

1 **Mixed cytomegalovirus genotypes in HIV positive mothers show compartmentalization and**  
2 **distinct patterns of transmission to infants.**

3

4 Juanita Pang<sup>1¶</sup>, Jennifer A. Slyker<sup>2¶</sup>, Sunando Roy<sup>1</sup>, Josephine Bryant<sup>1</sup>, Claire Atkinson<sup>3</sup>, Juliana  
5 Cudini<sup>1</sup>, Carey Farquhar<sup>4</sup>, Paul Griffiths<sup>3</sup>, James Kiarie<sup>5</sup>, Sofia Morfopoulou<sup>1</sup>, Alison C. Roxby<sup>4</sup>,  
6 Helena Tutil<sup>1</sup>, Rachel Williams<sup>1</sup>, Soren Gantt<sup>6</sup>, Richard A. Goldstein<sup>1&</sup>, Judith Breuer<sup>7&</sup>

7

8 <sup>1</sup>Division of Infection and Immunity, University College London, Cruciform Building, Gower St,  
9 London, WC1E 6BT

10 <sup>2</sup>Departments of Global Health and Epidemiology, University of Washington, Seattle WA, USA

11 <sup>3</sup>Institute of Immunology and Transplantation, Division of Infection and Immunity, University  
12 College London, Royal Free Campus

13 <sup>4</sup>Departments of Global Health, Epidemiology, Medicine (Div. Allergy and Infectious Diseases),  
14 University of Washington, Seattle WA, USA

15 <sup>5</sup>University of Nairobi, Department of Obstetrics and Gynaecology, Kenya, World Health  
16 Organization

17 <sup>6</sup>Research Centre of the Sainte-Justine University Hospital, Department of Microbiology,  
18 Infectious Diseases and Immunology, University of Montréal QC, Canada

19 <sup>7</sup> Department of Infection, Immunity and Inflammation, UCL Great Ormond Street Institute of  
20 Child Health, University College London, London, United Kingdom

21

22 \*Corresponding author. Email: j.breuer@ucl.ac.uk

23 ¶These authors contributed equally to this work.

24 &These authors also contributed equally to this work.

25

26 **Abstract**

27 Cytomegalovirus (CMV) is the commonest cause of congenital infection (cCMVi) and particularly  
28 so among infants born to HIV-infected women. Studies of cCMVi pathogenesis are complicated  
29 by the presence of multiple infecting maternal CMV strains, especially in HIV-positive women,  
30 and the large, recombinant CMV genome. Using newly developed tools to reconstruct CMV  
31 haplotypes, we demonstrate anatomic CMV compartmentalization in five HIV-infected mothers  
32 and identify the possibility of congenitally transmitted genotypes in three of their infants. A single  
33 CMV strain was transmitted in each congenitally infected case, and all were closely related to  
34 those that predominate in the cognate maternal cervix. Compared to non-transmitted strains,  
35 these congenitally transmitted CMV strains showed statistically significant similarities in 19 genes  
36 associated with tissue-tropism and immunomodulation. In all infants, incident superinfections  
37 with distinct strains from breast milk were captured during follow-up. The results represent  
38 potentially important new insights into the virologic determinants of early CMV infection.

39

## 40 **Introduction**

41 Human cytomegalovirus (CMV) is the commonest infectious cause of congenitally-acquired  
42 disability [1]. Between 0.2% and 2% of all live births have congenital CMV infection (cCMVi), and  
43 of these an estimated 15%-20% develop permanent sequelae ranging from sensorineural hearing  
44 loss to severe neurocognitive impairment [2, 3]. Maternal coinfection with HIV, even when  
45 mitigated by antiretroviral treatment, is associated with higher CMV viral loads in plasma, saliva,  
46 cervix and breast milk, and a greater risk of both congenital and postnatal CMV transmission [4-  
47 7]. Numerous studies have highlighted the negative health impacts of CMV on both HIV-infected  
48 and HIV-exposed uninfected (HEU) infants and children [8-10].

49  
50 Primary maternal CMV infection during pregnancy confers a 30%-40% risk of transmission to the  
51 fetus [11]. Pre-existing maternal CMV immunity appears to reduce the risk of cCMVi, though it is  
52 clearly imperfect [12]. Over two-thirds of infants with cCMVi are born to seropositive women,  
53 which constitute 88.4% of women in the Kenyan community from whom these study participants  
54 were drawn [13]. Moreover, the overall risk of cCMVi is directly proportional to the maternal  
55 seroprevalence in a population [14]. Increasing evidence points to the importance of maternal  
56 CMV reinfection with new antigenic strains during pregnancy as a major risk factor for non-  
57 primary cCMVi [12, 15]. Evidence that household children may be a source of maternal  
58 reinfection provides additional support for this hypothesis [16, 17].

59  
60 The CMV genome is the largest of the human herpesviruses. Regions of extensive sequence  
61 variability together with high levels of recombination between different strains results in high

62 diversity for a DNA virus [18-20]. Individuals are often infected with multiple CMV strains. We  
63 have recently demonstrated that separate CMV haplotypes can be resolved from high-  
64 throughput sequencing (HTS) data [21]. This advance, by enabling tracking of individual genomes  
65 within mixed CMV infections, has already revealed the impact of mutation, recombination and  
66 selection in shaping the course of infection [21]. Here we apply these methods to CMV genomes  
67 sequenced from samples from five HIV-infected Kenyan women and their infants that were  
68 collected between 1993 and 1998 originally for studies of maternal-infant HIV transmission [7].  
69 By reconstructing genome-wide haplotypes from these longitudinal samples, we are able to  
70 examine the diversity of CMV shed by HIV-infected women and the specific genotypes that are  
71 transmitted in congenital and postnatal infections, and to reconstruct the likely chronology with  
72 which specific CMV variants were transmitted from mothers to infants.

73

## 74 **Results**

### 75 **Participant characteristics, sampling, depth of sequencing**

76 Details of the study cohort, follow-up, sample collection, and HIV and CMV infection status and  
77 transmission have been previously described [22-24]. Sufficient residual sample was available  
78 from the five families analysed here. To maximise the chance of recovering near full genomes,  
79 we selected samples reported in the original publication [23] to have  $> 10^3$  copies/ml, as this is  
80 the limit at which we generally can generate whole genomes from blood. Of the five mother-  
81 infant pairs analysed, four infants were HIV-exposed uninfected (HEU) (Infants 22, 123, 41, 14),  
82 and one was HIV-infected (Infant 12).

83

### 84 **CMV viral loads and sequencing**

85 Cervical, breast milk, and infant blood CMV viral loads, Mother blood plasma HIV viral loads, and  
86 time of sample collection for the five mother-infant pairs studied are shown in Figure 1. The  
87 percentage of genome coverage and mean read depths are shown in Table 1. While breast milk  
88 samples had greater than 70% coverage at depths of 10x or more, the cervical and infant samples  
89 were of generally of lower depth, likely due to degradation of DNA due to the age and handling  
90 of the samples; genome coverage and mean de-duplicated read depth were directly related to  
91 actual CMV genome copy number present in the input material (Figure 1 – Figure Supplement 1).  
92 For all subsequent analysis, we removed samples with genome coverage of less than 20%.  
93 Fourteen of the remaining 20 cervical and baby samples had genome coverage above 70% and  
94 read depths of greater than 10x (Table 1).

95

## 96 **CMV genome sequence relatedness and diversity**

97 We used multidimensional scaling to cluster CMV genomic sequences by nucleotide similarity  
98 (Figure 2), as use of phylogenetic trees is problematic due to the high levels of CMV  
99 recombination. Sequences from families 12, 14 and 41 all clustered by family. Families 22 and  
100 123 clustered in two distinct spaces, suggesting infection with more than one strain. In all five  
101 cases, the first sample from each infant (indicated by an arrow) clustered most closely with that  
102 of its mother, indicating the likelihood of recent maternal-infant transmission.

103

104 To further investigate the possibility of mixed infections, we calculated the within-sample  
105 nucleotide diversity, a metric that we have shown previously can be used as a proxy for the  
106 likelihood of mixed strain infections [21]. It has previously been reported that a nucleotide diversity  
107 of 0.005 or above is likely to indicate a mixed infection [21]. Figure 2 – Figure Supplement 1  
108 shows that almost all the breast milk samples were highly diverse and therefore likely to contain  
109 multiple virus strains, a finding consistent with previous analyses of breast milk from HIV-infected  
110 women [25]. In contrast, the cervical and infant samples, with the exception of one cervical  
111 sample from family 12, showed lower diversity. We used subsampling to demonstrate that  
112 computed nucleotide diversities are robust down to sequencing depths of >10 (Figure 2 – Figure  
113 Supplement 2). Low diversity was also observed in cervical and blood spots with higher coverage  
114 and read depths (Table 1).

115

## 116 **Reconstruction of individual haplotypes reveals CMV compartmentalization**

117 To resolve the individual viral sequences (haplotypes) within each sample, we used our previously  
118 described method HaROLD [26]. Figure 3 shows that haplotypes for each sample tended to  
119 cluster by family group albeit with clear evidence of distinct clusters even within a family e.g.  
120 family 22.

121  
122 The presence of mixed infections within a single family was supported by data showing that a  
123 subset of the sequence haplotypes within each family had pairwise distances as great as those  
124 between unrelated GenBank sequences (Figure 3 – Figure Supplement 1). Within-family  
125 phylogenetic analysis (Figure 3 – Figure Supplement 2) shows distinct clusters of the  
126 phylogenetically related sequence haplotypes recovered from breast milk, cervix and baby, likely  
127 to represent variants forming distinct viral strains (Figure 3 – Figure Supplement 2). Based on the  
128 distribution of pairwise distances (see Methods, Figure 3 – Figure Supplement 3), we clustered  
129 similar haplotypes together into strains henceforth termed genotypes, so that all members of a  
130 cluster have a pairwise evolutionary distance with all other members less than 0.017, resulting in  
131 26 clusters which we refer to as genotypes. In no cases did haplotypes from different families  
132 fulfil our clustering criterion confirming that haplotypes were not shared between unrelated  
133 families.

134  
135 For ease of reference, genotypes were coloured differently, with the genotype predominating in  
136 the first cervical sample of each family coloured red (Figure 3 – Figure Supplement 2). Other  
137 genotypes were coloured by their phylogenetic and pairwise distances from this genotype (Figure

138 3 – Figure Supplement 2). From our data, we identified a total of 26 genotypes with between 3  
139 and 9 genotypes for each family (Figure 3 – Figure Supplement 2).

140

141 To elucidate the relationship between maternal and infant genotypes, we plotted the abundance  
142 of each within a sample over time (Figure 4). All five mothers were infected with multiple  
143 genotypes in breast milk. In many cases genotypes within a single maternal sample were as  
144 genetically distant as unrelated database sequences, suggesting the presence of multiple distinct  
145 CMV strains (Figure 3 – Figure Supplement 2, Figure 4). Relative genotype abundances present in  
146 breast milk changed over time. One unique genotype appeared in the breast milk of mother 22  
147 at 6 weeks, disappearing from a subsequent sample (Figure 4). This genotype was genetically  
148 distinct not only from other genotypes in family 22 but from genotypes in all other families,  
149 reducing the likelihood that it was a contaminant and may therefore have represented a new  
150 reinfection or reactivation of pre-existing latent infection. All cervical samples showed a single  
151 dominant genotype (Figure 4), including mother 12, whose sample was more diverse and found  
152 to contain low levels of other genotypes. Overall, the data point to compartmentalization of CMV  
153 populations between cervix and breast milk.

154

### 155 **Transmission bottlenecks**

156 CMV genomes from individual infant blood spots also showed lower diversity (Figure 2 – Figure  
157 Supplement 1), and predominance of one genotype (Figure 4), including in samples with good  
158 sequence read depth e.g. Baby12 DEL and 9M, Baby14 6W,14W and 6M, Baby22 14W, Baby123  
159 10W and 12M, (Table 1), indicating the likelihood of a bottleneck in mother-to-child

160 transmission. Two infants (families 12 and 123 Figure 1) who tested positive at birth were first  
161 infected with the genotype present in the greatest abundance in the cervix (Figure 4 and Figure  
162 3 – Figure Supplement 2). The same pattern was found in a third infant (family 22) whose first  
163 sample at two weeks of age tested positive (Figure 2, **Figure 3 – Figure Supplement 2, and Figure**  
164 **4**). Interestingly, all three of these congenitally infected infants were subsequently re-infected  
165 with distinct genotypes present in breast milk (Figure 4). Two infants with initially two (family 14)  
166 and three (family 41) negative tests from birth onwards, first became positive at 6 and 10 weeks  
167 respectively. The genotypes detected in the blood spots from both of these infants were present  
168 in breast milk and differed from the most abundant genotype in cervix (Figure 4).

169

#### 170 **Subsampling to control for the impact of read depths**

171 To determine the degree to which results were affected by the quality of sequence, we  
172 subsampled reads of different samples to show that sample diversity calculations are robust at  
173 read depths of  $\geq 5$  (Figure 2 – Figure Supplement 2); eight of the 18 blood spots and four of seven  
174 cervical samples had mean read depth  $\geq 10$  (Table 1) and all except one were of low diversity  
175 (Figure 2 – Figure Supplement 1). To determine the extent to which read depth affected  
176 haplotype frequencies, the 12-month breastmilk sample from mother 12, which had a mean read  
177 depth of 779.72 and five haplotypes (Figure 3 – Figure Supplement 2), was subsampled down to  
178 mean read depth  $< 4$  (Figure 4 – Figure Supplement 1). All of the haplotypes in this sample were  
179 present for read depths of 22 or more, with three haplotypes identified even at the lowest read  
180 depth. Nine out of ten cervical and blood spot samples from four families with read depths of  
181  $> 20$  (Table 1), had either single genotypes or multiple closely related variants (Figure 4)

182 supporting previous conclusions around compartmentalization and transmission bottlenecks  
183 [27].

184

### 185 **Genotype compartmentalization**

186 Given the observation of multiple haplotypes in each of the mother-baby pairs, we can ask whether  
187 certain genotypes are more likely than others to be found in different compartments, and whether  
188 there are common characteristics of the genotypes observed in similar compartments in different  
189 individuals. In order to address this question, we considered all possible subsets of between two  
190 and five genotypes where each genotype was derived from a different mother-baby pair. We then  
191 used fixation index (FST) to compare the genetic similarities of all of the genotypes in this set  
192 relative to the remaining genotypes. P-values and false discovery rates for each pair were  
193 calculated using non-parametric bootstrapping. In order to compare various subsets, we  
194 computed a confidence weighted sum of FST (cwsFST) values for each subset. The distribution of  
195 cwsFST values is shown in Figure 5 – Figure Supplement 1. As can be seen, there are a large  
196 number of subsets with significant cwsFST values, far in excess of what is observed for scrambled  
197 sequences (black line).

198

199 The sum weighted FST value for the subset of five genotypes that predominated in the cervical  
200 samples was not significantly different from other subsets, suggesting overall, that genotypes  
201 that predominated in the cervix of these women were less closely related than most other  
202 comparisons (**Figure 5 – Figure Supplement 1**, black arrow). Intriguingly, however, the subset of  
203 cervical genotypes from mother-baby pairs 12, 22, and 123 had a sum weighted FST with a value

204 greater than 99.6% of the other subsets (**Figure 5 – Figure Supplement 1**, blue arrow), indicating  
205 a strong signal of inter-patient viral convergence. These genotypes were from the three mother-  
206 baby pairs with proven congenital infection based on first detection of CMV in the baby at  $\leq 2$   
207 weeks of age, and in whom the baby's genotype was identical to that predominating in cervix. In  
208 contrast, the predominant cervical genotypes from mothers 14 and 41 showed low levels of  
209 relatedness (**Figure 5 – Figure Supplement 1**, red arrow). The infant strains from 14 and 41 were  
210 most closely related to those from their mothers' breast milk (**Figure 3 – Figure Supplement 2**  
211 **and Figure 4**).

212

213 The FST analysis identified 19 genes as likely to be contributing to the genetic similarity between  
214 congenitally transmitted genotypes from mothers 12, 22, 123 (FDR < 0.05) (Figure 5). The  
215 comparison between these congenitally-transmitted and other genotypes generally yielded the  
216 same genes when the pairwise difference was varied to cluster haplotypes into more or fewer  
217 genotypes (**Figure 5 – Figure Supplement 2**), suggesting that this finding is not an artefact of  
218 decisions about haplotype clustering.

