# Genomic and geographical structure in Human Cytomegalovirus

Oscar J. Charles[1,+], Cristina Venturini[1,+,*], Soren Gantt[2], Claire Atkinson[3], Paul Griffiths[3], Richard A. Goldstein[4] and Judith Breuer[1,4]


1. Department of Infection, Immunity and Inflammation, UCL Great Ormond Street Institute of Child Health, London, UK
2. Research Centre of the Sainte-Justine University Hospital and Department of Microbiology, Infectious Diseases and Immunology, University of Montréal, Canada
3. Division of Infection and Immunity, Institute for Immunity and Transplantation, University College London, London, UK
4. Division of Infection and Immunity, University College London, The Cruciform Building, Gower St, London, UK
5. Great Ormond Street Hospital for Children NHS Foundation Trust, London, UK


[+]These authors contributed equally to this work.

[*]Corresponding author:

Cristina Venturini

c.venturini@ucl.ac.uk

**Classification**

Biological Sciences

Evolution

**Keywords**

Human cytomegalovirus, genotyping, hypervariability, genomics, Hidden Markov Models, phylogeography, molecular population genetics

36

37

38

39 **Abstract**

40 Human cytomegalovirus (CMV) has infected humans since the origin of our species and currently

41 infects most of the world's population. Variability between CMV genomes is the highest of any

42 human herpesvirus, yet large portions of the genome are conserved. Here we show that the genome

43 encodes 74 regions of relatively high variability each with 2-8 alleles. We then identified two

44 patterns in the CMV genome. Conserved parts of the genome and a minority (32) of variable regions

45 show geographic population structure with evidence for African or European clustering, although

46 hybrid strains are present. We find no evidence that geographic segregation has been driven by host

47 immune pressure affecting known antigenic sites. Forty-two variable regions show no geographical

48 structure, with similar allele distributions across different continental populations. These "non-

49 geographical" regions are significantly enriched for genes encoding immunomodulatory functions

50 suggesting a core functional importance. We hypothesise that at least two CMV founder populations

51 account for the geographical differences that are largely seen in the conserved portions of the

52 genome, although the timing of separation and direction of spread between the two is not clear. In

53 contrast, the similar allele frequencies among 42 variable regions of the genome, irrespective of

54 geographical origin, is indicative of a second evolutionary process, namely balancing selection that

55 may preserve properties critical to CMV biological function.  Given that genetic differences between

56 CMV viruses are postulated to alter immunogenicity and potentially function, understanding these

57 two evolutionary processes could contribute important information for the development of globally

58 effective vaccines and the identification of novel drug targets.

59

60 **Significance statement**

61

62 CMV genome diversity is higher than other human herpesvirus, and recombination is pervasive.

63 Here, using Hidden Markov modelling, we describe 74 multi-allelic regions, with the remaining 86%

64 of the genome showing lower variability, albeit with single nucleotide polymorphisms. We

65 demonstrate for the first time that CMV diversity is influenced by two distinct evolutionary forces. A

66 founder effect results in geographical segregation affecting the regions of low variability and 32

67 variable regions. In contrast, the 42 remaining regions, which are enriched for immunomodulatory

68 functions, show so-called balancing selection, resulting in maintenance of equal allele frequencies

69    irrespective of geography. These new insights into CMV evolution are likely to provide insights in

70    virus biology and inform the development of drugs and global vaccines.

71

72    **Introduction**

73    Human cytomegalovirus (CMV) is a member of the *Betaherpesvirinae* that infects circa 66% to 90%

74    of adults in any given country (1). Like all human herpesviruses, CMV is a linear double-stranded

75    DNA virus that causes lifelong latent infection by establishing latency in long-lived cell populations,

76    and periodically reactivates resulting in lytic viral replication (2, 3). CMV causes significant burden of

77    disease in those with compromised immune systems (4),  and is also the most common infectious

78    cause of congenital disability worldwide (5). Because of this, developing a vaccine is a high public

79    health priority (6).

80    At approximately 236 kb, CMV has the largest genome of all human herpesviruses (7) and the

81    highest level of genetic diversity of all the known human herpesviruses (8, 9). The virus is known to

82    readily undergo recombination (10, 11), and co-infection is frequently observed, especially in

83    individuals with weakened immune systems (12). Most of the observed CMV diversity occurs in

84    discrete hypervariable regions where sequences cluster into genotypes also known as alleles (13,

85    14). In some hypervariable genes, e.g. UL55 (glycoprotein B, gB), alleles have been defined (15) and

86    used alone or in combination with other multi-allelic regions to genotype CMV and identify mixed

87    infections (16, 17). However, attempts to correlate individual alleles with transmission and

88    pathogenesis have so far been unclear or contradictory (14, 18–24).

89    Despite the efforts to define CMV diversity, our understanding of the evolutionary history that led to

90    its present variation and whether this variability follows a global pattern is limited. CMV, like other

91    herpesviruses,  shows remarkable species-specificity, which results from long-term co-evolution and

92    adaptation to the host  (25, 26). In contrast, unlike other herpesviruses, human CMV (HCMV)  whole

93    genomes from clinical samples show little evidence of geographical or other  population structure,

94    with the exception of two Asian genomes that have been shown to cluster phylogenetically (10, 11).

95    This, together with the absence of data from ancient CMV genomes, makes for uncertainty as to

96    when and how current HCMV diversity evolved.

97    To better understand CMV genomic structure and how it has evolved, we employed Hidden Markov

98    Model (HMM) clustering to delineate the hypervariable and conserved regions across a global and

99    diverse dataset of published and unpublished CMV genomes, which together represent unrelated

100   clinical and low-passage strains (27–29). HMM is able to determine the number of sequence clusters

101   (i.e. alleles) that best explain the diversity across CMV genomes and to identify regions where

102   multiple alleles are present. Using the outputs of the model we describe precise co-ordinates for

103     regions of multi-allelic variability, some of which are novel. We also show, for the first time, that

104     CMV genomes do display geographic population structure and that this is particularly clear within

105     the relatively conserved monoallelic regions of the genomes. We highlight examples of where our

106     approach can help to provide insights for further research into questions of viral pathogenesis and

107     ancient evolution.

108

109     **Results**

110     **CMV diversity and population structure is determined by 74 discrete multi-allelic regions**

111     We compiled a set of 259 CMV whole genome sequences which had been collected worldwide

112     (Supplementary figure 1, full details are shown in Dataset S1). 233 sequences were retrieved from

113     GenBank (30) with available metadata for 106. Of these 106, 35 were from patients who were

114     immunocompromised through HIV, organ, or bone marrow transplant and 71 were from immune-

115     competent individuals (of whom 60 were congenitally infected babies). Short read data were

116     available for 17 samples allowing us to check for mixed infections. Nine samples contained only a

117     single CMV strain, of which five were from HIV-positive and CMV-positive mothers and four were

118     immuno-competent children (primary infection or sibling). Eight samples showed evidence of

119     multiple strains for which we reconstructed the haplotypes (17 sequences in total) with validated

120     methods (8, 31). All eight samples came from HIV-positive mothers with CMV infection (8, 32).

121     Most genomes were derived from virus obtained in Europe and the Middle East (n=216, including

122     Belgium, Czech Republic, France, Germany, Greece, Italy, Netherlands, United Kingdom, and Israel).

123     As Israeli genomes appeared as if European in analyses, and the migratory history of the country is

124     complex, we have labelled Israeli sampled genomes as European for simplicity. 30 sequences

125     (including the reconstructed haplotypes) were from samples collected in Africa (Zambia, Kenya and

126     Uganda) (17) (32–34). The remaining sequences were collected in other parts of the world: 11 from

127     the United States (America); 2 from Asia (China and South Korea) and 1 from Australia (Oceania).

128     To characterise CMV genomic diversity we first aligned the genomes and then calculated

129     heterozygosity along a sliding window of 50 base pairs (bp) (Track1 in red, Figure 1). This revealed

130     that while most of the genome is highly conserved there are regions of significant nucleotide

131     diversity, some of which have previously been described (35, 36). Aligning HCMV genomes in these

132     hypervariable regions is challenging and leads to numerous gaps in alignments that are often

133     ignored in phylogenetic and genetic distance calculations. To deal with the complexity of CMV

134     genomes alignment in these regions, we developed a sequence clustering method based on Hidden

135     Markov Model called "HMMcluster". This approach is unaffected by gaps and it groups together

136     sequences based on the statistical likelihood that they came from the same underlying source. The

137      genome-wide implementation has two steps: in the first (1) a fixed window is defined (in this case

138      200bp, chosen as it returned regions no smaller than 20bp) and the model calculates the minimum

139      number of sequence clusters (alleles) that best explain the diversity within the window (37). The

140      second step (2) concatenates contiguous windows of sequence variation, and it refines the

141      coordinates for a "variable region". Variable regions with multiple alleles are defined as "multi-

142      allelic".

