in women [4], with  
56 more than 0.5 million new cases and 0.2 million deaths occur worldwide annually [5].

57

58 Distinct HPV types are defined based on the L1 open reading frame (ORF) genetic  
59 sequence, with dissimilarity at least 10% to all other characterized viruses as a  
60 “novel” HPV type. Isolates of the same HPV type are referred to variant lineages and  
61 sublineages based on the complete genome nucleotide sequences differing  
62 approximately 1%-10% and 0.5%-1%, respectively [6]. HPV16, 18, 52 and 58 have  
63 been found to be the top four prevalent high-risk HPVs among Chinese females [7–  
64 10]. So far, multiple variant lineages and sublineages have been designated to  
65 HPV16 (A1-A4, B1-B4, C1-C4, D1-D4) [6,11,12], HPV18 (A1-A5, B1-B3, C) [6],  
66 HPV52 (A1-A2, B1-B2, C1-C2, D) [6] and HPV58 (A1-A3, B1-B2, C, D1-D2) [6]. HPV  
67 variants have different phenotypic characteristics including carcinogenicity and  
68 ethnic dispersion. For example, HPV16 A3 and A4 were linked with higher cancer  
69 risks in Asian populations comparing to the European prototype A1, while A4 and D  
70 displayed higher carcinogenesis in North Americans than A1 [12,13]. Moreover,  
71 HPV16 A1 and A2 tended to cause higher cancer risks in white Americans, and D2  
72 and D3 in Latino Americans [11]. Similarly, HPV18 B/C, HPV52 B and HPV58 A3

73 might be linked with higher cancer risks than the other variants of the same type [14–  
74 18].

75

76 Since infections with different HPV variants herald cancer risks differently,  
77 identification of lineage/sublineages will provide instructions on the triage and  
78 screening frequency of HPV-infected individuals. Herein, we used the publicly  
79 available genomes of HPV16, 18, 52 and 58 to pinpoint the marker genes and  
80 signature nucleotide positions for the determination of HPV variants, which would  
81 facilitate cost-effective diagnostic detections of HPV lineages and/or sublineages.  
82 Genomes of HPV16, 18, 52 and 58 types obtained from Chinese females were  
83 utilized for variant classification accuracy using the signature genes and sites.

84

## 85 **Methods**

### 86 **Data preparation**

87 Genome sequences for HPV16 (n=3,718), 18 (n=129), 52 (n=91) and 58 (n=172)  
88 were downloaded from NCBI nucleotide dataset by keyword search (Keyword:  
89 txid333760 for HPV16, txid333761 for HPV18, txid10618 for HPV52, txid10598 for  
90 HPV58; Species: Viruses; Molecular types: genomic DNA/RNA; Sequence type:  
91 Nucleotide; Release Date: from 0000/01/01 to 2019/07/25; Sequence length: from  
92 7,000 to 8,500; accessed on 25 July 2019). Reference genomes [6,19] were  
93 retrieved from PaVE [20]. Only unique genomes with over 95% of the HPV genome  
94 length (after excluding ambiguous sites) were selected. Moreover, all genes (E1, E2,  
95 E4, E5, E6, E7, L1 and L2) and LCR (long control region) sequences of the selected  
96 genomes had >70% coverage of the corresponding gene/region complete length.  
97 The sequences were aligned with MAFFT v7.427 and manually scrutinized and

98 edited with BioEdit v7.0.5. After exclusion of highly similar sequences, a total of  
99 2,695 genomes (HPV16, n=2,385; HPV18, n=99; HPV52, n=77; HPV58, n=134; see  
100 **Supplementary Table 1**) were retained for downstream analysis.

101

## 102 **Phylogeny reconstruction and data visualization**

103 Nucleotide substitution model test was conducted with IQ-TREE ModelFinder [21]  
104 and the best model identified was subsequently used for Maximum Likelihood (ML)  
105 phylogeny construction for HPV16, 18, 52 and 58 using the aforementioned  
106 datasets, with 1,000 ultrafast bootstrap pseudo-replications [22]. Lineage and  
107 sublineage assignments were conducted according to the classification criteria  
108 proposed by Burk et al. [6]. Representative phylogenies of HPV16, 18, 52 and 58  
109 were reconstructed with mean intra- and inter-group sequence distances calculated  
110 using R package *seqinr* and in-house R scripts. Phylogeny and the associated data  
111 were visualized with *ggtree* package in R [23].

112

## 113 **Pairwise distance matrix comparison**

114 Pairwise p-distances for HPV sequences were calculated with R package *ape*, with  
115 gaps deleted in a pairwise manner. Correlations between the DNA distance matrices  
116 of partial genomic sequences and the full genomes were analyzed using Mantel test  
117 with the R package *vegan* [24].

118

## 119 **Determination of signature nucleotide positions**

120 All sequences over 95% of the reference complete genomes were used to identify  
121 lineage- and sublineage-specific signature nucleotide positions as marker sites for  
122 variant classification. Positions with over 10% gaps or displaying a conservation rate

123 of over 98% across the alignments were excluded. All signature nucleotide positions  
124 were highly conserved in 99% of the sequences in the corresponding  
125 lineage/sublineage or genetic cluster. The HPV reference genomes used were  
126 downloaded from NCBI except for HPV16: HPV16, K02718 (downloaded from PaVE  
127 [20]); HPV18, AY262282; HPV52, X74481; HPV58, D90400. The distribution of the  
128 signature positions along the HPV genomes were further summarized based on a  
129 sliding window size of 1000bp and a step size of 500bp.

130

### 131 **HPV-positive cervical samples from Chinese women**

132 Exfoliated cervical cells were obtained from women participating in National Cervical  
133 Cancer Screening Program in Eastern China, including Anhui, Jiangsu, Shandong  
134 and Guangdong provinces. HPV DNA detections were conducted with BGI SeqHPV  
135 Kit (BGI-Shenzhen, China) [25,26]. The majority of subjects involving in this large-  
136 scale screening program displayed no clinical illness or slight inflammation in  
137 histopathological examination. Only samples from participants who consented to  
138 donate their residual sample for microbial investigation were selected, with their  
139 personal data anonymized. A total of 347 participants were recruited in this study. All  
140 the participants aged from 30 to 67 years old, with a median age of 48.

141

### 142 **HPV genome sequencing and assembly**

143 Probes covering the complete genomes of 18 HPV types (HPV6, 11, 16, 18, 31, 33,  
144 35, 39, 45, 51, 52, 56, 58, 59, 66, 68, 69 and 82) were designed by MyGenostics  
145 [27]. The extracted DNA was sheared to fragments around 250 bp in length, adapter-  
146 added, hybridized with probes at 65°C for 24h, and washed to remove uncaptured  
147 fragments. The DNA library was sequenced using a BGISEQ-500 platform for

148 paired-end 100bp reads (BGI-Shenzhen, China). Following demultiplexing, raw  
149 reads were trimmed with fastp [28] and deduplicated with BMap  
150 (<https://sourceforge.net/projects/bbmap/>). Quality-filtered reads were mapped to HPV  
151 reference genomes with BWA alignment tool [29]. Reads with both ends aligned to  
152 HPV16, 18, 52 and 58 were extracted and subject to *de novo* assembly with SPAdes  
153 3.12.0 [30].