219

220

## 221 **Discussion**

222 We used next generation sequencing and haplotype reconstruction of individual CMV genomes,  
223 obtained from samples of HIV-infected women and their infants, to identify mixed infections,  
224 compartmentalization and distinct viral-genotype associations with transmission of CMV from

225 mother-to-infant. Breast milk CMV showed high nucleotide diversity and, as has been previously  
226 reported [25], contained a mixture of viral genotypes, some of which were as genetically distant  
227 from each other as unrelated GenBank sequences and can therefore be considered distinct viral  
228 strains. Cervical samples were of low nucleotide diversity and dominated by a single viral  
229 genotype that was, with one exception, present in lower abundance in breast milk. Our data fit  
230 with most but not all [28] previous reports of CMV within-host compartmentalization based on  
231 genotyping of subgenomic fragments [29-32]. We found little evidence for widespread new  
232 superinfecting or reactivating viruses in these mothers. In line with findings from the  
233 immunosuppressed RhCMV monkey model of congenital infection, cCMVi [33] genotypes  
234 (strains) comprised families of closely related haplotypes. However, unlike the finding for  
235 congenitally transmitted gB and gL RhCMV variants, even where we found transmission of one  
236 genotype, maternal and infant haplotypes were not completely identical either in early,  
237 potentially congenital CMV infections, or in postnatally transmitted viruses from breastmilk.  
238 Neither were haplotypes sampled at different times from maternal breast milk conserved,  
239 suggesting a measure of de novo mutation in this patient group, in line with previous findings  
240 [20].

241  
242 Our method of reconstructing viral haplotypes in serial samples provides insights into the natural  
243 history of CMV infection. While all mothers had mixtures of genotypes in breastmilk, the  
244 proportions changed over time for some (family 22 and 41) and remained more stable in others.  
245 Whether expanding genotypes in mothers 22 and 41 had been recently acquired is not known  
246 but would be consistent with incident reinfection. In contrast, all infants were initially infected

247 with a single genotype (Figure 4), supporting a bottleneck to CMV transmission [21, 33,  
248 34]. Apparent reinfection by viruses present in breast milk occurred in all four infants with  
249 multiple samples (Figure 4). We posit that the appearance of a new strain in an infant sampled  
250 from birth can confidently be interpreted as a newly acquired exogenous virus rather than  
251 reactivation of a previously undetected one. In all cases, the reinfecting strains were genetically  
252 distant from and replaced the previously dominant strain (Figure 4). Taken together with the rise  
253 and fall of infant CMV viral loads over time (Figure 1), this pattern is consistent with immunity  
254 against the infants' first CMV strain not being protective against reinfection with antigenically  
255 distinct strains, a concept that can be further tested. Of note, reinfection with the closely related  
256 strains also appears to occur readily with both human CMV and in animal models [16, 35].  
257 Repeated reinfection with distinct strains may explain the high genetic variability observed  
258 between sequential samples in early sequencing studies of CMV genomes from congenitally-  
259 infected infants [19, 32].

260• Those infants who tested positive at <3 weeks from birth were congenitally infected by  
261 definition[15]. In contrast, we cannot formally rule out cCMVi in the two others who were  
262 classified as having post-natal infection, since sensitivity of PCR detection of CMV DNA in new-  
263 born blood spots is only approximately 84% [36], and new-born saliva or urine were not available.  
264 However, this is unlikely given that only a small minority of infants have cCMVi, even among  
265 those born to HIV-infected women. Furthermore, it is striking that genotypes in babies with  
266 proven cCMVi were highly similar to maternal cervical genotypes, while those with negative tests  
267 for the first six weeks of life were not, and the strains detected later in the blood of these two  
268 infants were most similar to those in their mothers' breast milk.

269 While it has previously been noted that a severe genetic bottleneck occurs during CMV  
270 transmission from mother to fetus or infant [20, 32, 37], it remains unknown whether CMV  
271 transmitted/founder virus populations share genotypic features that confer a fitness advantage  
272 for establishing an initial infection, such as seen in HIV [38]. Notwithstanding the apparent  
273 dominance of one genotype in each of the cervical samples, our analysis did not show evidence  
274 for inter-patient convergence of cervical genotypes per se. Rather the three cervical genotypes  
275 that were detected in babies 12, 22 and 123, who were infected at birth showed a higher level of  
276 genetic similarity than over 99.6% of other subset comparisons and much greater than would be  
277 expected by chance (black line) (Figure 5 – Figure Supplement 1). Nineteen genes (Figure 5, Table  
278 2) had particularly high ( $p < 0.01$ ) similarity scores. Twelve of the 19 genes with the highest  
279 similarity scores (Figure 5) are part of the highly diverse RL11 gene family. Uniquely, RL11 genes  
280 form an island of linkage within the otherwise highly recombinant CMV genome [18]. Phylogeny  
281 of primate CMV RL11 complexes recapitulates the evolutionary history of the cognate host,  
282 suggesting it to be a potential driver of CMV co-evolution and speciation [18]. It is intriguing that  
283 RL11 family proteins influence tissue tropism [34] or are immunomodulatory [34, 39-43].  
284 Together with its functional properties (Table 2) and extreme diversity [18], the possibility that  
285 within-species CMV RL11 gene-family variation may also influence within-host viral adaption to  
286 different compartments and/or transplacental transmission presents a tractable hypothesis that  
287 can now be tested. cCMVi is thought to occur primarily through maternal viremia followed by  
288 replication in placental cytotrophoblasts resulting in spread to the fetus [44]. The three mothers  
289 who transmitted their viruses congenitally had higher cervical viral loads than mothers whose  
290 babies become infected post-partum (Figure 1). Analysis of data from the whole cohort of

291 mothers confirmed that women who transmitted CMV *in utero* had mean cervical CMV vial loads  
292 at 38 weeks that were 0.83 log<sub>10</sub> copies/ml (SD=1.0, p=0.02) higher than women who did not  
293 transmit CMV *in utero* (data not shown) [23]. We therefore speculate that virus sampled in the  
294 cervix is representative of CMV populations that infect and cross the placenta, and that a possible  
295 explanation for our findings is that the properties that promote replication to higher titers in  
296 genital tissue may also predispose to transplacental infection.

297  
298 Other genes with high similarity (FST) scores include US27, which codes for a G-protein-coupled  
299 receptor (GPCR) homologue that modulates signalling of the CXCR4 chemokine and may have a  
300 role during viral entry and egress [45], and US26 whose function is unknown. Less marked but  
301 still significantly different from non-congenitally transmitted strains, UL40 protein [46]  
302 modulates NK cell function. NK cells are the most abundant lymphocytes in placental tissue [44],  
303 while UL50 is also immunomodulatory [47, 48]. Finally, UL74, coding for glycoprotein O, which is  
304 highly significantly similar in all bar one comparisons (Figure 5 – Figure Supplement 2), is part of  
305 the glycoprotein complex which is critical for tropism and entry into both fibroblasts and  
306 epithelial cells [49]. Of interest, gB and gL, which showed considerable diversity in the congenital  
307 RhCMV model were, as might be expected, not represented among the genes sharing significant  
308 genetic similarity in our analysis. One possibility that would unite our findings and those of the  
309 congenital RhCMV model is that CMV transmission bottlenecks are agnostic of variation in genes  
310 not implicated in transmission.

311

312 Being born to an HIV-infected women is a major risk factor for cCMVi as well as long term CMV-  
313 related complications, whether or not the child acquires HIV [8, 9]. We show here that,  
314 irrespective of the route of first infection, HIV-exposed uninfected (HEU) children frequently  
315 acquire repeated infections with different CMV viruses within the first year of life. Preliminary  
316 evidence suggests that breast milk of HIV-uninfected women may have lower CMV viral loads  
317 and carry fewer strains [50]. If this is true, the possibility that HEU, as well as HIV-infected, infants  
318 are exposed to greater numbers of CMV strains during infancy as compared with HIV-uninfected  
319 infants, may provide an explanation for their worse clinical outcomes, a hypothesis that can now  
320 be tested in prospective studies. Similarly, these methods promise to be invaluable for studying  
321 the role of maternal CMV reinfection during pregnancy, a question of central importance in the  
322 field [12].

323  
324 This study potentially provides several new insights into the pathogenesis of CMV infection.  
325 However, the study is limited by the small number of subjects, the fact that all women were HIV-  
326 1 infected and the lack of samples and data to absolutely confirm the route of CMV acquisition  
327 by these infants. Because we were only able to analyse maternal breast milk, cervical samples  
328 and infant blood, and only intermittently, it is possible that some transmitted viral variants were  
329 not captured. Some, particularly cervical and blood spot samples, had low CMV viral loads and,  
330 as a result, suboptimal genome coverage. Mapping data confirmed that in these cases sequence  
331 loss was random, excluding the possibility of systematic bias. To further address this potential  
332 bias, we subsampled samples with good coverage to identify read-depth thresholds above which  
333 the diversity estimation is robust and haplotype frequency to 5% and above is preserved (**Figure**

334 **2 – Figure Supplement 2 and Figure 4 – Figure Supplement 1).** Analysis of only those samples  
335 with read depths above the identified thresholds supported our overall conclusions. The quality  
336 of the sequence and the numbers of samples allowed for conclusions to be drawn at gene level  
337 only and precluded robust identification of putative motifs or single nucleotide polymorphisms  
338 associated with biological differences.

339

340 In summary, by reconstructing the individual CMV haplotypes we found evidence for mixed CMV  
341 infection in HIV-infected women, and compartmentalization of viral strains between cervical and  
342 breast milk. Infants appeared usually to acquire one virus genotype initially, indicating a  
343 transmission bottleneck, though subsequent reinfection with a second virus from maternal  
344 breast milk was common. We also found that viruses transmitted congenitally resembled the  
345 virus genotypes that were present at highest abundance in cervix, and shared genetic features  
346 that distinguished them from CMV strains predominating in breast milk and in the cervixes of  
347 women whose infants were apparently first infected post-partum. These data provide new  
348 testable insights into the pathogenesis of CMV transmission from mothers to their infants, as  
349 well as tools to unravel the importance of viral diversity for reinfection and congenital  
350 transmission, questions that are central to the development of a vaccine to prevent the global  
351 burden of disease due to CMV.

352

## 353 **Materials and Methods**

354 Samples were approved for research by the Institutional Review Board of the University of  
355 Washington and the Ethics and Research Committee of Kenyatta National Hospital IRB  
356 NCT00530777 and sequenced under the ULCP Biobank REC approval. Approval for use of  
357 anonymised residual diagnostic specimens were obtained through the University College  
358 London/University College London Hospitals (UCL/UCLH) Pathogen Biobank National Research  
359 Ethics Service Committee London Fulham (Research Ethics Committee reference: 12/LO/1089).  
360 Informed patient consent was not required.

361

### 362 **Patient specimens**

363 Mother-child pairs were selected from a randomized, placebo-controlled trial to determine the  
364 impact of twice-daily valacyclovir (500 mg) on breast milk HIV RNA viral load in HIV-1/HSV-2 co-  
365 infected women (NCT 00530777). Trial design, participant characteristics, and follow-up have  
366 been reported elsewhere, [22-24] and the University of Washington Institutional Review Board  
367 and Kenyatta National Hospital Research and Ethics Committee approved the research. Women  
368 received short course antiretrovirals for prevention of mother-to-child HIV transmission, but no  
369 women or infants received combination antiretroviral therapy, as the study was conducted  
370 before recommendations for universal treatment. All women were HIV-1, HSV-2 and CMV co-  
371 infected. For this CMV genomics study, we selected 5 mother-infant pairs from the placebo arm,  
372 who had well-defined timing of infant CMV infection. All infants were HIV-exposed, and one was  
373 HIV-infected. Women had cervical swabs and blood specimens collected at 34- and 38-weeks  
374 gestation. Maternal blood and infant dried blood spots were collected delivery, then postpartum

375 at 2, 6, 10, 14, 24, 36, and 52 weeks. Breast milk was collected at all times after delivery. Blood  
376 plasma, cervical swabs, and breast milk supernatant (whey) were cryopreserved at  $-80^{\circ}\text{C}$  for the  
377 study of HIV and other co-infections.

378

#### 379 **DNA extraction and CMV DNA measurement**

380 Viral nucleic acids were extracted from blood plasma, dried blood spots, breast milk supernatant  
381 and cervical swabs as previously described using the Qiagen UltraSens Viral Nucleic Acid  
382 extraction kit [23]. Quantitative real-time PCR was used to measure CMV DNA levels in these  
383 specimens [23].

384

#### 385 **Sure-select sequencing**

386 Hybridization and library preparation were performed as previously described [51]. Briefly,  
387 extracted DNA was sheared by acoustic sonication (Covaris e220, Covaris Inc.). DNA fragments  
388 underwent end-repair, A'-tailing, and (Illumina) adaptor ligation. DNA libraries were hybridised  
389 with biotinylated 120-mer custom RNA baits designed using all available CMV full genomes in  
390 Genbank for 16-24 hrs at  $65^{\circ}\text{C}$  and subsequently bound to MyOne™ Streptavidin T1 Dynabeads™  
391 (ThermoFisher Scientific). Following washing, libraries were amplified (18 cycles) to generate  
392 sufficient input material for Illumina sequencing. Paired-end sequencing was performed on an  
393 Illumina MiSeq using the 500 cycle v2 Reagent Kit (Illumina, MS-102-2003). Samples were  
394 sequenced in four different batches by family group.

395

396 Reads generated were quality checked and mapped to the Merlin Reference sequence followed  
397 by removal of duplicates using the CLC Genomics Workbench ver. 10.1. Consensus sequence was  
398 extracted with a minimum coverage of 2X. All consensus sequences along with other Genbank  
399 reference sequences were aligned using MAFFT 7.212 [52] and refined by manual editing.

400

#### 401 **Clustering**

402 Pairwise distances between sequences were calculated using the dist.dna function from R  
403 package Ape v.5.3 [53]. Sequences were clustered using multidimensional scaling as  
404 implemented by the cmdscale function from R package Stats v.3.6 [54].

405

#### 406 **Nucleotide diversity**

407 Nucleotide diversity was calculated by fitting the observed variant frequency spectrum to the  
408 mixture of two distributions, one representing sequencing errors (represented by a Beta  
409 distribution), the other representing true diversity (represented by a four-dimensional Dirichlet  
410 distribution plus delta function, the latter representing invariant sites). The parameters for these  
411 two distributions were optimized by maximizing the log likelihood. This framework allows all of  
412 the sequencing data to be used and does not require pre-filtering the data to remove sites with  
413 low read depth or few variants resulting in the favourable robustness to read depth, as shown in  
414 Figure 2 – Figure Supplement 2. Software is available for download at GitHub Repository,  
415 <https://github.com/ucl-pathgenomics/NucleotideDiversity>.

416

#### 417 **Haplotype reconstruction**

418 Haplotype reconstruction was accomplished using HaROLD with default settings [26]. Details of  
419 this procedure are described in the associated publications. In brief, HaROLD employs a two-step  
420 process. The first step is based on the assumption that there are a limited number of haplotypes  
421 that are the same for all of the samples from a given mother/ child data set, so that the  
422 differences in the frequencies of polymorphisms represent different mixtures of these  
423 haplotypes. By taking advantage of the co-variation of variant frequencies, HaROLD creates a set  
424 of haplotypes for each of the data sets, optimized so that linear combinations of these haplotypes  
425 can best account for the observed variant frequencies. The number of haplotypes is chosen to  
426 maximize the log likelihood of the observed frequencies. The second step involves relaxing the  
427 assumption of constant haplotypes, with each sample treated individually. For each sample,  
428 reads are assigned probabilistically to the various haplotypes generated by the first step. These  
429 haplotype sequences and frequencies are then adjusted based on the assigned reads. The reads  
430 are then re-assigned to these adjusted haplotypes, and the procedure is repeated until  
431 convergence. During this process, haplotypes can be merged if that decreases the Akaike  
432 Information Criterion (AIC) [55]. This procedure results in a set of haplotypes for each sample,  
433 loosely based on the haplotypes derived from the first step.

434

#### 435 **Haplotype trees**

436 Maximum Likelihood trees of the haplotypes from each family were computed using RaxML  
437 v8.2.10, implementing the GTR model, with 1000 bootstrap replicates [56].

438

#### 439 **Haplotype clustering**

440 The haplotypes for each mother/baby data set were divided into genotypes. We calculated the  
441 pairwise evolutionary distance (the sum of distances on the evolutionary tree between the  
442 haplotypes and their latest common ancestor) for all pairs of haplotypes in each family. As shown  
443 in Figure 3 – Figure Supplement 3, the observed distribution of such pairwise distances fits the  
444 sum of a Gamma distribution (69.3%,  $\alpha = 19.5$ ,  $\beta = 0.0015$ ) and an exponential distribution  
445 (30.7%,  $\text{mean} = 0.01$ ), indicative of two classes of relationships – pairs of sequences that are  
446 highly similar, modelled by the exponential, representing small accumulated variations, and pairs  
447 that are more distinct, represented by the Gamma distribution. We chose the crossing point of  
448 these two distributions, at a cut-off distance of 0.017, as differentiating small variations from  
449 larger differences (Figure 3 – Figure Supplement 3). We then grouped the haplotypes into clusters  
450 so that all members of a cluster have a pairwise evolutionary distance with all other members  
451 less than 0.017, resulting in 26 clusters which we refer to as genotypes. We used these groups to  
452 assign colours to the different haplotype-clusters (genotypes) in Figure 4 and Figure 3 – Figure  
453 Supplement 2.