143      From this we identified 74 discrete variable regions for which the model provided statistical support

144      for the presence of more than one allele (Track 3 in blue - Figure 1, Table S1, Table S2). These multi-

145      allelic regions range in size from 26 to 4760 nucleotides, encompassing 14% of the genome (Figure 1,

146      Table S1). The remaining 86% of the genome was found to be highly conserved with no statistical

147      support for multiple alleles. Whilst multi-allelic regions generally occur in regions of high nucleotide

148      diversity, some, generally smaller in size and with fewer segregating Single Nucleotide

149      Polymorphisms (SNPs), were found in regions of otherwise lower nucleotide diversity. The 74 multi-

150      allelic regions were not constrained to known coding regions and encompass the previously

151      documented 12 hypervariable genes (7). Fifty-two of the regions were each contained entirely

152      within a single gene, 17 crossed gene boundaries by either spanning the gap between two genes or

153      extending beyond a gene terminus into non-coding regions, and 5 regions were entirely in non-

154      coding or otherwise unassigned portions of the genome.

155

156      To characterise the evolutionary relationships between the alleles within each multi-allelic region we

157      constructed maximum likelihood phylogenetic trees for each of the multi-allelic regions, excluding

158      10 with evidence of recombination breakpoints (Table 1 and Supplementary figure 2). For the

159      remaining 64/74 we observed that multi-allelic regions consisted of well separated clades

160      representing HMMcluster supported alleles, whereas reconstructed phylogenies of similarly sized

161      conserved regions tended to show poorly resolved clades with low bootstrap support for key central

162      nodes (Figure 2). Genetic variation in multi-allelic regions were as much as an order of magnitude

163      greater than comparably sized conserved genome portions (Figure 2). Well separated clades with

164      restricted recombination have been identified previously in CMV hypervariable regions and are

165      thought to represent an inability of homologous strands to anneal, supporting the development of

166      population structure within those loci (11).

167

168      **European and African whole genomes display geographical population structure**

169      To examine the relationship between CMV strains we first analysed whole genomes by the

170      geographical region (continent) in which they were sampled. Geographical segregation of genomes

171   is well described for other herpesviruses (38–41). Because CMV is known to be highly recombinant
172   such that accurately reconstructing phylogenetic relationships and tree topology become
173   problematic, we initially used Multi Dimensions Scaling (MDS) to analyse the genomes (Figure 3) as
174   dimensionality reduction techniques simply represent "closeness" of sequences and have been able
175   to derive representations of genetic data resembling their geographic ancestry (42, 43). The
176   resulting clustering pattern showed geographic segregation of whole genomes in the second most
177   important dimension (component 2, Supplementary Figure 3) with those sampled in Africa clustering
178   away from most sequences sampled in Europe. Two sequences sampled in Asia were located
179   amongst the European sequences (Figure 3). MDS components 1, 3 and 4 which represent large
180   variance in the data were uninformative with respect to geographical segregation of whole
181   genomes.
182   This geographical split was corroborated by a phylogenetic network tree, which allows
183   representations of shared recombination events as network splits (Supplementary Figure 4).
184   Sequences sampled in the Americas (all from the USA) were distributed throughout the plot. The
185   single sequence from Oceania (Australia) appeared to resemble sequences sampled in Europe. A
186   minority of sequences sampled in Europe clustered with sequences sampled in Africa, but not vice
187   versa (Figure 3). To examine this in more detail, we made use of a separate set of CMV genomes
188   from a cohort of CMV seropositive solid organ transplant recipients (n=11, Dataset S2) with available
189   self-reported ethnicity data, all of whom were sampled in the UK (44). Performing MDS of these
190   sequences along with the European, African, and Asian sampled CMV showed that CMV from
191   transplant recipients of self-reported African and Afro-Caribbean origin but sampled in the UK
192   clustered predominantly with strains sampled in Africa, and separately from UK patients of non-
193   African origin (Supplementary Figure 5). This suggests that sequences cluster by host ancestry rather
194   than simply by sampling location. The few European sampled strains that clustered with the bulk of
195   African samples sequences are therefore likely to be from individuals of African ancestry. Our
196   findings contrast with previously published results that failed to identify CMV population structure
197   related to geographical origin in partial genomes (45, 46) or a subset of highly heterogenous whole
198   genomes (11), although the two published Asian sequences have been noted as tightly clustering in
199   an analysis of whole genomes (10).
200
201   **Conserved and multi-allelic genomic regions show distinct patterns of phylogeography**
202   Because of potential differences in evolutionary histories between conserved and multi-allelic
203   regions we next analysed these separately for geographic structure. MDS of the concatenated
204   conserved genomes (concatenome) showed more distinct separation between viruses of African and

205 European origin than the whole genome sequences, as well as clearer segregation of the Asian CMV
206 viruses (Figure 3; Supplementary Figures 4 and 5). In contrast, multi-allelic regions when
207 concatenated appeared to show minimal geographic clustering. We quantified the differences
208 between African and European populations in conserved regions using Fixation Index (Fst) which
209 compares diversity within and between different populations. As Fst can be biased if sample sizes
210 vary between populations (47), we randomly chose 30 African-sampled and 30 European-sampled
211 sequences. The continent labels were then randomly scrambled to generate a null hypothesis Fst
212 and both steps were repeated 10,000 times to obtain true and null Fst distributions (Supplementary
213 Figure 6).
214 The results showed that the conserved regions in the CMV genome encode clear geographic
215 (continental) differences, with a mean Fst of 0.21 (a 423% increase on the mean null Fst). Fst of the
216 concatenated multi-allelic regions was relatively weaker at 0.097, which is only a 194% increase on
217 the mean null Fst. While both Fst values were significantly different to their respective null
218 distributions (both Mann Whitney test p-values < 0.0005.), geographic signal appears to be more
219 enriched (423% vs 194%) for the conserved regions of the CMV genome.
220
221 **Admixture model-based estimation of ancestry supports continental population structure**
222 To further examine the segregation of conserved genomic sequences by continent and to use a
223 model based approach to complement the visual MDS and network phylogeny, we undertook an
224 admixture analysis which attempts to infer the ancestral lineages and the contributions from each
225 that gave rise to a set of modern sequences (Figure 4) (48). Admixture analyses, like dimensionality
226 reduction, can also be skewed by large sample sizes differences between groups and by
227 heterologous sampling methodology between groups (49). To account for the former, we randomly
228 subsampled to generate more proportionate sample sizes per continent (30 European, 30 African
229 and both Asian sequences), then calculated the Cross Validation Error (CVE) for K=1 through 10. This
230 was repeated for 1000 random sample draws.
231 Two lineages (K=2), one African and one Eurasian were found to be the consensus result from the 10
232 lowest error models (Figure 4A and Supplementary figure 7A). CMV sequences sampled in Africa
233 from individuals of African ancestry, were clearly identified by the model as being part of this African
234 lineage, with no hint of admixing and vice versa for white Europeans. Few sequences sampled in
235 Europe were assigned strongly to the African lineage but in the MDS these are likely explained as
236 being sampled from ethnic African individuals in Europe (Figure 3A, 4B; Supplementary Figure 5).
237 Many sequences also appeared as admixes and correspond well with those sequences that lie
238 between the clusters of African and European in the MDS.

239     For a minority of subsampling replicates there was a random skew toward sampling more Asian-like

240     CMV sequences and in three replicates the optimal model was K=3 clusters (Supplementary figure

241     7). Although the current data is best explained by two ancestral genomes (K=2), the few  Asian

242     sequences (n=2), although clustering together will not change the admixture result for most

243     replicates (49).  Notwithstanding, there was also clustering of two Asian CMV conserved

244     concatenomes from UK patients of self-reported Asian ethnicity (Supplementary Figure 5 and

245     Dataset S2) with the two Asian-origin GenBank sequences making it likely that, with more genomes

246     three ancestral sequences (K=3) may turn out to better represent the data.

247     To overcome the complexity of sample origin not necessarily reflecting virus ancestry (i.e. human

248     migration) and to allow better delineation of population differences between African and European

249     CMV lineages, we limited further analyses to "archetypal" sequences from each continent.

250     Archetypal sequences were defined as those with >90% admixture proportion to a single ancestral

251     cluster using the conserved concatenome data, in the lowest error k=3 model. Using this definition,

252     we identified 42 strains (12 of which were sampled in Europe) of archetypal African ancestry and 129

253     of archetypal European viral ancestry, which includes the reference strain Merlin. Four CMV

254     sequences were identified as of archetypal Asian ancestry (2 of which were sampled in Europe from

255     subjects with self-reported Asian ancestry).