154

### 155 **Verification of signature genes and nucleotide positions**

156 HPV16 (n=120), 18 (n=48), 52 (n=120) and 58 (n=93) genomes from Chinese  
157 females generated in this work were used to verify the hierarchical signature  
158 nucleotide positions for variant lineage/sublineage classification. The designations  
159 were conducted based on ML tree topologies and signature nucleotide position  
160 mapping algorithms. In brief, the genomes combining with the reference sequences  
161 **(Supplementary Table 2)** were aligned with MAFFT v7.427 and constructed for ML  
162 trees using IQ-TREE with 1,000 ultrafast bootstrap implementations [21,22]. The  
163 phylogenetic trees inferred from the ORFs and LCR were reconstructed following the  
164 same procedure, and the subsequent classification results were compared against  
165 those defined by complete genome data. R scripts were developed in-house to map  
166 the sequences against the signature nucleotide position sets to define variant  
167 lineages and sublineages.

168

### 169 **Ethical statement**

170 This study was reviewed and approved by the Institutional Review Board of Beijing  
171 Genomics Institute, Shenzhen, China (BGI-R071-1-T1 & BGI-R071-1-T2). All the

172 participants consented to the donation of their exfoliated cell samples and  
173 anonymized associated data for research purposes.

174

#### 175 **Data availability**

176 The data that support the findings of this study have been deposited into CNSA  
177 (CNGB Sequence Archive) of CNGBdb with accession number CNP0001117  
178 (<https://db.cngb.org/cnsa/>).

179

## 180 **Results**

### 181 **Signature genes for sublineage/lineage identification**

182 In order to pinpoint marker genes for HPV16, 18, 52 and 58 variant classification, we  
183 compared the consistency between the complete genomes, individual ORF/region  
184 (E1, E2, E4, E5, E6, E7, L1, L2, LCR), and the concatenated regions (4R:  
185 E1+E2+L1+L2, 5R: E1+E2+L1+L2+LCR, 8R: E1+E2+E4+E5+E6+E7+L1+L2, and  
186 9R: E1+E2+E4+E5+E6+E7+L1+L2+LCR) regarding to sequence diversity and  
187 lineage/sublineage assignment.

188

189 **Sequence diversity comparison.** All the partial genomic regions showed a positive  
190 correlation with full genome sequences based on pairwise distance matrix  
191 comparison ( $p < 0.05$ ) but yielded different correlation coefficient values. The distance  
192 matrices of the concatenated partial genomes usually display a high correlation with  
193 that of the full genome, with correlation coefficients ranging from 0.957 to 0.996  
194 (**Figure 1**). For individual genomic regions, E2 of HPV16, L2 of HPV18, and E1 of  
195 HPV58 displayed the highest correlation with its full genome. For HPV52, the



196 distance matrices of L1 and LCR displayed similar correlation with that of full  
197 genome depending on the nucleotide substitution model utilized.

198

199 **Lineage/sublineage classification.** Phylogeny of each genomic region for  
200 lineage/sublineage classification was reconstructed with the same setting for each  
201 HPV type. The assignments based on partial genomes (e.g., 4R, 5R, 8R and 9R)  
202 were 100% consistent with those by complete genomes for HPV18, 52 and 58  
203 ( $n < 150$ ). For HPV16, the classification consistencies of partial genomes were around  
204 99.8%, which may be due to the relatively abundant dataset ( $n = 2385$ ) and the  
205 variation in the non-coding regions of this type (**Figure 2**). We also found high  
206 accuracy of variant classification inferred from distinct ORF/region including HPV16  
207 E2 (98.49%), HPV18 L2 (100%), HPV52 L1 (100%) and LCR (100%), and HPV58 L2  
208 (100%), LCR (100%) and E1 (100%). The phylogenetic topologies inferred from E4,  
209 E5, E6 and E7 genes were severely distorted, given the limited variation these genes  
210 contain. Although some short ORFs, e.g., HPV18 E4 and E5, HPV52 E6 and E7,  
211 achieved lineage/sublineage classification results with high accuracies in our study,  
212 the assignment process relied heavily on the abundance of references and  
213 experience, and may be prone to human subjectivity and visual error. Therefore, the  
214 short ORFs (E4, E5, E6, E7) may be not suitable for phylogenetic classification.

215

#### 216 **Signature nucleotide positions for lineage/sublineage assignment**

217 Since HPVs have a relatively low mutation rate [31] and previous reports have  
218 shown certain positions to be population- or lineage- specific, we sought to identify  
219 signature nucleotide positions with lineage and/or sublineage fixation across the  
220 complete genome. Because no exclusive single nucleotide positions were

221 determinable for HPV16, 18, 52 and 58 variants, we used tree topologies and a  
222 hierarchical manner to identify lineage- and/or sublineage-specific nucleotide sites  
223 and patterns.

224

225 Certain sublineages, including HPV16 A1-3, B2-4, C1-4 and D1-4, HPV18 A1-4 and  
226 B1-3, and HPV52 A1-2 were merged together because of the limited sequence  
227 variation or inadequate distances these clusters contained (**Supplementary Figure**  
228 **1**). A total of 79, 133, 161 and 123 signature nucleotide positions were characterized  
229 for HPV16, 18, 52 and 58, respectively (**Figure 4** and **Supplementary Table 3**). At  
230 lineage level of HPV16, 35 sites were able to discriminate lineage A from lineage  
231 B/C/D, followed by 13 sites further discriminating B from C/D, and 17 discriminating  
232 C from D. At the sublineage level of HPV16, we found 12 positions for the  
233 discrimination of A1-3 and A4, and 2 for B1 and B2-4. The hierarchical structure of  
234 signature nucleotide positions for HPV18, 52, and 58 can be interpreted similarly  
235 (**Figure 3**).

236

237 Using a sliding window size of 1000bp and a step size of 500bp, we explored the  
238 distribution of the signature sites along the viral genomes. Results showed that the  
239 1001<sup>st</sup>-2000<sup>th</sup> genomic positions of HPV16 (the 5' terminal region of E1), the 3501<sup>st</sup>-  
240 4500<sup>th</sup> of HPV18 (flanking E2, E4, E5 and L2), the 6501<sup>st</sup>-7500<sup>th</sup> of HPV52 (the 3'  
241 terminal region of L1) and the 7001<sup>st</sup>-7824<sup>th</sup> of HPV58 (LCR) may contain the most  
242 sufficient signature sites to distinguish all the hierarchical levels for each HPV type.  
243 (**Figure 4**).