454

455 We used  $F_{ST}$  to identify sequence characteristics associated with sets of genotypes. Consensus  
456 sequences were constructed for each genotype.  $F_{ST}$  values, representing the genetic difference  
457 between a subset of genotypes and the other genotypes, were calculated for each gene. P-values  
458 and corresponding false discovery rates were estimated by non-parametric bootstrapping,  
459 through scrambling the bases at each position amongst the clusters. The results are shown for  
460 the 26 genotypes obtained with a cut-off distance of 0.017; changing this cut-off resulted in

461 increased or decreased numbers of genotypes, but yielded similar results, especially for the more  
462 confident identifications (Figure 5 – Figure Supplement 2).

463

#### 464 **Evaluating the similarity between subsets of genotypes**

465 We use  $F_{ST}$  values to identify similarities between individual genes from subsets of genotypes  
466 compared with the other genotypes. In order to compare the magnitude of the similarities of  
467 different subsets, we would like to take the sum of the  $F_{ST}$  values for all genes where the  
468 similarities are real and not the result of random associations. As we cannot definitively identify  
469 these genes, we instead consider the sum of the  $F_{ST}$  values for all genes weighted by our  
470 confidence that the  $F_{ST}$  value is significant, represented as one minus the false discovery rate.

471

#### 472 **Acknowledgments**

473 We acknowledge the support of the MRC/NIHR UCLH/UCL Biomedical Research Centre funded  
474 Pathogen Genomics Unit. This work was funded by EUFP7 grant 304875 (PI Breuer), Wellcome  
475 Trust grant 204870 (PI Griffiths), NIH National Institute of Allergy and Infectious Diseases grant  
476 AI087369 (PI Slyker), AI027757 (PI Slyker, Holmes), AI076105 and K24 AI087399 (Farquhar),  
477 National Institute of Child Health and Human Development HD057773–01, HD054314 (Farquhar).  
478 JP is funded by a Rosetrees Trust PhD Studentship M876. SM and J Bryant are funded by Henry  
479 Wellcome fellowships. J Breuer receives funding from the UCL/UCLH NIHR Biomedical Research  
480 Centre.

481

#### 482 **Data availability**

483 Sequence reads have been deposited in NCBI Sequence Read Archive under BioProject ID  
484 PRJNA605798.

485

486 All software used are available for download at GitHub Repository, [https://github.com/ucl-](https://github.com/ucl-pathgenomics/NucleotideDiversity)  
487 [pathgenomics/NucleotideDiversity](https://github.com/ucl-pathgenomics/HAROLD) and <https://github.com/ucl-pathgenomics/HAROLD>.

488

489 **References**

- 490 1. Morton CC, Nance WE. Newborn hearing screening--a silent revolution. *N Engl J Med.* 2006;354(20):2151-  
491 64. Epub 2006/05/19. doi: 10.1056/NEJMra050700. PubMed PMID: 16707752.
- 492 2. Boppana SB, Ross SA, Fowler KB. Congenital cytomegalovirus infection: clinical outcome. *Clin Infect Dis.*  
493 2013;57 Suppl 4:S178-81. Epub 2013/12/07. doi: 10.1093/cid/cit629. PubMed PMID: 24257422.
- 494 3. Dollard SC, Grosse SD, Ross DS. New estimates of the prevalence of neurological and sensory sequelae and  
495 mortality associated with congenital cytomegalovirus infection. *Rev Med Virol.* 2007;17(5):355-63. Epub 2007/06/02.  
496 doi: 10.1002/rmv.544. PubMed PMID: 17542052.
- 497 4. Gantt S, Orem J, Krantz EM, Morrow RA, Selke S, Huang ML, et al. Prospective Characterization of the  
498 Risk Factors for Transmission and Symptoms of Primary Human Herpesvirus Infections Among Ugandan Infants. *J*  
499 *Infect Dis.* 2016;214(1):36-44. doi: 10.1093/infdis/jiw076. PubMed PMID: 26917575; PubMed Central PMCID:  
500 PMC4907408.
- 501 5. Gantt S, Leister E, Jacobsen DL, Boucoiran I, Huang ML, Jerome KR, et al. Risk of congenital  
502 cytomegalovirus infection among HIV-exposed uninfected infants is not decreased by maternal nelfinavir use during  
503 pregnancy. *J Med Virol.* 2016;88(6):1051-8. doi: 10.1002/jmv.24420. PubMed PMID: 26519647; PubMed Central  
504 PMCID: PMC4818099.
- 505 6. Slyker JA, Richardson B, Chung MH, Atkinson C, Asbjornsdottir KH, Lehman DA, et al. Maternal Highly  
506 Active Antiretroviral Therapy Reduces Vertical Cytomegalovirus Transmission But Does Not Reduce Breast Milk  
507 Cytomegalovirus Levels. *AIDS Res Hum Retroviruses.* 2017;33(4):332-8. Epub 2016/11/01. doi:  
508 10.1089/AID.2016.0121. PubMed PMID: 27796131; PubMed Central PMCID: PMC4818099.
- 509 7. Richardson BA, John-Stewart G, Atkinson C, Nduati R, Asbjornsdottir K, Boeckh M, et al. Vertical  
510 Cytomegalovirus Transmission From HIV-Infected Women Randomized to Formula-Feed or Breastfeed Their  
511 Infants. *J Infect Dis.* 2016;213(6):992-8. doi: 10.1093/infdis/jiv515. PubMed PMID: 26518046; PubMed Central  
512 PMCID: PMC4760415.
- 513 8. Garcia-Knight MA, Nduati E, Hassan AS, Nkumama I, Etyang TJ, Hajj NJ, et al. Cytomegalovirus viraemia  
514 is associated with poor growth and T-cell activation with an increased burden in HIV-exposed uninfected infants.  
515 *AIDS.* 2017;31(13):1809-18. Epub 2017/06/14. doi: 10.1097/QAD.0000000000001568. PubMed PMID: 28609400;  
516 PubMed Central PMCID: PMC5538302.

- 517 9. Gompels UA, Larke N, Sanz-Ramos M, Bates M, Musonda K, Manno D, et al. Human cytomegalovirus  
518 infant infection adversely affects growth and development in maternally HIV-exposed and unexposed infants in  
519 Zambia. *Clin Infect Dis*. 2012;54(3):434-42. doi: 10.1093/cid/cir837. PubMed PMID: 22247303; PubMed Central  
520 PMCID: PMC3258277.
- 521 10. Hsiao NY, Zampoli M, Morrow B, Zar HJ, Hardie D. Cytomegalovirus viraemia in HIV exposed and infected  
522 infants: prevalence and clinical utility for diagnosing CMV pneumonia. *J Clin Virol*. 2013;58(1):74-8. Epub  
523 2013/06/04. doi: 10.1016/j.jcv.2013.05.002. PubMed PMID: 23727304.
- 524 11. Kenneson A, Cannon MJ. Review and meta-analysis of the epidemiology of congenital cytomegalovirus  
525 (CMV) infection. *Rev Med Virol*. 2007;17(4):253-76. Epub 2007/06/21. doi: 10.1002/rmv.535. PubMed PMID:  
526 17579921.
- 527 12. Britt WJ. Congenital Human Cytomegalovirus Infection and the Enigma of Maternal Immunity. *J Virol*.  
528 2017;91(15). doi: 10.1128/JVI.02392-16. PubMed PMID: 28490582; PubMed Central PMCID: PMC5512250.
- 529 13. Maingi Z, Nyamache AK. Seroprevalence of Cytomegalo Virus (CMV) among pregnant women in Thika,  
530 Kenya. *BMC Res Notes*. 2014;7:794. Epub 2014/11/14. doi: 10.1186/1756-0500-7-794. PubMed PMID: 25392013;  
531 PubMed Central PMCID: PMC4247150.
- 532 14. de Vries JJ, van Zwet EW, Dekker FW, Kroes AC, Verkerk PH, Vossen AC. The apparent paradox of  
533 maternal seropositivity as a risk factor for congenital cytomegalovirus infection: a population-based prediction model.  
534 *Rev Med Virol*. 2013;23(4):241-9. doi: 10.1002/rmv.1744. PubMed PMID: 23559569.
- 535 15. Boppana SB, Fowler KB, Britt WJ, Stagno S, Pass RF. Symptomatic congenital cytomegalovirus infection  
536 in infants born to mothers with preexisting immunity to cytomegalovirus. *Pediatrics*. 1999;104(1 Pt 1):55-60. Epub  
537 1999/07/02. doi: 10.1542/peds.104.1.55. PubMed PMID: 10390260.
- 538 16. Boucoiran I, Mayer BT, Krantz EM, Marchant A, Pati S, Boppana S, et al. Nonprimary Maternal  
539 Cytomegalovirus Infection After Viral Shedding in Infants. *Pediatr Infect Dis J*. 2018;37(7):627-31. doi:  
540 10.1097/INF.0000000000001877. PubMed PMID: 29889809.
- 541 17. Barbosa NG, Yamamoto AY, Duarte G, Aragon DC, Fowler KB, Boppana S, et al. Cytomegalovirus  
542 Shedding in Seropositive Pregnant Women From a High-Seroprevalence Population: The Brazilian Cytomegalovirus  
543 Hearing and Maternal Secondary Infection Study. *Clin Infect Dis*. 2018;67(5):743-50. Epub 2018/03/01. doi:  
544 10.1093/cid/ciy166. PubMed PMID: 29490030; PubMed Central PMCID: PMC6094000.

- 545 18. Lassalle F, Depledge DP, Reeves MB, Brown AC, Christiansen MT, Tutill HJ, et al. Islands of linkage in an  
546 ocean of pervasive recombination reveals two-speed evolution of human cytomegalovirus genomes. *Virus Evol.*  
547 2016;2(1):vew017. Epub 2016/06/15. doi: 10.1093/ve/vew017. PubMed PMID: 30288299; PubMed Central PMCID:  
548 PMC6167919.
- 549 19. Pokalyuk C, Renzette N, Irwin KK, Pfeifer SP, Gibson L, Britt WJ, et al. Characterizing human  
550 cytomegalovirus reinfection in congenitally infected infants: an evolutionary perspective. *Mol Ecol.* 2017;26(7):1980-  
551 90. Epub 2016/12/19. doi: 10.1111/mec.13953. PubMed PMID: 27988973.
- 552 20. Sackman AM, Pfeifer SP, Kowalik TF, Jensen JD. On the Demographic and Selective Forces Shaping  
553 Patterns of Human Cytomegalovirus Variation within Hosts. *Pathogens.* 2018;7(1). Epub 2018/02/01. doi:  
554 10.3390/pathogens7010016. PubMed PMID: 29382090; PubMed Central PMCID: PMC6167919.
- 555 21. Cudini J, Roy S, Houldcroft CJ, Bryant JM, Depledge DP, Tutill H, et al. Human cytomegalovirus haplotype  
556 reconstruction reveals high diversity due to superinfection and evidence of within-host recombination. *Proc Natl Acad*  
557 *Sci U S A.* 2019;116(12):5693-8. Epub 2019/03/02. doi: 10.1073/pnas.1818130116. PubMed PMID: 30819890;  
558 PubMed Central PMCID: PMC6431178.
- 559 22. Drake AL, Roxby AC, Ongecha-Owuor F, Kiarie J, John-Stewart G, Wald A, et al. Valacyclovir suppressive  
560 therapy reduces plasma and breast milk HIV-1 RNA levels during pregnancy and postpartum: a randomized trial. *J*  
561 *Infect Dis.* 2012;205(3):366-75. Epub 2011/12/08. doi: 10.1093/infdis/jir766. PubMed PMID: 22147786; PubMed  
562 Central PMCID: PMC3256951.
- 563 23. Roxby AC, Atkinson C, Asbjornsdottir K, Farquhar C, Kiarie JN, Drake AL, et al. Maternal valacyclovir and  
564 infant cytomegalovirus acquisition: a randomized controlled trial among HIV-infected women. *PLoS One.*  
565 2014;9(2):e87855. doi: 10.1371/journal.pone.0087855. PubMed PMID: 24504006; PubMed Central PMCID:  
566 PMC3913686.
- 567 24. Slyker J, Farquhar C, Atkinson C, Asbjornsdottir K, Roxby A, Drake A, et al. Compartmentalized  
568 cytomegalovirus replication and transmission in the setting of maternal HIV-1 infection. *Clin Infect Dis.*  
569 2014;58(4):564-72. doi: 10.1093/cid/cit727. PubMed PMID: 24192386; PubMed Central PMCID:  
570 PMC3905754.
- 571 25. Suarez NM, Musonda KG, Escriva E, Njenga M, Agbueze A, Camiolo S, et al. Multiple-Strain Infections of  
572 Human Cytomegalovirus With High Genomic Diversity Are Common in Breast Milk From Human

- 573 Immunodeficiency Virus-Infected Women in Zambia. *J Infect Dis.* 2019;220(5):792-801. Epub 2019/05/06. doi:  
574 10.1093/infdis/jiz209. PubMed PMID: 31050737; PubMed Central PMCID: PMC6667993.
- 575 26. Goldstein RA, Tamuri AU, Roy S, Breuer J. Haplotype assignment of virus NGS data using co-variation of  
576 variant frequencies. *bioRxiv.* 2018. doi: <https://doi.org/10.1101/444877>.
- 577 27. Renzette N, Bhattacharjee B, Jensen JD, Gibson L, Kowalik TF. Extensive genome-wide variability of human  
578 cytomegalovirus in congenitally infected infants. *PLoS Pathog.* 2011;7(5):e1001344. Epub 2011/06/01. doi:  
579 10.1371/journal.ppat.1001344. PubMed PMID: 21625576; PubMed Central PMCID: PMC6667993.
- 580 28. Puchhammer-Stockl E, Gorzer I, Zoufaly A, Jaksch P, Bauer CC, Klepetko W, et al. Emergence of multiple  
581 cytomegalovirus strains in blood and lung of lung transplant recipients. *Transplantation.* 2006;81(2):187-94. Epub  
582 2006/01/27. doi: 10.1097/01.tp.0000194858.50812.cb. PubMed PMID: 16436961.
- 583 29. Hage E, Wilkie GS, Linnenweber-Held S, Dhingra A, Suarez NM, Schmidt JJ, et al. Characterization of  
584 Human Cytomegalovirus Genome Diversity in Immunocompromised Hosts by Whole-Genome Sequencing Directly  
585 From Clinical Specimens. *J Infect Dis.* 2017;215(11):1673-83. Epub 2017/04/04. doi: 10.1093/infdis/jix157. PubMed  
586 PMID: 28368496.
- 587 30. Kadambari S, Atkinson C, Luck S, Macartney M, Conibear T, Harrison I, et al. Characterising variation in  
588 five genetic loci of cytomegalovirus during treatment for congenital infection. *J Med Virol.* 2017;89(3):502-7. Epub  
589 2016/08/04. doi: 10.1002/jmv.24654. PubMed PMID: 27486960.
- 590 31. Ross SA, Novak Z, Pati S, Patro RK, Blumenthal J, Danthuluri VR, et al. Mixed infection and strain diversity  
591 in congenital cytomegalovirus infection. *J Infect Dis.* 2011;204(7):1003-7. Epub 2011/09/02. doi:  
592 10.1093/infdis/jir457. PubMed PMID: 21881114; PubMed Central PMCID: PMC3164425.
- 593 32. Renzette N, Gibson L, Bhattacharjee B, Fisher D, Schleiss MR, Jensen JD, et al. Rapid intrahost evolution  
594 of human cytomegalovirus is shaped by demography and positive selection. *PLoS Genet.* 2013;9(9):e1003735. doi:  
595 10.1371/journal.pgen.1003735. PubMed PMID: 24086142; PubMed Central PMCID: PMC3784496.
- 596 33. Vera Cruz D, Nelson CS, Tran D, Barry PA, Kaur A, Koelle K, et al. Intrahost cytomegalovirus population  
597 genetics following antibody pretreatment in a monkey model of congenital transmission. *PLoS Pathog.*  
598 2020;16(2):e1007968. Epub 2020/02/15. doi: 10.1371/journal.ppat.1007968. PubMed PMID: 32059027.
- 599 34. Stanton RJ, Baluchova K, Dargan DJ, Cunningham C, Sheehy O, Seirafian S, et al. Reconstruction of the  
600 complete human cytomegalovirus genome in a BAC reveals RL13 to be a potent inhibitor of replication. *J Clin Invest.*