256

**257     Geographic signal is consistent across the conserved regions and sequences between continental**

**258     clusters have random patterns ruling out recent recombination**

259     To understand how the geographic signal is ordered across the viral genome we identified the loci

260     within the conserved concatenome with high (>0.5) Fst values, i.e. which showed the most

261     distinction between African and European strains. To provide a visual representation of the

262     continent of origin for each strain, we coloured high Fst value African nucleotides present at

263     consensus (>50%) in archetypal African strains red and high Fst value European nucleotides present

264     at consensus (>50%) in archetypal European strains, blue.  Sites which were not present at

265     consensus in ether archetypal African or European strains were coloured (black) (Supplementary

266     Figure 8). In the archetypal African or European subset, sequences were overwhelmingly

267     represented by the base most common to their continent. Moreover, where an African consensus

268     base appears in an archetypal European strain, or vice versa this appeared to be largely random.

269     Sequences not archetypal of these three continents had complex, although apparently non-random

270     patterns of admixture.

271

**272     Multi-allelic genomic regions represent a mixture of geographic relationships**

273    We next examined the apparent but weaker geographical segregation of the 74 multi-allelic regions
274    between African and European populations. As the multi-allelic regions each contain variable
275    numbers of alleles and are of different lengths, we considered each of the 74 regions separately.  We
276    performed chi squared tests with allele distributions to evaluate whether multi-allelic regions
277    segregated with archetypal African or archetypal European strains. The results revealed strong
278    geographic distribution of alleles for 32 of the 74 regions (Table S1, we calculated false discovery
279    rate, FDR, using the Benjamini-Hochberg procedure and  we set the FDR threshold at 0.05) (50),
280    while five showed more moderate geographic differences  (0.05 < FDR < 0.3; Table S1) with the rest
281    showing no evidence of  geographical segregation.  Colouring African dominant alleles red and
282    European-dominant alleles blue (as defined above) for the 32 multi-allelic regions that segregated
283    geographically we observe a similar pattern to that seen for the conserved regions, with African
284    strains largely red and European strains largely blue (Supplementary Figure 8, C and D).

285

286    **Host Immune mediated selection is not obviously responsible for driving geographic segregation**
287    **of African and European strains**
288    To assess whether the geographic population structure observed in conserved regions of the CMV
289    genome was due to differences in the selective pressures exerted by different host populations, we
290    looked at the 440 most geographically informative sites (Fst > 0.5) within the archetypal subsets of
291    European and African strains. We asked how many of these sites resulted in a different continental
292    consensus amino acid between African and European sequences and, if so, whether the change
293    occurred within known B and T cell epitopes, as recorded in the Immune Epitope DataBase (IEDB)
294    (51). Only 15% (64) of the 440 sites encoded nonsynonymous changes. Of these, 16% (10 of 64) lay
295    within known epitopes, compared with 14% (52 of 376) of synonymous sites, providing little
296    evidence that geographic population structure in CMV is driven by continentally unique host
297    immune pressure.

298

299    **Alleles in immunomodulatory genes tend to maintain similar diversity across continents**
300    We next tested whether certain gene functions were over-represented in three classes of genomic
301    regions: conserved, geographically segregating multi-allelic, and non-geographically segregating
302    multi-allelic regions. To do this we identified the genes lying within each region class and annotated
303    their function, using three predefined functional groups from a published gene-ontology [Latency,
304    Tropism and Immunomodulation] (52). We then determined whether key functional groups were
305    significantly over/under-represented by region class (Supplementary Figure 9). We found the 74
306    multi-allelic regions together were significantly enriched for genes encoding immunomodulatory

307 functions with the proportion of genes with immunomodulatory function of 0.371 (37.1%) where for
308 all genes this proportion is 0.206 (FDR= 0.0096). This enrichment was clearer still if only the 42 multi-
309 allelic regions with no evidence of geographical segregation were considered (39.1%, FDR=0.014).
310 No other comparison was significant.
311

312 **Geographic and multi-allelic genetic differences may impact biological function**
313 One interesting application of studying multi-allelic regions is the possible association between
314 different alleles, and function. To illustrate this, we investigated glycoprotein B (gB, UL55) as an
315 example. The reference strain Towne has been extensively used to develop vaccines and to study
316 CMV immunogenicity and function (53) and its gB protein variant has been the basis for many
317 vaccine candidates (54). From the admixture analysis we identified Towne as being 79% African,
318 while two other common reference strains AD169 and Merlin were 75% and 99% European
319 respectively. gB (UL55) sequence is composed of 3 separated multi-allelic regions (22, 23 and 24 in
320 Table S1) which are linked (Monte Carlo chi square p<0.001) in 12 possible combinations or
321 haplotypes. In addition to showing geographically informative genetic differences from European
322 strains in its conserved regions, Towne gB also differs in the sequence of its multi-allelic regions
323 sharing the same 22,23,24 (UL55) haplotype as only 4% of the 259 CMV genomes sequenced here
324 (Supplementary Figure 10). By contrast, the Merlin multi-allelic region 22,23,24 (UL55) haplotype is
325 shared with 32% and AD169 with 20% of the 259 genomes. Genetic differences in gB (UL55) have
326 been mooted to underlie observed differences in cross-strain neutralisation by antibodies raised
327 against one strain (55, 56). In general, immunotherapies and drug development targeting CMV that
328 rely on alleles that differ across geographic isolates, may now require further investigation as to
329 whether treatment effect will be advantageous to only certain human populations.
330

331 **Discussion**
332 We have characterised the variability present in whole CMV genome sequences, including several
333 known to be of African origin identifying  both known hypervariable regions and over 40 which have
334 not previously identified (10, 11, 13, 16, 57).  Altogether we describe 74 hypervariable regions
335 comprising 14% of the genome, all of which are multi-allelic. The remaining 86% of the genome is
336 monoallelic and at least ten times less variable than the multi-allelic regions. While previous reports
337 have identified around 30 multi-allelic hypervariable regions, they identified them by the gene in
338 which they were located, despite in many cases, much of the gene concerned not being
339 hypervariable and thus subject to different constraints on recombination and diversity. In contrast,
340 our data precisely delineates the nucleotide coordinates of variable and conserved regions, using an

341 unbiassed and consistent assignment model. This has allowed the identification of 17 multi-allelic

342 regions that cross gene boundaries and 5 are entirely in non-coding or in otherwise unassigned. For

343 each of the 74 regions, our statistical approach also returns the number of alleles that most

344 parsimoniously fit the data, leading to between 2 and 8 alleles per region (Table S1). In some cases,

345 for example  UL146 and UL144 , the number of alleles we identify differs from previous numbers

346 reported  (58, 59). However, previous estimates were determined largely by visual inspection of

347 whole-gene phylogenetic trees, a  process which resulted in differences in reported allele numbers

348 not only from us but between other authors (for example for UL55), with, on occasion some alleles

349 not being reliably distinguishable  (16, 60–62). The objective mathematical approach we have

350 adopted provides clear and reproducible multi-allelic region boundaries and allele numbers,

351 properties which will have advantages for standardising genotyping nomenclature.

352 In contrast to previous reports (10, 11) we observe clear geographical segregation of CMV genomes

353 with evidence for African, European, mixed and possibly Asian genotypes. From an evolutionary

354 viewpoint, geographical segregation of genomes is well described for other human herpesviruses

355 including herpes simplex virus (HSV-1), varicella zoster virus (VZV) and Epstein Barr virus (EBV) (38–

356 41, 63) and is therefore not surprising for CMV. Our data suggest that most geographically

357 informative SNPs in CMV are in the monoallelic  genomic portions and like the rest of the genome

358 are under purifying selection (10).   Unlike EBV, in which host immune selection drives local

359 adaptation of virus to different human host populations, shaping the pattern of genetic diversity

360 (64), we find no evidence that selection within immunogenic regions of the genome are a dominant

361 driver of the observed genetic geographical differences. Instead, we postulate that genetic drift and

362 bottleneck events such as founder viruses are plausible explanations for the population structure

363 observed in CMV. If this is the case, there remains some difficulty in establishing the direction and

364 date of split for European and African CMV populations, due to the low association between

365 sampling date and distance from phylogenetic tree root and a mutation rate that has proven difficult

366 to determine for double-stranded DNA viruses (65).  The difficulty in determining mutations rates is

367 likely to be, in part, a result of CMV's longstanding free recombination within these geographically

368 isolated pockets. As 13 multi-allelic regions were found to contain alleles unique to Europe while the

369 opposite was not seen for African viruses, this could be taken as evidence supporting a European

370 origin of CMV, where Africa has restricted diversity. However, this is likely more simply explained as

371 an artefact of the differences in size and sampling heterogeneity of our available genomes.