244

245 **Genetic diversity of verification dataset**

246 To verify the signature of marker genes and sites in variant classification, we  
247 accessed HPV16, 18, 52 and 58 genomes from Chinese women who participated in  
248 the National Cervical Cancer Screening Program. Most of the participants were from  
249 Eastern China. All sequences are >95% coverage in size of the complete genomes.  
250 Based on tree topologies and distance threshold, 116 out of 117 HPV16 sequences  
251 were unambiguously assigned as lineages A (A1-2=25, A3=24, A4=66) and D  
252 (D3=1) (**Figure 5, Supplementary Table 4**). The high prevalence of HPV16 A4  
253 variants in Eastern Chinese women was consistent with previous reports in Eastern  
254 Asian [18]. HPV16 A4 was also linked with higher cancer risks in Asian populations  
255 than other variants [12,13]. All HPV18 genomes belonged to lineage A (A1=37,  
256 A3=1, A4=10) (**Figure 5, Supplementary Table 4**), which were consistent with  
257 previous reports on HPV18 sublineage distribution in China [32,33]. The majority of  
258 HPV52 and HPV58 sequences were lineages B (111/117, 94.9%) and A (88/89,  
259 98.9%), respectively, that were further divided into HPV52 B2 (n=111) and HPV58  
260 A1 (n=51), A2 (n=25) and A3 (n=12) (**Figure 5, Supplementary Table 4**). Rest of  
261 HPV52 sequences belonged to A1 (n=2), C2 (n=3) and D (n=1), and HPV58 to B2  
262 (n=1). HPV52 lineages B and C were common in Asian countries, with B the most  
263 prevalent in China [34–36]. However, the cancer risk of HPV52 B variants was  
264 reported to be lower than lineage C [16]. Lineage A of HPV58 was found to be  
265 globally distributed and was the most prevalent in Asian females [16,33,35,37,38].  
266 Moreover, HPV58 sublineage A3 might pose higher cancer risk than other variants  
267 [17,18]. The prevalence of HPV variants associated with higher cancer risk (e.g.,  
268 HPV16 A4, HPV58 A3) in Eastern China called for continuous surveillance on female  
269 populations in this region.  
270

271 **Verification of HPV lineage classifications with signature regions and**  
272 **nucleotide positions**

273 Using the assembled HPV sequences from Eastern Chinese women, we further  
274 confirmed that HPV16 E2, HPV18 L2, HPV52 L1 and LCR, and HPV58 LCR  
275 contained sufficient variation information to assign the target genes/region with  
276 proper lineages and/or sublineages, consistent with the classification by the  
277 complete genomes (**Table 1, Supplementary Table 4**). In addition, HPV58 E1 and  
278 L2 reached 98.88% accuracy in variant classification, except for one sequence with  
279 ambiguous assignment. Classification results by signature nucleotide positions also  
280 showed high consistency with those by the complete genomes, with 100% and >98.9%  
281 accuracies in variant classification at lineage and sublineage levels, respectively. It's  
282 worth noting that one HPV58 A2 sequence based on the complete genome  
283 assignment was grouped to A1 since only one variation was discriminative for A1, A2  
284 and A3. Hence, mutations at certain individual position may affect the accuracy of  
285 sublineage assignment.

286

287 **Discussion**

288 It has been recommended to use the complete genomes to identify HPV variant  
289 lineages and sublineages [6]. However, ambiguous assignment may arise when the  
290 complete genomic sequences are not available in clinical settings or in developing  
291 areas. Lineage fixation of genetic changes in one gene/region highly correlated with  
292 other changes within genomes from the same lineage and sublineage is observed  
293 throughout HPVs and may represent adapted variations in natural selection with  
294 different phenotypic characteristics and carcinogenicity. In this study, we sought to  
295 identify genes and sites with lineage- and sublineage-specific significance for HPV

296 variant characterization in large-scale screening program. Our data indicated that  
297 marker genes (e.g., HPV16 E2, HPV18 L2) as well as signature nucleotide positions  
298 proved high accuracy in lineage/sublineage classification, which was further verified  
299 by assembled HPV sequences from Eastern Chinese women.

300

301 Hotspots of the genetic signature position may be chosen as the target regions for  
302 the developing of rapid classification assays, such as the 1001<sup>st</sup>-2000<sup>th</sup> genomic  
303 positions of HPV16, the 3501<sup>st</sup>-4500<sup>th</sup> of HPV18, the 6501<sup>st</sup>-7500<sup>th</sup> of HPV52 and the  
304 7001<sup>st</sup>-7824<sup>th</sup> of HPV58 (**Figure 3, Figure 4, Supplementary Table 3**). Because  
305 most detection methods, such as qPCR, utilize short genomic regions as detection  
306 targets, we recommend using multiple genetic regions to achieve optimal detection  
307 accuracy. Due to the uneven distribution of the genetic signature positions, single  
308 regions may not be able to provide high-resolution classification. For example, while  
309 the 1001<sup>st</sup>-2000<sup>th</sup> positions of HPV16 may be able to distinguish all the classification  
310 levels, this region contains much less information for the separation of Ax  
311 sublineages than the 3001<sup>th</sup>-4000<sup>th</sup> region. Therefore, selection of multiple genetic  
312 regions pertinent to the local HPV diversity and detection methods such as multiplex-  
313 PCR may facilitate the cost-effective classification of HPVs to lineage/sublineage  
314 levels, which would help promote large-scale epidemiological study on the  
315 carcinogenesis of HPV variants.

316

317 This study has several limitations: 1) Though we have already gathered all available  
318 genomes from public database, the global diversity of HPVs may not be  
319 comprehensively covered. The scarcity of high-quality full genomes for certain  
320 lineages and sublineages, such as HPV16 lineage C and D, hinders the exploration

321 of cluster-specific genetic signatures. 2) Due to the limited sample sizes of this study,  
322 the genetic diversities of high-risk HPVs in China remain to be explored. 3) Because  
323 the clinical information of the surveyed sequences was not available in this work, the  
324 association between disease statuses and specific variants/variations in China  
325 remains elusive. Such defects call for genetic studies on HPV-infected individuals  
326 from more diverse geographic regions with sufficient clinical records to enhance our  
327 understandings of specific variants that may pose significant effects on cervical  
328 health.

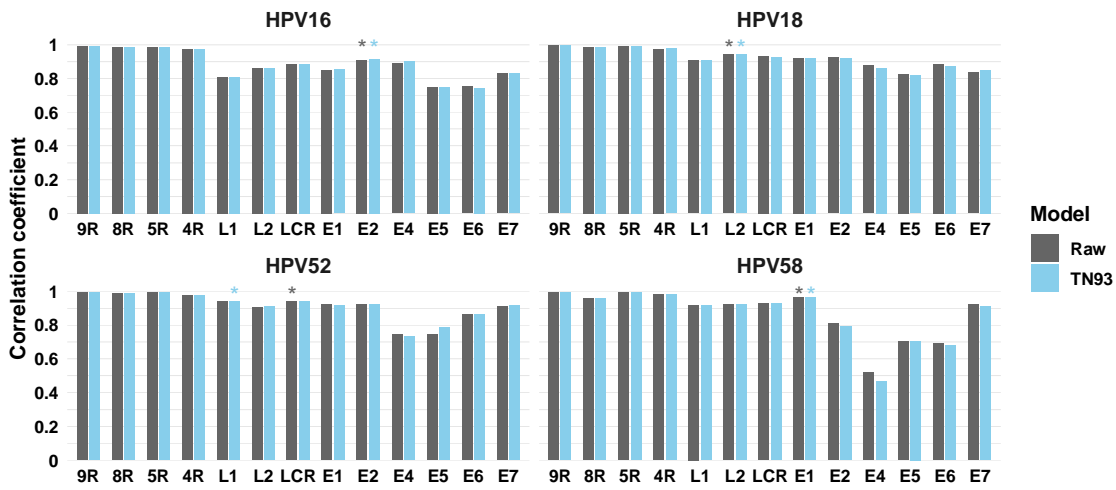
329

330 The sublineage distributions of the four HPV types in China generated in this study  
331 were consistent with previous reports [16,33–40]. However, with the intensification of  
332 globalization, the genetic diversity of HPV in China and other parts of the planet  
333 remains dynamic and requires continuous surveillance. The genetic signatures  
334 characterized by this study may provide valuable references for the design of cheap  
335 and fast detection assay to classify the four high-risk HPV types in Eastern China.  
336 Nevertheless, high-quality genomes were still scarce for many HPV types except  
337 HPV16. With the extensive application of whole genome sequencing in HPV  
338 research, the classification power using signature regions or nucleotide positions  
339 could be further increased.