- 601 2010;120(9):3191-208. Epub 2010/08/04. doi: 10.1172/JCI42955. PubMed PMID: 20679731; PubMed Central  
602 PMCID: PMCPMC2929729.
- 603 35. Hansen SG, Powers CJ, Richards R, Ventura AB, Ford JC, Siess D, et al. Evasion of CD8+ T cells is critical  
604 for superinfection by cytomegalovirus. *Science*. 2010;328(5974):102-6. Epub 2010/04/03. doi:  
605 10.1126/science.1185350. PubMed PMID: 20360110; PubMed Central PMCID: PMCPMC2883175.
- 606 36. Wang L, Xu X, Zhang H, Qian J, Zhu J. Dried blood spots PCR assays to screen congenital cytomegalovirus  
607 infection: a meta-analysis. *Viol J*. 2015;12:60. doi: 10.1186/s12985-015-0281-9. PubMed PMID: 25889596; PubMed  
608 Central PMCID: PMCPMC4408583.
- 609 37. Mayer BT, Krantz EM, Swan D, Ferrenberg J, Simmons K, Selke S, et al. Transient Oral Human  
610 Cytomegalovirus Infections Indicate Inefficient Viral Spread from Very Few Initially Infected Cells. *J Virol*.  
611 2017;91(12). doi: 10.1128/JVI.00380-17. PubMed PMID: 28381570; PubMed Central PMCID: PMCPMC5446638.
- 612 38. Joseph SB, Swanstrom R, Kashuba AD, Cohen MS. Bottlenecks in HIV-1 transmission: insights from the  
613 study of founder viruses. *Nat Rev Microbiol*. 2015;13(7):414-25. doi: 10.1038/nrmicro3471. PubMed PMID:  
614 26052661; PubMed Central PMCID: PMCPMC4793885.
- 615 39. Cortese M, Calo S, D'Aurizio R, Lilja A, Pacchiani N, Merola M. Recombinant human cytomegalovirus  
616 (HCMV) RL13 binds human immunoglobulin G Fc. *PLoS One*. 2012;7(11):e50166. Epub 2012/12/12. doi:  
617 10.1371/journal.pone.0050166. PubMed PMID: 23226246; PubMed Central PMCID: PMCPMC3511460.
- 618 40. Van Damme E, Van Loock M. Functional annotation of human cytomegalovirus gene products: an update.  
619 *Front Microbiol*. 2014;5:218. Epub 2014/06/07. doi: 10.3389/fmicb.2014.00218. PubMed PMID: 24904534; PubMed  
620 Central PMCID: PMCPMC4032930.
- 621 41. Perez-Carmona N, Martinez-Vicente P, Farre D, Gabaev I, Messerle M, Engel P, et al. A Prominent Role of  
622 the Human Cytomegalovirus UL8 Glycoprotein in Restraining Proinflammatory Cytokine Production by Myeloid  
623 Cells at Late Times during Infection. *J Virol*. 2018;92(9). Epub 2018/02/23. doi: 10.1128/JVI.02229-17. PubMed  
624 PMID: 29467314; PubMed Central PMCID: PMCPMC5899185.
- 625 42. Bruno L, Cortese M, Monda G, Gentile M, Calo S, Schiavetti F, et al. Human cytomegalovirus pUL10  
626 interacts with leukocytes and impairs TCR-mediated T-cell activation. *Immunol Cell Biol*. 2016;94(9):849-60. Epub  
627 2016/10/19. doi: 10.1038/icb.2016.49. PubMed PMID: 27192938.

- 628 43. Gabaev I, Elbasani E, Ameres S, Steinbruck L, Stanton R, Doring M, et al. Expression of the human  
629 cytomegalovirus UL11 glycoprotein in viral infection and evaluation of its effect on virus-specific CD8 T cells. *J*  
630 *Virol.* 2014;88(24):14326-39. Epub 2014/10/03. doi: 10.1128/JVI.01691-14. PubMed PMID: 25275132; PubMed  
631 Central PMCID: PMC4249143.
- 632 44. Pereira L, Tabata T, Petitt M, Fang-Hoover J. Congenital cytomegalovirus infection undermines early  
633 development and functions of the human placenta. *Placenta.* 2017;59 Suppl 1:S8-S16. Epub 2017/05/10. doi:  
634 10.1016/j.placenta.2017.04.020. PubMed PMID: 28477968.
- 635 45. Frank T, Niemann I, Reichel A, Stamminger T. Emerging roles of cytomegalovirus-encoded G protein-  
636 coupled receptors during lytic and latent infection. *Med Microbiol Immunol.* 2019;208(3-4):447-56. Epub 2019/03/23.  
637 doi: 10.1007/s00430-019-00595-9. PubMed PMID: 30900091.
- 638 46. Heatley SL, Pietra G, Lin J, Widjaja JM, Harpur CM, Lester S, et al. Polymorphism in human  
639 cytomegalovirus UL40 impacts on recognition of human leukocyte antigen-E (HLA-E) by natural killer cells. *J Biol*  
640 *Chem.* 2013;288(12):8679-90. Epub 2013/01/22. doi: 10.1074/jbc.M112.409672. PubMed PMID: 23335510;  
641 PubMed Central PMCID: PMC3605686.
- 642 47. Lee MK, Kim YJ, Kim YE, Han TH, Milbradt J, Marschall M, et al. Transmembrane Protein pUL50 of  
643 Human Cytomegalovirus Inhibits ISGylation by Downregulating UBE1L. *J Virol.* 2018;92(15). Epub 2018/05/11.  
644 doi: 10.1128/JVI.00462-18. PubMed PMID: 29743376; PubMed Central PMCID: PMC6052311.
- 645 48. DeRussy BM, Boland MT, Tandon R. Human Cytomegalovirus pUL93 Links Nucleocapsid Maturation and  
646 Nuclear Egress. *J Virol.* 2016;90(16):7109-17. Epub 2016/05/27. doi: 10.1128/JVI.00728-16. PubMed PMID:  
647 27226374; PubMed Central PMCID: PMC4984640.
- 648 49. Wu Y, Prager A, Boos S, Resch M, Brizic I, Mach M, et al. Human cytomegalovirus glycoprotein complex  
649 gH/gL/gO uses PDGFR-alpha as a key for entry. *PLoS Pathog.* 2017;13(4):e1006281. Epub 2017/04/14. doi:  
650 10.1371/journal.ppat.1006281. PubMed PMID: 28403202; PubMed Central PMCID: PMC5389851.
- 651 50. Arcangeletti MC, Vasile Simone R, Rodighiero I, De Conto F, Medici MC, Martorana D, et al. Combined  
652 genetic variants of human cytomegalovirus envelope glycoproteins as congenital infection markers. *Virol J.*  
653 2015;12:202. Epub 2015/11/28. doi: 10.1186/s12985-015-0428-8. PubMed PMID: 26611326; PubMed Central  
654 PMCID: PMC4662005.

- 655 51. Houldcroft CJ, Bryant JM, Depledge DP, Margetts BK, Simmonds J, Nicolaou S, et al. Detection of Low  
656 Frequency Multi-Drug Resistance and Novel Putative Maribavir Resistance in Immunocompromised Pediatric  
657 Patients with Cytomegalovirus. *Front Microbiol.* 2016;7:1317. Epub 2016/09/27. doi: 10.3389/fmicb.2016.01317.  
658 PubMed PMID: 27667983; PubMed Central PMCID: PMC5016526.
- 659 52. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in  
660 performance and usability. *Mol Biol Evol.* 2013;30(4):772-80. Epub 2013/01/19. doi: 10.1093/molbev/mst010.  
661 PubMed PMID: 23329690; PubMed Central PMCID: PMC3603318.
- 662 53. Paradis E, Schliep K. ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R.  
663 *Bioinformatics.* 2019;35(3):526-8. Epub 2018/07/18. doi: 10.1093/bioinformatics/bty633. PubMed PMID: 30016406.
- 664 54. Team RC. A language and environment for statistical computing. R Foundation for Statistical Computing.  
665 Vienna, Austria2012.
- 666 55. Akaike H. Information theory and an extension of the maximum likelihood principle. 2nd International  
667 Symposium on Information Theory (BN Petrov and F Cs ä ki, eds); Akademiai Ki à do, Budapest1973.
- 668 56. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.  
669 *Bioinformatics.* 2014;30(9):1312-3. Epub 2014/01/24. doi: 10.1093/bioinformatics/btu033. PubMed PMID:  
670 24451623; PubMed Central PMCID: PMC3998144.

671

672

673

674

## 675 **Figures Legends**

676 **Figure 1.** CMV viral loads of longitudinal samples for each family from breast milk (red), baby  
677 blood spots (green) and cervix (blue), and HIV viral loads from mother's blood plasma. Vertical  
678 line indicates date of delivery. Horizontal line indicates minimum threshold of detection. Red  
679 circles indicate the samples that were submitted for whole genome sequencing.

680

681 **Figure 1 – Figure Supplement 1.** Scatter plots showing relationship between input viral load and  
682 (A) mean read depth and (B) genome coverage respectively.

683

684 **Figure 2.** Multidimensional scaling showing clustering of consensus genome sequences for each  
685 sample by family. Arrows indicate that the first baby blood spot clusters with their own maternal  
686 sequences in all cases.

687

688 **Figure 2 – Figure Supplement 1.** Within sample nucleotide diversity shown by family (colour)  
689 and sample type (symbol). BM; breast milk, CV; cervix, BS; baby blood spot. The figure shows  
690 that most cervical and blood spot samples are of low diversity, while most breast milk samples are  
691 of high diversity. Diversity of breast milk versus cervix;  $p=1.619e-07$  and versus baby blood spot;  
692  $p=9.69e-6$  (Mann-Whitney test).

693

694 **Figure 2 – Figure Supplement 2.** Effect of down-sampling on estimated diversity. Samples tested  
695 include family 14: 14W BS (green squares), family 41: 14W BM (blue dots), family 14: 6W BM  
696 (green triangles), family 12: 12M BM (maroon diamonds) all of which had initial read depths of

697 150 or more. The estimated diversity is relatively insensitive to read depth; in particular, down-  
698 sampling of high read-depth samples shows no tendency of the analysis to underestimate the  
699 diversity of low read-depth samples. This indicates that the low diversity observed in many of the  
700 CV and BS samples is not an artefact but is rather consistent with the presence of significant  
701 bottlenecks.

702

703 **Figure 3.** Multidimensional scaling showing clustering of haplotype sequences by family. Colours  
704 indicate the families; shapes indicate the types of sample.

705

706 **Figure 3 – Figure Supplement 1.** Pairwise differences between haplotypes within a family.  
707 Distances are compared with random GenBank sequences and sequences previously analyzed by  
708 the same pipeline and reported [21]. Higher values are similar to those seen between unrelated  
709 database sequences and indicate the presence of distinct strains.

710

711 **Figure 3 – Figure Supplement 2.** Maximum-likelihood phylogenetic tree to show haplotypes  
712 clusters (genotypes). By convention, the genotype most prevalent in cervix was coloured red for  
713 each family. Genotypes were designated where a distinct cluster of related haplotypes (pairwise  
714 distance  $\leq 0.017$ ) occurred with a bootstrap value of 100 (see methods and supplementary figure  
715 9). The genotype containing the most abundant haplotype present in the cervix is coloured red for  
716 each family. Thereafter sequences that are genetically closest to the red genotype are coloured  
717 magenta. Genotypes that are as distant from the cervical genotype as unrelated GenBank sequences  
718 are coloured shades of green, blue and purple. The number of clusters between 18 and 34 did not

719 affect subsequent conclusions about genetic similarity between cervical versus other strains (see  
720 Figure 5 – Figure Supplement 2).

721  
722 **Figure 3 – Figure Supplement 3** Distribution of pairwise evolutionary distances for haplotypes  
723 within families. Black, observed distribution of pairwise evolutionary distances; green, gamma  
724 distribution; blue, exponential distribution; orange, sum of Gamma distribution plus Exponential  
725 Distribution. The chosen cut-off distance to differentiate small variations from large differences  
726 is the crossing point of the two distributions, at 0.017.

727  
728 **Figure 4.** Abundance of haplotypes within each sample plotted for breast milk (BM), Cervix (CV)  
729 and Blood spots (BS). The timing of sampling is shown along the x axis. For ease of reference, the  
730 genotype containing the most abundant haplotype present in the cervix is coloured red for each  
731 family. Thereafter sequences that are genetically closest to the red genotype (**Figure 3 – Figure**  
732 **Supplement 2**) are coloured magenta. Genotypes that are as distant from the cervical genotype as  
733 unrelated GenBank sequences are coloured shades of green, blue and purple. Single variants are  
734 coloured in shades of the nearest genotype.

735  
736 **Figure 4 – Figure Supplement 1.** Boxplot showing number of haplotypes reconstructed in  
737 relation to read depth. Analysis was performed on the 12-month breastmilk sample from family  
738 12.

739

740 **Figure 5.** The magnitude of FST values plotted for each gene (x axis). P values, adjusted with false  
741 discovery rate are shown in Red for  $p < 0.01$ , Grey for  $p > 0.05$  and turquoise for  $p = 0.01-0.05$ .

742

743 **Figure 5 – Figure Supplement 1.** Distribution of confidence-weighted sums of FST (cwsFST) values  
744 for all subsets of two (cyan), three (purple), four (green) and five (magenta) genotypes from  
745 different mother-baby pairs. For comparison, we also show the distribution obtained when the  
746 genotype sequences corresponding to each mother-baby pair are scrambled (black line). Arrows  
747 mark the values for the five genotypes that predominated in the cervical samples (black), the  
748 three predominant genotypes from cervical samples for mother-baby pairs 12, 22, and 123  
749 (blue), and the two predominant genotypes from cervical samples for mother-baby pairs 14 and  
750 41 (red).

751

752 **Figure 5 – Figure Supplement 2.** Heatmap showing genes identified as significant in FST analysis  
753 are robust to changes in the number of clusters. Colors indicated the false discovery rate value, red  
754 =  $< 0.001$ ; magenta =  $0.001-0.01$ ; pink =  $0.01-0.05$ ; purple =  $0.05-0.1$ ; blue =  $0.1-0.2$ ; grey =  $> 0.2$ .

755

756

## 757 **Tables**

758 **Table 1.** Sequencing characteristics for samples from each family. OTR: on target read; %  
759 Genome: % of genome coverage; % Dup: % of duplicated reads. Samples with genome coverages  
760 too low to be included in any analysis are shaded in grey. Cervical or baby samples with good  
761 coverage and read depth are highlighted in yellow.

Sample	%OTR	%Genome	%Dup	Mean Depth	Viral Load
<b>Family 12</b>					
Breast milk 2W	26.41	99	29.49	224.45	1235136.63
Breast milk 6W	68.99	99	13.84	578.56	14926741
Breast milk 14W	76.4	99	5.02	683.04	7309960
Breast milk 6M	77.47	99	8.07	730.04	10876521
Breast milk 12M	77.81	99	7.68	779.72	6135712.5
Cervix 38W Pregnant	14.73	99	47.56	325.97	95842
Baby Delivery	1.35	76	82.27	31.86	27393.9395
Baby 6W	0.02	2	81.79	0.29	4067.86694
Baby 10W	0.1	12	77.77	2.63	1959.9679
Baby 9M	1.1	78	79.41	28.53	2501.75195
<b>Family 14</b>					
Breast milk 2W	13.54	98	65.41	101.66	232442.219
Breast milk 6W	60.32	98	49.85	656.47	20485190
Breast milk 14W	11.15	97	65.77	80.09	345851.781
Cervix 38W Pregnant	0.22	63	56.04	4.34	1377
Baby 6W	1.4	91	69.35	21.35	55400.7148
Baby 14W	3.33	96	78.59	113.92	3960.64233
Baby 6M	0.34	66	74.11	11.42	154.414169
Baby 12M	0.02	7	75.97	0.75	3054.47485
<b>Family 22</b>					
Breast milk 2W	6.08	96	34.22	54.34	55000.2891
Breast milk 6W	43.18	98	44.57	352.49	107861.141
Breast milk 14W	6.4	97	44.41	38.3	56883.9805
Cervix 34W Pregnant	0.16	46	54.95	2.97	1125
Cervix 38W Pregnant	0.16	67	47.91	4.14	1377
Baby 2W	0.01	1	46.34	0.03	1703.49292
Baby 6W	0.08	1	43.61	0.03	22082.6465
Baby 14W	2.29	92	79.42	46.53	10962.7197
Baby 6M	0.3	33	79.36	5.98	2124.86548
Baby 9M	0.22	25	79.33	5.01	82937.5
<b>Family 41</b>					
Breast milk 2W	43.33	98	60.89	224.53	7163743
Breast milk 6W	37.05	98	61.89	289.61	323325.531
Breast milk 14W	48.15	98	68.02	438.05	2697832.75
Cervix 38W Pregnant	0.61	91	47.53	12.6	122
Baby 14W	0.12	32	74.47	4.67	1848.62402