372 The finding of African-clustering CMV strains in patients self-reporting as being of Afro-Caribbean

373 ethnicity, many of whom presumably have not lived in Africa can potentially be explained by early

374 acquisition of CMV from family members and by assortative mating of racial groups. A similar effect

375  has been observed for VZV in subjects of Afro-Caribbean origin growing up in the UK (41). In African

376  countries most children are CMV positive by their first birthday (33, 66, 67). Early Infection, most

377  likely acquired from maternal or sibling transmission (68) may explain why subjects of African origin

378  living in Europe test positive for strains that cluster with known African strains. This would date the

379  split of African and European CMV strains to at least 500 years ago, the time at which the first

380  African slaves were transported to Europe (69). Although the actual separation is likely to be much

381  older given even the highest estimates of CMV's mutation rate (70), or borrowing from rates

382  presumably more accurately estimated for a similar virus, HSV1, for which there are ancient

383  genomes available (71). With greater mixing of populations and the pervasive genome-wide

384  recombination that occurs in CMV, we see evidence for increasing numbers of hybrid strains

385  including some of the reference strains for example Towne and AD159 (Figure 4). This serves to

386  further muddy insights into the phylogeographic origins of currently circulating strains of CMV and

387  elucidating these undoubtedly requires additional and more granular worldwide sampling, as well as

388  the inclusion of ancient CMV genetic material if these can be found.

389  While in 32 of the multi-allelic regions, the alleles, like the conserved regions, are different or

390  differently distributed between geographical regions, most notably countries with predominantly

391  European or African populations, from which most strains originate, the majority (42) show no

392  evidence of geographical distribution, but instead appear to maintain the full allele palette in these

393  genome portions across both continents and at similar frequencies. This effect has previously been

394  observed in RNA viruses (72, 73). Initial studies into CMV virion envelope complexes, linking variants

395  to function, have reported that allele differences can modulate virus cell tropism (74, 75). Similar

396  examples from hepatitis C virus and human immunodeficiency virus (HIV) have shown genotypic

397  differences to affect viral compartmentalisation (76), while certain hepatitis B virus genotypes

398  appear to be associated with chronic infection (77). From our analyses, the non-geographically

399  segregated multi-allelic regions were significantly enriched for genes encoding immunomodulatory

400  functions. For example, region 6 which shows no geographical segregation encodes a portion of the

401  non-recombinant haplotype RL11D block, RL11, RL12, RL13, UL1, UL2 and UL4, all of which are

402  proven or predicted to be virion membrane glycoproteins (52). In addition, RL11D (region 6)

403  variability has been suggested to be critical for the adaptation of CMV to different primate species

404  (11). Of interest, many of the variable multi-allelic regions correspond to regions previously

405  identified as being in local linkage disequilibrium and thus not affected by the pervasive

406  recombination occurring throughout the more conserved regions of the genome (11). Taken

407  together, our data strongly support the likelihood that CMV genome is the result of two distinct

408  evolutionary forces, genetic drift occurring in segregated viral populations and so called frequency-

409     dependent balancing selection (78), a form of adaptation that maintains pre-existing diversity in the

410     face of genetic drift.

411     This granularity of CMV genome analysis allows deeper insights into how genome might be related

412     to function. Following Towne gB + MF59 adjuvant vaccination, antibody titres to the antigenic AD2

413     region have been shown to correlate with better protection against post renal transplant CMV

414     viremia (54–56). Baraniak and colleagues also showed that only ~50% of individuals vaccinated with

415     Towne gB had detectable AD2 antibody response against gB peptides in an assay derived from the

416     AD169 laboratory strain. The AD169 gB AD2 allele (region 24, UL55) differs from that found in Towne

417     gB (Supplementary figure 8). Since AD2 antibodies are not broadly reactive (54), there is some

418     question about what the Towne gB vaccine antibodies are recognising within the AD169 gB AD2

419     peptides. Multi-allelic region 24 (gB/UL55 AD2) also segregates differently between African (Towne-

420     like) and European-like (AD169 and Merlin) viral populations. This together with the finding that the

421     Towne gB conserved region carries predominantly African-segregating SNPS raises the possibility

422     that a vaccine based on Towne gB sequence might not confer cross-protective immunity against

423     European strains.

424     Notably, the CMV pentamer complex, which is being developed as part of an alternative multi-

425     antigenic vaccine using the Merlin strain, contains no multi-allelic regions and may therefore be

426     more tractable than gB (79). However, geographically related differences are still present and

427     potentially need evaluation. For example, Q35K which segregates with African strains and L40P

428     which segregates with European strains are both present within a known B cell neutralising epitope

429     in the pentameric complex UL130 protein (IEDB ID: 142031).

430     These analyses are subject to limitations, the clearest of which is the potential biases related to the

431     available sequences and the samples from which they were derived. First, the representation of

432     genomes was heavily skewed towards Europe. Second, the European samples were typically

433     collected from unrelated patients for clinical purposes, whereas African samples were obtained from

434     study participants in southern and eastern Africa (Zambia, Uganda, and Kenya), many of whom were

435     HIV co-infected. African-clustering sequences that were sampled in Europe likely reflect divergent

436     sampling of unrelated individuals and may to some extent mitigate the bias of African strains.

437     However, we had only two Asian strains collected in the early 2000s and no strains from south

438     America or many other parts of the world. Until these gaps are closed our conclusions must remain

439     incomplete. Seventeen African genomes were reconstructed from samples containing mixed CMV

440     infections where generating consensus genomes has typically been challenging. However, HaRold,

441     the program we used to reconstruct haplotypes (31) performs with high accuracy in validation

442     exercises using simulated and real mixtures of CMV genomes containing known sequences. Lack of

443 homology between some alleles within a multi-allelic region could be a limitation on constructing

444 alignments for HMM.  However, even for the most divergent alleles, the bordering 5' and 3'

445 sequences are identical, a factor that mitigates this potential constraint. Finally, when considering B

446 and T cell CMV epitopes, we are limited by epitopes included in the IEDB database which are largely

447 generated for European strains.  However, since most of the viruses analysed here were European,

448 the conclusions that most nonsynonymous differences from African strains do not lie within

449 epitopes is likely to be true.

450

451 **Conclusion**

452 Our findings provide several new insights into the genomic landscape of CMV.  First, we identify and

453 precisely delineate 74 discrete variable regions which consist of multiple alleles, showing that the

454 rest of the genome (86%) is monoallelic. We identify for the first time, that CMV genomic evolution

455 is shaped by two distinct processes: likely genetic drift occurring within geographically distinct

456 populations and balancing selection which counteracts genetic drift to maintain similar diversity in

457 variable multi-allelic regions irrespective of geographical location.  We identify that variable regions

458 under balancing selection are enriched for key CMV properties, highlighting that better

459 characterisation of diversity in these regions is likely to be important for understanding CMV biology

460 and control.  Our results provide a genomic roadmap to enable studies of how variation across the

461 CMV genome interacts to cause clinical disease. At the same time, the data raise questions about

462 how geographical differences arose and the direction of spread from one region to another.  The

463 answers to these questions will require further sampling of geographically diverse whole CMV

464 genomes, CMV sequence data from ancient samples, or both. Finally, the data highlight that the

465 geographical and allelic differences between proteins being trialled as potential vaccines needs to be

466 considered when designing vaccines. Our findings raise the possibility that vaccines based on strain-

467 specific gB or other viral antigens may fail to induce sufficient cross-protection globally against

468 circulating variants.

469

475

476 **Conflicts of interest**

479

480  **Code availability**

481  The HMMcluster program is freely available under the MIT license, at https://github.com/ucl-

482  pathgenomics/hmmcluster

483

484  **Data availability**

485  The accession identifiers (GenBank or SRA accessions) for the 259 genomes included in the main

486  analysis and the genomes of those patients with self-reported ancestry are available in Dataset S1

487  and Dataset S2. The full alignment and full resolution phylogenetic trees are also available at

488  https://github.com/ucl-pathgenomics/HCMV_resources_public.

489

490  **Tables**

491

492 **Table 1. Recombination in multi-allelic regions**.

493

| Region | Genes | Number of breakpoints | NC_006273.2 coordinates of breakpoints (bp) |
|---|---|---|---|
| **2** | **RL5A; RL6** | **2** | **5907, 6290** |
| **6** | **RL11; RL12; RL13; UL1; UL2; UL4** | **2** | **12606, 13305** |
| **8** | **UL10; UL11; UL6; UL7; UL8; UL9** | **6** | **15751, 15860, 16390, 17553, 18419, 18672** |
| 10 | UL20; UL21A | 2 | 26468, 26732 |
| 28 | UL73; UL74 | 2 | 107640, 107848 |
| 42 | UL116 | 1 | 166458 |
| 43 | UL119; UL120; UL121 | 6 | 169045, 169283, 169475, 169525, 169782, 170046 |
| 50 | UL144 | 2 | 182464, 182551 |
| 51 | UL150A | 1 | 183170 |
| 72 | TRS1 | 4 | 231683, 231783, 232227, 232415 |

494 Each multi-allelic region was assessed for evidence of recombination. Firstly, each region was examined using a

495 set of seven recombination methods implemented in RDP5. We then visually investigated the phylogenies

496 either side of predicted breakpoints of those multi-allelic regions with evidence of recombination (a region was

497 considered to have evidence of recombination if at least 5 methods in RDP5 were significant). We removed

498 breakpoints that could be explained by sub-clade structure. We underline and highlight in bold those regions

499 where most sequences in the multiple-sequence alignment showed the recombination breakpoints.