340

341

342 **Figures and tables**



343

344 **Figure 1: Correlation of the genetic distance matrix between full genome and**

345 **partial genomic sequences of HPV16, 18, 52 and 58.** Raw pairwise p-distances

346 and genetic distances based on the TN93 model for HPV sequences were calculated

347 using R package *ape*, with gaps deleted in a pairwise manner. Correlations between

348 the DNA distance matrices of partial genomic sequences and the full genomes were

349 explored using Mantel test with Spearman correlation method. Asterisks indicate the

350 individual gene that showed the highest correlation with full genome under the

351 corresponding substitution model for each HPV type. Abbreviations: FG, full

352 genome; LCR, long control region; 4R, partial genome concatenated with 4 genetic

353 regions: E1+E2+L2+L1; 5R, partial genome concatenated with 5 genetic regions:

354 E1+E2+L2+L1+LCR; 8R, partial genome concatenated with 8 genetic regions:

355 E6+E7+E4+E5+E1+E2+L1+L2; 9R, partial genome concatenated with 9 genetic

356 regions: E6+E7+E4+E5+E1+E2+L1+L2+LCR.

357

358

359

Sequence	Sublineage				Lineage			
	HPV16	HPV18	HPV52	HPV58	HPV16	HPV18	HPV52	HPV58
9R	99.79	100	100	100	100	100	100	100
8R	99.75	100	100	100	100	100	100	100
5R	99.83	100	100	100	100	100	100	100
4R	99.79	100	100	100	100	100	100	100
L1	86.71	98.99	100	94.78	99.96	100	100	97.76
L2	93.71	100	96.1	100	99.92	100	100	100
LCR	93.54	96.97	100	100	99.83	100	100	100
E1	97.36	98.99	98.7	100	99.87	100	100	100
E2	98.49	98.99	98.7	92.54	99.96	100	100	97.76
E4	94.88	84.85	88.31	87.31	99.29	100	88.31	96.27
E5	77.15	62.63	93.51	82.84	95.3	100	96.1	93.28
E6	78.32	82.83	93.51	77.61	99.54	98.99	100	97.76
E7	84.32	83.84	94.81	91.04	97.65	95.96	100	97.76

360

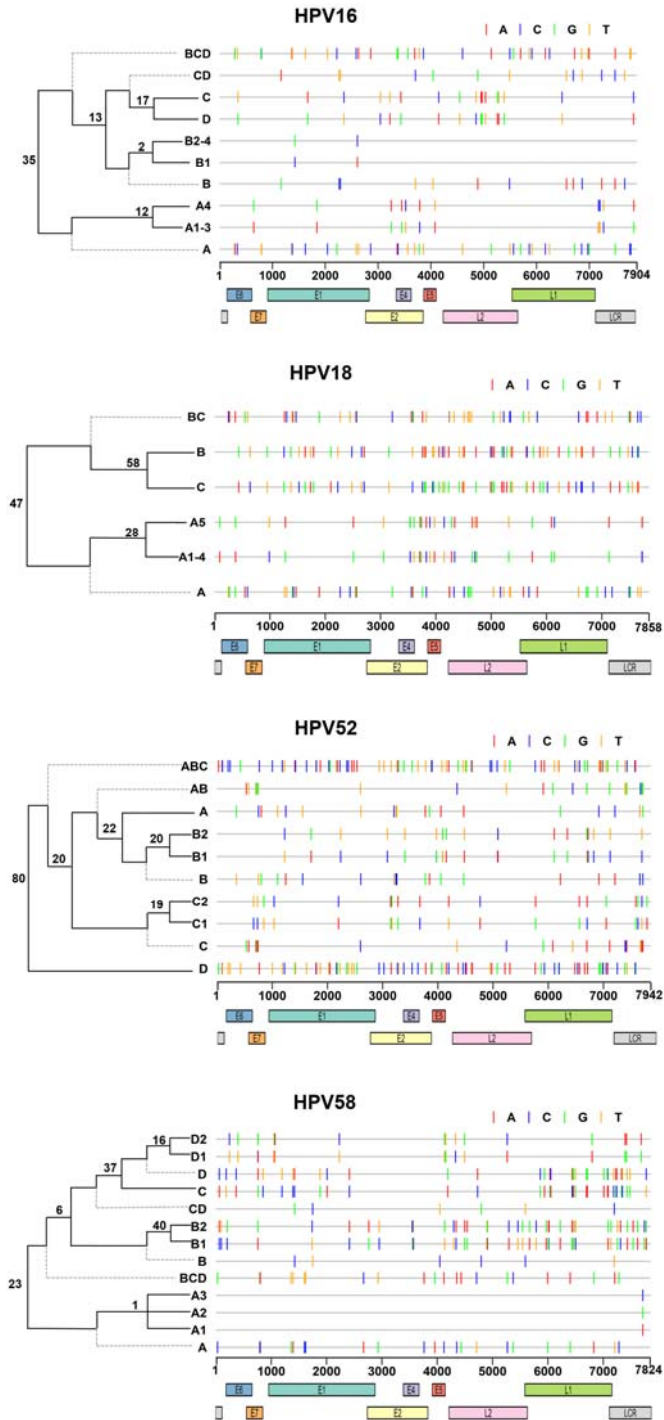
361 **Figure 2: Accuracies of lineage/sublineage assignments using different part of**  
 362 **the genomes.** Values indicate the percentage consistency in lineage/sublineage  
 363 assignment using the corresponding genomic regions comparing against the  
 364 assignment results using full genomes. Cell colors: yellow, 100%; red: 95~100%;  
 365 white, <95%. Abbreviations are the same as **Figure 1.**