Family 123					
Breast milk 2W	16.11	98	60.11	117.25	518071.875
Breast milk 6W	16.96	98	64.77	107.35	262400.719
Breast milk 14W	13.95	98	64.01	122.08	518071.875
Breast milk 6M	15.81	98	63.07	101.92	678250.313
Cervix 34W Pregnant	2.45	97	49.46	41.91	7931
Cervix 38W Pregnant	1.36	96	49.61	28.07	4326
Baby Delivery	0.21	84	10.93	6.1	939.190735
Baby 10W	2.19	91	78.64	43.96	93297.3047
Baby 6M	0.13	20	77.67	3.1	5428.83545
Baby 12M	1.36	85	80.13	40.56	6205.88281

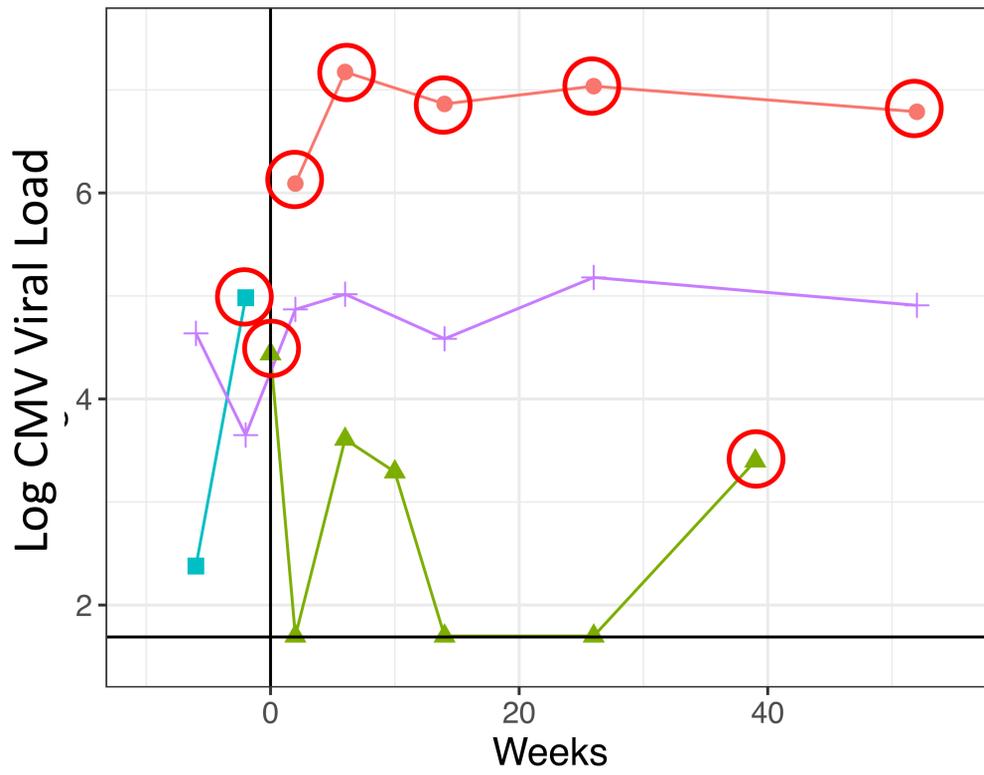
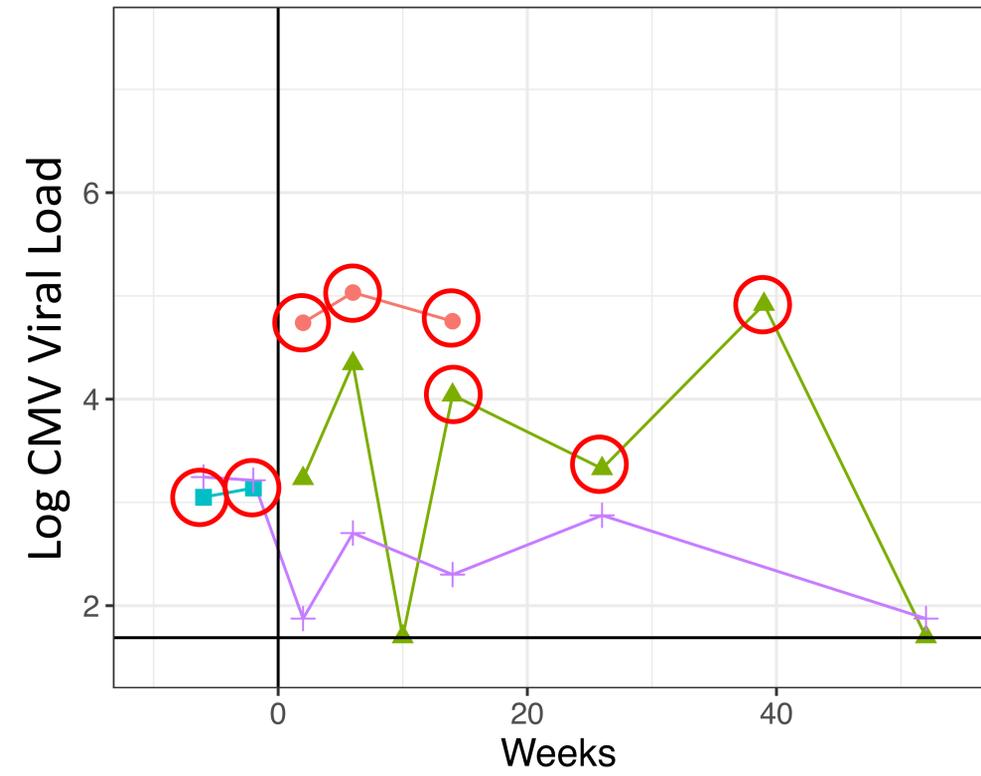
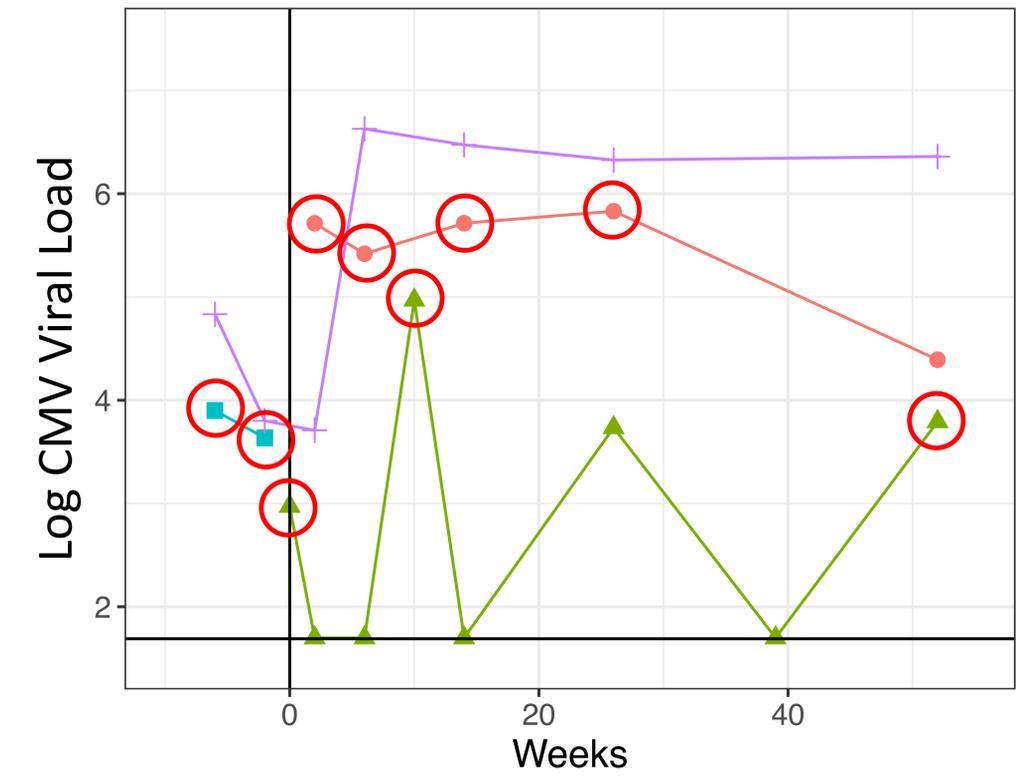
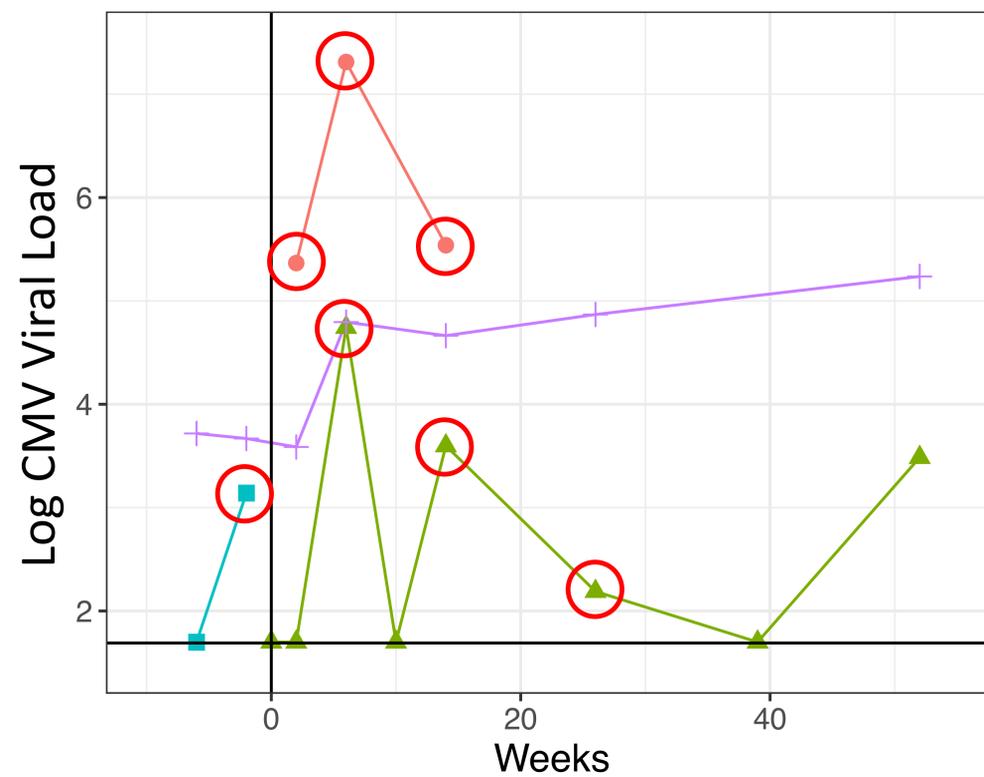
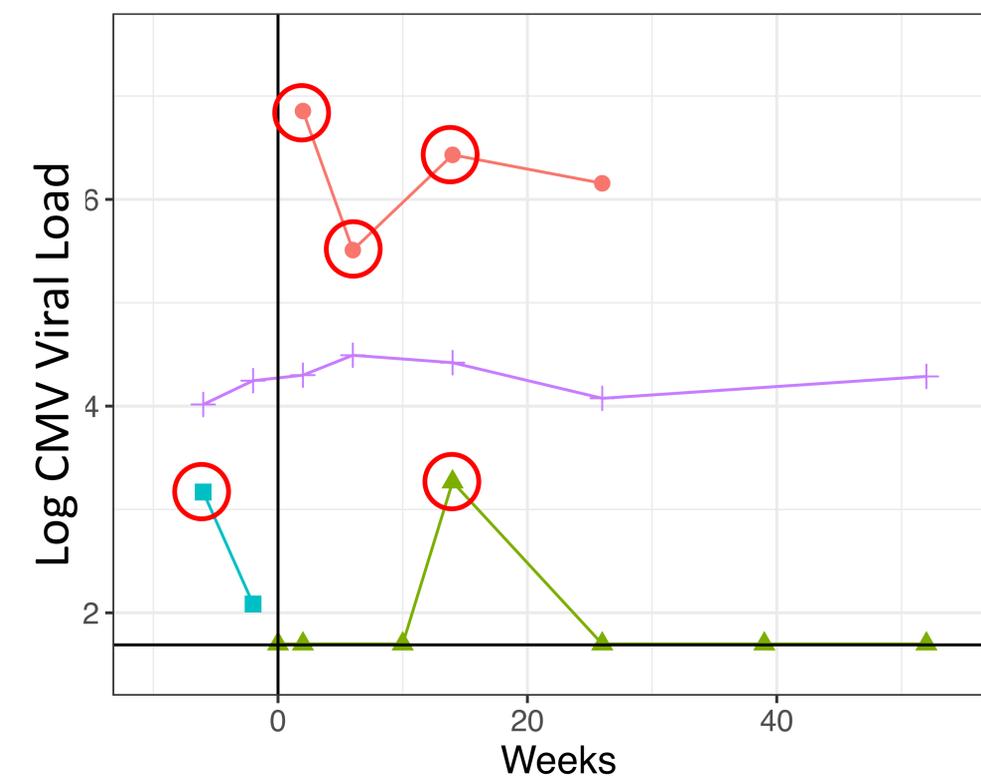
762

763 **Table 2.** Open Reading Frames (ORFs) identified by FST as being significantly more similar in  
764 strains transmitted prenatally. LD: Found to contain one of 33 hotspots of genetic linkage  
765 disequilibrium [18].