500

501 **Materials and methods**

502

503 **Data retrieval**

504 A python script using Biopython (80), specifically the Entrez module, was used to access the SRA and

505 NCBI nucleotide databases for sequence information, and extract country and continent assignment

506 for sequences.

507

508 **Sequence assembly**

509 SRA sequences for Zambian CMV genomes were downloaded using the SRA toolkit and assembled

510 using an in house de novo assembly pipeline, which involves contig generation, optimal reference

511 identification, scaffolding on to the reference sequence, and subsequent iterative mapping of NGS

512 reads on the genome scaffold. These were then subject to haplotype reconstruction and relevant

513 consensus sequences determined.

514 Ugandan sequences and historical clinical sequences of known ethnicity were also de novo

515 assembled. Kenyan sequence data were assembled to a reference sequence using an in-house

516 pipeline using the strain. Sequence positions with less than 10 read depth were labelled as n.

517

518 **Haplotype reconstruction**

519 Possible mixed infections were investigated with HaROLD, a tool for reconstructing haplotypes using

520 co-varying variant frequencies in a probabilistic framework (31). HaROLD takes the bam files

521 obtained from the assembly step (as explained above) and then reconstructs the optimal number of

522 haplotypes for each sample. Haplotypes' sequences are then checked by reconstructing

523 phylogenetic trees and are considered distinct if they have >2000 bp differences.  We reconstructed

524 a total of 17 haplotypes from 8 samples (1-3 haplotypes per sample). In line with the approach taken

525 in de novo assembly, we ignored haplotypes with an average read depth of less than 10 bases

526 (haplotype frequency * mean read depth).

527

528 **Multiple sequence alignment**

529 Multiple sequence alignments were obtained using MAFFT v7 (81), particularly variable sections

530 were re-aligned using MUSCLE (82) and finished manually. Sequence alignments were viewed in the

531 lightweight alignment viewer AliView (83). Alignments relative to a reference strain were only used

532 to generate the heterozygosity per reference position calculation, these were generated using

533 MAFFT with the " --add – keeplength" options, which allowed SNPs to be called based on differences

534 to the reference Merlin (Refseq accession: NC_006273.2).

535

536 **Measures of sequence diversity**

537 Heterozygosity was generated using an in-house R function, using the following calculation.

538 h is heterozygosity for a given polymorphic site with I alleles, such that the sum of all allele

539 frequencies p equals 1. N is the number of sequences in the sample. Summing over all segregating

540 sites S in an alignment, we get sum of site heterozygosity $\pi$.

$$h = \frac{n}{n-1} \left(1 - \sum_{j=1}^{S} p_i^2\right) \pi = \sum_{j=1}^{S} h_j$$

541

542 **Multi-Dimensional Scaling**

543    Pairwise distances were calculated using the dist.dna() with the nucleotide-nucleotide substitution

544    matrix "TN93" (84) and with pairwise deletion by way of the R package Ape v.5.4 (85). Multi-

545    dimensional scaling much like principal component analysis (PCA) is a method to attempt to simplify

546    complex data into a more interpretable format, by reducing dimensionality of data whilst retaining

547    most of the variation. In a genomics context we can use this on pairwise distance matrices, where

548    each dimension is a sequence with data points of n-1 sequences pairwise distance. This allows us to

549    observer patterns of population structure as "clusters". MDS was implemented using the cmdscale()

550    function with pairwise deletion in R (86).

551

552    **Phylogenetic Reconstruction**

553    Phylogenetic relationships of multi-allelic and example conserved regions in Figure 2 were

554    constructed from nucleic acid sequences in IQ-TREE (87). Using a Maximum Likelihood GTR

555    substitution model with a discrete Gamma heterogeneity model (88) and 1000 rounds of

556    bootstrapping . We attempted to root the CMV homologous genome to the most suitable ex-CMV

557    taxa Chimpanzee CMV (accession AF480884), however this outgroup was too far removed with

558    distance >5 so unrooted trees were preferred. Trees were visualised using Figtree (89).

559    For the whole genome and concatenome phylogenetic analysis where CMV is known to recombine

560    freely, a Neighbor-net split phylogenetic network analysis was undertaken using Splitstree version

561    4.1.5 (90). Non-default options chosen were HKY85 distance matrix, with equal site rate variation.

562    Both terminal repeat regions were trimmed from alignments (although they had negligible impact)

563    before analysis.

564

565    **"$F_{st}$" F statistics**

566    For calculating a $F_{st}$ like statistic from sequence data, we can use the sum of site heterozygosity's

567    across a locus to produce $\gamma_{st}$. Where $\pi_T$ is calculated as above using all samples in an alignment, $\pi_S$ is

568    an average of the same calculation for each sub population separately. (91).

$$\gamma_{ST} = \frac{\pi_T - \pi_S}{\pi_T}$$

569    Only sites with greater than 5% minor allele frequency were considered. To account for uneven

570    African and European populations, either when defined by sampling location or when considering

571    those sequences that are archetypally (90%> in admixture analyses) African or European, we

572    repeatedly subsampled the European population to be equivalent to the number of African

573    sequences (1000*) and took the mean of the site Fst's.

574    When we multiply bootstrap sampled the 30 African and 30 European sequences, the mean number

575    of pairwise differences for sequences within each population were determined, as well as the mean

576    number of pairwise differences across all sequences. This can be used to estimate Fst in an efficient

577    manner for multiple bootstraps (92).

578

579    **Chi-squared analysis of allele proportions**

580    For each region the allele assignments from HMMcluster were grouped by origin into African and

581    European allele frequencies as observed in the admixture archetypal strains, which we tested for

582    significant differences using a chi squared test of independence with the base R function chisq.test

583    over allele frequencies. From the European data we generated expected allele frequencies which

584    were compared against observed allele frequencies from Africa. A Benjamin-Hochberg adjusted

585    False Discovery Rate (FDR) < 0.05 we determined as inferring significant distribution deviation where

586    we assigned the multi-allelic region as "geographic", otherwise they were labelled "pervasive".

587

588    **HMMcluster – Sequence clustering by Optimal Hidden Markov Models**

589    We implemented a maximum likelihood allele assignment model in java on a Hidden Markov Models

590    statistical framework. Briefly, this approach considers the genomic alignment as a set of contiguous

591    blocks, within each block the model instantiates by perfectly representing each sequence as its own

592    HMM, this results in the highest Likelihood but with an excessive number of parameters. Then the

593    model considers the optimal way to combine HMM's to keep the highest likelihood i.e. 259 models

594    to 258, and continues to iterate with a greedy stepwise algorithm until only a single HMM is

595    reached. A single HMM most poorly represents each sequence, such that the likelihood is lowest,

596    but it uses the minimal number of parameters. To balance this likelihood and parameter problem, a

597    typical approach is to appeal to the Akaike Information Criterion (AIC), and we use this here to

598    identify the most parsimonious representation of a given genomic segment.

599    Consider that we have a set of sequences $\{x_{1j}, x_{2j} \dots x_{nj}\}$ where $x_{ij}$ is the base found at position $i$ in

600    sequence $j$, where there are $N$ positions in sequence $(1 \leq i \leq N)$ and $M$ sequences $(1 \leq j \leq M)$.

601    Each Hidden Markov Model, not considering insertion and deletion states (of which we ignore

602    insertion states) is defined as a series of match states which are represented by the probability of

603    the emissions from that state. That is the Hidden Markov Model is defined by where $p_i(x)$ is the

604    probability that match state $i$ emits base $x$, and $\sum_k p_j(x_k) = 1$.

$$\left\{ \begin{pmatrix} p_1(x_1) \\ p_1(x_2) \\ \vdots \end{pmatrix}, \begin{pmatrix} p_2(x_1) \\ p_2(x_2) \\ \vdots \end{pmatrix}, \begin{pmatrix} p_3(x_1) \\ p_3(x_2) \\ \vdots \end{pmatrix} \dots \right\}$$

605    In this case, the probability that sequence $j$ would arise from this hidden Markov model is equal to

606    $\prod_i p_i(x_{ij})$ or the log likelihood is given by $\sum_i \log p_i(x_{ij})$. The total log likelihood for the set of $M$

607    sequences is then equal to $\sum_j \sum_i \log p_i(x_{ij}) = \sum_i \sum_j \log p_i(x_{ij})$.