366

367

368

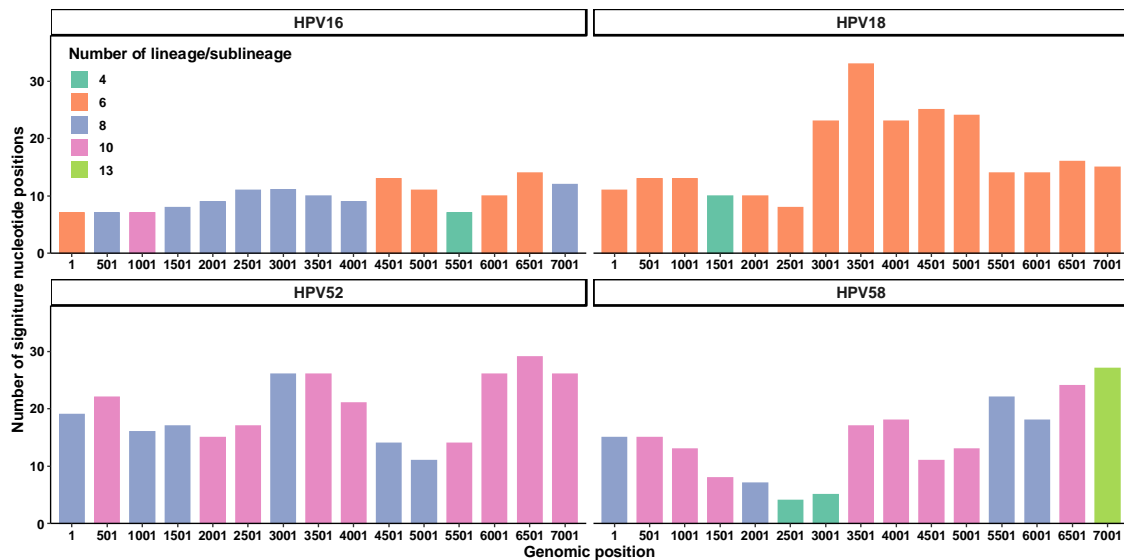




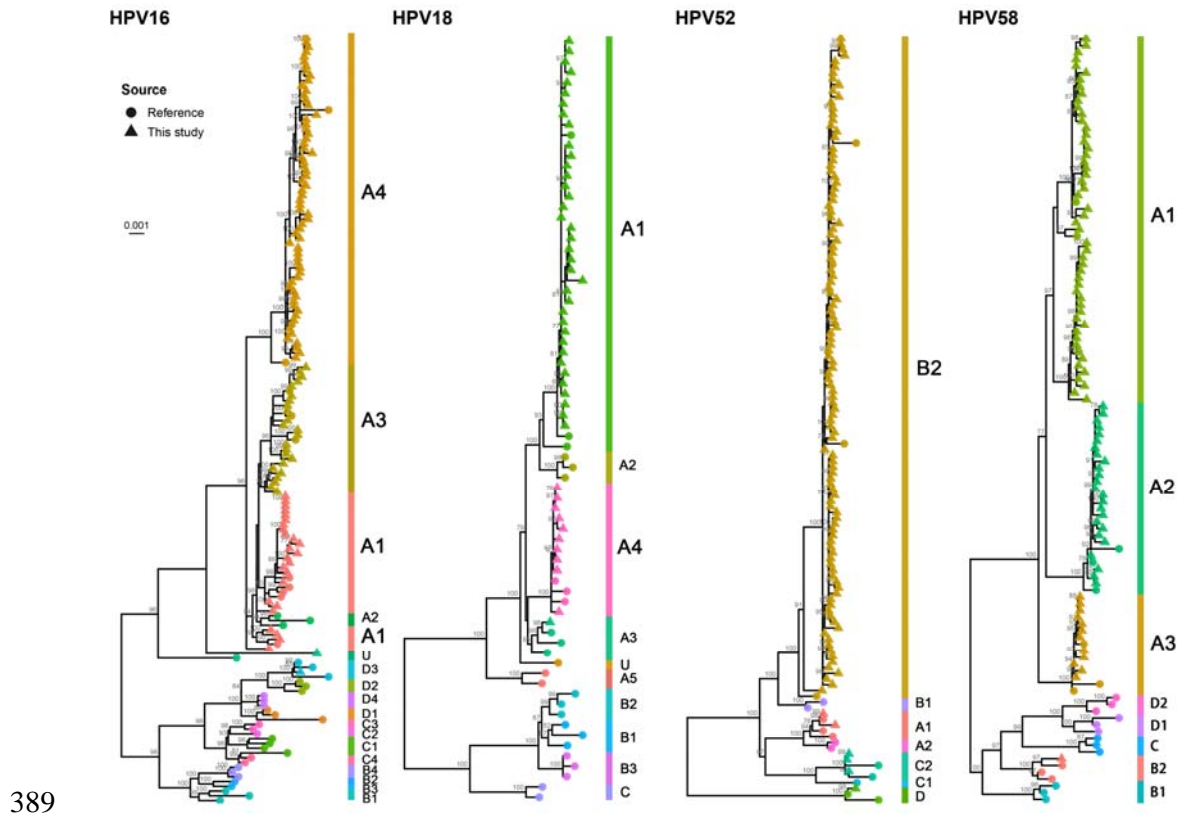
369

370 **Figure 3. Hierarchical discrimination of HPV lineages and sublineages with**  
371 **signature nucleotide positions.** The number of nucleotide positions determined for  
372 discriminating variant lineages and/or sublineages are displayed on the common  
373 node of two lineage/sublineage level. For example, in the panel of HPV16, 35 is the

374 number of signature sites to distinguish lineage A and BCD, and 12 to distinguish  
375 A1-3 and A4. The hierarchical structure was not scaled to genetic distances. The  
376 signature positions were color-labelled for each classification level based on the  
377 genomic organization of the reference genome of each HPV type, with the ORF/LCR  
378 locations indicated by color bars.  
379



380 **Figure 4. Genetic hotspots for signature nucleotide positions of HPV16, 18, 52**  
381 **and 58 types.** Using a sliding window size of 1000bp and a step size of 500bp,  
382 distribution of the signature nucleotide positions identified for the four HPV types  
383 were summarized. The X axis indicates the start position for each sliding window.  
384 The Y axis indicates the number of distinct signature positions for each window  
385 based on the corresponding reference genomes. Color code indicates the number of  
386 lineage/sublineage identifiable within each sliding window. The color legend applies  
387 to all types.  
388



390 **Figure 5. ML phylogenies of HPV16, 18, 52 and 58 sequenced in this study.** The  
391 ML phylogenies were constructed with full genomes using IQ-TREE, implementing  
392 1000 ultra-fast bootstrap tests. Bootstrap values of over 70 were displayed on nodes.  
393 Tips labelled by solid circles are reference genomes obtained from NCBI, while  
394 those labelled by triangles were sequenced in this study. Lineage/sublineages are  
395 indicated by colors and the name for each lineage/sublineage is presented with the  
396 same color. U: Undefined.

397

398 **Table 1: Accuracies of lineage/sublineage assignments using E1, E2, L1, L2**  
399 **and LCR with verification dataset.** Values indicate the percentage consistency in  
400 lineage/sublineage assignment using the corresponding sequences comparing  
401 against the assignment results using full genomes.