766

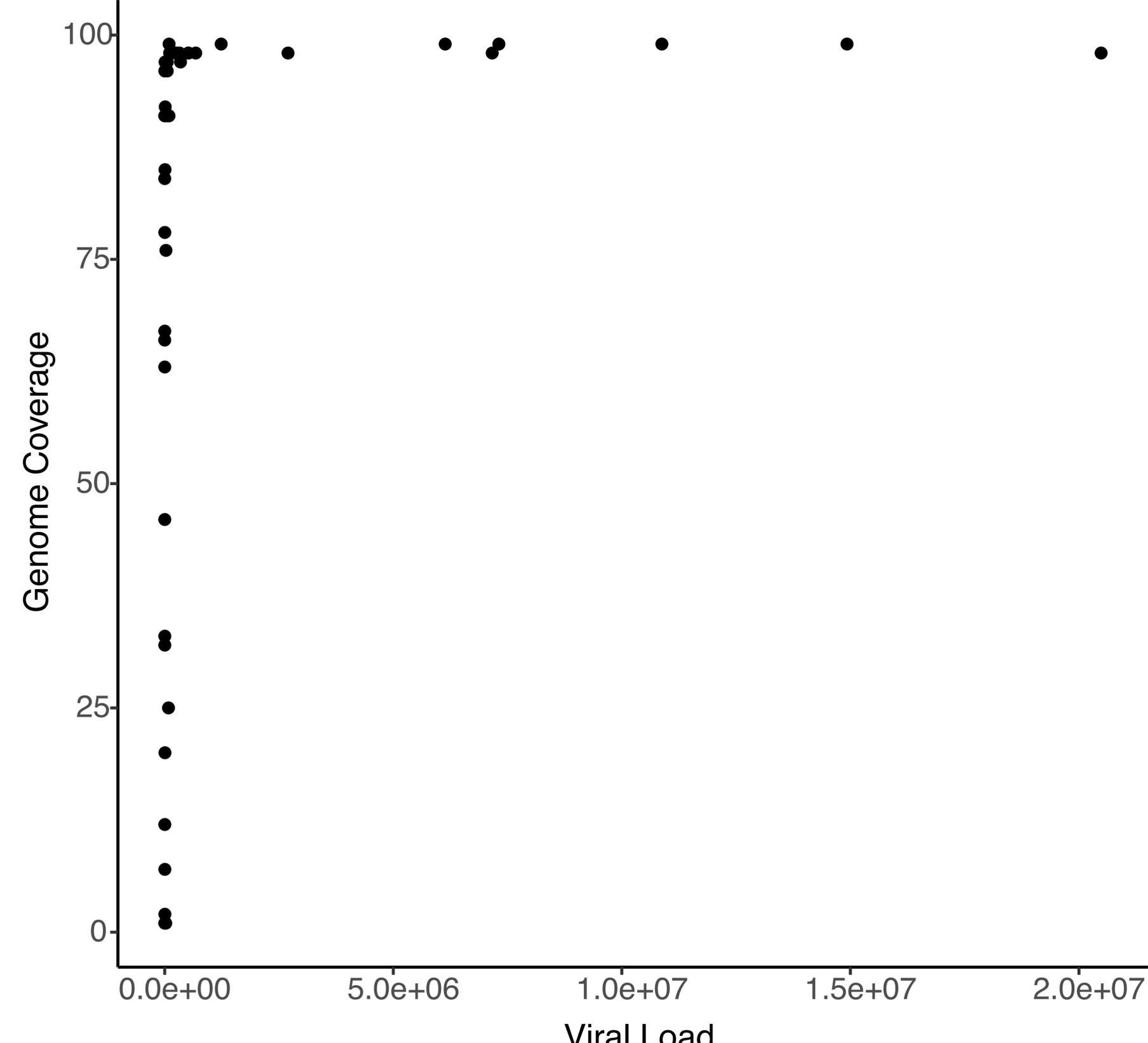
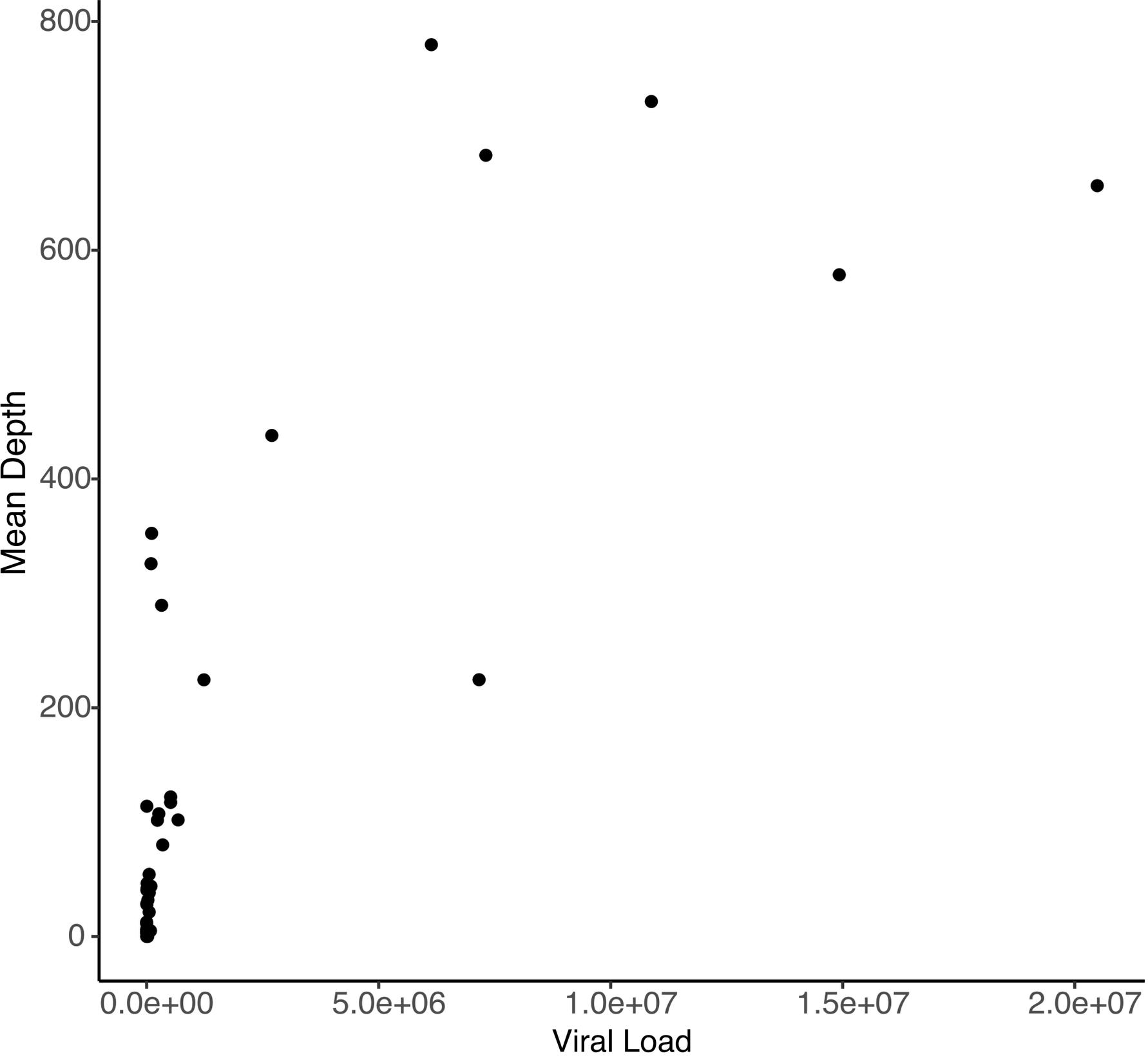
767

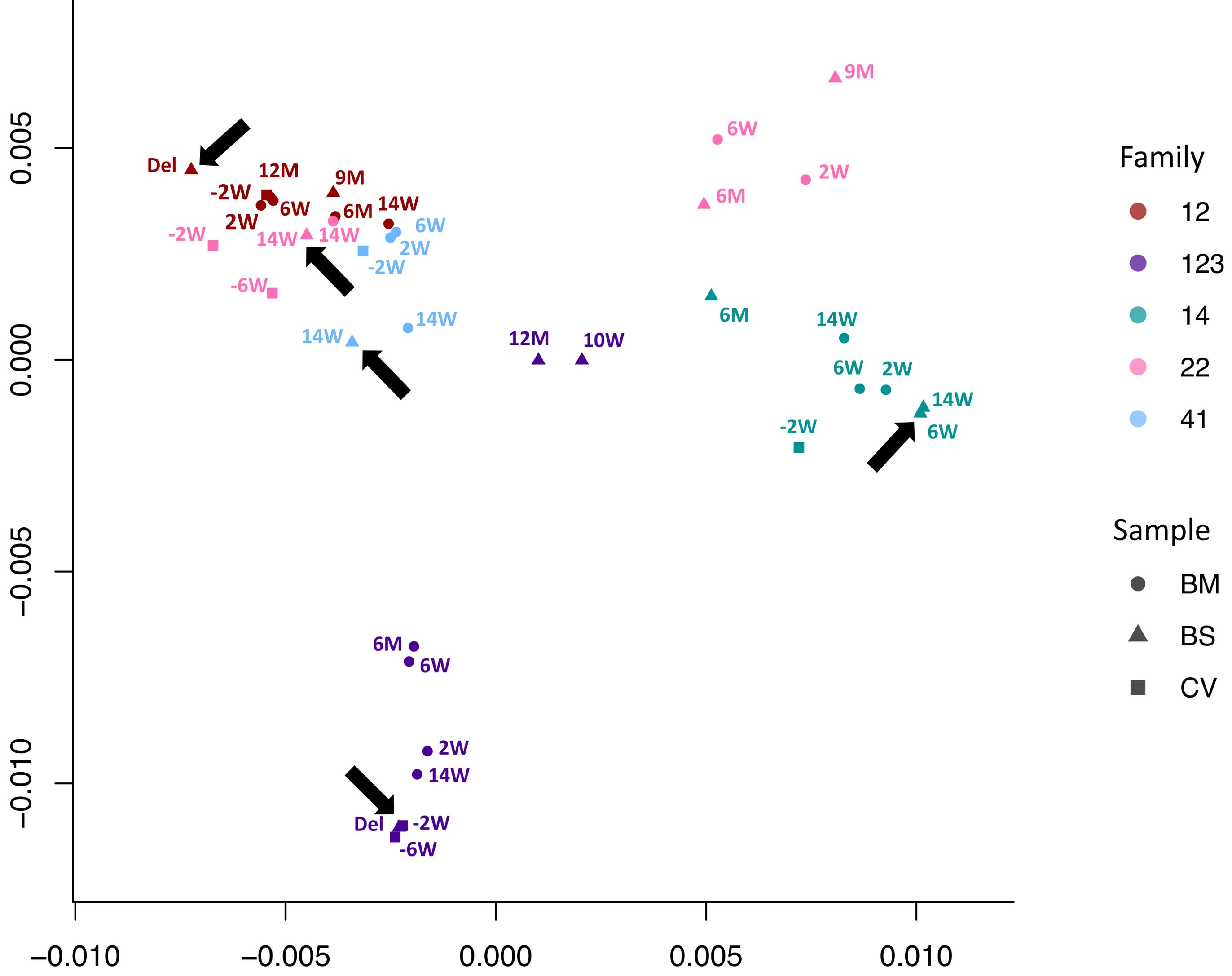
ORF	LD	FUNCTION
UL10	Y	Putative membrane glycoprotein, Immunosuppressive impairs T cell function [42]
UL11	Y	Membrane glycoprotein modulation of T cell signalling/function [43, 50]
UL13		Unknown function
UL4	Y	Putative membrane glycoprotein [40]
UL5		Putative membrane glycoprotein [40]
UL6	Y	Putative membrane glycoprotein [40]
UL7	Y	Membrane glycoprotein, modulates chemo-and/or cytokine signalling function [41]
UL8	Y	Transmembrane glycoprotein. Inhibits proinflammatory cytokines [41]
US26		Unknown function
US27	Y	Membrane glycoprotein Activates CXCR4 signalling to increase HCMV replication [45]
UL150A		Fibroblast and Epithelial cell entry [51]
UL2		Putative membrane glycoprotein [40]
RL11	Y	Membrane glycoprotein. Binds IgG Fc domain involved in immune regulation [40]
UL147		$\alpha$ -chemokine homologue [52, 53]
UL40		Control of NK recognition [46]
RL13	Y	Glycoprotein, repression of replication, bind IgG domain immune regulation [34, 39]
RL10		Membrane glycoprotein
UL57		Ss DNA binding protein [40]
UL50		Nuclear Egress complex. Reduces interferon mediated antiviral effect [48]

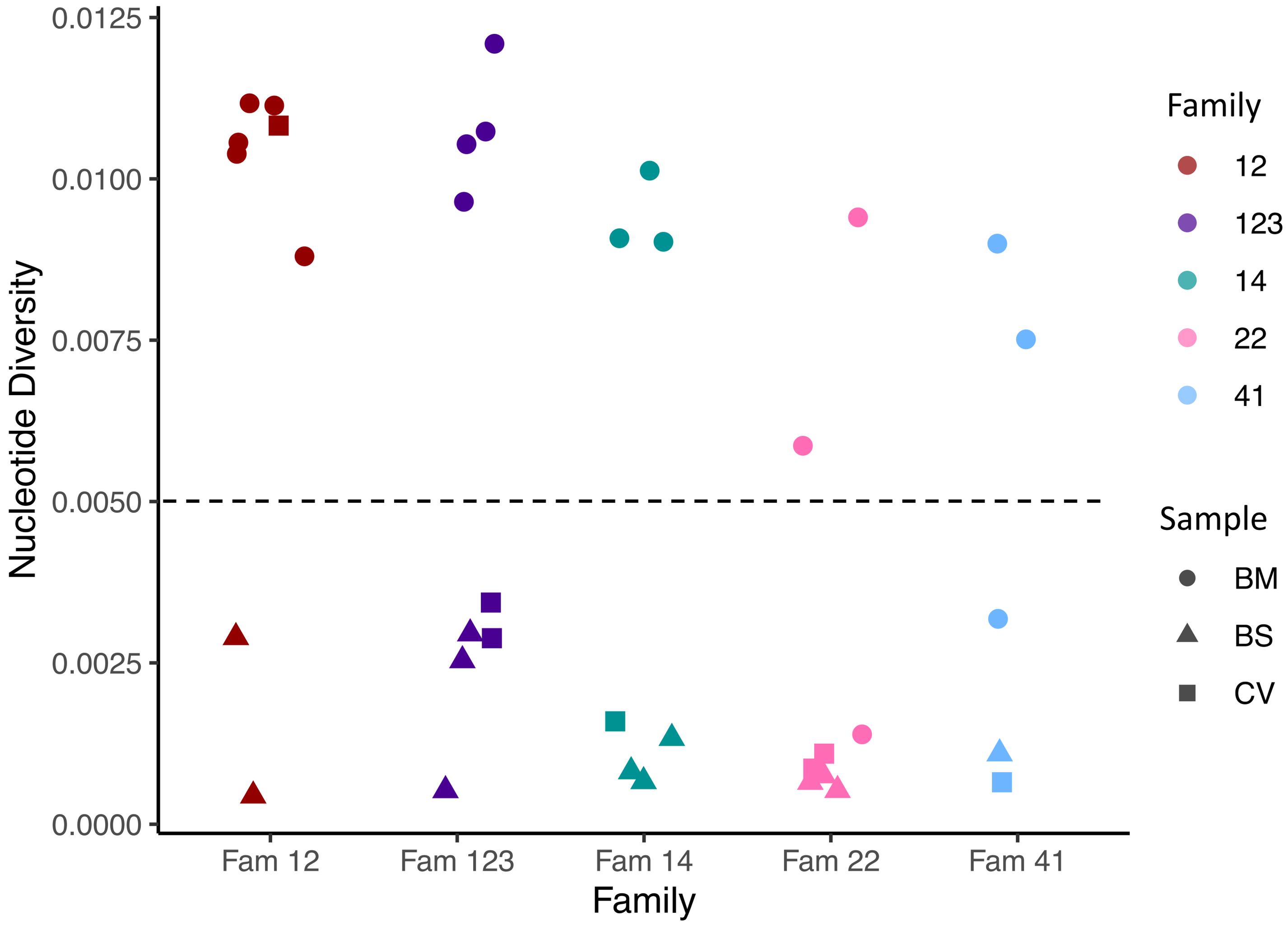
**Family 12****Family 22****Family 123****Family 14****Family 41**