608  If we consider a given location $i$ and imagine that at this site $m_{i1}$ of the sequences have base $x_1$, $m_{i2}$

609  of the sequences have base $x_2$, etc, with $\sum_k m_{ik} = M$, then we can sum over identities of bases

610  rather than sum over sequences, and the log likelihood becomes $\sum_i \sum_k m_{ik} \log p_i(x_k)$. It turns out,

611  not surprisingly, that the best (i.e. maximum likelihood) values for $p_i(x_k)$ is equal to the fraction of

612  the sequences that have base $x_k$ at that position, that is, $\hat{p}_i(x_k) = \frac{m_{ik}}{M}$. Substituting this in yields the

613  highest likelihood of the set of sequences is given by $\sum_i \sum_k m_{ik} \log \frac{m_{ik}}{M}$.

614  We implemented the model to look for evidence of population structure initially within 200bp

615  genome slices. Identified loci were refined by Maximum likelihood, and any overlapping regions

616  concatenated, and again start / stop positions refined by maximum likelihood.

617

618  **Recombination analysis**

619  Genome sequences were examined for evidence of systematic recombination events using the

620  Recombination Detection Program (RDP) version RDP5.5 with the maximum likelihood tree option

621  (93). The RDP software includes a suite of recombination-detecting algorithms where we used

622  seven, namely phylogenetic (RDP, BOOTSCAN and SISCAN) and substitution (GENECONV, MAXCHI,

623  CHIMAERA, and 3-SEQ) methods to generate evidence of recombination. Using a Bonferroni

624  corrected P-value cut-off of $\leqslant 0.05$ significant scores with 5 or more of the seven algorithms, found

625  in a group of 4 or more non-haplotype sequences were considered significant and the phylogeny

626  either side examined to determine if it was a true significant.

627

628  **Population structure**

629  Population structure was analysed in an unsupervised fashion with Admixture 1.3.0 (48). Alignments

630  were converted to VCF format using snp-sites (94), and sites with minimum allele frequency < 5%

631  were trimmed. Sequences were randomly subsampled to generate more proportionate sample sizes

632  per continent (30 European, 30 African and both (2) Asian) 1000 time. For each of the 1000 sample

633  draws admixture was run for a k ranging from 1 to 10 with 20-fold cross validation. As recommended

634  in the admixture manual we thinned the markers according to the observed sample correlation

635  coefficients using the plink argument "--indep-pairwise 50 10 0.1". Analyses were visualised in R.

636

637  **Identifying Epitopes**

638  Known CMV B and T cell epitopes were downloaded from IEDB (51) then mapped to Merlin

639  reference strain genomic co-ordinates by tblastn (95). Predicted epitopes were ignored. cmvdrg (96)

640  was used to identify which variants are synonymous or nonsynonymous when translated. Sites with

641  less than 10% variants in either African or European populations were ignored.  As sites with low

642 variability can still exhibit high Fst values, we limited the analysis to sites where the consensus base

643 was different between the archetypal African and European sequences, this removed 11% of sites.

644 This allowed variant sites to be analysed together from geographical, immune, and protein effect

645 frames of reference.

646

647 **References**

648 1. M. Zuhair, *et al.*, Estimation of the worldwide seroprevalence of cytomegalovirus: A
649 systematic review and meta-analysis. *Rev. Med. Virol.* **29**, e2034 (2019).

650 2. J. I. Cohen, Herpesvirus latency. *J. Clin. Invest.* **130**, 3361 (2020).

651 3. A. J. Davison, *et al.*, The Order Herpesvirales. *Arch. Virol.* **154**, 171–177 (2009).

652 4. S.-Y. Cho, D.-G. Lee, H.-J. Kim, Cytomegalovirus Infections after Hematopoietic Stem Cell
653 Transplantation: Current Status and Future Immunotherapy. *Int. J. Mol. Sci.* **20** (2019).

654 5. C. Marsico, D. W. Kimberlin, Congenital Cytomegalovirus infection: advances and challenges
655 in diagnosis, prevention and treatment. *Ital. J. Pediatr.* **43**, 38 (2017).

656 6. Institute of Medicine (US) Committee to Study Priorities for Vaccine Development, *Vaccines*
657 *for the 21st Century: A Tool for Decisionmaking*, K. R. Stratton, J. S. Durch, R. S. Lawrence, Eds.
658 (National Academies Press (US), 2000) (June 27, 2022).

659 7. A. Dolan, *et al.*, Genetic content of wild-type human cytomegalovirus. *J. Gen. Virol.* **85**, 1301–
660 1312 (2004).

661 8. J. Cudini, *et al.*, Human cytomegalovirus haplotype reconstruction reveals high diversity due
662 to superinfection and evidence of within-host recombination. *Proc. Natl. Acad. Sci.* **116**,
663 5693–5698 (2019).

664 9. N. Renzette, B. Bhattacharjee, J. D. Jensen, L. Gibson, T. F. Kowalik, Extensive Genome-Wide
665 Variability of Human Cytomegalovirus in Congenitally Infected Infants. *PLoS Pathog.* **7** (2011).

666 10. S. Sijmons, *et al.*, High-Throughput Analysis of Human Cytomegalovirus Genome Diversity
667 Highlights the Widespread Occurrence of Gene-Disrupting Mutations and Pervasive
668 Recombination. *J. Virol.* **89**, 7673–7695 (2015).

669 11. F. Lassalle, *et al.*, Islands of linkage in an ocean of pervasive recombination reveals two-speed
670 evolution of human cytomegalovirus genomes. *Virus Evol.* **2** (2016).

671 12. S. Chou, Reactivation and Recombination of Multiple Cytomegalovirus Strains from Individual
672 Organ Donors. *J. Infect. Dis.* **160**, 11–15 (1989).

673 13. E. Puchhammer-Stöckl, I. Görzer, Human cytomegalovirus: an enormous variety of strains and
674 their possible clinical significance in the human host. *Future Virol.* **6**, 259–271 (2011).

675 14. E. Puchhammer-Stöckl, I. Görzer, Cytomegalovirus and Epstein-Barr virus subtypes—The
676 search for clinical significance. *J. Clin. Virol.* **36**, 239–248 (2006).

677  15.  U. Meyer-König, K. Ebert, B. Schrage, S. Pollak, F. T. Hufert, Simultaneous infection of healthy
678      people with multiple human cytomegalovirus strains. *The Lancet* **352**, 1280–1281 (1998).

679  16.  N. M. Suárez, *et al.*, Human Cytomegalovirus Genomes Sequenced Directly From Clinical
680      Material: Variation, Multiple-Strain Infection, Recombination, and Gene Loss. *J. Infect. Dis.*
681      **220**, 781–791 (2019).

682  17.  N. M. Suárez, *et al.*, Multiple-Strain Infections of Human Cytomegalovirus With High Genomic
683      Diversity Are Common in Breast Milk From Human Immunodeficiency Virus–Infected Women
684      in Zambia. *J. Infect. Dis.* **220**, 792–801 (2019).

685  18.  D. H. Shepp, *et al.*, Cytomegalovirus Glycoprotein B Groups Associated with Retinitis in AIDS.
686      *J. Infect. Dis.* **174**, 184–187 (1996).

687  19.  V. C. Emery, *et al.*, Differential decay kinetics of human cytomegalovirus glycoprotein B
688      genotypes following antiviral chemotherapy. *J. Clin. Virol.* **54**, 56–60 (2012).

689  20.  E. Paradowska, *et al.*, Distribution of cytomegalovirus gN variants and associated clinical
690      sequelae in infants. *J. Clin. Virol.* **58**, 271–275 (2013).

691  21.  S. Pignatelli, *et al.*, Cytomegalovirus gN Genotypes Distribution among Congenitally Infected
692      Newborns and Their Relationship with Symptoms at Birth and Sequelae. *Clin. Infect. Dis.* **51**,
693      33–41 (2010).

694  22.  E. Paradowska, *et al.*, Cytomegalovirus alpha-chemokine genotypes are associated with
695      clinical manifestations in children with congenital or postnatal infections. *Virology* **462–463**,
696      207–217 (2014).

697  23.  R. He, *et al.*, Sequence Variability of Human Cytomegalovirus UL146 and UL147 Genes in Low-
698      Passage Clinical Isolates. *Intervirology* **49**, 215–223 (2006).

699  24.  S. K. Pati, *et al.*, Genotypic Diversity and Mixed Infection in Newborn Disease and Hearing
700      Loss in Congenital Cytomegalovirus Infection. *Pediatr. Infect. Dis. J.* **32**, 1050–1054 (2013).

701  25.  A. Mozzi, *et al.*, Past and ongoing adaptation of human cytomegalovirus to its host. *PLOS
702      Pathog.* **16**, e1008476 (2020).

703  26.  A. J. Davison, *et al.*, The human cytomegalovirus genome revisited: comparison with the
704      chimpanzee cytomegalovirus genome. *J. Gen. Virol.* **84**, 1053–1053 (2003).