Region	Sublineage				Lineage			
	HPV16	HPV18	HPV52	HPV58	HPV16	HPV18	HPV52	HPV58
E1	78.63	100	99.15	98.88	100	100	100	98.88
E2	100	100	99.15	73.03	100	100	100	73.03
L1	48.72	100	100	97.75	100	100	100	97.75
L2	49.57	100	97.44	98.88	100	100	100	98.88
LCR	94.02	100	100	98.88	100	100	100	100

402

403

404

405

## 406 **Supplementary Data**

407 **Supplementary Table 1: Genomes of HPV16, 18, 52, 58 downloaded from**  
408 **public database.**

409 **Supplementary Table 2: Reference HPV complete genomes selected for the**  
410 **four HPV types.**

411 **Supplementary Table 3: Hierarchical SNPs for lineage/sublineage**  
412 **classification of HPV16, 18, 52 and 58.**

413 **Supplementary Table 4: Classification of the HPV16, 18, 52, 58 genomes**  
414 **generated by this study.**

415 **Supplementary Figure 1: Representative phylogenies of HPV16, 18, 52 and 58.**

416 Using publicly available complete genomes, ML trees were reconstructed for the four  
417 HPV types and each sequence was classified to sublineage level based on tree  
418 topologies. The representative phylogenies were reconstructed based on mean intra-  
419 and inter-group percentage differences of clades, which are simultaneously  
420 displayed on the right panel of each tree. The number of HPV genomes in each  
421 clade is displayed in tip labels. The sizes of black circles on tree tips indicate the  
422 relative sequence abundance of the corresponding clade. Abbreviations: Ax and Cx,  
423 the strains belonged to lineage A and C, but could not be assigned to any existing  
424 sublineages; U, (i.e., Undefined), the strains were not assigned to any existing  
425 lineage. Number of genomes used: HPV16, n=2,385; HPV18, n=99; HPV52, n=77;  
426 HPV58, n=134.

427

428

429 **Notes**

430 **Author contributions**

431 J.L., Z.O. and Z.C. designed the study. J.L., N.L., L.L., S.Z. and X.W. coordinated  
432 sample collection. H.L., W.L., G.H., C.G. and P.R. conducted viral genome  
433 sequencing. Z.O. and W.L. conducted data analysis. Z.O. wrote the manuscript. Z.C.,  
434 J.L., Y.Z. and L.L. provided critical revision of the manuscript.

435

436 **Declarations of interest**

437 The authors declare no conflict of interest.

438

439 **Funding**

440 This work received funding by Guangdong Provincial Key Laboratory of Genome  
441 Read and Write (2017B030301011) and Shenzhen Engineering Laboratory for  
442 Innovative Molecular Diagnostics (DRC-SZ[2016]884).

443

444 **Information of corresponding author**

445 Junhua Li; PhD; Tel: 0086-13929566296; E-mail: [lijunhua@genomics.cn](mailto:lijunhua@genomics.cn); Affiliations:  
446 Shenzhen Key Laboratory of Unknown Pathogen Identification, BGI-Shenzhen,  
447 Shenzhen 518083, China & School of Biology and Biological Engineering, South  
448 China University of Technology, Guangzhou, China; Address: China National  
449 GeneBank, Dapeng New District, Shenzhen 518120, China.

450

451 **Related Meetings**

452 Part of the lineage/sublineage classification results has been presented in EUROGIN  
453 2019 (December 4-7, 2019, Monaco). Poster Number: #0397; Title: Preliminary

454 Analysis on the Genetic Diversities of High-risk Human Papillomaviruses in Chinese  
455 Women.

456

#### 457 **Acknowledgments**

458 We warmly thank Dr. Houshun Zhu, Mr. Xianchu Cai, Ms. Jieyao Yu, Ms. Wei Zhou,  
459 Mr. Hongcheng Zhou, Miss. Yi Wang and Miss. Di Wu for their assistance in viral  
460 genome sequencing and dataset preparation. We thank China National GeneBank  
461 for providing sequencing service for this project. Thanks also to all the authors for  
462 contributing HPV genomes to NCBI GenBank. Last but not the least, the authors are  
463 grateful to the inspirational communication from Miss Shanshan Mo, Miss Feiyun Ou,  
464 Mr. Geer Xi and Miss Xiaobai Zhong.

465

#### 466 **References**

- 467 1. hpvcenter – International Human Papillomavirus Reference Center [Internet]. [cited  
468 2020 Jul 14]. Available from: <https://www.hpvcenter.se/>
- 469 2. Sanjose S de, Quint WG, Alemany L, et al. Human papillomavirus genotype attribution  
470 in invasive cervical cancer: a retrospective cross-sectional worldwide study. *The Lancet*  
471 *Oncology*. **2010**; 11(11):1048–1056.
- 472 3. IARC Monographs on the Identification of Carcinogenic Hazards to Humans –  
473 INTERNATIONAL AGENCY FOR RESEARCH ON CANCER [Internet]. [cited 2020  
474 Jul 14]. Available from: <https://monographs.iarc.fr/>
- 475 4. Arbyn M, Weiderpass E, Bruni L, et al. Estimates of incidence and mortality of cervical  
476 cancer in 2018: a worldwide analysis. *The Lancet Global Health*. **2020**; 8(2):e191–e203.
- 477 5. Cancer Statistics Review, 1975-2017 - SEER Statistics [Internet]. [cited 2020 Jul 14].  
478 Available from: [https://seer.cancer.gov/csr/1975\\_2017/](https://seer.cancer.gov/csr/1975_2017/)
- 479 6. Burk RD, Harari A, Chen Z. Human papillomavirus genome variants. *Virology*. **2013**;  
480 445(1–2):232–243.
- 481 7. Long W, Yang Z, Li X, et al. HPV-16, HPV-58, and HPV-33 are the most carcinogenic  
482 HPV genotypes in Southwestern China and their viral loads are associated with severity  
483 of premalignant lesions in the cervix. *Viol J*. **2018**; 15(1):94.