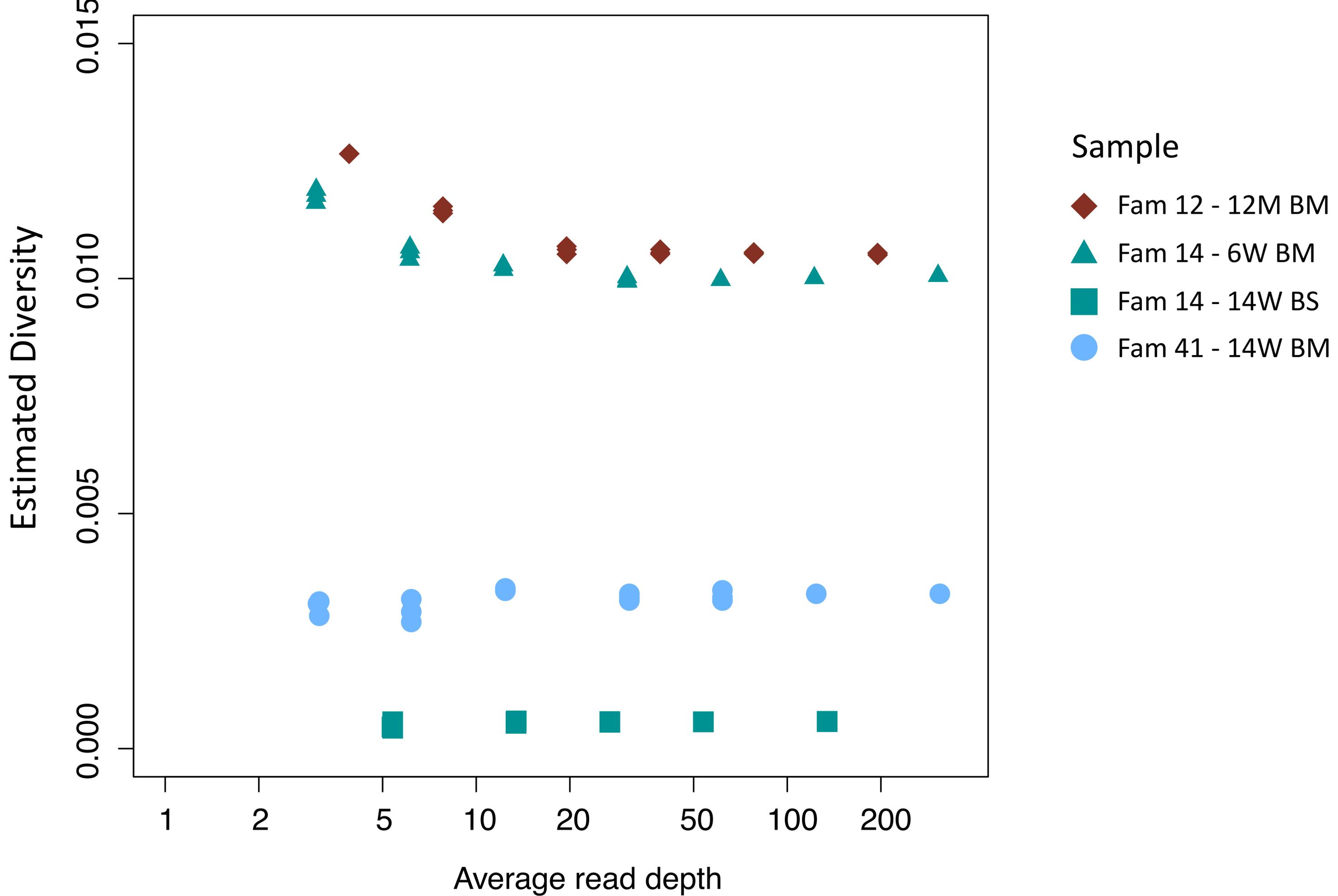
## Samples

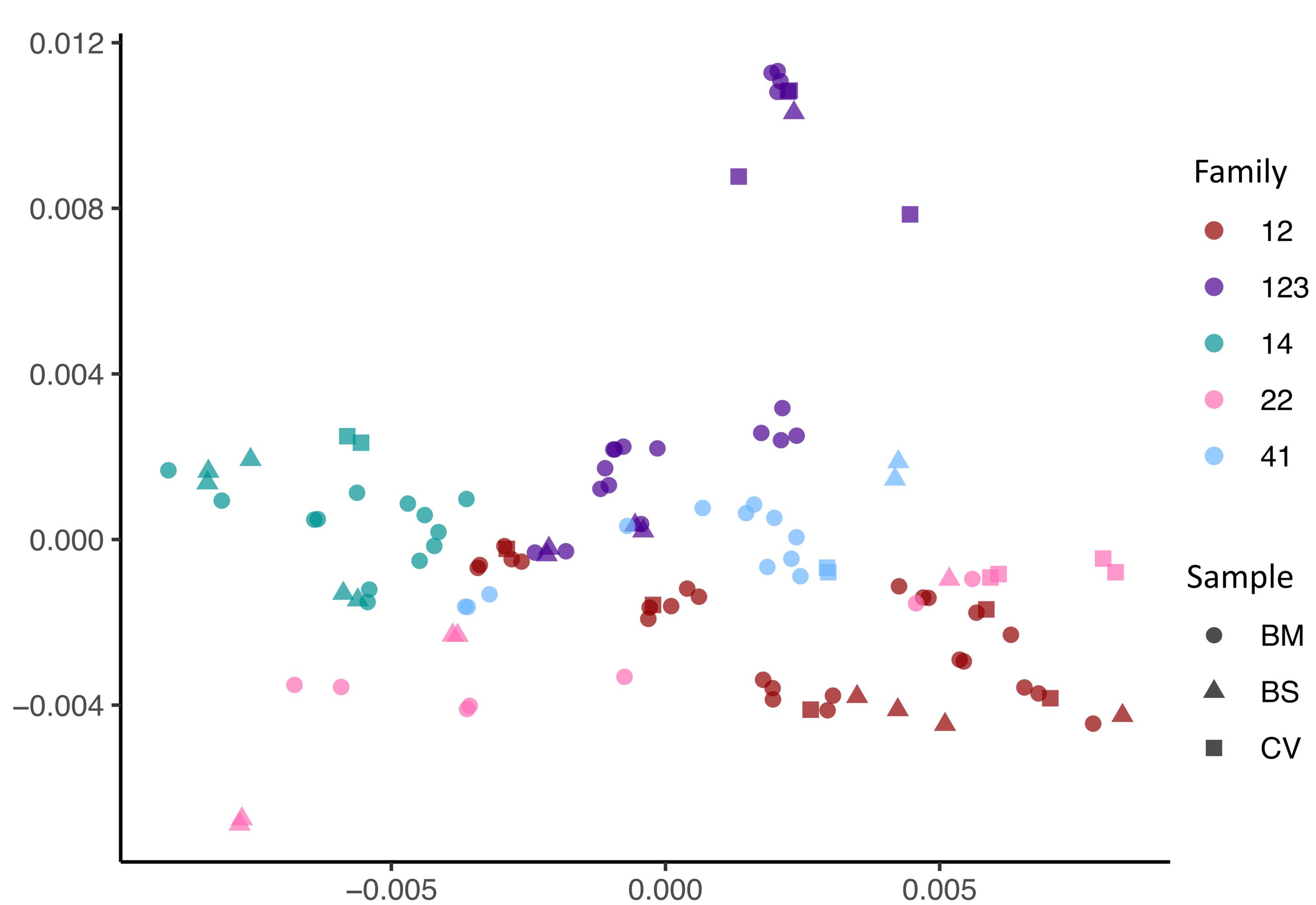
- Breast milk (BM)
- ▲ Baby blood spots (BS)
- Cervix (CV)
- Viral Sequence Obtained
- + HIV Plasma Viral Load