705  27.  S. Sijmons, *et al.*, A Method Enabling High-Throughput Sequencing of Human
706      Cytomegalovirus Complete Genomes from Clinical Isolates. *PLOS ONE* **9**, e95501 (2014).

707  28.  D. J. Dargan, *et al.*, Sequential mutations associated with adaptation of human
708      cytomegalovirus to growth in cell culture. *J. Gen. Virol.* **91**, 1535–1546 (2010).

709  29.  G. S. Jung, *et al.*, Full genome sequencing and analysis of human cytomegalovirus strain JHC
710      isolated from a Korean patient. *Virus Res.* **156**, 113–120 (2011).

711  30. ,   Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*
712      **46**, D8–D13 (2018).

713  31.  C. Venturini, *et al.*, Haplotype assignment of longitudinal viral deep-sequencing data using co-
714       variation of variant frequencies. *Virus Evol.*, veac093 (2022).

715  32.  J. Pang, *et al.*, Mixed cytomegalovirus genotypes in HIV-positive mothers show
716       compartmentalization and distinct patterns of transmission to infants. *eLife* **9**, e63199 (2020).

717  33.  S. Gantt, *et al.*, Prospective Characterization of the Risk Factors for Transmission and
718       Symptoms of Primary Human Herpesvirus Infections Among Ugandan Infants. *J. Infect. Dis.*
719       **214**, 36–44 (2016).

720  34.  D. P. Depledge, *et al.*, Deep Sequencing of Viral Genomes Provides Insight into the Evolution
721       and Pathogenesis of Varicella Zoster Virus and Its Vaccine in Humans. *Mol. Biol. Evol.* **31**, 397–
722       409 (2014).

723  35.  N. M. Suárez, *et al.*, Whole-Genome Approach to Assessing Human Cytomegalovirus
724       Dynamics in Transplant Patients Undergoing Antiviral Therapy. *Front. Cell. Infect. Microbiol.*
725       **10** (2020).

726  36.  C. Mattick, *et al.*, Linkage of human cytomegalovirus glycoprotein gO variant groups identified
727       from worldwide clinical isolates with gN genotypes, implications for disease associations and
728       evidence for N-terminal sites of positive selection. *Virology* **318**, 582–597 (2004).

729  37.  P. Smyth, "Clustering Sequences with Hidden Markov Models" in *Advances in Neural
730       Information Processing Systems 9*, M. C. Mozer, M. I. Jordan, T. Petsche, Eds. (MIT Press,
731       1997), pp. 648–654.

732  38.  S. Correia, *et al.*, Sequence Variation of Epstein-Barr Virus: Viral Types, Geography, Codon
733       Usage, and Diseases. *J. Virol.* **92**, e01132-18 (2018).

734  39.  C. Grose, Pangaea and the Out-of-Africa Model of Varicella-Zoster Virus Evolution and
735       Phylogeography. *J. Virol.* **86**, 9558–9565 (2012).

736  40.  M. L. Szpara, *et al.*, Evolution and Diversity in Human Herpes Simplex Virus Genomes. *J. Virol.*
737       **88**, 1209–1227 (2014).

738  41.  M. Quinlivan, *et al.*, The Molecular Epidemiology of Varicella-Zoster Virus: Evidence for
739       Geographic Segregation. *J. Infect. Dis.* **186**, 888–894 (2002).

740  42.  J. Novembre, *et al.*, Genes mirror geography within Europe. *Nature* **456**, 98–101 (2008).

741  43.  D. Gurdasani, *et al.*, The African Genome Variation Project shapes medical genetics in Africa.
742       *Nature* **517**, 327–332 (2015).

743  44.  S. F. Atabani, *et al.*, Cytomegalovirus Replication Kinetics in Solid Organ Transplant Recipients
744       Managed by Preemptive Therapy. *Am. J. Transplant.* **12**, 2457–2464 (2012).

745  45.  A. J. Bradley, *et al.*, Genotypic analysis of two hypervariable human cytomegalovirus genes. *J.
746       Med. Virol.* **80**, 1615–1623 (2008).

747  46.  S. Pignatelli, *et al.*, Human cytomegalovirus glycoprotein N (gpUL73-gN) genomic variants:
748       identification of a novel subgroup, geographical distribution and evidence of positive
749       selective pressure. *J. Gen. Virol.* **84**, 647–655 (2003).

750  47.  S. Shringarpure, E. P. Xing, Effects of Sample Selection Bias on the Accuracy of Population
751       Structure and Ancestry Inference. *G3 GenesGenomesGenetics* **4**, 901–911 (2014).

752  48.  D. H. Alexander, J. Novembre, K. Lange, Fast model-based estimation of ancestry in unrelated
753       individuals. *Genome Res.* **19**, 1655–1664 (2009).

754  49.  D. J. Lawson, L. van Dorp, D. Falush, A tutorial on how not to over-interpret STRUCTURE and
755       ADMIXTURE bar plots. *Nat. Commun.* **9**, 3258 (2018).

756  50.  Y. Benjamini, Y. Hochberg, Controlling the False Discovery Rate: A Practical and Powerful
757       Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).

758  51.  R. Vita, *et al.*, The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res.* **47**,
759       D339–D343 (2019).

760  52.  E. Van Damme, M. Van Loock, Functional annotation of human cytomegalovirus gene
761       products: an update. *Front. Microbiol.* **5** (2014).

762  53.  G. W. G. Wilkinson, *et al.*, Human cytomegalovirus: taking the strain. *Med. Microbiol.
763       Immunol. (Berl.)* **204**, 273–284 (2015).

764  54.  M. R. Schleiss, Recombinant cytomegalovirus glycoprotein B vaccine: Rethinking the
765       immunological basis of protection. *Proc. Natl. Acad. Sci.* **115**, 6110–6112 (2018).

766  55.  I. Baraniak, *et al.*, Protection from cytomegalovirus viremia following glycoprotein B
767       vaccination is not dependent on neutralizing antibodies. *Proc. Natl. Acad. Sci.* **115**, 6273–6278
768       (2018).

769  56.  C. S. Nelson, *et al.*, HCMV glycoprotein B subunit vaccine efficacy mediated by
770       nonneutralizing antibody effector functions. *Proc. Natl. Acad. Sci.* **115**, 6267–6272 (2018).

771  57.  M. Foglierini, J. Marcandalli, L. Perez, HCMV Envelope Glycoprotein Diversity Demystified.
772       *Front. Microbiol.* **10** (2019).

773  58.  G. Guo, *et al.*, Polymorphisms and features of cytomegalovirus UL144 and UL146 in
774       congenitally infected neonates with hepatic involvement. *PLoS ONE* **12**, e0171959 (2017).

775  59.  C. Berg, *et al.*, The frequency of cytomegalovirus non-ELR UL146 genotypes in neonates with
776       congenital CMV disease is comparable to strains in the background population. *BMC Infect.
777       Dis.* **21**, 386 (2021).

778  60.  J. J. C. de Vries, *et al.*, Rapid Genotyping of Cytomegalovirus in Dried Blood Spots by Multiplex
779       Real-Time PCR Assays Targeting the Envelope Glycoprotein gB and gH Genes. *J. Clin.
780       Microbiol.* **50**, 232–237 (2012).

781  61.  H.-Y. Wang, *et al.*, Common Polymorphisms in the Glycoproteins of Human Cytomegalovirus
782       and Associated Strain-Specific Immunity. *Viruses* **13**, 1106 (2021).

783  62.  F. Zavaglio, *et al.*, Detection of Genotype-Specific Antibody Responses to Glycoproteins B and
784       H in Primary and Non-Primary Human Cytomegalovirus Infections by Peptide-Based ELISA.
785       *Viruses* **13**, 399 (2021).

786  63.  W. Barrett-Muir, *et al.*, Genetic variation of varicella-zoster virus: evidence for geographical
787       separation of strains. *J. Med. Virol.* **70 Suppl 1**, S42-47 (2003).

788  64.  F. Wegner, F. Lassalle, D. P. Depledge, F. Balloux, J. Breuer, Coevolution of Sites under
789       Immune Selection Shapes Epstein–Barr Virus Population Structure. *Mol. Biol. Evol.* **36**, 2512–
790       2521 (2019).

791  65.  C. Firth, *et al.*, Using Time-Structured Data to Estimate Evolutionary Rates of Double-Stranded
792       DNA Viruses. *Mol. Biol. Evol.* **27**, 2038–2051 (2010).

793  66.  M. Bates, A. B. Brantsaeter, Human cytomegalovirus (CMV) in Africa: a neglected but
794       important pathogen. *J. Virus Erad.* **2**, 136–142 (2016).

795  67.  U. A. Gompels, *et al.*, Human Cytomegalovirus Infant Infection Adversely Affects Growth and
796       Development in Maternally HIV-Exposed and Unexposed Infants in Zambia. *Clin. Infect. Dis.*
797       **54**, 434–442 (2012).