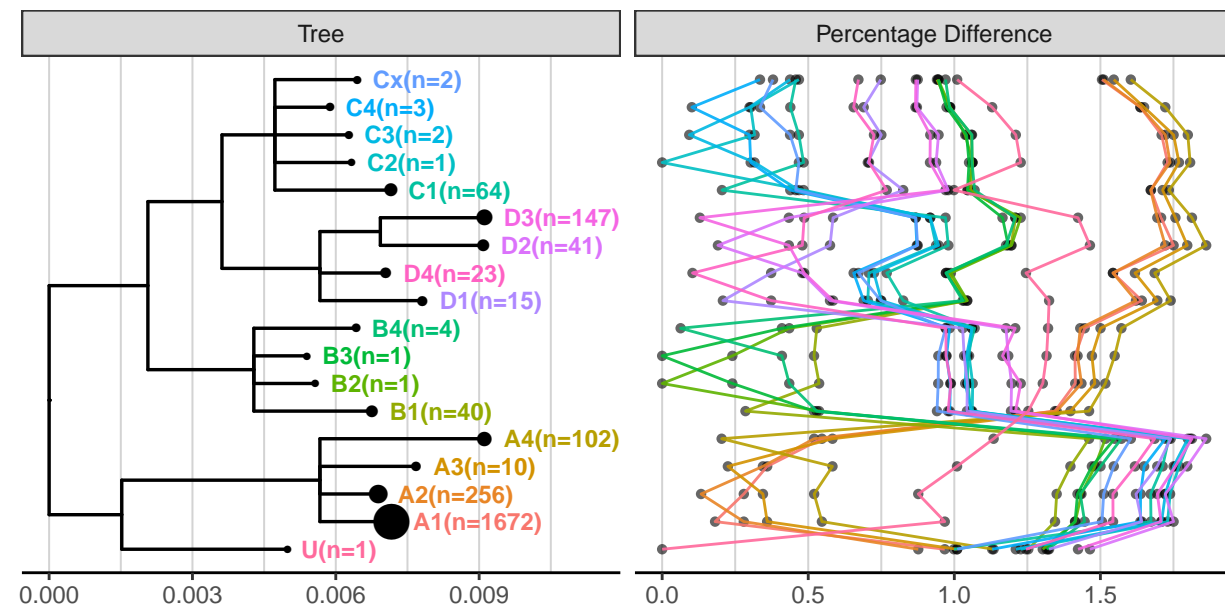
- 484 8. Lu J, Shen G, Li Q, Chen X, Ma C, Zhu T. Genotype distribution characteristics of  
485 multiple human papillomavirus in women from the Taihu River Basin, on the coast of  
486 eastern China. *BMC Infect Dis.* **2017**; 17(1):226.
- 487 9. Wang R, Guo X, Wisman GBeaA, et al. Nationwide prevalence of human  
488 papillomavirus infection and viral genotype distribution in 37 cities in China. *BMC*  
489 *Infect Dis.* **2015**; 15(1):257.
- 490 10. Cao D, Zhang S, Zhang Q, et al. Prevalence of high-risk human papillomavirus  
491 infection among women in Shaanxi province of China: A hospital-based investigation. *J*  
492 *Med Virol.* **2017**; 89(7):1281–1286.
- 493 11. Mirabello L, Yeager M, Cullen M, et al. HPV16 Sublineage Associations With  
494 Histology-Specific Cancer Risk Using HPV Whole-Genome Sequences in 3200 Women.  
495 *JNCI J Natl Cancer Inst.* **2016**; 108(9):djwt100.
- 496 12. Clifford GM, Tenet V, Georges D, et al. Human papillomavirus 16 sub-lineage dispersal  
497 and cervical cancer risk worldwide: Whole viral genome sequences from 7116 HPV16-  
498 positive women. *Papillomavirus Research.* **2019**; 7:67–74.
- 499 13. Hang D, Yin Y, Han J, et al. Analysis of human papillomavirus 16 variants and risk for  
500 cervical cancer in Chinese population. *Virology.* **2016**; 488:156–161.
- 501 14. Sichero L, Ferreira S, Trottier H, et al. High grade cervical lesions are caused  
502 preferentially by non-European variants of HPVs 16 and 18. *Int J Cancer.* **2007**;  
503 120(8):1763–1768.
- 504 15. Chen AA, Gheit T, Franceschi S, Tommasino M, Clifford GM. Human Papillomavirus  
505 18 Genetic Variation and Cervical Cancer Risk Worldwide. Banks L, editor. *J Virol.*  
506 **2015**; 89(20):10680–10687.
- 507 16. Chang Y-J, Chen H-C, Lee B-H, et al. Unique variants of human papillomavirus  
508 genotypes 52 and 58 and risk of cervical neoplasia. *Int J Cancer.* **2011**; 129(4):965–973.
- 509 17. Chen Z, Schiffman M, Herrero R, et al. Evolution and Taxonomic Classification of  
510 Human Papillomavirus 16 (HPV16)-Related Variant Genomes: HPV31, HPV33,  
511 HPV35, HPV52, HPV58 and HPV67. Chan KYK, editor. *PLoS ONE.* **2011**;  
512 6(5):e20183.
- 513 18. Chan PKS. Association of Human Papillomavirus Type 58 Variant With the Risk of  
514 Cervical Cancer. *CancerSpectrum Knowledge Environment.* **2002**; 94(16):1249–1253.
- 515 19. Mirabello L, Clarke M, Nelson C, et al. The Intersection of HPV Epidemiology,  
516 Genomics and Mechanistic Studies of HPV-Mediated Carcinogenesis. *Viruses.* **2018**;  
517 10(2):80.
- 518 20. Van Doorslaer K, Li Z, Xirasagar S, et al. The Papillomavirus Episteme: a major update  
519 to the papillomavirus sequence database. *Nucleic Acids Res.* **2017**; 45(D1):D499–D506.
- 520 21. Nguyen L-T, Schmidt HA, Haeseler A von, Minh BQ. IQ-TREE: A Fast and Effective  
521 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular*  
522 *Biology and Evolution.* **2015**; 32(1):268–274.