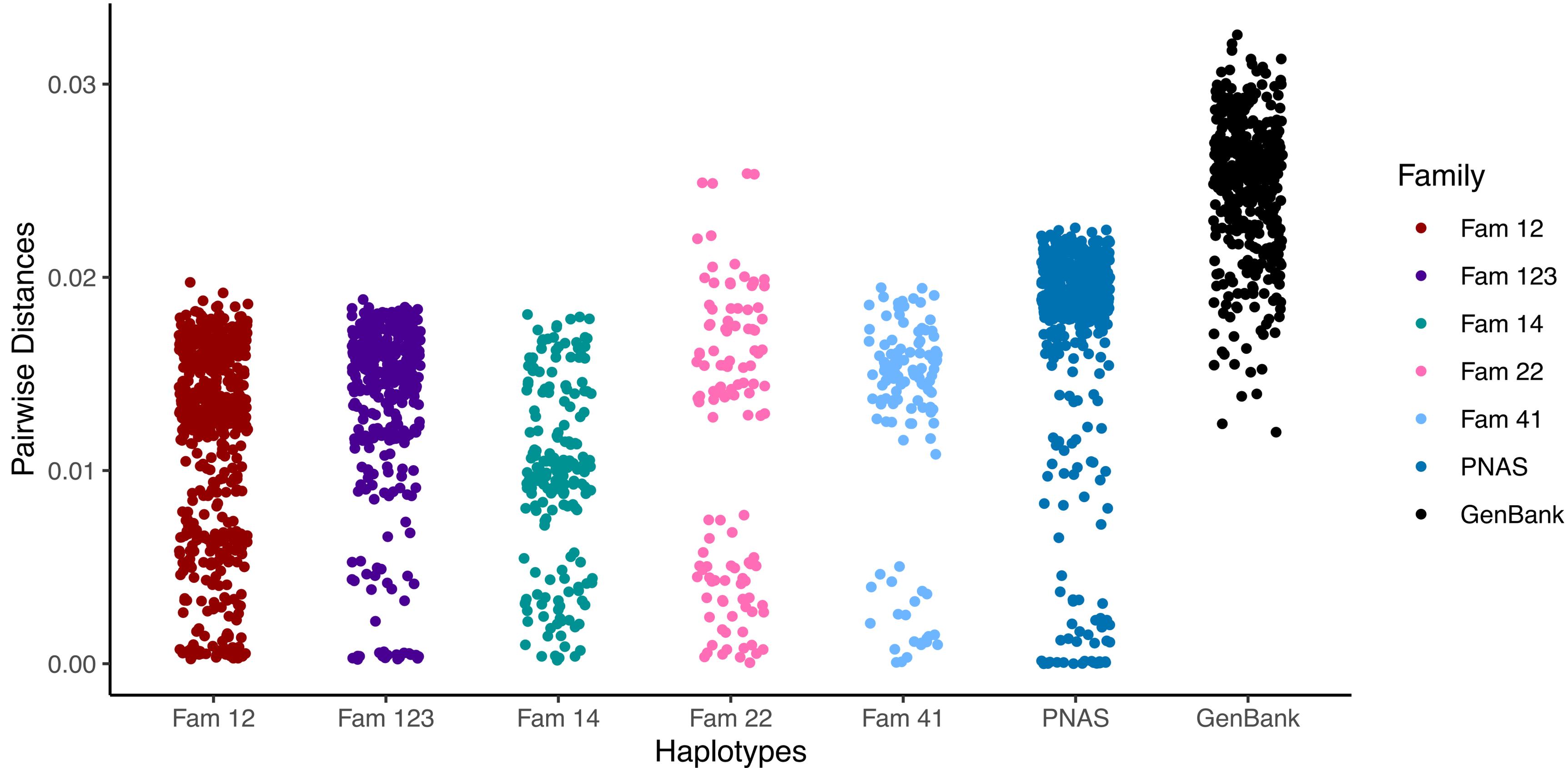




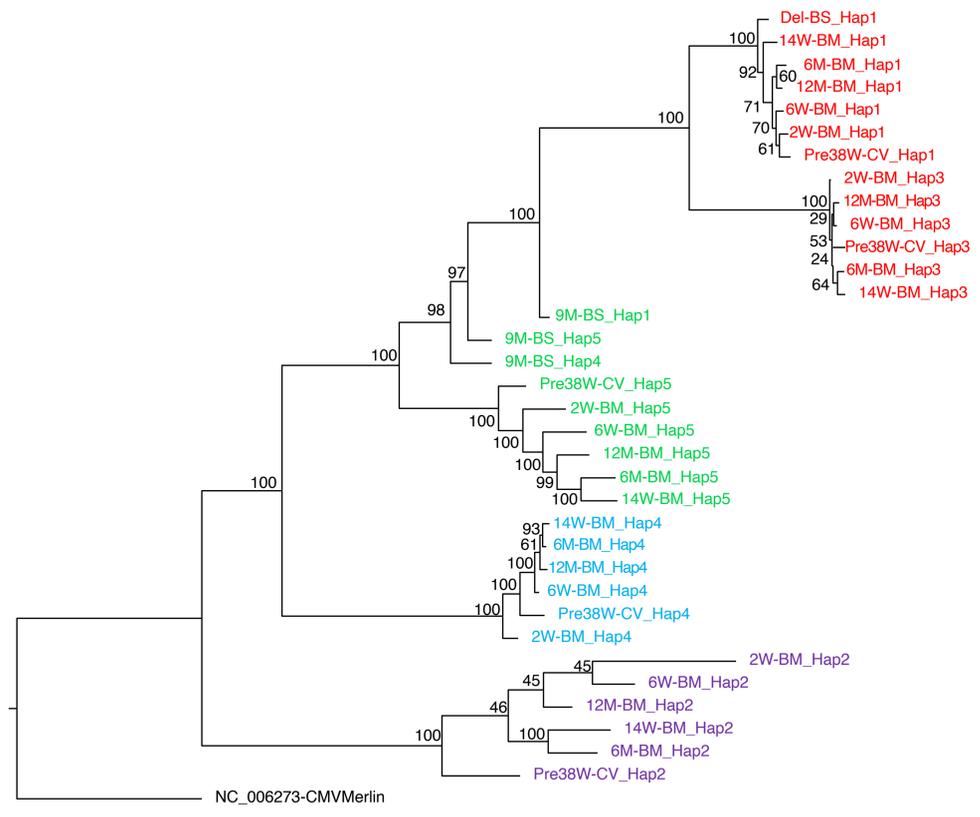




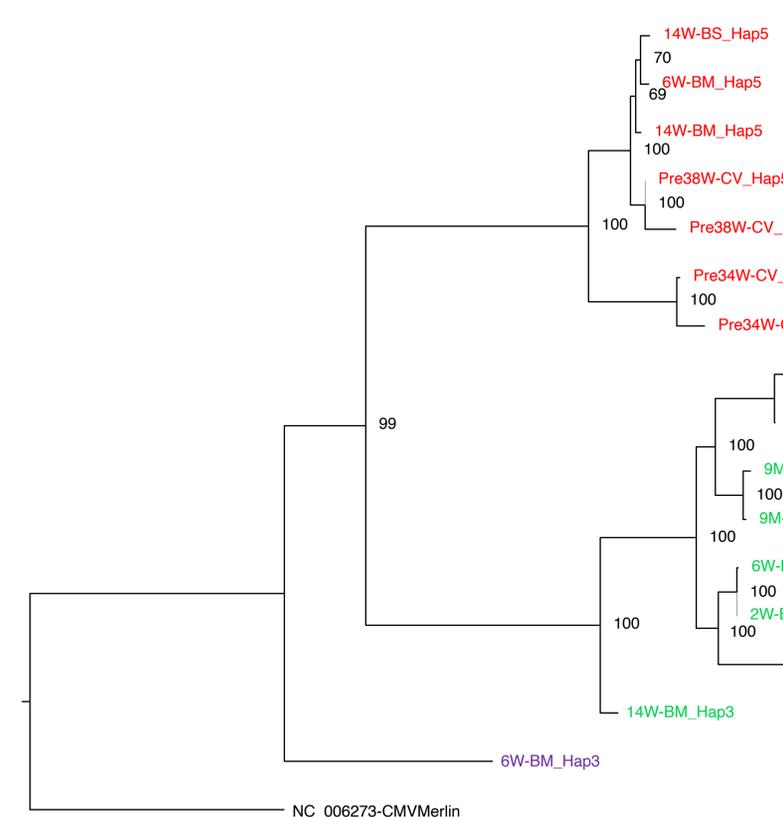




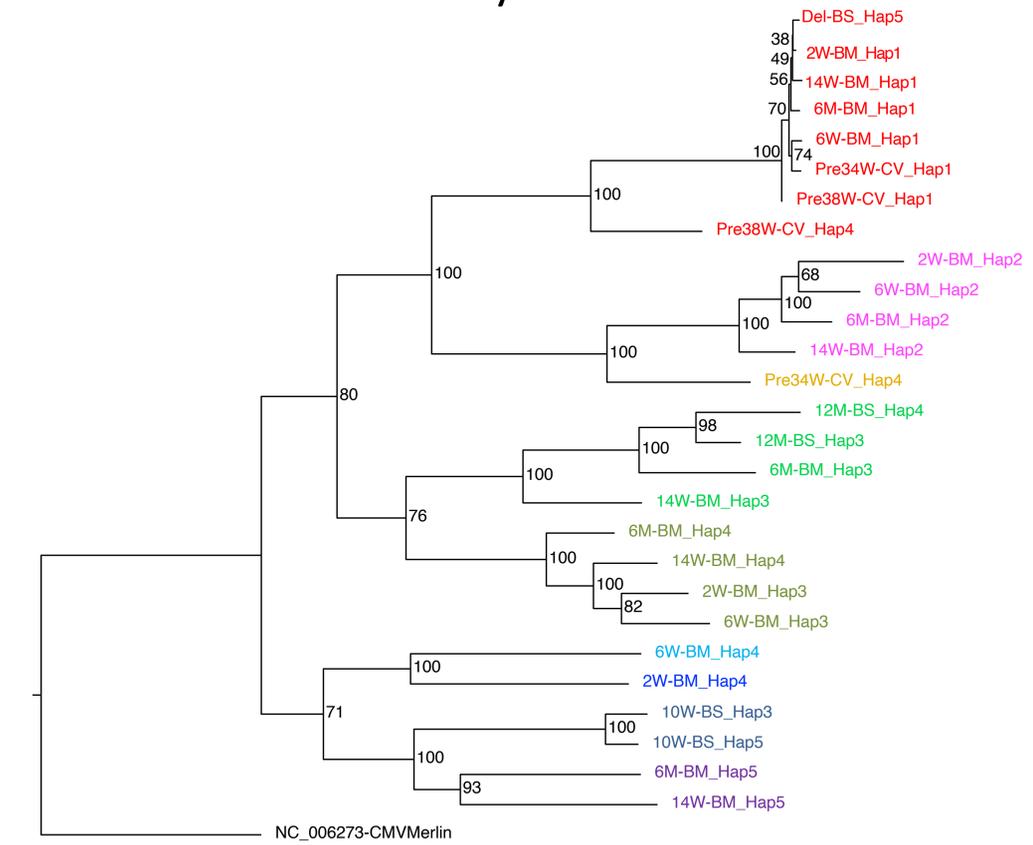
Family 12



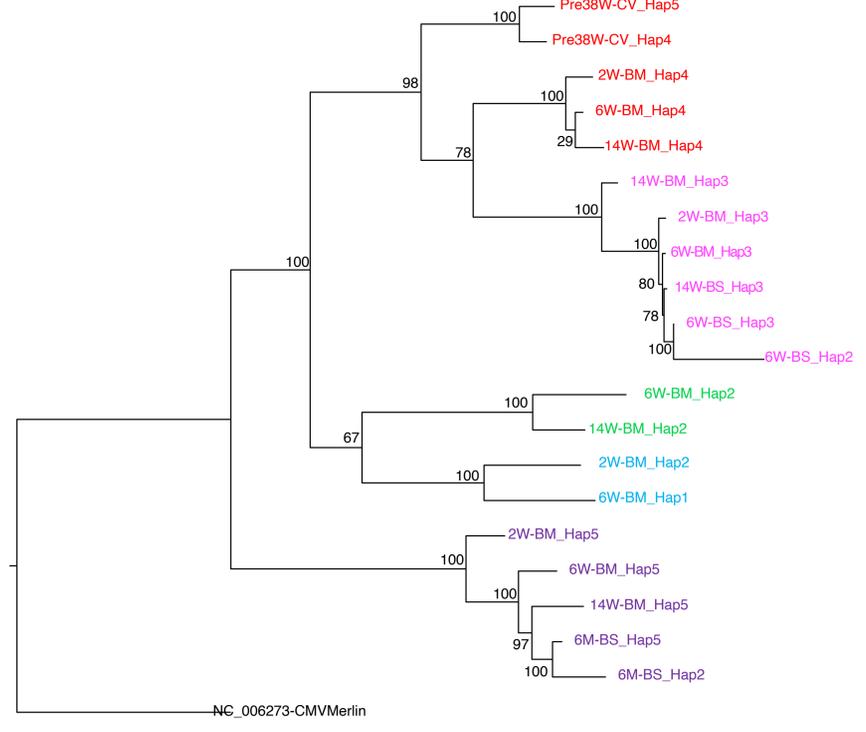
Family 22



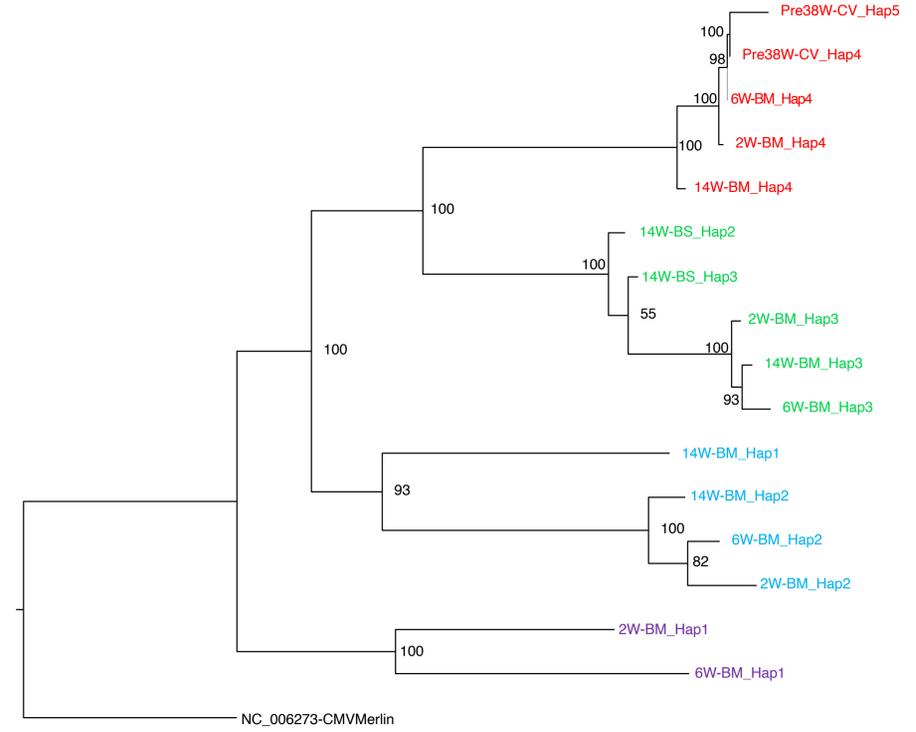
Family 123

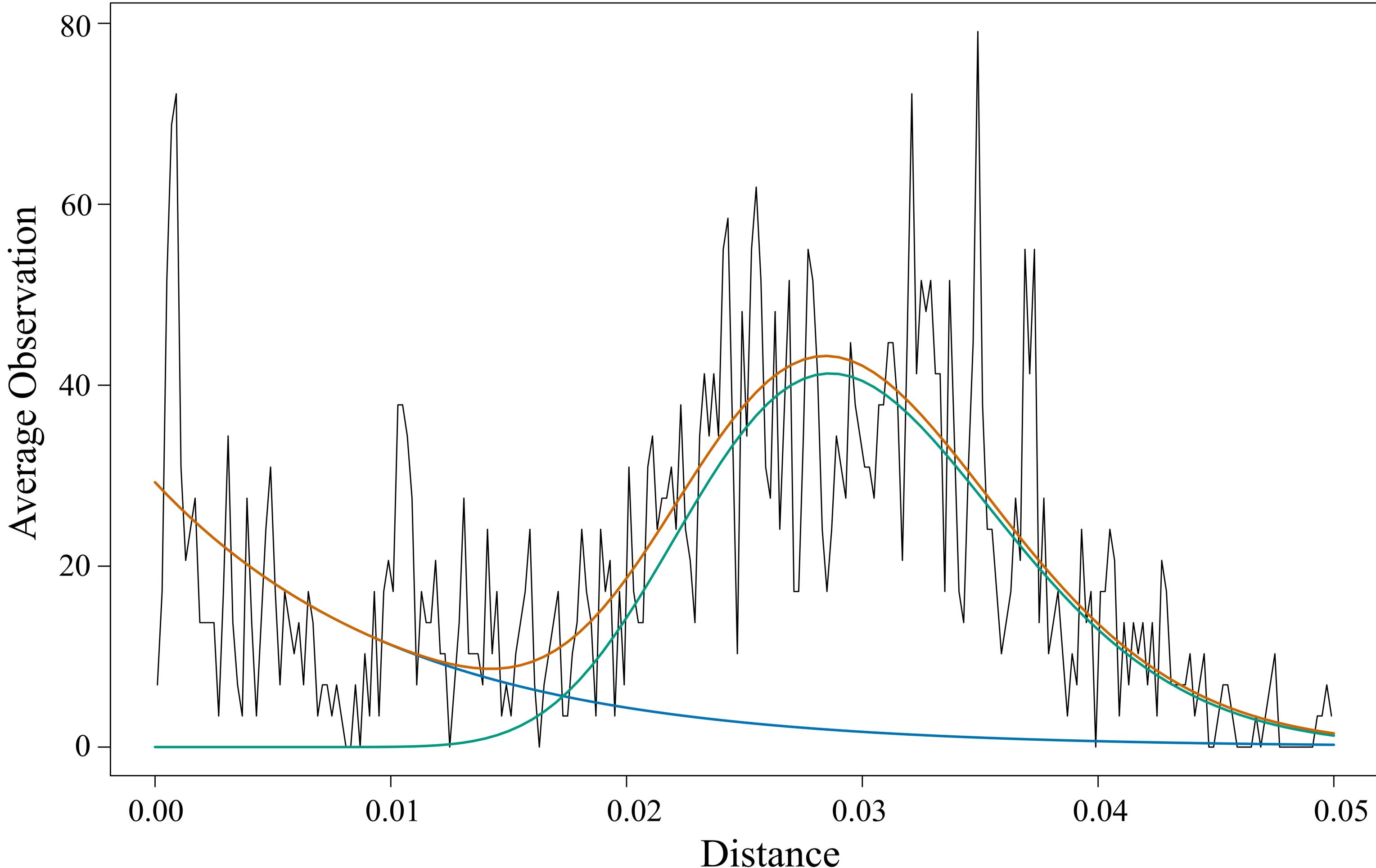


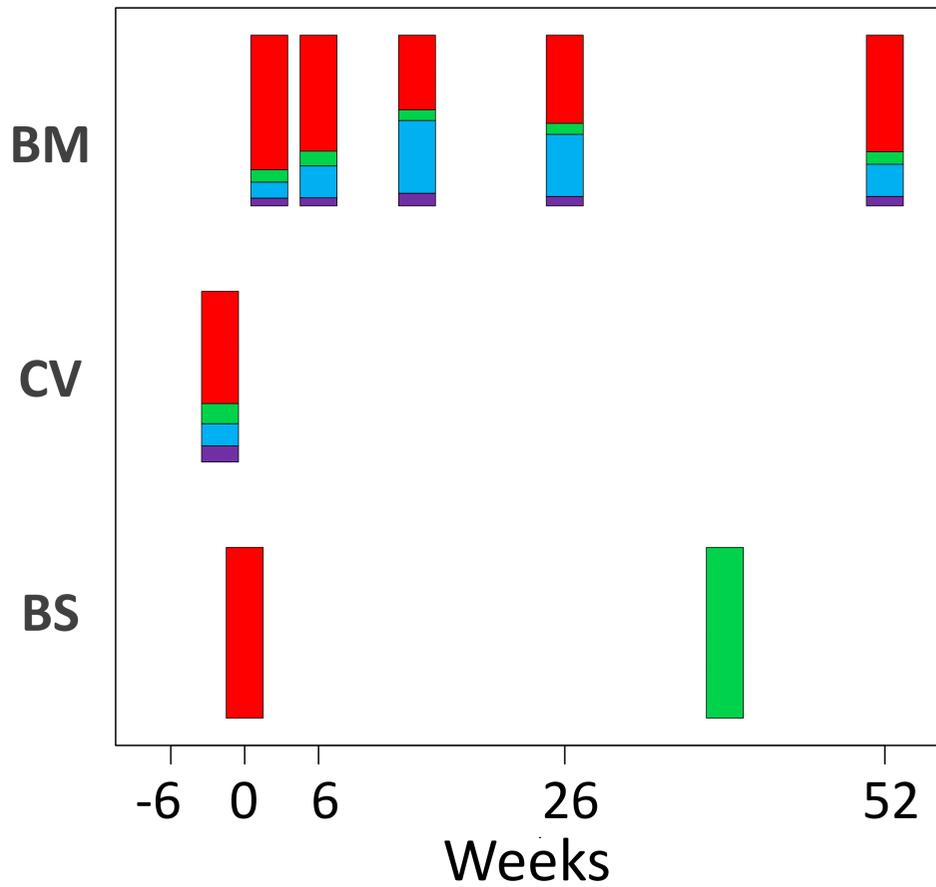
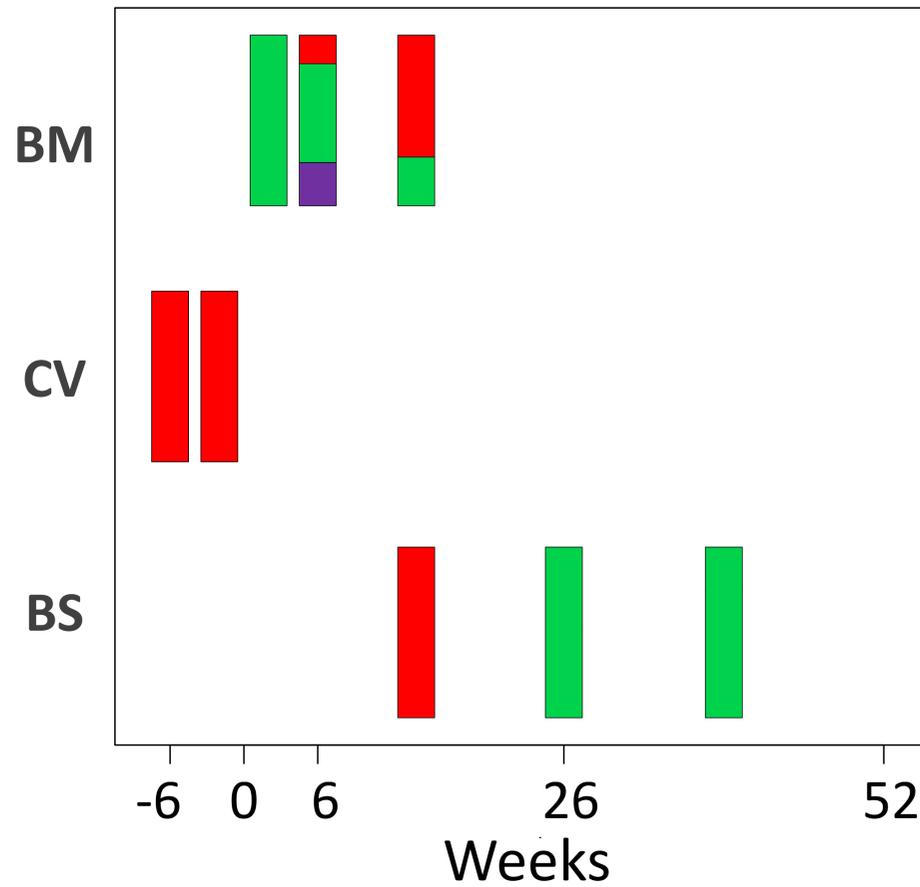
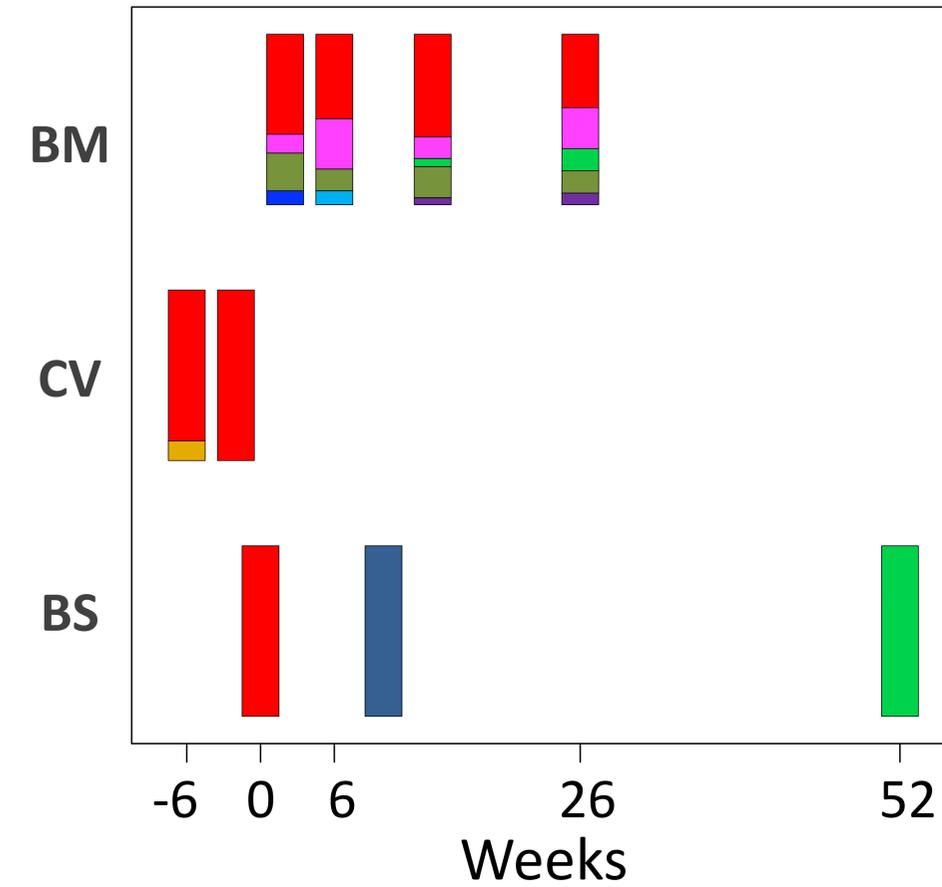
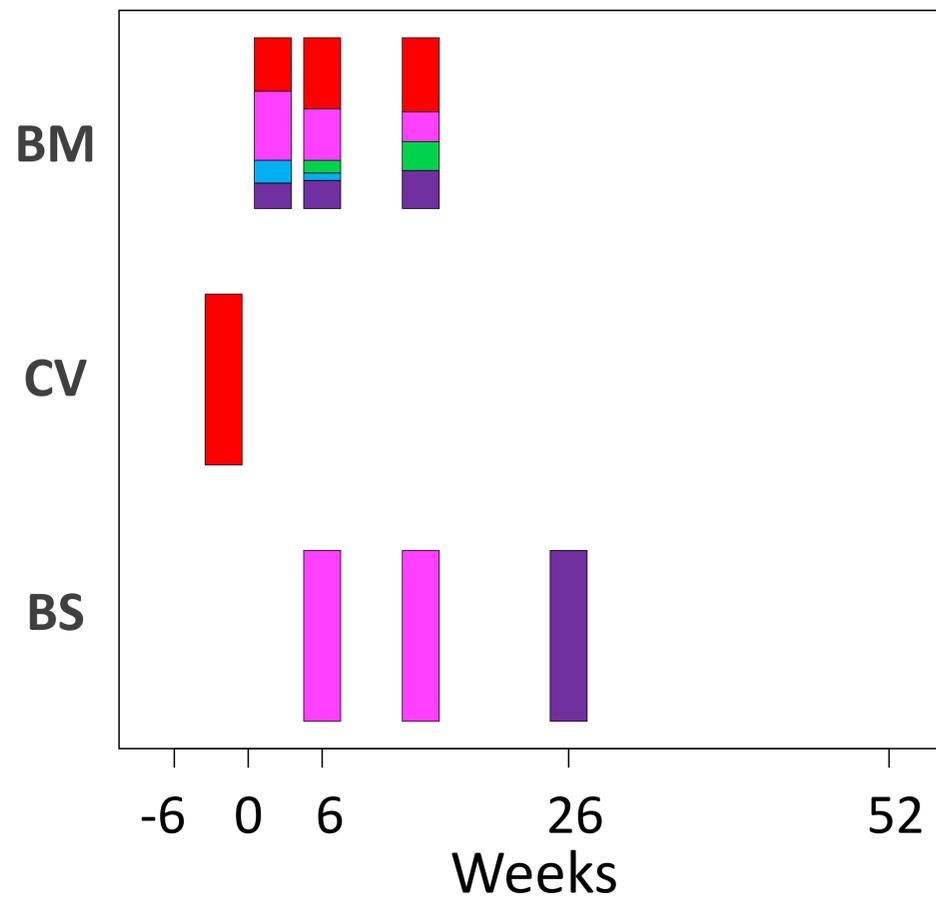
Family 14



Family 41





**Family 12****Family 22****Family 123****Family 14****Family 41**