798  68.  R. F. Pass, B. Anderson, Mother-to-Child Transmission of Cytomegalovirus and Prevention of
799       Congenital Infection. *J. Pediatr. Infect. Dis. Soc.* **3**, S2–S6 (2014).

800  69.  F. Ribeiro da Silva, The slave trade and the development of the Atlantic Africa port system,
801       1400s–1800s. *Int. J. Marit. Hist.* **29**, 138–154 (2017).

802  70.  N. Renzette, *et al.*, Limits and patterns of cytomegalovirus genomic diversity in humans. *Proc.*
803       *Natl. Acad. Sci.* **112**, E4120–E4128 (2015).

804  71.  M. Guellil, *et al.*, "Ancient herpes simplex 1 genomes reveal recent viral structure in Eurasia"
805       (Genomics, 2022) https:/doi.org/10.1101/2022.01.19.476912 (January 27, 2022).

806  72.  J. B. Plotkin, J. Dushoff, Codon bias and frequency-dependent selection on the hemagglutinin
807       epitopes of influenza A virus. *Proc. Natl. Acad. Sci.* **100**, 7152–7157 (2003).

808  73.  S. F. Elena, R. Miralles, A. Moya, Frequency-Dependent Selection in a Mammalian RNA Virus.
809       *Evolution* **51**, 984–987 (1997).

810  74.  M. Zhou, J.-M. Lanchy, B. J. Ryckman, Human Cytomegalovirus gH/gL/gO Promotes the Fusion
811       Step of Entry into All Cell Types, whereas gH/gL/UL128-131 Broadens Virus Tropism through a
812       Distinct Mechanism. *J. Virol.* **89**, 8999–9009 (2015).

813  75.  J. Kalser, *et al.*, Differences in Growth Properties among Two Human Cytomegalovirus
814       Glycoprotein O Genotypes. *Front. Microbiol.* **8**, 1609 (2017).

815  76.  S. L. Fishman, A. D. Branch, The Quasispecies Nature and Biological Implications of the
816       Hepatitis C Virus. *Infect. Genet. Evol. J. Mol. Epidemiol. Evol. Genet. Infect. Dis.* **9**, 1158–1167
817       (2009).

818  77.  A. M. Di Bisceglie, *et al.*, AGE, RACE AND VIRAL GENOTYPE ARE ASSOCIATED WITH THE
819       PREVALENCE OF HEPATITIS B E ANTIGEN IN CHILDREN AND ADULTS WITH CHRONIC
820       HEPATITIS B. *J. Viral Hepat.* **26**, 856–865 (2019).

821  78.  L. Ségurel, Z. Gao, M. Przeworski, Ancestry runs deeper than blood: The evolutionary history
822       of ABO points to cryptic variation of functional importance. *Bioessays* **35**, 862–867 (2013).

823  79.  S. John, *et al.*, Multi-antigenic human cytomegalovirus mRNA vaccines that elicit potent
824       humoral and cell-mediated immunity. *Vaccine* **36**, 1689–1699 (2018).

825  80.  P. J. A. Cock, *et al.*, Biopython: freely available Python tools for computational molecular
826       biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).

827  81.  K. Katoh, D. M. Standley, MAFFT Multiple Sequence Alignment Software Version 7:
828       Improvements in Performance and Usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).

829  82.  R. C. Edgar, MUSCLE: a multiple sequence alignment method with reduced time and space
830       complexity. *BMC Bioinformatics* **5**, 113 (2004).

831  83.  A. Larsson, AliView: a fast and lightweight alignment viewer and editor for large datasets.
832       *Bioinformatics* **30**, 3276–3278 (2014).

833  84.  K. Tamura, M. Nei, Estimation of the number of nucleotide substitutions in the control region
834       of mitochondrial DNA in humans and chimpanzees. *Mol. Biol. Evol.* **10**, 512–526 (1993).

835  85.  E. Paradis, J. Claude, K. Strimmer, APE: Analyses of Phylogenetics and Evolution in R language.
836       *Bioinformatics* **20**, 289–290 (2004).

837  86.  R Core Team, R: A language and environment for statistical computing. R Foundation for
838       Statistical Computing, Vienna, Austria. (2014).

839  87.  L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, B. Q. Minh, IQ-TREE: A Fast and Effective
840       Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol. Evol.* **32**,
841       268–274 (2015).

842  88.  Z. Yang, Maximum likelihood phylogenetic estimation from DNA sequences with variable
843       rates over sites: Approximate methods. *J. Mol. Evol.* **39**, 306–314 (1994).

844  89.  A. Rambaut, FigTree v1.4.2, A Graphical Viewer of Phylogenetic Trees (2018).

845  90.  D. H. Huson, SplitsTree: analyzing and visualizing evolutionary data. *Bioinformatics* **14**, 68–73
846       (1998).

847  91.  R. Chakraborty, M. Nei, Genetic differentiation of quantitative characters between
848       populations or species: I. Mutation and random genetic drift. *Genet. Res.* **39**, 303–314 (1982).

849  92.  R. R. Hudson, M. Slatkin, W. P. Maddison, Estimation of Levels of Gene Flow from DNA
850       Sequence Data. *Genetics* **132**, 583–589 (1992).

851  93.  D. P. Martin, B. Murrell, M. Golden, A. Khoosal, B. Muhire, RDP4: Detection and analysis of
852       recombination patterns in virus genomes. *Virus Evol.* **1** (2015).

853  94.  A. J. Page, *et al.*, SNP-sites: rapid efficient extraction of SNPs from multi-FASTA alignments.
854       *Microb. Genomics* **2** (2016).

855  95.  C. Camacho, *et al.*, BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).

856  96.  O. J. Charles, C. Venturini, J. Breuer, "cmvdrg - An R package for Human Cytomegalovirus
857       antiviral Drug Resistance Genotyping" (Bioinformatics, 2020)
858       https:/doi.org/10.1101/2020.05.15.097907 (November 13, 2020).

859

860 **Figure Legends**

861

862 **Figure 1. Circular genome map showing nucleotide diversity and multi-allelic regions.** Tracks

863 numbered from in to out: (Track1 - red) Barplot of nucleotide diversity (calculated as heterozygosity)

864 is shown as bars of Heterozygosity (red) along a 50bp moving average; (Track2)  Open Reading

865 Frames in the CMV genome are coloured by gene family as defined by bottom legend and in (26);

866 (Track3 – blue)  Multi-allelic regions as defined using HMMcluster are highlighted in translucent

867 blue; (Track4) ORF names. We also show a representative multi-sequence alignment for conserved

868 (left) and multi-allelic (right) regions.

869

870 **Figure 2. Unrooted Maximum Likelihood phylogenetic trees of representative multi-allelic and**

871 **conserved regions.** Tips were grouped if within 5% of the maximum taxa distance and are shown as

872 triangles where size indicates the number of grouped sequences and colour represents different

873 allele from HMMcluster. Small hard to see fans have been blown up and are represented by fans

874 within circles. Nodes with bootstrap support >90% are shown as red diamonds. **Note:** Scale bars

875 differ for each figure. **A)** Multi-allelic region 2 (RL5A RL6) (5 alleles). **B)** Multi-allelic region 30 (UL75)

876 (2 alleles). **C)** Example conserved region (UL105) (1 allele) of comparable alignment length.

877 Variability of C is much less than A and B with no support for HMM derived clusters**. D)** A and C when

878 drawn to the scale of B, the example conserved region tree becomes difficult to see at this

879 representation reflecting the relatively minor variation it encodes. Sequences with greater than 15%

880 ambiguous bases were removed before phylogenetic reconstruction.

881

882 **Figure 3 Multi-Dimensional Scaling of all CMV genomes.** The figure shows multidimensional scale

883 analysis for all CMV strains analysed (n=259): each dot represents a CMV strain, and the colour

884 indicates the continent of isolation (Europe includes European and Middle eastern genomes). The

885 analysis was done in three scenarios: A) whole genome, B) conserved regions (conserved

886 concatenome), C) multi-allelic regions (multi-allelic concatenome). This analysis shows an overall

887 trend for geographical segregation for the whole genome (A) and the conserved regions (C), but not

888 for the multi-allelic regions (B).

889

890 **Figure 4. Admixture analysis in CMV's conserved regions.** Admixture inferred ancestral lineages

891 reconstructed from CMV conserved concatenomes support evidence of geographic segregation in

892 CMV. The admixtures derived from a representative K=2 model of 62 sequences, were projected to

893 the remaining 197 sequences. A) This plot shows admixture proportions for whole CMV dataset

894    (n=259 strains) grouped by continent. B) The red and blue cluster components were used to colour

895    sequences in the conserved concatenome MDS. Select common reference strains have been

896    labelled.