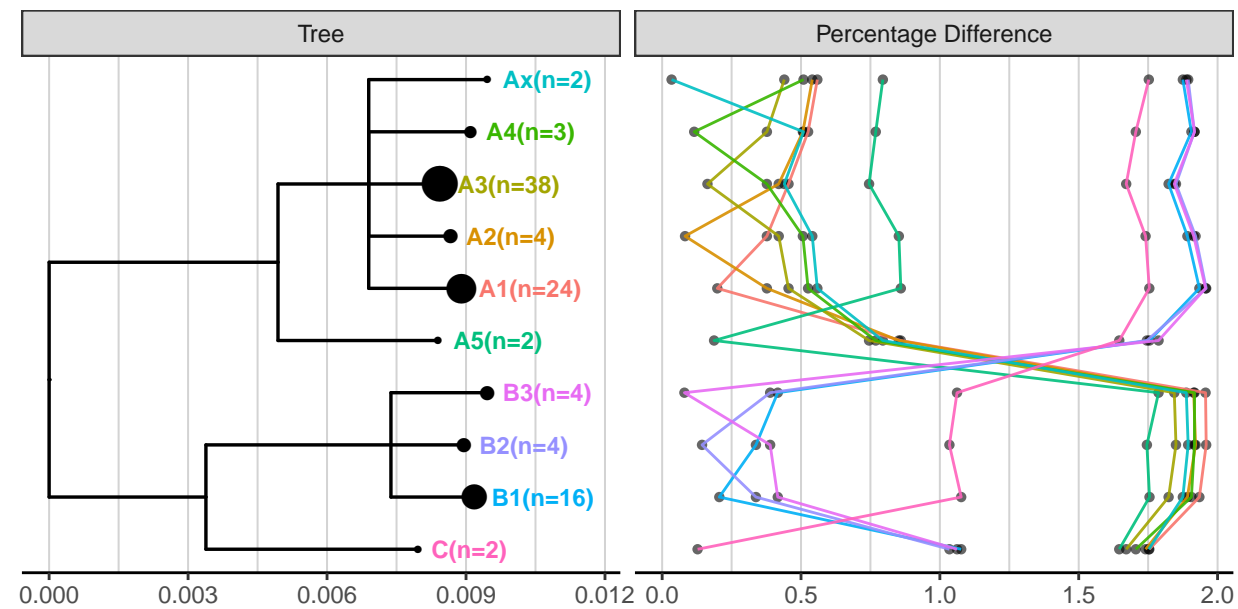
- 523 22. Hoang DT, Chernomor O, Haeseler A von, Minh BQ, Vinh LS. UFBoot2: Improving  
524 the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution*. **2018**;  
525 35(2):518–522.
- 526 23. Yu G. Using ggtree to Visualize Data on Tree-Like Structures. *Current Protocols in*  
527 *Bioinformatics* [Internet]. **2020** [cited 2020 Mar 6]; 69(1). Available from:  
528 <https://onlinelibrary.wiley.com/doi/abs/10.1002/cpbi.96>
- 529 24. Dixon P. VEGAN, a package of R functions for community ecology. *Journal of*  
530 *Vegetation Science*. **2003**; 14(6):927–930.
- 531 25. Yi X, Zou J, Xu J, et al. Development and Validation of a New HPV Genotyping Assay  
532 Based on Next-Generation Sequencing. *American Journal of Clinical Pathology*. **2014**;  
533 141(6):796–804.
- 534 26. Yi X, Li J, Yu S, et al. A New PCR-Based Mass Spectrometry System for High-Risk  
535 HPV, Part I. *American Journal of Clinical Pathology*. **2011**; 136(6):913–919.
- 536 27. Hu Z, Zhu D, Wang W, et al. Genome-wide profiling of HPV integration in cervical  
537 cancer identifies clustered genomic hot spots and a potential microhomology-mediated  
538 integration mechanism. *Nat Genet*. **2015**; 47(2):158–163.
- 539 28. Chen S, Zhou Y, Chen Y, Gu J. fastp: an ultra-fast all-in-one FASTQ preprocessor.  
540 *Bioinformatics*. **2018**; 34(17):i884–i890.
- 541 29. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
542 transform. *Bioinformatics*. **2009**; 25(14):1754–1760.
- 543 30. Bankevich A, Nurk S, Antipov D, et al. SPAdes: A New Genome Assembly Algorithm  
544 and Its Applications to Single-Cell Sequencing. *Journal of Computational Biology*.  
545 **2012**; 19(5):455–477.
- 546 31. Rector A, Lemey P, Tachezy R, et al. Ancient papillomavirus-host co-speciation in  
547 Felidae. *Genome Biol*. **2007**; 8(4):R57.
- 548 32. Xu H-H, Zheng L-Z, Lin A-F, Dong S-S, Chai Z-Y, Yan W-H. Human papillomavirus  
549 (HPV) 18 genetic variants and cervical cancer risk in Taizhou area, China. *Gene*. **2018**;  
550 647:192–197.
- 551 33. Liu Y, Pan Y, Gao W, Ke Y, Lu Z. Whole-Genome Analysis of Human Papillomavirus  
552 Types 16, 18, and 58 Isolated from Cervical Precancer and Cancer Samples in Chinese  
553 Women. *Sci Rep*. **2017**; 7(1):263.
- 554 34. Zhang C, Park J-S, Grce M, et al. Geographical Distribution and Risk Association of  
555 Human Papillomavirus Genotype 52-Variant Lineages. *J Infect Dis*. **2014**;  
556 210(10):1600–1604.
- 557 35. Tenjimbayashi Y, Onuki M, Hirose Y, et al. Whole-genome analysis of human  
558 papillomavirus genotypes 52 and 58 isolated from Japanese women with cervical  
559 intraepithelial neoplasia and invasive cervical cancer. *Infect Agents Cancer*. **2017**;  
560 12(1):44.

- 561 36. Choi YJ, Ki EY, Zhang C, et al. Analysis of Sequence Variation and Risk Association  
562 of Human Papillomavirus 52 Variants Circulating in Korea. Liu X, editor. PLoS ONE.  
563 **2016**; 11(12):e0168178.
- 564 37. Chan PKS, Luk ACS, Park J-S, et al. Identification of Human Papillomavirus Type 58  
565 Lineages and the Distribution Worldwide. *The Journal of Infectious Diseases*. **2011**;  
566 203(11):1565–1573.
- 567 38. Liu J, Lu Z, Wang G, et al. Variations of human papillomavirus type 58 E6, E7, L1  
568 genes and long control region in strains from women with cervical lesions in Liaoning  
569 province, China. *Infection, Genetics and Evolution*. **2012**; 12(7):1466–1472.
- 570 39. Sun M, Gao L, Liu Y, et al. Whole Genome Sequencing and Evolutionary Analysis of  
571 Human Papillomavirus Type 16 in Central China. Zheng Z-M, editor. PLoS ONE. **2012**;  
572 7(5):e36577.
- 573 40. Guo Y, Hu J, Zhu L, et al. Physical Status and Variant Analysis of Human  
574 Papillomavirus 16 in Women from Shanghai. *Gynecol Obstet Invest*. **2016**; 81(1):61–70.
- 575

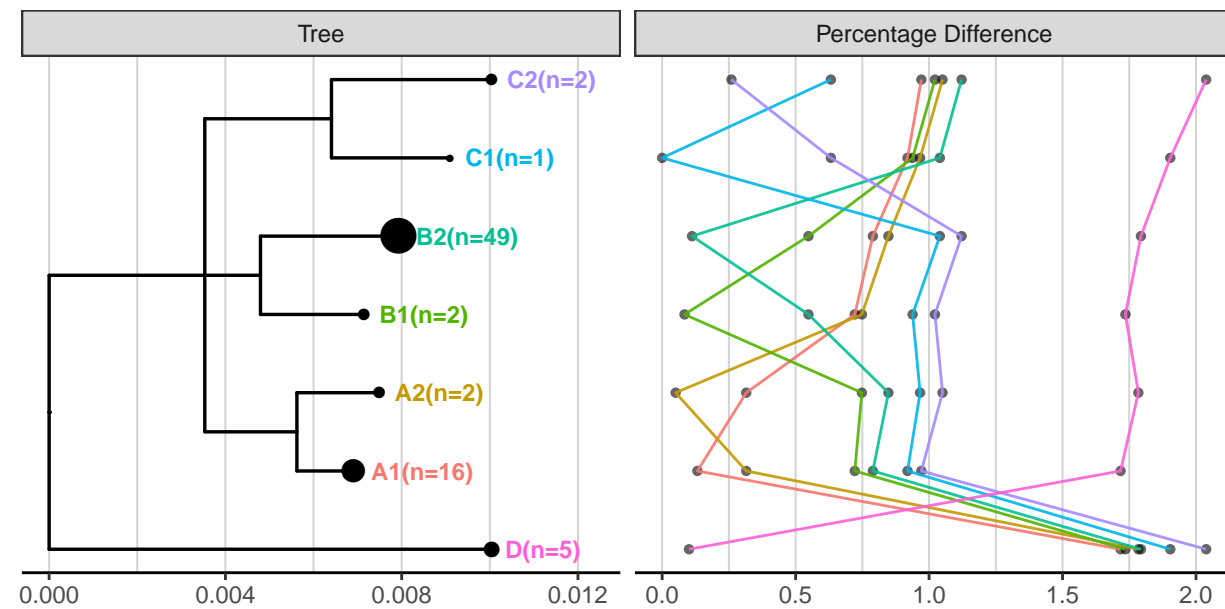
HPV16 (n=2385)



HPV18 (n=99)



HPV52 (n=77)



HPV58 (n=134)

